

2ij. 1

**UNIVERSIDAD NACIONAL AUTONOMA DE MEXICO**

**ESCUELA NACIONAL DE ESTUDIOS PROFESIONALES  
A C A T L A N**



**ANALISIS ROBUSTO DE MORTALIDAD  
EN MEXICO DURANTE EL PERIODO 1922-1982.**

**Tesis Profesional**

Que para obtener el Título de  
**A C T U A R I O**

presenta

**VIVIANE DENISE BURGUNDER SANTOSCOY**

Asesor: **DR. JAIME B. CURTS GARCIA**

**EDO. DE MEX.**

**TESIS CON  
FALLA DE ORIGEN**

**1988**



## **UNAM – Dirección General de Bibliotecas Tesis Digitales Restricciones de uso**

### **DERECHOS RESERVADOS © PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL**

Todo el material contenido en esta tesis está protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

## RESUMEN

La mortalidad es una de las variables demográficas que determinan el volumen, el crecimiento y la estructura por edades y sexo de la población. Su incidencia es el reflejo de las condiciones sociales, económicas y culturales entre los habitantes del país.

Este estudio fue diseñado para:

1. Analizar las diferencias de la esperanza de vida dentro y entre las diferentes entidades federativas de la República Mexicana, en 1975.
2. Analizar el comportamiento secuencial de la mortalidad en México de 1922 a 1982.

El análisis de los datos se basó principalmente en las ideas del análisis exploratorio de datos utilizando fundamentalmente los procesos iterativos de mediana pulida y de suavización de datos.

Los resultados más significativos fueron:

1. Existen grandes disparidades entre los diferentes estados:

a) Menor esperanza de vida en los Estados de Oaxaca, Puebla, Hidalgo y Chiapas.

b) Mayor esperanza de vida en los Estados de Quintana Roo, Nuevo León, Durango y Tamaulipas.

2. Prevalece en México el patrón universal de la mortalidad:

a) Mayor esperanza de vida en el sexo femenino.

b) Altas tasas de mortalidad en el primer año de vida, que descienden rápidamente hasta llegar a valores mínimos entre los 10 y 15 años, y que incrementan de nuevo más lentamente.

3. Existen cambios en la mortalidad con respecto al tiempo:

a) De 1922 a 1942: inestable

b) De 1943 a 1947: descenso

c) De 1948 a 1951: constante

d) De 1952 a 1965: descenso

e) De 1966 a 1969: constante

f) De 1970 a 1982: descenso.

Finalmente, los resultados de este estudio sugieren que se investigue, con mayor cuidado, sobre las causas de las altas de mortalidad en algunas entidades rezagadas. Asimismo, se sugiere que se estudie más la disminución del ritmo de descenso en la mortalidad de los últimos años, a fin de que los avances futuros sean más sustanciales.

## INDICE

"Ser consciente de la propia ignorancia  
es un gran paso hacia el saber."

B. Disraeli

## INDICE

INTRODUCCION.....	1
Antecedentes del estudio	
Propósitos de estudio	
Justificación del estudio	
Estructura organizativa	
<b>I- REVISION DE LA LITERATURA.....</b>	<b>10</b>
<u>A- El análisis tradicional de la mortalidad.....</u>	<u>11</u>
1. El análisis demográfico.	
2. Panorama del estudio de la mortalidad en México.	
<u>B- El análisis robusto de datos.....</u>	<u>21</u>
1. Filosofía del "Análisis Exploratorio de Datos"	
2. Herramientas del EDA aplicadas a la demografía.	
<b>II- DESCRIPCION DE TECNICAS ROBUSTAS.....</b>	<b>31</b>
<u>A- Suavización no lineal.....</u>	<u>33</u>
1. Secuencias de datos y resúmenes suavizadores.	
2. Suavizadores no lineales	
a. Definiciones	
b. Unidades elementales	
3. Suavizadores Compuestos	
<u>B- Análisis de tablas de doble-entrada por medianas.....</u>	<u>52</u>
1. La tabla de doble-entrada.	
2. La mediana pulida.	
<u>C- Descripción de la noestacionariedad.....</u>	<u>61</u>
1. Diagnóstico de la noestacionariedad	
2. Datos suavizados en series de tiempo.	
3. Secuencias estructurales	
<u>D- Revelación gráfica de los datos.....</u>	<u>71</u>
1. El diagrama de Tallo-y-Hoja.	
2. El diagrama de caja.	
<b>III- DISEÑO EXPERIMENTAL.....</b>	<b>78</b>
<u>A- Fuentes de información.....</u>	<u>79</u>
<u>B- Los datos demográficos.....</u>	<u>80</u>
1. Tasa de mortalidad 1982-1982.	
2. Esperanza de vida, por edad exacta, según Entidad Federativa y según sexo, 1975	

<b>IV- ANALISIS E INTERPRETACION DE LOS RESULTADOS.....</b>	<b>89</b>
<u><b>A- La serie de tiempo: tasa de mortalidad.....</b></u>	<b>90</b>
1. Presentación de los resultados.	
2. Discusión de los resultados.	
<u><b>B- La tabla de doble-entrada: esperanza de vida.....</b></u>	<b>107</b>
1. Presentación de los resultados.	
2. Discusión de los resultados.	
<b>CONCLUSION.....</b>	<b>168</b>
Conclusiones y recomendaciones.	
<b>BIBLIOGRAFIA.....</b>	<b>174</b>

## INTRODUCCION

"Seguramente la estadística no puede ser ignorada por ningún investigador, aún cuando no tenga ocasión de emplearla en todos sus detalles y ramificaciones."

Oscar Bernard

## INTRODUCCION

### ANTECEDENTES DEL ESTUDIO.

La Demografía y la Estadística son disciplinas estrechamente relacionadas. Una y otra han contribuido enormemente a su mutuo desarrollo. La demografía, por su parte, ha impulsado el desarrollo y perfeccionamiento de diversas técnicas estadísticas al plantearle problemas específicos en el estudio de las características demográficas, en tanto que la estadística ha enriquecido el conjunto de los métodos y técnicas que el demógrafo utiliza en el estudio de la población.

A propósito de la relación entre la demografía y la estadística, R. Hom Chanda (1983) ha comentado que: "... se consigna el descubrimiento de que la sangre circula por el cuerpo humano, en razón de lo cual el corazón palpita. En el caso de la estadística y la demografía, se da un acontecimiento de este tipo, a partir del cual los estadísticos dicen que nació la estadística como ciencia, y los demógrafos manifiestan que se inició la demografía científica. Este origen común es la publicación, en 1662, de la obra de John Graunt, *Natural and Political Observations, Mentionned in a Following Index and Made Upon the Bills of Mortality.*"

En términos generales, la estadística se utiliza en las tres fases de la investigación demográfica: 1) la generación de información; 2) el análisis



demográfico; 3) la investigación causal.

La generación de datos, como su nombre lo indica, corresponde a la recopilación de datos demográficos a través de las fuentes clásicas: censos, encuestas, registro civil y estadísticas vitales. Aquí, la estadística no sólo proporciona las bases teóricas para el diseño de encuestas, sino también métodos efectivos de verificación y control de calidad de los datos demográficos.

La estadística es ampliamente utilizada en la segunda fase identificada con el tratamiento de datos numéricos de observación. En las tareas de evaluación, ajuste y corrección de información, se utilizan los métodos de regresión y algunos modelos estadísticos como el análisis de factores y el de componentes principales. Así, por ejemplo, la estimación de los parámetros para las tasas de mortalidad y la esperanza de vida es a través de técnicas de regresión y otros métodos del análisis confirmatorio de datos.

Finalmente, la investigación causal, que en sí constituye el fin último de la investigación demográfica, trata de examinar las causas de los fenómenos, empleando tanto las diversas técnicas estadísticas de verificación de hipótesis, como los métodos de regresión, los análisis de factores, discriminantes y componentes principales, entre otros.

Pero es preciso ver que, frecuentemente, esas investigaciones presentan algunos errores comunes. A menudo, la literatura demográfica "abunda en ingenuidades, en errores, bobadas en última instancia, atribuibles a la mala

comprensión de los datos estadísticos." (Pressat, 1970).

Un catálogo de los errores de apreciación provocados por un análisis erróneo de los datos numéricos resultaría muy valioso. Serviría para convercarse del carácter indispensable y a veces demasiado sutil del análisis demográfico. Pero una presentación suficientemente nutrida de análisis erróneos sólo podría resultar demasiado laboriosa como paso previo al examen de los métodos del análisis demográfico.

Por esta razón, se podría mencionar sólo los casos más comunes que son los análisis en los que se utiliza erróneamente la ecuación de regresión o el coeficiente de correlación para la evaluación del trabajo de investigación. Curtis (1984) afirmó que "lo anterior puede deberse a dos factores: 1) Por utilizar referencias que no contienen los avances que existen con respecto a cierta técnica o utilizar. 2) Porque existe confusión en los principios o supuestos en los que se fundamentan las técnicas cuantitativas a utilizar."

Sin embargo, se ha dado recientemente una atención especial al papel y desarrollo que juegan los métodos robustos en la estadística. El análisis exploratorio de datos es una colección de métodos gráfico-numéricos previo al clásico análisis estadístico minimizando los supuestos probabilísticos que tradicionalmente se asumen con respecto al comportamiento o la distribución de los datos. El enfoque exploratorio para el análisis de datos trata de maximizar todo lo que se pueda aprender acerca de los datos.

El desarrollo en este campo del análisis de datos se debe en gran parte a J. Tukey, quien publicó en el año de 1977 su libro llamado "Análisis Exploratorio de Datos" o EDA (siglas en Inglés Exploratory Data Analysis).

El análisis exploratorio de datos tiene como finalidad la descripción de la estructura subyacente a los datos a partir de la ecuación

$$\text{DATO} = \text{MODELO} + \text{RESIDUO}$$

y la manifestación de la forma en que un conjunto de datos se desvía de un modelo particular.

#### PROPOSITOS DEL ESTUDIO.

Este estudio está enfocada al análisis del comportamiento secuencial de tasas de mortalidad y de la esperanza de vida utilizando técnicas estadísticas robustas y del análisis exploratorio de datos.

Con el objeto de aportar elementos para el conocimiento del fenómeno de la mortalidad en México, se describen los cambios en la mortalidad de 1922 a 1982 y a fin de completar el estudio de estos cambios, se planteó la necesidad de observar el comportamiento de un indicador en las entidades federativas. Se trata entonces de estudiar lo que ocurre en las treinta y dos entidades con la esperanza de vida por sexo y en diferentes edades, en 1975.

Sin embargo, el trabajar solamente con indicadores demográficos de mortalidad

hace que se cubra una parcialidad del fenómeno, ya que no se investiga sobre los determinantes sociales, económicos y culturales de la mortalidad, ni sobre las causas de los cambios.

En esta dirección cabe señalar que esta investigación solo cubre un espacio limitado, como es el del análisis cuantitativo de los cambios de la mortalidad. Este intento constituye una parte fundamental y complementaria al de otros estudios que componen el tema de la mortalidad y que se han mencionado dentro de este estudio: análisis de las causas de muerte y su relación con factores sociales y económicos.

En base a lo anterior, cabe preguntarse en qué años se presentan tasas brutas de mortalidad altas. ¿Se atribuyen, acaso, descensos inesperados y rápidos en años específicos? ¿Por qué razones? ¿Cómo se podría interpretar la relación entre las variables -la entidad federativa, la edad y el sexo- y la variable de respuesta -la esperanza de vida-?

De esta manera, se intenta no sólo describir las diferencias dentro y entre entidades federativas por grupos, sino llevar a cabo un análisis del significado de estas diferencias y del comportamiento en el interior de las diferentes entidades. El trabajo incluye un bosquejo general de las técnicas del análisis exploratorio utilizadas para dicho fin.

## JUSTIFICACION DEL ESTUDIO

La medición de la magnitud relativa del acortamiento de las defunciones humanas ha sido un aspecto tradicional y necesariamente incluido en las investigaciones que procuran indagar el nivel de vida y las condiciones de salud de las sociedades donde éstas se producen. Por lo mismo, se han desarrollado y utilizado técnicas e indicadores que permiten cuantificar, además de la incidencia general de la mortalidad, las causas mórbidas de la misma y sus diferentes y cambiantes niveles conforme a las características biológicas, culturales, sociales y económicas de los diversos grupos de población.

Para poder cumplir con el objetivo del trabajo, sobre el análisis cuantitativo de los cambios de la mortalidad, se eligieron los indicadores de mortalidad que reflejan mejor estos cambios: la esperanza de vida por edades, por sexo y por entidad federativa en 1975; y la tasa bruta de mortalidad de 1982 a 1982.

La razón de seleccionar estos indicadores es porque la esperanza de vida, sobre todo al nacer, es un indicador sumamente sensible a los cambios del desarrollo social y económico e incluso se toma como un buen elemento para medir el grado de desarrollo de un país o de una región; la esperanza de vida en edades activas permite observar tanto las disminuciones de los riesgos que causan la muerte de los trabajadores, como el patrón de la mortalidad diferencial por sexo; la esperanza de vida en la vejez sirve para diferenciar los efectos de cambios en la mortalidad de población adulta según las transformaciones sociales y

económicas de cada entidad federativa; y por último la tasa bruta de mortalidad de 1922 a 1982, se eligió por ser el indicador resumen de las condiciones de mortalidad en cada momento de la estimación.

No obstante la importancia y la utilidad que tiene el análisis exploratorio de datos aplicada a la Demografía, son pocos los trabajos publicados, consistiendo éstos principalmente en libros recientes, y estas técnicas no se han llevado a cabo más que en Suecia, por Breckenridge (1983) quien ha analizado particularmente problemas relacionados con índices de fertilidad.

#### ESTRUCTURA ORGANIZATIVA.

En el primer capítulo, se expone un panorama del análisis tradicional de datos relacionados con la mortalidad. Se discute el problema de los casos extremos cuando se utiliza la técnica de mínimos cuadrados y, se hace referencia a la calidad de ajuste de una ecuación de regresión. En seguida, se abordan los componentes básicos de la filosofía del análisis exploratorio de datos y se comentan algunas aplicaciones del análisis robusto de datos demográficos.

El segundo capítulo aborda una descripción de las técnicas robustas, enfatizando los procesos iterativos llamados "mediana pulida" y "suavización no-lineal". En ambos casos, se describe el problema de la no-estacionaridad por lo que se explican las técnicas de transformación de datos.

En el tercer capítulo, se presentan la población y las variables del diseño experimental. Asimismo, se examina la estructura de los datos.

La información que sirve para este estudio proviene del trabajo sobre los cuadros de los cuadernos de estadísticas vitales de la SPP.

En el cuarto capítulo, se presentan el análisis e interpretación de los datos. Se reseña la situación de los datos demográficos y se expone una serie de cuadros y diagramas derivados de los cálculos de las técnicas exploratorias.

En vista de que las técnicas robustas han crecido con la aparición de diversos algoritmos computarizados, los datos de este trabajo se sistematizaron haciendo uso del programa de computadora STATGRAF.

Finalmente, se enumeran las conclusiones y recomendaciones que este trabajo pueda aportar para investigaciones futuras.

## CAPITULO I

### REVISION DE LA LITERATURA

"La demografía es el análisis estadístico de la existencia humana; por lo tanto, la investigación del estudio y los movimientos de la población, la genealogía, la eugenesia, la antropometría y la patología, concebido siempre en las condiciones de un posible examen cuantitativo."

Whipple



## I- REVISION DE LA LITERATURA

En este primer capítulo, se pretende analizar el estudio tradicional de datos demográficos, mencionando los métodos estadísticos clásicos. A parte, la primera parte de este capítulo ofrece una visión panorámica del estudio de la mortalidad en México. La segunda parte trata de dar un enfoque de la filosofía del Análisis Exploratorio de Datos y sus aplicaciones a la Demografía.

### A- EL ANÁLISIS TRADICIONAL DE LA MORTALIDAD

#### 1- El análisis demográfico.

El análisis demográfico se identifica considerablemente con el procesamiento de datos numéricos de observación. Por sus conceptos y sus modos de descripción, el análisis puede orientar por buen camino la investigación causal.

Como bien señalan Rodríguez et al. (1987), "las estadísticas demográficas son registros cuantitativos de las características biológicas y sociales de una población, referidas a un momento dado o, bien, a cambios continuos que ocurren durante un periodo de tiempo", y por ello, no es sorprendente que el objeto de estudio de las investigaciones demográficas se presten a ser contabilizadas: se está obligado a censar ya sea acontecimientos (nacimientos, divorcios, defunciones, ...), personas, según diversos estados; o bien, el tiempo tratase ésta medida por el calendario o de las diversas duraciones transcurridas desde

tal o cual suceso inicial: edad, duración de la viudez, ...

Como varias de esas medidas puedan intervenir simultáneamente, es factible que los cuadros estadísticos resultantes sean extremadamente complejos. Puede decirse que el aprendizaje empieza por el contacto con los datos cifrados, en sus aspectos multiformes; ningún trabajo fructífero será posible mientras no se alcance una familiaridad suficiente con estos datos.

A menudo, los investigadores tienen problemas para entender el conjunto de datos con sólo observarlo, porque pareciera que el conjunto de datos sólo es una lista de números. Uno de las razones por las que existen métodos estadísticos es justamente hacer que todas las características importantes que un conjunto de datos pueda contener se vuelvan claras y aparentes. Para reafirmar lo anterior, basta con citar a Siegel (1988) quien afirma que: "la estadística es el arte de hacer inferencias y sacar conclusiones a partir de datos imperfectos. Los valores de datos son imperfectos, en la medida en que comunican la información útil pero no cuentan la historia completa."

Conviene señalar, como lo hacen Rodríguez et al. (1987), que "la materia prima del análisis demográfico está constituido por los datos estadísticos", ... expresados generalmente, en números absolutos; los cuales son insuficientes para hacer comparaciones en el tiempo y en el espacio.

Por ello, una de las funciones más importantes del demógrafo es la de transformar esos datos en bruto arrojados por la observación en números

relativos; es decir, cuando las frecuencias absolutas se relacionan entre sí y se expresan en términos de razones o tasas. (La diferencia estriba según el tipo de cifras que se emplean). En aras de este estudio, se ejemplifica cuatro formas de expresar la tasa de mortalidad:

Tasa Específica de Mortalidad Según la Edad = (Número de defunciones de un grupo específico de edad / Población de dicho grupo específico de edad) \* 1000

Se considera que esta tasa es de gran utilidad, ya que sirve para medir el riesgo de muerte para cada edad o grupo de edades que se seleccione. Además, se puede calcular por edad y por sexo y referir a determinada área geográfica.

Tasa de Mortalidad Infantil = (Número de defunciones de menores de un año de edad / Número de nacidos vivos) \* 1000

Esta tasa se refiere, específicamente, a las muertes de 0 a 1 año de edad, y es igual a la tasa específica de mortalidad por edad que se refiera a ese mismo intervalo de edad.

Tasa Media de Mortalidad =  $(D^{-1} + D^0 + D^{+1}) / 3N_0$ .

En esta tasa,  $D^{-1}$ ,  $D^0$  y  $D^{+1}$  se consideran como las defunciones de tres años consecutivos;  $N_0$  es la población media de todo el intervalo de tiempo. Se considera que ésta es aplicable a la mortalidad general; pero también a la mortalidad de un grupo específico de población.

Tasa de Mortalidad por causa específica =  $(N_0 \text{ de defunciones por causa específica ocurridas durante el año} / \text{Población total a mitad del año}) \times 1000$

En rigor, esta tasa se puede considerar como "bruta" para un grupo específico, ya que, en el denominador interviene toda la población, está expuesta o no, al riesgo de la enfermedad específica que se está midiendo.

Nótese la insistencia en la importancia y las dificultades del trabajo que debe realizarse para pasar de datos en bruto a datos convenientemente elaborados. Los datos estadísticos recopilados no poseen un significado simple e inmediato; debe efectuarse un trabajo de análisis basado en conceptos y métodos específicos.

Por lo general, la importancia de esa forma de análisis, así como sus dificultades y el carácter específico de sus métodos, son desconocidos, lo cual es fuente de innumerables equivocaciones de las que constantemente da cuenta la literatura demográfica; ignorarlo equivale a resignarse a disertar sobre las apariencias, principalmente sobre las que cubren las manifestaciones en bruto de los fenómenos demográficos.

Por otro lado, para comprender, prever y controlar mejor los fenómenos, la búsqueda de sus causas es el fin fundamental de la demografía. Un buen análisis demográfico puede encaminar seriamente al investigador hacia la investigación causal. Sin embargo, los errores que se cometen, frecuentemente, son con respecto al análisis de un cierto problema por medio de una técnica cuantitativa "que se

crea es la más adecuada para el caso estudiado." (Curts, 1984).

Por ejemplo, si se considera una variable demográfica y una variable que caracterice un modo de pertenencia social, geográfica o de otro tipo; y si el objetivo de la investigación es elaborar un modelo matemático que explique el fenómeno bajo estudio, se propone un modelo de regresión lineal simple

$$y_i = B_0 + B_1x_i + e_i$$

donde  $y_i$  = variable dependiente

$x_i$  = variable independiente

$B_0$  = constante

$B_1$  = pendiente (razón de cambio de  $y_i$  con respecto a  $x_i$ )

$e_i$  = error estadístico

$i = 1, 2, 3, \dots, n$

$n$  = tamaño de la muestra.

El modelo obtenido no refleja la situación real del fenómeno, por ello el término  $e_i$  representa el error que se contempla con respecto a que "todos los valores observados caigan exactamente sobre una línea recta." Además, para poder estimar los parámetros  $B_0$  y  $B_1$  por mínimos cuadrados, es necesario seguir ciertos supuestos. Se debe tomar en cuenta que esta técnica es bastante sensible bajo la presencia de casos aberrantes. En otras palabras, los valores estimados de los parámetros pueden estar fuertemente influenciados por los casos aberrantes, y no por la mayoría del conjunto de los datos (Curts, 1984). Por lo que se puede caracterizar a los mínimos cuadrados como una técnica confirmatoria poco resistente.

Otro error que se encuentra a menudo es el de afirmar, tan sólo por el coeficiente de correlación, que el hecho de que los cambios de una segunda variable se vinculen a los cambios de la primera implica que exista entre ambas una relación de causa a efecto.

Por ejemplo, no se puede aceptar confiadamente que la disminución de la mortalidad esté relacionada con el aumento de población urbana.

Sin que jamás se insista suficientemente en este aspecto, se puede decir de una manera más sucinta que correlación no implica causalidad. De esta manera, no se está seguro de haber aislado la causa o las causas fundamentales del fenómeno: siempre puede existir algún factor, de acción más profunda, del que nada se había sospechado a priori.

Como un último ejemplo, si el experimentador desea encontrar la hipótesis nula ( $H_0: \mu_A = \mu_B$ ) haciendo uso de la prueba de "t-student" deberá recordar lo siguiente. Los métodos clásicos de inferencia que utilizan los estadísticos como  $\chi^2$ , t o F están basados en la hipótesis que las poblaciones de donde se extraen las muestras tengan una distribución Normal o Gaussiana.

Asimismo, se tiene que tener mucho cuidado en toda prueba de hipótesis porque los prejuicios que se tengan pueden afectar los resultados. Y los supuestos deben ser naturalmente tomados en cuenta cuando se interpretan los resultados del análisis. Este es uno de los primeros aspectos por observar cuando se critica un estudio o cuando se trata de entender porque "si se le dan a dos estadísticos el mismo problema se llegará a tres conclusiones diferentes".

## 2- Panorama del estudio de la mortalidad en México.

A continuación, se presenta una recopilación de estudios sobre la mortalidad en México para ofrecer una visión panorámica de los avances en su conocimiento, de las insuficiencias y limitaciones en su documentación y análisis y de las tareas y líneas de investigación que se emprendan actualmente.

La mortalidad es una de las variables demográficas que determinan el volumen, el crecimiento y la estructura por edades y sexo de la población. Su incidencia es el reflejo de las condiciones sociales, económicas y culturales que prevalecen entre los habitantes del país. La disminución de los niveles de mortalidad en cierta área responden al estancamiento o desarrollo de los servicios directa o indirectamente asociados con la salud, como infraestructura hospitalaria, servicios médicos y paramédicos, campañas de vacunación, servicios de agua potable y drenaje, alimentación, vivienda, etc... Es también esencial en investigaciones detalladas, como la elaboración de proyecciones de población, cálculo de probabilidades de contraer matrimonio, etc...

Con estas bases se fundamenta la necesidad e importancia de efectuar estudios sobre el nivel y la evolución de la mortalidad en las diferentes zonas de la República, así como investigaciones donde se vincule el comportamiento de esta variable demográfica con el desarrollo socioeconómico del país a distintos niveles geográficos. Para llevar a cabo tales estudios e investigaciones, que permitan, entre otras cosas, la elaboración de diagnósticos y la determinación de planes y programas tendientes a mejorar el nivel de vida de los mexicanos, se requiere información estadística adecuada, confiable y detallada.

En México, las estadísticas sobre defunciones generadas están basadas en la información que se recaba sobre este hecho vital en las Oficinas del Registro Civil. (Se entiende por defunción la desaparición permanente de todo signo de vida, cualquiera que sea el tiempo transcurrido desde el nacimiento con vida.)

Sin embargo, se ha podido observar que dicha información adolece de exactitud y precisión y las correcciones sobre estos datos son difíciles de emprender. A pesar del cuidado puesto para asegurar la calidad de la información recogida por enumeración y registro, "las tabulaciones finales muestran a veces indicios obvios de errores en la información básica; más o menudo sucede que los errores sólo son inferidos" (Spiegelman, 1955).

Por ejemplo, se sospecha que las defunciones no son informadas tan íntegramente en las áreas rurales como en las urbanas. Por lo tanto, la existencia de tasas de mortalidad muy bajas en algunas localidades puede indicar la persistencia de omisión en el registro de los decesos.

A pesar de ello, los registros de defunciones en México se consideran completos. El criterio de evaluación de "completo" proviene de la apreciación de las estadísticas vitales que lleva a cabo la ONU, donde se considera "completo" aquel registro que abarque por lo menos el 90% de los sucesos que se intentan medir. (Composartego y Juárez, 1977).

El primer acercamiento al fenómeno de la mortalidad en México está basado en dos indicadores: la tasa bruta de mortalidad y la esperanza de vida al nacer. El



primero, el más sencillo, es un buen indicador de la mortalidad en el tiempo, siempre que la estructura por edad de la población no haya sufrido cambios considerables. El segundo indicador, libre de cualquier tipo de estructura, refuerza el análisis de la evolución de la mortalidad en México.

La búsqueda de documentos recopilados por Alameda Bay (1975), comprende un registro de los estudios sobre la mortalidad en México en el período de 1900 a 1979 en las más importantes publicaciones periódicas de los investigadores de la mortalidad en la población mexicana y de los ficheros hemerobibliográficos de diversas instituciones académicas y asistenciales.

Los criterios adoptados por los autores de esta recopilación de estudios sobre la mortandad han sido el rigor del contenido médico y estadístico y la contribución a la descripción y análisis de la mortalidad en la población mexicana. Los temas de estudio sobre la mortalidad realizados en México, abarcan los siguientes puntos:

- La descripción y análisis de la mortalidad general y de la mortalidad por grupos de edad, sobre todo lo que concierne a la infantil, a la perinatal y a la materna.
- Los estudios de mortalidad por causas específicas.
- Los estudios sobre la esperanza de vida: tendencias y diferenciales.
- Los estudios sobre la distribución geográfica de la mortalidad mediante comparaciones internacionales, estudios migratorios, diferencias regionales y urbano-rurales.

- Los estudios sobre la mortalidad por características de la población: sexo, edad, grupo étnico, ocupación, clase social, estado marital y cohorte de nacimientos.
- Los estudios sobre la calidad de la información de la mortalidad.

Por lo general, estos análisis que se han efectuado a partir de técnicas del análisis clásico estadístico, se han basado exclusivamente en la evaluación de índices de correlación entre algún indicador de la mortalidad y varios agentes que se consideran explicativos.

## B- EL ANALISIS ROBUSTO DE DATOS.

### 1- Filosofía del análisis exploratorio de datos.

John Tukey ha realizado numerosos métodos novedosos para el análisis de datos. En 1977, escribió un libro intitulado "Exploratory Data Analysis" (técnica denominada EDA) donde hace énfasis en la exploración de los datos por métodos gráficos previo al análisis estadístico tradicional. Desde entonces, se han expuesto textos más recientes sobre la filosofía del análisis de datos lo que ha impulsado a otros autores como lo son Breckenridge (1983), Velleman y Hoaglin (1981), entre otros a desarrollar estos métodos. Particularmente, se ha prestado atención al papel que juegan los métodos gráficos en la estadística. Chambers, Cleveland, Kleiner y Tukey (1983) comentan en su libro que no existe ninguna herramienta estadística que sea tan poderosa como una gráfica bien seleccionada. El objetivo de un análisis gráfico es mostrar los datos generados por una investigación de manera clara y precisa. Una gráfica debe ser el retrato de la información cuantitativa de los datos experimentales.

Por otra parte, Curtis (1986) ha mencionado que: "La visualización de los datos permite al investigador penetrar en su estructura, minimizando los supuestos probabilísticos que tradicionalmente se asumen con respecto a su comportamiento o distribución. Lo anterior equivale a proporcionarle al investigador "una lente de aumento" que le permita:

- a) Exhibir características o patrones ocultos dentro de los datos.
- b) Resaltar con claridad la tendencia que conforman los datos.
- c) Proponer hipótesis o modelos acerca del comportamiento de los datos."

Este mismo autor ha comentado que diversos "lentes de aumento" han sido desarrollados en los últimos años para examinar de manera preliminar los datos experimentales; constituyendo éstas las herramientas fundamentales del Análisis Exploratorio de Datos.

Para lograr lo anterior, Curtis y Rodríguez (1988) mencionan que "un buen analista de datos debe ser escéptico ante cualesquier medida que resuma a uno o varios lotes de datos; es decir no debe conformarse o concluir algún aspecto acerca de los datos tan sólo porque el resumen numérico -la media aritmética, el coeficiente de correlación, etc.- sea significativo. Lo anterior es sumamente importante, ya que los resúmenes numéricos tienen la propiedad de encubrir fácilmente ciertos aspectos informativos de los mismos datos". Estos mismos autores agregan que el análisis exploratorio de datos se basa en un segundo principio: criterio abierto. Este se refiere a "la buena voluntad del analista, de siempre buscar en sus datos evidencias sobre características ocultas o patrones inesperados en los datos."

El conocimiento que ya se tenga acerca de los datos puede ayudar a ver diferentes alternativas al comienzo de la exploración y a cada paso posterior. Para empezar, ¿Las cantidades, tasas o diferencias serán la manera más apropiada para examinar los datos? ¿Los modelos serán más aparentes si los datos están

ordenados primero de alguna manera? ¿Qué otros datos pueden ocultar información valiosa? ¿Las distribuciones de frecuencia o las distribuciones acumulativas harán mejorar nuestros propósitos?

La exploración comienza sin supuestos detallados acerca de los datos o de los modelos estimados. El análisis revela una creciente descripción completa de los modelos moviendo repetidamente los datos (a menudo empezando por remover la mediana) y después se observan los residuos. Dos preguntas deben ser formuladas a cada paso: "¿Qué tan lejos se ha llegado? ¿Se puede ir más allá en el desarrollo de una descripción útil de lo que está pasando en los datos?". Y se prosigue con esta "disección" iterativa (como lo define Breckenridge, 1983) hasta que la cantidad de movimientos regulares adicionales en un paso sea insignificante. La descripción final apropiada es simplemente la combinación (a veces la suma) de todos los ajustes sucesivos. Algunas veces, por supuesto, la penetración a una etapa posterior lleva a regresar y a manejar una etapa anterior diferente.

Puesto que los datos son casi siempre el resultado de mediciones indirectas e imperfectas, y como es posible caer en el mejor ajuste, es decir, "la mejor aproximación a la curva" (Tukey, 1977); los residuos de tamaño variables son el resultado esperado de la exploración. Una característica distintiva del EDA es que nada es descartado. El análisis no se detiene en un modelo y una exposición del porcentaje de la variación total de los datos explicada por la descripción. En su lugar, los residuos no son solamente examinados para modelos posteriores en cada etapa del análisis sino que son incluso retenidos y examinados en detalle al final para ver donde y en que dirección los datos parten de la descripción.

asentada; por ejemplo, ¿en años específicos?, ¿en edades específicas?

Por otra parte, los conceptos de resistencia, residuos, re-expresión y revelación, mejor conocidos como las 4-Rs, son los componentes básicos del EDA.

El concepto de resistencia se refiere a la capacidad de un índice de localización (basado en procedimientos de ordenación y conteo) de suavizar los efectos de los cambios bruscos o arbitrarios al que un lote de datos puede estar sujeto. La media aritmética, por ejemplo, es un resumen numérico "poco resistente" bajo la presencia de casos extremos, ya que su cómputo depende de los puntajes individuales que constituyen un lote de datos. Para fines exploratorios, la mediana es, por su característica resistente, ideal para representar la mayoría de los números que integran un lote de datos.

El examen minucioso de los residuos constituye un instrumento poderoso para descubrir propiedades inherentes del lote de datos. La ecuación fundamental del análisis de datos:

$$\text{DATO} = \text{MODELO} + \text{RESIDUO}$$

considera al término RESIDUO como una medida de la variación de cada elemento  $Y_i$  con respecto al modelo o resumen numérico que se haya seleccionado para estudiar el lote de datos.

Puesto que los datos no siempre están disponibles en una forma fácil de analizar, es necesario buscar la representación más conveniente para un lote de datos. La resolución a este problema consiste en re-expresar matemáticamente la

escala original de los datos en otra que facilite su interpretación.

El concepto de revelación gráfica de los datos consiste en la inspección visual de los mismos.

Es menester insistir sobre el segundo componente. La examinación de los residuos guía a la vez la interpretación y un nuevo análisis exploratorio. A menudo esto lleva a la identificación de cambios mayores en los modelos considerados. Incluso se puede sumar a los efectos entendidos de eventos singulares. Por ejemplo, en el contexto de la tendencia de una población en el patrón de edad de fertilidad, ¿Qué efectos tiene una guerra o un período de creciente emigración en los patrones de natalidad?. La examinación de residuos puede ayudar en la identificación de desviaciones debida a errores o clasificaciones mal elaboradas y entonces estimar las correcciones apropiadas. En un EDA completo, como en todo proceso, se necesita repetición, particularmente después de haber "afinado" la expresión de los datos.

## 2- Herramientas del EDA aplicadas a la demografía.

Uno de los trabajos pioneros del EDA aplicados a la demografía fué el estudio de Breckenridge (1983) quien estudió los patrones de edad de fertilidad aplicando estos modelos. Dicho estudio respondió a una necesidad específica de hacer comparaciones de la fertilidad a través del tiempo y del espacio, necesidad de encontrar una medida de resumen de cambio que diera más información a cerca de

las dinámicas de cambio que lo que pueden dar las tasas de fertilidad arrojados en bruto. El método del EDA de Tukey se ha mostrado muy efectivo en otras áreas para detectar modelos de distribución, errores de dato y casos aberrantes. La demografía empezó a tomar ventaja sobre las fuerzas de este enfoque con la descripción de los conjuntos de datos complejos: aún con el crecimiento de disponibilidad de datos y el comportamiento individual. Las medidas de las consecuencias del comportamiento individual continúan teniendo un lugar central en el análisis demográfico por varias buenas razones:

- 1- Para identificar las tendencias y los cambios de la población.
- 2- Para proporcionar medidas generales para medir el comportamiento de la población.
- 3- Para permitir la comparación de modelos actuales con modelos históricos o con modelos de aquellas poblaciones para las cuales no es posible obtener información detallada.
- 4- Para resumir el comportamiento de la fertilidad y usarlo en proyecciones de población y en relaciones existentes entre un cambio demográfico y un cambio socio-económico.

Si se consideran sólo las medidas de fertilidad total no se puede responder de lleno a estos propósitos, porque patrones diferentes de edad en la fertilidad pueden resultar de la misma tasa de fertilidad total.

Los esfuerzos por modelar los patrones de edad de fertilidad encontraron dos dificultades: o no describieron bien un rango de escalas de fertilidad



suficientemente amplia o no proporcionaron descripciones estables, interpretables de cambio de fertilidad a través del tiempo para las poblaciones actuales. El éxito reportado en el estudio hecho por Breckenridge sobre los modelos de fertilidad en series de tiempo largos fue el resultado de la flexibilidad del comportamiento de datos en un estudio del EDR, con métodos robusto-resistentes de parámetro de estimación. Como lo dice Breckenridge en su libro "Age, Time and Fertility", el proceso introduce una doble etapa de modelaje -una futura innovación para el modelaje demográfico-.

La primera etapa trae la mayor variabilidad sistemática en un pequeño número de parámetros. Así el modelo crece fuera de los datos en vez de estar impuesto dentro de los datos. La segunda etapa estandariza la forma de los parámetros ajustados. Existen numerosas expresiones algebraicas equivalentes a los parámetros; hay que seleccionar una forma estándar guiada demográficamente para describir apropiadamente los cambios, identificar las tendencias y soluciones y comparar una serie de tiempo de fertilidad con otra.

En demografía, los estudios de las tasas se llevan a cabo en perspectivas de cohortes y en diagramas de Lexis. Esto está ilustrado en el libro de Breckenridge. Aunque la contribución primaria de su trabajo fue para modelar el cambio demográfico en series de tiempo en lugar de haber sido para estudiar la historia de la fertilidad sueca.

La investigación continua de otros caminos para expresar cambios o hechos demográficos en pequeño número de parámetros muy significantes tiene su origen en

la necesidad de simplificar comparaciones a través del tiempo y del espacio. Una medida de resumen, seleccionada con cuidado, puede también introducir una asociación del cambio demográfico con el cambio social y económico.

El observar los datos de diferentes maneras y después de nuevo, de otras maneras, ha sido demostrado repetidamente. Los problemas relacionados con datos incompletos o erróneos continúan siendo tema de estudio para nuevas metodologías.

En la apreciación de tendencias en tecnología que inciden sobre investigación demográfica, Breckenridge (1983) menciona que Winsborough (1978) hace énfasis no sólo en los avances en computación, el procesamiento de datos y los sistemas retroactivos sino también en el desarrollo de nuevos métodos de construcción del modelo, análisis de datos y estimación de parámetros. Él se refiere a los nuevos métodos en términos de nuevos estilos analíticos para manipular el volumen creciente y la complejidad de datos demográficos. También, señala que el trabajo del EDA de Tukey y Mosteller (1983), es un estudio cuya importancia para la demografía debe ser tratada desde ahora. Breckenridge afirma que "la utilidad del EDA de Tukey y Mosteller (1983) aplicado a datos defectuosos y de muy alta calidad, demuestra el valor de este enfoque flexible, de datos "guiados", afrontando una necesidad específica en análisis demográfico -la necesidad de una estable e interpretable descripción de cambio en el modelo de edad de fertilidad a través del tiempo".

Las metas alcanzadas por estos análisis exploratorios fueron útiles:

- 1- Para extraer modelos de largas series de tiempo en edad específica de fertilidad,
- 2- Para expresar estos modelos en un pequeño número de parámetros que revelen no sólo las diferencias en los modelos sino también las dinámicas de cambio en una forma demográficamente entendible, y
- 3- Para describir, suficientemente bien, las tendencias y los cambios a fin de que las desviaciones de las tendencias puedan ser examinadas en relación a eventos singulares -guerras, períodos de peligro económico, años de alta emigración- que sería expuesto a tener largos, aunque temporales, efectos en la fertilidad.

Los esfuerzos por modelar un patrón de la edad de fertilidad, no tuvieron éxito para proveer tales descripciones de experiencia de la población actual sueca. Un enfoque hacia el modelaje de fertilidad se ha esforzado por expresar la clara función de maternidad en términos de distribuciones estadísticas específicas, así como una curva normal de logaritmo, una función Beta o una función Gompertz. Breckenridge (1983) menciona que Keyfitz (1977) observa que los tests de varias distribuciones propuestas ofrecen poca satisfacción, y que Brass (1974) además de Keyfitz enfatiza que la mayoría de las distribuciones simples están lejos de describir, con exactitud o precisión un amplio rango de valores de la fertilidad.

Mientras que los esfuerzos previos probaron cómo los datos de la fertilidad de edad específica ajustan externamente patrones derivados, el enfoque del EDR permite al modelo de patrones de cambio de fertilidad tener su origen en los

datos de series de tiempo. Las primeras medidas de éxito de los procedimientos del EDA son las descripciones más adecuadas desarrolladas para las secuencias de datos de "diagrama de Lexis" y de "cohorte", y el retrato coherente de cambio demográfico expresado en los parámetros ajustados.

**CAPITULO II**  
**DESCRIPCION DE TECNICAS ROBUSTAS**

**"Un estudio de la historia de la opinión debe preceder necesariamente a la emancipación de la mente. No puedo concebir que cosa vuelve más conservador a un hombre: no conocer sino el presente, o no conocer sino el pasado."**

**Keynes**

## II - DESCRIPCION DE TECNICAS ROBUSTAS

El objetivo de este capítulo es presentar un marco teórico de dos de las técnicas del análisis exploratorio de datos, específicamente la suavización no-lineal y la mediana pulida. Ambas técnicas consisten en procesos iterativos a partir de la ecuación:

$$D = M + R$$

donde

$D$  = cualesquier dato  $Y_i$  del lote de datos en estudio.

$M$  = modelo o resumen numérico que representa en forma sintética al lote de datos en estudio.

$R$  = residuo o error que resulta de la diferencia entre  $Y_i$  y el modelo  $M$  o resumen numérico contemplado.

Y la iteración se hace al rededor de  $M$ , con el fin de "amortiguar" los cambios que el lote de datos pueda sufrir. De manera esquemática, este proceso se puede describir como:

$$\text{DATO} = \boxed{\text{MODELO}} + \text{RESIDUO}$$

donde  $\boxed{\phantom{x}}$  indica el proceso iterativo.

## A- SUAVIZADORES MULTINIALES.

### 1- Secuencias de datos u resúmenes suavizados.

En las ciencias sociales, principalmente en la demografía, se analizan datos en los que el orden de la secuencia es importante, por ejemplo, las series de tiempo; u otro tipo de variables (no necesariamente numéricas) que proporcionan también una secuencia en los datos como lo son las variables por regiones geográficas. Para describir el patrón general de la secuencia de datos, se ordenan los datos y se emplea una de las técnicas de suavización de datos; las medias corridas (también llamadas promedios móviles) es la técnica más usual. Así, por ejemplo, las tendencias demográficas a largo plazo y estacionales serán más fáciles de ver en una secuencia suave que en los datos crudos.

Cuando los valores de  $x$  son igualmente espaciados, se puede preguntar como la variable  $y$  cambia suavemente de punto en punto a lo largo del eje  $x$ , refiriéndose a los valores de  $y$  como una secuencia de datos. Cuando la secuencia proviene de registrar un valor por cada tiempo de intervalo sucesivo, los valores de  $y$  son conocidos como una serie de tiempo. Sin embargo, el orden de los valores de los datos en una secuencia no necesita estar definido por el tiempo. Las secuencias son una forma especializada de los datos  $(x, y)$  en donde los valores de  $x$  son importantes por el orden que especifican -en tiempo, en espacio, o en cualquier otro:

Para esto, se ha tratado siempre de encontrar una línea resistente como una simple descripción de la relación  $y$ -versus- $x$  y separar los datos en:

$$\text{DATO} = \text{MODELO} + \text{RESIDUO}$$

Tal separación puede ser útil aún cuando el ajuste no está descrito por una fórmula. Todo lo que se requiere es que el ajuste sea una simple descripción bien-estructura de los datos. Lo que se intenta en un ajuste simple son las curvas suaves. Cuando se trabaja a mano, se puede graficar la secuencia de los valores de  $y$  contra los valores de  $x$  correspondientes. Con esta curva se trataría de capturar el comportamiento de la secuencia de datos, a gran escala -esto es, donde la secuencia sube, donde baja, y si muestra regularidades o ciclos. Las fluctuaciones a pequeña-escala, tales como valores de datos aislados o casos aberrantes, oscilaciones rápidamente cambiantes, aparecerían después en los residuos.

Sin embargo, si se quiere un ajuste simple para ser reproducible o producido por computador, se deben de definir operaciones. Estas operaciones suavizadoras resumen usualmente, traslapando segmentos de la secuencia definidos por  $t$ - por ejemplo, los cinco primeros valores de datos, luego del segundo al sexto, y así sucesivamente. Puesto que los segmentos resumidos se traslapan, los resúmenes cambian suavemente. Por consiguiente, se le conoce como "smooth" (parte suave), en contraste a los residuos llamados "rough" (parte rugosa). El "smooth" y el "rough", como los valores de datos, son secuencias ordenadas por  $t$ . El objetivo de los métodos de datos suavizados consiste en separar la secuencia de datos en dos partes: la parte suave que varía lentamente y la parte gruesa que varía rápidamente:

$$\text{DATA} = \text{SMOOTH} + \text{ROUGH}$$

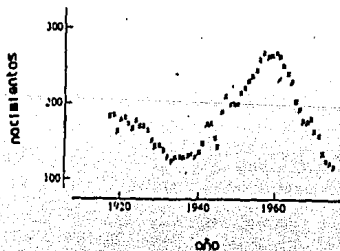
$$(1)^k$$



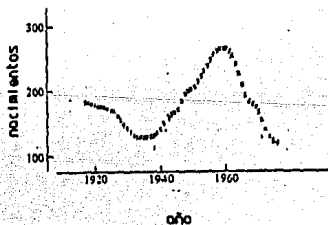
Una gran ventaja que presentan las secuencias suaves es que los datos no están restringidos por supuestos paramétricos que se asumen tradicionalmente.

Por otro lado, hay que notar que (como en las rectas ajustadas) se puede estar más interesado en los residuos, o rough, que en el modelo ajustado, o smooth. El hecho de denominar a estas técnicas "suavizadores de datos" puede fomentar a algunos análisis a olvidar la importancia de la parte rugosa, puesto que al escuchar el término "suavización de datos" se presta demasiada atención en la secuencia suave. Sin embargo, algunos aspectos de la secuencia rugosa, como los valores individuales, tienen un interés mayor que los de la secuencia suave. Esto es evidente puesto que los "rough" son simplemente los residuos de los "smooth".

Las gráficas (tomadas de Velleman y Hooglin (1981) que se presentan a continuación son muy ilustrativas al respecto. Muestran el número de nacimientos por 10,000 mujeres de 23 años, en Estados Unidos entre 1917 y 1975 y la secuencia suavizada de estos datos por 4253h, doble. Se observa una tendencia clara: la tasa de natalidad disminuye durante la depresión, aumenta a partir de la segunda Guerra Mundial y disminuye de nuevo después de 1960.

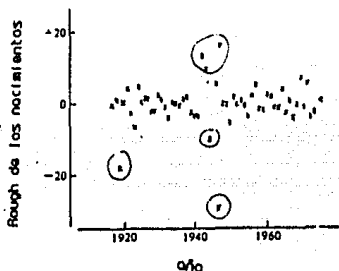


(a) Datos



(b) Smooth por 4253h, doble

Mientras que la gráfica presentada a continuación contiene información más interesante. Muestra la secuencia rugosa para estos mismos datos. La natalidad fue inestable en los años 20's, Irregular durante la segunda Guerra Mundial y nuevamente inestable en los años 60's. En los demás periodos, se presentaron sólo ligeros cambios.



(c) Rough de los nacimientos

Por otra parte, los métodos tradicionales de datos suavizados dependen de las medias móviles, que sustituyen cada valor,  $y_t$ , por un promedio pasado de  $y_t$  y los valores anteriores y posteriores a éstos en la secuencia de datos. Se seleccionan los pesos de acuerdo a la frecuencia de los datos. Para las secuencias que están formadas por frecuencias altas y bajas, los suavizadores que usan medias móviles pueden ser designados para separar los patrones de baja frecuencia de los de alta frecuencia. A veces se quisiera preservar el componente de alta frecuencia o remover el componente de baja frecuencia. Por ejemplo, una observación extrema aislada (llamada data-pico) contamina no solamente el valor suave en ese punto sino cualquier valor suave en cualquier lado en el que esté involucrado el promedio. Es decir, los datos-picos influyen en todas las frecuencias. Para eliminarlos, un suavizador debe rechazar su componente de baja frecuencia así como su componente de alta frecuencia.

Por otro lado, Tukey (1977) afirma que los suavizadores de datos basados en medianas móviles proporcionan mayor resistencia para los datos-pícos que los basados en medias móviles. Lo anterior es evidente ya que Mallows y Velleman (1980) han seguido haciendo investigaciones teóricas y han mejorado y propuesto nuevas técnicas. Recientemente se han hecho programas de cómputo para algunas de estas técnicas (Velleman y Hoaglin, 1981) que han sido incorporadas dentro de los paquetes estadísticos estándares. A continuación, se presentan algunos aspectos teóricos sobre los suavizadores no lineales.

## 2- Suavizadores no-lineales.

La propiedad fundamental de una secuencia suave es que cada valor de dato es semejante a sus vecinos; los cambios no son repentinos. Un camino sencillo para lograr esto es reemplazar cada valor de  $y$  con la mediana de tres valores de  $y$ ; el mismo, su predecesor y su sucesor. El valor de  $y$  que está fuera del paso con sus vecinos será sustituido por uno u otro de ellos, el más cercano.

### a) Definiciones.

Sea  $(y_t)$  la secuencia de datos originales donde el índice  $t$  indica las observaciones igualmente espaciadas en el tiempo (se pueden tener otras variables relacionadas con un índice que no sea el tiempo y el espaciado igual no es un requerimiento estricto). Un suavizador de datos,  $S_{ij}$ , descompone una secuencia de datos,  $(y_t)$ , en la suma de una secuencia suave llamada smooth,  $(z_t)$ , y una secuencia residual llamada rough,  $(R_t)$ . Los suavizadores producen cada valor,  $z_t$ , operando en un corto segmento de la secuencia de datos y se mueven o se recorren a lo largo de la secuencia. El número de puntos de datos en el segmento es el "span" (amplitud) del suavizador.

Los suavizadores de datos basados en medianas móviles sacrifican alguna simplicidad matemática de los métodos de suavización tradicionales. Los suavizadores lineales, tales como las medias móviles, satisfacen

$$S_a[ax + by] = aS_a[x] + bS_a[y] \quad (2)^*$$

pero esta condición obliga al suavizador a responder en los datos-picos. Los

medias móviles son miembros de la clase de los suavizadores no lineales que satisfacen una forma incompleta de (2)<sup>o</sup>, específicamente,

$$S_1[Y_t + c] = aS_1[Y_t] + c \quad (3)^*$$

Las medias móviles se combinan casi siempre con otros suavizadores lineales para mejorar la ejecución. Los suavizadores compuestos están contruidos usando dos operaciones: la resuavización y la "reroughing". La resuavización aplica un segundo suavizador,  $S_2$ , a la secuencia suave producida por  $S_1$ :

$$[Z_t] = S_2[S_1[Y_t]] \quad (4)^*$$

El "reroughing" agrega la contribución de los residuos de la primera suavización a la primera secuencia suave:

$$[Z_t] = S_1[Y_t] + S_2[Y_t - S_1[Y_t]] \quad (5)^*$$

Generalmente,  $S_1 = S_2$  en (5)<sup>o</sup>, esta combinación se refiere a un  $S_1$  doble. Sea  $R_0[Y_t] = [Y_t] - S_1[Y_t]$  el rough, entonces un estado equivalente a (5)<sup>o</sup> es

$$R_0[Z_t] = R_0[R_0[Y_t]] \quad (5')^*$$

Esta fórmula similar a (4)<sup>o</sup> expresa el fenómeno "reroughing". Los términos doble y "reroughing" fueron propuestos por Tukey (1971, 1977).

#### b) Unidades elementales de los suavizadores.

El suavizador no lineal más sencillo es la mediana móvil de la "amplitud" (span, en inglés)  $v$ , que es definido como sigue. (El número de valores de datos resumidos por cada mediana es conocido como el "span" del suavizador.) Para un span impar,  $v=2u+1$ ,

$$z_t = \text{med}(y_{t-u}, \dots, y_t, \dots, y_{t+u}) \quad (6)^*$$

Para un span par,  $v=2u$ , el resultado de cada operación de la mediana es naturalmente localizada al centro del span -entre los dos valores originales  $t$ -:

$$z_{t-1/2} = \text{med}(y_{t-u}, \dots, y_t, \dots, y_{t+u-1}) \quad (7)^{\dagger}$$

Para volver a centrar la secuencia de valores en número igual a los valores originales de  $t$ , se puede aplicar un segundo suavizador de span-par (a menudo con un span diferente) usando la definición alternativa:

$$z_{t+1/2} = \text{med}(y_{t-u+1}, \dots, y_t, \dots, y_{t+u}) \quad (7'')^{\ddagger}$$

(7) y (7'') pueden ser aplicadas en otro orden. Aplicando uno se obtienen los valores  $t=1/2, 3/2, 5/2, \dots$ ; Aplicando el otro al resultado de la primera suavización se obtienen de nuevo los valores  $t=1, 2, 3, \dots$ . Al volver a centrar, cuando el segundo suavizador es una mediana móvil de span-2 (equivalente a una media móvil de span-2), se pueden combinar las dos etapas como:

$$z_t = 1/2 \text{med}(y_{t-u}, \dots, y_t, \dots, y_{t+u-1}) + 1/2 \text{med}(y_{t-u+1}, \dots, y_t, \dots, y_{t+u})$$

que utiliza  $v+1$  valores de datos consecutivos pero da el primer y el último  $1/2$  peso.

Para las medianas móviles de span impar se selecciona el valor del suavizador en algún punto  $t$  desde el conjunto de valores de datos cercano a  $t$ . Se obtienen secuencias monótonas invariables con saltos repentinos resistentes a los picos.

Puesto que las medianas de tres no pueden ser correctas para dos casos aberrantes o datos discrepantes (outliers, en inglés), se pueden tomar más puntos para la mediana; se puede basar cada mediana en el smooth sobre cinco puntos en vez de tres, observando los dos puntos más cercanos y los dos puntos más lejanos al valor de  $y$  que está siendo modificado. Estos métodos son ejemplos de los

suavizadores de mediana-corrída, así llamados porque se corre a lo largo de la secuencia de datos y se encuentra la mediana de los tres o los cinco valores de datos más cercanos a cada punto.

Para medianas de tres, el valor de dato inicial en la secuencia plantea un problema cuando éste no está a la mitad de los tres valores de datos. Por ahora, éste se copia sin ninguna modificación. Es evidente, lo mismo sucede para el valor de dato final, y éste también se copia para el "smooth". Para las medianas de cinco, los dos valores de datos de cada extremo de la secuencia son difíciles de suavizar. Se copian los dos valores finales, pero se usa una mediana de tres para suavizar el segundo y el penúltimo valor.

Cada uno de estos suavizadores de medianas-corrídas pueden ser calculados fácilmente a mano, pero ambas son justamente pesadas en sus efectos en secuencias de datos. Las medianas-corrídas de cuatro valores de datos consecutivos son levemente suaves. A diferencia de los suavizadores que seleccionan el valor de dato medio de tres o cinco, una mediana de cuatro valores ignora el mayor y el menor valor en cada segmento de cuatro y promedia los dos valores de en medio. Nótese que los valores seleccionados para promediar están a la mitad de la medida en el sentido que los valores de  $y$  caen en medio de los demás valores de  $y$ . No necesitan ser la mitad de los dos valores acordando el orden definido por  $t$ - claro, no necesitan ser puntos de datos consecutivos en la secuencia.

Cuando se utilizan medianas-corrídas de peso-igual, se deben promediar los

valores de  $t$  también. La mediana de un segmento de peso-impair de la secuencia de datos es naturalmente inscrita a la mitad del valor de  $t$ , así se registra la mediana en la laguna entre las dos mitades de los valores de  $t$ . Un par de medianas flanquea cada valor de  $t$ . Se puede alinear un nuevo valor de  $y$  con un valor de  $t$  original promediando las medianas corridas en otro lado. Es posible ilustrar esta operación como:

valores de datos:     ...  $y_5$     $y_6$     $y_7$     $y_8$     $y_9$     $y_{10}$  ...  
 suavizados por 4's: ...  $z_{4.5}$   $z_{5.5}$   $z_{6.5}$   $z_{7.5}$   $z_{8.5}$   $z_{9.5}$  ...  
 recentrados por pares: ...  $z_5$     $z_6$     $z_7$     $z_8$     $z_9$     $z_{10}$  ...

Después de haber suavizado el resto de la secuencia, se puede modificar los valores finales en vez de copiarlos. Más adelante, se propone un método para el tratamiento de los valores finales de la secuencia.

Evidentemente, la medida de recentrado es una mediana corrida de dos porque la mediana de dos números es también su promedio. Algebráicamente, una mediana corrida de cuatro, recentrada con una mediana corrida de dos, reemplaza el valor de dato  $y_t$  por:

$$z_t = 1/2(\text{med}(y_{t-2}, y_{t-1}, y_t, y_{t+1}) + \text{med}(y_{t-1}, y_t, y_{t+1}, y_{t+2})).$$

Esta ecuación utiliza cinco valores de datos,  $y_{t-2}$  hasta  $y_{t+2}$ , pero el primer y último valor aparecen en uno de los dos segmentos cuyas medianas son promediadas, y así tienen casi la mitad del efecto de cualquiera de los demás puntos.

Hasta ahora, se han examinado suavizadores con  $\text{span}^2$  de 2, 3, 4 y 5. Los



medianas suavizadoras con grandes "spans" pueden resistir más en los extremos. Así, una mediana de "span-2" no afectará ningún punto extraordinario. Un suavizador con una mediana de "span-3" o un "span-4" no afectará los extremos. Una mediana de "span-3" seguirá un par extremo, pero una mediana de "span-4" reducirá fuertemente el tamaño de 2-puntos pico a la mitad. Una mediana de "span-5" será completamente resistente en 2-puntos pico. Las secuencias suaves resultantes son similares pero difieren en la medida que las secuencias suavizadas con medianas de cinco son más suaves pero menos parecidas a la secuencia de datos original.

#### NOTACION

El símbolo para una mediana móvil es el dígito correspondiente a su "span". Estas medianas móviles de tres están denotadas 3. A fin de proporcionar una notación compacta para las operaciones elementales de suavización, se hace referencia a ellas con nombres de carácter único. El nombre para una mediana corrida es el dígito simple correspondiente a su "span", tal como 3 o 5. Cuando una mediana corrida de "span-4" es seguida por una operación de par-promediado para recentrar los resultados, se utiliza la notación 42. El nombre de dos-dígitos es apropiado porque involucra dos operaciones. En realidad, unas combinaciones poco sofisticadas insertan otras operaciones elementales entre un 4 y un 2. Ya que es raro usar medianas móviles de más de 7 puntos, hay una pequeña posibilidad de confundir 42 con una mediana móvil de 42 valores de datos. La concatenación de nombres de un-carácter se explica en seguida, donde se combinan operaciones elementales de suavización en orden para obtener una mejor representación.

### "HANNING"

Se puede querer obtener una operación de suavización aún más suave que 42. Para esto, se utiliza un promedio pesado corrido que es útil para suavizar secuencias de datos sustituyendo cada valor de dato por un promedio de los valores de datos alrededor de éste. A veces los valores de datos son multiplicados en cada operación promediada por pesos. Así, por ejemplo, se puede reemplazar  $y_t$  por

$$z_t = 1/4y_{t-1} + 1/2y_t + 1/4y_{t+1}.$$

Un suavizador lineal simple se incorpora casi siempre dentro los suavizadores no lineales compuestos para suavizar los bordes del rough que las medianas móviles pueden dejar. Para retener la resistencia en los picos, las operaciones lineales serían aplicadas solamente después de que una operación no lineal haya removido los valores de las afueras. Se utilizan las medias móviles de 'span' corto cuando se desea evitar una gran cantidad de pasos. Un número ilimitado de promedios pesados corridos son posibles, pero uno mismo se limita a esta fórmula particular para explorar más los datos. El más usual es la media móvil pesada en tres puntos  $1/4, 1/2, 1/4$  (equivalente a 2 resuavizada por 2) (la suma de estos pesos debe ser 1). Este suavizador es llamado hanning, por Julius von Hann; quien advocó su uso y es denotado por H. Ningún promedio pesado corrido estará afectado por un extremo, así se usan generalmente tales suavizadores sólo después de que los extremos hayan sido suavizados por un suavizador de mediana corrida.

### SUAVIZANDO LOS PUNTOS FINALES

Hasta ahora se ha dado poco de como suavizar los valores inicial y final de la secuencia de datos. No se pueden suavizar estos valores de la misma manera que se suavizan los demás porque no están rodeados por otros valores. Con un suavizador de "span-largo" como 5, se puede impedir el problema encontrando medianas de "span-corto" cerca de los puntos-finales.

Para los primeros y últimos puntos de la secuencia, el "span" de la mediana móvil debe ser menor, por ejemplo, para medianas móviles de 5:

$$z_3 = \text{med}(v_1, v_2, v_3, v_4, v_5)$$

$$z_2 = \text{med}(v_1, v_2, v_3)$$

$$z_1 = v_1$$

Para los primeros y últimos valores, denotada por  $\epsilon$ , la fórmula siguiente puede ser aplicada:

$$z_1 = \text{med}(3z_2 - 2z_3, v_1, z_2) \quad (8)^\dagger$$

Hay que notar que:

$$3z_2 - 2z_3 = z_2 - 2(z_3 - z_2) = z_0$$

predice el valor  $z_0$  por una extrapolación lineal desde los dos valores smooth más a la izquierda,  $z_2$  y  $z_3$ , hasta  $t=0$ . Este valor sirve para suavizar enseguida el valor de dato  $t=1$ . Similarmente, del otro lado,  $\epsilon$  utiliza los dos valores smooth más a la derecha.

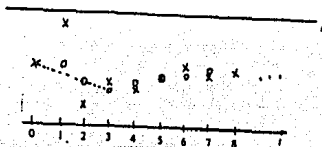
Así, para medianas corridas de cinco, se toman medianas de tres para el

segundo y el penúltimo valor:

$$z_2 = \text{med}(v_1, v_2, v_3)$$

$$z_{n-1} = \text{med}(v_{n-2}, v_{n-1}, v_n)$$

Los valores finales,  $z_1$  y  $z_n$ , requieren de un enfoque diferente. Hasta aquí, se ha dicho que se copie tal cual el valor final -esto es, para usar los valores finales sin cambiarlos. Se puede mejorar esto extrapolando los valores suavizados cerca del final. Primero se estima qué próximo valor pasa, el valor final pudo haber sido. No se puede usar el mismo valor como valor final en esta estimación porque no se ha suavizado esto todavía. Un enfoque simple es encontrar la línea recta que pasa a través de los segundos y los terceros valores suavizados desde el final y trazar el punto estimado en esta línea en el valor- $t$  que habría ocupado (ver gráfica 1).



Gráfica 1

Extrapolación del punto final

- x puntos de los datos
- o valores suavizados
- x valor extrapolado

Para valores igualmente espaciados con espaciado-t  $D_t$ , la línea en el final bajo tiene por pendiente:

$$(z_3 - z_2)/D_t$$

Se están extrapolando dos intervalos-t más allá de  $z_2$ , así el valor estimado es

$$\begin{aligned} \hat{v}_0 &= z_2 - 2D_t(z_3 - z_2)/D_t \\ &= 3z_2 - 2z_3 \end{aligned}$$

donde las  $z$ 's son los valores ya estimados. Similarmente, para el punto final se estima el punto siguiente como:

$$\hat{v}_{n+1} = 3z_{n-1} - 2z_n.$$

Después, se encuentra la mediana del punto extrapolado, el punto-final observado, y el punto suavizado siguiente al final:

$$z_1 = \text{med}(\hat{v}_0, v_1, z_2)$$

$$z_n = \text{med}(\hat{v}_{n+1}, v_n, z_{n-1}).$$

Así, si se denota a esta operación por  $\epsilon$ , se pudo usar 4253EH, doble.

El suavizador 42 tiene un problema de valor-final adicional porque necesita volver a centrar el resultado de la primer suavización. Cuando se suaviza por medianas corridas de cuatro, se obtiene una secuencia un punto más largo que la secuencia de datos originales. Se puede denotar esta secuencia larga por  $z_{1/2}, z_{1-1/2}, \dots, z_{n+1/2}$ . Aquí, se copian los valores finales:  $z_{1/2} = v_1$ ,  $z_{n+1/2} = v_n$ . Los próximos valores de cada final son medianas de dos:  $z_{1-1/2} = \text{med}(v_1, v_2)$ ,  $z_{n-1/2} = \text{med}(v_{n-1}, v_n)$ . La re-centrada subsiguiente por medianas corridas de 2 devuelve la secuencia a su largo original. De nuevo, los valores finales son copiados:

$z_1 = z_{1/p} (= y_1)$ ,  $z_n = z_{n+1/p} (= y_n)$ . Todos los demás valores son promedios de valores adyacentes; por ejemplo,  $z_2 = \text{med}(z_{1-1/p}, z_{2-1/p}) = (z_{1-1/p} + z_{2-1/p})/2$ .

### 3- Suavizadores Compuestos.

Mientras que las medianas corridas simples suavizan una secuencia de datos y pueden resistir valores de datos extraordinarios ocasionales, las secuencias suaves que ellas producen pueden describir los datos sólo crudamente o toscamente; Se puede demostrar en la descripción -obteniendo los suavizadores de datos cuyas secuencias suaves son casi iguales a los datos sin perder su suavidad- hasta una combinación juiciosa de los procedimientos de suavización.

Se han propuesto ya varias combinaciones diferentes de operaciones de suavización no lineal elementales. El suavizador compuesto 53H, doble fué propuesto por Tukey (1971) y usado con éxito por Beaton y Tukey (1974). Cleveland, Dunn, y Terpenning (1979) utilizan suavizadores relacionados con el proceso de ajuste estacional para series de tiempo en problemas de Economía y Demografía. Valleman (1975, 1980) propuso el suavizador 4253H, doble por su mejor funcionamiento en la práctica y su estabilidad en presencia de secuencias de datos altamente estructuradas. Esto empieza con una mediana corrida de 4 recentrado con una mediana corrida de 2 (equivalente a un promedio móvil de 2). Después, el "smooth" es resuavizado por medianas de 5, por medianas de 3, y finalmente por "hanning". Se aplica después, el proceso entero a los residuos y los residuos suavizados resultantes se agregan al "smooth": ésta es la "doble" operación.

Ahora se presenta una pequeña teoría para guiar la elección y la combinación de los suavizadores no lineales, propuestas por Velleman (1982). En 1980, proporcionó una guía y reportó los resultados de experimentos con varias alternativas. Estos experimentos y consideraciones teóricas implican la preferencia de 4253H, doble entre los suavizadores compuestos investigados hasta la fecha para uso general.

### Resuavización

El aplicar un suavizador a los resultados del suavizador previo es conocido como resuavización. Como con el nombre 42, se denotan tales series de operaciones concatenando sus nombres de un carácter. Si se trabaja a mano, se puede escoger a usar sólo medianas corridas de 3 y resuavizar repetidamente hasta que la nueva resuavización no cambie más los resultados. Se denota esta combinación repetida por 3R.

### "Roughing"

Generalmente, los suavizadores de medianas-corridas suavizan demasiado una secuencia de datos; Resueven patrones interesantes. Una operación complementaria puede ser utilizada para recobrar patrones suaves a partir de los residuos -esto es, a partir de la parte llamada "rough" en la fórmula

$$\text{Data} = \text{smooth} + \text{rough}$$

Se suaviza la secuencia "rough" y se agrega el resultado a la secuencia "smooth". Se espera que los patrones que hayan sido suavizados lejos por el primer paso de la resuavización pueda ser recobrada desde el "rough" y usada para hacer que el "smooth" se parezca a la secuencia de datos original. Por analogía a la resuavización, esta operación se llama "reroughing".

Se usa casi siempre el mismo suavizador en la suavización y la "reroughing", y a éste se le conoce como suavizador doble.

El "reroughing" es un ejemplo de una operación que funda en las técnicas exploratorias que "pulén" un ajuste. En una línea resistente, el paso "reroughing" ajusta una línea a los residuos, agregando esta línea al ajuste.

#### 4253H

Los suavizadores compuestos combinan casi siempre diversos suavizadores elementales por resuavización y "reroughing". Los primeros pasos en un suavizador compuesto se concentran en la protección desde los extremos en la secuencia de datos. Los pasos tardíos de la resuavización pueden entonces emplear un promedio pesado corrido. Curiosamente, las medianas corridas de 3 ó 5 pueden alterar algunas secuencias rápidamente oscilantes extrañamente. Por ejemplo, la secuencia infinita  $\dots, +1, -1, +1, -1, +1, -1, \dots$  no es modificada del todo por una mediana corrida de "span=5", aunque la secuencia oscila rápidamente. Más extraño aún, una mediana corrida de "span=3" invertirá la secuencia, como si cada valor estuviera multiplicado por  $-1$ . Así, las medianas corridas de "span-par" son



preferidas a veces -especialmente cuando una computadora está disponible para hacer el cómputo de todos los promedios que ellas requieren.

Consideraciones similares surgen en la "reroughing" porque el "rough", por su diseño, contiene picos reflejando los extremos presentes en los datos originales y oscilan generalmente rápidamente. Además, los suavizadores aplicados al "rough" deben ser también resistentes a estas características.

Una combinación de suavizadores que parece mejorar perfectamente es 4253H. Empezamos con una mediana corrida de cuatro, 4 recentrada por 2. Esto después resuaviza por 5, por 3, y finalmente -ahora que los extremos han sido suavizados fuera- por H. El resultado de esta suavización es a menudo "reroughed" -o pulido- por residuos de cómputo, aplicándoles el mismo suavizador y agregando el resultado al "smooth" del primer paso. Esto produce el suavizador completo, 4253H, doble.

## R- ANALISIS DE TABLAS DE DOBLE-ENTRADA POR MEDIANAS.

La tabla de doble-entrada es un tipo de estructura de datos útil por sus diversos campos de aplicación a diferentes niveles. Cada uno de los dos factores varía regular y separadamente uno del otro, y el valor de la variable respuesta es observada en cada combinación. Por ejemplo, el índice de mortalidad infantil puede variar en función de la región geográfica y del nivel de instrucción de los padres.

Se puede utilizar la mediana como una simple moda iterativa para asegurar un análisis resistente para los datos presentados de esta forma. A partir de la técnica exploratoria es más fácil ver alguna dependencia adicional entre la variable respuesta y los dos factores, además de que los residuales proporcionan directamente más información cuando provienen de un análisis resistente. En este trabajo, se discute el uso de las medianas y otros métodos del análisis de tablas de doble-entrada.

A continuación, se hace una pequeña revisión de la estructura natural de una tabla de doble-entrada, y después se explica e ilustra el procedimiento básico iterativo llamado "mediana pulida".

### 1- La Tabla De Doble-Entrada.

La tabla de dos-entradas es un conjunto de datos en el que las observaciones se escriben:

$$v_{ij} \quad (i=1, \dots, I; \quad j=1, \dots, J) \quad (11)''$$

y se colocan en una tabla rectangular.

	1	2	...	J
1	$v_{11}$	$v_{12}$	...	$v_{1J}$
2	$v_{21}$	$v_{22}$	...	$v_{2J}$
.	.	.	.	.
.	.	.	.	.
.	.	.	.	.
I	$v_{I1}$	$v_{I2}$	...	$v_{IJ}$

De esta manera, se pueden apreciar tres variables en la tabla: el factor renglón con  $I$  niveles; el factor columna con  $J$  niveles; y la respuesta  $y$ , en la que se tienen  $I \cdot J$  observaciones, una para cada combinación entre el renglón y la columna. (Sólo  $y$  tiene que ser numérico, aunque el renglón y/o la columna puedan ser etiquetados numéricamente ...)

#### Los modelos aditivos.

La relación entre la variable respuesta y los dos factores es fácil de resumir e interpretar si la contribución conjunta es la suma de la contribución separada de cada factor.

El modelo aditivo simple puede ser escrito formalmente como:

$$V_{ij} = \mu + A_i + B_j + e_{ij} \quad (2)^\circ$$

En este modelo,  $\mu$  es un valor típico (promedio) de toda la tabla; se le denomina "valor común". La contribución incrementada del nivel  $i$  del factor renglón, relativa al valor promedio, es  $A_i$  - a "efecto renglón". Similarmente, la contribución incrementada del nivel  $j$  del factor columna es  $B_j$  - a "efecto columna". Finalmente,  $e_{ij}$  representa la desviación de  $V_{ij}$  para el modelo aditivo de  $\mu + A_i + B_j$ ; se considera  $e_{ij}$  como una fluctuación aleatoria.

Todo análisis de tablas de doble-entrada correspondiente a la ecuación (2)<sup>o</sup> descompone el orden del dato observado  $V_{ij}$  en cuatro términos adicionales: uno que es constante en la tabla completa, uno que es constante por renglones, otro que es constante por columnas, y otro que contiene a los residuos. Se pueden especificar varias formas diferentes para lograr esto; cada una permite una técnica de análisis. Aquí, se trata de las técnicas de análisis resistentes, entonces los errores aislados en un pequeño número de celdas no afectan el valor común, los efectos-renglón, o los efectos-columna, y, por consecuencia, se verán reflejados en los residuos.

## 2- La mediana cuilida.

Para obtener el valor aditivo en la forma de la ecuación (2), se puede operar

iterativamente con los datos de la tabla, calculando y restando las medianas por rengiones y las medianas por columnas. Por ejemplo, se puede empezar por los rengiones, calculando la mediana de cada renglón y luego restando este valor a cada observación del renglón. Después, se puede continuar con las columnas de la tabla resultante, encontrando la mediana de cada una y restándola a los valores de la columna correspondiente. Por supuesto, si el renglón o la columna tiene una mediana igual a cero, entonces no habrá ningún cambio en ese renglón o en esa columna. Por principio, se puede continuar repitiendo el proceso de resta de las medianas -llamado "mediana pulida"- hasta que todos los rengiones y columnas tengan mediana cero. Esto quiere decir que, habiendo empezado por los rengiones y luego puliendo las columnas, es necesario verificar los rengiones y pulir de nuevo lo que ahora tenga mediana diferente de cero. Y así el proceso iterativo podría continuar.

*Una presentación formal.*

Para formalizar este proceso, se introduce la siguiente notación algebraica.

El ajuste aditivo análogo a (2) es

$$v_{ij} = m + a_i + b_j + e_{ij}. \quad (3)$$

Se denota por

$$v_{ij} = m(n) + a_i(n) + b_j(n) + e_{ij}(n). \quad (4)$$

el ajuste y los residuos al cabo de  $n$  iteraciones.

Puesto que se está describiendo un proceso iterativo, se plantean dos condiciones iniciales, antes de la primera iteración:

$$\begin{aligned} a(0) &= 0 \\ a_j(0) &= 0 \quad j=1, \dots, I, \quad (5) \\ b_j(0) &= 0 \quad j=1, \dots, J. \end{aligned}$$

Como se mencionó anteriormente, en esta presentación se empieza a trabajar por los renglones y luego por las columnas. (Si se trabajara a la inversa: primero las columnas y luego los renglones, el resultado no necesariamente es el mismo, pero la diferencia entre las dos soluciones es generalmente poca importante).

En seguida, se da una lista de nueve ecuaciones, numeradas de (6a) a (6i). Parecen ser complicadas por la notación matemática, pero no resulta difícil comprender su estructura. Las tres primeras ecuaciones describen la mediana pulida de los renglones, incluyendo el renglón de los efectos-columna y el nuevo valor de los residuos. Las tres ecuaciones siguientes dan el pulido correspondiente por columnas. Las tres últimas ecuaciones proporcionan el valor común, los efectos-renglón y los efectos-columna. Las nueve ecuaciones preparan conjuntamente la tabla para la próxima iteración.

El símbolo  $\Delta$  representa un cambio, y se asume que  $n$  es 0, inicialmente. Los

pasos en una iteración son los siguientes:

**Regiones:**

$$\Delta a_j^{(n)} = \text{med}\{e_j^{(n-1)} | j = 1, \dots, J\}; \quad I = 1, \dots, I; \quad (6a)^o$$

$$\Delta m_i^{(n)} = \text{med}\{b_j^{(n-1)} | j = 1, \dots, J\}; \quad (6b)^o$$

$$d_{ij}^{(n)} = e_j^{(n-1)} - \Delta a_j^{(n)}; \quad j = 1, \dots, J; \quad I = 1, \dots, I; \quad (6c)^o$$

**Columnas:**

$$\Delta b_j^{(n)} = \text{med}\{d_{ij}^{(n)} | I = 1, \dots, I\}; \quad j = 1, \dots, J; \quad (6d)^o$$

$$\Delta m_i^{(n)} = \text{med}\{a_j^{(n-1)} + \Delta a_j^{(n)} | I = 1, \dots, I\}; \quad (6e)^o$$

$$e_{ij}^{(n)} = d_{ij}^{(n)} - \Delta b_j^{(n)}; \quad I = 1, \dots, I; \quad j = 1, \dots, J; \quad (6f)^o$$

**Valor Cédula y Efectos:**

$$m_i^{(n)} = m_i^{(n-1)} + \Delta m_i^{(n)} + \Delta m_i^{(n)}; \quad (6g)^o$$

$$a_j^{(n)} = a_j^{(n-1)} + \Delta a_j^{(n)} - \Delta m_i^{(n)}; \quad I = 1, \dots, I; \quad (6h)^o$$

$$b_j^{(n)} = b_j^{(n-1)} - \Delta m_i^{(n)} + \Delta b_j^{(n)}; \quad j = 1, \dots, J. \quad (6i)^o$$

En la práctica, los pasos dados por las nueve ecuaciones son muy confiables para llevar a cabo a mano en pequeñas tablas de doble-entrada si se observa cuidadosamente el cálculo de los residuos  $e_{ij}^{(n-1)}$ . Para facilitar la contabilidad, se dan a continuación dos tablas de los pasos por regiones y por columnas. Las tablas Ia y Ib muestran las versiones esquemáticas de estas tablas.

Tabla 1a

Mediana pulida por renglón  
en la iteración n

i	j			D. med.	prev.
	1	...	J		
1	$a_{11}^{(n-1)}$	...	$a_{1J}^{(n-1)}$	$\{\Delta a_1^{(n)}\}$	$a_1^{(n-1)}$
⋮	⋮	⋮	⋮	⋮	⋮
J	$a_{J1}^{(n-1)}$	...	$a_{JJ}^{(n-1)}$	$\{\Delta a_J^{(n)}\}$	$a_J^{(n-1)}$
prev.	$b_1^{(n-1)}$	...	$b_J^{(n-1)}$	$\{\Delta m_J^{(n)}\}$	$m^{(n-1)}$

Tabla 1b

Mediana pulida por columna  
en la iteración n  
(Resultado inmediato después  
de la tabla 1a)

i	j			prev.
	1	...	J	
1	$d_{11}^{(n)}$	...	$d_{1J}^{(n)}$	$a_1^{(n-1)} + \Delta a_1^{(n)}$
⋮	⋮	⋮	⋮	⋮
J	$d_{J1}^{(n)}$	...	$d_{JJ}^{(n)}$	$a_J^{(n-1)} + \Delta a_J^{(n)}$
D. med.	$\{\Delta b_1^{(n)}\}$	...	$\{\Delta b_J^{(n)}\}$	$\{\Delta m_0^{(n)}\}$
prev.	$b_1^{(n-1)} - \Delta m_1^{(n)}$	...	$b_J^{(n-1)} - \Delta m_J^{(n)}$	$m^{(n-1)} + \Delta m_0^{(n)}$

\* tomado de Emerson y Hoaglin (1983)



Estudiando estas dos tablas, se ve que la tabla la muestra los resultados de las ecuaciones  $(6a)^2$  y  $(6b)^2$ , mientras que la tabla lb muestra los resultados de las ecuaciones  $(6c)^0$ ,  $(6d)^n$  y  $(6e)^2$ . Continuando con las ecuaciones  $(6f)^0$  a  $(6i)^2$  se llega a la tabla la en la que  $n+1$  reemplaza a  $n$  y la columna denominada "med" omitida.

La columna y el renglón denominados "prev" en las tablas la y lb forman parte del ajuste previo que, en el último paso, se transforman en efectos-columna y efectos-renglón. Como lo indican las ecuaciones  $(6b)^2$  y  $(6e)^2$ , se opera paralelamente con los renglones y las columnas de los residuos (parciales). Cuando se trabaja a mano, esto es conveniente, pero puede dejar un paso del ajuste para las  $b_j(n)$  para ser hecho después de terminar las iteraciones. Las ecuaciones  $(6e)^2$  y  $(6h)$  centran el  $b_j(n)$  de modo que su mediana sea 0.

A veces el significado de un conjunto de versiones (o ambas) es tal que "mediana = 0" no sea el centrado natural. Cuando algún otro centrado es más natural, se usa éste.

El orden de ejecución de las ecuaciones  $(6b)^2$  y  $(6i)^2$  significa que la versión de la ecuación  $(6b)^2$  puede tener que seguir la ecuación  $(6i)$  y también cambiar  $b_j(n)$  y  $m(n)$  si se tiene el mismo centrado para  $b_j(n)$ . Se podría conservar equivalentemente todos los cálculos centrados para la última etapa, y el algoritmo de cómputo presentado por Velleman y Hoaglin (1981) sigue este enfoque.

Utilizando  $n$  en las iteraciones, las fórmulas pueden hacer aparecer lo que se contempla en un gran número de iteraciones, pero ésta no es así. La versión "estandar" de este análisis basado en medianas utiliza  $n=2$ . En algunos casos, no todas las medianas de los ranglones y de las columnas de  $e_{ij}(2)$  pueden ser 0, pero esto no es de gran importancia. A menudo sucede que las medianas por renglón y las medianas por columna de los residuos se vuelven cero antes de esta etapa.

El uso de la técnica de la mediana pulida para modelos apropiados de las tablas de doble-entrada es un desarrollo relativamente reciente. Junto con otras técnicas exploratorias que resisten la influencia de los pocos valores de datos "malos", la mediana pulida fue descrita por Tukey (1970). Bood (1950) y Brown y Bood (1951) utilizan el mismo proceso como parte de un procedimiento para probar las interacciones en las tablas de doble-entrada con varias observaciones por celda. Existe literatura más reciente, incluyendo los libros de Tukey (1977), Mosteller y Tukey (1977), y Velleman y Hooglin (1981), Hooglin, Mosteller y Tukey (1983), donde se describe la técnica de la mediana pulida y en la que se ilustra ésta usando conjuntos de datos reales.

Como ya se mencionó, la mediana pulida para el análisis de tablas de doble-entrada es una técnica resistente en la que los errores aislados se ven reflejados en los residuos. Por esta razón, se adopta la mediana pulida para el análisis exploratorio en vez de utilizar las medias para un análisis correspondiente que forma las bases para el análisis de varianza clásico, en el que se asume que los  $e_{ij}$  son fluctuaciones o errores independientes con una distribución Gaussiana teniendo media cero y varianza común.

### C- DESCRIPCIÓN DE LA NO-ESTACIONARIDAD.

Después de haber definido los suavizadores no lineales y comprendido sus propiedades, se presentan a continuación dos nuevas técnicas empleando suavizadores no lineales. La primera de éstas provee un diagnóstico del "rough" en donde resulta mejor analizar una secuencia de datos después de una re-expresión de los datos. La segunda combina los métodos de suavización no lineales con otros métodos exploratorios para investigar secuencias de datos estructuradas con ciclos conocidos. (¿Cómo detectar si el modelo es no-aditivo?)

#### 1- Diagnóstico de la no-estacionaridad.

Cuando una secuencia de datos exhibe una fuerte tendencia o en otras circunstancias tiene segmentos con valores generalmente altos y segmentos con valores generalmente bajos, no es poco común encontrar que la variabilidad de valores de datos individuales (como se indicó en la variabilidad de sus vecinos inmediatos en la secuencia) cambia en esta forma, esto se conoce formalmente como varianza no-estacionaria y es aproximadamente equivalente a la forma más sencilla de heteroscedasticidad en regresión. La forma más común de varianza no-estacionaria ocurre cuando los valores de datos más grandes exhiben proporcionalmente la variabilidad más grande.

A menudo, la variabilidad puede ser estabilizada re-expresando los datos. Las re-expresiones más comunes son potencias enteras, raíces o raíces recíprocas y

logaritmos modificados para preservar el orden entre los valores positivos. Esto es para  $0 < a < b$  y la re-expresión  $f$ ,  $f(a) < f(b)$ . Una forma de la familia de potencias es:

$$F_p(y) = \begin{cases} y^p & p > 0 \\ \log(y) & p = 0 \\ -y^p & p < 0 \end{cases}$$

Puede consultarse a Stoto y Emerson (1984) y Berenson, Lavine y Goldstein (1983) para un análisis más profundo de la escala de transformación. Los efectos que estas re-expresiones tienen en los datos cambian sistemáticamente con  $p$ . Por ejemplo,  $\log(y)$  ( $p=0$ ) altera los datos de la misma manera que  $\sqrt{y}$  ( $p=1/2$ ) pero a un grado mayor. Por esto, Tukey (1977) lo denomina "escala de potencias". La consistencia con la cual los patrones en los datos cambian a través de diferentes potencias permite una selección sistemática de potencias apropiadas en una gran variedad de situaciones.

La re-expresión de datos utilizando la escala de potencias puede mejorar la simetría, estabilizar la varianza a través de los grupos, mejorar la linealidad de una relación  $x$ - $y$ , y mejorar la aditividad de una tabla de doble-entrada. Una ventaja práctica del EDA es que provee métodos sistemáticos, resistentes a los extremos, para seleccionar las re-expresiones que simplifican los modelos en los datos. Ver Tukey (1977), Leinhardt y Bosserman (1979) y Velleman y Hoaglin (1981) para discusiones y ejemplos. Ver Mosteller y Tukey (1977) para las referencias ahí expuestas para el marco teórico y discusiones filosóficas.

Un modelo natural para ser usado (cuando se presentan ciclos tales como picos o valles con duraciones de por lo menos 4 puntos) es el modelo aditivo usado en análisis de doble-entrada de varianza:

$$y_{ij} = \text{común} + \text{renglón } i + \text{columna } j + e_{ij} \quad (1)$$

### 3- Diagnóstico de la amplitud del ciclo no-estacionario.

Cuando la amplitud del ciclo cambia, el modelo aditivo simple (1) no se ajusta bien a los datos. Sin embargo, cuando los cambios en la amplitud del ciclo están relacionados a los cambios en el nivel de la secuencia de datos (por ejemplo, la amplitud del ciclo mayor cuando los valores de datos son mayores y menor amplitud cuando los valores de datos son menores) una re-expresión bien-escogida puede casi siempre simplificar el patrón. A menudo es posible diagnosticar el tipo de no-aditividad producido en una tabla de doble-entrada por esta clase de no-estacionaridad en un análisis de varianza con la prueba de un-grado-de-libertad usado por Tukey.

Para usar el diagnóstico exploratorio de no-aditividad, hay que calcular los valores de la comparación:

$$\text{com}_{ij} = (\text{renglón } i + \text{columna } j) / \text{común} \quad (2)$$

Luego graficar los residuales  $R_{ij}$  contra  $\text{Com}_{ij}$  y encontrar la pendiente,  $b$ , de  $R_{ij}$  sobre  $\text{Com}_{ij}$ . Es probable que la potencia  $1-b$  mejore la aditividad de la tabla, en este caso estabilizando las amplitudes de los ciclos.

En los trabajos del EDA no se discutió el mejoramiento de la estacionariedad de varianza en una secuencia de datos. El método propuesto por Velleman está conceptualmente relacionado con el método discutido en Tukey (1977) y Lainhardt y Wasserman (1979) para identificar una transformación para estabilizar la variabilidad a través de los grupos. En ese método, el log. de una medida de variabilidad, o amplitud, para cada grupo está graficado contra el log. de la mediana de ese grupo y se ajusta una recta a la exposición resultante. Una recta con una pendiente  $b$  sugiere que la potencia  $1-b$  en la escala de potencias hace que las variabilidades entre los grupos sean más parecidas. Por ejemplo, la amplitud es proporcional al nivel (como se midió aquí por medio de la mediana) entonces

$$\log(\text{valor absoluto del residuo}) = \log(\text{dato suavizado}) + \log(\text{proporcionalidad de la constante}) \quad (3)$$

Si la recta tiene una pendiente de 1 entonces  $p=0$  por la re-expresión del logaritmo. Si la re-expresión nula ( $p=1$ ) es necesaria, la pendiente va a estar cerca de 0 y (3) se transforma en:

$$\log(\text{valor absoluto del residuo}) = \text{constante} \quad (3')$$

La relación entre pendiente y potencia es exacta sólo en  $p=0$  y  $p=1$ , y es una buena aproximación cerca de estas potencias.

Para una secuencia de datos, cualquier suavizador que satisface la ecuación (3)\* -no importa que tan complejo- se comporta localmente como un estimador del la posición o nivel. Los valores del residuo,  $r_t$ , son desviaciones de la posición local especificada por el "smooth". Las medidas de escala natural local son

análogas a la desviación estándar,  $(S_m[r]t/2)^{1/2}$ , y al valor absoluto residual de la mediana,  $(S_m[r]t/3)$ . De los dos,  $S_m[r]t/3$  se parece más a otras medidas exploratorias: esto es simple y, cuando se utiliza un suavizador no-lineal apropiado, éste es resistente a los efectos de los valores extraordinarios en el residuo. (Esta resistencia es importante; un suavizador resistente causa valores de datos extraordinarios para transformarse en grandes residuos, entonces los picos tienen que estar expuestos en el residuo).

El método de diagnóstico gráfico  $\log(S_m[r]t/3)$  versus  $\log(S_m[y]t)$ . Al comparar estas dos medidas en grupos, la recta de pendiente b hace que la potencia  $(1-b)$  mejore la estacionaridad de la varianza.

Los análisis exploratorios emplean generalmente métodos de rectas-ajustadas resistentes a los extremos ocasionales. Mientras que los valores de la ordenada y de la abscisa habrán indicado picos suavizados fuera de ellos, esto garantiza solamente la suavización a lo largo del orden de la secuencia. La exposición del diagnóstico ignora este orden y puede incluir puntos ocasionales lejos del modo completo. La técnica exploratoria conocida como la recta resistente es una elección posible del método de ajuste. (Ver Tukey, 1977, Leinhardt y Wasserman, 1979, y Velleman y Hoaglin, 1981, para el método y ejemplos. Ver Johnstone y Velleman, 1981, para la historia, teoría, y resultados de distribuciones.) La escala de ajuste por mínimos-cuadrados puede también ser utilizada -especialmente si los residuos de la recta resistente revela valores no-extraordinarios.

La secuencia rugosa  $[rt]$  puede incluir valores nulos. Sin embargo, es menos

probable que el valor absoluto del residuo suavizado tenga ceros (y no puede tener valores negativos). Cuando los ceros ocurren, los puntos correspondientes pueden ser omitidos al encontrar los logaritmos y se pueden agregar operaciones subsiguientes a una pequeña constante (1/6 generalmente) a todos los valores absolutos de los residuos antes de tomar los logaritmos. Si el residuo es enteramente o predominantemente cero, la secuencia de datos originales era bastante suave, entonces es una pequeña varianza residual por estabilizar.

## 2- Secuencias estructuradas.

Una diferencia primaria entre la suavización de datos y los métodos más comunes de la exploración y el análisis de datos es que los suavizadores no-lineales de datos no requiere ni de un modelo ni de una estructura a priori en los datos más allá del ordenamiento de la secuencia. Esta ausencia de supuestos detallados puede ser una desventaja en el análisis confirmatorio, en el que un modelo paramétrico puede proveer hipótesis probadas y bases para hacer predicciones. En exploración, sin embargo, la libertad de supuestos acerca de los datos es una ventaja.

Consecuentemente, estarían incluidos otros supuestos estrictos gradualmente y con algunas trepidaciones. Una clase de estructura que puede a menudo ser asumida con una pequeña pérdida de generalidades es la repetición cíclica de patrones (sin especificar los patrones ellos-mismos). A menudo el periodo del ciclo es comúnmente conocido o tiene que ser confirmado previamente. Por ejemplo, los datos mensuales presentan casi siempre un ciclo anual; datos diarios, un ciclo



semanal; datos por hora, un ciclo diario. A menudo, esta estructura ayuda a separar los datos en patrón y residuo. Un ejemplo común es el ajuste estacionario de las series de tiempo en Economía y en Demografía. Los métodos de ajuste estacionario ya han sido desarrollados y extraordinariamente sofisticados. El programa SABL (Cleveland, Dunn, y Terpenning, 1979), el cual combina los suavizadores no-lineales como los discutidos aquí, con otras técnicas, es un buen ejemplo. Los análisis exploratorios de las secuencias no requieren de este grado de sofisticación (ni de la inversión de recursos computacionales). Algunos métodos más elementales, menos eficientes estadísticamente, y menos caros son generalmente suficientes. Los métodos exploratorios sacrifican casi siempre la eficiencia estadística (medida, por ejemplo, en términos de varianza relativa) para obtener resistencia y estar libre de supuestos restrictivos. Una regla común es que los buenos métodos exploratorios necesitan ser sólo 50% eficientes en relación al método óptimo. (Es evidente ya que el método óptimo depende del conocimiento de la población ni disponible para ni asumido por el método exploratorio). En realidad, los métodos exploratorios son, por lo general, mucho más eficientes que los estadísticas "óptimas" clásicas de datos "chatarra".

Algunos métodos de ajuste estacionarios requieren secuencias de datos largas para establecer patrones estacionales, otros se enfocan a especificar la extensión de los períodos (expresado como números de puntos). Los métodos exploratorios discutidos aquí pueden ser aplicados a secuencias relativamente cortas (cualquiera que tenga más de tres ciclos) y ajustan a algún período mayor de 7 puntos aproximadamente.

Si la secuencia del dato crudo es muy irregular, la secuencia debe ser suavizada y después la secuencia suavizada debe ser analizada en la tabla de doble-entrada. Si el dato crudo demuestra un claro comportamiento periódico, o si el patrón periódico no es suave, sería más acertado analizar los datos crudos en una tabla de doble-entrada.

Una vez que la secuencia suave haya sido analizada, una recta de mínimos-cuadrados puede ser apropiada para encontrar la pendiente. Si se usan los datos, una línea resistente permite juzgar de una mejor manera. En cualquier caso, es importante examinar la gráfica. Las pendientes fuertes y consistentes arguyen a favor de la re-expresión más fuertemente que los puntos borrosos; a estos últimos les puede suceder tener una pendiente distinta de cero.

Si, por ejemplo, los residuos son generalmente más negativos en los dos esquinas de la tabla; el extremo alto izquierdo y el extremo bajo derecho y más positivos arriba a la derecha y abajo a la izquierda - la clase de modelo de "silla de montar" o menudo indica la necesidad de una re-expresión. Entonces habría que graficar los residuos contra los valores de comparación, calcular la pendiente de esta gráfica e indicar una re-expresión para estabilizar la amplitud del ciclo con éxito. Los patrones en el renglón y los efectos columna de tablas más grandes pueden ser más claras graficando los efectos contra el número del renglón o la columna y suavizando.

Es probable que la amplitud del ciclo no-estacionario sea prominente cuando

la tendencia a largo plazo es escalonada y las series son muy largas (a fin de que la magnitud de los valores de datos cambie substancialmente). Cuando la tendencia es casi lineal (como se asume a menudo), el modelo aditivo simple está bien ajustado. En los datos mensuales con ciclo anual, por ejemplo, los efectos mes cambian linealmente con una pendiente correspondiente a la de la tendencia; los efectos año que estiman un nivel central para cada periodo anual son también lineales. Si la tendencia no es lineal el modelo aditivo no puede ajustarla bien. Es probable que esto suceda si los datos cubren un ancho rango y requieren una re-expresión. Cleveland, Dunn, y Terpenning (1979) mencionan brevemente este problema como una razón para emplear una técnica iterativa más compleja para diagnosticar la no-estacionaridad del ciclo en SABL. Se pueden hacer varios ajustes a los métodos descritos aquí para encontrar este problema.

Si la amplitud de este patrón periódico es poca relativo al cambio de largo plazo en el nivel, es posible utilizar las medias estándares exploratorias de la no-linealidad de estimación en las relaciones de  $x$ - $y$ . Para hacer que la tendencia a largo plazo sea casi lineal, la transformación de la secuencia de datos va a mejorar el ajuste del modelo aditivo y estabilizará la amplitud del ciclo así como la variabilidad de los residuos.

Alternativamente, el modelo aditivo (1) puede ser desarrollado para ajustar una pendiente para cada año. El modelo resultante es

$$Y_t = \text{común} + \text{año } i + \text{edad } j + b_{ij} + \text{residuo } t$$

$$i=0, \dots, I; j=1, \dots, J; t=Ji+j \quad (4)$$

Se calcula, para cada año, la pendiente  $b_i$ . La ecuación (4) es equivalente a ajustar una descripción de un "lotajucioso-lineal" de la tendencia. Si los datos son suficientemente irregulares, si se espera que la tendencia sea sistemática, y si se dispone de varios años, entonces la secuencia de los  $b_i$ 's puede ser suavizada.

Este método puede resultar más laborioso y más caro que lo que usualmente se espera de los métodos exploratorios. Este requiere seguramente de una computadora y, al nivel de programas publicados (Valleman y Hoaglin, 1981) y programación adicional.

## D- REVELACION GRAFICA DE LOS DATOS.

Una de las aportaciones que ha impulsado fuertemente al análisis exploratorio de datos, ha sido el desarrollo de métodos visuales para inspeccionar lotes de datos. Una gráfica bien seleccionada puede transmitir sintéticamente un gran volumen de información cuantitativa, sin eliminar detalles o variaciones azarosas que puedan contener los datos.

### 1- El diagrama de Tallo-y-hoja.

En la opinión de Wainer y Theissen (1981), el diagrama de "Tallo-y-Hoja" es el instrumento más importante para el análisis de lotes de datos desde que apareció la clásica prueba de "t-student". El diagrama fué ideado por Tukey (1977) con el objeto de comunicar simultáneamente los valores numéricos de un lote de datos con la forma natural de su distribución. De hecho, el diagrama de tallo-y-hoja puede definirse como un híbrido que combina los aspectos visuales de un histograma con la información numérica que proporciona una tabla de distribución de frecuencias. Para su uso y construcción, el lector interesado puede consultar a Curtis, 1986.

Para explicar su uso y construcción, conviene considerar el siguiente ejemplo. Considérese la siguiente tabla, tomada de Rodríguez et al. (1987):

Fondos para Salubridad, Bienestar y Asistencia.Gasto Federal y cambio social. En México.

Año	Porcentaje ejercido	Año	Porcentaje ejercido
1940	6.4	1952	2.5
1941	6.5	1953	3.2
1942	6.4	1954	2.7
1943	5.8	1955	2.8
1944	4.7	1956	2.9
1945	4.9	1957	3.3
1946	3.4	1958	3.3
1947	4.9	1959	3.4
1948	4.1	1960	3.5
1949	3.3	1961	3.9
1950	3.8	1962	4.0
1951	3.1	1963	3.3

Para construir la forma más simple de un diagrama de tallo-y-hoja se procede de la siguiente manera:

1. Ordenar el lote de datos en magnitud creciente.
2. Escoger un par conveniente de dígitos que permita fraccionar en dos partes el lote de datos. Por ejemplo, el primer valor correspondiente es 6.4, al cual puede ser fraccionado en su parte entera (6) y su parte decimal (4).
3. Formar el tallo y la hoja con las fracciones respectivas. Este proceso se ilustra como:

Valor del dato: 6.4

Tallo: 6

Hoja: 4

4. Construir el tallo escribiendo verticalmente los dígitos enteros entre el 2 y el 6, asociando a cada uno su hoja respectiva. Los dígitos del tallo están separados de los dígitos de la hoja por medio de una línea vertical.

El diagrama completo se muestra en la siguiente figura. Nótese que el diagrama contiene el recordatorio que el valor de los datos está en unidades de 0.1. Cabe mencionar que, como en este caso, el formato del diagrama puede aparecer demasiado tupido, teniendo muchas hojas por línea, se "alarga" el diagrama, al duplicar los dígitos. Las hojas cuyos valores estén dentro del intervalo  $[0,4)$  se colocan en la línea que contenga el círculo "a"; las hojas cuyos valores estén en el intervalo  $[5,9]$  se colocan en la línea que contenga la estrella "\*".

Unidad = 0.1

112 representa 1.2

20	5789
30	12333344
35	589
40	01
45	799
50	
55	8
60	44
65	5

El ejemplo analizado permite afirmar que el diagrama de tallo-y-hoja es una forma efectiva y flexible para observar la estructura general de los datos. La apariencia general del diagrama se asemeja a la de un histograma; con la salvedad que las hojas preservan toda la información numérica del lote de datos. En términos generales, un diagrama de esta naturaleza hace visible las siguientes características:

1. Muestra el rango de valores que los datos cubren.
2. Determina donde se concentra la mayoría de los datos.
3. Describe la simetría del conjunto de datos.
4. Identifica si existen "huecos" en la distribución de los datos.
5. Señala aquellos valores que claramente se desvían del conjunto de datos.



## 2- El Diagrama de Caja.

El diagrama de caja fue ideado con objeto de resumir gráficamente la distribución de los datos. Por su naturaleza visual, las distancias representadas numericamente en el diagrama de tallo-y-hoja están comunicadas de manera más efectiva en un diagrama de "caja."

Para dibujar el diagrama de caja correspondiente a los datos del ejemplo anterior, se debe ordenar los datos de menor a mayor, calcular la mediana y su rango:

2.5	3.5
2.7	3.8
2.8	3.9
2.9	4.0
3.1	4.1
3.2	4.7
3.3	4.9
3.3	4.9
3.5	5.8
3.3	6.4
3.4	6.4
3.4	6.5

El rango de la mediana es  $(N+1)/2$ , donde  $N$  es el número de datos. En este ejemplo,  $N$  es igual a 24, por lo tanto el rango es 12.5. Entonces la mediana es la observación que se encuentra en el centro  $(3.4+3.5)/2 = 3.45$ .

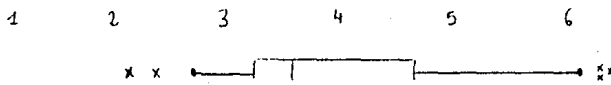
Un resumen importante es el de los extremos, es decir los valores mayor y menor del conjunto de datos. Aquí, los valores extremos son 2.5 y 6.5 e indican que tan dispersos están los datos.

Otro par de resúmenes importantes son los "bisogras", que encierran la mitad central de los datos. El rango de las bisogras es  $(n+1)/2$ , donde  $n$  es la mitad de  $N$ . Por lo tanto, el rango es  $(n+1)/2 = (12+1)/2 = 6.5$ . Por lo tanto, la primera "bisogra" es el valor de dato que se encuentra entre la sexta y la séptima observación  $(3.2+3.3)/2 = 3.25$  y la segunda "bisogra" es  $(4.7+4.9)/2 = 4.8$  (el promedio entre la sexta y la séptima observación de mayor a menor).

Nótese que se ha obtenido, por el momento, un resumen de "cinco números": la mediana, las bisogras y los extremos.

Esto es (2.5, 3.25, 3.45, 4.8, 6.5).

Como lo ilustra la siguiente figura, la línea vertical al centro de la caja representa la localización de la mediana de la distribución. La línea vertical del lado izquierdo es la localización de la primera bisogra y la del lado derecho de la caja es la segunda bisogra.



De esta manera, se hace visible la simetría en las orillas de la parte central, las observaciones adyacentes y las observaciones fuera del patrón. El lector que desee mayor información sobre el uso y la construcción del diagrama de caja, pueda consultar a Tukey (1977).

Es importante señalar que el diagrama de caja no contiene información numérica: no se muestra todo lo que muestra el diagrama de tallo-y-hoja, pero provee una mejor imagen de las colas de la distribución. Por lo tanto, los diagramas de caja son particularmente útiles para identificar y comunicar las características mayores de las distribuciones cuyas colas se desvían de las de la distribución normal.

### CAPITULO III

#### DISEÑO EXPERIMENTAL

"El centro de la actividad cognoscitiva de los seres humanos son las hipótesis y no los datos."

Mario Bunge '

### III- DISEÑO EXPERIMENTAL

A continuación, se hace una descripción de los casos de estudio: la población, las variables, las hipótesis y los modelos apropiados para el análisis de los datos.

#### A- FUENTES DE INFORMACION

El censo de población constituye la fuente primaria de los datos demográficos básicos de cualquier país y es la base de un programa de recolección de información, con una periodicidad de 10 años.

Las posibilidades de utilización de la información censal dependen, en gran medida, del tipo, amplitud y comparabilidad de los datos básicos, así como del conocimiento de los procedimientos de recolección, sistematización y presentación de la información; ya sea para fines administrativos, para la planificación económica y social o para investigaciones que tomen en cuenta las características de la población.

Entre las estadísticas demográficas y sociales, las más interesantes para este estudio son las "Estadísticas Vitales", que son aquellas que se elaboran a partir de los datos del registro civil y que, fundamentalmente, registran hechos.

SECRETARÍA DE ECONOMÍA  
ESTADÍSTICA  
CENSO DE PUEBLO Y VIVIENDAS  
1970

## B- LOS DATOS DEMOGRAFICOS

La información proporcionada por el censo o el registro civil, expresada en números absolutos, es de gran interés y utilidad para fines administrativos, para programas de salud pública, para edificación de escuelas, para estimaciones de la población, etc; pero es insuficiente, sin embargo, para hacer comparaciones en el tiempo y en el espacio.

Por esta razón, Rodríguez (1987) aclara que una de las primeras fases del análisis de datos demográficos se inicia a partir del cálculo de los números "relativos"; es decir, cuando las frecuencias absolutas provenientes tanto del censo como de las estadísticas vitales puede relacionarse entre sí, convenientemente. Por ejemplo, si se relacionan grandes grupos de edad con el total de la población, se obtiene un número "relativo", que informa sobre la estructura de ésta.

No obstante, el propósito en este trabajo no consiste en el procedimiento de obtención de los datos, sino en la aplicación de técnicas robustas de análisis de datos que se obtuvieron a partir de información censal. Cabe agregar, que se incluye una breve descripción de los datos a examinar.

### 1- TASA DE MORTALIDAD GENERAL DE 1922-1982.

En el estudio de la mortandad, las tasas de mortalidad general son los datos más utilizados, porque no sólo miden la calidad de la vida sino que se utilizan como indicadores para evaluar la situación de la salud en el país.

La tasa bruta de mortalidad, relación entre las defunciones de un año y el promedio de población para ese mismo lapso, es el índice más simple; el examen de sus variaciones en el seno de una misma población proporciona, por sí sola, indicaciones interesantes.

La tabla 4 muestra la tasa de mortalidad por mil anual, de 1922 a 1982, en México. Los datos fueron tomados de los cuadros "Movimientos de la Población" de los Cuadernos de Estadísticas Vitales de la SPP de 1960, 1970, 1980 y 1985.

Unidades de análisis: Individuos de la población.

Variabla independiente: X = tiempo de 1922 a 1982 año por año.

Variabla dependiente: Y = tasa de mortalidad por mil.

Elementos lógicos: la tasa de mortalidad disminuye conforme avanza el tiempo.

La disminución de la mortalidad es un fenómeno lento, colectivo, del que es difícil percatorse a nivel individual. Lleva consigo toda la evolución demográfica desde hace más de medio-siglo. Una inspección somera a la tabla 4 revela cambios en la mortalidad. Entre 1922 y 1982, la mortalidad ha disminuido constantemente.

Tabla 4<sup>a</sup> Tasa Bruta de Mortalidad  
México, 1922-1982

Edad	mortalidad
22	25.3
23	24.4
24	25.6
25	26.5
26	26.9
27	26
28	25.5
29	26.8
30	26.4
31	25.9
32	26.1
33	25.7
34	25.8
35	25.4
36	25.5
37	26.4
38	22.9
39	23
40	23.2
41	22.1
42	22.8
43	22.4
44	20.6
45	19.5
46	19.4
47	16.6
48	16.9
49	17.9
50	16.2
51	17.3
52	15
53	15.9
54	13.2
55	13.5
56	11.7
57	12.7
58	12
59	11.4
60	11.2
61	10.4
62	10.5
63	10.4
64	9.9
65	9.5
66	9.6
67	9.2
68	9.6
69	9.4
70	9.6
71	8.7
72	8.8
73	8.2
74	7.5
75	7.2
76	7.5
77	7
78	6.8
79	6.4
80	6.2
81	5.5
82	5.7



Entre 1922 y 1933, la tasa bruta de mortalidad oscilaba entre 24 y 26.8 defunciones por mil habitantes. A partir de 1933, como puede verse en el cuadro, la trayectoria descendente de la mortalidad aparece claramente: de 1933 a 1943, la tasa bruta de mortalidad se redujo poco a poco: hasta llegar a 22.4 por mil. Y en 30 años, la tasa se redujo en un 42% al pasar de 26 a 15 por mil en 1952. Mientras que en los 30 años siguientes de 1954 a 1982, la reducción fue del 58.5%. Si se separan los periodos por cada 10 años, se observa que el ritmo de decrecimiento tiende a ser más lento, a partir de 1960.

¿En qué años, se presentan picos? ¿Se atribuyen, acaso, descansos inesperados y rápidos en años específicos? Habría que profundizar la cuestión para comprender la reducción importante del ritmo de aumento de la vida durante estos años y la aparición de cierta reactivación durante tales otros años. Estas cuestiones se podrían aclarar, analizando los datos trabajados en vez de analizar los datos arrojados en bruto.

Puesto que la variable  $X$  tiene sus valores igualmente espaciados, cabe preguntarse qué cambios sufre la variable  $Y$ , si se suaviza esta secuencia de datos. Y como aquí la secuencia proviene del registro de un valor por intervalo de tiempo, los valores de  $Y$  forman una serie de tiempo. Para describir el patrón general de esta serie, la técnica exploratoria más usual es la suavización de datos por medianas corridas. Así, la tendencia demográfica de la tasa de mortalidad general es más fácil estudiarla a partir de los datos suavizados que de los datos crudos.

2- Esperanza de vida, por edad exacta, según Entidad Federativa, y según sexo, 1975.

La esperanza de vida a la edad  $X$  es el número de años que les queda por vivir, de promedio, a las personas de edad  $X$ , en las condiciones de mortalidad descritas por la tabla de mortalidad; la esperanza de vida al nacer, en particular, es la duración media de vida de individuos sometidos, a partir de su nacimiento, a las condiciones de mortalidad de la tabla; también se llama vida media.

De la información que se presenta en las tablas 5a y 5b se desprenden varias observaciones:

Unidades de análisis: Individuos de la población.

Variablas independientes:  $X$  = Edad exacta (cada 5 años).

EF = Entidad Federativa.

S = Sexo (variable dicotómica).

Variable dependiente:  $Y$  = Esperanza de vida.

Elementos lógicos: - A menor nivel de vida, menor esperanza de vida.

- Al sexo femenino, esperanza de vida mayor.

- A menor edad, mayor esperanza de vida.

ESTADÍSTICA DE FUMOS, POR EDAD (SEXOS), SEGUN CUANTOS FUMAN (1975)  
 Sexes (1975)

Estados	0	1 año	5 años	10 años	15 años	20 años	25 años	30 años	35 años	40 años	45 años	50 años	55 años	60 años	65 años	70 años	75 años	80 años
Total	62.94	62.22	62.45	57.81	55.96	60.54	61.25	60.81	60.75	61.26	61.70	61.87	60.33	60.53	59.66	58.63	57.93	57.82
1. Aguascaltecas	63.90	64.22	63.10	60.40	57.81	61.17	64.01	60.62	60.20	61.20	61.10	61.11	60.33	60.30	59.10	58.10	57.10	57.10
2. Baja California	61.64	60.13	60.69	57.00	53.87	60.10	60.00	59.60	59.60	59.60	59.60	59.60	59.60	59.60	59.60	59.60	59.60	59.60
3. Baja California Sur	60.04	60.00	61.87	60.00	60.00	61.31	60.00	61.20	61.37	61.34	60.62	60.61	60.30	59.10	58.10	57.10	56.10	56.10
4. Campeche	60.95	60.20	61.11	60.20	58.10	61.10	60.00	61.11	57.75	58.00	59.81	60.00	61.00	61.00	61.00	61.00	61.00	61.00
5. Coahuila	61.83	61.91	62.11	57.20	52.11	60.32	61.00	60.00	58.16	60.00	60.70	61.00	60.15	59.17	58.17	57.17	56.17	56.17
6. Colima	60.30	60.91	60.27	57.00	52.11	61.11	61.00	60.50	58.12	60.00	60.91	61.00	60.00	59.10	58.10	57.10	56.10	56.10
7. Chiapas	60.63	62.20	60.13	56.97	51.30	57.83	61.04	61.63	60.51	60.64	60.45	61.00	60.10	59.10	58.10	57.10	56.10	56.10
8. Chihuahua	64.72	67.25	64.10	59.20	56.36	60.27	61.63	61.00	57.13	60.00	60.60	61.00	60.00	59.00	58.00	57.00	56.00	56.00
9. Distrito Federal	65.91	67.52	63.91	59.10	54.17	61.17	61.00	60.40	58.73	60.00	61.67	61.00	60.00	59.00	58.00	57.00	56.00	56.00
10. Durango	60.63	60.00	61.60	59.13	55.12	61.10	60.10	61.00	57.00	55.17	59.11	59.00	58.00	57.00	56.00	55.00	54.00	54.00
11. Guanajuato	61.57	60.70	62.90	60.31	51.00	59.20	61.50	60.21	60.00	60.00	60.10	61.00	60.00	59.00	58.00	57.00	56.00	56.00
12. Guerrero	61.71	60.20	62.94	61.42	58.00	60.10	60.00	60.11	58.97	58.00	60.20	60.00	59.00	58.00	57.00	56.00	55.00	55.00
13. Hidalgo	60.87	61.20	60.10	58.00	53.76	59.87	60.00	57.00	57.00	57.00	58.00	58.00	58.00	58.00	58.00	58.00	58.00	58.00
14. Jalisco	63.54	60.54	61.20	58.53	54.27	60.72	60.00	60.67	58.63	58.61	58.20	58.20	58.20	58.20	58.20	58.20	58.20	58.20
15. Jalapa	60.10	60.11	62.11	57.00	51.62	60.17	61.17	60.00	58.16	60.00	61.00	61.00	61.00	61.00	61.00	61.00	61.00	61.00
16. Michoacán	61.64	60.00	62.11	57.15	51.92	57.92	61.70	60.70	58.20	60.00	60.10	60.10	60.10	60.10	60.10	60.10	60.10	60.10
17. Morelos	61.00	60.20	61.00	60.21	51.00	58.10	62.00	59.83	58.10	60.00	57.87	58.00	58.10	58.10	58.10	58.10	58.10	58.10
18. Nayarit	62.27	60.00	62.11	60.00	56.10	60.10	61.11	60.00	58.20	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00
19. Nuevo León	60.00	61.20	61.21	60.00	54.20	60.00	60.00	60.00	58.21	58.11	58.10	58.10	58.10	58.10	58.10	58.10	58.10	58.10
20. Oaxaca	60.00	57.67	60.62	60.00	49.00	61.10	60.10	60.00	58.20	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00
21. Puebla	61.00	60.72	61.13	58.00	53.00	61.22	60.12	60.00	58.00	60.00	60.53	60.00	59.00	58.00	57.00	56.00	55.00	55.00
22. Querétaro	61.12	64.00	60.00	57.87	51.17	60.64	61.00	60.00	58.00	60.00	60.52	61.00	60.00	59.00	58.00	57.00	56.00	56.00
23. Quintana Roo	60.44	60.64	61.00	60.25	57.00	53.10	61.12	61.00	60.70	60.00	60.12	61.00	60.00	59.00	58.00	57.00	56.00	56.00
24. San Luis Potosí	62.23	64.00	61.11	57.00	52.00	60.44	61.00	60.00	58.00	60.00	60.34	61.17	60.10	59.10	58.10	57.10	56.10	56.10
25. Sinaloa	60.77	61.10	61.11	60.00	54.00	60.00	61.00	61.00	57.11	58.00	60.01	61.00	61.00	61.00	61.00	61.00	61.00	61.00
26. Sonora	61.61	60.00	60.73	57.87	51.11	60.70	60.70	60.00	58.69	60.00	61.00	61.00	61.00	61.00	61.00	61.00	61.00	61.00
27. Tamaulipas	61.12	61.62	60.70	58.20	51.17	60.20	61.00	61.00	57.52	58.00	59.40	61.17	61.00	61.00	61.00	61.00	61.00	61.00
28. Tlaxcala	60.00	60.00	61.00	60.00	50.10	60.10	60.00	60.00	58.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00	60.00
29. Veracruz	61.03	60.00	61.00	60.10	51.10	60.77	61.00	60.00	58.11	60.00	60.10	61.00	60.00	59.00	58.00	57.00	56.00	56.00
30. Yucatán	61.00	60.27	61.43	60.63	51.14	61.63	61.12	60.00	58.10	60.00	61.00	61.00	60.00	59.00	58.00	57.00	56.00	56.00
31. Zacatecas	60.30	61.17	61.00	60.00	50.00	60.00	60.00	61.00	57.50	58.00	59.13	60.23	60.60	60.00	59.00	58.00	57.00	57.00
32. Baja Verapaz	64.77	61.00	62.11	61.40	50.00	58.00	61.63	61.63	58.00	60.00	60.65	60.66	61.54	60.53	59.60	58.60	57.60	57.60

tomado de las Estadísticas Vitales.  
 1966-1975  
 SSP

TABLA 5a

ESTADÍSTICA DE ODM. POR EDAD SEXO Y ESTADO CIVIL  
 SEXO Femenino  
 1975

Entidad	0	5 años	10 años	15 años	20 años	25 años	30 años	35 años	40 años	45 años	50 años	55 años	60 años	65 años	70 años	75 años	80 años	
Total	67.87	69.00	66.90	62.30	57.60	52.70	48.20	43.67	39.70	34.30	30.64	26.95	23.36	19.80	15.80	11.80	6.77	2.83
Aguascalientes	60.34	70.42	67.52	60.70	50.00	42.30	40.70	40.23	39.70	35.30	30.90	26.60	22.20	18.10	14.10	10.70	7.30	3.33
Baja California	60.62	70.32	66.70	60.00	57.10	50.30	47.00	43.00	38.40	33.00	29.51	25.12	21.00	17.10	13.70	10.70	6.67	2.70
Baja California Sur	70.00	75.00	69.50	64.71	59.00	54.00	50.00	46.20	42.50	39.00	35.50	32.20	29.00	25.50	22.00	18.50	15.00	6.70
Belize	69.00	71.00	69.00	63.00	56.00	50.17	45.00	40.20	36.20	32.00	28.07	24.27	20.80	17.77	14.00	11.17	8.02	4.83
Campeche	66.50	70.20	67.00	62.21	57.30	52.61	47.00	42.00	38.71	34.00	30.00	26.50	23.50	20.00	16.90	13.70	10.00	5.91
Colima	67.27	66.74	64.42	61.00	56.30	50.21	47.54	43.00	38.20	33.90	29.50	25.50	21.50	18.42	14.12	10.70	7.10	3.13
Chiapas	64.29	61.00	54.00	48.21	40.71	36.67	34.00	33.00	32.73	31.90	31.03	29.90	28.90	27.00	24.70	21.51	18.04	6.75
Chihuahua	69.40	71.61	68.51	62.71	56.00	50.12	45.00	40.70	36.17	31.65	27.31	23.07	19.00	15.30	12.20	9.20	6.27	3.17
Distrito Federal	70.00	72.00	69.70	63.00	56.00	50.00	45.00	40.71	36.00	31.42	26.90	22.00	17.50	13.00	10.10	7.00	5.00	4.21
Durango	71.20	71.43	69.00	64.20	59.00	54.00	49.00	45.17	40.50	36.00	31.04	27.20	23.12	19.10	14.40	10.00	6.00	3.21
Guatemala	66.30	70.10	67.50	62.70	58.00	53.27	48.61	44.00	39.51	35.10	30.70	26.20	22.20	18.10	14.10	10.00	6.12	3.04
Herrera	60.70	64.00	60.20	62.10	57.41	52.20	47.54	43.00	38.42	34.00	30.00	27.20	23.00	19.00	15.00	11.00	7.12	4.14
Michoacán	63.00	64.00	62.30	57.91	50.17	46.67	44.00	40.00	36.20	31.00	27.00	24.10	20.41	16.90	13.50	10.10	7.12	3.13
Morelos	60.20	70.00	67.63	61.00	53.10	47.42	42.20	38.00	34.00	30.00	26.70	23.00	19.20	15.10	11.51	8.50	5.95	3.05
Nayarit	66.71	70.07	67.00	62.47	58.00	53.32	48.77	44.21	39.60	35.11	30.90	26.90	22.90	19.20	15.00	11.00	7.10	4.10
Nuevo León	69.00	70.20	67.20	62.47	57.64	52.00	47.00	43.70	39.30	35.00	30.00	26.50	22.30	18.10	14.10	10.10	6.12	3.04
Oaxaca	69.60	70.70	67.63	62.00	56.10	51.71	46.00	41.20	36.50	31.00	26.90	22.20	17.70	13.00	10.00	6.00	4.00	3.00
Puebla	73.70	73.20	69.57	65.10	60.21	55.43	50.00	45.90	41.30	36.52	32.13	27.67	23.14	18.50	13.70	9.10	5.57	2.50
Quercuaro	60.13	63.20	60.70	56.70	52.10	47.73	43.00	38.50	34.00	31.00	26.00	22.20	18.50	14.71	11.00	7.10	3.10	0.10
Quintana Roo	63.70	66.50	64.51	60.00	53.00	50.00	46.00	42.00	38.11	33.50	29.70	25.50	21.70	18.20	15.00	11.77	8.75	3.04
Guerrero	60.61	64.63	61.63	61.10	56.17	51.70	47.20	42.00	38.63	34.42	30.20	26.13	22.00	17.90	13.27	9.10	5.07	2.01
Veracruz	73.62	74.61	71.11	66.20	61.20	56.43	51.70	47.10	42.62	38.73	34.90	30.00	25.60	21.20	17.10	13.20	9.11	5.00
San Luis Potosí	60.30	66.42	60.00	61.20	56.00	51.00	47.00	42.82	38.20	34.00	29.90	26.00	22.00	18.00	14.00	10.00	6.00	3.00
Tlaxcala	72.10	72.11	68.04	64.00	60.00	54.90	50.00	46.00	41.71	36.13	31.00	27.00	22.70	18.11	13.61	9.57	5.60	2.47
Tamaulipas	60.00	70.70	67.13	61.30	54.00	50.20	46.00	41.00	36.74	34.20	29.30	25.00	21.53	17.90	14.00	10.00	6.10	3.10
Tlaxcala	60.11	65.00	60.00	61.67	57.20	53.00	48.00	43.10	38.20	33.04	28.40	23.00	18.50	14.00	10.00	6.00	3.00	0.10
Tlaxcala	71.54	71.03	66.00	61.20	56.00	50.00	45.00	40.00	35.00	30.11	25.20	20.20	15.20	10.20	6.20	2.20	0.20	0.10
Tlaxcala	60.20	70.57	67.27	62.00	56.21	50.00	45.00	40.52	36.13	31.00	26.00	21.20	16.50	12.00	8.00	4.00	2.00	0.10
Veracruz	67.30	66.00	60.27	61.00	56.00	51.20	46.20	41.50	36.70	31.00	26.00	21.00	16.00	11.00	6.00	1.00	0.00	0.00
Yucatán	60.60	70.00	67.00	62.20	57.43	52.70	48.00	43.00	38.20	34.00	29.00	24.70	20.70	17.00	13.00	9.00	5.00	2.00
Zacatecas	66.70	71.20	68.20	62.57	56.70	51.20	46.00	41.00	36.11	31.27	27.50	23.40	19.20	15.20	11.00	7.00	3.00	0.10

\* Tomado de las Estadísticas Vitales  
 1966-1975  
 SPP

TABLA 5b

En primer término, se observan diferencias entre la esperanza de vida femenina y la masculina. La esperanza de vida al nacer es, en 1975, de 62.94 años para los hombres y de 67.87 para las mujeres. Se observa que la esperanza de vida femenina es mayor que la masculina, en casi todas las edades y que el descenso de esta última es más rápida que el de aquellas. Este fenómeno se conoce como la sobremortalidad masculina en todas las edades.

No obstante la correspondencia entre la esperanza de vida por edad a los diferentes niveles, el descenso de la esperanza de vida no tiene la misma intensidad en todas las edades. Esto se puede observar calculando la ganancia (o pérdida) media en años, por edad. En México, para cualquier época y a diferentes niveles, el patrón universal de la mortalidad ha prevalecido: elevadas tasas de mortalidad en el primer año de vida que descienden rápidamente hasta llegar a valores mínimos entre los 10 y 15 años; a partir de estas edades las tasas se vuelven a incrementar, de manera lenta en un principio y a mayor ritmo después de los 60 años de edad aproximadamente, conforme la edad aumenta.

Y por lo que se refiere a la esperanza de vida por regiones, se aprecia que existen disparidades en los niveles de mortalidad entre las diferentes entidades federativas del país. Tanto en hombres como en mujeres y para casi todas las edades, los valores superiores corresponden a los Estados de Durango, Nuevo León, Quintana Roo, Sinaloa y Tamaulipas y los inferiores a los de Chiapas, Hidalgo, Oaxaca y Puebla.

Sin embargo, ¿cómo se puede interpretar, con mayor claridad, la relación entre estas tres variables (entidad federativa, edad y sexo) y la variable de respuesta (esperanza de vida)? Una posible explicación es a través de un análisis de tablas de doble-entrada, por la técnica denominada "mediana pulida".

#### CAPITULO IV

#### ANALISIS DE RESULTADOS

"Frecuentemente, cuando se termina una pieza de investigación, se ve que han surgido nuevos problemas, nuevos temas y nuevas cuestiones como resultado de los que originalmente se habían considerado en el trabajo de investigación".

Pauline Young

## **IV - ANALISIS E INTERPRETACION DE LOS RESULTADOS**

En este apartado, se reseña la situación de los datos demográficos y se expone una serie de cuadros y diagramas derivados de los cálculos de las técnicas explicadas anteriormente. La presentación de estos cuadros conlleva el propósito de proporcionar, de manera resumida, un panorama general de la mortalidad en las diferentes entidades y en los períodos que abarcan las tablas de mortalidad presentadas en el capítulo anterior. Este mismo grupo de cuadros constituye el material básico utilizado en esta investigación para validar los métodos exploratorios.

A continuación, se presentan los aspectos más relevantes de los resultados de la investigación realizada, es decir, se efectúa básicamente una interpretación conjunta y breve de los mismos, con el doble objeto de resaltar aquellos aspectos comunes y de no extender esta ponencia.

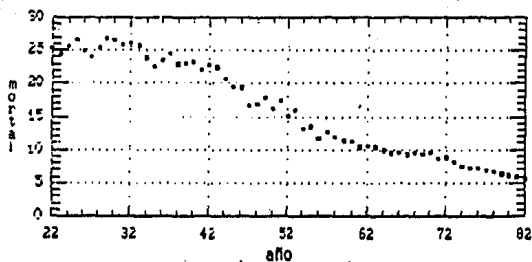
### **A- LA TASA DE MORTALIDAD EN MEXICO 1922-1982.**

#### **1- Presentación de los resultados.**

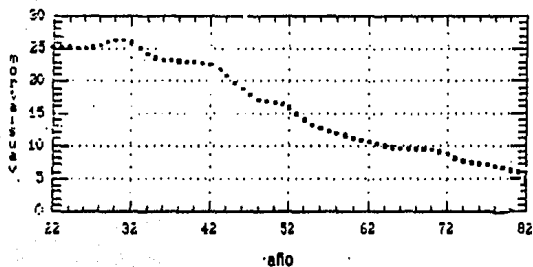
La tabla 4 y la gráfica 2 muestran los datos de tasas de mortalidad de 1922 a 1982.



Gráfica 2  
Tasa bruta de mortalidad



Gráfica 3  
Tasa bruta de mortalidad suavizada



Gráfica 4  
Residuos de la Suavización

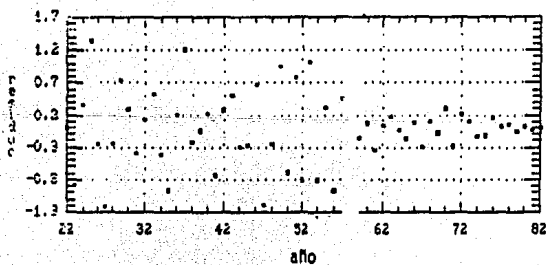


Tabla 6

SUAVIZACION DE LA TASA DE MORTALIDAD POR 4253hd.  
1922-1982

año	Yt	Zt	Rt
22	25.3	25.3	0
23	24.4	25.25	-0.85
24	25.6	25.25	0.37
25	26.5	25.16	1.34
26	24.9	25.14	-0.24
27	24	25.21	-1.21
28	25.5	25.55	-0.25
29	26.8	26.06	0.74
30	26.6	26.31	0.29
31	25.9	26.29	-0.39
32	26.1	25.97	0.13
33	25.7	25.17	0.55
34	23.8	24.21	-0.41
35	22.6	25.57	-0.97
36	23.5	23.3	0.2
37	24.4	23.2	1.2
38	22.9	23.13	-0.23
39	23	23.05	-0.05
40	23.2	22.97	0.23
41	22.1	22.83	-0.73
42	22.8	22.52	0.26
43	22.4	21.9	0.5
44	20.6	20.89	-0.29
45	19.5	19.77	-0.27
46	19.4	18.73	0.67
47	16.6	17.79	-1.19
48	16.9	17.14	-0.24
49	17.9	16.95	0.95
50	16.2	16.88	-0.68
51	17.3	16.51	0.79
52	15	15.8	-0.8
53	15.9	14.87	1.03
54	13.1	13.9	-0.8
55	15.5	13.18	0.32
56	11.7	12.67	-0.97
57	12.7	12.24	0.46
58	12	11.89	0.11
59	11.4	11.55	-0.15
60	11.2	11.13	0.07
61	10.4	10.74	-0.34
62	10.5	10.46	0.04
63	10.4	10.23	0.17
64	9.5	9.94	-0.04
65	9.5	9.66	-0.16
66	9.6	9.32	0.08
67	9.2	9.49	-0.29
68	9.6	9.5	0.1
69	9.4	9.47	-0.07
70	9.6	9.29	0.31
71	8.7	8.96	-0.28
72	8.8	8.59	0.21
73	8.2	8.1	0.1
74	7.5	7.63	-0.13
75	7.2	7.32	-0.12
76	7.3	7.15	0.15
77	7	6.98	0.02
78	6.6	6.74	0.06
79	6.4	6.46	-0.06
80	6.2	6.16	0.02
81	5.5	5.93	-0.02
82	5.7	5.7	0

donde

Yt = Tasa de mortalidad

Zt = Tasa de mortalidad Suavizada

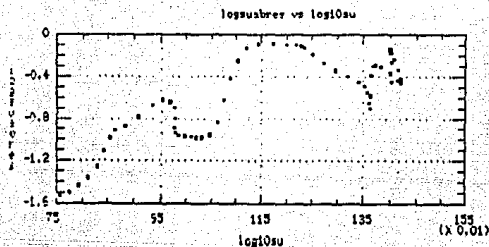
Rt = Residuo de la suavización

Tabla. 7  
DIAGNOSTICO DE LA NO-ESTACIONARIDAD

año	$S_m(Rt)$	$Log(Zt)$	$Log(S_m(Rt))$
22	0.3591	1.40312	-0.444785
23	0.5543	1.40295	-0.256255
24	0.6751	1.40192	-0.170632
25	0.7214	1.40071	-0.141824
26	0.7262	1.40037	-0.138944
27	0.6816	1.40157	-0.16647
28	0.5766	1.40739	-0.237622
29	0.4594	1.41597	-0.337809
30	0.3732	1.42012	-0.428058
31	0.3456	1.41979	-0.463947
32	0.3677	1.41447	-0.434506
33	0.4298	1.40088	-0.366734
34	0.4878	1.38359	-0.311758
35	0.5135	1.37236	-0.289121
36	0.5044	1.36736	-0.297255
37	0.4686	1.36549	-0.388702
38	0.2639	1.36418	-0.578561
39	0.1394	1.36267	-0.700275
40	0.2261	1.36116	-0.645699
41	0.2835	1.35851	-0.546835
42	0.3251	1.35257	-0.487983
43	0.3561	1.34644	-0.448428
44	0.4035	1.31994	-0.394156
45	0.4537	1.29601	-0.343231
46	0.531	1.27254	-0.274905
47	0.6541	1.25018	-0.184356
48	0.7502	1.23401	-0.124823
49	0.7803	1.22917	-0.107738
50	0.7888	1.22737	-0.103033
51	0.7955	1.21775	-0.0933598
52	0.8002	1.19866	-0.0968015
53	0.8087	1.17231	-0.0922126
54	0.8088	1.14301	-0.0921589
55	0.7314	1.11992	-0.135845
56	0.5563	1.10278	-0.254691
57	0.3765	1.08778	-0.423774
58	0.2379	1.07518	-0.623606
59	0.1445	1.06258	-0.840132
60	0.1102	1.0465	-0.957818
61	0.1047	1.031	-0.980053
62	0.1038	1.01953	-0.983803
63	0.106	1.00988	-0.974694
64	0.1073	0.997386	-0.9694
65	0.1092	0.984977	-0.961777
66	0.1163	0.978637	-0.93442
67	0.1292	0.977266	-0.888737
68	0.1578	0.977724	-0.801893
69	0.2009	0.97633	-0.69702
70	0.2286	0.968016	-0.640924
71	0.2342	0.953276	-0.630413
72	0.2135	0.933993	-0.670602
73	0.1657	0.908465	-0.786677
74	0.1323	0.882525	-0.87844
75	0.1216	0.864511	-0.915066
76	0.1047	0.854306	-0.980053
77	0.0778	0.843855	-1.10902
78	0.0552	0.82866	-1.25806
79	0.0436	0.810233	-1.36051
80	0.037	0.790988	-1.4318
81	0.0319	0.773055	-1.48621
82	0.03	0.755875	-1.52248

donde  
 $S_m(|Rt|)$  = Suavización del valor  
absoluto del residuo  
 $Log(S_m(|Rt|))$  = Logaritmo en base  
10 de  $Zt$   
 $Log(S_m(|Rt|))$  = Logaritmo en base  
10 de  $S_m(|Rt|)$   
año = 19..

Gráfica 5  
Gráfica del Diagnóstico



Cuadro 1  
ANALISIS DE REGRESION

independent variable	coefficient	std. error	t-value	sig.level
INSTANT	-2.239685	0.173582	-12.8731	0.0000
logsu	1.433265	0.147662	9.7064	0.0000

--R<sup>2</sup>. (Adj.) = 0.6084 SE= 0.245088 MAE= 0.200457 DurWat= 0.117  
 --R<sup>2</sup> previous: 0.0000 0.000000 0.000000 0.0000  
 --: observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 2  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	5.65922	1	5.65922	94.2136	.0000
Error	3.54401	59	0.0600680		
Total (Corrected)	9.20324	60			

--R-squared = 0.614917

--R-squared (Adj. for d.f.) = 0.50639

Std. error of est. = 0.245086  
 Durbin-Watson statistic = 0.117003

Tabla 8  
SUAVIZACION DE (-tasa de mortalidad)<sup>-0.5</sup> POR 4253ehd  
1922-1982

año	transmort	transuav	transres
22	-0.198811	-0.1968	-1.0693E-5
23	-0.202444	-0.1989	-1.5440E-3
24	-0.197442	-0.1991	1.45765E-3
25	-0.194257	-0.1994	5.14283E-3
26	-0.200401	-0.1995	-9.0120E-4
27	-0.204124	-0.1992	-4.9241E-3
28	-0.198611	-0.1979	-9.1069E-4
29	-0.193167	-0.1959	2.73315E-3
30	-0.193692	-0.1945	1.06521E-3
31	-0.196494	-0.195	-1.4943E-3
32	-0.19574	-0.1963	5.59927E-4
33	-0.197257	-0.1994	2.14254E-3
34	-0.20456	-0.2034	-1.5800E-3
35	-0.210352	-0.206	-4.3515E-3
36	-0.202644	-0.2071	8.15751E-4
37	-0.202444	-0.2076	5.15592E-3
38	-0.206569	-0.2079	-1.0659E-3
39	-0.206514	-0.2064	-2.1441E-4
40	-0.207614	-0.2086	9.863E-4
41	-0.212714	-0.2092	-3.5178E-3
42	-0.209427	-0.2106	1.17305E-3
43	-0.211289	-0.2138	2.51144E-3
44	-0.220324	-0.219	-1.3263E-3
45	-0.226455	-0.225	-1.4554E-3
46	-0.227038	-0.2314	4.3617E-3
47	-0.24544	-0.2375	-7.9403E-3
48	-0.243252	-0.2416	-1.65211E-3
49	-0.23636	-0.2429	6.54027E-3
50	-0.248452	-0.2434	-5.052E-3
51	-0.246424	-0.2463	5.67648E-3
52	-0.258199	-0.252	-6.1563E-3
53	-0.250785	-0.2601	9.31507E-3
54	-0.276289	-0.2689	-7.3694E-3
55	-0.272166	-0.2758	3.63447E-3
56	-0.294352	-0.281	-0.0113527
57	-0.280607	-0.2856	5.19323E-3
58	-0.286675	-0.2899	1.22467E-3
59	-0.296174	-0.2942	-1.9744E-3
60	-0.298807	-0.2998	9.92846E-4
61	-0.310087	-0.3052	-4.66648E-3
62	-0.308607	-0.3091	4.932E-4
63	-0.310067	-0.3126	-2.51316E-3
64	-0.317821	-0.3172	-6.2086E-4
65	-0.324443	-0.3217	-2.7428E-3
66	-0.322749	-0.324	1.25159E-3
67	-0.32969	-0.3244	-5.0962E-3
68	-0.322749	-0.3244	1.65139E-3
69	-0.326164	-0.3248	-1.3640E-3
70	-0.324749	-0.3278	5.05139E-3
71	-0.339032	-0.3335	-5.5317E-3
72	-0.3371	-0.3413	4.20007E-3
73	-0.349215	-0.3516	2.38485E-3
74	-0.365146	-0.3622	-2.9463E-3
75	-0.372678	-0.3695	-3.178E-3
76	-0.370517	-0.3738	3.66359E-3
77	-0.377364	-0.3784	4.35527E-4
78	-0.363482	-0.3651	1.61751E-3
79	-0.395285	-0.3936	-1.6847E-3
80	-0.40161	-0.4024	7.90356E-4
81	-0.411693	-0.4106	-8.5546E-4
82	-0.418854	-0.4169	4.60917E-3

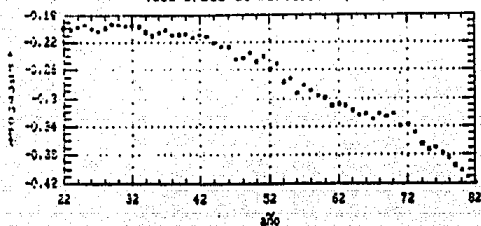
Donde

transmort = Tasa de mortalidad transformada

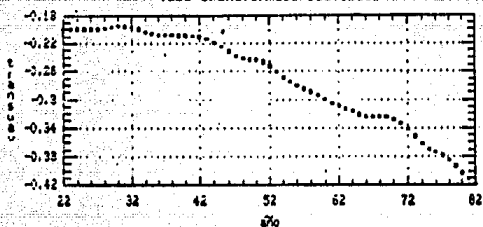
transuav = Tasa de mortalidad transformada suavizada por 4253ehd

transres = Residuo de la suavización

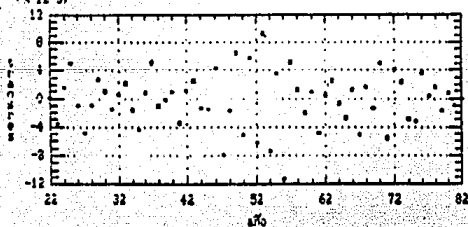
Gráfica 6  
Tasa bruta de mortalidad transformada



Gráfica 7  
Tasa transformada suavizada



(1E-3) Residuos de la Suavización



Gráfica 8

Una inspección cuidadosa a la gráfica 2 revela cambios en los datos como la tasa de mortalidad.

La tabla 6 ilustra los resultados de la suavización de estos datos por la técnica 4253EH, doble. En dicha tabla se muestran los valores suavizados ( $Z_t$ ) y los residuos ( $R_t$ ). La gráfica 3 muestra una tendencia a la baja en la mortalidad a lo largo de los años y la gráfica 4 indica que la distribución es aleatoria de los residuos al rededor del cero, pero con dos componentes de varianza. De 1922 a 1957, los residuos varían aproximadamente entre -1.3 y 1.7; Mientras que en el período 1958-1982, la variación cambia, los residuos fluctúan entre -0.3 y 0.3. Estos gráficos indican un tipo diferente de no-estacionaridad que se diagnostica a continuación.

Con el objeto de detectar la no-estacionaridad de estos resultados, en la tabla 7, se presentan los resultados obtenidos del cómputo de los logaritmos en base 10 de los datos suavizados y de los datos derivados de la suavización de los valores absolutos de los residuos. La gráfica 5 representa la curva  $\log(Sa|rt/1)$  versus  $\log(Z_t)$ .

Para determinar la pendiente de esta curva, se llevó a cabo un análisis de regresión completo (Ver cuadros 1 y 2). El valor de la pendiente (1.43) ajustada por mínimos cuadrados sugiere que  $p = -0.43 = -0.5 = -1/2$ , que es aproximadamente equivalente a la re-expresión  $-1/\sqrt{t}$ . Después de haber suavizado los datos transformados bajo esta re-expresión (tabla 8), se graficaron los resultados bajo

esta transformación (gráficas 6, 7 y 8). Al compararla con las de los datos no transformados con (gráficas 2, 3 y 4), se desprende una observación evidente: el claro mejoramiento en la regularidad de los datos.

Finalmente, con el objeto de comparar las dos tendencias de disminución de la tasa de mortalidad transformada suavizada, se calcularon los pendientes (-0.005 y -0.007 respectivamente), para los períodos 1952-1965 y 1970-1982.

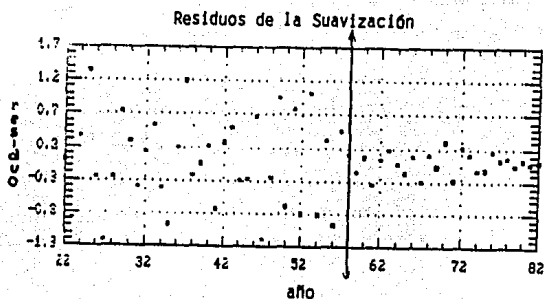
## 2- Interpretación de los resultados.

Los autores Martínez y Bustamante, (Almada Bay, 1975) señalan también que la erradicación de algunas enfermedades como el cólera, la peste bubónica, la viruela y la fiebre amarilla urbana; el abatimiento del tifo, la fiebre recurrente y el paludismo y la marcada reducción de las tasas de defunción por difteria, tos ferina y sarampión conforman un panorama optimista. Entre 1922 y 1930 la tasa bruta de mortalidad oscilaba entre 25 y 26 defunciones por mil habitantes. A partir de 1930, como puede verse en la gráfica 2, la trayectoria descendente de la mortalidad aparece con claridad: en 30 años, de 1930 a 1960, la tasa bruta de mortalidad se redujo en un 60% al pasar de 26.6 a 11.2 por mil. La tasa promedio entre 1965 y 1967 es de 9.4, que representa el 36.7% de la de 1930-34 y solo el 28.4% de la que prevalecía a principios de siglo.

A partir de la gráfica 4, se puede observar indirectamente un cambio en la mortalidad a partir de 1958. Aunque las condiciones de la población son difíciles de evaluar, dadas las estadísticas con que se cuenta; es posible formarse una idea de cuáles son los factores explicativos de este cambio, comparando las



condiciones de vida en 1930, 1940 y 1960.



Gráfica 4

Entre 1930 y 1960, existieron en México cambios importantes referentes a la mortalidad, misma época en que el país sufrió una serie de cambios económicos. El producto interno bruto per capita pasó de 1087 pesos anuales en 1935-1939, a 2134 en 1960-1964, lo que significa un porcentaje de cambio del 96.32%. La producción agrícola, pasó de 134.9 en 1940, tomando 1900=100, a 429.5 en 1960 y la producción manufacturera, con la misma base (1900=100) de 358.7 en 1940, a 1662.7 en 1960. Estas cifras dan un claro ejemplo del desarrollo económico, que ha existido en nuestro país. Aunque ha habido un desarrollo, éste no ha traído para

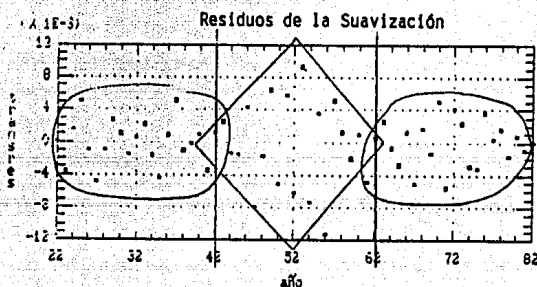
el grueso de la población mejoras substanciales, con lo que se podría plantear que el desarrollo económico no representa un factor explicativo importante del descenso de la mortalidad, aunque no hay que olvidar que sí ha ayudado al descenso de los fallecimientos.

En 1940, sólo un 50.2% de la población usaba zapatos, lo que denota que por lo menos una mitad de la población tenía una situación social deficiente. Este porcentaje aumentó para 1960 a 62.3%, quedando todavía un 40% de la población en condiciones sociales dudosas. Este modesto aumento contrasta con los logros en materia económica. El porcentaje de la población que no sabía leer y escribir, era en 1940 de 58%, pasando a ser el 38% en 1960.

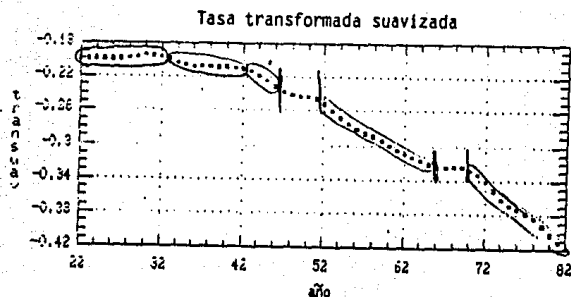
Otro indicador que muestra el pésimo estado de la sociedad mexicana durante ese período, es el porcentaje de viviendas sin drenaje. En 1940, el 86.5% de ellas no contaba con este servicio y en 1960 era el 71.1%. Estas cifras muestran que aunque en el período considerado existieron mejoramientos en las condiciones sociales, éstas no alcanzaron niveles satisfactorios. Siendo por lo tanto muy arriesgado plantear la hipótesis de que el cambio en la mortalidad se debió al mejoramiento de las condiciones de la población. Puesto que se pueda observar que existe una ligera baja de la mortalidad en el período presidencial de Avila Camacho.

Otra observación digna de mencionar es la presentada en la gráfica 8, en la que el comportamiento de los residuos presenta posiblemente un patrón oculto. La distribución parece indicar que existen tres comportamientos con respecto al

tiempo: 1922-1942, 1942-1962 y 1962-1982.



Por otra parte, como se aprecia en la gráfica 7, la disminución de la mortalidad está marcada por cuatro épocas bien distintas: la primera de 1922 a 1940, etapa "reformista", en ella empiezan a aplicarse la Reforma Agraria, renace la educación y se fortalecen las organizaciones obreras; la segunda es llamada de Estabilidad Política y de Avance Económico, que abarcaría de 1940 a 1954, ésta se acepta, como su nombre lo indica, como una de desarrollo económico. De 1955 a 1970: es época de gran auge tecnológico, y de 1970 a 1982 es una época de gran auge de tecnología médica.



Gráfica 7

En dicha gráfica, se señalan claramente períodos en que la tasa transformada suavizada se mantiene constante, otros en los que baja ligeramente y otros dos períodos en los que el descenso es más pronunciado. Habría entonces que profundizar la cuestión, llevando a cabo una descripción más detallada del marco histórico.

La tasa transformada suavizada se mantiene constante de 1922 a 1933, desciende ligeramente manteniéndose de nuevo constante durante el período 1934-1942. Aquí cabe señalar que hasta 1930-1935, la mortalidad representada por niveles altos y fluctuantes se relaciona con un período de guerras, epidemias y

de muy pequeño, casi nulo, desarrollo económico además de que en este período, las condiciones sociales de la población se mantuvieron por lo general estables y a niveles ínfimos. Sin embargo, en 1934, con el ascenso al poder del General Lázaro Cárdenas, las prestaciones sociales se agudizaron, manifestándose una gran agitación social.

En general, la reducción importante del ritmo de aumento de la vida durante los años 1922-1940 y la aparición de cierta reactivación a partir de 1940 aproximadamente deben relacionarse con los adelantos de la inmunización contra la viruela, de la alimentación y del nivel cultural de la población.

Por ello, se puede decir que el descenso claro de la mortalidad se da a partir de 1942, época en que principia el desarrollo de la economía mexicana. De 1946 a 1952, bajo el gobierno de Miguel Alemán, en realidad, el país vivió una época de crecimiento, pero crear riqueza no implicaba mejorar las condiciones en que habitaba la población. La mortalidad no cambió aunque está claro que continuaron los avances científicos, administrativos y culturales: la industrialización y la urbanización contrarrestaban los adelantos médicos y sanitarios.

Entre 1940 y 1950, la mortalidad se redujo de 23.2 a 16.5 por mil, y en el decenio siguiente se obtienen un descenso un poco menos grande: hasta obtener una tasa de 11.2 por mil. A partir de 1960, el ritmo de incremento tiende a ser más lento.

En los años 40's, la Administración Pública empezó a difundir en la población los principios de higiene y los conocimientos biológicos aportados por los descubrimientos de Pasteur. Se puede decir, que en el decenio de los cuarenta, empezaron en México las Campañas de Vacunación. En este tiempo, se usaron por primera vez en forma masiva, la vacuna BCG contra la tuberculosis, la de Salk y Sabine, contra la poliomielitis y otras contra la difteria y la varicela. También en estos años (1947), tuvo éxito el uso del DDT, que influyó notablemente en la reducción de muchas enfermedades, como el paludismo y otras infecciosas. Otro aspecto es que, el DDT logró mejorar los rendimientos agrícolas. Como se ve, una posible explicación del cambio en la mortalidad, podría ser la importación de tecnología médica.

De 1952 a 1965, los cambios en los niveles de la mortalidad se hacen más importantes. Este periodo corresponde a los períodos presidenciales de Ruiz Cortines y López Mateos. En contraste, los niveles de mortalidad a partir de esta fecha tienden a estabilizarse de 1966 a 1969, período correspondiente a la presencia de Díaz Ordaz en el poder. Aunque cabe aclarar que las condiciones de vida en esa época mejoran relativamente: la proporción de viviendas con drenaje aumento a un 41.15% en 1970, la proporción de personas que usan zapatos a un 80.1% y las personas alfabetas a 76.3%, y aunque los porcentajes alcanzados no son motivo de orgullo, si representan aumentos considerables.

Las implicaciones políticas de las conclusiones anteriores obligan a examinarlas con atención y a evaluar con gran cuidado la importancia relativa de

los factores socioeconómicos versus los otros más directamente ligados a la salud.

Los programas de prevención de la salud, el mejoramiento de las instalaciones médicas para hacerlas accesibles a los grupos de alta mortalidad que todavía no tienen acceso a ellas, y los mejoramientos de las condiciones ambientales que afectan los estándares de vida se ven seriamente obstaculizados por las actuales tendencias en la distribución y redistribución de la población. La distribución desigual de la salud y de los servicios sociales básicos a lo largo del territorio nacional, su concentración en la capital y otros núcleos urbanos, el poco acceso que a ellos tiene la población rural dispersa, son hechos bien conocidos. La velocidad de la urbanización y de la concentración urbana hace muy difícil que el gobierno pueda cumplir con las necesidades urbanas crecientes de vivienda, abastecimiento de agua potable, servicios de salud, instalaciones de alcantarillado. La consecuencia de esto es que grandes segmentos de la población urbana están viviendo ahora bajo condiciones altamente favorables para la aparición de enfermedades infecciosas y parasitarias, situación que se agrava en muchos casos por la desigual distribución de los servicios médicos dentro de los límites de la ciudad.

En resumen, desde un punto de vista político, parece clara ahora que los futuros mejoramientos en la esperanza de vida al nacer, así como los mejoramientos más específicos en la mortalidad infantil y en la niñez requieren que los programas de salud sean incluidos en políticas socioeconómicas y de población más amplias, tendientes a elevar el nivel de vida de los grupos de

mortalidad alta tanto en áreas rurales como urbanas. Por el contrario, no parece razonable esperar que los programas aislados de salud conduzcan en el futuro, y mientras la importancia de diversos tipos de causas de muerte no se altere, a descensos constantes y significativos de la mortalidad.



## 8- LA ESPERANZA DE VIDA EN 1975.

### 1- Presentación de los resultados.

Ante las dificultades que plantea el análisis de dos variables (esperanza de vida femenina y masculina) para 32 objetos, con el agregado de la dimensión edad, se siguió el siguiente camino.

Las tablas 5a y 5b son unas tablas de doble-entrada de  $32 \times 18$  que facilitan el análisis de similitudes de la esperanza de vida por entidad federativa y por edad.

Aunque éste es un análisis de doble-entrada de un modelo de varianza, un ANOVA estándar no es muy aconsejable; puesto que esta técnica es poco resistente a los casos extremos. (Curts, 1984)

Un modelo natural para usar en tales circunstancias es el modelo aditivo usado en análisis de varianza de doble-entrada:

$$Y = \mu + e.f. + edad + e$$

Aquí, el término común  $\mu$  es un nivel general de la esperanza de vida de toda la tabla. Los términos del renglón e.f. y de la columna edad son valores que miden la desviación de cada entidad federativa (renglón) y por edad (columna) para el nivel general.

Por lo tanto, a partir de los indicadores de la esperanza de vida presentados en las tablas 5a y 5b, se operaron estos datos por el método de la mediana pulida

de un análisis de tablas de doble-entrada. Las tablas 9a y 9b muestran los resultados después de 10 iteraciones del proceso del "pulimento" de las tablas 5. Se dispone de los residuos del ajuste para una futura inspección; seguramente un diagrama de tallo-y-hoja es de gran utilidad para observar la forma general de los residuos y un diagrama de caja para mostrar el rango de valores que los residuos cubren y determinar dónde se concentran la mayoría de éstos.

Para tratar de recobrar los ciclos regionales indicados en las gráficas 9a y 9b, la mediana pulida ajustó iterativamente los renglones y las columnas de la tabla original hasta obtener una mediana igual a cero, restando las medianas por renglón y por columna. El resultado es una descripción de las tablas de acuerdo a las ecuaciones (6).

Los efectos renglón y los efectos columna resultantes (tablas 10a y 10b) se presentan en los cuadros 7 y 8. Los efectos-renglón muestran los cambios de los datos, región-por-región. Similarmente, los efectos-columna describen patrones edad-por-edad.

Y a fin de visualizar los efectos de una manera más clara, se dibujaron los diagramas de tallo-y-hoja en espejo, por sexo. (Figuras 1 y 2). Al examinar dichas figuras, se puede apreciar que ambos diagramas muestran que las distribuciones de los efectos son simétricas. Y la figura 1 contiene unos casos extremos que habría que estudiar detenidamente.

MEIANA POLIDA DE LA ESPERANZA DE VIDA.

Sexo Masculino.

E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-0.237136	-0.244472	-0.440613	3.72962E+4	0.9594328	0.115401	0.650620	1.026001	0.0994691	0.0410433	7.23441E+3	-8.32778E+4	-0.226793	-0.189799	-0.430752	-0.245841	-0.027799	-0.146641
2	3.465044	6.274364	-4.44181E-3	-0.159428	-0.210249	-0.1042	0.0167023	0.0307998	8.60792E-3	-1.78728	-0.342967	-0.200454	-0.204594	-0.0496001	0.9572466	0.216958	0.4744	0.283558
3	1.136E+5	1.74133	1.61554	1.47254	1.60442	1.20278	0.853644	0.497784	0.134832	-0.221774	-0.505782	-0.45367	-0.72941	-0.211616	-0.143749	-0.496536	-0.864616	-0.829546
4	6.347884	0.745143	1.69434	0.371469	0.340423	0.536597	0.3675	0.301597	0.139443	-1.17394	-0.211749	-0.251857	-0.485797	-0.227803	-1.12994	-1.460424	-1.5028	-1.25464
5	0.347662	0.440541	0.240154	0.647147	0.0162268	0.0423598	-0.0167023	-0.0726048	-0.254537	-0.302183	-0.545571	-0.456459	-0.389599	-0.444005	-0.284158	0.0815324	0.521955	0.510154
6	-0.156425	-0.269526	0.274277	0.36129	0.34025	0.284518	0.207421	0.174518	0.195946	0.0419404	-0.0418442	-0.0493256	-0.175876	-0.237482	-0.500035	-0.824536	-0.512892	0.504277
7	-1.27018	-1.25044	-1.10564	-0.844831	-0.679731	-0.654073	-0.5825	-0.476073	-0.3907547	0.0504753	0.148213	0.420143	0.414203	0.482197	0.490044	0.423753	0.427197	1.13436
8	0.277528	0.198784	0.348	0.115014	0.0840732	0.170242	0.141144	0.155242	0.113289	-0.0445162	-0.224425	-0.344412	-0.452133	-0.424158	-0.676212	-1.1804	-1.19914	-0.852
9	1.4024	1.10448	0.724652	0.495729	0.428765	0.339933	0.260476	-0.150647	-0.452018	0.685464	-0.713413	-0.54152	-0.397461	-0.189467	-0.10162	-0.104092	0.275513	-0.117308
10	1.54975	0.374029	0.17244	0.0422236	1.31331E-3	-2.5181E-3	-0.0361158	-0.0715183	1.10553	-0.15424	0.145115	0.151628	-4.9124E-3	-0.0969183	-0.419072	-0.43336	-0.301918	-1.75476
11	-1.94625	-0.46277	6.4619E-3	-0.0163446	0.0225131	0.0468034	-0.9504141	-0.104217	-0.0884879	0.0414257	0.146217	0.37223	0.506289	0.534283	0.49213	-0.0421583	-0.150737	-2.55821E-3
12	0.014964	-1.20543	-1.04644	-1.05963	-1.04037	-1.0042	-0.772398	-0.38972	-0.1532E-3	0.221442	0.574746	0.949344	1.04941	1.1116	0.82247	0.144758	-0.0236001	3.5584E-3
13	-1.44814	-2.37948	-1.75064	-1.44163	-1.50437	-1.3184	-0.9475	-0.643403	-0.115355	0.117747	0.514241	1.03714	1.3872	1.4872	1.97504	2.05074	3.3022	2.45914
14	-0.619427	-0.624172	0.195043	0.184279	0.251138	0.307307	0.208209	0.112037	3.54671E+4	-0.0474312	-0.0914598	5.52774E+4	4.31244E-3	0.142907	-0.0504446	-0.313333	-0.281093	-0.314735
15	-3.7118	-0.550232	0.0481901	3.90348E-3	-0.0450346	-0.0179458	-0.0338483	-0.0179458	-0.14582	0.0704146	-7.23441E+3	0.224476	0.428738	0.496732	0.804579	1.09029	1.19173	0.53889
16	-0.040694	-1.03679	-1.07357	-1.11056	-1.04503	-1.07313	-1.00333	-0.890333	-0.182285	0.0704146	0.244414	0.449213	0.562273	0.630267	0.648114	0.513825	0.525247	0.202425
17	0.0191249	-1.25479	-1.31737	-1.39054	-1.43151	-1.25233	-0.618023	-0.224285	-0.0499977	0.149401	0.518413	0.592473	0.710247	1.06411	1.15793	1.40247	0.712425	0.712425
18	1.149578	0.141043	0.0807147	2.39027E-3	0.0231499	2.5181E-3	-0.106379	0.0715183	0.0455644	0.05344E-3	-0.0493461	-0.0219555	0.0112142	-0.151892	-0.408123	-0.404482	-0.489753	
19	0.44444	1.74014	1.51954	1.32637	1.28547	1.2414	1.0525	0.694597	0.254645	-0.2494	0.636749	-1.064557	-1.1108	-1.2618	-1.7494	-2.02124	-1.4478	-0.365444
20	-0.21114	-0.80564	-0.846044	-2.32323	-2.05457	-1.9484	-1.3175	-0.953403	-0.453353	0.447077	1.17343	1.39714	2.3642	2.8572	3.51504	3.70074	3.5822	3.13974
21	-0.04014	-2.41484	-1.80564	-1.42463	-1.49357	-1.3134	-1.0923	-0.708403	-0.280353	0.447077	1.17343	1.39714	2.3642	2.8572	3.51504	3.70074	3.5822	3.13974
22	-0.04197	-1.29445	-0.521442	-0.524428	-0.465349	-0.3732	-0.31874	-0.3842	-0.116132	0.147442	0.444414	0.544246	0.504406	0.44644	0.724247	0.709758	0.4814	0.118958
23	1.54781	1.014768	0.118	0.6820125	0.6340732	0.460242	0.441144	0.375242	0.23239	-0.048144	-0.24915	-0.344412	-0.452133	-1.409614	-2.40071	-3.4504	-4.19914	-3.472
24	-1.21857	-1.30545	-0.82442	-0.744248	-0.610348	-0.5442	-0.4605298	-0.4292	-0.121152	0.174442	0.317454	0.505346	0.561404	0.7314	0.689247	0.454958	0.4544	0.113359
25	0.36047	0.173446	-0.105242	-0.211238	-0.243144	-0.0449994	0.032036	0.108	0.144044	-0.0255576	-0.0753462	-0.0346237	-0.079594	0.0404601	0.0844499	0.6421502	0.0236001	-0.475242
26	0.84244	0.614601	0.42461	-0.72914E+4	-1.3132E-3	0.124855	0.165758	0.0998555	-0.0420948	0.23575	0.113311	0.61559	-0.447339	-0.389545	-0.401658	-0.205107	0.165455	0.222413
27	-0.45123	-1.07503	-0.70141	-0.51428	-0.405248	-0.1891	0.0816629	0.4759	0.863448	1.04404	1.14433	1.27445	1.07851	0.8645	0.434347	-0.0795414	-0.2283	-0.301341
28	-1.33147	0.443483	0.444254	0.523169	0.300429	0.204597	0.1173	0.205397	0.114445	-0.107744	-0.111749	-0.319957	-0.343757	-0.18745	-0.219956	-0.214215	-0.242803	-0.95484
29	-2.50414	-0.844537	-0.19944	-0.148131	-0.235751	-0.223403	-0.1425	-0.0784028	3.54457E+4	-0.05944E-3	0.148213	0.501415	0.744023	0.852197	0.700044	0.825735	0.847797	0.148454
30	-0.613904	-1.19179	-0.634573	-0.415641	-0.744501	-0.480313	-0.52443	-0.355333	-0.237285	0.0150193	0.241401	0.633213	0.817243	0.975247	0.373114	0.088425	0.470267	0.215445
31	0.65244	0.151643	0.74054	0.441349	0.380429	0.584397	0.3673	0.291597	0.119463	-0.15714	-0.201749	-0.329837	-0.455737	-1.0174	-1.21994	-1.75424	-1.7428	-1.75424
32	-0.34054	0.360774	-2.19001	0.157	0.510599	0.492228	0.42313	0.717228	0.37574	0.15747	-0.0434412	-0.6442262	-0.250167	-0.312172	-0.844326	-1.27861	-1.48717	-1.96001

Tabla 9a

MEDIANA PULIDA DE LA ESPERANZA DE VIDA.

Sexo Masculino.

E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-0.217174	3.2441672	-6.460013	3.7298214	-0.0594226	0.117501	0.0563025	0.120601	0.0186491	0.0618037	7.286438-3	-0.327774	-0.226793	-0.189795	-0.430952	-0.305861	-0.093799	-0.144641
2	0.603264	6.373266	-6.441871	-0.159428	-0.210269	-0.1424	0.0167023	0.0370798	8.849218-3	-0.178750	-0.262567	-0.300454	-0.306934	-0.0496001	0.0576166	0.231658	0.4764	0.283558
3	1.33653	6.74133	1.61056	1.43756	1.46641	1.20278	-0.853486	0.497597	0.149432	-0.231776	-0.305782	-0.49347	-0.73961	-0.211616	-0.143769	-0.639058	-0.856618	-0.839458
4	0.342484	0.763163	1.43438	0.0871871	0.0182289	0.0623352	0.303297	0.301597	0.139465	-0.211769	-0.259387	-0.485797	-0.812005	-0.284158	0.0813536	0.312319	0.510194	0.510194
5	0.347664	0.665461	0.260156	0.0716271	0.0182289	0.0623352	0.303297	0.301597	0.139465	-0.211769	-0.259387	-0.485797	-0.812005	-0.284158	0.0813536	0.312319	0.510194	0.510194
6	-0.159475	-0.136946	0.274277	0.381229	0.360253	0.286318	0.207421	0.213318	0.153566	0.0419606	-0.0418844	-0.0493236	-0.175876	-0.237482	-0.540025	-0.816124	-0.212862	0.266277
7	-1.33945	-2.25646	-1.105664	-0.848821	-0.479571	-0.433403	-0.2822	-0.159009	-0.0902587	0.0502155	0.148326	0.420143	0.414829	0.472197	0.460064	0.423255	0.627197	1.13164
8	0.217558	0.198786	0.248	0.115014	0.0940763	0.170212	0.161144	0.157242	0.11329	-0.0491444	-0.231625	-0.364112	-0.452153	-0.468158	-0.676312	-1.1806	-1.19516	-0.852
9	1.2224	1.01048	0.723693	0.459703	0.428763	0.324932	0.193836	-0.150067	-0.452018	-0.684562	-0.734322	-0.512152	-0.397611	-0.189467	-0.10162	0.104098	0.275313	-0.117308
10	1.24753	0.376216	0.17521	0.0422236	1.113118-3	-2.51818-3	-0.0316158	-0.0135183	0.12053	1.252626	0.183119	0.151028	-0.91264-3	-0.0889189	-0.419072	-0.43316	-0.201310	-1.05679
11	-0.36003	-0.46277	6.41819-3	-0.0184616	0.025131	0.0484034	-0.0504161	-0.106717	-0.092405	-0.0461257	0.146217	0.37223	0.506289	0.536283	0.49213	-0.041858	-0.130517	0.558218-3
12	0.0126444	-1.20565	-1.04644	-1.03963	-1.00037	-1.00621	-0.773298	-0.5992	-1.1528-3	0.271442	0.367436	0.193364	1.09161	1.1114	0.829267	0.184958	-0.0236001	3.598218-3
13	-1.48114	-2.27586	-1.75044	-1.64342	-1.50457	-1.3184	-1.0475	-0.643602	-0.115755	0.117509	0.50444	1.03514	1.3892	0.8712	1.97534	2.05076	2.3042	2.36936
14	-0.413627	-0.0241672	0.195043	0.180079	0.201138	0.307307	0.408809	0.412207	1.366718-4	-0.0472312	-0.0913589	5.327761-4	4.912444-3	0.142807	-0.312325	-0.281093	-0.316935	0.23869
15	-1.17171	-0.520222	0.088801	5.90368E-3	-0.0450246	-0.0386881	-0.0179558	-0.0228483	-0.16382	-0.074664	-7.2346E-3	0.246576	0.428728	0.496722	0.806379	1.09029	1.19173	0.23869
16	-0.084664	-1.03679	-1.07757	-1.11056	-1.0715	-1.00323	-0.81643	-0.90733	-1.16285	-0.701019	-0.466463	0.482213	0.526273	0.450267	0.688114	0.523825	0.525267	0.202425
17	0.0519256	-1.23879	-1.23757	-1.39056	-1.4515	-1.25223	-1.01443	-0.800358	-0.322828	-0.226987	0.186401	0.516133	0.762273	0.710267	1.00811	1.15783	1.20327	0.712625
18	1.44578	0.140385	0.082747	2.9027E-3	0.0267499	2.9181E-3	-0.106378	0.0175183	0.0435444	-0.23841-3	-0.3654981	-0.2293355	0.0101262	-0.191882	-0.726215	-0.803223	-0.604882	-0.489223
19	0.41646	1.74614	1.51916	1.32637	1.28561	1.2416	1.0325	0.696397	0.256465	-0.154781	-0.436769	-0.806437	-1.11108	-1.2828	-1.78446	-2.02916	-1.4478	-0.38344
20	-4.31114	-6.96786	-2.86664	-2.32263	-2.03457	-1.5404	-1.7179	-0.93403	-0.245255	0.607029	1.17442	1.31514	2.3642	2.9372	3.15004	3.5074	3.7622	3.13936
21	-3.98918	-2.81886	-1.90561	-1.49537	-1.14937	-1.1736	-1.0235	-0.709403	-0.280352	0.488029	0.749251	1.46034	1.9362	2.4522	2.76064	2.92576	2.8712	2.56636
22	-2.40153	-1.39665	-0.521462	-0.324228	-0.465348	-0.7732	-0.81828	-0.3862	-0.116152	0.176142	0.444316	0.586166	0.500406	0.8884	0.726247	0.709708	0.6814	0.18939
23	1.50751	0.158768	0.111	0.0501305	0.0367732	0.160242	0.481144	0.275262	0.12329	-0.024124	-0.40145	-0.166212	-0.852153	-1.40416	-2.40631	-3.6326	-4.19918	-1.472
24	-1.21939	-1.30565	-0.828462	-0.749428	-0.610368	-0.5442	-0.601298	-0.4392	0.121152	0.173142	0.287436	0.501346	0.561606	0.7316	0.689267	0.449458	0.0764	0.113559
25	2.39017	0.171564	-0.107442	-0.212228	-0.241181	-0.0683999	0.0329026	0.109	0.186048	-0.0135757	-0.073462	-0.0363317	-0.079394	0.0496001	0.0864693	0.443583	0.0236001	-0.478142
26	0.821142	0.11363	0.125417	0.7338E-4	1.3121E-3	1.54825	0.185758	0.0989751	-0.0420948	-0.13763	-0.403156	-0.411959	-0.467579	-0.395959	-0.401648	-0.205947	-0.164555	0.221613
27	-0.81813	-1.70555	-0.701341	-0.516228	-0.405248	-0.1891	0.081629	0.4795	0.613948	1.04834	1.11153	1.27445	1.07851	0.8865	0.438467	-0.079914	-0.2285	0.30141
28	1.53216	0.483163	0.448356	0.328369	0.300643	0.206397	0.2375	0.201597	0.156475	-0.107181	-0.111769	-0.19357	-0.265797	-0.197803	-0.279956	-0.212425	-0.22820	-0.85644
29	-2.5011	-0.584837	-0.115644	-0.138631	-0.239571	-0.223603	-0.1525	-0.0784028	-3.567E-4	2.39986E-3	0.148431	0.240143	0.746203	0.821197	0.700044	0.832755	0.847197	0.14036
30	-0.019044	-1.19179	-0.832755	-0.815961	-0.746301	-0.680232	-0.52942	-0.355313	-0.237885	-0.0151093	0.221401	0.32313	0.887273	0.973287	0.373114	0.818825	0.870267	0.716245
31	0.430466	0.513163	0.704756	0.643249	0.606419	0.646397	0.3875	0.281597	0.119645	-0.117861	-0.201769	-0.22957	-0.635797	-1.0178	-1.11996	-1.79426	-1.7428	-1.25344
32	-0.240504	0.360776	-0.23601	0.297	0.516059	0.682228	0.62113	0.717218	0.15276	0.15767	0.0426182	-0.0442162	-0.250167	-0.512172	-0.846236	-1.23841	-1.08717	-1.76001

Tabla 9a

MEJORA POR LÍNEA DE LA ESPERANZA DE VIDA.

Sexo Femenino.

E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-0.432706	0.018761	-0.002961	-0.074779	0.014838	0.000422	0.147131	0.233523	0.294364	0.230496	0.214564	0.099934	-0.019099	-0.145036	-0.148273	-0.172172	-0.160997	-0.155337
2	0.374534	0.012210	0.115355	0.071071	0.064136	1.122512	-1.340081	-0.064912	-0.111955	-0.225005	-0.212955	-0.126264	-0.212955	-0.212955	-0.121205	0.313128	0.441325	0.415464
3	1.27381	1.27167	0.930740	0.07042	0.71153	0.457136	0.453064	0.230137	0.111258	-0.116112	-0.140742	-0.273353	-0.288322	-0.270361	-0.236634	-0.315450	-0.433262	-0.436223
4	0.634749	0.300111	0.081409	0.161361	-0.037900	0.090138	0.043799	0.912973	0.030801	0.029709	-0.015649	-0.060099	-0.14471	-0.287739	-0.417739	-0.779953	-1.21462	-1.29442
5	0.634749	0.42441	0.081409	0.161361	-0.037900	0.090138	0.043799	-0.023199	-0.108211	-0.117001	-0.19523	-0.277601	-0.324411	-0.277601	0.382309	0.315483	1.09748	0.442119
6	-0.146607	-0.70041	0.701914	-0.13542	0.131472	0.117254	0.057984	0.040294	0.021302	-0.359929	-0.064619	-0.13223	-0.11824	-0.45924	-0.144511	-0.223336	0.400078	0.701399
7	-0.246419	0.24025	-0.72922	-1.27990	-1.06492	-0.82864	-0.464136	-0.197812	0.237796	-1.915914	0.012378	0.146447	0.421418	1.32146	1.72136	1.75454	1.72474	1.51118
8	0.123604	0.26147	0.230634	0.16279	0.151426	0.147014	0.117713	0.050117	-0.208634	-0.170313	-0.188664	-0.203474	-0.050034	-0.302118	-0.134794	6.420119	-3.02418	-0.439453
9	1.20043	1.29011	0.677730	0.357044	0.448134	0.44778	0.260209	0.224081	0.379022	-0.23546	-0.27409	-0.191738	-0.111952	0.420119	-0.044811	0	0	0
10	1.55381	0.81646	0.430740	0.30042	0.25123	0.257136	-0.113064	0.110317	0.041278	-0.040132	-0.148742	-0.24293	-0.229382	-0.24293	-0.418634	-0.32454	-0.372262	-0.728223
11	-0.206172	-0.34311	-0.607011	-0.607011	0.050790	0.062349	0.062053	0.050799	0.050799	0.050799	0.040189	-0.041317	0.040189	-0.115111	-0.219415	-0.39437	-0.24064	-0.079621
12	-0.062190	-1.29423	-0.935419	-0.805737	-0.754448	-0.410611	-0.432133	-0.34574	-1.729116	0.173811	0.112461	0.200619	0.359621	0.359621	0.465364	0.31954	0.030742	0.175721
13	-1.73119	-2.20253	-1.79423	-1.52454	-1.43242	-1.16796	-0.801136	-0.536743	-0.163743	0.164006	0.444238	1.07145	1.51662	1.79442	2.15434	2.01954	2.16174	2.83618
14	-0.34781	0.22065	0.22913	0.108002	0.107932	0.145311	0.062344	0.038619	-0.070360	-0.11211	-0.100136	-0.074970	0	0	0.0797456	-0.027075	-0.014879	-0.41142
15	-0.06192	-0.73646	-0.134900	-0.173316	-0.014186	-0.1186	-0.041871	0.301773	-0.024479	0.0146723	0.013521	0.200912	0.419882	0.645882	0.645882	0.778004	0.631001	-0.358812
16	0.605452	-0.071077	-0.222649	-0.23279	-0.251806	-0.20622	-1.18991	-0.103119	-0.097812	-0.114652	0.117602	0.167292	0.164262	0.042424	0.100600	-0.016811	3.282071	0.127421
17	0.813006	-0.22053	-0.475252	-0.44536	-0.41848	-0.132064	-0.264136	-0.295913	-0.148742	-0.080192	0.0812378	0.216447	0.361418	0.431418	0.51246	0.57462	0.614738	0.011748
18	1.62407	0.44173	0.331009	0.330481	0.391811	0.167397	0.044136	0.964294	0.045134	0.5399232	-0.048406	-0.193091	-0.31812	-0.30812	-0.544379	-0.955197	-0.973	-0.348562
19	1.48061	1.09447	0.777640	0.43542	0.57452	0.562136	0.448064	0.315237	0.062378	-0.095192	-0.293742	-0.448323	-0.483382	-0.442382	-0.46436	-0.360456	-0.609262	-1.10382
20	-0.91419	-0.20212	-0.34323	-2.23950	-2.20643	-1.97286	-1.46416	-0.849763	-0.258712	0.239808	0.811254	1.28645	1.58162	1.11162	2.41136	2.30434	2.71674	2.87118
21	-0.33119	-0.20212	-0.74623	-1.44454	-1.31415	-1.07386	-0.701136	-0.247463	-0.073712	0.071808	0.246250	0.361647	0.726418	1.15642	1.53634	1.52571	1.51174	1.74618
22	-1.54843	-0.70047	-0.25437	-0.21165	-1.20053	-0.894913	0.631793	0.318132	0.531533	0.237093	0.349132	0.164264	0.949193	0.895913	-0.207415	-0.32564	-1.17337	0.419072
23	1.09129	0.74952	0.296232	0.257904	0.129034	0.044202	1.340841	0.017713	-1.237712	-0.022706	-0.012278	0.044136	-0.297514	-0.23049	-0.001152	-1.44737	-1.93778	-2.26134
24	-1.93119	-0.86252	-1.67423	-1.34434	-1.15345	-1.01796	-0.861136	-0.487463	-0.127372	0.174006	0.496258	0.781647	0.876418	0.736418	0.736418	0.419543	0.341738	0.24177
25	2.22203	0.31606	0.079643	-1.361312	-2.317912	-0.043263	0.040203	0.108136	0.049163	0.0402623	-0.000527	-0.133134	-0.240143	-0.180163	-0.170419	-0.11741	-0.224043	-0.740053
26	0.82151	0.82004	0.180925	0.097970	0.091974	-1.464412	-0.026737	-0.026737	-0.113564	-0.248618	-0.29944	-0.287973	0.207415	0.295004	0.207415	0.719451	0.174118	0.807353
27	-1.24928	-0.37518	-0.72908	-0.20216	-0.571186	-0.3104	-0.264871	-1.984214	-0.240522	0.239012	0.570522	0.79312	0.99082	0.84082	0.59628	-0.043922	4.411814	-0.043922
28	1.44634	0.74003	0.413281	0.31255	0.264085	0.204971	0.206199	0.102774	0.712722	-0.187457	-0.256207	-0.230818	-0.125847	-0.095471	-0.144101	-2.92311	-0.170727	-0.90480
29	-0.26093	-0.37327	-0.263991	-0.312319	-0.233189	-0.207643	-1.106874	-0.034016	0.074314	0.0320494	0.064314	0.241909	0.65688	0.51688	0.246203	0.469603	0.352	0.28162
30	-0.326492	-1.47323	-1.11995	-1.01028	-0.94913	-0.93364	-0.436462	-0.240462	-0.049413	0.129109	0.070543	0.20548	0.290919	0.410919	0.460464	0.232642	0.146039	0.170477
31	0.462184	0.41942	-0.701914	-0.613701	-0.800643	0.032578	0.246412	-0.844129	-0.620360	-0.11811	-0.21036	-0.144371	0	0	0	0.069194	1.212511	0.12312
32	-0.67951	-0.41951	-0.049429	0.031012	0.217112	-0.034121	-0.074535	-1.080612	-0.21994	0.24649	0.20994	0.39529	0.4403	0.2203	-0.109954	-0.211774	-0.27479	-0.450411

Tabla 9b



Figura 4a

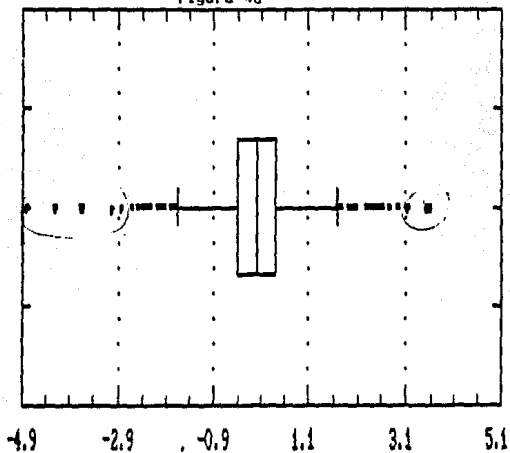


DIAGRAMA DE CAJA

Vector de Residuos - Sexo Masculino

Figura 4b

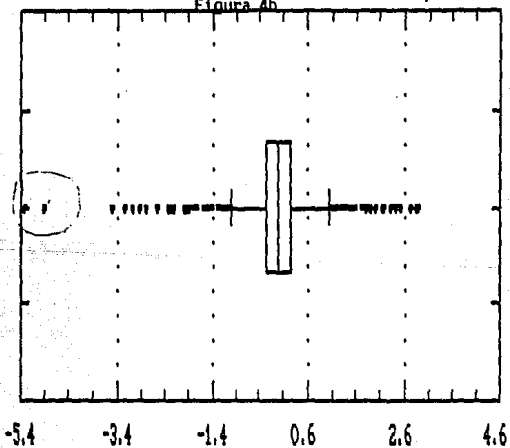


DIAGRAMA DE CAJA

Vector de Residuos - Sexo Femenino

Cuadro 7  
EFECTOS BENGLON

ENT. FED.	SEXO	
	MASC.	FEM.
AGS.	0.15	-0.18
B.C.	-0.30	-1.13
B.C.S.	1.10	0.81
CAMP.	1.46	0.70
COAH.	-0.70	-0.82
COL.	-1.01	-1.34
CHIS.	-1.52	-1.94
CH.	0.89	0.64
D.F.	0.26	0.30
DGO.	1.54	0.87
GTO.	-0.02	-0.24
GRU.	0.15	-0.03
HGO.	-2.62	-3.60
JAL.	0.40	-0.15
MEX.	-0.50	0.02
NICH.	-0.08	-0.30
NOR.	-0.64	0.10
NAV.	0.57	0.10
N.L.	-0.11	1.44
OAX.	-3.48	-3.73
PUC.	-2.82	-1.48
QRO.	0.05	-1.57
Q.R.	4.32	2.96
S.L.P.	-0.00	-0.48
SIN.	0.86	1.01
SON.	-0.37	-0.70
TAB.	0.53	-0.03
TAMPS.	1.36	1.32
TLAX.	-0.01	0.38
VER.	-0.70	-0.33
YUC.	1.29	-0.74
ZAC.	2.56	0.51

Cuadro 8

EFECTOS COLUMNA  
EDAD

SEXO	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
Sexo Masc.	29.4	32.1	28.9	24.3	19.5	14.9	10.6	6.3	2.0	-2.0	-6.1	-10.1	-13.9	-17.5	-20.7	-23.6	-26.4	-28
Sexo Fem.	31.3	33.3	30.3	25.5	20.7	15.9	11.3	6.7	2.2	-2.2	-6.5	-10.8	-15.0	-18.9	-22.5	-25.8	-28.8	-31



## VALORES DE COMPARACION - Sexo masculino

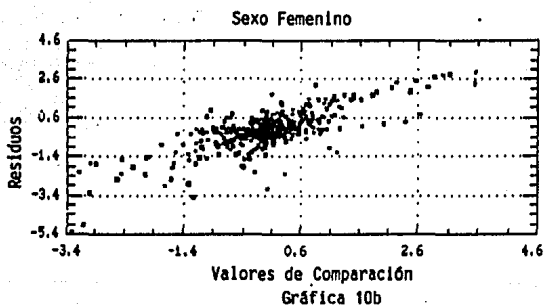
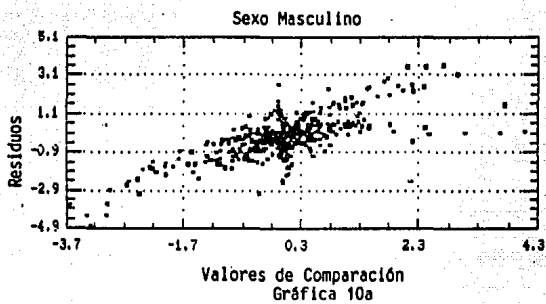
E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	0.11615	0.12328	0.12836	0.13168	0.043479	0.043809	0.045033	0.028132	0.705222	-0.45047	-0.033734	-0.042178	-0.0391664	-0.074505	-0.088621	-0.101132	-0.113008	-0.122949
4	-0.41637	-0.28613	-0.25736	-0.21678	-0.17912	-0.13062	-0.0941726	-0.0558594	-0.0161256	0.016416	0.0941288	0.0957464	0.123161	0.136051	0.149483	0.160687	0.155594	0.226179
8	0.94983	0.10359	0.39043	0.78122	0.82823	0.68076	0.340185	0.301795	0.0653117	-0.068139	-0.195211	-0.184925	-0.466721	-0.563718	-0.661781	-0.710179	-0.75105	-0.72341
9	1.12726	1.36537	1.32203	1.07949	0.831878	0.620637	0.281154	0.0987476	-0.064611	-0.258087	-0.432162	-0.591607	-0.718406	-0.802321	-0.882121	-1.00779	-1.12474	-1.22518
3	-0.56482	-0.44973	-0.394917	-0.31368	-0.23934	-0.202212	-0.213858	-0.181852	-0.0481184	0.0042151	0.128266	0.252121	0.280287	0.15432	0.419173	0.478453	0.520018	0.58176
6	-0.61844	-0.4566	-0.635464	-0.72152	-0.580128	-0.49002	-0.316201	-0.165312	-0.0405078	0.0682214	0.180518	0.301496	0.412658	0.506457	0.415432	0.702946	0.786046	0.874726
7	-1.31165	-1.48427	-1.25127	-1.08475	-0.87178	-0.647181	-0.421106	-0.200041	-0.0909186	0.081493	0.271235	0.453623	0.62052	0.761135	0.725376	1.02626	1.16111	1.1882
8	0.716267	0.824613	0.754788	0.626264	0.509128	0.389977	0.27959	0.181648	0.093182	-0.236181	-0.138399	-0.1468	-0.36632	-0.472284	-0.560894	-0.617528	-0.630174	-0.750494
9	0.21173	0.26183	0.217561	0.182382	0.148758	0.112798	0.0794564	0.0471389	0.0130032	-0.015007	-0.0487626	-0.0765211	-0.104368	-0.131681	-0.155737	-0.177744	-0.198821	-0.218171
10	1.32423	1.45233	1.30749	1.05837	0.802225	0.672555	0.478063	0.285522	0.121816	-0.078762	-0.438714	-0.627839	-0.782129	-0.938394	-1.04951	-1.19376	-1.30043	-1.40043
11	-0.016749	-0.0160461	-0.0146532	-0.0121618	-0.758972	-7.46772	-5.28432	-3.13652	-1.01766	1.11112	2.02705	3.07011	4.36651	5.75671	6.61037	0.616823	0.012201	0.013752
12	0.125993	0.127151	0.123468	0.103781	0.032864	0.0637328	0.0491623	0.0267767	0.653366	-6.81723	-0.0238442	-0.0421819	-0.0518876	-0.0486811	-0.100993	-0.129234	-0.128234	-0.12801
13	-2.31396	-2.44632	-2.21861	-1.84363	-1.49732	-1.19428	-0.811113	-0.481132	-0.154202	0.15322	0.494164	0.77316	1.04528	1.34608	1.58181	1.61887	1.02115	2.20649
14	0.34169	0.386286	0.29237	0.257191	0.21135	0.176855	0.125169	0.076354	0.0241046	-0.019311	-0.078172	-0.130169	-0.184391	-0.207413	-0.165318	0.36579	0.38473	0.426482
15	-0.61616	-0.46918	-0.422955	-0.35516	-0.285589	-0.218239	-0.156443	-0.091292	-0.029781	0.0002917	0.0844776	0.148191	0.204104	0.256159	0.30112	0.365919	0.38473	0.426482
16	-0.217012	-0.0781183	-0.070247	-0.0590773	-0.0474837	-0.0361357	-0.0237122	-0.015215	-0.75151	0.00617	0.0167773	0.0468225	0.032165	0.0464072	0.0503973	0.053266	0.046289	0.049592
17	-0.35941	-0.40446	-0.346509	-0.457618	-0.367452	-0.281327	-0.199082	-0.116088	-0.0363185	0.068276	0.116416	0.191031	0.281464	0.29295	0.190211	0.465393	0.49048	0.47928
18	0.01717	-0.335263	0.48188	0.404811	0.322389	0.24938	0.178183	0.106304	0.0339282	-0.0321519	-0.101237	-0.188461	-0.27123	-0.331954	-0.345312	-0.19471	-0.460768	-0.541545
19	-0.308228	-0.10248	-0.0321278	-0.0782134	-0.0623809	-0.0481177	-0.0340094	-0.0201968	-6.957118	6.51513	0.0195689	0.0212726	0.0447188	0.0564227	0.0667388	0.0761171	0.051823	0.0926252
20	-0.66494	-2.12322	-2.94466	-4.67539	-1.94381	-1.25225	-1.07734	-0.695849	-0.207675	0.20617	0.819181	1.01328	1.41495	1.81928	1.78728	2.71421	3.95468	2.18734
21	-2.44937	-2.63716	-2.35264	-2.00347	-1.61512	-1.22529	-0.876581	-0.518787	-0.168426	0.117465	0.502639	0.819143	1.04663	1.48929	1.71421	2.95468	1.18734	2.57113
22	0.44937	0.0477553	0.0477553	0.061456	0.0390521	0.0212316	0.0137313	0.113181	3.1029532	-3.1029532	-0.010213	-0.010213	-0.0200481	-0.0200481	-0.0200481	-0.0200481	-0.0200481	-0.0200481
23	1.72716	0.46642	1.57511	3.07383	4.07682	3.8964	1.37395	0.791558	0.457438	-0.150618	-0.746187	-1.18729	-1.79706	-2.21851	-2.62223	-2.93209	-3.34693	-3.52523
24	-1.81848	-3.29365	-3.54811	-3.77712	-3.1538	-1.81211	-1.29588	-0.746811	-0.289544	0.137114	0.447468	0.742411	1.201651	2.107281	2.51851	2.89044	3.24817	3.52523
25	0.741901	0.40674	0.727129	0.610731	0.490493	0.258417	0.258417	0.157672	0.0511901	0.0506483	-0.15277	-0.150668	-0.34911	-0.49048	0.286464	0.286157	0.286157	0.286157
26	-0.32355	-0.310163	-0.316041	-0.165495	-0.211393	-0.163293	-0.115535	-0.0456064	-0.0222236	0.0041165	0.0646935	0.110879	0.151759	0.191478	0.226464	0.258157	0.286157	0.286157
27	0.438593	0.499182	0.449181	0.377503	0.303423	0.232188	0.164303	0.091581	0.0164629	-0.014682	-0.094484	-0.151639	-0.215787	-0.271213	-0.325013	-0.37521	-0.420947	-0.461513
28	1.07683	1.17521	1.14578	0.969256	0.776709	0.576223	0.415698	0.24881	0.0807326	-0.0662777	-0.141094	-0.206515	-0.256747	-0.302103	-0.34221	-0.37521	-0.420947	-0.461513
29	-0.016968	-0.0160927	-0.0124867	-0.0106577	-0.564812	-6.25351	-4.82952	-2.70181	0.932714	1.678274	2.465881	3.45078	4.61021	7.66612	9.09781	0.0162776	0.0116043	0.0116183
30	0.407163	0.456579	0.251074	0.196341	-0.20997	-0.205199	-0.216109	-0.151809	-0.0416174	0.0413555	0.142003	0.20737	0.26326	0.32811	0.432584	0.48369	0.50646	0.587864
31	1.11898	1.21806	1.09653	0.92177	0.740041	0.564573	0.400923	0.238113	0.0773085	-0.076742	-0.200419	-0.344711	-0.508552	-0.684362	-0.785811	-0.84919	-1.003	-1.07064
32	2.31419	2.60793	2.14718	1.82055	1.44428	1.11776	0.79257	0.46994	0.15251	-0.15163	-0.45394	-0.760317	-1.06064	-1.311	-1.5526	-1.7727	-1.96226	-2.15344

Tabla 11a

## VALORES DE COMPARACION - Sexo Femenino

E.F.	0	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80	
1	-0.149309	-0.158911	-0.168657	-0.178519	-0.096493	-0.0761011	-0.0539779	-0.0319992	-0.0104208	0.01051814	0.0311803	0.0515593	0.071545	0.0904595	0.107776	0.123261	0.137223	0.149441
2	-0.147663	-0.100416	-0.112097	-0.077054	-0.426773	-0.481027	-0.361169	-0.202884	-0.0479163	0.04444	0.156771	0.325711	0.532228	0.371458	0.478713	0.777731	0.767731	0.964611
3	0.646733	0.727776	0.437941	0.553156	0.430173	0.364609	0.265817	0.137313	0.0476932	-0.0478716	-0.141763	-0.234696	-0.325857	-0.411916	-0.489053	-0.560202	-0.626992	-0.690504
4	0.346778	0.61824	0.565090	0.476288	0.394395	0.326118	0.210034	0.119713	0.0405786	-0.0405002	-0.131131	-0.250206	-0.31839	-0.39391	-0.477113	-0.478768	-0.533949	-0.581375
5	-0.61722	-0.725562	-0.635566	-0.556231	-0.431203	-0.367446	-0.266452	-0.1610104	-0.0476127	-0.0749394	-0.162136	-0.232735	-0.326464	-0.429539	-0.540112	-0.617188	-0.626538	0.682180
6	-1.11501	-1.78603	-1.07897	-1.91931	-0.736249	-0.596413	-0.431711	-0.239006	-0.0776464	-0.0726536	-0.232317	-0.364663	-0.526463	-0.670550	-0.819018	-0.915018	-1.02494	-1.143719
7	-1.61909	-1.48003	-1.56357	-1.31391	-1.06702	-0.827703	-0.584216	-0.349593	-0.11887	-0.113771	0.238966	0.357762	0.776187	0.976697	1.16282	1.33177	1.48327	1.613719
8	0.457470	0.483943	0.439387	0.371302	0.301004	0.231737	0.164386	0.0974501	0.0317973	-0.0261106	-0.094635	-0.158927	-0.217682	-0.275423	-0.327022	-0.371079	-0.417496	-0.459013
9	0.214463	0.264856	0.260585	0.203001	0.146411	0.128748	0.0489971	0.0328293	0.0178672	-0.0176057	-0.0318453	-0.082819	-0.119156	-0.150623	-0.178299	-0.204918	-0.232334	-0.261883
10	0.736137	0.77709	0.70641	0.596053	0.483283	0.372349	0.239558	0.15668	0.0909961	-0.0206007	-0.13221	-0.251366	-0.381611	-0.521258	-0.670485	-0.810485	-0.950485	-1.090485
11	-0.204362	-0.217504	-0.197731	-0.166832	-0.135283	-0.104181	-0.0738806	-0.0437979	-0.014273	0.0143563	0.044264	0.0705291	0.0973243	0.125784	0.164948	0.166409	0.178719	0.204501
12	-0.046664	-0.0262501	-0.0238423	-0.0201347	-0.016327	-0.012371	-0.91643	-5.20287	-1.72225	1.736118	3.142728	8.31278	0.016181	0.0169393	0.0177372	0.0202349	0.0246675	0.0266680
13	-2.37716	-1.39078	-2.39054	-2.46742	-1.38939	-1.32023	-1.09382	-0.84621	-0.207383	0.211023	0.425064	1.03645	1.81573	2.159	2.47056	2.75518	3.00001	3.20001
14	-0.14513	-0.12678	-0.148135	-0.104913	-0.0825046	-0.0635029	-0.046607	-0.0376248	-0.0379738	-0.0473174	-0.0467966	-0.044523	-0.043184	0.0778403	0.0761224	0.105906	0.118118	0.126402
15	0.164653	0.071949	0.0181873	0.0137682	0.0111645	0.990778	3.099138	3.61455	1.17798	-1.18738	1.51813	5.82038	0.561415	-0.0121537	0.0212168	-0.0120882	-0.0155501	-0.0164769
16	-0.232607	-0.268533	-0.264108	-0.205376	-0.167028	-0.128598	-0.0912138	-0.0560756	-0.0176216	0.0177441	0.0526049	0.087076	0.130299	0.152621	0.181448	0.211884	0.231884	0.251884
17	0.0472218	0.0482859	0.0842873	0.0712045	0.037739	0.044456	0.0312323	0.018493	0.091728	-1.14034	-0.061853	-0.0301019	-0.0417965	-0.052832	-0.0622735	-0.0718771	-0.0801816	-0.0872813
18	0.0876039	0.0925399	0.0843786	0.0710266	0.0379947	0.0433449	0.0318123	0.0164663	0.0785113	-1.28495	-0.0181299	-0.0300167	-0.04169	-0.0527	-0.0625192	-0.0718972	-0.0799611	-0.0870656
19	1.20188	1.37917	1.14223	0.981171	0.755821	0.412587	0.496303	0.273762	0.089917	-0.0461112	-0.250587	-0.414152	-0.571914	-0.728025	-0.88454	-0.99039	-1.10459	-1.20271
20	-1.11464	-3.3107	-1.00957	-2.51912	-2.05119	-1.35947	-1.12406	-0.616462	-0.217256	0.218769	0.648557	1.61255	1.49523	1.82419	2.21704	2.56341	2.65665	3.11278
21	-1.21216	-1.15481	-1.19601	-1.00917	-0.81527	-0.65645	-0.44653	-0.264538	-0.043373	0.0872358	0.257308	0.42643	0.592317	0.74871	0.869007	1.0167	1.13612	1.21484
22	-1.31276	-1.23822	-1.27103	-1.07349	-0.86969	-0.6969	-0.47942	-0.281538	-0.0915641	0.6294859	0.272309	0.453397	0.629511	0.79376	0.946783	1.08242	1.2074	1.31484
23	2.44896	3.44772	3.24673	3.01525	1.63439	0.892569	0.529133	0.172435	-0.171811	-0.316761	-0.654078	-1.18503	-1.49549	-1.77955	-2.03459	-2.26909	-2.47063	-2.67063
24	-0.401244	-0.427228	-0.388488	-0.327564	-0.216602	-0.204501	-0.082991	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829	-0.082829
25	0.840164	0.849194	0.813162	0.685078	0.556171	0.428223	0.303733	0.180094	0.584786	-0.0591464	-0.175117	-0.289356	-0.402285	-0.508906	-0.606208	-0.692336	-0.771255	-0.840739
26	-0.340571	-0.626421	-0.363645	-0.465087	-0.289421	-0.219518	-0.126214	-0.011045	0.0434467	0.122716	0.201127	0.280228	0.356505	0.432273	0.495023	0.549037	0.594937	0.632937
27	-0.043572	-0.309536	-0.018128	-0.023711	-0.019232	-0.0148101	-0.0105047	-0.22716	-0.02948	0.0455783	0.038228	0.0100522	0.0131336	0.0176005	0.0208166	0.0239452	0.026705	0.029077
28	1.10747	1.37473	1.04778	0.91008	0.78066	0.56237	0.399064	0.265521	0.077088	-0.077028	-0.206127	-0.369655	-0.526889	-0.669566	-0.791287	-0.903978	-1.0148	-1.10931
29	0.316418	0.338777	0.306328	0.258673	0.205973	0.161176	0.118463	0.0479358	0.0251181	-0.0222594	-0.0462013	-0.070787	-0.101716	-0.131781	-0.162758	-0.194513	-0.226787	-0.258432
30	-1.17576	-0.297552	-0.370489	-0.228233	-0.150272	-0.142496	-0.101071	-0.059917	-0.0192529	-0.0198419	-0.0282897	0.0764661	0.137166	0.185267	0.201057	0.210289	0.216543	0.219765
31	-0.441138	-0.461028	-0.400356	-0.320773	-0.211118	-0.165267	-0.119753	-0.13312	-0.0423818	0.047276	0.165504	0.211466	0.297623	0.376254	0.444494	0.511883	0.576993	0.621563
32	0.422236	0.465462	0.460581	0.364753	0.279256	0.212246	0.132671	0.0905064	0.028486	-0.037227	-0.0480484	-0.143763	-0.202157	-0.255798	-0.307102	-0.360003	-0.388119	-0.422593

Tabla 11b



## ANALISIS DE REGRESION

Sexo Masculino

independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.044356	0.029305	-1.5136	0.1307
HIVIVIANE.VECCOMMAS	0.790516	0.033952	23.2831	0.0000

R-sq. (ADJ.) = 0.4848 SE= 0.703319 MAE= 0.492598 DurbinWat= 0.529  
 Previously: 0.0000 0.000000 0.000000 0.0000  
 576 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

## Cuadro 10a

## ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	268.156	1	268.156	542.104	.0000
Error	283.934	574	0.494658		
Total (Corr.)	552.090	575			

R-squared = 0.485711

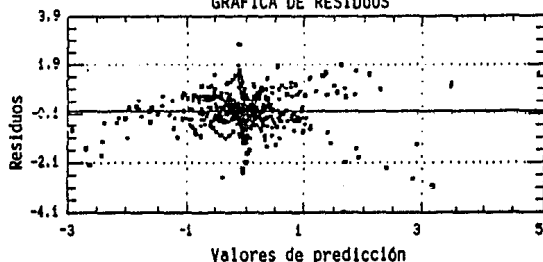
Std. error of est. = 0.703319

R-squared (Adj. for d.f.) = 0.484815

Durbin-Watson statistic = 0.528965

Gráfica 11a

GRAFICA DE RESIDUOS



Gráfica 12a

Gráfica: Probabilidad Normal

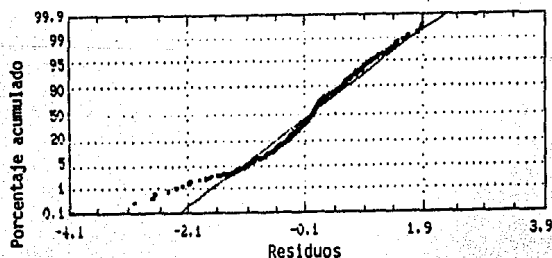


Tabla 12a  
Observaciones influenciadas:

Obs. Number	Std. Residual	Leverage	Mahalanobis Dist.	DFITS
217	0.24722	0.01371	6.97811	0.02914
218	-0.41165	0.01592	8.28587	-0.05235
219	0.06784	0.01323	6.69824	0.00786
233	1.06177	0.01131	5.56700	0.11355
234	0.81475	0.01306	6.59328	0.09371
325	3.94604	0.00176	0.01277	0.16561
343	-2.56702	0.02284	12.4195	-0.39553
344	-3.16349	0.02674	14.7735	-0.52438
345	-0.69973	0.02200	11.9167	-0.10496
346	-0.46186	0.01604	8.36149	-0.05898
347	-0.62502	0.01098	5.37672	-0.06587
355	1.85110	0.00639	2.69053	0.14839
356	2.25065	0.00914	4.29826	0.21622
357	2.71879	0.01210	6.03430	0.30094
358	2.38774	0.01525	7.88915	0.29711
359	2.15381	0.01863	9.90086	0.29679
360	1.24679	0.02172	11.7452	0.18577
361	-1.56714	0.01565	8.12754	-0.19760
362	-0.67410	0.01822	9.65444	-0.09183
363	0.18560	0.01510	7.80092	0.02298
364	0.00606	0.01117	5.48548	0.00064
374	1.93171	0.00661	2.82344	0.15762
375	2.07545	0.00856	3.96017	0.19290
376	1.47863	0.01064	5.17215	0.15331
377	1.27712	0.01287	6.48361	0.14581
378	0.95683	0.01490	7.68332	0.11767
397	-1.94320	0.03418	19.3182	-0.36558
398	-4.43900	0.04019	23.0389	-0.90635
399	-3.99860	0.03290	18.5263	-0.73747
400	-3.33970	0.02372	12.9479	-0.52057
401	-2.70164	0.01593	8.29515	-0.34376
410	0.56176	0.01322	6.68884	0.06501
411	-0.41435	0.01779	9.39869	-0.05577
412	-1.75804	0.02265	12.3030	-0.26762
413	-2.18322	0.02788	15.4639	-0.36973
414	-1.08536	0.03264	18.3720	-0.19938
433	2.60196	0.00301	0.73483	0.14297
559	-2.94548	0.01311	6.62559	-0.33946
560	-2.12398	0.01821	7.88924	-0.28399
561	-5.98118	0.01266	6.35942	-0.67718
575	0.17740	0.01092	5.33675	0.01864
576	-0.30272	0.01259	6.31950	-0.03418

Number of flagged observations (high leverage or DFITS) = 42



Cuadro 9b  
ANALISIS DE REGRESION  
Sexo Femenino

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.024448	0.023237	-1.0521	0.2932
MIVIVIANE.VECOMFEM	0.802562	0.030537	26.2815	0.0000
R-SQ. (ADJ.) = 0.5454 SE = 0.357676 MAE = 0.380127 DurWat = 0.626				
Previously: 0.0000 0.000000 0.000000 0.000				
576 observations fitted, forecast(s) computed for 0 missing val. of dep. var.				

Cuadro 10b  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	214.815	1	214.815	690.717	.0000
Error	178.516	574	0.311003		
Total (Corr.)	393.330	575			

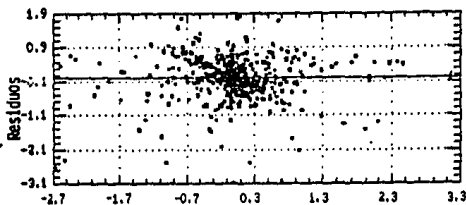
R-squared = 0.546144

Std. error of est. = 0.557676

R-squared (Adj. for d.f.) = 0.545353

Durbin-Watson statistic = 0.626303

Gráfica 11b  
GRAFICA DE RESIDUOS



Gráfica 12b  
Gráfica: Probabilidad Normal

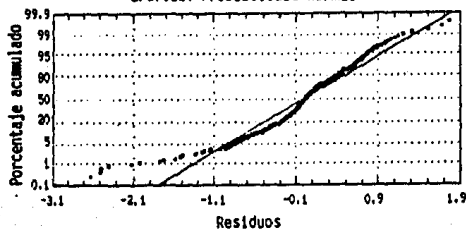


Tabla 12b  
Observaciones Influenciales

Obs. Number	Std. Residual	Leverage	Mahalanobis Dist.	DFITS
17	3.18690	0.00438	1.52435	0.21140
20	2.72351	0.00473	1.72856	0.18772
73	2.17351	0.00311	0.79077	0.12134
110	-2.76251	0.01055	5.12193	-0.28525
217	1.27296	0.02859	15.8929	0.21837
218	0.69541	0.03216	18.0732	0.12676
219	1.03264	0.02687	14.8486	0.17158
220	0.84022	0.01962	10.4861	0.11885
221	0.33148	0.01348	6.84495	0.03875
230	0.58348	0.01168	5.78766	0.06344
231	0.81399	0.01575	8.18417	0.10295
232	0.11082	0.02012	10.7873	0.01588
233	-0.00924	0.02459	13.4724	-0.00147
234	0.82391	0.02882	16.0363	0.14194
307	3.21820	0.00176	0.01479	0.13520
343	-4.43272	0.03065	17.1489	-0.78817
344	-4.88880	0.03449	19.5070	-0.92401
345	-1.47431	0.02879	16.0199	-0.25386
346	-0.90072	0.02099	11.3070	-0.13188
347	-1.14066	0.01438	7.37745	-0.13779
356	0.76459	0.01244	6.23437	0.08583
357	1.15855	0.01682	8.81894	0.15151
358	0.85467	0.02152	11.6283	0.12679
359	0.81167	0.02634	14.5282	0.13349
360	0.72362	0.03089	17.2992	0.12920
397	-0.11898	0.02010	10.7731	-0.01704
398	-2.39577	0.02253	12.2304	-0.36370
399	-2.72047	0.01892	10.0738	-0.37783
400	-2.42142	0.01398	7.14238	-0.28837
409	1.75310	0.00589	2.40470	0.13498
410	1.79296	0.00839	3.85982	0.16495
411	1.16937	0.01113	5.46229	0.12406
412	0.37760	0.01408	7.19957	0.04515
413	-0.16309	0.01710	8.98737	-0.02151
414	-0.45986	0.01996	10.6901	-0.06562
433	3.02560	0.00388	1.23787	0.18884
451	2.37984	0.00276	0.58770	0.12509

Number of flagged observations (high leverage or DFITS) = 37





Los residuos obtenidos en las tablas 9a y 9b fueron expuestos en los diagramas de tallo-y-hoja (figuras 3a y 3b) y en los diagramas de caja (figuras 4a y 4b) los cuales confirman los picos acentuados de las regiones en las que hay mayor y menor esperanza de vida.

Para diagnosticar la no-aditividad del modelo, se calcularon los valores de comparación como se describieron en el segundo capítulo (tablas 11a y 11b), y se graficaron los residuos contra estos valores. (gráficas 10a y 10b).

Detectando aquellos casos aberrantes o influenciados que pudieran sesgar la estimación de los parámetros, se calculó la pendiente B de los  $R_{ij}$  sobre  $C_{ij}$ , por mínimos cuadrados. Para ello, se llevó a cabo el análisis de regresión por mínimos cuadrados (cuadros 9a y 9b), el análisis de varianza (cuadros 10a y 10b), la gráfica sobre el supuesto de la normalidad de los residuos (gráficas 11a y 11b) y la gráfica de los residuos observados versus los residuos estimados (gráficas 12a y 12b).

Al observar las gráficas 12, se detecta la presencia de residuos llamados "casos aberrantes" que son observaciones cuyos valores exceden los valores de las demás observaciones en una gran medida, tal vez tres o cuatro desviaciones estándar más allá del valor medio de las observaciones. Pero, ¿Qué tanto afectan al valor de la pendiente? ¿Deberían descartarse estas observaciones corriendo una regresión con las restantes observaciones?

Primero para detectar cuáles son esas observaciones, se calcularon los

residuos estandarizados, los puntos de nivel (leverage points, en inglés), la distancia de Mahalanobis y la D-ajustado (D-fits, en inglés). Se observaron 42 casos influenciales o aberrantes, en el caso del sexo masculino (tabla 12a) y 37 en el caso del sexo femenino (tabla 12b). Para detectar, los casos influenciales fueron calculados utilizando la fórmula siguiente (Weisberg, 1980):

$$H_i = 1/n + (X_i - \bar{X})^2 / \text{Sum}(X_j - \bar{X})$$

$$\text{En este caso, } Y_i = B_0 + B_1X_i + e_i$$

con  $p = 1$  (variable independiente) y  $p' = 2$  (parámetros  $B_0$  y  $B_1$ )

por lo que  $\text{Sum } H_i = 2$  y  $H_i = p'/n = 2/576 = 0.00347$

Entonces si  $H_i > 2H_i = 0.00694$  entonces se sospechó de un caso inflencial.

Después se hizo un diagrama de tallo-y-hoja y el diagrama de caja correspondiente a los residuos estandarizados derivados de la regresión (figuras 5a, 5b y 6a, 6b; respectivamente). Al igual que la gráfica de los "leverage" contra los valores comparativos calculados anteriormente, Como se puede observar, estas dos gráficas (13a y 13b) describen una parábola.

En seguida, se corrió de nuevo una regresión de los residuos (cuadros 11a y 11b; 12a y 12b) pero esta vez extrayendo las 42 y 37 observaciones. Por el hecho de que la técnica de mínimos<sup>2</sup> es poco resistente a los casos aberrantes. Se observó que la pendiente cambia en un 12% por abajo de la pendiente total, la desviación estándar en un 10% aproximadamente y la R cuadrada en un 40%.

Cuadro 11b  
ANALISIS DE REGRESION  
Sin datos influenciados - Sexo Masculino

independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.027668	0.025683	-1.0773	0.2818
-economas	0.740293	0.046493	15.9227	0.0000
r-sq. (ADJ.) = 0.3215 SE= 0.593499 MAE= 0.433692 DurWat= 0.580				
Previously:	0.0000	0.000000	0.000000	0.000
334 observations fitted, forecast(s) computed for 0 missing val. of dep. var.				

Cuadro 12a  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	89.3040	1	89.3040	253.531	.0000
Error	187.392	532	0.352241		
Total (corr.)	276.696	533			

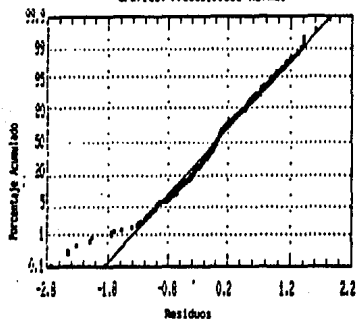
r-squared = 0.322751

r-squared (Adj. for d.f.) = 0.321478

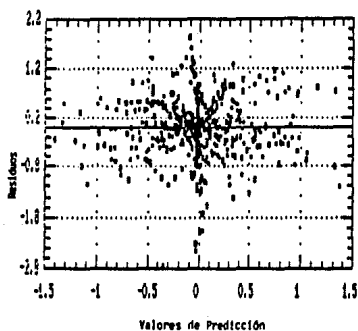
Std. error of est. = 0.593499

Durbin-Watson statistic = 0.580191

Gráfica 14a  
Gráfica: Probabilidad Normal



Gráfica 15a  
GRAFICA DE RESIDUOS



Cuadro 11b  
ANALISIS DE REGRESION  
Sin datos influenciados - Sexo Femenino

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.033028	0.021671	-1.5240	0.1281
veccomfen	0.732944	0.04242	17.2781	0.0000

R-SQ. (ADJ.) = 0.3561 SE= 0.503034 MAE= 0.345513 DurbWat= 0.671

539 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 12b  
ANALISIS DE VARIANZA

source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	75.5420	1	75.5420	298.534	.0000
Error	135.884	537	0.253043		
Total (Corr.)	211.426	538			

R-squared = 0.357298

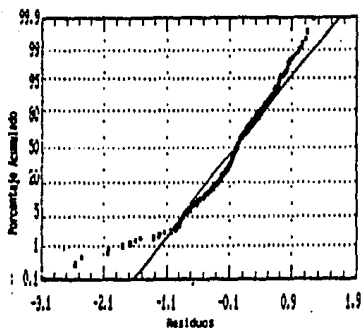
Std. error of est. = 0.503034

R-squared (Adj. for d.f.) = 0.356101

Durbin-Watson statistic = 0.670963

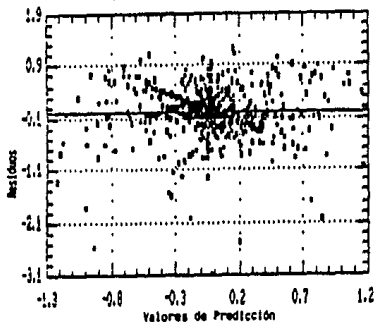
Gráfica 14b

Gráfica: Probabilidad Normal



Gráfica 15b

GRAFICA DE RESIDUOS



: (Esperanza de vida)<sup>0.25</sup>

Valor típico: 2.4156

Sexo Masculino

S.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	2.61177	2.65107	2.61943	2.76548	2.70647	2.64939	2.58786	2.52156	2.45576	2.38116	2.40691	2.31579	2.1187	2.01721	1.9032	1.7766	1.64756	1.52188
2	2.61429	2.65167	2.61329	2.75826	2.69906	2.64204	2.58076	2.51411	2.44831	2.37351	2.40077	2.29977	2.10277	2.00177	1.88782	1.7612	1.63156	1.50588
3	2.61681	2.65419	2.61181	2.75674	2.69754	2.64052	2.57924	2.51259	2.44679	2.37200	2.40000	2.29900	2.10200	2.00100	1.88705	1.76045	1.63079	1.50511
4	2.61933	2.65671	2.61433	2.75867	2.69947	2.64245	2.58117	2.51452	2.44872	2.37393	2.40193	2.30093	2.10393	2.00293	1.88898	1.76238	1.63272	1.50704
5	2.62185	2.65923	2.61685	2.76062	2.70142	2.64440	2.58312	2.51647	2.45067	2.37588	2.40388	2.30288	2.10588	2.00488	1.89093	1.76433	1.63467	1.50899
6	2.62437	2.66175	2.61937	2.76248	2.70328	2.64626	2.58498	2.51833	2.45253	2.37774	2.40574	2.30474	2.10774	2.00674	1.89279	1.76619	1.63653	1.51085
7	2.62689	2.66427	2.62189	2.76429	2.70509	2.64807	2.58680	2.52015	2.45435	2.37956	2.40756	2.30656	2.10956	2.00856	1.89461	1.76801	1.63835	1.51267
8	2.62941	2.66679	2.62441	2.76610	2.70690	2.64988	2.58861	2.52196	2.45616	2.38137	2.40937	2.30837	2.11137	2.01037	1.89642	1.76982	1.64016	1.51499
9	2.63193	2.66931	2.62693	2.76791	2.70871	2.65169	2.59041	2.52376	2.45796	2.38317	2.41117	2.31017	2.11317	2.01217	1.89822	1.77162	1.64196	1.51681
10	2.63445	2.67183	2.62945	2.76972	2.71052	2.65350	2.59222	2.52557	2.45977	2.38498	2.41298	2.31198	2.11498	2.01398	1.90003	1.77343	1.64377	1.51866
11	2.63697	2.67435	2.63197	2.77153	2.71233	2.65531	2.59404	2.52739	2.46159	2.38680	2.41480	2.31380	2.11680	2.01580	1.90185	1.77483	1.64475	1.51955
12	2.63949	2.67687	2.63449	2.77334	2.71414	2.65712	2.59577	2.52912	2.46332	2.38853	2.41653	2.31553	2.11853	2.01753	1.90359	1.77581	1.64573	1.52044
13	2.64201	2.67939	2.63701	2.77515	2.71595	2.65893	2.59750	2.53085	2.46505	2.39026	2.41826	2.31726	2.12026	2.01926	1.90464	1.77679	1.64671	1.52133
14	2.64453	2.68191	2.63953	2.77696	2.71776	2.66074	2.59933	2.53268	2.46688	2.39209	2.42009	2.31909	2.12209	2.02109	1.90562	1.77777	1.64769	1.52222
15	2.64705	2.68443	2.64205	2.77877	2.71957	2.66255	2.59986	2.53450	2.46869	2.39390	2.42190	2.32090	2.12509	2.02409	1.90660	1.77875	1.64867	1.52311
16	2.64957	2.68695	2.64457	2.78058	2.72138	2.66436	2.60139	2.53631	2.47049	2.39570	2.42390	2.32290	2.12709	2.02609	1.90758	1.77973	1.64965	1.52400
17	2.65209	2.68947	2.64709	2.78239	2.72319	2.66617	2.60292	2.53812	2.47220	2.39741	2.42591	2.32491	2.13029	2.02929	1.90856	1.78071	1.65063	1.52489
18	2.65461	2.69199	2.64961	2.78420	2.72500	2.66798	2.60445	2.54033	2.47439	2.39960	2.42792	2.32692	2.13368	2.03268	1.90954	1.78169	1.65161	1.52578
19	2.65713	2.69451	2.65213	2.78601	2.72681	2.66979	2.60598	2.54214	2.47620	2.40141	2.43093	2.32893	2.13647	2.03547	1.91052	1.78267	1.65259	1.52667
20	2.65965	2.69703	2.65465	2.78782	2.72862	2.67160	2.60761	2.54395	2.47807	2.40328	2.43295	2.33095	2.13936	2.03836	1.91150	1.78365	1.65357	1.52756
21	2.66217	2.69955	2.65717	2.78963	2.73043	2.67341	2.60922	2.54576	2.47998	2.40519	2.43496	2.33296	2.14174	2.04074	1.91248	1.78463	1.65455	1.52845
22	2.66469	2.70207	2.65969	2.79144	2.73224	2.67522	2.61083	2.54757	2.48179	2.40700	2.43697	2.33497	2.14458	2.04358	1.91346	1.78561	1.65553	1.52934
23	2.66721	2.70459	2.66221	2.79325	2.73405	2.67703	2.61244	2.54938	2.48360	2.40881	2.43894	2.33694	2.14740	2.04650	1.91444	1.78659	1.65651	1.53023
24	2.66973	2.70711	2.66473	2.79506	2.73586	2.67884	2.61365	2.55119	2.48541	2.41062	2.44091	2.33891	2.15022	2.04950	1.91542	1.78757	1.65749	1.53112
25	2.67225	2.70963	2.66725	2.79687	2.73767	2.68065	2.61486	2.55296	2.48722	2.41243	2.44292	2.34092	2.15304	2.05104	1.91640	1.78855	1.65847	1.53201
26	2.67477	2.71215	2.66977	2.79868	2.73948	2.68246	2.61607	2.55477	2.48903	2.41424	2.44491	2.34291	2.15586	2.05306	1.91738	1.78953	1.65945	1.53290
27	2.67729	2.71467	2.67229	2.80049	2.74129	2.68427	2.61728	2.55658	2.49084	2.41605	2.44692	2.34492	2.15868	2.05506	1.91836	1.79051	1.66043	1.53379
28	2.67981	2.71719	2.67481	2.80230	2.74310	2.68608	2.61849	2.55839	2.49265	2.41786	2.44881	2.34681	2.16150	2.05706	1.91934	1.79149	1.66141	1.53468
29	2.68233	2.71971	2.67733	2.80411	2.74491	2.68789	2.61970	2.56020	2.49446	2.41907	2.45068	2.34868	2.16432	2.05906	1.92032	1.79247	1.66239	1.53557
30	2.68485	2.72223	2.67985	2.80592	2.74672	2.68970	2.62091	2.56201	2.49627	2.42088	2.45249	2.35049	2.16714	2.06106	1.92130	1.79345	1.66337	1.53646
31	2.68737	2.72475	2.68237	2.80773	2.74853	2.69151	2.62212	2.56382	2.49808	2.42269	2.45430	2.35230	2.16996	2.06306	1.92228	1.79443	1.66435	1.53735
32	2.68989	2.72727	2.68489	2.80954	2.75034	2.69332	2.62333	2.56563	2.49989	2.42450	2.45611	2.35411	2.17278	2.06506	1.92326	1.79541	1.66533	1.53824
33	2.69241	2.72979	2.68741	2.81135	2.75215	2.69513	2.62454	2.56744	2.50170	2.42631	2.45772	2.35572	2.17560	2.06706	1.92424	1.79639	1.66631	1.53913
34	2.69493	2.73231	2.68993	2.81316	2.75396	2.69694	2.62575	2.56925	2.50351	2.42812	2.45933	2.35733	2.17842	2.06906	1.92522	1.79737	1.66729	1.54002
35	2.69745	2.73483	2.69245	2.81497	2.75577	2.69875	2.62696	2.57106	2.50532	2.43033	2.46094	2.35934	2.18124	2.07106	1.92620	1.79835	1.66827	1.54091
36	2.69997	2.73735	2.69497	2.81678	2.75758	2.70056	2.62817	2.57287	2.50713	2.43214	2.46255	2.36135	2.18406	2.07306	1.92718	1.79933	1.66925	1.54180
37	2.70249	2.73987	2.69749	2.81859	2.75939	2.70237	2.62938	2.57468	2.50894	2.43395	2.46416	2.36316	2.18688	2.07506	1.92816	1.80031	1.67023	1.54269
38	2.70501	2.74239	2.69999	2.82040	2.76120	2.70418	2.63059	2.57649	2.51075	2.43576	2.46597	2.36497	2.18970	2.07706	1.92914	1.80129	1.67121	1.54358
39	2.70753	2.74491	2.70251	2.82221	2.76301	2.70599	2.63180	2.57830	2.51256	2.43757	2.46818	2.36618	2.19252	2.07906	1.93012	1.80227	1.67219	1.54447
40	2.71005	2.74743	2.70503	2.82402	2.76482	2.70780	2.63301	2.58011	2.51437	2.43938	2.47039	2.36839	2.19534	2.08106	1.93110	1.80325	1.67317	1.54536
41	2.71257	2.74995	2.70755	2.82583	2.76663	2.70961	2.63422	2.58192	2.51618	2.44119	2.47260	2.37060	2.19816	2.08306	1.93208	1.80423	1.67415	1.54625
42	2.71509	2.75247	2.71007	2.82764	2.76844	2.71142	2.63543	2.58373	2.51799	2.44290	2.47481	2.37261	2.20098	2.08506	1.93306	1.80521	1.67513	1.54714
43	2.71761	2.75499	2.71259	2.82945	2.77025	2.71323	2.63664	2.58554	2.51980	2.44471	2.47702	2.37442	2.20380	2.08706	1.93404	1.80619	1.67611	1.54803
44	2.72013	2.75751	2.71511	2.83126	2.77206	2.71504	2.63785	2.58735	2.52161	2.44652	2.47923	2.37623	2.20662	2.08906	1.93502	1.80717	1.67709	1.54892
45	2.72265	2.76003	2.71763	2.83307	2.77387	2.71685	2.63906	2.58916	2.52342	2.44833	2.48144	2.37804	2.20944	2.09106	1.93600	1.80815	1.67807	1.54981
46	2.72517	2.76255	2.72015	2.83488	2.77568	2.71866	2.64027	2.59097	2.52523	2.45014	2.48365	2.38025	2.21226	2.09306	1.93698	1.80913	1.67905	1.55070
47	2.72769	2.76507	2.72267	2.83669	2.77749	2.72047	2.64148	2.59278	2.52704	2.45195	2.48586	2.38206	2.21508	2.09506	1.93796	1.81011	1.68003	1.55159
48	2.73021	2.76759	2.72519	2.83850	2.77930	2.72228	2.64269	2.59459	2.52885	2.45376	2.48807	2.38387	2.21790	2.09706	1.93894	1.81109	1.68101	1.55248
49	2.73273	2.77011	2.72771	2.84031	2.78111	2.72409	2.64390	2.59640	2.53066	2.45557	2.49028	2.38568	2.22072	2.09906	1.93992	1.81207	1.68199	1.55337
50	2.73525	2.77263	2.73023	2.84212	2.78292	2.72590	2.64511	2.59821	2.53247	2.45738	2.49249	2.38749	2.22354	2.10106	1.94090	1.81305	1.68297	1.55426
51	2.73777	2.77515	2.73275	2.84393	2.78473	2.72771	2.64632	2.60002	2.53428	2.45919	2.49470	2.38930	2.22636	2.10306	1.94188	1.81403	1.68395	1.55515
52	2.74029	2.77767	2.73527	2.84574	2.78654	2.72952	2.64753	2.60183	2.53609	2.46090	2.49691	2.39111	2.22918	2.10506	1.94286	1.81501	1.68493	1.55604
53	2.74281	2.78019	2.73779	2.84755	2.78835	2.73133	2.64874	2.60364	2.53790	2.46271	2.49912	2.39292	2.23200	2.10706	1.94384	1.81599	1.68591	1.5569

(Esperanza de vida)<sup>0.25</sup>  
Sexo Femenino

Valor típico: 2.47347

S.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	2.8731	2.84989	2.86643	2.84396	2.86003	2.89281	2.64721	2.37857	2.31511	2.41836	2.39942	2.37102	2.34799	2.06282	1.93912	1.85031	1.81021	1.53288
2	2.67412	2.67571	2.65256	2.66067	2.74326	2.64551	2.52747	2.56009	2.48964	2.41274	2.35773	2.32679	2.31627	2.32551	2.32457	2.00944	1.68364	1.56717
3	2.96358	2.92331	2.87723	2.83214	2.78084	2.72203	2.65792	2.59418	2.52605	2.44419	2.37577	2.28414	2.18018	2.08440	1.95113	1.84602	1.73728	1.59303
4	2.69463	2.51551	2.42221	2.37712	2.31294	2.23194	2.15248	2.09772	2.02707	2.04941	2.07036	2.16232	2.18722	2.08143	1.96427	1.82816	1.68004	1.54704
5	2.67282	2.67433	2.66114	2.66444	2.75119	2.61219	2.63122	2.56447	2.49336	2.41476	2.33507	2.25027	2.15238	2.05649	1.93271	1.83001	1.73052	1.58374
6	2.81648	2.6134	2.56275	2.50575	2.43746	2.36806	2.28756	2.23511	2.16648	2.07267	2.01152	1.94007	1.88381	2.03123	2.10467	1.78259	1.64877	1.53348
7	2.61918	2.64941	2.61211	2.70055	2.62587	2.64001	2.60178	2.56464	2.4784	2.40229	2.3312	2.23005	2.15992	2.19992	2.05719	1.86191	1.70948	1.56422
8	2.84122	2.7601	2.6772	2.62544	2.57028	2.71211	2.6318	2.56683	2.51513	2.44251	2.44941	2.37913	2.38912	2.08940	1.97949	1.87043	1.71643	1.5824
9	2.67725	2.61679	2.57551	2.6271	2.77125	2.71119	2.6510	2.59594	2.51887	2.43554	2.35885	2.27272	2.18611	2.07122	1.97321	1.87088	1.72019	1.57980
10	2.91481	2.71729	2.66218	2.61804	2.77314	2.71351	2.65728	2.57446	2.54856	2.45088	2.37627	2.28621	2.19303	2.09126	1.98194	1.86374	1.73201	1.58498
11	2.65272	2.64677	2.63815	2.61894	2.73967	2.7018	2.64649	2.57599	2.50713	2.43253	2.35894	2.26529	2.17226	2.04391	1.95178	1.82323	1.68959	1.55634
12	2.67516	2.66631	2.65876	2.60421	2.73241	2.69638	2.63423	2.57419	2.50912	2.43256	2.36605	2.28415	2.19113	2.08449	1.97777	1.87771	1.71842	1.57914
13	2.66172	2.6359	2.61034	2.7388	2.70023	2.64139	2.59018	2.51581	2.44776	2.37364	2.29461	2.21658	2.12355	2.03122	1.91825	1.78233	1.64429	1.52161
14	2.67487	2.60125	2.66612	2.61798	2.7821	2.7025	2.64234	2.57819	2.50602	2.43287	2.35368	2.26781	2.17320	2.07827	1.96311	1.82958	1.68574	1.54574
15	2.64711	2.63913	2.63793	2.61386	2.74284	2.70211	2.64844	2.57856	2.50751	2.43187	2.35247	2.26739	2.18003	2.05256	1.96773	1.83610	1.70244	1.56882
16	2.66885	2.67632	2.61664	2.61117	2.75528	2.67123	2.63623	2.57328	2.50675	2.43188	2.3558	2.26974	2.17842	2.07645	1.96772	1.83610	1.70244	1.56882
17	2.61222	2.63933	2.61227	2.61372	2.75814	2.70071	2.63932	2.57395	2.50887	2.43459	2.35827	2.28141	2.18875	2.08544	1.98777	1.84441	1.70442	1.56742
18	2.69771	2.66419	2.67263	2.62246	2.74444	2.70729	2.64944	2.57974	2.51251	2.43677	2.35904	2.27225	2.18444	2.07664	1.97174	1.81845	1.67385	1.53388
19	2.60731	2.62621	2.67222	2.64025	2.73858	2.71858	2.66014	2.60344	2.53372	2.45957	2.38627	2.31532	2.26036	2.10194	1.97211	1.83164	1.71577	1.58101
20	2.74444	2.60513	2.73216	2.74654	2.68741	2.68864	2.56764	2.50871	2.44221	2.37474	2.30185	2.21849	2.12677	2.01612	1.92248	1.79483	1.64434	1.51801
21	2.66235	2.62565	2.63605	2.71818	2.72821	2.67153	2.61228	2.55147	2.48412	2.41692	2.34649	2.26952	2.18566	2.0866	1.98727	1.85122	1.71577	1.55454
22	2.66663	2.67815	2.66051	2.73721	2.74129	2.68123	2.62181	2.55971	2.49505	2.42674	2.35716	2.28632	2.19622	2.09652	1.9121	1.73112	1.58312	1.51088
23	2.6322	2.67605	2.70391	2.63244	2.73788	2.74609	2.68212	2.62056	2.55507	2.48627	2.41676	2.33618	2.24594	2.14708	2.04041	1.90047	1.76441	1.61418
24	2.65412	2.67828	2.68042	2.74218	2.74228	2.68415	2.62158	2.56402	2.49318	2.41674	2.33906	2.24767	2.15878	2.06111	1.95519	1.82558	1.72123	1.56419
25	2.66663	2.67815	2.66051	2.73721	2.74129	2.68123	2.62181	2.55971	2.49505	2.42674	2.35716	2.28632	2.19622	2.09652	1.9121	1.73112	1.58312	1.51088
26	2.6322	2.67605	2.70391	2.63244	2.73788	2.74609	2.68212	2.62056	2.55507	2.48627	2.41676	2.33618	2.24594	2.14708	2.04041	1.90047	1.76441	1.61418
27	2.65412	2.67828	2.68042	2.74218	2.74228	2.68415	2.62158	2.56402	2.49318	2.41674	2.33906	2.24767	2.15878	2.06111	1.95519	1.82558	1.72123	1.56419
28	2.66663	2.67815	2.66051	2.73721	2.74129	2.68123	2.62181	2.55971	2.49505	2.42674	2.35716	2.28632	2.19622	2.09652	1.9121	1.73112	1.58312	1.51088
29	2.6322	2.67605	2.70391	2.63244	2.73788	2.74609	2.68212	2.62056	2.55507	2.48627	2.41676	2.33618	2.24594	2.14708	2.04041	1.90047	1.76441	1.61418
30	2.65412	2.67828	2.68042	2.74218	2.74228	2.68415	2.62158	2.56402	2.49318	2.41674	2.33906	2.24767	2.15878	2.06111	1.95519	1.82558	1.72123	1.56419
31	2.66663	2.67815	2.66051	2.73721	2.74129	2.68123	2.62181	2.55971	2.49505	2.42674	2.35716	2.28632	2.19622	2.09652	1.9121	1.73112	1.58312	1.51088
32	2.6322	2.67605	2.70391	2.63244	2.73788	2.74609	2.68212	2.62056	2.55507	2.48627	2.41676	2.33618	2.24594	2.14708	2.04041	1.90047	1.76441	1.61418

Tabla 13b

MEDIANA PALIDA DE LOS DATOS TRANSFORMADOS

Sexo Masculino

E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-1.5144E-3	5.29250E-3	1.13919E-3	1.81610E-3	3.02772E-3	3.40041E-3	2.61768E-3	1.33266E-3	-1.50772E-3	-6.10411E-4	-6.13422E-4	-1.0916E-3	-0.0100379	-0.0151126	-0.021878	-0.0172991	-6.20311E-3	-5.42161E-3
2	6.26751E-3	0.0104961	1.19922E-3	1.97321E-3	6.14942E-4	1.34399E-3	2.51877E-3	1.01927E-4	-6.4012E-4	-5.2310E-3	-0.0103621	-0.015421	-0.016434	-0.0176316	-0.0137201	-6.41601E-3	6.92301E-3	-4.1336E-3
3	7.63115E-3	0.0141451	7.74693E-3	6.89321E-3	0.0112764	0.0104764	7.7937E-3	2.20712E-3	-1.2237E-3	-1.2441E-3	-0.0114149	-0.0114421	-0.0163099	7.1221E-3	1.34511E-3	-0.21142E-3	-0.18110E-3	-7.8911E-3
4	1.48127E-3	6.4246E-3	0.017061	1.31212E-3	2.2214E-3	3.79941E-3	1.41871E-3	1.81544E-3	1.4784E-3	-1.2142E-3	-0.012771E-3	-0.14121E-3	-0.0109966	-0.0137982	-0.0176491	-0.0176491	-0.0134474	-0.0169607
5	0.015153	0.0187444	0.0113256	0.0115587	0.0100771	0.0100611	6.10724E-3	4.0791E-3	-1.1394E-4	-6.7621E-3	-0.0164482	-0.0164562	-0.0113006	-0.0067262	-0.0290923	-0.0144864	-1.3372E-3	1.73361E-4
6	6.48194E-3	5.36633E-3	0.0107203	0.0104764	0.0100075	5.29931E-3	7.51272E-3	4.11289E-3	2.09417E-3	-2.1847E-3	-6.6054E-3	-0.0110496	-0.0204496	-0.0236758	-0.0304283	-0.0162026	-0.077674E-3	-0.0272953
7	-4.3744E-3	-5.2999E-3	1.24622E-3	6.40472E-3	1.11932E-3	2.12328E-3	3.20594E-3	2.30975E-3	1.96447E-3	2.11697E-3	-8.9917E-3	-0.01787E-3	-5.13700E-3	-0.0134976	-5.3660E-3	-0.025331E-3	-0.0277493	-1.9316E-3
8	-6.1231E-4	-0.1759E-3	1.90209E-3	5.4221E-4	2.84952E-3	2.73221E-3	1.96095E-3	2.30975E-3	1.99591E-3	4.9731E-3	-1.8184E-3	-5.7254E-3	-7.8771E-3	-0.0170956	-0.02324E-3	-0.02324E-3	-0.0179499	-0.0161373
9	9.48014E-3	0.0123536	5.19352E-3	6.29251E-3	1.78864E-3	1.04845E-3	6.41821E-4	-5.7111E-3	-0.0107067	-0.20619E-3	-0.0179495	-0.017236	-0.0137379	-0.0113381	-0.0114951	1.88441E-3	0.0140184	-1.43941E-3
10	5.05719E-3	-2.1617E-3	-6.6172E-3	-6.4959E-3	-8.1697E-3	-6.4931E-3	-4.6468E-3	-0.02731E-3	-2.1719E-3	1.3170E-3	2.2721E-3	6.12531E-3	3.38222E-3	1.7471E-3	3.1045E-3	9.11531E-3	0.017133	5.9364E-4
11	-0.0280709	-1.0734E-3	1.17721E-3	2.85122E-4	7.3475E-4	1.1779E-4	-4.4132E-4	-3.9409E-3	-1.8148E-3	-2.110E-3	6.73551E-3	2.1141E-3	3.39409E-3	2.02742E-3	3.41944E-3	-0.0103686	-0.018773	-6.4549E-3
12	4.18732E-3	-0.13081E-3	-0.010694	-0.0106937	-0.0112980	-0.011561	-6.2964E-3	-8.2897E-3	-6.2964E-3	6.20771E-3	3.0147E-3	0.01991	0.0212185	0.0218226	0.0202668	3.51451E-3	-1.8152E-3	5.0194E-3
13	-1.7764E-4	-1.0241E-4	-5.4937E-5	-1.7812E-3	-3.2301E-3	-4.8901E-3	-6.4184E-3	-6.728E-3	-3.26045E-3	-6.0164E-3	-6.8467E-3	6.5077E-3	3.59375E-3	7.6461E-3	0.012124	0.018487	0.0234476	0.0194264
14	-4.8502E-3	4.91548E-3	1.4313E-3	6.06061E-3	1.23907E-3	8.49364E-3	3.39588E-3	2.47545E-3	1.2094E-3	1.19351E-3	-1.6661E-4	-8.2160E-3	2.0878E-4	2.2848E-4	2.8478E-3	1.4879E-3	-5.80075E-3	-3.2285E-3
15	-0.0165564	-0.01744E-3	2.7993E-3	1.96422E-3	5.02952E-4	3.0203E-3	-8.7901E-3	-8.2451E-3	-5.7208E-3	-4.1843E-3	-6.5224E-3	-2.7441E-3	-1.7901E-3	-1.3094E-3	6.09921E-3	0.02391E-3	0.0264783	4.92959E-3
16	-4.1870E-3	-4.8974E-3	-1.9058E-3	-0.010482	-0.0112734	-0.0115281	-6.8932E-3	-8.1334E-3	-5.7208E-3	2.60581E-3	3.03541E-3	6.13721E-3	0.01694E-3	6.0949E-3	0.012123	0.0134671	0.0176049	6.28695E-3
17	7.04413E-3	-1.8772E-3	-7.4507E-3	-9.2679E-3	-0.011849	-0.0110486	-6.5628E-3	-6.4820E-3	-4.1919E-3	-1.8746E-3	1.8534E-3	5.7409E-3	3.80892E-3	3.2773E-3	0.0182289	0.027944	0.023237	0.0164264
18	0.013849	6.70639E-3	1.86645E-3	-1.4131E-3	5.1681E-3	6.64958E-3	-1.1221E-4	-0.0029E-3	1.1713E-3	2.13791E-3	-1.0564E-4	-6.344E-4	2.09944E-4	-9.2461E-3	-0.021724	-0.025944	-0.018246	-6.8031E-3
19	0.028191	0.024144	0.0301899	0.0194491	0.0152248	0.0200944	0.0182124	0.0115481	6.8923E-3	-6.87E-3	-0.0141733	-0.0227264	-0.0235472	-0.0150707	-0.0293627	-6.10139E-3	-0.110249	-0.016478
20	-0.048424	-0.0335E-3	-6.8015E-4	-1.9131E-3	-3.6420E-3	-5.1537E-3	-6.9748E-3	-6.2401E-3	-5.3415E-3	6.7181E-3	0.0131232	0.020367	0.0269273	0.0365262	0.0395078	0.01826E-3	0.0267953	0.0267953
21	-0.0149424	-0.0335E-3	-6.8015E-4	-1.9131E-3	-3.6420E-3	-5.1537E-3	-6.9748E-3	-6.2401E-3	-5.3415E-3	6.8292E-3	7.9125E-3	0.0144149	0.0210645	0.0276767	0.01826E-3	0.031204	0.029819	0.0292326
22	-0.041943	-0.03040E-3	-4.4077E-3	-6.2381E-3	-4.1228E-3	-9.3777E-3	-8.5964E-3	-6.4493E-3	-2.351E-3	5.1844E-3	-1.0564E-4	-6.344E-4	2.09944E-4	-9.2461E-3	-0.021724	-0.025944	-0.018246	-6.8031E-3
23	1.50741E-3	-0.0281E-3	-0.0161865	-1.0451E-3	-4.3815E-3	-1.1994E-3	6.49041E-3	7.82467E-3	0.018518	-0.0270E-3	0.016096	0.0137941	0.0171896	5.6413E-3	-4.1352E-3	-0.0199961	-0.0484392	-0.014445
24	-0.0130727	-0.01511E-3	-1.2174E-3	-4.4520E-3	-5.4183E-3	-5.3904E-3	-6.4489E-3	-6.9773E-3	-2.3240E-3	2.2979E-3	5.04754E-3	7.5646E-3	9.02444E-3	0.0101552	0.0101551	0.0164974	0.0170749	5.2376E-3
25	0.017713	2.35637E-3	-4.6181E-3	-4.7740E-3	-4.7740E-3	-1.9475E-3	8.9501E-3	-1.0192E-3	1.9314E-3	-2.8744E-3	-0.0741E-4	-6.0278E-3	-2.2895E-3	7.41955E-3	0.01826E-3	0.0260749	1.4184E-3	
26	0.146159	0.017417	0.0109957	8.4843E-3	6.41593E-3	9.7932E-3	0.021148	6.32474E-3	1.6194E-3	-1.7151E-3	-0.0110761	-0.0167658	-0.0219701	-0.0219498	-0.0267906	-0.02041E-3	-1.4891E-3	-1.2194E-3
27	-0.0197811	-0.012819	-0.0125476	-0.0106685	-6.8410E-3	-5.8483E-3	-1.4872E-3	2.297E-3	0.0107294	-0.04143E-3	0.0131005	0.0137482	0.0226595	0.0175267	0.0114998	2.91805E-3	-3.2582E-3	6.12144E-3
28	6.81937E-3	6.89344E-3	-1.7913E-3	-1.2812E-3	-6.1164E-3	-1.6201E-3	1.49301E-3	3.2644E-3	6.39502E-3	-1.4603E-3	-5.8000E-3	-1.9319E-3	-2.4480E-3	-5.3927E-3	4.13461E-3	0.0152326	0.0268748	4.9289E-3
29	-0.0285266	-5.2499E-3	-2.3604E-3	-1.41158E-3	-1.3273E-3	-1.9900E-3	-3.9960E-4	-2.1971E-3	-1.1341E-3	-1.7646E-3	1.14074E-3	7.3771E-3	0.0128974	0.0121556	0.0137909	0.0264565	0.0137749	6.02321E-3
30	6.82415E-3	-8.4072E-4	-1.9131E-3	-1.6190E-3	-2.0822E-3	-2.322E-3	-1.6949E-3	-2.8140E-3	-1.2753E-3	6.12391E-3	2.7847E-3	6.12391E-3	0.0134289	0.0109957E-3	-7.2153E-3	0.0156921	0.0162291	0.0121934
31	-3.8438E-4	-0.01602E-3	1.46011E-3	6.36411E-3	6.9608E-3	3.9791E-3	5.36429E-3	3.45644E-3	2.8808E-3	-2.4636E-3	1.05417E-3	-2.8152E-3	-9.7031E-3	-0.0237903	-0.0275077	-0.0407931	-0.0254261	-0.021437
32	-0.0210178	-0.0306132	0.0400631	-9.5479E-3	-7.9371E-3	-5.9003E-3	-1.7487E-3	-1.49361E-3	1.81071E-3	6.81351E-3	-2.6467E-3	1.44217E-3	-4.5943E-3	-4.5943E-3	-4.5943E-3	-4.5943E-3	-4.5943E-3	-4.5943E-3

Tabla 14a



MEDIANA PULIDA DE LOS DATOS TRANSFORMADOS

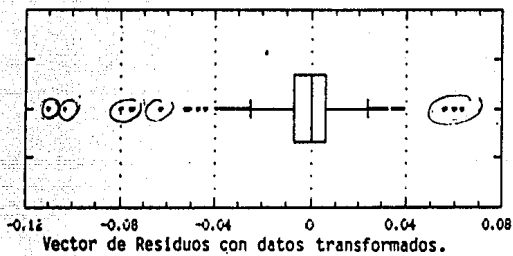
Sexo Femenino

E.F.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-7.16117	-3.65704	-3.78823	-3.71828	-3.60659	-3.31787	-3.77629	-3.11352	-3.55781	-2.10195	-1.33548	-1.79172	-5.36412	-0.01501	-0.03147	-0.04149	-0.03921	-0.01495
2	-0.61246	-0.61639	0.01345	0.01521	0.01749	0.25932	-0.07412	5.47094	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073	-2.14073
3	5.85615	-7.46713	5.26652	3.90121	5.53167	3.37915	3.75711	-1.97117	-1.17714	-4.56112	-2.87212	-5.44672	-3.63942	-2.17412	-4.37712	-6.14291	-1.20312	-0.01495
4	-4.14612	-1.73612	-1.71212	-5.71622	-3.07812	4.32512	5.27542	-7.52904	-5.07721	-1.13641	-1.04001	-0.02212	-2.10512	-7.41112	-0.01612	-0.07502	-0.04107	-0.01495
5	6.33154	6.49272	3.33321	2.20958	0.00921	-6.49094	-4.36461	-3.29122	-5.75812	-2.16541	-0.01121	-0.01402	-0.01372	-0.01372	-0.01372	-0.01372	-0.01372	-0.01372
6	3.01621	-3.15312	5.95612	-0.01622	8.30521	3.12842	8.00212	3.71212	3.20712	-1.26912	-5.14212	-0.01921	-0.01921	-0.01921	-0.01921	-0.01921	-0.01921	-0.01921
7	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312
8	-2.35121	-2.51642	1.53321	4.17432	-0.79312	3.79382	-6.27921	-3.42312	-1.75202	-2.17812	-7.65541	-3.32742	3.63452	-1.36122	-1.71212	8.04412	0.01604	0.01604
9	8.46504	3.74194	4.41822	3.44322	3.90462	3.40462	2.61602	3.24612	-2.44122	-5.13112	-7.65541	-0.53212	-6.54892	-6.87612	-4.71212	1.52612	6.70461	-3.34412
10	5.36122	3.04372	7.70742	6.11521	1.49522	-1.27021	3.11312	6.64421	6.43512	-4.14812	-1.25742	-2.01972	-4.19121	-5.58712	-2.77072	-4.42121	-1.54612	-1.01312
11	-0.02242	-0.27122	3.27322	3.44622	3.30142	3.90192	3.41312	3.87302	1.53921	1.97702	6.31654	-4.99521	-1.93002	-0.01031	-0.01074	-0.02502	-0.02542	-0.01254
12	-1.25172	-0.01201	-0.20921	-7.10642	-7.36212	-4.42302	-4.48272	-2.90212	-1.58172	2.34112	7.11607	0.01284	0.01284	6.50262	0.01784	5.42812	1.54612	3.01342
13	1.68744	-1.10712	3.12412	0.01582	-1.07412	-4.10382	-1.33222	-1.56612	-1.43712	-1.23521	-2.27142	3.18312	6.45192	-3.34121	-7.79212	-4.07312	-0.01272	0.01372
14	-1.78162	3.75494	3.10472	6.99522	3.54462	6.19122	3.22512	3.62502	1.16021	2.46721	-2.12021	-9.38421	-3.98121	-0.34121	-7.79212	-4.07312	-1.59072	-1.05042
15	-0.01482	-7.23012	-5.05612	-5.36222	2.24312	-9.37412	3.62312	-5.12552	-1.21812	-4.63212	-5.41121	4.71912	7.55872	0.01612	0.01794	0.01612	0.02712	0.03712
16	6.79472	3.64612	3.35112	-1.18672	-3.43712	-2.00021	-3.50212	-4.56712	2.09942	-1.12942	4.23712	-1.53572	1.20612	-6.02212	-4.81512	-6.01512	-0.01622	-1.22342
17	6.85612	-1.25972	-4.07772	-3.96302	-4.31912	-2.29602	-3.32312	-3.51512	-1.49042	-1.43512	1.47712	4.10272	7.80212	9.08742	0.01512	0.02512	0.04942	0.01912
18	0.01642	6.79172	6.10772	6.65412	4.00702	3.18252	1.60312	1.81542	8.48464	-1.84912	-1.85342	-4.40212	-4.63972	-0.01794	-0.02034	-0.04074	-0.04674	-0.01374
19	0.01421	1.94921	-4.04972	-1.32482	4.78922	1.62492	1.75712	1.02321	-8.85721	-1.94312	-2.47202	-4.72412	-1.32142	-1.34212	-2.89742	6.07952	0.01624	0.01624
20	-0.01231	-0.03242	-0.01442	-7.19712	-9.40712	-9.82962	-0.04812	-5.01712	-1.20982	2.53942	6.46121	3.10291	8.45592	7.36472	0.01624	8.90012	0.01642	0.01624
21	-0.02192	-0.01802	-1.53612	-7.02112	-7.58712	-5.96812	-3.79322	-1.22432	-1.60412	1.84442	1.44912	1.77212	8.94312	0.01624	0.01624	0.01624	0.01624	0.01624
22	-0.01942	-1.45722	1.42122	3.00432	1.78704	-1.06412	-1.44812	-1.49012	2.93542	4.47712	7.80521	6.85412	2.81731	-0.01031	-0.04512	-0.01242	-0.15642	-0.04122
23	6.12674	-4.71121	-0.01662	-0.01642	-0.48772	-1.75172	-2.24312	-1.51192	2.14542	1.65542	0.74892	0.01104	0.01942	0.01147	0.01755	6.00712	5.35912	-5.44412
24	-0.01602	-0.01942	-0.01142	-1.76412	-0.01612	-9.74472	-0.07072	-5.43112	-8.63121	2.46542	6.07742	0.01344	0.01937	0.01147	0.01310	9.02112	0.01624	0.01624
25	0.01942	-3.46122	-4.04321	-4.10321	-1.81112	-1.09112	-1.09772	-1.21282	-1.21282	9.49912	-2.12021	3.44621	-2.30021	3.44421	3.68112	0.01122	0.01942	3.74122
26	0.01547	0.01547	0.01134	0.01134	0.01134	0.01134	0.01134	0.01134	0.01134	-8.12742	-4.13642	-6.94342	-5.41121	-0.01312	-0.01312	-0.01312	-0.01312	-0.01312
27	-6.81312	-9.19042	-7.05212	-4.22612	-4.86462	-3.27742	-1.54842	-1.64472	2.75792	3.17452	5.97992	0.01474	0.01914	0.01914	0.01914	0.01914	0.01914	0.01914
28	3.56494	-1.31402	-2.10402	-1.98712	-1.91412	-4.75224	-1.42512	-1.46542	-1.24212	-2.58021	-1.99902	0.02231	5.04674	-1.73972	0.01291	0.01554	0.01942	0.01942
29	-0.01822	-1.30202	-3.06412	-4.30822	-5.44112	-2.44422	-7.74612	-1.04542	-2.07821	1.60212	1.55492	5.90961	0.01192	0.01374	0.01244	0.01794	0.02442	3.65572
30	-2.07122	-3.39212	-6.46412	-3.13642	-1.63172	-3.23942	-1.44912	-4.06121	-1.67942	6.44121	-2.94972	8.31921	5.64921	6.46012	8.51921	3.64612	-5.39121	-1.44972
31	3.07012	7.16164	7.43112	1.97412	3.94874	6.25492	3.05072	3.17702	-1.67512	-1.67512	-4.74972	-4.31102	-5.94302	-0.01121	-0.01921	-0.01921	-0.01794	-5.28342
32	-9.16442	-1.48712	-2.73112	-5.87912	-6.30452	-4.86712	-2.18112	-5.12552	3.83034	7.04892	7.29512	8.84192	-0.01121	2.44212	1.79412	-2.50112	-2.92972	-2.84012

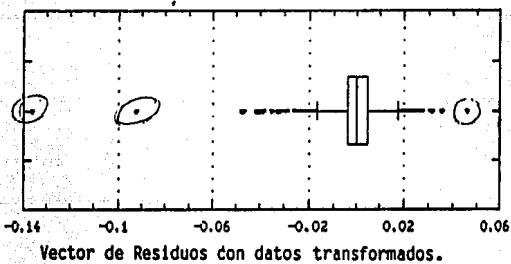
Tabla 14b



Figura 10a  
DIAGRAMA DE CAJA



Sexo Masculino  
DIAGRAMA DE CAJA



Sexo Femenino  
Figura 10 b

VALORES DE COMPARACION DE LOS DATOS TRANSFORMADOS  
 Sexo Masculino

C.E.	..8..	..9..	..10..	..15..	..20..	..25..	..30..	..40..	..45..	..50..	..55..	..60..	..65..	..70..	..75..	..80..		
1	6.789041	-5.037648	-4.482284	-4.063718	-3.779414	2.677792	-1.963818	-1.230674	-4.217742	-1.177218	-2.347218	-3.460184	-4.543118	-5.602118	-7.185648	-8.755418	-1.047118	
2	-4.993118	-7.294318	-4.460218	-5.939418	-4.941118	-1.514818	-2.045118	-1.799518	-4.140318	6.102018	1.957048	1.431548	4.984418	6.641918	8.481218	1.050518	1.280018	1.305218
3	1.804378	0.618758	2.728518	3.288058	2.685468	2.187468	1.355958	3.778318	3.340058	-2.254618	-1.165818	-2.708918	-6.609718	-6.161048	-5.702718	-7.351718	-9.254818	
4	4.376118	1.531918	6.157018	3.620018	3.019878	3.191218	1.753648	1.099918	2.746874	-2.728218	-1.945518	-2.077118	-1.044188	-0.053918	-5.144718	-6.400218	-7.022718	-9.254818
5	-2.450418	-2.808918	-2.811718	-2.245218	-1.485718	-1.454118	-1.794118	-4.463218	-1.250118	2.377218	7.465118	1.303118	3.905418	2.521918	3.283318	4.007318	4.483418	3.842618
6	-7.471518	-2.964018	-2.728118	-3.261018	-1.964818	-1.558718	-1.411118	-7.156418	2.649618	4.245518	7.785618	1.345318	1.902018	3.641118	1.374718	4.177308	5.000018	6.047918
7	-4.930918	-5.372118	-4.891618	-4.154818	-3.222118	-2.721218	-2.045018	-1.292018	-4.292118	4.945118	1.292968	2.446718	3.257918	4.725678	6.044618	7.490118	9.126498	0.020397
8	2.198448	3.523018	2.240018	2.026428	1.635778	1.357198	3.790018	6.178998	2.101968	-2.081618	-1.172818	-1.700718	-2.264218	-2.994618	-3.544418	-4.267218	-3.233718	
9	1.280548	1.753178	1.245118	1.086158	0.970544	1.138428	5.487028	1.290078	1.118548	-1.135548	-2.902018	-4.775148	-9.141748	-1.253518	-1.253518	-1.721118	-2.304048	2.797918
10	5.082018	3.364048	4.961818	4.312018	3.580958	3.042778	2.025318	1.206488	4.477578	-4.429718	-1.431718	-2.493818	-3.419488	-4.822318	-4.160218	-7.621018	-9.294918	-0.011171
11	3.749318	3.961438	4.318081	3.190448	2.645214	2.107918	1.544974	0.487158	3.116948	3.294648	-1.054118	-1.847418	-2.683018	-3.574218	-4.547018	-5.654318	-6.873218	-6.243218
12	4.365474	4.827104	4.478548	3.877418	3.225684	2.555874	1.873118	1.174548	4.021564	-2.902218	-1.278118	-2.240118	-3.253118	-4.538118	-6.875918	-8.256418	-9.593418	0.017303
13	-0.021218	-0.467228	-7.818228	-6.406498	-5.442718	-4.466158	-3.286198	-2.059798	6.911814	2.247798	3.932698	3.712378	3.511898	3.742328	0.012092	0.014695	0.017303	0.017303
14	9.732594	-1.021818	3.528218	6.258218	6.970118	3.443118	3.909518	6.201584	8.563225	-6.161718	-2.722118	-4.771118	-6.930218	-9.234818	-1.179528	-1.462468	-1.773718	-2.118418
15	-8.218418	-6.845158	-4.096318	-4.973218	-5.801318	-4.594118	-1.348818	-2.112328	-7.122618	1.162218	2.298448	4.028918	6.821118	7.749098	4.964018	4.233318	1.502838	1.757478
16	-1.081618	-1.142518	-1.0914	-0.177218	-1.434918	-5.08158	-4.431818	-2.7408	-9.518718	8.425718	3.027218	5.303128	7.701818	1.024918	1.622118	1.977818	2.345918	4.017303
17	-2.05748	-1.172218	-2.01648	-1.745278	-1.452318	-1.150418	-8.437718	-5.288218	-1.81048	1.730218	5.754638	1.008428	1.443038	3.192228	3.4914818	3.0874918	3.7422718	4.939718
18	1.639818	1.732228	1.602118	1.791618	1.157528	3.197918	6.721718	4.214718	1.443128	-1.429518	-4.394638	-6.038848	-1.674478	-1.555818	-1.767318	-2.440318	-2.992518	-2.748418
19	3.757974	-2.975218	-1.684718	-1.780118	-2.453718	-2.102118	-1.561118	-6.642318	2.208128	2.276518	1.031548	1.843118	4.497178	4.547418	4.755494	6.370218	6.870218	6.022484
20	-0.010734	-0.011393	-0.010304	-0.10818	-0.377318	-0.400148	-4.400118	-2.759518	-9.444418	5.354918	3.002418	5.262318	7.647488	0.010184	0.010094	0.016106	0.019619	0.023473
21	-6.401618	-0.875318	-1.136118	-1.128718	-5.935418	-4.688718	-3.443918	-7.293918	7.218184	6.249938	6.187228	3.982588	3.797198	0.010206	0.012807	0.015314	0.018373	0.018373
22	2.784538	2.961478	-2.740218	-2.262728	-1.965384	-1.557318	-1.141428	7.157118	-2.450634	-2.444418	-7.788918	-1.343518	-1.992018	-2.643118	-3.748718	-4.178918	-5.071018	-4.090118
23	0.010360	0.011632	0.012504	0.066018	7.458494	3.909448	4.231618	3.715948	3.793218	-9.287718	-2.555118	-5.180118	-7.574218	-0.010024	-0.012802	-0.015314	-0.018373	-0.021104
24	6.817318	7.202188	-6.441718	5.7857718	4.8123718	3.819418	2.7859718	1.752418	0.020448	-9.942118	-1.907118	-3.242718	-4.852418	-4.469518	-2.831818	-1.023118	-1.244318	-1.493118
25	2.3747918	2.718998	3.233418	2.184748	1.8175218	1.1418	1.025418	6.4179718	2.263648	-2.218348	-7.201418	-1.242218	-1.873418	-2.443118	-3.160418	-3.864118	-4.703618	-5.613618
26	-1.7408	-1.86448	-1.713218	-1.483218	-1.233918	-1.1733918	-7.116518	-4.493118	-1.538418	1.523494	4.099448	6.367448	1.246778	1.658678	2.119378	2.629448	3.194588	3.823498
27	2.232318	2.348818	1.719718	1.864678	1.567918	1.183948	3.114498	5.712118	1.954828	-1.373718	-4.231618	-1.018718	-1.583318	-1.079018	-1.649618	-2.231618	-4.066518	-4.183618
28	0.017318	2.143618	1.931118	3.606448	2.832718	2.344718	1.644718	0.022948	1.537948	3.500918	-1.124418	-1.964318	-2.840578	-3.811718	-4.864418	-6.024818	-7.344418	-7.781818
29	1.829048	1.932138	1.712518	1.551918	1.811118	1.022918	4.675318	5.7012218	0.609718	-1.573198	-5.115918	-8.944518	-1.302418	-1.735218	-2.164418	-2.747418	-3.445618	-4.000218
30	-2.284218	-2.412918	-2.238618	-1.998218	-1.612418	-1.277518	-5.363518	-3.871118	-2.010218	1.996478	6.388998	1.119918	1.628548	3.167418	2.759348	3.428058	4.178498	4.951818
31	3.564232	3.769418	4.971718	3.027788	2.518848	1.959448	1.463618	3.171818	2.160228	-2.109718	-6.980518	-1.749218	-2.540118	-3.385918	-4.342518	-5.353118	-6.522118	-7.801818
32	6.1570218	0.616718	0.794318	8.5213218	3.757918	4.361494	3.343618	2.096918	7.1786318	-7.108618	-2.261318	-3.998318	-5.808118	-7.719618	-9.897318	-0.012246	-0.014911	-0.017862

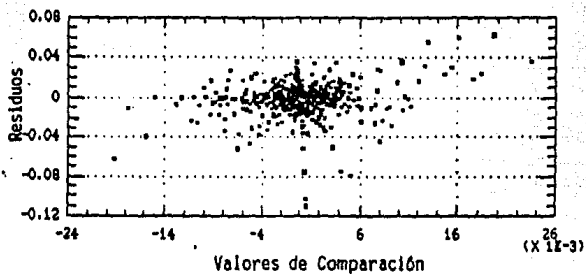
Tabla 16a

VALORES DE COMPARACION DE LOS DATOS TRANSFORMADOS  
Sexo Femenino

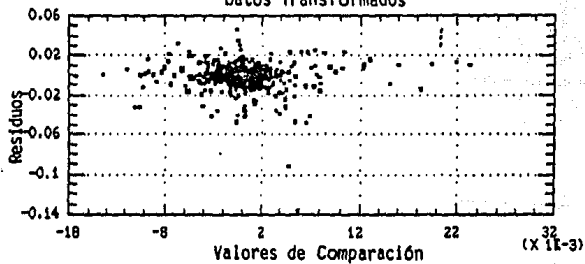
E.T.	0	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
1	-1.1816E-4	-1.2346E-4	-1.1622E-4	-9.2426E-5	-8.3254E-5	-6.8774E-5	-6.8712E-5	-6.0159E-5	-1.0448E-5	1.0584E-5	3.3218E-5	3.7894E-5	8.3415E-5	1.1657E-4	1.6794E-4	1.8002E-4	3.1889E-4	2.4156E-4
2	-9.1327E-4	-8.7794E-4	-3.5114E-3	-2.0309E-2	-1.3538E-2	-8.0378E-2	-1.4757E-2	-9.2717E-2	-3.2181E-4	2.3936E-4	1.0224E-3	1.7827E-3	2.6737E-3	2.5323E-3	4.6467E-3	3.3451E-3	6.7292E-3	8.9405E-3
3	-4.4932E-3	-3.3787E-3	-2.3074E-3	-6.0958E-3	-1.7318E-2	-1.8265E-2	-1.0011E-2	6.1941E-3	2.1513E-2	-2.1757E-4	-6.4917E-4	-1.1914E-4	-1.7594E-4	-8.3991E-5	3.4918E-5	-6.9302E-5	-8.2554E-5	-1.2554E-5
4	1.1292E-4	1.8078E-4	1.4724E-4	1.4328E-4	1.2484E-4	9.7035E-5	7.1312E-5	4.4152E-5	1.5219E-5	-1.5453E-4	-6.8479E-4	-8.4984E-4	-1.2504E-3	-1.6773E-3	-2.1261E-3	-2.4351E-3	-3.2064E-3	-3.8212E-3
5	-1.3551E-4	-1.4296E-4	-1.5075E-4	-1.3079E-4	-1.0918E-4	-8.7412E-5	-4.4291E-4	-2.9804E-4	-1.3914E-4	1.2968E-4	6.3986E-4	7.6518E-4	1.1273E-3	1.3122E-3	1.9248E-3	2.3741E-3	2.8891E-3	1.9323E-3
6	-3.6761E-4	-3.8016E-4	-3.318E-4	-3.0567E-4	-2.5642E-4	-2.0413E-4	-1.5001E-4	-2.2942E-4	-2.2241E-4	3.2597E-4	1.0241E-3	1.7893E-3	2.4209E-3	3.5296E-3	4.4979E-3	5.3649E-3	6.7197E-3	8.0521E-3
7	-5.4044E-4	-5.4561E-4	-5.8161E-4	-4.7072E-4	-3.9467E-4	-3.1431E-4	-2.3019E-4	-1.4204E-4	-6.9649E-5	3.0196E-4	1.3751E-3	2.5805E-3	4.0512E-3	5.4944E-3	6.9221E-3	8.5184E-3	1.0103E-2	6.6124E-2
8	1.4014E-4	1.4725E-4	1.5079E-4	1.3647E-4	1.1328E-4	9.9015E-5	6.4014E-5	4.0072E-5	1.4166E-5	-1.4342E-4	-6.8022E-4	-9.8739E-4	-1.1572E-3	-1.5527E-3	-1.9737E-3	-2.4197E-3	-2.9644E-3	-3.5447E-3
9	1.3755E-4	1.3287E-4	1.2285E-4	1.0718E-4	8.9680E-5	7.1534E-5	5.2575E-5	3.2528E-5	1.1294E-5	-1.1423E-4	-5.5895E-4	-8.2582E-4	-9.2194E-4	-1.2361E-3	-1.9559E-3	-2.3624E-3	-2.8231E-3	-3.4231E-3
10	2.2932E-4	2.5007E-4	2.7131E-4	2.0078E-4	1.4859E-4	1.3421E-4	9.9647E-4	6.1074E-4	2.1198E-4	-1.1432E-4	-6.7389E-4	-1.1742E-3	-1.7197E-3	-2.3801E-3	-2.9559E-3	-3.6329E-3	-4.4328E-3	-5.2963E-3
11	-7.2601E-4	-7.5641E-4	-6.9987E-4	-6.081E-4	-5.1103E-4	-4.0092E-4	-2.9641E-4	-1.8479E-4	-4.4141E-5	6.8487E-5	2.0374E-4	3.5392E-4	5.2177E-4	7.0259E-4	8.9487E-4	1.1031E-3	1.3424E-3	1.6027E-3
12	-4.4023E-4	-4.4624E-4	-5.7994E-4	-5.2079E-4	-4.3704E-4	-3.479E-4	-2.5371E-4	-1.5822E-4	-9.4731E-5	3.5537E-4	1.4755E-3	2.9487E-3	4.4893E-3	6.0147E-3	7.6654E-3	9.4505E-3	1.1490E-2	1.3797E-2
13	-9.2015E-4	-8.0103E-4	-6.5727E-4	-6.3174E-4	-6.3975E-4	-5.5564E-4	-4.0825E-4	-3.5274E-4	-8.7731E-5	8.8439E-4	2.7462E-3	4.8762E-3	7.1566E-3	9.6024E-3	0.0121E-2	0.0150E-2	0.0184E-2	0.0219E-2
14	-7.7546E-4	-8.1044E-4	-1.0777E-3	-6.5143E-4	-5.4452E-4	-4.2504E-4	-3.1974E-4	-1.9797E-4	-4.8714E-5	6.9472E-5	2.1827E-4	3.8042E-4	5.6064E-4	7.5218E-4	9.5801E-4	1.1174E-3	1.4748E-3	1.7149E-3
15	1.9077E-3	1.9934E-3	1.8491E-3	1.6022E-3	1.3640E-3	1.0499E-3	7.8641E-4	6.4062E-4	1.8897E-4	-1.7086E-5	-5.3402E-5	-9.361E-5	-1.3789E-4	-1.8697E-4	-2.2861E-4	-2.7944E-4	-3.5138E-4	-4.2274E-4
16	-7.9747E-4	-8.3317E-4	-1.7081E-3	-1.7081E-3	-1.6183E-3	-4.4726E-4	-2.2874E-4	-2.0233E-4	-7.0644E-5	7.1439E-5	2.2464E-4	3.9155E-4	5.7644E-4	7.7247E-4	9.8494E-4	1.2149E-3	1.4772E-3	1.7452E-3
17	2.8816E-4	3.4381E-4	3.18E-4	2.743E-4	3.1877E-4	1.8431E-4	1.3561E-4	8.3966E-5	3.5147E-5	-2.8645E-5	-5.2575E-5	-1.6143E-4	-2.378E-4	-3.187E-4	-4.6432E-4	-5.0121E-4	-6.9841E-4	-7.821E-4
18	3.4443E-4	3.5997E-4	3.2284E-4	2.8929E-4	2.4268E-4	1.9318E-4	1.4192E-4	8.7912E-5	3.2513E-5	-1.0849E-5	-9.4242E-5	-1.691E-4	-2.4899E-4	-3.3399E-4	-4.2542E-4	-5.1477E-4	-6.3805E-4	-7.643E-4
19	4.1085E-4	4.2935E-4	3.7151E-4	3.4504E-4	2.8944E-4	2.3042E-4	1.6934E-4	1.0446E-4	3.4316E-5	-3.6797E-4	-1.1561E-3	-2.014E-3	-2.9698E-3	-3.9937E-3	-5.0744E-3	-6.2594E-3	-7.6104E-3	-9.094E-3
20	-0.0150E-2	-0.0197E-2	-0.0191E-2	-8.8153E-3	-7.2995E-3	-5.8896E-3	-4.228E-3	-2.4401E-3	-9.5025E-4	9.4052E-4	2.3549E-3	3.1504E-3	3.9804E-3	0.0118E-2	0.0126E-2	0.0159E-2	0.0194E-2	0.0233E-2
21	-4.8031E-3	-2.4639E-3	3.9441E-3	3.4248E-3	-2.8747E-3	-2.796E-3	-1.482E-3	-1.0411E-3	-1.615E-4	3.4544E-4	1.1493E-3	2.0021E-3	2.8493E-3	3.9631E-3	5.0364E-3	6.2163E-3	7.5524E-3	8.0164E-3
22	-1.222E-3	2.3483E-3	-1.0978E-3	-2.4917E-3	-2.2584E-3	-1.795E-3	-1.3112E-3	-8.14E-4	-2.829E-4	2.8705E-4	9.0188E-4	1.3749E-3	2.1167E-3	3.1274E-3	3.9584E-3	4.6428E-3	5.7349E-3	7.243E-3
23	-1.8142E-3	8.1454E-4	7.3524E-4	6.3422E-4	5.5052E-4	4.3823E-4	3.221E-4	1.9942E-4	6.9217E-5	-6.991E-4	-2.1987E-3	-3.4341E-3	-5.6479E-3	-7.9523E-3	-9.8501E-3	-1.1419E-2	-1.3444E-2	-1.6173E-2
24	-1.4931E-3	-1.3429E-3	-1.2421E-3	-1.0733E-3	-9.0542E-4	-7.7074E-4	-6.2976E-4	-4.3799E-4	-1.1844E-4	1.2909E-4	3.6142E-4	6.3039E-4	9.2893E-4	1.2446E-3	1.5874E-3	1.9548E-3	2.3802E-3	2.8433E-3
25	2.8778E-3	2.9551E-3	2.7343E-3	2.3749E-3	1.9926E-3	1.5840E-3	1.1673E-3	7.2149E-4	2.5059E-4	-2.5237E-4	-9.3933E-4	-1.8164E-3	-2.6461E-3	-3.7495E-3	-5.0545E-3	-6.4061E-3	-7.8281E-3	-9.2364E-3
26	-5.5095E-3	-4.6222E-3	-1.4273E-3	-8.1674E-3	-1.7488E-3	-1.4078E-3	-1.0245E-3	-6.4057E-4	-2.2313E-4	2.2474E-4	7.0474E-4	1.2312E-3	1.8119E-3	2.4332E-3	3.0994E-3	3.8223E-3	4.6446E-3	5.5505E-3
27	8.8001E-4	8.5498E-4	7.9187E-4	6.4873E-4	5.7775E-4	4.3964E-4	3.2805E-4	2.0921E-4	7.2644E-5	-1.3447E-4	-8.3076E-4	-8.024E-4	-9.927E-4	-7.9551E-4	-1.0128E-3	-1.2491E-3	-1.5191E-3	-1.8152E-3
28	3.6312E-3	3.7844E-3	3.9064E-3	3.0419E-3	2.3514E-3	2.0310E-3	1.4728E-3	9.2449E-4	3.2801E-4	-2.2431E-4	-1.093E-3	-1.7749E-3	-2.6174E-3	-3.5112E-3	-4.4724E-3	-5.5174E-3	-6.7081E-3	-8.0188E-3
29	1.1871E-3	1.24054E-3	1.14741E-3	9.3694E-4	6.3493E-4	4.6791E-4	3.1934E-4	1.8279E-4	1.0514E-4	-1.0432E-4	-3.3404E-4	-5.323E-4	-7.151E-4	-9.151E-4	-1.1480E-3	-1.4805E-3	-1.9895E-3	-2.6174E-3
30	-1.2314E-3	-1.3745E-3	-1.2893E-3	-1.1207E-3	-9.4054E-4	-7.4841E-4	-5.956E-4	-4.4058E-4	-1.1921E-4	1.1375E-4	3.7930E-4	6.5474E-4	9.6452E-4	1.2978E-3	1.6481E-3	2.021E-3	2.4747E-3	2.9174E-3
31	-2.3514E-3	-2.3231E-3	-1.1744E-3	-1.491E-3	-1.3864E-3	-1.2438E-3	-9.282E-4	-6.7448E-4	-1.9545E-4	2.0144E-4	6.1361E-4	1.0404E-3	1.6275E-3	2.1832E-3	2.7809E-3	3.4044E-3	4.1702E-3	4.9644E-3
32	1.2322E-3	1.3817E-3	1.3704E-3	1.1104E-3	9.2108E-4	7.4164E-4	6.4564E-4	3.7481E-4	1.1712E-4	-1.1442E-4	-7.207E-4	-1.1048E-3	-9.3574E-3	-1.282E-2	-1.633E-2	-2.044E-2	-2.4492E-2	-2.947E-2

Tabla 16b

Gráfica 16a  
Sexo Masculino  
Datos Transformados



Gráfica 16b  
Sexo Femenino  
Datos Transformados



Cuadro 13a  
ANALISIS DE REGRESION  
con datos transformados  
SEXO MASCULINO

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.000917	0.000656	-1.3967	0.1630
VECCOMMAS4	0.90939	0.1507	6.0345	0.0000

R-SQ. (ADJ.) = 0.0580 SE= 0.015750 MAE= 0.010092 DurWat= 0.381  
 Previously: 0.0000 0.000000 0.000000 0.000

576 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 14a  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	0.00903350	1	0.00903350	36.4146	.0000
Error	0.142394	574	0.000248073		
Total (Corr.)	0.151428	575			

R-squared = 0.0596556  
 R-squared (Adj. for d.f.) = 0.0580174

Std. error of est. = 0.0157503  
 Durbin-Watson statistic = 0.381114

Cuadro 13b  
ANALISIS DE REGRESION  
con datos transformados  
SEXO FEMENINO

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.000775	0.000538	-1.4407	0.1502
VECCOMPFA	-0.195333	0.143153	-1.3645	0.1729
R-SQ. (ADJ.) = 0.0015 SE= 0.012899 MAE= 0.007569 DurWat= 0.511				

576 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 14b  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	0.000309781	1	0.000309781	1.86187	.1729
Error	0.0955032	574	0.000166382		
Total (Corr.)	0.0958130	575			

R-squared = 3.23319E-3

R-squared (Adj. for d.f.) = 1.49666E-3

Std. error of est. = 0.0128989

Durbin-Watson statistic = 0.510701



Tabla 17a  
 Datos influenciales  
 Sexo Masculino

Obs. Number	Std. Residual	Leverage	Mahalanobis Dist.	DFITS
126	-0.67328	0.01272	6.39823	-0.07643
179	1.69226	0.00958	4.53348	0.16643
180	0.74230	0.01297	6.54521	0.08510
231	0.45844	0.01046	5.06850	0.04713
232	0.43775	0.01509	7.79664	0.05419
233	0.72140	0.02154	11.6385	0.10704
234	0.28424	0.03004	16.7803	0.05002
343	-0.90673	0.01221	6.09610	-0.10080
344	-0.78634	0.01343	6.81355	-0.09173
345	0.18972	0.01179	5.85214	0.02073
355	1.46014	0.00714	3.12918	0.12382
356	1.77229	0.01131	5.56559	0.18952
357	2.90030	0.01732	9.12019	0.38507
358	2.97522	0.02561	14.0877	0.48234
359	2.98545	0.03715	21.1480	0.58641
360	1.03481	0.05236	30.7185	0.24325
374	1.34853	0.00761	3.40392	0.11810
375	1.92927	0.01130	5.56207	0.20625
376	1.34431	0.01638	8.56043	0.17348
377	1.07724	0.02345	12.7875	0.16694
378	0.52510	0.03278	18.4539	0.09666
397	-0.45865	0.01203	5.99297	-0.05062
398	-1.16764	0.01322	6.69280	-0.13516
399	-1.29851	0.01163	5.75486	-0.14085
409	1.62979	0.00687	2.96980	0.13551
410	1.25703	0.01087	5.30843	0.13176
411	0.36514	0.01666	8.72586	0.04753
412	-1.56147	0.02465	13.5056	-0.24821
413	-2.87268	0.03578	20.3008	-0.55337
414	0.03051	0.05047	29.5095	0.00703
503	2.20341	0.00662	2.82932	0.17993
573	0.65184	0.01061	5.15869	0.06751
574	0.76815	0.01537	7.96073	0.09597
575	0.95128	0.02200	11.9127	0.14267
576	0.43335	0.03075	17.2099	0.07718

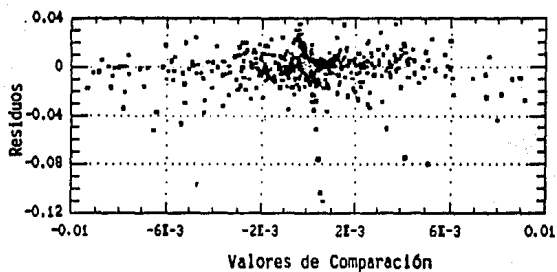
Number of flagged observations (high leverage or DFITS) = 35

Tabla 17b  
 Datos Influenciales  
 Sexo Femenino

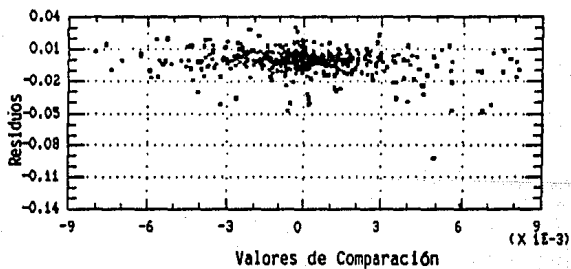
Obs. Number	Std. Residual	Leverage	Mahalanobis Dist.	DFITS
123	2.03085	0.00742	3.29531	0.17564
124	2.00951	0.01045	3.06491	0.20653
125	2.16751	0.01469	7.56001	0.26467
126	1.11583	0.02031	10.9002	0.16065
217	0.02558	0.01413	7.22628	0.00306
218	-0.08837	0.01525	7.89211	-0.01100
219	0.15822	0.01332	6.75196	0.01839
220	0.16866	0.01052	5.10278	0.01739
230	0.43871	0.01280	6.44242	0.04995
231	0.88758	0.01979	10.5877	0.12610
232	-0.41128	0.02931	16.3326	-0.07146
233	-0.78212	0.04262	24.5546	-0.16502
234	1.50350	0.06025	35.7993	0.38068
304	2.84150	0.00178	0.02793	0.12014
305	3.65659	0.00180	0.03839	0.15539
342	-0.05596	0.01221	6.09541	-0.00622
343	-2.65724	0.01565	8.12575	-0.33502
344	-2.64843	0.01691	8.87609	-0.34737
345	-0.98825	0.01474	7.59119	-0.12089
346	-0.65042	0.01159	5.73367	-0.07044
356	0.77839	0.01419	7.26463	0.09339
357	1.51826	0.02205	11.9458	0.22800
358	1.08801	0.03277	18.4466	0.20025
359	1.27305	0.04774	27.7774	0.28504
360	1.25053	0.06756	40.5923	0.33662
375	1.97176	0.00471	1.71773	0.13563
376	2.05308	0.00630	2.64342	0.16353
377	2.20872	0.00854	3.94557	0.20498
378	0.80723	0.01150	5.68205	0.08708
410	1.60202	0.00904	4.23991	0.15304
411	1.22484	0.01351	6.86187	0.14333
412	0.34838	0.01956	10.4543	0.04921
413	-0.11915	0.02799	15.5311	-0.02022
414	-0.24863	0.03912	22.3701	-0.05017
502	1.88967	0.00566	2.26824	0.14255
503	2.44729	0.00749	3.33282	0.21258

Number of flagged observations (high leverage or DFITS) = 36

Gráfica 17a  
Sexo masc.



Gráfica 17b  
Sexo fem.



Cuadro 15a  
SEXO MASCULINO  
ANALISIS DE REGRESION  
con datos transformados  
sin puntos influyentes

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.001677	0.000645	-2.5977	0.0096
veccommas4	0.030967	0.226692	0.1366	0.8914

R-SQ. (ADJ.) = 0.0000    SE=    0.014996    MAE=    0.009522    DurbinWat= 0.421  
 Previously: 0.0000    0.000000    0.000000

541 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 16a  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	0.00000419626	1	0.00000419626	0.0186607	.8929
Error	0.121206	539	0.000224871		
Total (Corr.)	0.121210	540			

R-squared = 3.46198E-5  
 R-squared (Adj. for d.f.) = 0

Std. error of est. = 0.0149957  
 Durbin-Watson statistic = 0.420518

Cuadro 15b  
SEXO FEMENINO  
ANALISIS DE REGRESION  
con datos transformados  
sin puntos influenciados

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.001438	0.000515	-2.7900	0.0055
veocomfem	-1.261879	0.215998	-5.8421	0.0000

R-SQ. (ADJ.) = 0.0579 SE= 0.011973 MAE= 0.007201 DurWat= 0.601  
 Previously: 0.0000 0.014996 0.009522 0.421  
 540 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 16b  
ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	0.00489253	1	0.00489253	34.1300	.0000
Error	0.0771221	538	0.000143350		
Total (Corr.)	0.0820147	539			

R-squared = 0.0596543  
 R-squared (Adj. for d.f.) = 0.0579065

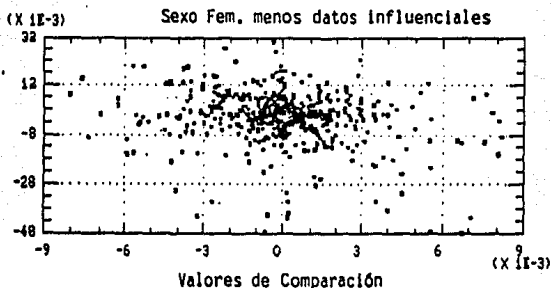
Std. error of est. = 0.0119729  
 Durbin-Watson statistic = 0.601368

Tabla 18

## DATOS INFLUENCIALES- Sexo Femenino.

Obs. Number	Std. Residual	Leverage	Mahalanobis Dist.	DFITS
19	1.72405	0.00623	2.36296	0.13652
20	1.65124	0.00663	2.58179	0.13495
34	-1.54234	0.01203	5.52964	-0.17021
35	-0.42486	0.01690	8.21386	-0.05570
36	0.58841	0.02333	11.8042	0.09094
54	-0.69172	0.01144	5.20350	-0.07441
89	2.67377	0.00463	1.49687	0.18242
106	-4.20968	0.01207	5.55073	-0.46532
107	-4.19428	0.01695	8.24533	-0.55080
108	-0.79384	0.02341	11.8497	-0.12290
109	-1.93946	0.01227	5.66170	-0.21619
110	-2.02048	0.01323	6.18916	-0.23397
111	-0.77747	0.01159	5.28597	-0.08418
122	1.88508	0.01167	5.32933	0.20482
176	-0.25226	0.00364	0.96076	-0.01525
213	1.91175	0.00186	0.00112	0.08254
214	0.65766	0.00186	0.00168	0.02840
221	-0.52354	0.01800	8.82796	-0.07089
256	0.01876	0.00187	0.00473	0.00081
325	0.87059	0.00747	3.03429	0.07551
326	0.68101	0.00798	3.31518	0.06109
327	0.44526	0.00710	2.83404	0.03765
328	0.39469	0.00582	2.13715	0.03019
334	-0.12253	0.00190	0.02384	-0.00535
335	-0.34509	0.00230	0.23767	-0.01657
349	-1.08878	0.00807	3.36228	-0.09820
365	-0.91038	0.00460	1.47676	-0.06186
366	-0.69339	0.00339	0.93445	-0.04164
367	-0.42401	0.00279	0.50402	-0.02245
368	-0.11051	0.00222	0.19274	-0.00521
369	0.05487	0.00190	0.02300	0.00240
370	0.15871	0.00190	0.02417	0.00693
378	-32768.0	-32768.0	-32768.0	-32768.0
397	0.88093	0.02213	11.1325	0.13253
413	-32768.0	-32768.0	-32768.0	-32768.0
414	-32768.0	-32768.0	-32768.0	-32768.0
432	0.62516	0.00455	1.45102	0.04226
466	-1.23764	0.00672	2.62579	-0.10177
483	1.42040	0.00186	0.00174	0.06135

Number of flagged observations (high leverage or DFITS) = 39



Gráfica 18

Cuadro 17

## ANALISIS DE REGRESION

Independent variable	coefficient	std. error	t-value	sig.level
CONSTANT	-0.001043	0.00043	-2.4268	0.0156
VECcomfem4	-0.874359	0.181627	-4.8140	0.0000

R-SQ. (ADJ.) = 0.0397    SI=    0.009971    MAE=    0.006672    DurbinWat= 0.588  
 Previously:    0.0000    0.000000    0.000000    0.000000    0.000  
 538 observations fitted, forecast(s) computed for 0 missing val. of dep. var.

Cuadro 18

## ANALISIS DE VARIANZA

Source	Sum of Squares	DF	Mean Square	F-Ratio	P-value
Model	0.00230392	1	0.00230392	23.1749	.0000
Error	0.0532862	536	0.0000994145		
Total (Corr.)	0.0555901	537			

R-squared = 0.0414449  
 R-squared (Adj. for d.f.) = 0.0396565

Std. error of est. = 9.97068E-3  
 Durbin-Watson statistic = 0.587526

Cuadro 19  
RELACION ENTRE LOS COEFICIENTES DE LA REGRESION  
"Valores de Comparación vs Residuos".

SEXO MASCULINO	(1)	(2)	(3)	(4)		
	$n_1^*$	$n_2^*$	$n_3^*$	$n_4^*$	(2)/(1)	(4)/(3)
$B_0$	576	534	576	541	0.62	0.93
$B_1$					0.94	0.93
e.s. $B_0$					0.88	0.98
e.s. $B_1$					1.37	1.50
t $B_0$					0.71	0.86
t $B_1$					0.68	0.02
SE					0.84	0.95
$R^2$					0.67	0.00
F					0.47	0.00

SEXO FEMENINO	(1)	(2)	(3)	(4)	(5)		
	$n_1^*$	$n_2^*$	$n_3^*$	$n_4^*$	$n_5^*$	(3)/(1)	(5)/(4)
$B_0$	576	539	576	540	538	1.35	1.83
$B_1$						0.91	0.69
e.s. $B_0$						0.93	0.83
e.s. $B_1$						1.39	0.84
t $B_0$						1.45	0.87
t $B_1$						0.66	0.82
SE						0.90	0.83
$R^2$						0.65	0.69
F						0.43	0.68

- (1) Regresión a partir de los datos crudos.  
 (2) Regresión a partir de los datos crudos menos los casos influyentes.  
 (3) Regresión a partir de los datos transformados.  
 (4) Regresión a partir de los datos transformados menos los influyentes.  
 (5) Regresión a partir de los datos de (4) menos dos puntos extremos en el caso del sexo femenino.



Se ve que la normalidad (gráficos 14a y 14b) mejoró bastante en el caso de los datos masculinos, mientras que en el otro caso empeoró.

La pendiente es de 0.79 (0.74) y 0.80 (0.73) para los datos completos (incompletos) masculinos y femeninos respectivamente. Es decir una pendiente aproximada en todos los casos a 0.765, por lo que  $p = 1-b = 1-0.765 = 0.235$  indica una reexpresión cercana a la raíz cuarta, lo que estabilizará la amplitud del ciclo.

Finalmente, se construyó una nueva tabla a partir de los datos transformados (tablas 13a y 13b), que al operarla por el proceso de la mediana pulida, se generó una tabla de residuos después de 10 iteraciones (tablas 14a y 14b). Se volvieron a calcular los valores comparativos de estos datos (tablas 16a y 16b) después del cálculo de los efectos (tablas 15) y los diagramas de tallo-y-hoja y de caja de los residuos obtenidos (figuras 9 y 10). Se graficaron los valores comparativos contra los residuos (gráficos 16). La gráfica presenta una tendencia lineal "casi horizontal" cuya pendiente es 0.09 (0.03) y -0.19 (-1.26) para el caso masculino y femenino respectivamente, para los datos transformados, (sin datos influyentes) [cuadros 13 y 14; gráficos 17 y tablas 17].

Para el caso del sexo masculino, las pendientes parecen indicar que la transformación que se hizo fue la más conveniente. Sin embargo, en el caso del sexo femenino, la gráfica 18 ilustra la misma gráfica 17b menos dos puntos aberrantes muy significativos (tabla 18). La pendiente varía de -1.26 a -0.87

(cuadros 17 y 18), es decir en un 70% solo por el simple hecho de no tomar estos dos puntos; lo que sugiere que la transformación es también aceptable. El cuadro 19 permite ver las diferencias entre los diferentes parámetros obtenidos del análisis de regresión, para los datos con y sin datos influenciales, antes y después de la transformación.

Finalmente, se presenta una tabla de los datos originales, marcados por los casos extremos (tablas 19a y 19b). Y como un resumen visual, para poder apreciar la utilidad de las técnicas aquí expuestas, se graficaron los residuos de los datos transformados contra la entidad federativa (con la edad) (gráficas 19a y 19b).

Como se observó a lo largo de esta presentación, se está siempre tan interesado en la secuencia residual, como en la secuencia suave. El residuo puede revelar casos aberrantes, tan bien como porciones de la secuencia que parece estar sujeta a grandes fluctuaciones.

Ahora se procede a dar una breve y posible interpretación demográfica a los casos comunes y a los casos extremos que se encontraron en esta segunda parte del último capítulo. Se debe mencionar que son grandes las lagunas que se tienen en el conocimiento tanto de los factores sociales, económicos, culturales, médicos, que están afectando directa e indirectamente a la mortalidad, como el nivel real de esta incidencia. Muy poco se conoce sobre la problemática causal, sobre la esperanza de vida y las probabilidades de sobrevivencia de las poblaciones que pertenecen a uno y a otro grupo social.

## 2- Interpretación de los resultados.

De esta manera, se puede señalar que al observar las gráficas de los datos crudos (gráficas 9a y 9b), se ve que en México ha prevalecido el patrón universal: elevadas tasas en el primer año de vida que desciendan rápidamente hasta alcanzar valores mínimos en el grupo de 10 a 14 años. A partir de esta edad, las tasas se incrementan, primero lentamente y a mayor ritmo después.

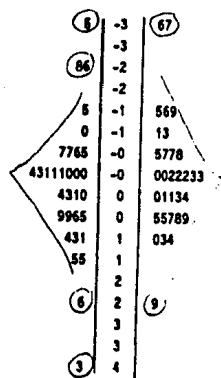
El análisis de los cambios en la mortalidad confirma una de las hipótesis de Almada Bay (1975), la que expresa que "el más profundo cambio en salud y en los patrones de enfermedades, se da en la niñez", ya que de todas las edades, las que cambiaron en mayor proporción fueron las de los niños.

En este estudio, se encontró que la tasa de mortalidad en los grupos infantiles y de edad avanzada (Ver Gráficas 19a y 19b), es siempre más elevada que entre los adultos. Como ya se dijo, los cambios registrados en la esperanza de vida no fueron uniformes en todas las edades, lo principal de este hecho consiste en los diferentes riesgos de muerte que se presentan en las distintas edades. Estos riesgos se refieren a las causas de muerte, que en un sentido amplio podrían dividirse en factores biológicos y factores sociales.

Asimismo, respecto de la mortalidad por sexos, se encontró que ésta es más elevada entre los hombres que entre las mujeres para casi todas las grupos de edad. Es necesario recalcar que México no solamente sigue el patrón universal de mortalidad diferencial por edad, sino que también sigue el patrón universal de

mortalidad por sexo.

A fin de observar los datos que se desvían claramente del conjunto de los efectos, se construyó el diagrama de tallo-y-hoja. Así, se puede comparar simultáneamente los dos grupos de efectos en un diagrama en espejo (por sexo). Se observa que, para ambos sexos, los Estados de Hidalgo y Oaxaca son los que presentan un efecto-ranglón mayor en valor absoluto indicando una esperanza de vida menor y el Estado de Quintana-Roo la mayor. También, se nota que el Estado de Puebla y el de Zacatecas muestran para el sexo masculino efectos absolutos grandes con una mortalidad grande y pequeña respectivamente.



masc. fem.  
• Efectos Ranglón

DIAGRAMA DE TALLO-y-HOJA  
EN ESPEJO

Unidad=0.1 112 representa 1.2

Figura 1

Ranglón	masc.	fem.
-3	9641	369
-2	740	159
-1	62	27
0	62	27
1	951	1
2	994	067
3	2	013

masc. fem.  
Efectos Columna

DIAGRAMA DE TALLO-y-HOJA  
EN ESPEJO

Unidad=1 112 representa 12

Figura 2

Por lo que respecta a los efectos-columna, el diagrama anterior no muestra cambios muy importantes, se nota que existe una consistencia en los datos; los efectos varían de mayor a menor conforme avanza la edad, excepto el hecho de que se presenta un "brinco" en la segunda columna; esto se debe a que el factor edad varía de 5 en 5 años cada columna, mientras que el cambio de la primera a la segunda es de un año (mortalidad infantil). Si se compara por sexo, vemos que para todas las edades los datos femeninos presentan efectos-columna mayores en valor absoluto que los datos masculinos; sobretodo en las edades 0, 65, 70, 75 y 80 años. Esto se debe probablemente a que existe una sobremortalidad masculina.

Dos posibles causas desempeñan cierto papel en este retraso relativo: se trata del alcoholismo y de los accidentes. Pueden existir otras, en particular el tabaco, pero estas dos causas están bien identificadas.

Por otra parte, los puntos influyentes (tablas 19a y 19b) indican aquellos valores de datos que se encuentran fuera de la tendencia general, son puntos que bien vale la pena de hacer resaltar.

Presentan una diferencia muy significativa los siguientes Estados: para ambos sexos, Hidalgo, Oaxaca y Quintana Roo en los primeros años de vida y en edades

avanzadas y Sinaloa a la primera edad; para el caso del sexo masculino, el Estado de Nuevo León a la edad 0 y el Estado de Zacatecas en los primeros años de vida y en edades avanzadas; para el caso del sexo femenino, los Estados de Baja California, Coahuila, Chiapas, Nayarit y Sonora a edades muy tempranas.

También existen diferencias notables en cuanto a la esperanza de vida masculina en los Estados de Chiapas a la edad 30, Durango a las edades 75 y 80 y Tamaulipas a la edad 75; mientras que en el caso de las mujeres las diferencias se presentan principalmente en Chiapas en edades avanzadas, Nuevo León en las primeras edades y a las edades 40, 45 y 80, Puebla en casi todas las edades, Morelos, Querétaro y Tamaulipas a los 70 y 75 años.

Es aceptado que la población no está constituida por individuos biológicamente homogéneos. De suerte que "la acción de los factores determinantes produce resultados diferentes según sean distintas las características presentes en los individuos afectados, las características biológicas establecen diferencias en el estado de la mortalidad, ya que los factores determinantes inciden de distinta manera sobre los individuos, de acuerdo a su sexo y con su edad". Los factores biológicos fijan límites a los niveles inferiores de mortalidad y los factores sociales aumentan este nivel en relación a la situación del medio externo. En otras palabras, la mortalidad tiene dos componentes primarios: uno endógeno (causas intrínsecas al organismo) y otro exógeno (causas de muerte exterior al organismo). Esta parte puede considerarse controlable por el avance del desarrollo social de las comunidades. Prescindiendo de los factores biológicos que influyen en la mortalidad los cuales se reflejan en la edad, los

diferentes niveles de mortalidad puedan explicarse en función de caracteres socioeconómicos.

Es importante pues conocer las diferencias en el interior de la República en términos de sus distintas regiones geográficas. En cuanto a la esperanza de vida no se cuenta con cifras suficientemente confiables para años recientes de las entidades federativas del país; sin embargo, se ha realizado una clasificación de dichas entidades en cinco categorías, de acuerdo con sus niveles de mortalidad, mediante el cálculo de los efectos región obtenidos en los cuadros y el cálculo de la distancia de Mahalanobis para detectar los puntos influyentes.

Estos grupos han sido conformados como sigue:

Grupo I: Nivel de mortalidad muy bajo:

Quintana Roo.

Grupo II: Nivel de mortalidad bajo:

Baja California Sur, Campeche, Chihuahua, Durango,

Sinaloa y Tamaulipas.

Grupo III: Nivel de mortalidad medio:

Aguascalientes, Baja California, Distrito Federal, Guanajuato,

Gerrero, Jalisco, México, Michoacán, Morelos, Nayarit,

Nuevo León, Querétaro, San Luis Potosí, Tabasco,

Tlaxcala, Veracruz, Yucatán y Zacatecas.

Grupo IV: Nivel de mortalidad alto:

Chiapas, Coahuila, Colima, Puebla y Sonora.

Grupo V: Nivel de mortalidad muy alto:

Hidalgo, Oaxaca.

Se puede observar que las entidades con bajos niveles de mortalidad se encuentran en el noreste y noroeste, con excepción de los Estados de Campeche y Quintana Roo que se encuentran al sureste de la República, en correspondencia con las regiones en donde se observan los más altos niveles socioeconómicos. La mayoría de los estados del tercer grupo, de mortalidad media, se encuentran en el centro del país. Los clasificados en los grupos IV y V, de alta mortalidad, se localizan en el sur del país, y se caracterizan, casi todos, por ser los estados mayormente marginados del desarrollo económico del país.

Hay que mencionar que existen algunas diferencias por sexo dentro de los estados del grupo III. En el caso del sexo femenino, el estado de Nuevo León pasa a formar parte del grupo II, es decir tiene en general una esperanza de vida femenina significativamente mayor que la masculina; y los efectos (figura 1) correspondientes a los estados de Baja California y Querétaro indican que el nivel de mortalidad de las mujeres en estos estados es menor que el nivel general de la población femenina. En el caso del sexo masculino, el estado de Zacatecas pasa al grupo I, los estados de Tabasco y Nayarit pasan al grupo II y los estados de Morelos y Veracruz al grupo IV, con respecto al resto de la población



masculina y el estado de Puebla que se consideró como entidad con un nivel de mortalidad bajo pasa a formar parte dentro del grupo cuyo nivel es muy alto. Finalmente, en el estado de Yucatan es donde se encuentra mayor diferencia por sexo, puesto que en el caso de las mujeres, esta entidad tiende a estar considerada entre los estados cuyo nivel de mortalidad es alto, y en el caso de los hombres éste podría estar considerado en el grupo cuyo nivel de mortalidad es bajo.

Es evidente que pueden quedar ocultos ante el análisis de las cifras globales, algunos contrastes de marginación y bienestar coexistentes en las grandes áreas urbanas que probablemente constituyen las unidades geográficas. Aún cuando en las zonas metropolitanas se está incrementando la mortalidad por algunas causas que son propias de las grandes ciudades (accidentes automovilísticos o de trabajo, violencias, enfermedades mentales, sobretudo en edad activa), éstas muestran una mejor posición con respecto al promedio del país.

Los estados del grupo IV y V son las entidades económica y socialmente menos desarrolladas del país. No es nuestro objetivo resumirlos, sin embargo, se mencionan algunos aspectos que inciden directa y negativamente en la mortalidad y la morbilidad, y por los cuales se piensa que aún siendo altos los niveles de mortalidad para estos estados -medidos a través de indicadores en los que se han utilizado registros de defunciones-, subestiman los valores reales de este fenómeno. Estos aspectos saltan a la vista de cualquier observador que se interna en estos estados.

Tal vez lo primero que llama la atención al recorrer estos estados sea la gran cantidad de pequeñas comunidades esparcidas a lo largo de su abrupto territorio. El sello característico de estas poblaciones lo constituye su aislamiento, pobreza y falta de servicios municipales (falta de agua potable, servicio de alcantarillado insuficiente, presencia de animales en las viviendas).

Por otro lado, la actividad principal de la población económicamente activa se constituye por labores agropecuarias. La falta de recursos económicos de los habitantes de estos estados es notoria, sobre todo en el medio rural.

Los bajos niveles educativos, por su parte, además de relacionarse en forma directa con la productividad, el ingreso, el subempleo y otros, inciden negativamente en las condiciones de la salud, sobretodo en un lugar como el estado de Oaxaca en donde el problema educativo es alarmante. Según datos del Censo de 1970, en Oaxaca, el 49.1% de las mujeres de diez y más años y el 34.6% de los hombres en las mismas edades eran analfabetas, y de la población de seis años y más de edad, el 51.2% no había recibido ni un año de instrucción formal.

Otros aspectos que cabe mencionar son los étnico-culturales y las condiciones de vivienda y alimentación. Entre las costumbres de la población rural que más importancia tienen para la salud, se encuentran los hábitos alimenticios (para Oaxaca, los porcentajes de población que durante la semana anterior al censo declararon no haber consumido carne, huevos, leche, pescada y pan de trigo, fueron respectivamente 24.4, 26.8, 68.3, 64.1 y 25.8%), la forma de atención a

los enfermos y algunas festividades que tienen su origen en la religión.

Las condiciones de la vivienda dependen mucho de la zona geográfica y el clima prevalente. En las zonas montañosas, por ejemplo, donde las temperaturas son bajas, se acostumbra construir con paredes de adobe y con techos de palma, teje o madera, viviendas de un solo cuarto, sin ventanas; en donde vive la familia junto con los animales que poseen, duermen en patates y se cocina sobre el suelo. Estas condiciones crean además del acocimiento, situaciones propicias para la contaminación del aire y la transmisión de enfermedades.

Finalmente, se aborda uno de los aspectos que tienen mayor relación con la mortalidad: los servicios médico-asistenciales. Los índices de habitantes por médico y personas por cama para todo el estado de Hidalgo, por ejemplo, fueron respectivamente 4376 y 1893, en 1980.

Mientras que los estados definidos por los grupos I y II están entre las entidades que poseen los mayores ingresos per capita. Además del Producto Nacional Bruto que durante mucho tiempo era considerado como punto de comparación, se incluyen otras variables que pretenden ilustrar el proceso de industrialización, la modernización agrícola, la capacidad productiva por habitante, la infraestructura socioeconómica y las condiciones sociales de la entidad.

En resumen, se puede decir que entre las entidades que conforman la República, existen diferencias en cuanto a niveles de mortalidad en los distintos grupos de

edad, tanto para hombres como para mujeres. El valor nacional es la resultante de valores muy dispares y por lo tanto es un indicador parcial de la realidad nacional en la cual prevalecen las desigualdades. Esta situación se mantiene a lo largo del tiempo, si bien parece haber una tendencia a concentrarse alrededor del valor nacional con algunas entidades acentuadamente rezagadas.

ESPANOL DE VIDA, POR CADA CATEGORIA, SEGUN ENTORNO PRECISADO  
 SEGUN ANÁLISIS  
 1979

Valor típico: 34.093

Entidad	0	1 año	2 años	3 años	4 años	5 años	6 años	7 años	8 años	9 años	10 años	11 años	12 años	13 años	14 años	15 años	16 años	17 años	18 años	19 años	20 años		
Total	42.96	45.37	42.16	37.21	31.06	40.56	44.25	40.25	35.75	30.76	27.20	23.59	20.22	16.62	13.60	10.56	7.93	5.42					
Agrupaciones:																							
Baja California	43.00	46.32	43.10	36.10	33.31	45.27	44.00	40.52	35.20	30.20	26.10	22.10	18.11	14.10	10.10	7.10	4.10	1.10					
Baja California Sur	42.65	44.10	42.60	37.00	33.87	40.60	44.00	40.00	35.53	31.90	27.80	23.00	18.10	13.50	9.50	6.43	3.00	0.00					
Camacho	46.96	46.20	46.11	40.27	34.52	31.00	40.00	42.11	37.70	33.20	29.27	24.00	19.00	14.23	10.53	6.23	2.04	0.04					
Casullo	43.43	43.91	42.31	37.29	32.91	40.23	42.23	38.00	33.10	28.00	24.20	20.00	15.10	10.10	6.10	2.10	0.10	0.10					
Colima	42.30	44.53	42.22	37.09	32.37	41.23	43.00	39.50	35.27	31.00	26.50	22.00	18.00	13.00	8.00	4.00	0.00	0.00					
Coltaco	40.43	42.91	42.22	36.97	31.36	41.62	42.50	39.62	34.51	30.04	25.43	21.00	16.00	11.00	6.00	1.00	0.00	0.00					
Coltepec	44.12	47.25	44.10	38.25	34.36	40.77	43.62	41.00	37.13	33.00	28.00	24.00	19.00	14.00	9.00	4.00	0.00	0.00					
Alta y Baja California	45.41	47.52	43.24	39.20	34.27	41.17	43.61	40.40	35.71	31.04	27.47	23.70	19.90	15.67	11.56	6.60	2.21	0.21					
Guaymas	44.65	44.00	44.60	38.23	34.13	41.65	44.16	41.30	37.90	34.77	30.20	26.20	21.70	16.00	10.53	5.57	0.57	0.57					
Guaymas	43.17	45.74	42.90	38.21	33.00	41.10	43.10	40.22	35.00	30.00	25.00	20.00	15.00	10.00	5.00	0.00	0.00	0.00					
Guaymas	42.71	44.00	42.94	37.13	32.60	40.35	41.00	40.11	36.27	32.40	28.20	24.00	19.00	14.00	9.00	4.00	0.00	0.00					
Guaymas	40.27	41.25	40.53	34.00	31.24	41.67	40.00	37.00	33.20	29.57	25.24	21.20	16.00	11.00	6.00	1.00	0.00	0.00					
Jalisco	42.34	44.54	42.50	36.23	34.21	40.23	41.26	40.67	36.53	32.43	28.20	24.20	19.20	14.20	9.20	4.20	0.20	0.20					
México	39.90	41.11	42.81	37.43	33.60	40.43	41.10	39.00	35.10	31.00	27.00	23.20	19.00	14.00	9.00	4.00	0.00	0.00					
Micahocán	42.40	45.04	41.00	37.15	32.43	41.00	42.70	39.70	34.20	29.00	24.00	19.00	14.00	9.00	4.00	0.00	0.00	0.00					
Moravia	42.00	44.26	41.00	36.21	31.49	41.21	42.20	39.63	34.30	30.00	25.00	20.00	15.00	10.00	5.00	0.00	0.00	0.00					
Nayarit	40.17	40.00	41.21	36.00	34.14	41.50	41.11	40.20	36.70	32.64	28.50	24.31	20.01	15.00	10.00	5.00	0.00	0.00					
Nuevo León	40.44	42.70	41.21	34.50	34.70	40.54	40.50	40.00	36.27	32.71	27.20	23.00	18.01	13.24	8.20	3.20	0.20	0.20					
Oaxaca	39.00	37.67	40.62	34.64	31.91	41.54	40.45	39.20	34.21	29.63	25.20	20.62	16.00	11.42	6.42	1.42	0.42	0.42					
Panama	37.00	40.72	40.33	33.00	34.27	40.73	40.00	37.00	33.00	29.50	26.00	22.00	18.00	14.00	10.00	6.00	2.00	0.00					
Queretaro	41.52	44.00	42.00	37.07	31.17	40.64	41.30	40.00	36.00	32.00	28.00	24.00	19.00	14.00	9.00	4.00	0.00	0.00					
Queretaro	40.44	40.64	37.40	32.75	31.91	41.50	40.43	39.00	34.20	29.30	24.20	19.20	14.20	9.20	4.20	0.20	0.20	0.20					
San Luis Potosí	40.33	44.00	41.13	37.00	34.57	40.44	41.00	39.00	34.20	29.20	24.20	19.20	14.20	9.20	4.20	0.20	0.20	0.20					
Sonora	40.22	42.10	42.71	36.00	34.00	40.00	41.00	39.10	34.20	29.10	24.10	19.10	14.10	9.10	4.10	0.10	0.10	0.10					
Sonora	40.21	44.20	44.71	37.07	32.21	40.70	40.40	38.50	33.20	28.20	23.00	18.00	13.00	8.00	3.00	0.00	0.00	0.00					
Tampico	41.17	43.42	42.79	36.20	33.21	41.21	41.10	41.00	37.50	33.63	29.00	24.33	19.00	14.20	9.20	4.20	0.20	0.20					
Tampico	40.40	40.00	40.00	36.00	34.20	40.20	40.20	37.00	33.00	29.00	25.00	21.00	17.00	13.00	9.00	5.00	1.00	0.00					
Tampico	41.63	40.00	42.62	36.10	33.23	40.20	40.20	38.10	33.00	28.00	23.00	18.00	13.00	8.00	3.00	0.00	0.00	0.00					
Tampico	40.00	40.27	41.43	36.23	32.43	41.63	41.63	39.20	34.10	29.00	24.00	19.00	14.00	9.00	4.00	0.00	0.00	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50	5.50	1.50	0.00					
Tampico	40.20	40.27	41.00	36.00	34.00	40.20	40.20	37.50	33.50	29.50	25.50	21.50	17.50	13.50	9.50</								

ESTADÍSTICA DE VOTOS POR CADA CANTÓN, SEGUNDO CONTORNO ELECTORAL  
 SECONDO CONTO  
 1978

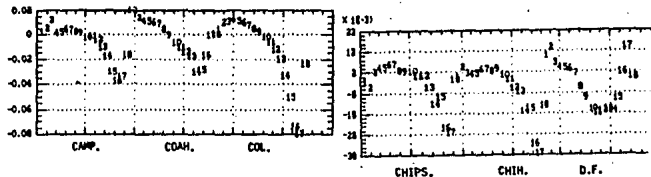
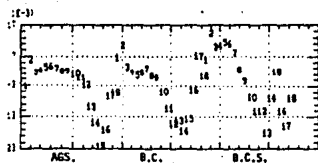
Valor típico: 37.4787

CANTÓN	0	1 año	2 años	3 años	4 años	5 años	6 años	7 años	8 años	9 años	10 años	11 años	12 años	13 años	14 años	15 años	16 años	17 años	18 años
Total	67.67	67.64	66.96	67.29	67.40	67.70	66.70	67.67	69.29	66.96	66.64	66.70	66.45	66.10	66.45	66.60	66.60	66.60	66.60
Aguascaltecos	66.54	66.52	67.52	66.79	66.60	67.29	66.76	66.23	66.70	66.36	66.90	66.60	66.60	66.20	66.10	66.10	66.10	66.10	66.10
Deje California	66.62	66.71	66.70	66.70	67.12	66.90	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06	67.06
Deje California Sur	66.00	67.71	66.50	66.71	66.00	66.00	66.70	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50
Comacho	66.62	66.62	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60
Cachilla	66.59	66.70	67.02	66.71	67.70	67.61	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94	67.94
Colono	67.27	66.76	66.42	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90
Chopan	66.29	65.90	66.00	66.01	66.71	66.47	66.70	67.00	67.70	67.24	67.43	67.43	67.43	67.43	67.43	67.43	67.43	67.43	67.43
Cachucha	66.00	67.51	66.53	67.71	66.99	66.12	66.45	66.70	66.11	65.63	66.11	66.11	66.11	66.11	66.11	66.11	66.11	66.11	66.11
Distrito Federal	66.10	67.36	66.75	67.00	66.90	66.10	66.45	66.21	66.00	65.42	66.90	66.00	66.00	66.00	66.00	66.00	66.00	66.00	66.00
Guaranga	67.12	67.49	67.05	66.20	66.32	66.56	66.66	67.17	66.62	66.67	66.66	66.66	66.66	66.66	66.66	66.66	66.66	66.66	66.66
Guacajala	66.32	66.10	67.50	66.70	66.00	67.27	66.53	66.63	66.63	66.63	66.63	66.63	66.63	66.63	66.63	66.63	66.63	66.63	66.63
Guaymas	66.70	67.46	66.70	66.10	67.41	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71	66.71
Huachuco	67.41	66.70	67.30	67.31	67.17	66.47	66.12	66.96	66.71	66.63	67.02	66.19	66.19	66.19	66.19	66.19	66.19	66.19	66.19
Jaltaco	66.27	66.84	67.23	67.06	66.23	67.42	66.72	66.04	66.64	66.61	66.70	66.66	66.66	66.66	66.66	66.66	66.66	66.66	66.66
México	66.71	66.80	67.06	66.60	66.60	67.12	66.77	66.71	66.64	66.31	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90	66.90
Wigman	66.07	66.70	67.23	66.47	67.04	66.90	66.70	67.70	66.36	66.00	66.00	66.34	66.34	66.34	66.34	66.34	66.34	66.34	66.34
Morales	66.93	66.56	67.04	66.56	67.00	67.10	66.54	66.50	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60
Muyil	66.70	66.10	66.10	66.46	66.50	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49	66.49
San Juan	67.70	67.27	66.57	66.10	66.71	66.41	66.00	66.94	66.50	66.62	66.10	66.57	66.54	66.54	66.54	66.54	66.54	66.54	66.54
Quana	66.12	67.71	66.10	66.70	66.70	67.70	66.60	66.61	66.60	66.61	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60
Popala	67.20	66.50	66.51	66.00	67.00	66.91	66.61	66.30	66.31	66.32	66.74	66.30	66.30	66.30	66.30	66.30	66.30	66.30	66.30
Quetzaco	66.61	66.71	66.66	66.13	66.47	66.16	66.75	66.23	66.53	66.49	66.13	66.13	66.13	66.13	66.13	66.13	66.13	66.13	66.13
Quetzaco Am	67.62	66.51	67.11	66.20	67.20	66.43	67.71	67.10	67.60	66.71	67.00	67.00	67.00	67.00	67.00	67.00	67.00	67.00	67.00
San Luis Potosí	66.30	66.42	66.50	66.10	66.06	66.91	67.49	67.27	66.66	66.67	66.66	66.66	66.66	66.66	66.66	66.66	66.66	66.66	66.66
Sancti	67.10	67.11	66.64	66.12	66.70	66.94	66.00	66.70	66.72	66.23	66.00	66.00	66.00	66.00	66.00	66.00	66.00	66.00	66.00
Sanora	66.69	66.70	67.11	66.34	67.00	66.70	66.60	67.00	66.74	66.75	66.96	66.99	66.99	66.99	66.99	66.99	66.99	66.99	66.99
Tancitaro	66.31	67.62	66.90	66.17	67.70	66.96	66.16	66.19	66.37	66.60	66.19	66.19	66.19	66.19	66.19	66.19	66.19	66.19	66.19
Tancitaro	67.54	67.52	67.49	66.46	67.70	66.92	66.77	66.61	66.70	66.41	66.00	66.77	66.68	66.70	66.70	66.70	66.70	66.70	66.70
Tancitaro	66.69	66.57	67.77	67.00	66.72	66.60	66.63	66.12	66.60	66.62	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60	66.60
Tancitaro	67.00	66.99	66.27	66.00	66.64	66.20	67.02	67.02	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50	66.50
Tancitaro	66.67	66.66	67.05	66.20	67.43	66.71	66.66	67.12	66.30	66.00	66.00	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70
Tancitaro	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70	66.70

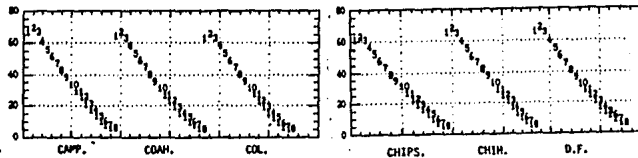
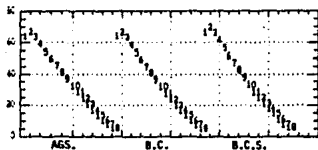
\* Datos Influenciales: Primera Regresión  
 + Datos Influenciales: Segunda Regresión

Tabla 15b

GRAFICA 19a  
RESIDUOS DE LA MEDIANA PULIDA DE (ESPERANZA DE VIDA)<sup>0.25</sup>

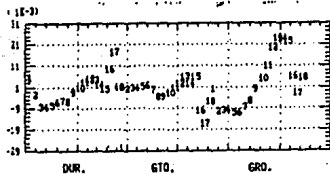


GRAFICA 9a  
DATOS CRUDOS DE LA ESPERANZA DE VIDA - SEXO MASCULINO

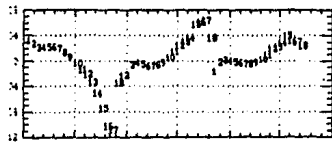


\* Codificación de los puntos:

- |             |              |
|-------------|--------------|
| 1 = edad 0  | 10 = edad 40 |
| 2 = edad 1  | 11 = edad 45 |
| 3 = edad 5  | 12 = edad 50 |
| 4 = edad 10 | 13 = edad 55 |
| 5 = edad 15 | 14 = edad 60 |
| 6 = edad 20 | 15 = edad 65 |
| 7 = edad 25 | 16 = edad 70 |
| 8 = edad 30 | 17 = edad 75 |
| 9 = edad 35 | 18 = edad 80 |

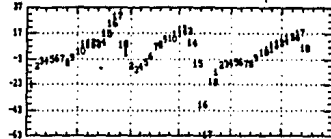






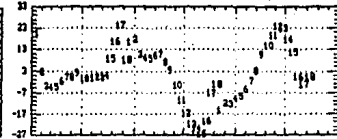
N.L. QAX. PUE.

(X 10-3)



QRO. Q.R. S.L.P.

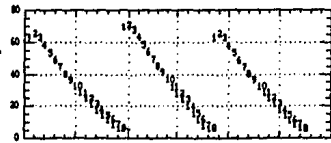
(X 10-3)



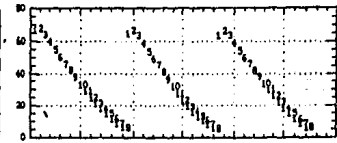
SIN. SOM. TAB.



N.L. QAX. PUE.

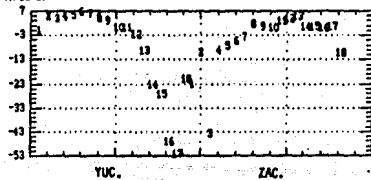


QRO. Q.R. S.L.P.

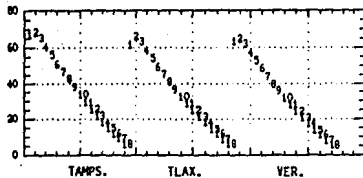
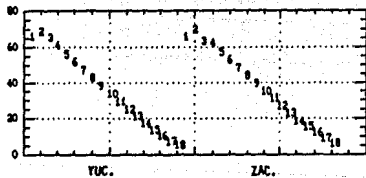
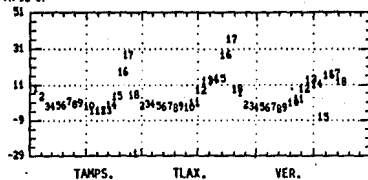


SIN. SOM. TAB.

(X 18-3)

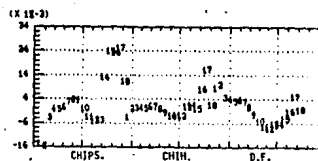
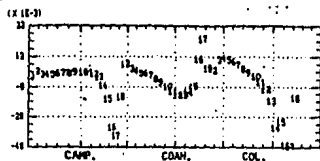
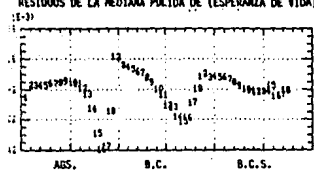


(X 18-3)



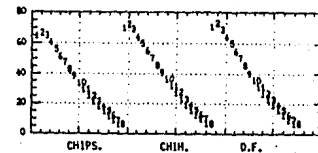
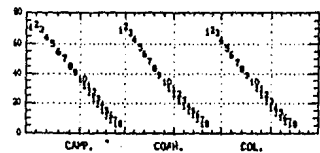
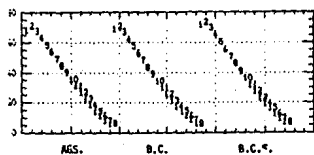
GRAFICA 19b

RESIDUOS DE LA MEDIANA PULIDA DE (ESPERANZA DE VIDA)<sup>0.25</sup>

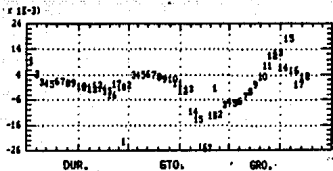


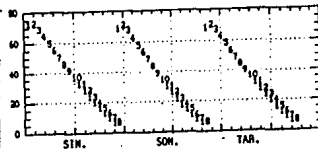
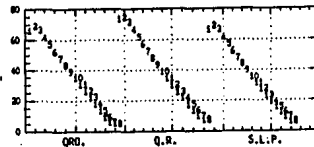
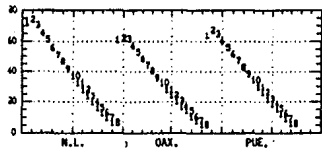
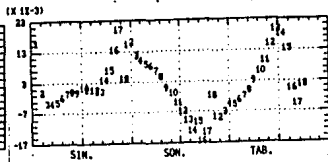
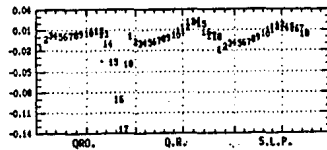
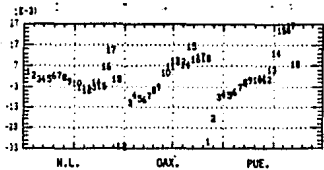
GRAFICA 5b

DATOS CRUOSOS DE LA ESPERANZA DE VIDA - SEXO FEMENINO

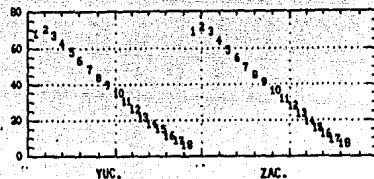
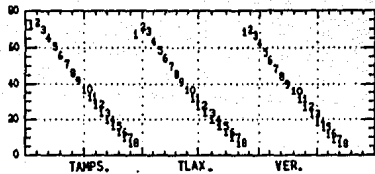
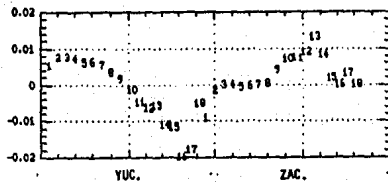
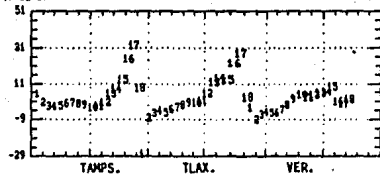


- \* Codificación de los puntos:
- |             |              |
|-------------|--------------|
| 1 = edad 0  | 10 = edad 40 |
| 2 = edad 1  | 11 = edad 45 |
| 3 = edad 5  | 12 = edad 50 |
| 4 = edad 10 | 13 = edad 55 |
| 5 = edad 15 | 14 = edad 60 |
| 6 = edad 20 | 15 = edad 65 |
| 7 = edad 25 | 16 = edad 70 |
| 8 = edad 30 | 17 = edad 75 |
| 9 = edad 35 | 18 = edad 80 |





(X 1E-9)



**CONCLUSION**

**"Es una necesidad todo intento de conclusión".**

**Gustave Flaubert**

## CONCLUSIONES Y RECOMENDACIONES

El análisis exploratorio de datos es un conjunto de técnicas gráficas y numéricas basadas principalmente en el ordenamiento de los datos. La finalidad de este análisis consiste en describir la estructura subyacente a los datos a partir de la ecuación  $\text{DATO} = \text{MODELO} + \text{RESIDUO}$  (donde el modelo es una función generalmente desconocida) y poner de manifiesto la forma en que conjunto de datos se desvía de un modelo (o índice de localización) particular. Es de hacer notar que las técnicas derivadas del análisis exploratorio permiten estudiar a los datos de interés de una manera un tanto informal; es decir, sin hacer suposiciones "a priori" sobre el comportamiento probabilístico de los mismos.

Las consideraciones anteriores permiten concluir que el enfoque exploratorio del análisis de datos está basado en dos principios generales: escepticismo y criterio abierto. Curtis (1986) explica que: "El escepticismo estadístico significa que el investigador debe estar siempre consciente que los resúmenes numéricos (como la media aritmética) o las pruebas confirmatorias (como la prueba de t) representan una forma efectiva de eliminar detalles y variaciones azarosas en los datos. El principio de criterio abierto se refiere a la buena voluntad de buscar suposiciones no razonables o patrones inesperados en el lote de datos."

Por otro lado, el estudio de los casos aberrantes es una parte importante del análisis exploratorio de datos. Curtis y Rodríguez (1988) comentan, a propósito de



dicho estudio, que "no es un tópico nuevo en la literatura estadística (por ejemplo, véase Peirce, 1852); sin embargo, es un tema que ha generado diversas actitudes a lo largo del tiempo. Algunas de estas conductas van, desde no atreverse a dudar de la calidad que un lote de datos pueda tener, hasta aquellas que indiscriminadamente desechan las observaciones que no se ajustan al patrón que el investigador observa o pretende esperar. Afortunadamente, hoy día se tienen suficientes bases teóricas y prácticas que permiten manejar en forma adecuada los casos aberrantes de un lote de datos. El uso de estadísticas resistentes como la mediana es un ejemplo del criterio abierto que un buen analista debe de poseer para identificar y valorar adecuadamente estos casos."

En este sentido, los resultados fueron satisfactorios, puesto que se encontró una gran "consistencia" en el comportamiento de la mortalidad, tanto entre grupos de edad y entre sexos, como en el tiempo y en el espacio: algunos estados y ciertos periodos presentan conductas que se apartan de las normas habituales. Se consideró que éstas puedan, tal vez, deberse a particularidades de dichos estados más que a problemas de información.

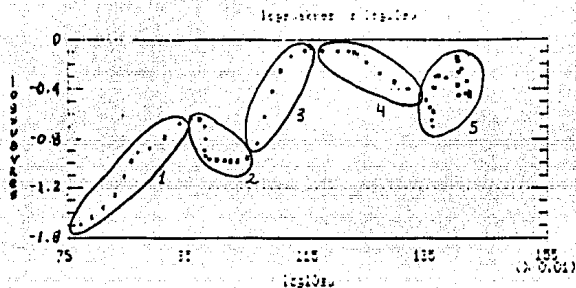
De esta manera, se puede señalar que existen fuertes diferencias en cuanto a niveles de mortalidad en los distintos grupos de edad, tanto para hombres como para mujeres. El valor nacional es la resultante de valores muy dispares y por lo tanto es un indicador parcial de la realidad nacional en la cual prevalecen las desigualdades. Esta situación se mantiene a lo largo del tiempo, si bien parece haber una tendencia a concentrarse alrededor del valor nacional con algunas entidades acentuadamente rezagadas.

Indudablemente que han sido grandes los avances que México ha logrado en el abatimiento de los niveles de la mortalidad; sin embargo, dentro de esta trayectoria general descendente, se observó respecto a los últimos años estudiados una disminución del ritmo de descenso, lo que hace suponer, como muchos autores apuntan, que los avances futuros serán más lentos y probablemente condicionados a sustanciales adelantos en el campo del desarrollo socioeconómico.

Por otro lado, los niveles alcanzados por otros países de mayor desarrollo señalan metas a las que es posible aspirar con las presentes técnicas médicas y con un desarrollo socioeconómico que permita su aplicación a la totalidad de la población.

Así pues, se trató de efectuar una breve interpretación del análisis que tomara en cuenta las dimensiones espacial y temporal. Sin embargo, se pueden explotar más los cuadros y las gráficas aquí presentadas. Por ejemplo, el lector podría preguntarse por qué los valores de  $\log(Z_t)$  versus los valores de  $\log(S_m | R_t)$  de la gráfica 5 (Diagnóstico de la no-estacionaridad de la serie de tiempo) marcan épocas bien definidas. Véase la gráfica siguiente. El primer conjunto de datos corresponde al período 1971-1982, el segundo: 1960-1970, el tercero: 1954-1959, el cuarto: 1942-1953 y el quinto: 1941-1922. Ésta y otras observaciones puede ser objeto de un futuro análisis.

Gráfica 5  
Gráfica del Diagnóstico



En síntesis, Curtis y Rodríguez (1988), explican que "el análisis exploratorio de datos es análogo al trabajo que efectúa un detective: se desea poner al descubierto lo que los datos tratan de decir". Y se puede agregar que tanto la pictografía de datos como los residuos son para el analista de datos, "lo que una lente de aumento, una prueba química para manchas de sangre, o un dispositivo auditivo representan para un novelista de historias detectivescas".

¡INTERES ante!  
¡Muy INTERES ante!



## BIBLIOGRAFIA

"Todo hombre debería de leer sólo  
aquello a que le lleva su inclinación;  
ya que lo que lee como obligación  
de poco le aprovechará."

Samuel Johnson

## BIBLIOGRAFIA

- Almada B., I. (1980).  
La Mortalidad en México. 1982-1975.  
I.B.S.S.
- Berenson, M., Levine, D. & Goldstein, M. (1983).  
Intermediate Statistical Methods and Applications.  
A Computer Package Approach.  
Prentice Hall.
- Breckenridge, M. B. (1983).  
Age, Time and Fertility.  
Applications of Exploratory Data Analysis.  
New-York: Academic Press.
- Comasortega C, S. & Juez C, R. (1977).  
Descripción y Análisis de la Mortalidad en México 1900-1973.  
Tesis. Facultad de Ciencias. U.N.A.M.
- Centro de Estudios Económicos y Demográficos. (1981).  
Dinámica de la Población de México.  
El Colegio de México.
- Chambers, J. M., Cleveland, W. S., Kline, B., & Tukey, P. R. (1983).  
Graphical methods for data analysis.  
Belmont, CA: Wadsworth International Group.
- Consejo Nacional de Población. (1984).  
Reunión Nacional Sobre Mortalidad y Políticas de Salud.
- Consejo Nacional de Población. (1984).  
Estado Actual del Conocimiento sobre los Niveles y Tendencias de la Mortalidad en México.
- Curts, J. B. (1984).  
Introducción al análisis de residuos en biología.  
Biótica, 9 (3), 271-278.
- Curts, J. B. (1985).  
Teaching College Biology Students The simple Linear Regression Model Using an Interactive Microcomputer Graphics Software Package.  
Dissertation Abstracts International Vol 46, Number 7, 1986
- Curts, J. B. (1986).  
Regresión lineal resistente en biología.  
Biótica.
- Curts, J. B., Alcantara, L. & Chiappa, X., (1987).  
Introducción al Análisis Exploratorio de Datos Multidimensionales.  
CIENCIAS, Revista de Difusión Nº. 11. U.N.A.M.

- Curts, J. B. & Silva, A. (1988).  
Métodos Cuantitativos en Psicología. (Capítulo 10).  
Editorial Trillas.
- González C., P. & Florescano, E. (1979).  
México, Hoy.  
Siglo XXI Editoras.
- Han Chanda, A. (1983).  
John Grant: Un origen de la Estadística, la Demografía,  
... y algo de la Actuaría.  
ACTUA. Enlace Informativo. Núm. 3.  
Colegio Nacional De Actuarios.
- Hartwig, F., & Dearing, B., (1979).  
Exploratory Data Analysis.  
Sage Publications. Beverly Hills. London.
- Hoaglin, D. C., Mosteller, F., & Tukey, J. W. (1983).  
Understanding robust and exploratory data analysis.  
New-York: John Wiley & Sons, Inc.
- Johnston, J., (1972).  
Econometric Methods.  
McGraw-Hill Book Company, New-York.
- Leinhardt, S.; Emerson, M.; Stoto, M. & Valleman, P. (1984).  
Sociological Methodology.  
Jossey-Bass. San Francisco.
- Martínez M., J. (1982).  
La revolución demográfica en México. 1970-1980.  
I.M.S.S.
- Ornelas, R. & Minujin Z., A. (1984).  
Los Factores del cambio demográfico en México.  
Instituto de Investigaciones Sociales. U.N.A.M.  
Siglo XXI.
- Pressat, Roland. (1970).  
El análisis demográfico.  
Métodos, resultados, aplicaciones.  
México: Fondo de Cultura Económica.
- Rodríguez, D. (1987).  
Temas Demográficos.  
Escuela Nacional de Trabajo Social.  
Instituto de Investigaciones Económicas.  
U.N.A.M.

Secretaría de Programación y Presupuesto. (1985).  
Estadísticos Vitales. 1901-1982.

Siegel, A. (1988).  
Statistics and Data Analysis.  
An Introduction.  
John Wiley & Sons.

Spiegelman, M. (1955).  
Introducción a la Demografía.  
Fondo de Cultura Económica. Ed. 1985.

Statistical Graphics Corporation. (1985-1986).  
STATGRAPHICS. Statistical Software Manual. Version 2.1

Tukey, J. W. (1977).  
Exploratory data analysis.  
Reading, MA: Addison-Wesley Publishing Company.

Unidad Académica de los Ciclos Profesional y  
de Posgrado del Colegio de Ciencias y Humanidades.  
Primer Foro Del Proyecto Académico  
Especialización en Estadística Aplicada. Sept. 1986.  
U.N.A.M.

Venables, P. F. & Ripley, D. C. (1981).  
Applications, basics and computing of exploratory data analysis.  
Boston, MA: Dux bury Press.

Wainer, H. & Thissen, D. (1981).  
Graphical data analysis.  
Annual Review of Psychology, 32, 191-241.

Weisberg, S. (1980).  
Applied Regression Analysis.  
John Wiley & Sons. New-York.