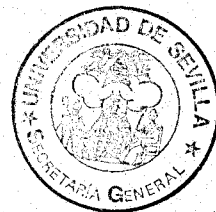


UNIVERSIDAD DE SEVILLA

Depositado en
de la
de esta Universidad desde el día
hasta el día

Sevilla de de 19
EL DIRECTOR DE



UNIVERSIDAD DE SEVILLA
FACULTAD DE MATEMATICAS

UNIVERSIDAD DE SEVILLA
FACULTAD DE MATEMATICAS
BIBLIOTECA

ANALISIS CUALITATIVO

DE DATOS ESTADISTICOS

UNIVERSIDAD DE SEVILLA
SECRETARIA GENERAL

Queda registrada esta Tesis Doctoral
al folio 133 número 23 del libro
correspondiente.

Sevilla, 22 MAYO 1967

El Jefe del Negociado de Tesis,

H. Y. J. de la Cruz Roldán

Memoria dirigida por:

Prof. Dr. D. Joaquin Muñoz Garcia

Prof. Dr. D. Antonio Pascual Acosta

Memoria presentada por:

Juan Luis Moreno Rebollo

Admito la envuelta de este
memoria

Juan Luis Moreno Rebollo

R 11192

LBS 602217

043

172

UNIVERSIDAD DE SEVILLA

FACULTAD DE

MATEMATICAS

ANALISIS CUALITATIVO

DE DATOS ESTADISTICOS

Visado en Sevilla a

21 de Mayo de 1987

Fdo. Prof. Dr. D. Joaquin Muñoz Garcia

Fdo. Prof. Dr. D. Antonio Pascual Acosta

Memoria que presenta

Juan L. Moreno Rebollo

para optar al grado de

Doctor en Ciencias

Matemáticas

Fdo. Juan Luis Moreno

Rebollo

INDICE

CAPITULO I .- OUTLIERS. PROBLEMATICA DE LAS TECNICAS DE IDENTIFICACION .	
-- Introduccion	3
1.1 Errores en las Observaciones Expe- rimentales	5
1.2 Outliers	8
1.3 Metodos de Analisis de Outliers ...	13
1.4 Metodos de Identificacion	17
1.5 Problematica de los Metodos de Identificacion	21

CAPITULO II .- APROXIMACION FORMAL AL ANALISIS CUALITATIVO DE DATOS.

-- Introduccion	32
2.1 Propiedades de la relacion criterio .	33
2.2 Funciones de C-similitud. Seleccion .	37
2.3 Invarianza	69

CAPITULO III .- ANALISIS CUALITATIVO DE DATOS .
TECNICAS DE IDENTIFICACION DE
OUTLIERS.

-- Introduccion	77
3.1 Estructura Basica. Aproximacion General	78
3.2 Criterio de Dispersion Central	83
3.3 Criterio de Dispersion Basado en el Recorrido	88
3.4 Criterio Natural (Origen Conocido).	93
3.5 Generalizacion al caso de mas de una Observacion	95
3.6 Criterio de Dispersion Central Multivariante	104

BIBLIOGRAFIA	109
--------------------	-----

INTRODUCCION

Hacia la segunda mitad del siglo XVIII , ciertos astrónomos aconsejan la eliminación de las observaciones extremas (valores máximo y mínimo) resultantes de sus cuantificaciones experimentales, argumentando para ello, que estas observaciones pueden perturbar las inferencias a realizar.

La idea básica de esta práctica experimental, puede considerarse como una primera aproximación a un problema que posteriormente ha adquirido una gran importancia dentro de la estadística, como es, el análisis de la calidad de las observaciones experimentales, lo cual es básico para garantizar la máxima fiabilidad de las conclusiones.

Sin embargo, el procedimiento empleado para la depuración de los datos, antes referido, no es el más correcto, pues se basa únicamente en la subjetividad del experimentador. Por esta razón, diversos autores se dedican a la búsqueda de procedimientos que corrijan en cierta medida la falta de objetividad, y no es hasta la primera mitad del siglo XX cuando aparecen métodos con cierto fundamento teórico. A partir de entonces, surgen numerosas técnicas, que presentan, entre otros, el inconveniente de carecer de una base formal común, y que dan lugar a lo que hoy en día se conoce como "Análisis Estadísticos de Outliers".

En el Capítulo I, se introduce un esquema que permite analizar de un modo genérico los distintos errores que pueden afectar a las observaciones experimentales. Asimismo, se destacan los rasgos e inconvenientes principales que presentan las distintas definiciones que sobre el término outlier han dado diversos autores, y se propone una nueva definición para el mismo. Por último, se tipifica la problemática que presentan las Técnicas de Identificación de Outliers.

Con el objeto de realizar un análisis cualitativo de los datos experimentales, en el Capítulo II, se propone un marco general -basado en la idea de que "el comportamiento anómalo de una observación depende en gran medida

de su relacion con las restantes observaciones obtenidas bajo condiciones similares y del objetivo que se persigue en la experimentacion"- , que permite obtener una funcion que cuantifica las desviaciones de las observaciones entre si, respecto de un criterio dependiente del analisis que se desea realizar con los datos.

El esquema propuesto en el Capitulo II, aplicado al caso particular del Analisis de Outliers, presenta, como se recoge en las conclusiones del Capitulo III, ciertas ventajas respecto del enfoque clasico, entre las que cabe destacar, que algunos estadisticos en los que se fundamentan las Tecnicas de Identificacion de Outliers, se obtienen a partir de una base formal comun.

UNIVERSIDAD DE VALPARAISO
FACULTAD DE MATEMATICAS
BIBLIOTECA

CAPITULO I

OUTLIERS .

PROBLEMATICA DE LAS TECNICAS DE IDENTIFICACION .

--- Introduccion .

1.1 Errores en la Observaciones Experimentales .

1.1.1 Variabilidad de la Fuente o Naturaleza .

1.1.2 Error del Medio .

1.1.3 Error del Experimentador .

1.2 Outliers .

1.3 Metodos de Analisis de Outliers .

1.3.1 Metodos de Acomodacion .

1.3.2 Metodos de Identificacion .

1.4 Metodos de Identificacion .

1.4.1 Metodos o Procedimientos Formales .

1.4.2 Metodos o Procedimientos no Formales o Graficos .

1.5 Problematica de las Tecnicas de Identificacion .

1.5.1 - De los Metodos Formales .

1.5.2 - De los Metodos no Formales o Graficos .

1.5.3 Efecto Enmascaramiento.

INTRODUCCION

El estudio de los fenomenos naturales lleva consigo un proceso de codificacion de los elementos de la naturaleza.

El conjunto de codigos resultantes da lugar a lo que se conoce como masa de datos, a partir de la cual se pretende obtener conclusiones sobre el fenomeno que se analiza. Sin embargo, los datos experimentales, no deben considerarse en ningun caso , como valores absolutamente fidedignos del experimento en cuestion, ya que pueden estar afectados por diversos errores o variabilidades. Por tanto, una de las principales etapas de cualquier analisis estadisco de datos, es la de investigar la calidad de las observaciones, ya que si no se detectan los

errores existentes en las mismas, los procedimientos estadísticos que se apliquen posteriormente pueden conducir a conclusiones inexactas.

Para abordar este problema, en este capítulo, se introduce un esquema que permite analizar, de un modo genérico, los diversos errores que pueden presentarse en las observaciones resultantes de una experimentación.

Los errores que se definen en el epígrafe 1.1, dan lugar a diferentes tipos de observaciones erróneas, sin embargo, se concluye que desde un punto de vista práctico, al menos inicialmente, resulta imposible tal diferenciación, por lo que se las denomina, globalmente, con el término anglosajón OUTLIER.

Se recogen los rasgos principales que caracterizan las definiciones que de outlier han dado diversos autores, indicándose los inconvenientes que presentan, por lo que se propone una nueva definición para este término.

Por último, se realiza un estudio general de las distintas técnicas existentes para abordar el problema de la presencia de outliers en una masa de datos experimentales, analizando con especial interés las Técnicas de Identificación, y la problemática actualmente asociada a las mismas.

1.1 ERRORES EN LAS OBSERVACIONES EXPERIMENTALES

La estadística se puede considerar como aquella ciencia que permite extraer información relevante y útil de un conjunto de datos experimentales.

Al analizar esta definición, se han de tener presentes las afirmaciones que hacen diversos autores, entre los que se pueden citar Anscombe (1960), Beckman-Cook (1983), en el sentido de que las observaciones resultantes de una experimentación no deben considerarse en ningún caso como valores absolutamente fidedignos del experimento en cuestión, ya que estas pueden estar afectadas por diversos errores o variabilidades. Por esta razón existen autores como Anscombe (1960), Hampel (1968), Grubbs (1969), Hawkins (1980), Beckman y Cook (1983), Barnett y Lewis (1984), Hampel et al (1986) que se preocupan del análisis de las diferentes fuentes de variabilidad o error que pueden influir en los valores muestrales. Sin embargo, cabe resaltar que entre estos análisis, unos son bastante incompletos y poco formales en sus consideraciones y otros exponen situaciones demasiado particulares.

Esto motiva que se introduzca un esquema que permita analizar de un modo más genérico y preciso, los posibles errores que se pueden presentar en una experimentación y la forma en que estos influyen en las observaciones de que dispone el estadístico.

La realizacion de un experimento lleva consigo una triplete de elementos, que se representara por (Ω, X, E) , donde Ω es la fuente o naturaleza que provee los elementos $w_i, i=1, \dots, n$ y X es el medio por el cual el experimentador, E , cuantifica o transforma los elementos $w_i, i=1, \dots, n$ en las observaciones $x_i, i=1, \dots, n$ de las que dispondra el estadístico para analizar el experimento.

Segun este esquema, es posible afirmar, que las variabilidades o errores que pueden influir en las observaciones se pueden clasificar en:

- Variabilidad de la fuente o naturaleza
- Error del medio
- Error del experimentador

1.1.1 VARIABILIDAD DE LA FUENTE O NATURALEZA

Es la que se manifiesta en las observaciones al recorrer las mismas el soporte de la poblacion bajo estudio. Tal variacion ha de ser considerada como un comportamiento natural de los datos.

1.1.2 ERROR DEL MEDIO.

Se produce cuando no se dispone de la tecnica adecuada, o cuando no existe un procedimiento que permita realizar de forma exacta la transformacion del espacio empirico al espacio matematico.

En este error se incluye , por ejemplo, el ocasionado por el redondeo forzoso que se ha de realizar cuando se trabaja con variables continuas.

1.1.3 ERROR DEL EXPERIMENTADOR

Este error es atribuible al propio experimentador y se puede presentar de distintas formas:

- Error de informacion: Se comete al suponer un modelo matematico no adecuado o preciso para la poblacion, o al admitir informacion o hipotesis iniciales erroneas.
- Error de planificacion: Por el cual el experimentador no llega a delimitar de forma correcta el espacio Ω y realiza el experimento sobre un espacio distinto $\Omega \cup \Omega'$.

En este tipo de error se contempla el ocasionado al tomar muestras sesgadas o al incluir individuos no pertenecientes a la poblacion bajo estudio.

- Error de realizacion: Se comete al llevar a cabo valoraciones erroneas de las transformaciones de los elementos de .

En este tipo de error se incluyen las transcripciones erroneas de datos , falsas lecturas, etc..

1.2 OUTLIERS

Las variabilidades o errores definidos en el epigrafe 1.1 se suelen presentar en los datos resultantes de una experimentacion por lo que en toda muestra cabe distinguir dos tipos de observaciones ,segun esten afectadas o no por los mismos . A las primeras , que constituyen el objeto de nuestro estudio ,se les ha denominado de muy diversas formas en la bibliografia existente sobre el tema.

Definiendolas segun el tipo de variabilidad o error que presentan, se adopta la siguiente terminologia:

Observacion atipica: Aquel valor $x \in X$ que presenta una gran variabilidad de tipo inherente.

Observacion Erronea: Aquella observacion que presenta un gran error de medio y/o un gran error del experimentador,el cual puede manifestarse a traves de uno o varios de los errores definidos anteriormente.

La posible presencia de estas observaciones en una experimentacion , hace que algunos autores , entre ellos Chatfield (1985), destaquen,entre las principales etapas que se han de seguir en un analisis estadistico de datos, la de investigar la estructura y calidad de las observa-

vaciones , así como la de realizar un examen inicial de las mismas , ya que si no se detectan los errores presentes en las observaciones y el estadístico actúa suponiendo la ausencia de ellos , los procedimientos estadísticos que se apliquen pueden conducir a conclusiones inexactas.

En el análisis a realizar sobre la calidad de los datos , sería deseable que se distinguiesen las observaciones atípicas de las erróneas , ya que las primeras carecerían de error . Sin embargo , desde un punto de vista práctico tal distinción es imposible ya que como muy bien dice Anscombe (1960) , esto es materia de conjetura.

Por tanto , aunque desde un punto de vista teórico se pueda admitir la diferencia entre los distintos tipos de observaciones definidas anteriormente , e incluso se podrían definir muchos más , según los errores citados , en la práctica solo se considerará de un único tipo de observaciones , para la que se adoptará el término anglosajón de OUTLIER.

Antes de definir observación outlier , se destacan y discuten los rasgos que comúnmente se presentan en las definiciones más utilizadas sobre la misma.

Generalmente en las definiciones dadas se caracteriza el outlier mediante su desviación del resto de las

observaciones muestrales, Miller (1981) , y en ocasiones aparecen además algunas de las siguientes peculiaridades:

a) Intervienen los diversos tipos de error que pueden afectar a las observaciones . En este sentido se pueden citar las dadas por Grubbs (1950) , Kendall y Buckland (1957) , Anscombe (1960) , Hawkins (1980), Beckman y Cook (1983).

b) Admiten la subjetividad del experimentador para identificarlas . Entre estas se encuentran las de Grubbs (1950) , Kendal y Buckland (1957) , Ferguson (1961), Grubbs (1969), Hawkins (1980), Beckman y Cook (1983), Barnett y Lewis (1984).

La idea de caracterizar al outlier mediante su desviación del resto de los elementos que componen la muestra parece evidente, ya que , como indican Beckman y Cook (1983):

"... la fiabilidad probable de una observación se refleja mediante su relación con otras observaciones que se obtienen bajo condiciones similares..."

Con respecto a la primera de las peculiaridades , parece necesario que los errores aparezcan en la definición, de forma implícita o explícita , sin embargo , una definición formal de este concepto no debería contener referencia alguna a la influencia subjetiva del experi-

mentador , ya que la presencia de terminos como "despertar sospechas" , "parecer inconsistente" ,... se han de considerar como una reminiscencia de los albores de este problema . Basta citar , a modo de ejemplo , a Boscovich (1755) , quien para determinar la elipticidad de la tierra rechazaba las dos observaciones extremas resultantes de la experimentacion , sin criterio alguno , unicamente por ser la menor y la mayor de las observaciones.

El problema de la subjetividad se pone de manifiesto en trabajos , como el Collet y Lewis (1976) , donde se afirma:

"... existen comportamientos diferentes entre los individuos en su reaccion a valores sorprendentes e incluso entre los juicios realizados por un mismo individuo en distintas ocasiones...".

Es mas, en las "situaciones mas estructuradas", segun la clasificacion realizada por Barnett y Lewis (1984), o cuando se observan datos multivariantes , resulta en general imposible la influencia subjetiva del experimentador sobre la apreciacion de observaciones outliers.

Esto hace que la definicion que se adopte sobre este termino evite cualquier grado de subjetividad. De hecho, existen autores que ponen en entredicho la influencia del experimentador . Asi Beckman y Cook (1983) muestran

su preocupacion hacia definiciones que dependan de una nocion tan subjetiva como "parecer sorprendente".

Por ultimo, la definicion tambien deberia reflejar el hecho de que una observacion sera outlier para un determinado modelo, metodo o criterio , baste citar para ello lo apuntado por Gnanadesikan y Kettering (1972):

"... la complejidad del caso multivariante sugiere que seria conveniente investigar procedimientos que ofrezcan proteccion contra cualquier tipo de outlier , aunque una aproximacion mas razonable seria la de construir procedimientos que permitan proteger los datos ante situaciones especificas..."

Ademas, un outlier para modelo , metodo o criterio no lo es necesariamente para otro.

Asi pues , teniendo en cuenta las consideraciones realizadas se propone la siguiente definicion para el termino outlier:

OUTLIER es aquella observacion que siendo atipica y/o erronea , se desvia marcadamente del comportamiento general de la masa de datos , respecto del criterio (modelo, metodo o situacion) que se desea analizar en los mismos.

Con el termino anglosajon de INLIERS se denominaran a las observaciones no caracterizadas como outliers.

De lo anteriormente expuesto se deduce que es necesario estudiar la presencia de outliers dentro de la muestra, ya que estas observaciones pueden distorsionar, de forma considerable , estimaciones , tests u otro tipo de analisis estadistico que se desee realizar sobre los datos ; incluso, existen bastantes situaciones, donde la caracterizacion de estos es mas importante que los errores que ocasionan.

A continuacion se estudian las diferentes formas de abordar la posible presencia de estas observaciones en una masa de datos.

1.3 METODOS DE ANALISIS DE OUTLIERS.

Los diversos metodos para analizar la presencia de outliers en un conjunto de datos son clasificados por Beckman y Cook (1983) en:

- Metodos de Acomodacion.
- Metodos de Identificacion.

La utilizacion de un metodo u otro dependera , del objetivo o de la informacion que se desee obtener de la muestra resultante de la experimentacion.

1.3.1 METODOS DE ACOMODACION.

Dentro de estos metodos se engloban dos formas distintas de abordar el problema:

- a) Dotar a la fuente o naturaleza de una estructura o modelo matematico que explique todos los valores muestrales.

Este tipo de procedimiento puede dar lugar a la determinacion de una estructura o modelo matematico erroneo , en caso de que las observaciones outliers fuesen producidas por efecto de la variabilidad inherente.

- b) Utilizar procedimientos estadisticos que no sean influenciados por la presencia de outliers (Metodos Robustos).

Desde esta perspectiva se presupone que todos los outliers son observaciones erroneas, sin admitir la posibilidad de que existan observaciones con cierto grado de variabilidad inherente, las cuales desde un punto de vista teorico han de ser consideradas.

Asi pues, puede considerarse que la mayor problematica que plantean estos metodos esta contenida en la conjetura de Anscombe (1960), citada previamente.

Ademas estos procedimientos presentan el inconveniente de que requieren de una gran informacion para poder

cumplir con sus objetivos , informacion generalmente no disponible.

1.3.2 METODOS DE IDENTIFICACION.

Dentro de estos metodos , se engloban todos aquellos procedimientos que se caracterizan por las dos etapas siguientes:

I. Determinacion de la observacion u observaciones outliers.

II. Analisis pormenorizado de cada outlier , y en conjunto de todos ellos , con el objetivo de determinar las variabilidades o errores presentes , lo que permitiria llegar a una o varias de las siguientes conclusiones:

a) Mantener o rechazar las observaciones detectadas como miembros del conjunto inicial de datos.

b) Confirmar el modelo matematico caso que este hubiese sido supuesto inicialmente o determinar las desviaciones que se presenten sobre las hipotesis realizadas sobre la poblacion.

c) Identificar rasgos de interes practico que en un principio son desconocidos por el investigador.

d) Determinar comportamientos singulares en la experimentación.

Las dos últimas situaciones originan en muchos casos que los outliers tengan más importancia que la muestra en sí.

No obstante, se ha de reconocer, que a veces se tienen experimentos en los que no se llega a conclusiones como las expuestas, bien por la conjetura de Anscombe (1960), o bien porque no es posible el análisis de los outliers, lo cual puede ocurrir por muy diversas razones, por ejemplo, en aquellos experimentos donde los elementos del espacio empírico no vuelven a su estado inicial tras la experimentación. En estos casos, se tiende a rechazar las observaciones determinadas en la Etapa I, fundamentándose para ello en las afirmaciones realizadas al respecto por diversos autores, entre ellos Hampel et al (1986):

"... el peligro causado por el outlier es mucho mayor que el peligro de pérdida de eficiencia que motiva su rechazo en caso de que sea inlier...".

Esta argumentación también sería válida en la situación descrita en segundo lugar dentro de los Métodos de Acomodación.

El objetivo fundamental del estudio es dar un procedimiento de tipo general para la construcción de técnicas de identificación de outliers, por lo que a continuación se estudian las diferentes técnicas que se encuadran dentro de los Métodos de Identificación de outliers, y los principales inconvenientes que presentan estos métodos actualmente.

1.4 METODOS DE IDENTIFICACION.

En estas técnicas se han distinguido dos etapas diferentes, sin embargo, se analizará únicamente la primera de ellas, aquella en la que se determinan las observaciones outliers, ya que la segunda es más apropiada para el experimentador que para el estadístico.

Los procedimientos existentes para determinar observaciones outliers se pueden clasificar en dos grandes grupos:

- Métodos o Procedimientos Formales
- Métodos o Procedimientos no Formales o Gráficos

1.4.1 METODOS O PROCEDIMIENTOS FORMALES.

Estos métodos se caracterizan fundamentalmente, por la realización de un contraste de hipótesis, donde la hipótesis nula o hipótesis de trabajo supone que la muestra obtenida de la experimentación no contiene

outlier alguno . Para aceptar o rechazar esta hipotesis se construye una funcion test , que permite , en base al nivel del mismo , decidir sobre la presencia o no de outliers en los datos.

No obstante , para una misma hipotesis de trabajo , suelen presentarse situaciones en las que es necesario elegir entre diversas funciones test , por lo que , al igual que en la teoria clasica de contrastes de hipotesis , es necesario introducir algunas medidas que permitan elegir entre dichas funciones . Asi , se define la potencia del test, de forma similar a como se hace en la teoria clasica de Neymann - Pearson , y por tanto a cada hipotesis de trabajo se le ha de asociar una hipotesis alternativa que fija el modelo que provoca la presencia de outliers en la muestra.

Barnett y Lewis (1984) , estudian diferentes tipos de hipotesis alternativas en caso de que la hipotesis nula fije la pertenencia de las observaciones a un modelo distribucional.

En estos contrastes de hipotesis ,ademas de la potencia , aparecen otras medidas para la comparacion de las funciones test . Asi David y Paulson (1965) definen diversos criterios para el caso de una observacion outlier , los cuales han sido analizados por diversos autores , entre ellos Hawkins (1978,1980) , quien indica la importancia y utilidad de algunos de ellos .

Beckman y Cook (1983) , sugieren los siguientes parametros , para la comparacion de tests para la identificacion de outliers , que son validos para el caso en que exista mas de una de estas observaciones:

Sea Z el numero de outliers correctamente identificados , sea Y el numero de no outliers declarados como tales (numero de falsos positivos) , y sea k el numero de outliers existentes en la muestra.

- i.) La probabilidad de que no existan falsos positivos y todos los outliers sean detectados :

$$P(Z = k ; Y = 0)$$

- ii.) La probabilidad de ningun falso positivo y como minimo un outlier:

$$P(Z > 0 ; Y = 0)$$

- iii.) La proporcion esperada de outliers detectados:

$$E(Z/n)$$

donde n es el numero de elementos que componen la muestra.

- iv.) La probabilidad de al menos un falso positivo :

$$P(Y > 0)$$

- v.) El numero esperado de falsos positivos:

$$E(Y)$$

- vi.) La proporcion entre el valor esperado de falsos positivos y el valor esperado de todas las observa-

ciones declaradas outliers:

$$E(Y)/(E(Z)+E(Y))$$

Pero al igual que ocurre en la teoría clásica de los tests, no siempre es posible encontrar funciones con propiedades óptimas globales sobre los parámetros definidos, por esta razón se suelen buscar tests con condiciones de optimalidad local o insesgadez o que satisfagan ciertas condiciones de invarianza o bien se utilizan procedimientos que tengan buenas propiedades, como ocurre con el principio de razón de verosimilitud.

Los Métodos Formales así expuestos, presentan el inconveniente de que el rechazo de la hipótesis nula no determina que observación u observaciones son outliers, ya que la hipótesis alternativa únicamente formula el modelo que seguirían estas observaciones, caso de existir. Esto hace que se introduzcan las denominadas Hipótesis o Modelos Etiquetados y las Hipótesis o Modelos no Etiquetados.

HIPOTESIS O MODELOS ETIQUETADOS.

Son aquellos en que la hipótesis alternativa especifica las observaciones que se considerarían outliers en caso de rechazar la hipótesis nula.

HIPOTESIS O MODELOS NO ETIQUETADOS.

Son los que no expresan de forma explícita las observaciones que se admitirán como outliers en caso de rechazar la hipótesis de trabajo. En estos se suele indicar el número de outliers a detectar, o una cota superior del mismo. Para la determinación de esta cota se utilizan, entre otros, los denominados procedimientos de saltos.

1.4.2 METODOS O PROCEDIMIENTOS NO FORMALES O GRAFICOS .

Estos métodos tienen como objetivo la representación de los datos en un sistema de ejes cartesianos bi - o - tridimensional, utilizando procedimientos muy diversos. El más usual consiste en transformar cada observación mediante una función con valores en la recta real.

La función o estadístico se elige de acuerdo con el rasgo o propiedad que se desee analizar en los datos experimentales, considerándose outliers aquellas observaciones cuya representación gráfica se desvía relevantemente de la masa de datos.

1.5 PROBLEMATICA DE LOS METODOS DE IDENTIFICACION.

En este apartado se analizan los inconvenientes, teóricos y prácticos, que presentan los dos métodos englobados en este tipo de técnicas.

1.5.1 -DE LOS METODOS FORMALES.

En estos se presentan muy diversos problemas , que se analizaran atendiendo a los elementos caracteristicos de los mismos.

- RESPECTO DE LA HIPOTESIS ALTERNATIVA . Como ya se ha indicado , esta hipotesis fija el modelo que explica los outliers en caso de que estos existan.

Esto resulta muy dificil por las razones que se exponen a continuacion:

1.) Porque no todos los outliers han de proceder de un mismo modelo , ya que esto seria equivalente a afirmar que todos ellos presentan el mismo tipo de error.

Esta problematica ha sido abordada por diversos autores , para el caso particular de modelos poblacionales univariantes , entre ellos Kale (1976), Rosado (1984), los cuales contrastan la hipotesis nula frente a una hipotesis alternativa en la que recogen distintos modelos , todos ellos con la misma forma funcional , aunque de distintos parametros.

Rosado (1984) , reconoce que resulta dificultoso realizar generalizaciones en todos los sentidos posibles.

En el caso de modelos poblacionales multivariantes , o de situaciones mas estructuradas, la problematica que se plantea es aun mayor ya que en el primero de estos casos se han estudiado pocos modelos alternativos y en situaciones mas estructuradas , a la diversidad de modelos existentes , se les une el conjunto de condiciones iniciales que lleva consigo cada uno de ellos.

- ii.) Aun en el caso en que todos los outliers pudiesen explicarse por un unico modelo , cabria preguntarse por la unicidad de este.
- iii.) La incertidumbre existente sobre la eleccion de la hipotesis alternativa reduce la utilidad de cualquier propiedad de optimalidad para los tests , ya que estas dependen de la forma que se adopte para dicha hipotesis.

Estos problemas hacen que la mayor parte de las tecnicas formales esten referidas unicamente a la hipotesis nula.

- RESPECTO A LA HIPOTESIS ALTERNATIVA ETIQUETADA. Los modelos etiquetados obligan a determinar , a priori , las posibles observaciones outliers. Esto puede considerarse como una reminiscencia del caracter subjetivo con que se trataba esta teoria en sus inicios.

En este sentido se puede citar a Rosado (1984) , el cual indica que una de las principales aportaciones del metodo GAN (Generativo con Alternativa Natural) que el estudia radica en el hecho de que las observaciones consideradas outliers se obtienen a posteriori , con lo que se consigue retirar toda subjetividad en la seleccion de observaciones a testar como posibles outliers.

- RESPECTO A LA HIPOTESIS NULA. El problema fundamental es , al igual que en en el caso de la hipotesis alternativa, fijarla ya que generalmente no se tiene un conocimiento suficientemente preciso sobre el modelo o estructura matematica que siguen los datos.

En este sentido se puede citar la afirmacion de Barnett (1984):

"... un outlier para una distribucion normal puede no serlo para una distribucion Cauchy..."

Por tanto es necesario tener gran certeza sobre el modelo que se establece en esta hipótesis ya que se podría producir el denominado error de información.

- RESPECTO AL ESTADISTICO O TEST. La problemática que se plantea en este punto está muy relacionada, e incluso podría incluirse en el apartado anterior, ya que la distribución del estadístico a que da lugar la función test, se calcula basándose en el modelo que se propone en la hipótesis nula.

El mayor inconveniente radica en el hecho de que generalmente no es posible calcular la distribución exacta del estadístico, por lo que debe aproximarse, utilizándose para ello diversos métodos, siendo los más comunes:

- a) Método Conservativo: La aproximación de la distribución del estadístico da lugar a una región crítica conservativa, lo que presenta el inconveniente de que observaciones outliers pueden ser declaradas como inliers.
- b) Método Asintótico: En este caso se considera la distribución asintótica del estadístico.

Mas en esta situacion los outliers pierden parte de su importancia ya que estas observaciones cuando mas interesa ser estudiadas es en el caso de muestras de tamano reducido , que es donde ejercen su mayor efecto perturbador.

1.5.2 - DE LOS METODOS NO FORMALES O GRAFICOS .

Estas tecnicas presentan el inconveniente de que no existe una estructura teorica comun a todas ellas , ya que lo unico que se exige a las funciones que se utilizan en estos procedimientos es que tengan un fundamento relacionado con el problema que se desea analizar.

El problema fundamental , en casi todas las tecnicas que se engloban dentro de este apartado , es que no indican ni direccion ni magnitud de la desviacion del outlier respecto de la masa de datos.

1.5.3 EFECTO ENMASCARAMIENTO.

En distintos epigrafes se ha indicado que pueden presentarse situaciones en las que se desee analizar la presencia de mas de un outlier en la masa de datos , lo que da lugar a lo que algunos autores denominan problema de multiples outliers.

Este problema es abordado de formas muy diversas desde el punto de vista de las técnicas formales y presenta inconvenientes adicionales a los ya destacados para las mismas.

La primera forma de abordar el problema de múltiples outliers consiste en la aplicación sucesiva de procedimientos de identificación de un outlier, eliminando de la masa de datos la observación declarada como tal y aplicando la misma técnica de forma sucesiva a la muestra restringida hasta aceptar la hipótesis de trabajo.

Este procedimiento plantea como inconveniente lo que se denomina efecto de enmascaramiento de tipo A: considerar inliers a ciertas observaciones outliers.

El efecto de enmascaramiento de tipo A fue citado por primera vez por Pearson y Chandrasekar (1936), quienes afirman sobre su criterio:

"... este criterio parece eficiente únicamente en caso de que exista una única observación outlier..."

Este efecto causa la disminución de la potencia del método considerado, y presenta, además, un inconveniente de tipo metodológico, ya que el nivel de significación referido a la hipótesis de trabajo (todas las observaciones son inliers), no será el mismo que se le asigne al conjunto de observaciones declaradas como outliers.

Estos inconvenientes , hacen que algunos autores propongan otros tipos de procedimientos para abordar el problema de la presencia de multiples outliers.

Asi , se sugiere un procedimiento iterativo , que en una etapa preliminar selecciona las posibles observaciones outliers , de donde se extrae la observacion que menos sospechas despierta . A esta observacion se le aplica una tecnica de identificacion de un outlier en la muestra restringida (omitiendo los posibles outliers). En caso de aceptar la hipotesis nula , se incluiria esta observacion en la muestra restringida y se procederia a la seleccion de una nueva observacion de entre los restantes posibles outliers , aplicandose de nuevo el procedimiento anterior. En el momento en que se rechace la hipotesis nula, se concluye que la observacion que estaba siendo considerada , junto a las que aun no lo habian sido , son outliers.

Este procedimiento evita el problema del enmascaramiento de tipo A , pero presenta el inconveniente de que es necesario, en cada etapa, seleccionar la observacion a estudiar , lo cual quizas sea factible en el caso de datos univariantes pero no asi en el caso multivariante o en las denominadas situaciones mas estructuradas.

Otros tipos de procedimientos que evitan el efecto de enmascaramiento de tipo A , son las denominadas tecnicas de bloques , las que se aplican para la determinacion de

un numero fijo de outliers.

Desde un punto de vista metodologico , estos ultimos procedimientos , necesitan fijar inicialmente el numero de observaciones a testar como outliers , lo cual puede realizarse mediante los denominados procedimientos de saltos.

Fieller (1976) , indica que estos procedimientos para la identificacion de multiples outliers pueden presentar el denominado efecto de enmascaramiento de tipo B: considerar outliers observaciones que no lo son.

Ademas, tanto la determinacion del numero de outliers a testar , como el efecto de enmascaramiento de tipo B , afectan a las probabilidades de error de tipo I y tipo II.

CAPITULO II

APROXIMACION FORMAL AL ANALISIS CUALITATIVO

DE DATOS

--- Introduccion .

2.1 Propiedades de la relacion criterio .

2.2 Funciones de C-similitud . Seleccion .

2.3 Invarianza.

INTRODUCCION

En el capítulo anterior, se pone de manifiesto, que una de las etapas fundamentales de cualquier análisis estadístico, es el estudio de la calidad de los datos experimentales, con el fin de garantizar la máxima fiabilidad en las conclusiones que se obtengan del mismo.

Además, se ha reseñado, que el comportamiento anómalo de una observación, depende en gran medida de su relación con las restantes observaciones obtenidas bajo condiciones similares, y del análisis que se pretende llevar a cabo sobre las mismas.

Para responder a estas cuestiones, en este capítulo, se plantea el análisis del comportamiento de los elementos de un conjunto -expresado mediante relaciones de

preferencia entre los mismos- respecto de un criterio establecido, que depende del estudio que se desee realizar.

Por ello, este capítulo se centra en el estudio de la existencia de una función, bajo diversas restricciones, basadas en las relaciones definidas entre los elementos del conjunto inicial, proponiéndose un procedimiento de selección en aquellos casos en que no quede garantizada la unicidad de la misma.

Por último, se aborda el problema de la invarianza, con el objeto de que los resultados obtenidos no dependan del sistema de coordenadas elegido para representar las observaciones experimentales.

2.1 PROPIEDADES DE LA RELACION CRITERIO

Dado $X \subseteq \mathbb{R}^k$, sea C un criterio que permite comparar los elementos de X , mediante una relación binaria C_X , para la que dados $x, x' \in X$, $x C_X x'$ se interpreta como x es "al menos tan preferido como" x' , respecto del criterio C , considerados como elementos de X .

A partir de las relaciones binarias C_X , se consideran las siguientes definiciones:

Definición 2.1.1.

Dado $X \subseteq \mathbb{R}^k$, $x, x' \in X$, se dirá que x es "preferido" a x' respecto del criterio C , considerados como elementos de X , $x + C_X x'$, si :

$$x + C_X x' \langle == \rangle x C_X x' \text{ y } x' C_X x \quad (2.1)$$

Definición 2.1.2.

Dado $X \subseteq \mathbb{R}^k$, $x, x' \in X$, se dirá que x y x' son "equivalentes" respecto del criterio C , considerados como elementos de X , $x \sim C_X x'$, si :

$$x \sim C_X x' \langle == \rangle x C_X x' \text{ y } x' C_X x \quad (2.2)$$

Definición 2.1.3.

Un criterio C se dirá que es evaluable sobre $X \subseteq \mathbb{R}^k$, si C_X es fuertemente completa :

$$\forall x, x' \in X, == \rangle x C_X x' \text{ o } x' C_X x$$

Las relaciones C_X permiten introducir los conjuntos que se recogen en la siguiente

Definición 2.1.4.

Sea C un criterio evaluable sobre $X \subseteq \mathbb{R}^k$.

Dado $x \in X$, se define:

i) El conjunto C_X -dominante de x , $DC_X(x)$, como :

$$DC_X(x) = \{x' \in X / x' C_X x\}$$

ii) El conjunto C_X -dominado de x , $dc_X(x)$, como :

$$dc_X(x) = \{x' \in X / x C_X x'\}$$

Analogamente, se definen los conjuntos C_X -dominante

y C_X -dominado estrictos, que se representaran mediante $D(+C_X(x))$ y $d(+C_X(x))$, respectivamente.

Como generalizacion se introduce:

Definicion 2.1.5.

Sea C un criterio evaluable sobre $X \subseteq R^k$.

Dado $X' \subseteq X$, se define:

i) El conjunto C_X -dominante de X' , $DC_X(X')$, como:

$$DC_X(X') = \{x \in X / x C_X x', \forall x' \in X'\}$$

ii) El conjunto C_X -dominado de X' , $dC_X(X')$, como:

$$dC_X(X') = \{x \in X / x' C_X x, \forall x' \in X'\}$$

Proposicion 2.1.1.

Dado un criterio C , evaluable, sobre $X \subseteq R^k$ y $X' \subseteq X$, se verifica:

a) $DC_X(X') = \bigcap_{x' \in X'} DC_X(x')$

b) $dC_X(X') = \bigcap_{x' \in X'} dC_X(x')$

Estas definiciones permiten introducir:

Definicion 2.1.6.

Dado un criterio evaluable C , sobre $X \subseteq R^k$ y $X' \subseteq X$, se dira que X' esta C_X -dominado superiormente (inferiormente) sii $DC_X(X') \neq \emptyset$ ($dC_X(X') \neq \emptyset$).

Si $DC_X(X') \cap X' \neq \emptyset$ ($dC_X(X') \cap X' \neq \emptyset$), se considera que X' esta propiamente C_X -dominado superiormente (inferiormente).

Se dira que X' esta C_X -dominado (propiamente) si lo esta superior (propiamente) e inferiormente (propiamente).

Definicion 2.1.7.

Sea C un criterio evaluable sobre $X \subseteq R^k$, y $X' \subseteq X$. Entonces :

- i) Si X' esta C_X -dominado superiormente, el elemento $x^* \in DC_X(X')$ es C_X -dominante minimal (C_X D m) de X' -
 sii $\forall x_D \in DC_X(X') \Rightarrow x_D C_X x^*$
- ii) Si X' esta C_X -dominado inferiormente, el elemento $x_* \in dC_X(X')$ es C_X -dominado maximal (C_X d M) de X'
 sii $\forall x_d \in dC_X(X') \Rightarrow x_* C_X x_d$

Proposicion 2.1.2.

Sea C un criterio evaluable sobre $X \subseteq R^k$ y sea $X' \subseteq X$. Entonces :

- i) si x_1^*, x_2^* son elementos (C_X D m) de $X' \Rightarrow x_1^* \sim_{C_X} x_2^*$
- ii) si x_*^1, x_*^2 son elementos (C_X d m) de $X' \Rightarrow x_*^1 \sim_{C_X} x_*^2$

El resultado se deduce de forma inmediata a partir de la definicion 2.1.7.

En muchas situaciones, las relaciones binarias, C_X , consideradas sobre los conjuntos X de R^k son transitivas, por lo que se da la siguiente

Definicion 2.1.8.

Un criterio C se dira T-evaluable sobre $X \subseteq R^k$, si la relacion binaria C_X define un orden debil en X .

Proposición 2.1.3.

Si C es un criterio T -evaluabile sobre X , la relación dada mediante (2.2) define una relación de equivalencia en dicho conjunto.

Así, si se representa por $\hat{x} = \{x' \in X / x' \sim_{C_X} x\}$, y por $\hat{X} = \{\hat{x} / x \in X\}$, la relación definida en este conjunto mediante:

$$\hat{x} C_X \hat{x}' \iff x C_X x' \text{ es un orden simple en } \hat{X}.$$

Proposición 2.1.4.

Sea C un criterio T -evaluabile sobre $X \subseteq \mathbb{R}^k$ y \hat{X} su conjunto cociente, entonces si existe una clase $(C_X D m)$ o $(C_X d M)$ de $\hat{X}' \subseteq \hat{X}$, es única.

Este resultado se deduce a partir de la proposición 2.1.2.

Proposición 2.1.5.

Si C es un criterio T -evaluabile sobre X , de modo que $\#(\hat{X})$ es finito, X está propiamente C_X -dominado.

2.2 FUNCIONES DE C-SIMILITUD . SELECCION

Las relaciones binarias C_X definidas sobre los subconjuntos de \mathbb{R}^k , únicamente determinan el sentido de preferencia existente entre los elementos de dichos conjuntos respecto del criterio en cuestión, sin embargo, no permiten evaluar si los elementos presentan un comporta-

miento similar frente al criterio bajo estudio. Por esta razón, se define a continuación una función (real), sobre el producto cartesiano $X \times X$.

Definición 2.2.1.

Sea C un criterio T -evaluable sobre $X \subseteq \mathbb{R}^k$.

Se dirá que una función

$$SC_X : X \times X \rightarrow \mathbb{R}$$

es una C -similitud sobre X si verifica:

S1. $SC_X(x, x') \geq 0 \quad \forall x, x' \in X$

S2. $SC_X(x, x') = 0$ si $x \sim_C x'$

S3. $SC_X(x, x') = SC_X(x', x) \quad \forall x, x' \in X$

S4. $\forall x, x', x'' \in X / x \sim_C x' \sim_C x''$ se verifica:

S41. $SC_X(x, x'') \leq SC_X(x, x') + SC_X(x', x'')$

S42. $SC_X(x, x'') \geq \max\{SC_X(x, x'), SC_X(x', x'')\}$

Esta función refleja la similaridad de cada par de elementos de X respecto de C y cumple las siguientes propiedades:

Propiedad 1.

SC_X verifica la desigualdad triangular.

Demostración:

Dados $x, x', x'' \in X$, se presenta una de las siguientes situaciones:

$x \sim_C x' \sim_C x''$, para la que se cumple la afirmación debido a (S41)

$x \sim_C x'' \sim_C x' \implies (S42) \implies SC_X(x, x'') \leq SC_X(x, x')$

$$x' \in C_X \text{ y } C_X x'' \implies (S42) \implies SC_X(x, x'') \leq SC_X(x', x'')$$

De donde se deduce :

$$SC_X(x, x'') \leq SC_X(x, x') + SC_X(x', x'')$$

para cualesquiera $x, x', x'' \in X$

Propiedad 2.

Dados $x, x' \in X$:

$$x \sim_{C_X} x' \iff SC_X(x, y) = SC_X(x', y) \quad \forall y \in X .$$

Demostracion:

Si $x \sim_{C_X} x'$, se tiene

$$SC_X(x, y) \leq SC_X(x, x') + SC_X(x', y) \leq SC_X(x', x) + SC_X(x, y) \implies$$

$$\implies SC_X(x, y) = SC_X(x', y)$$

El reciproco es evidente.

Propiedad 3.

La funcion $SC_{\hat{X}} : \hat{X} \times \hat{X} \longrightarrow \mathbb{R}$

definida como

$$SC_{\hat{X}}(\hat{x}, \hat{y}) = SC_X(x, y)$$

esta bien definida y $(\hat{X}, SC_{\hat{X}})$ es un espacio metrico.

Propiedad 4.

Si $X' \subseteq X$ es un conjunto propiamente C_X -dominado, $\forall x, x' \in X'$, se verifica:

$$SC_X(x, x') \leq SC_X(x_d, x_D) \quad \forall (x_d, x_D) \in dC_X(X') \times DC_X(X')$$

y ademas

$$SC_X(x_*, x^*) = \sup_{x, x' \in X'} \{ SC_X(x, x') \}$$

donde x^* , x_* , son los elementos $(C_X D m)$ y $(C_X d M)$ de X' , respectivamente.

En ocasiones puede ser de especial interes comparar la similitud de los elementos de X respecto de un valor fijado $x_0 \in X$, por lo que se introduce :

Definicion 2.2.2.

Sea C un criterio T -evaluabile sobre $X \subseteq R^k$ y SC_X una funcion de C -similitud sobre X . Dado $x_0 \in X$, se define :

$$SC_{X,x_0} : X \rightarrow R$$

$$SC_{X,x_0}(x) = SC_X(x, x_0)$$

Esta permite enunciar

Propiedad 5.

Si $X' \subseteq X$ es un conjunto C_X -dominado, se verifica:

$$SC_{X,x_D}(x) < SC_{X,x_D}(x') \implies x_D C_X x + C_X x' C_X x_d \implies$$

$$\implies SC_{X,x_d}(x) \geq SC_{X,x_d}(x')$$

$$\forall (x_d, x_D) \in dC_X(X') \times DC_X(X') \text{ y } \forall x, x' \in X$$

Existen situaciones en las que la funcion SC_X no recoge la informacion que proporciona la relacion binaria ya que la C -similitud de dos elementos cualesquiera no refleja de que modo estan relacionados los mismos respecto del criterio C , salvo que dichos elementos sean C_X -equivalentes.

Por esta razon se introduce la siguiente

Definicion 2.2.3.

Sea C un criterio T-evaluable sobre $X \subseteq R^k$.

Se dira que una funcion

$$SSC_X : X \times X \rightarrow R$$

es una C-similitud signada sobre X si verifica :

$$SS1. SSC_X(x, x') \geq 0 \text{ si } x C_X x'$$

$$SS2. SSC_X(x, x') = - SSC_X(x', x) \quad \forall x, x' \in X$$

$$SS3. \forall x, x', x'' \in X / x C_X x' \wedge C_X x''$$

$$SS31. SSC_X(x, x'') \leq SSC_X(x, x') + SSC_X(x', x'')$$

$$SS32. SSC_X(x, x'') \geq \max \{SSC_X(x, x'), SSC_X(x', x'')\}$$

La funcion SSC_X verifica las siguientes propiedades:

Propiedad 6.

Si SSC_X es una funcion de C-similitud signada sobre X, entonces dados $x, x' \in X$:

$$SSC_X(x, x') = 0 \text{ si } x \sim C_X x'.$$

Demostracion:

$$SSC_X(x, x') = 0 \iff SSC_X(x, x') = 0 \text{ y } SSC_X(x', x) = 0$$

$$\iff x C_X x' \text{ y } x' C_X x \iff x \sim C_X x'.$$

Propiedad 7.

Dados $x, x' \in X$. Entonces :

$$x \sim C_X x' \iff SSC_X(x, y) = SSC_X(x', y) \quad \forall y \in X.$$

Demostracion:

Si $x \sim_{C_X} x'$, los dos casos que se pueden presentar:

$$\begin{aligned}
 & x' \sim_{C_X} x \text{ y } C_X y \quad ; \quad y \in C_X \text{ y } \sim_{C_X} x' \\
 & x' \sim_{C_X} x \text{ y } C_X y \implies SSC_X(x', y) \leq SSC_X(x', x) + SSC_X(x, y) \leq \\
 & \leq SSC_X(x, x') + SSC_X(x', y) \quad \text{por tanto} \\
 & SSC_X(x, y) = SSC_X(x', y)
 \end{aligned}$$

Analogamente se analizaria el otro caso.

El reciproco es inmediato a partir de la propiedad 6.

Propiedad 8.

La funcion $SSC_{\hat{X}} : \hat{X} \times \hat{X} \rightarrow R$

$$SSC_{\hat{X}}(\hat{x}, \hat{y}) = SSC_X(x, y)$$

esta bien definida.

Propiedad 9.

Si $X' \subseteq X$, esta propiamente C_X -dominado, se verifica

$$SSC_X(x_d, x_D) \leq SSC_X(x, x') \leq SSC_X(x_D, x_d)$$

$$\forall x, x' \in X \quad \text{y} \quad \forall (x_d, x_D) \in dC_X(X') \times DC_X(X') .$$

Ademas

$$SSC_X(x^*, x_*) = \sup_{x, x' \in X'} \{SSC_X(x, x')\} = - SSC_X(x_*, x^*) = \dots$$

$$\dots = \inf_{x, x' \in X'} \{SSC_X(x, x')\}$$

siendo x^* , x_* , elementos $(C_X \text{ D } m)$ y $(C_X \text{ d } M)$, de X' respectivamente.

Las funciones introducidas en las definiciones 2.2.1 y 2.2.3 permiten definir una relacion binaria sobre $X \times X$, tal como se recoge a continuacion.

Definicion 2.2.4.

Dados $x_i \in X$, $i=1,2,3,4$, se define la relacion binaria C'_X sobre $X \times X$ como :

$$(x_1, x_2) C'_X (x_3, x_4) \text{ sii } SC_X(x_1, x_2) \leq SC_X(x_3, x_4) \quad (2.3)$$

Esta definicion refleja que el par (x_1, x_2) es "al menos tan similar" como el par (x_3, x_4) respecto del criterio C , considerados como elementos de X .

Esta relacion permite introducir las siguientes:

Definicion 2.2.5.

Dados $x_i \in X$, $i=1,2,3,4$, se dira que el par (x_1, x_2) es "mas similar" que (x_3, x_4) respecto del criterio C , considerados como elementos de X

$$(x_1, x_2) +C'_X (x_3, x_4), \text{ sii:} \quad (2.4)$$

$$(x_1, x_2) C'_X (x_3, x_4) \text{ y } (x_3, x_4) \not C'_X (x_1, x_2) \quad \langle == \rangle$$

$$\langle == \rangle SC_X(x_1, x_2) < SC_X(x_3, x_4)$$

Definicion 2.2.6.

Dados $x_i \in X$, $i=1,2,3,4$, se dira que el par (x_1, x_2) es "equivalente" al par (x_3, x_4) respecto del criterio C , considerados como elementos de X

$$(x_1, x_2) \sim C'_X (x_3, x_4), \text{ sii:} \quad (2.5)$$

$$(x_1, x_2) C'_X (x_3, x_4) \text{ y } (x_3, x_4) C'_X (x_1, x_2) \quad \langle == \rangle$$

$$\langle == \rangle SC_X(x_1, x_2) = SC_X(x_3, x_4)$$

Proposición 2.2.1.

- i) C'_X define un un orden debil sobre $X \times X$.
- ii) $+C'_X$ define un orden debil estricto sobre $X \times X$
- iii) $\sim C'_X$ define una relacion de equivalencia sobre $X \times X$.

Las relaciones dadas mediante (2.3), (2.4) y (2.5) se pueden definir de la misma forma a traves de una funcion de similitud signada , para la que se puede enunciar una proposicion analoga a la proposicion 2.2.1. Ademas , en esta situacion se verifica

Proposición 2.2.2.

Las relaciones C'_X , $+C'_X$, $\sim C'_X$, definidas a traves de SSC_X mediante (2.3) , (2.4) y (2.5) verifican el axioma de signo reverso, es decir:

$$\begin{aligned} (x_1, x_2) C'_X (x_3, x_4) &\langle == \rangle (x_4, x_3) C'_X (x_2, x_1) \\ (x_1, x_2) +C'_X (x_3, x_4) &\langle == \rangle (x_4, x_3) +C'_X (x_2, x_1) \\ (x_1, x_2) \sim C'_X (x_3, x_4) &\langle == \rangle (x_4, x_3) \sim C'_X (x_2, x_1) \end{aligned}$$

En este caso se ha de destacar que la relacion C'_X lleva implicita la relacion C_X , ya que :

$$\begin{aligned} (x, x) C'_X (x, y) &\langle == \rangle SSC_X(x, x) \leq SSC_X(x, y) \langle == \rangle \\ &\langle == \rangle SSC_X(x, y) \geq 0 \langle == \rangle x C_X y. \end{aligned}$$

Para la posterior aplicacion practica de esta teoria es deseable disponer de una funcion sobre X en lugar de sobre X x X que refleje el sentido de preferencia establecido en las relaciones C_X y C'_X .

Para ello se estudiara bajo que condiciones existe una funcion

$$IC_X : X \rightarrow R$$

que verifique

$$SSC_X(x,y) \leq SSC_X(z,w) \Leftrightarrow IC_X(x) - IC_X(y) \leq IC_X(z) - IC_X(w)$$

ya que entonces se tendria que

$$y C_X x \Leftrightarrow SSC_X(y,x) \geq SSC_X(x,x) \Leftrightarrow IC_X(y) \geq IC_X(x)$$

y por tanto se recoge la informacion que proporcionan las relaciones .

Las funciones de similitud y similitud signadas introducidas en las definiciones 2.2.1 y 2.2.3 , incluyen como casos particulares a las definidas por:

$$SC_X(x,x') = \begin{cases} k & \text{si } x +C_X x' \text{ o } x' +C_X x \\ 0 & \text{si } x \sim C_X x' \end{cases}$$

$$SSC_X(x,x') = \begin{cases} k & \text{si } x +C_X x' \\ 0 & \text{si } x \sim C_X x' \\ -k & \text{si } x' +C_X x \end{cases}$$

siendo k una constante positiva.

Estas funciones no recogerían de forma deseable toda la información que proporciona la relación binaria, por lo que se introduce la siguiente

Definición 2.2.7.

Una función de C-similitud (signada), se dirá aditiva, si cumple:

$$\forall x, x', x'' \in X / x C_X x' C_X x'' \implies \\ \implies SC_X(x, x'') = SC_X(x, x') + SC_X(x', x'')$$

Teorema 2.2.1.

Si SSC_X es una función de C-similitud signada aditiva sobre X, y $C'_X, +C'_X, \sim C'_X$, son las relaciones dadas por (2.3), (2.4), (2.5), a través de SSC_X , y se verifican las condiciones:

T1. $\forall x_i, x'_i \in X, i=1,2,3$ tales que

$$\left. \begin{array}{l} (x_1, x_2) C'_X (x'_1, x'_2) \\ (x_2, x_3) C'_X (x'_2, x'_3) \end{array} \right\} \implies (x_1, x_3) C'_X (x'_1, x'_3)$$

T2. $\forall y, x_i \in X, i=1,2,3,4$ tales que

$$\left. \begin{array}{l} (x_1, x_2) C'_X (x_3, x_4) \\ (x_3, x_4) C'_X (y, y) \end{array} \right\} \implies$$

$\implies \exists z_1, z_2 \in X$ de forma que

$$(x_1, z_1) \sim C'_X (x_3, x_4) \quad y \quad (z_2, x_2) \sim C'_X (x_3, x_4)$$

Entonces :

$$\exists IC_X : X \rightarrow R \text{ tal que } \forall x_i \in X, i=1,2,3,4 \quad (2.6)$$

$$(x_1, x_2) + C'_X(x_3, x_4) \Leftrightarrow IC_X(x_1) - IC_X(x_2) < IC_X(x_3) - IC_X(x_4)$$

y ademas se verifica que si

$$IC_X^* : X \rightarrow R$$

es otra funcion que cumple (2.6) :

$$IC_X^* = \alpha IC_X + \beta \quad / \alpha > 0$$

Demostracion:

De los resultados obtenidos en las Proposiciones 2.2.1 y 2.2.2 , y de T1, T2, se deduce que el par $(X, +C'_X)$, tiene estructura de diferencia algebraica, de lo que se obtiene el resultado propuesto a partir de Krantz et al (1971).

Corolario 2.2.1.

En las condiciones del teorema anterior se verifica:

- i) $(x_1, x_2) \sim C'_X(x_3, x_4)$ sii $IC_X(x_1) - IC_X(x_2) = IC_X(x_3) - IC_X(x_4)$
- ii) $(x_1, x_2) C'_X(x_3, x_4)$ sii $IC_X(x_1) - IC_X(x_2) \leq IC_X(x_3) - IC_X(x_4)$
- iii) $x C_X y$ sii $IC_X(x) \geq IC_X(y)$
- iv) $x + C_X y$ sii $IC_X(x) > IC_X(y)$
- v) $x \sim C_X y$ sii $IC_X(x) = IC_X(y)$

La importancia de estos resultados radica en el hecho, de que la funcion IC_X es unica salvo transformacion lineal, lo que significa que incorpora dos constantes (un cero arbitrario y una unidad arbitraria).

En el caso particular en que X estuviese C_X -dominado, se podría considerar:

$$IC_X(DC_X(X)) \text{ -----} \rightarrow 1$$

$$IC_X(dC_X(X)) \text{ -----} \rightarrow 0$$

con lo que $IC_X(x) \in [0,1] \quad \forall x \in X$, lo que se puede interpretar como un conjunto difuso sobre X , definido a través del criterio C .

El teorema da lugar a una condición suficiente para la existencia de una función verificando (2.6), sin embargo, cuando X es un subconjunto de R^k , con cardinal finito, es posible establecer un teorema que proporciona una condición necesaria y suficiente para la existencia de dicha función, para lo cual es necesario dar previamente los siguientes resultados.

Sea $X \subset R^k / X \subseteq \prod_{i=1}^k X_i ; X_i \subset R$ con $\#(X_i) < +\infty$, y sea $X^m = X \times \dots \times X$.

Definición 2.2.8.

En X^m , $m > 1$, dados $x^i, y^i \in X \quad i=1, \dots, m$ se define la relación E_m , por:

$(x^1, \dots, x^m) E_m (y^1, \dots, y^m)$ si (x^1, \dots, x^m) es una permutación de (y^1, \dots, y^m) para $i=1, \dots, k$.

La relación así definida es una relación de equivalencia.

Lema 2.2.1.

Sea C un criterio T -evaluabile en $X \subset \mathbb{R}^k$, y sea $X \subseteq \bigcap_{i=1}^k X_i$; $\#(X_i) < +\infty$ / $x_{ij} = x_{ij'}$, $\forall j \neq j'$; $i=1, \dots, k$.

Entonces :

(2.7) $\exists u_1, \dots, u_k$ / $u_i : X_i \rightarrow \mathbb{R}$, verificando

$$x + C_X y \implies \sum_{i=1}^k u_i(x_i) > \sum_{i=1}^k u_i(y_i) \quad \text{sii}$$

(2.8) $\forall x^1, \dots, x^m, y^1, \dots, y^m \in X$, $m = 2, 3, \dots$ se cumple :

$$[(x^1, \dots, x^m) E_m (y^1, \dots, y^m) / x^j + C_X y^j \quad \text{o}$$

$$x^j = y^j, j=1, \dots, m-1] \implies x^m + C_X y^m$$

Demostracion:

Suponiendo (2.7) y dado $(x^1, \dots, x^m) E_m (y^1, \dots, y^m) /$

$$x^j + C_X y^j \quad \text{o} \quad x^j = y^j, j < m \implies$$

$$\sum_{j=1}^{m-1} \sum_{i=1}^k u_i(x_i^j) \geq \sum_{j=1}^{m-1} \sum_{i=1}^k u_i(y_i^j) \quad (2.9)$$

Mas como $(x^1, \dots, x^m) E_m (y^1, \dots, y^m) \implies$

$$\implies \sum_{j=1}^m \sum_{i=1}^k u_i(x_i^j) = \sum_{j=1}^m \sum_{i=1}^k u_i(y_i^j) \quad (2.10)$$

Mediante (2.9) y (2.10) se deduce:

$$\sum_{i=1}^k u_i(x_i^m) \leq \sum_{i=1}^k u_i(y_i^m) \implies x^m + C_X y^m$$

Recíprocamente:

Dado $X_i = \{ x_{ij} \mid j=1, \dots, n_i \}$, $i = 1, \dots, k$ se considera el vector

$$U = \{ u_i(x_{ij}) \mid i=1, \dots, k ; j=1, \dots, n_i \}. \quad (2.11)$$

Sea $H = \# \{ (x, y) \mid x + C_X y, x, y \in X \}$, por lo que se puede expresar

$$\{ (x, y) \mid x + C_X y \} = \{ (x^1, y^1), \dots, (x^H, y^H) \}$$

Dado (x^h, y^h) , $h = 1, \dots, H$ se definen los vectores

$$a_x^h = \{ (a_x^h)_{ij} \mid i=1, \dots, k ; j=1, \dots, n_i \} \quad \text{donde}$$

$$(a_x^h)_{ij} = \begin{cases} 1 & \text{si } x_{ij} = x_i^h \\ 0 & \text{en otro caso} \end{cases}$$

y análogamente

$$a_y^h = \{ (a_y^h)_{ij} \mid i=1, \dots, k ; j=1, \dots, n_i \} \quad \text{donde}$$

$$(a_y^h)_{ij} = \begin{cases} 1 & \text{si } x_{ij} = y_i^h \\ 0 & \text{en otro caso} \end{cases}$$

(2.12)

$$\text{Sea } a_{xy}^h = a_x^h - a_y^h = \{ (a_{xy}^h)_{ij} \mid i=1, \dots, k ; j=1, \dots, n_i \}$$

por lo que

$$(a_{xy}^h)_{ij} = \begin{cases} 1 & \text{si } x_{ij} = x_i^h ; x_{ij} \neq y_i^h \\ -1 & \text{si } x_{ij} \neq x_i^h ; x_{ij} = y_i^h \\ 0 & \text{si } \begin{cases} (x_{ij} = x_i^h ; x_{ij} = y_i^h) \\ \text{ó} \\ (x_{ij} \neq x_i^h ; x_{ij} \neq y_i^h) \end{cases} \end{cases}$$

A partir de los vectores definidos en (2.11) y (2.12), se considera el sistema

$$U' a_{xy}^h > 0 \quad h = 1, \dots, H \quad (2.13)$$

por lo que si este sistema tiene solución en U , se verifica que

$$\exists u_i : X_i \rightarrow \mathbb{R}, \quad i=1, \dots, k$$

tales que

$$x + C_X y \implies \sum_{i=1}^k u_i(x_i) > \sum_{i=1}^k u_i(y_i)$$

En caso de que no existiese solución para dicho sistema, el Teorema de la Alternativa (Tucker 1956), asegura

$\exists r_1, \dots, r_H \geq 0$ y enteros, no todos nulos, tales que

$$\sum_{h=1}^H r_h (a_{xy}^h)_{ij} = 0 \quad i=1, \dots, k; \quad j=1, \dots, n_i.$$

Es decir :

$$\sum_{h=1}^H r_h (a_x^h)_{ij} = \sum_{h=1}^H r_h (a_y^h)_{ij} \quad i=1, \dots, k; \quad j=1, \dots, n_i$$

de donde se deduce que

$$(r_1 x^1, s, \dots, r_H x^H, s) E_{r_1 + \dots + r_H} (r_1 y^1, s, \dots, r_H y^H, s)$$

lo que contradice la hipótesis inicial ya que

$$x^i + C_X y^i \quad i = 1, \dots, H$$

Por tanto ha de existir solución para el sistema (2.13)

Lema 2.2.2.

Sea C un criterio T-evaluable en $X \subset \mathbb{R}^k$, y sea $X \subseteq \prod_{i=1}^k X_i$, con $\#(X_i) < +\infty$.

Entonces:

$$\exists u_1, \dots, u_k \quad / \quad u_i : X_i \rightarrow \mathbb{R}, \quad (2.14)$$

verificando

$$x +_{C_X} y \leq \sum_{i=1}^k u_i(x_i) \geq \sum_{i=1}^k u_i(y_i) \quad \text{sii}$$

$$\forall x^1, \dots, x^m, y^1, \dots, y^m \in X, \quad m=2,3,\dots \quad (2.15)$$

se cumple:

$$[(x^1, \dots, x^m) E_m (y^1, \dots, y^m) / x^j +_{C_X} y^j, j=1, \dots, m-1] \Leftrightarrow$$

$$\Leftrightarrow x^m +_{C_X} y^m$$

Demostracion:

Suponiendo (2.14) y dado $(x^1, \dots, x^m) E_m (y^1, \dots, y^m)$ tal que $x^j +_{C_X} y^j$, $j < m$, se tiene:

$$\sum_{j=1}^{m-1} \sum_{i=1}^k u_i(x_i^j) \geq \sum_{j=1}^{m-1} \sum_{i=1}^k u_i(y_i^j),$$

y ya que

$$(x^1, \dots, x^m) E_m (y^1, \dots, y^m) \Rightarrow$$

$$\Rightarrow \sum_{j=1}^m \sum_{i=1}^k u_i(x_i^j) = \sum_{j=1}^m \sum_{i=1}^k u_i(y_i^j)$$

Por tanto:

$$\sum_{i=1}^k u_i(x_i^m) \leq \sum_{i=1}^k u_i(y_i^m) \Leftrightarrow x^m +_{C_X} y^m$$

Recíprocamente

Dado $X_i = \{ x_{ij} \mid j = 1, \dots, n_i \}$ $i = 1, \dots, k$ se considera el vector

$$U = \{ u_i(x_{ij}) \mid i=1, \dots, k ; j=1, \dots, n_i \}.$$

Sea $H = \# \{ (x,y) / x +C_X y ; x,y \in X \}$ y considerese el conjunto de pares $\{ (x,y) / x \sim C_X y , x \neq y \}$ donde de los pares (x,y) , (y,x) solo se considera uno de ellos y sea $M-H$ el cardinal de este conjunto.

En estas condiciones se tienen los conjuntos :

$$\{ (x^i, y^i) \mid i=1, \dots, H / x^i +C_X y^i \}$$

$$\{ (x^i, y^i) \mid i=H+1, \dots, M / x^i \sim C_X y^i \}$$

Para los pares (x^i, y^i) , $h=1, \dots, M$, se definen los vectores

$$a_x^h , a_y^h , a_{xy}^h$$

de la misma forma que en el lema 2.2.1.

A partir de estos vectores se considera el sistema

$$U' a_{xy}^h > 0 \quad h=1, \dots, H$$

$$U' a_{xy}^h = 0 \quad h=H+1, \dots, M$$

Si este sistema tiene solución , se verifica :

$$\exists u_i : X_i \rightarrow R , i=1, \dots, k$$

tales que

$$x +C_X y \implies \sum_{i=1}^k u_i(x_i) > \sum_{i=1}^k u_i(y_i)$$

$$x \sim_{C_X} y \implies \sum_{i=1}^k u_i(x_i) = \sum_{i=1}^k u_i(y_i)$$

Si el sistema no admite solución, a partir del Teorema de la Alternativa (Tucker 1956), se deduce

$\exists r_h, h=1, \dots, M / r_h \geq 0, h=1, \dots, H$, no todos nulos de forma que

$$\sum_{h=1}^H r_h (a_{xy}^h)_{ij} = 0 \quad i=1, \dots, k, j=1, \dots, n_i$$

Debido a que los valores $(a_{xy}^h)_{ij} \in \{1, 0, -1\}$, puede considerarse que $r_h \in \mathbb{Z}$, además, pueden considerarse positivos, ya que si alguno de ellos no lo fuese, bastaría reemplazar el par (x, y) por el (y, x) y r_h por $-r_h$, con lo cual se llegaría a una contradicción similar a la del lema 2.2.1.

Teorema 2.2.2.

Sea C un criterio T -evaluable en $X \subset \mathbb{R}^k$ con $\#(X) < +\infty$, y sea SSC_X una función de C -similitud signada sobre X . Entonces

$$\exists IC_X: X \rightarrow \mathbb{R} \text{ verificando} \quad (2.16)$$

$$(x, y) + C'_X(z, w) \implies IC_X(x) - IC_X(y) < IC_X(z) - IC_X(w)$$

si

para $m=2, 3, \dots$ dados $x^1, y^1, z^1, w^1 \in X, i=1, \dots, m$

cumpliendo

$$I) x^1, \dots, x^m, w^1, \dots, w^m \text{ es una permutación de } y^1, \dots, y^m, z^1, \dots, z^m$$

$$\text{II) } (x^i, y^i) + C'_X(z^i, w^i) \quad i=1, \dots, m-1$$

se tiene :

$$(z^m, w^m) + C'_X(x^m, y^m)$$

Demostracion:

Si se verifica (2.16) :

$$\text{dado que } (x^j, y^j) + C'_X(z^j, w^j) \quad , \quad j < m$$

se tiene que:

$$IC_X(x^j) - IC_X(y^j) < IC_X(z^j) - IC_X(w^j) \quad , \quad j=1, \dots, m-1$$

y como ademas:

$$\sum_{j=1}^m [IC_X(x^j) - IC_X(y^j)] = \sum_{j=1}^m [IC_X(z^j) - IC_X(w^j)]$$

se obtiene que :

$$IC_X(x^m) - IC_X(y^m) > IC_X(z^m) - IC_X(w^m)$$

y por tanto

$$(x^m, y^m) + C'_X(z^m, w^m).$$

Reciprocamente :

Dado $X = \{x_1, \dots, x_N\}$, se considera el vector

$$U = \{ u(x_i) \quad , \quad i=1, \dots, N \}$$

y el conjunto

$$A = \{ [(x, y), (z, w)] / (x, y) + C'_X(z, w) \quad , \quad x, y, z, w \in X \}$$

Ya que $\#(X) < +\infty \implies \#(A) = H < +\infty$, por lo que se puede considerar

$$A = \{ [(x^h, y^h) \quad , \quad (z^h, w^h)] \quad , \quad h=1, \dots, H \}$$

Dado $(x^h, y^h), (z^h, w^h)$, se definen los vectores

$$(a_x^h) = \{ (a_x^h)_i, i=1, \dots, N \}$$

$$(a_y^h) = \{ (a_y^h)_i, i=1, \dots, N \}$$

$$(a_z^h) = \{ (a_z^h)_i, i=1, \dots, N \}$$

$$(a_w^h) = \{ (a_w^h)_i, i=1, \dots, N \}$$

donde

$$(a_x^h)_i = \begin{cases} 1 & \text{si } x^h = x_i \\ 0 & \text{si } x^h \neq x_i \end{cases}$$

$$(a_y^h)_i = \begin{cases} 1 & \text{si } y^h = x_i \\ 0 & \text{si } y^h \neq x_i \end{cases}$$

$$(a_z^h)_i = \begin{cases} 1 & \text{si } z^h = x_i \\ 0 & \text{si } z^h \neq x_i \end{cases}$$

$$(a_w^h)_i = \begin{cases} 1 & \text{si } w^h = x_i \\ 0 & \text{si } w^h \neq x_i \end{cases}$$

y

$$a_{xyzw}^h = a_y^h + a_z^h - a_x^h - a_w^h$$

A partir de estos vectores se considera el sistema

$$U' a_{xyzw}^h > 0, \quad h=1, \dots, H \quad (2.17)$$

Si este sistema admite solución, entonces:

$$\exists IC_X : X \rightarrow R$$

verificando:

$$(x,y) + C'_X(z,w) \implies IC_X(x) - IC_X(y) < IC_X(z) - IC_X(w)$$

para lo cual basta considerar $IC_X(x_1) = u(x_1)$.

Si no existiese solución para el sistema (2.17), por Tucker (1956), existirían $r_h \in Z^+$, no todos nulos, tales que:

$$\sum_{h=1}^H r_h (a_{xyzw}^h)_j = 0 \quad j=1, \dots, N$$

y de

$(x^h, y^h) + C'_X(z^h, w^h)$, se deduce que existe una secuencia:

$$(x^1, y^1) + C'_X(z^1, w^1), \dots, (x^m, y^m) + C'_X(z^m, w^m)$$

donde $x^1, \dots, x^m, w^1, \dots, w^m$ es una permutación de $y^1, \dots, y^m, z^1, \dots, z^m$.

Con lo cual, si $m > 1$, esta condición contradice la hipótesis de partida y si $m = 1$, entonces:

$$(x,y) + C'_X(x,y) \text{ ó } (x,x) + C'_X(y,y)$$

que también contradice la hipótesis inicial.

Por tanto ha de existir solución para el sistema (2.17).

Teorema 2.2.3.

Sea C un criterio T-evaluable en $X \subset R^k$ con $\#(X) < +\infty$, y sea SSC_X una función de C-similitud signada sobre X . Entonces

$$\exists IC_X: X \rightarrow R \quad \text{verificando} \quad (2.18)$$

$$(x,y) + C'_X(z,w) \iff IC_X(x) - IC_X(y) < IC_X(z) - IC_X(w)$$

si (2.19)

para $m=2,3,\dots$, dados $x^i, y^i, z^i, w^i \in X, i=1,\dots,m$

cumpliendo

I') $x^1, \dots, x^m, w^1, \dots, w^m$ es una permutacion de $y^1, \dots, y^m, z^1, \dots, z^m$

II') $(x^i, y^i) C'_X (z^i, w^i) \quad i=1, \dots, m-1$

se tiene :

$$(z^m, w^m) C'_X (x^m, y^m)$$

Demostracion:

Si se verifica (2.18) :

$$IC_X(x^j) - IC_X(y^j) \leq IC_X(z^j) - IC_X(w^j) \quad , \quad j=1, \dots, m-1$$

y como ademas:

$$\sum_{j=1}^m [IC_X(x^j) - IC_X(y^j)] = \sum_{j=1}^m [IC_X(z^j) - IC_X(w^j)]$$

se obtiene que :

$$IC_X(x^m) - IC_X(y^m) \geq IC_X(z^m) - IC_X(w^m)$$

y por tanto

$$(x^m, y^m) + C'_X (z^m, w^m).$$

Reciprocamente:

De la condicion (2.19) , en el caso bidimensional, se deduce:

$$\left. \begin{aligned} & [(x^1, y^1), \dots, (x^m, y^m)] E_m [(z^1, w^1), \dots, (z^m, w^m)] \\ & (x^j, y^j) C'_X (z^j, w^j) \quad , j=1, \dots, m-1 \end{aligned} \right\} \Rightarrow$$

$\Rightarrow (x^m, y^m) + C'_X (z^m, w^m)$ y por lema 2.2.2 :

$\exists u_1, u_2 : X \rightarrow R$ de modo que

$$(x, y) + C'_X (z, w) \Leftrightarrow u_1(x) + u_2(y) < u_1(z) + u_2(w)$$

y de (2.19)

$$(x, y) + C'_X (z, w) \Leftrightarrow (w, z) + C'_X (y, x)$$

por lo que

$$(x, y) + C'_X (z, w) \Leftrightarrow u_1(w) + u_2(z) < u_1(y) + u_2(x)$$

Por tanto definiendo

$$IC_X(x) = u_1(x) - u_2(x)$$

se tiene que

$$(x, y) + C'_X (z, w) \Leftrightarrow IC_X(x) - IC_X(y) < IC_X(z) - IC_X(w) .$$

En el caso en que no exista IC_X verificando las condiciones impuestas en los teoremas 2.2.1, 2.2.2, 2.2.3, se construira una funcion que refleje unicamente el sentido de preferencia establecido en la relacion C_X .

Definicion 2.2.9.

Sea C un criterio sobre $X \subseteq R^k$, que define una relacion binaria C_X . Se dira que una funcion

$$IC_X : X \rightarrow R , \text{ es :}$$

i) Estrictamente C_X -creciente si :

$$x + C_X y \Leftrightarrow IC_X(x) > IC_X(y)$$

ii) Estrictamente C_X -decreciente si :

$$x C_X y \Leftrightarrow IC_X(x) < IC_X(y)$$

A continuacion se recogen una serie de resultados que caracterizan distintas situaciones en las que existen funciones estrictamente C_X -crecientes (decrecientes).

Teorema 2.2.4.

Sea C un criterio T -evaluabile sobre $X \subset \mathbb{R}^k$, con $\#(X)$ finito. Entonces:

$$IC_X : X \rightarrow \mathbb{R}$$

estrictamente C_X -creciente.

Demostracion:

$$\text{Basta considerar } \#(d+C_X(x)) = IC_X(x)$$

Teorema 2.2.5.

Sea C un criterio T -evaluabile sobre $X \subset \mathbb{R}^k$, con $\#(X)$ finito, e IC_X , estrictamente C_X -creciente.

Entonces una funcion

$$f : X \rightarrow \mathbb{R}$$

es estrictamente C_X -creciente sii:

$$\exists \emptyset : IC_X(X) \rightarrow \mathbb{R}, \text{ estrictamente creciente / } \emptyset \cdot IC_X = f$$

Demostracion:

Si $f = \emptyset \cdot IC_X$, con \emptyset estrictamente creciente, entonces

$$x +C_X x' \Leftrightarrow IC_X(x) > IC_X(x') \Leftrightarrow \emptyset \cdot IC_X(x) > \emptyset \cdot IC_X(x')$$

Reciprocamente, dada $f : X \rightarrow \mathbb{R}$, estrictamente C_X -creciente se define

$$\emptyset : IC_X(X) \rightarrow \mathbb{R}$$

$$\emptyset[IC_X(x)] = f(x),$$

la cual es estrictamente creciente, ya que

$$\begin{aligned} IC_X(x) > IC_X(x') &\Leftrightarrow f(x) > f(x') \Leftrightarrow \\ &\Leftrightarrow \emptyset[IC_X(x)] > \emptyset[IC_X(x')] \end{aligned}$$

Teorema 2.2.6.

Si C es un criterio T-evaluable sobre $X \subset \mathbb{R}^k$, con X numerable, existe una función estrictamente C_X -creciente.

Demostración:

$$\text{Sea } r_{ij} = \begin{cases} 1 & \text{si } x_i +_{C_X} x_j \\ 0 & \text{en otro caso} \end{cases} \quad \text{con } i, j \in \mathbb{N}$$

y considerese $IC_X(x_i) = \sum_{j=1}^{\infty} (1/2)^j r_{ij}$

de donde se deduce que si $x_i, x_k \in X$

$$x_i +_{C_X} x_k \Leftrightarrow IC_X(x_i) > IC_X(x_k)$$

Teorema 2.2.7.

Sea C un criterio T-evaluable sobre $X \subset \mathbb{R}^k$, con X numerable, e IC_X estrictamente C_X -creciente.

Entonces una función

$$f : X \rightarrow \mathbb{R}, \text{ es estrictamente}$$

C_X -creciente, si:

$$\exists \emptyset : IC_X(X) \rightarrow \mathbb{R}, \text{ estrictamente creciente / } \emptyset \cdot IC_X = f$$

La demostración de este resultado es análoga a la del teorema 2.2.5.

Corolario 2.2.6.

Sea C un criterio T -evaluabile sobre $X \subseteq \mathbb{R}^k$.

Si X/\sim_{C_X} es numerable :

$$\exists IC_X : X \rightarrow \mathbb{R} ,$$

estrictamente C_X -creciente.

En el caso en que se desee analizar la existencia de una funcion estrictamente C_X -decreciente, los resultados obtenidos son igualmente validos , ya que si IC_X es estrictamente C_X -creciente se tendra que $-IC_X$ es estrictamente C_X -decreciente.

En algunos casos, puede ser de interes el estudio de funciones bajo condiciones mas debiles, por lo que se da la siguiente

Definicion 2.2.10.

Sea C un criterio sobre $X \subseteq \mathbb{R}^k$, que define una relacion binaria C_X . Se dira que una funcion

$$IC_X : X \rightarrow \mathbb{R} , \text{ es :}$$

- i) C_X -creciente si $x C_X y \implies IC_X(x) \geq IC_X(y)$
- ii) C_X -decreciente si $x C_X y \implies IC_X(x) \leq IC_X(y)$
- iii) C_X -orden convexa si
 $x C_X z C_X y \implies IC_X(z) \geq \min\{IC_X(x), IC_X(y)\}$

La existencia de los distintos tipos de funciones introducidas en la definicion 2.2.10 esta garantizada, ya

que las funciones constantes son C_X -crecientes (decrecientes) y C_X -orden convexas, siendo de interes los siguientes resultados

Proposicion 2.2.3.

Dados un criterio T-evaluable sobre $X \subseteq R^k$, y $\{ f_j \}_{j \in J}$, una familia de funciones reales definidas sobre X , sea $f = \inf_{j \in J} \{ f_j \}$

Entonces si :

- i) f_j es C_X -creciente (C_X -decreciente) $\forall j \in J \implies$
 f es C_X -creciente (C_X -decreciente)
- ii) f_j es C_X -orden convexa $\forall j \in J \implies f$ es C_X -orden convexa.

Definicion 2.2.11.

Sea C un criterio T-evaluable sobre $X \subseteq R^k$ y $f : X \rightarrow R$.

i) Se dira que $f^* : X \rightarrow R$ es el cierre C_X -creciente de f y se representara por $CC_X(f)$, si verifica:

- C1 . Es C_X -creciente
- C2 . $f^*(x) \geq f(x) \quad \forall x \in X$
- C3 . $\forall f'$, C_X -creciente / cumple C1,C2

$$f'(x) \geq f^*(x) \quad \forall x \in X$$

ii) Se dira que $f_* : X \rightarrow R$ es el cierre C_X -decreciente de f y se representara por $DC_X(f)$, si verifica:

D1 . Es C_X -decreciente

D2 . $f_*(x) \geq f(x) \quad \forall x \in X$

D3 . $\forall f''$, C_X -decreciente / cumple D1,D2

$$f''(x) \geq f_*(x) \quad \forall x \in X$$

iii) Se dira que $\Psi : X \rightarrow R$ es el cierre C_X -orden convexo de f , y se representara por $0cC_X(f)$, si verifica:

OC1 . Es C_X -o.convexa

OC2 . $\Psi(x) \geq f(x) \quad \forall x \in X$

OC3 . $\forall \Psi'$, C_X -o.convexa / cumple OC1,OC2

$$\Psi'(x) \geq \Psi(x) \quad \forall x \in X$$

Proposicion 2.2.4.

Dados un criterio T-evaluable sobre $X \subseteq R^k$ y una funcion

$f : X \rightarrow R$, se verifica:

$$i) CC_X(f)(x) = \sup \{ f(y) / x C_X y \} = \sup \{ f(y) / y \in dC_X(x) \}$$

$$ii) DC_X(f)(x) = \sup \{ f(y) / y C_X x \} = \sup \{ f(y) / y \in DC_X(x) \}$$

$$iii) 0cC_X(f)(x) = \sup \{ \min \{ f(x_1), f(x_2) \} / x_2 C_X x C_X x_1 \} \\ = \sup \{ \min \{ f(x_1), f(x_2) \} / (x_1, x_2) \in DC_X(x) \times dC_X(x) \}$$

Demostracion.

i) Sea $f_0 : X \rightarrow R$, con

$$f_0(x) = \sup \{ f(y) / x C_X y \}$$

Como se verifica $x C_X x \quad \forall x \implies f(x) \leq f_0(x)$
 y como ademas

$$x_2 C_X x_1 \implies dC_X(x_1) \subseteq dC_X(x_2) \implies f_0(x_1) \leq f_0(x_2)$$

se concluye que f_0 es C_X -creciente.

Ademas para cualquier $f' : X \rightarrow R$, C_X -creciente
 con $f(x) \leq f'(x) \quad \forall x$, se verifica :

$$x C_X y \implies f'(x) \geq f'(y) \geq f(y)$$

de donde se deduce que

$$f_0 = CC_X(f)$$

ii) Se demostraria de forma analoga.

iii) Sea $f_1 : X \rightarrow R$, con

$$f_1(x) = \sup \{ \min \{ f(x_1), f(x_2) \} / x_2 C_X x C_X x_1 \}$$

Evidentemente $f_1(x) \geq f(x) \quad \forall x$.

Sea $\theta \in R$, tal que

$$\theta < \min \{ f_1(x), f_1(y) \} \implies \exists x_1, x_2 \text{ tales que}$$

$$x_2 C_X x C_X x_1 \quad \text{con} \quad \min \{ f(x_1), f(x_2) \} > \theta$$

De forma analoga $\exists y_1, y_2$ tales que

$$y_2 C_X y C_X y_1 \quad \text{con} \quad \min \{ f(y_1), f(y_2) \} > \theta$$

Por tanto para

$$x, y, z \in X / y C_X z C_X x$$

se verifica

$$y_2 C_X y C_X z C_X x C_X x_1 \implies$$

$$f_1(z) \geq \min \{ f(x_1), f(y_2) \} > \theta$$

para cualquier $\theta / \theta < \min \{ f_1(x), f_1(y) \}$

obteniendo que

$$f_1(z) \geq \min \{ f_1(x), f_1(y) \} \implies f_1 \text{ es } C_X\text{-o.convexa.}$$

Para cualquier f'' C_X -o.convexa, verificando

$$f''(x) \geq f(x) \quad \forall x \in X,$$

$$\begin{aligned} x_2 C_X x C_X x_1 &\implies f''(x) \geq \min \{ f''(x_1), f''(x_2) \} \geq \\ &\geq \min \{ f(x_1), f(x_2) \} \implies f''(x) \geq f_1(x) \quad \forall x \in X, \end{aligned}$$

concluyendo pues que $f_1 = OcC_X(f)$.

Lema 2.2.3.

Dado un criterio T-evaluable, C , sobre $X \subseteq \mathbb{R}^k$, y f_1, f_2 , tales que f_1 es C_X -creciente y f_2 es C_X -decreciente, se verifica:

$$f = \min \{ f_1, f_2 \} \text{ es } C_X\text{-o.convexa.}$$

Demostración:

Si $x C_X z C_X y$ entonces

$$f_1(z) \geq f_1(y) \geq f(y) \quad ; \quad f_2(z) \geq f_2(x) \geq f(x)$$

y por tanto

$$f(z) \geq \min \{ f(x), f(y) \} \implies$$

es C_X -o.convexa.

Teorema 2.2.8.

$$OcC_X(f) = \min \{ CC_X(f), dC_X(f) \}$$

Demostración:

$\Psi = \min \{ CC_X(f), dC_X(f) \}$ es C_X -o.convexa por lema 2.2.3 y es evidente que $f(x) \leq \Psi(x) \quad \forall x \in X \implies$

$$\Rightarrow \text{Oc}C_X(f)(x) \leq \Psi(x) \quad \forall x \in X .$$

$$\text{Si } \exists x \in X / \text{Oc}C_X(f)(x) \neq \Psi(x) \Rightarrow$$

$$\Rightarrow \exists \theta / \text{Oc}C_X(f)(x) < \theta < \Psi(x) \leq \text{CC}_X(f)(x) \Rightarrow$$

$$\Rightarrow \exists y_1 / x C_X y_1 \quad \text{con } f(y_1) > \theta$$

y como $\Psi(x) \leq \text{d}C_X(f)(x) \Rightarrow \exists y_2 / y_2 C_X x$ con $f(y_2) > \theta$, se deduce :

$$y_2 C_X x C_X y_1 \Rightarrow \text{Oc}C_X(f)(x) \geq \min \{ f(y_2), f(y_1) \} > \theta$$

por lo que $\text{Oc}C_X(f) = \Psi$.

Teorema 2.2.9.

Sea C un criterio T-evaluable sobre $X \subseteq R^k$ y sea

$$f : X \rightarrow R .$$

Entonces f es C_X -o.convexa sii existen dos funciones f_1 , C_X -creciente , y f_2 , C_X -decreciente de forma que :

$$f = \min \{ f_1 , f_2 \} .$$

Demostracion:

$$\text{Si } f \text{ es } C_X\text{-o.convexa} \Rightarrow f = \text{Oc}C_X(f) \Rightarrow$$

$$\Rightarrow f = \min \{ \text{CC}_X(f), \text{d}C_X(f) \}$$

Reciprocamente , si $f = \min \{ f_1, f_2 \}$, con f_1 C_X -creciente y f_2 C_X -decreciente, del lema 2.2.3 , se deduce que f es C_X -o.convexa.

Para el analisis de las situaciones experimentales , el estadístico dispone de un conjunto finito de observaciones $X \subseteq R^k$, con $\#(X) = n$.

En este caso , la existencia de una funcion

$$IC_X : X \rightarrow R$$

verificando la condicion

$$(x,y) + C'_X(z,w) \Leftrightarrow IC_X(x) - IC_X(y) < IC_X(z) - IC_X(w)$$

viene determinada por los teoremas 2.2.1 y 2.2.3 .

En cualquier caso siempre existiran C_X -crecientes , C_X -decrecientes y C_X -o.convexas , y la existencia de funciones estrictamente C_X -crecientes (decrecientes) , queda garantizada por el teorema 2.2.4 , en el caso en que el criterio C sea T -evaluabile.

Sin embargo , debido a la falta de unicidad de estas funciones, sera necesario utilizar algun procedimiento de seleccion entre las mismas, mediante el cual se refleje como se desvian las observaciones del comportamiento general de la masa de datos, respecto del criterio bajo estudio.

El procedimiento de seleccion que se utiliza en el Capitulo III , viene expresado a traves de la optimizacion de un funcional real, que depende de las observaciones y del conjunto de funciones que se desea analizar.

Asi , si se denota por $\mathfrak{X} = \{ X \subset R^k / \#(X) = n \}$,

$$F : \mathfrak{X} \times \mathfrak{M} \rightarrow R$$

$$F(X, \mu) = F(x_1, \dots, x_n, \mu(x_1), \dots, \mu(x_n))$$

donde \mathfrak{M} representa al conjunto de funciones que se desea analizar.

2.3 INVARIANZA.

Admitiendo el caracter subjetivo del criterio C , y del funcional F , seria deseable que las conclusiones o resultados que se obtengan no dependan del sistema de coordenadas elegido para representar los elementos de \mathfrak{X} , por lo que a continuacion se estudia este problema.

En el estudio que se realiza, se supone que g es una funcion

$$g : R^k \rightarrow R^k \quad \text{uno a uno}$$

En estas condiciones se dan las siguientes

Definicion 2.3.1.

Un criterio C , T-evaluable sobre los elementos de \mathfrak{X} , se dira que es g -invariante si verifica

$$x C_X x' \Leftrightarrow g(x) C_{g(X)} g(x') \quad \forall x, x' \in X, \forall X \in \mathfrak{X}$$

Se denotara por

$$I_{\mathfrak{X}}(C) = \{ g / C \text{ es } g\text{-invariante} \}$$

Definicion 2.3.2.

Sea C un criterio T-evaluable sobre los elementos de \mathfrak{X} , y $SSC(\mathfrak{X}) = \{ SSC_X / X \in \mathfrak{X} \}$, donde SSC_X es una funcion de C-similitud signada sobre X .

$SSC(\mathfrak{X})$ es g -invariante, si

$$(x, y) C'_X (z, w) \Leftrightarrow (g(x), g(y)) C'_{g(X)} (g(z), g(w))$$

para cualesquiera $x, y, z, w \in X, \forall X \in \mathfrak{X}$.

Se representara por $\text{ISSC}(\mathfrak{X})$ al conjunto

$$\text{ISSC}(\mathfrak{X}) = \{ g / \text{SSC}(\mathfrak{X}) \text{ es } g\text{-invariante} \}$$

Proposicion 2.3.1.

$$\text{Si } g \in \text{ISSC}(\mathfrak{X}) \implies g \in I_{\mathfrak{X}}(C)$$

Demostracion:

$$\begin{aligned} x C_X y &\iff (x, x) C'_X(x, y) \iff \\ (g(x), g(x)) C'_{g(X)}(g(x), g(y)) &\iff g(x) C_{g(X)} g(y) \end{aligned}$$

Proposicion 2.3.2.

Si C es un criterio T-evaluable sobre los elementos de \mathfrak{X} y $\varphi \subset I_{\mathfrak{X}}(C)$, siempre existe un grupo G , cuya ley interna es la composicion de funciones, verificando

$$\varphi \subset G \subset I_{\mathfrak{X}}(C)$$

Demostracion:

En efecto, $\forall g, g' \in \varphi \implies g \cdot g', g^{-1} \in I_{\mathfrak{X}}(C)$
ya que $\forall x, x' \in X, \forall X \in \mathfrak{X}$:

$$\text{i) } x C_X x' \iff g'(x) C_{g'(X)} g'(x') \iff gg'(x) C_{gg'(X)} gg'(x')$$

$$\text{ii) } g^{-1}(x) C_{g^{-1}(X)} g^{-1}(x') \iff gg^{-1}(x) C_{gg^{-1}(X)} gg^{-1}(x') \iff$$

$$\iff x C_X x'$$

de donde se deduce de forma inmediata el resultado de la proposicion.

Proposicion 2.3.3.

Dado $SSC(\mathfrak{X})$ y Ψ un subconjunto de $ISSC(\mathfrak{X})$, existe un grupo G , cuya ley interna es la composicion de funciones, verificando

$$\Psi \subset G \subset ISSC(\mathfrak{X})$$

La demostracion es analoga a la de la proposicion 2.3.2.

Al igual que se han introducido los conceptos de de invarianza de un criterio y de una funcion de C-similitud signada, para que las relaciones de preferencia definidas no dependan del sistema de coordenadas elegido, tambien es deseable que las funciones IC_X , que caracterizan a las distintas relaciones, fuesen tambien independientes del sistema de coordenadas.

Con tal objetivo se introduce:

Definicion 2.3.3.

Sea C un criterio T -evaluabile sobre los elementos de \mathfrak{X} , e $IC(\mathfrak{X}) = \{ IC_X / X \in \mathfrak{X} \}$ donde

$$IC_X : X \rightarrow R$$

Se dira que $IC(\mathfrak{X})$ es g -invariante si verifica :

$$IC_X(x) = IC_{g(X)}(g(x)) \quad \forall x \in X, \quad \forall X \in \mathfrak{X}.$$

Se denotara por $IIC(\mathfrak{X})$ al conjunto de funciones

$$IIC(\mathfrak{X}) = \{ g / IC(\mathfrak{X}) \text{ es } g\text{-invariante} \}$$

Proposición 2.3.4.

Si C es un criterio T-evaluable sobre los elementos de \mathfrak{X} y $\Psi \subset \text{IIC}(\mathfrak{X})$, existe un grupo de transformaciones G , cuya ley interna es la composición de funciones, que verifica:

$$\Psi \subset G \subset \text{IIC}(\mathfrak{X})$$

Proposición 2.3.5.

Sea $\text{IC}(\mathfrak{X}) = \{ \text{IC}_X / X \in \mathfrak{X} \}$, de forma que IC_X es estrictamente C_X -creciente (C_X -decreciente).

Entonces:

$$g \in \text{IIC}(\mathfrak{X}) \implies g \in I_{\mathfrak{X}}(C).$$

Demostración:

$$\begin{aligned} x C_X x' &\iff \text{IC}_X(x) \geq \text{IC}_X(x') \iff \\ &\iff \text{IC}_{g(X)}(g(x)) \geq \text{IC}_{g(X)}(g(x')) \iff g(x) C_{g(X)} g(x') \end{aligned}$$

Proposición 2.3.6.

Dado $\text{SSC}(\mathfrak{X})$, sea $\text{IC}(\mathfrak{X})$ verificando:

$$(x, y) C'_X (z, w) \iff \text{IC}_X(x) - \text{IC}_X(y) < \text{IC}_X(z) - \text{IC}_X(w)$$

$$\forall x, y, z, w \in X ; \forall X \in \mathfrak{X}$$

$$\text{Entonces: } g \in \text{IIC}(\mathfrak{X}) \implies g \in \text{ISSC}(\mathfrak{X})$$

Ya que las funciones IC_X , que se utilizarán en el capítulo siguiente se obtienen mediante la optimización de un funcional F , es lógico estudiar bajo que condiciones dichas funciones son invariantes.

El teorema esta referido al caso de funciones estrictamente C_X -crecientes pero en el caso de estrictamente C_X -decrecientes se obtiene un resultado analogo.

Teorema 2.3.1.

Sea un criterio C , T -evaluabile sobre los elementos de \mathfrak{E} , y $g \in I_{\mathfrak{E}}(C)$, de forma que :

$F(g(X), \mu) = g^* F(X, \mu)$ con g^* estrictamente creciente.

Entonces $IC^*(\mathfrak{E}) = \{ IC_X^* / X \in \mathfrak{E} \}$ es g -invariante si IC_X^* se determina a traves de la maximizacion (minimizacion) del funcional F , sobre la clase de funciones estrictamente C_X -crecientes.

Demostracion:

Si se denota por CC_X a la clase de funciones estrictamente C_X -crecientes, se tiene

$$\left. \begin{array}{l} \text{Max } F(g(X), \mu) \\ \mu \in CC_{g(X)} \end{array} \right\} \equiv \begin{array}{l} \text{Max } F(g(x_1), \dots, g(x_n), \mu_1, \dots, \mu_n) \\ \text{sujeto a} \\ g(x_i) C_{g(X)} g(x_j) \Leftrightarrow \mu_i \geq \mu_j \end{array}$$

que debido a las hipotesis es equivalente a

$$\left. \begin{array}{l} \text{Max } g^* F(x_1, \dots, x_n, \mu_1, \dots, \mu_n) \\ \text{sujeto a} \\ x_i C_X x_j \Leftrightarrow \mu_i \geq \mu_j \end{array} \right\} \equiv \begin{array}{l} \text{Max } F(X, \mu) \\ \mu \in CC_X \end{array}$$

Teorema 2.3.2.

Sea C un criterio T-evaluable sobre los elementos de \mathfrak{X} , y $g \in \text{ISSC}(\mathfrak{X})$ de forma que

$$F(g(X), \mu) = g^* F(X, \mu)$$

con g^* estrictamente creciente.

Entonces $g \in \text{IIC}^*(\mathfrak{X})$, si IC_X^* se determina a través de

$$\text{Max } F(x_1, \dots, x_n, \mu_1, \dots, \mu_n)$$

sujeto a

$$(x_i, x_j) + C'_X(x_h, x_k) \Leftrightarrow \mu_i - \mu_j < \mu_h - \mu_k$$

La demostración de este resultado es analoga a la del Teorema 2.3.1

CAPITULO III

ANALISIS CUALITATIVO DE DATOS .

TECNICAS DE IDENTIFICACION DE OUTLIERS .

-- Introduccion

3.1 Estructura Basica. Aproximacion General.

3.2 Criterio de Dispersion Central

3.2.1 Outliers

3.2.2 Caso de Parametros Poblacionales Conocidos

3.3 Criterio de Dispersion Basado en el Recorrido

3.3.1 Outliers

3.3.2 Caso de Parametros Conocidos

3.4 Criterio Natural (Origen Conocido)

3.4.1 Outliers

3.5 Generalizacion al Caso de mas de una Observacion

3.5.1 Criterio de Dispersion Central

3.5.2 Criterio de Dispersion Basado en el Recorrido

3.5.3 Criterio Natural

3.6 Criterio de Dispersion Central Multivariante

INTRODUCCION

En el capitulo I se ha destacado la problematica que presentan, actualmente, las Tecnicas de Identificacion de Outliers.

En este capitulo, se aplican los resultados obtenidos en el capitulo anterior, con el objeto de obtener algunos de los estadisticos que son de especial interes dentro de las tecnicas de deteccion de outliers.

En el apartado de conclusiones, se refleja como con el planteamiento propuesto se corrigen algunos de los defectos que hoy en dia aun persisten en la teoria de Outliers.

3.1 ESTRUCTURA BASICA. APROXIMACION GENERAL.

Sea un conjunto $X \subset \mathbb{R}^k$ / $\#(X) = n$, en el que se desea comparar sus elementos a traves de un criterio C , el cual determina una relacion binaria C_X definida por:

$$x C_X x' \iff C_X(x) \leq C_X(x') \quad (3.1)$$

con $C_X : X \rightarrow \mathbb{R}$.

Esta relacion binaria da lugar a los conjuntos C_X -dominantes:

$$DC_X(x) = \{ x' \in X / C_X(x') \leq C_X(x) \} \quad x \in X$$

y C_X -dominados:

$$dC_X(x) = \{ x' \in X / C_X(x) \leq C_X(x') \} \quad x \in X$$

introducidos en la definicion 2.1.4.

Cualquier subconjunto $X' \subseteq X$, estara propiamente C_X -dominado , y los elementos $(C_X D m)$ y $(C_X d M)$ de X' , que se representaran por x^* y x_* , respectivamente, son aquellos que verifican :

$$\text{Min}_{x \in X'} C_X(x) = C_X(x^*)$$

$$\text{Max}_{x \in X'} C_X(x) = C_X(x_*)$$

Ademas el criterio C es T-evaluable sobre X , y las clases de equivalencia estan formadas por aquellos elementos sobre los que la funcion C_X es constante.

De forma natural; la relacion binaria dada en (3.1) , da lugar a las funciones:

$$\exists C_X : X \times X \rightarrow R \quad (3.2)$$

$$SC_X(x, x') = |C_X(x') - C_X(x)|$$

y

$$SSC_X : X \times X \rightarrow R \quad (3.3)$$

$$SSC_X(x, x') = C_X(x') - C_X(x)$$

Proposición 3.1.1.

Las funciones SC_X y SSC_X , dadas en (3.2) y (3.3), son funciones de C-similitud y C-similitud signada sobre X, respectivamente, que son aditivas, en el sentido de la definición 2.2.7, a las que se denominaran funciones de C-similitud y C-similitud signada naturales.

La demostración de esta proposición es evidente a partir de las definiciones dadas de SC_X y SSC_X .

De forma analoga a (3.2) y (3.3), pueden definirse las funciones de C-similitud y C-similitud signada, tipificadas como:

$$SC_X(x, x') = [|C_X(x') - C_X(x)|] / [C_X(x_*) - C_X(x^*)]$$

$$SSC_X(x, x') = [C_X(x') - C_X(x)] / [C(x_*) - C(x^*)]$$

La función de C-similitud tipificada puede interpretarse como una relación difusa definida en $X \times X$, a través del criterio C.

Proposición 3.1.2.

Dada la función de C-similitud signada natural, existe

$$IC_X : X \rightarrow R$$

verificando:

$$(x, y) C'_X (z, w) \Leftrightarrow IC_X(y) - IC_X(x) \leq IC_X(w) - IC_X(z) \quad (3.4)$$

siendo C'_X , la relación introducida en la definición 2.2.4.

El resultado de esta proposición es inmediato, ya que sería suficiente considerar $IC_X = C_X$.

Proposición 3.1.3.

Para la función de C-similitud signada natural, se verifica:

$\forall x_i, x'_i \in X, i=1,2,3$, tales que

$$\left. \begin{array}{l} (x_1, x_2) C'_X (x'_1, x'_2) \\ (x_2, x_3) C'_X (x'_2, x'_3) \end{array} \right\} \Rightarrow (x_1, x_3) C'_X (x'_1, x'_3)$$

que es la condición T1 del Teorema 2.2.1.

Sin embargo:

Dados $y, x_i \in X, i=1,2,3,4$, tales que

$$(x_1, x_2) C'_X (x_3, x_4) ; (x_3, x_4) C'_X (y, y)$$

no está garantizada la existencia de $z_1, z_2 \in X$, cumpliendo:

$$(x_1, z_1) \sim C'_X (x_3, x_4) ; (z_2, x_2) \sim C'_X (x_3, x_4)$$

por tanto, la condición T2 del Teorema 2.2.1, no siempre se verifica, por lo que no queda garantizada la uni-

cidad salvo transformaciones lineales de IC_X .

Debido a la falta de unicidad, es necesario considerar algun criterio de optimalidad que permita seleccionar entre las diversas funciones que verifican (3.4).

En los casos que se analizan posteriormente, se considera como funcion IC_X optima aquella que:

$$\text{Min } \sum_{i=1}^n \{ [IC_X(x_i)]^2 [KC_X - C_X(x_i)] + [1 - IC_X(x_i)]^2 C_X(x_i) \}$$

sujeto a

$$(x, y) C'_X(z, w) \Leftrightarrow IC_X(y) - IC_X(x) \leq IC_X(w) - IC_X(z)$$

siendo KC_X , una constante positiva que depende del conjunto X y del criterio C , que se consideran.

Este criterio de optimalidad propuesto puede considerarse como una modificacion del dado por Dunn (1974).

El optimo IC_X^* , sin restricciones, viene dado por:

$$IC_X^*(x) = C_X(x) / KC_X \quad x \in X$$

el cual verifica la restriccion impuesta en (3.4). En efecto:

$$(x, y) C'_X(z, w) \Leftrightarrow C_X(y) - C_X(x) \leq C_X(w) - C_X(z) \Leftrightarrow \\ IC_X^*(y) - IC_X^*(x) \leq IC_X^*(w) - IC_X^*(z)$$

En el caso en que $IC_X^*(x) \in [0, 1]$, esta funcion puede interpretarse como un conjunto difuso definido sobre X a traves del criterio C .

Aunque de forma general, no esta garantizada la unicidad (salvo transformaciones lineales) de la funcion IC_X , existen situaciones en las que se verifica la condicion T2 del Teorema 2.2.1, por ejemplo, cuando la funcion C_X , es tal que :

$$C_X(X) = \{ 0, 1, \dots, r ; r \leq n-1 \}$$

o de forma mas general

$$C_X(X) = \{ \delta \cdot i + K ; i = 0, \dots, r ; r \leq n-1 \}$$

en cuyo caso queda garantizada la unicidad salvo transformaciones lineales de IC_X .

Es de destacar que esta funcion no seria valida para los propositos de la memoria, ya que en esta situacion no seria posible determinar que elementos se desvian marcadamente del comportamiento general de la masa de datos, ya que los rangos son robustos frente a la presencia de outliers.

A continuacion, se van a obtener mediante la aplicacion del procedimiento estudiado, algunos estadisticos de gran relevancia dentro de las actuales tecnicas de identificacion de outliers.

3.2 CRITERIO DE DISPERSION CENTRAL

Sea $X \subset \mathbb{R}$ / $\#(X) = n$, sobre el que se desea analizar el comportamiento de sus elementos respecto del criterio C , que define la relacion binaria :

$x C_X x'$ sii x esta al menos tan proximo a \bar{x} como x' .

Esta relacion puede expresarse mediante

$$x C_X x' \iff (x - \bar{x})^2 \leq (x' - \bar{x})^2 \quad (3.5)$$

o bien por

$$x C_X x' \iff |x - \bar{x}| \leq |x' - \bar{x}| \quad (3.6)$$

Evidentemente, (3.5) y (3.6) definen la misma relacion binaria sobre el conjunto X .

En el analisis que se realiza a continuacion se trabaja con la forma (3.5).

En este caso los conjuntos C_X -dominantes y C_X -dominados son :

$$DC_X(x) = \{ x' \in X / (x' - \bar{x})^2 \leq (x - \bar{x})^2 \}$$

$$dC_X(x) = \{ x' \in X / (x - \bar{x})^2 \leq (x' - \bar{x})^2 \}$$

respectivamente.

Como se indico en el epigrafe 3.1, todo $X' \subseteq X$ esta propiamente C_X -dominado, y ademas, en este caso, x_* , el elemento (C_X d M) de X , verifica:

$$(x_* - \bar{x})^2 = \text{Max}_{x \in X} \{(x - \bar{x})^2\} = \text{Max} \{(x_{(1)} - \bar{x})^2, (x_{(n)} - \bar{x})^2\}$$

donde $x_{(1)}$, $x_{(n)}$, son, respectivamente, la menor y la mayor de las observaciones, por lo que el elemento (C_X d M) de X, sera $x_{(1)}$ ó $x_{(n)}$.

Como ya se indico en epigrafe 3.1, el criterio C es T-evaluable sobre X, y en este caso particular, las clases de equivalencia estan formadas por elementos equidistantes a la media.

Para esta situacion, las funciones de C-similitud y C-similitud signada, naturales, introducidas mediante (3.2) y (3.3) vienen dadas por:

$$SC_X(x, x') = |(x - \bar{x})^2 - (x' - \bar{x})^2|$$

$$SSC_X(x, x') = (x' - \bar{x})^2 - (x - \bar{x})^2$$

Como quedo reflejado en el epigrafe anterior, en general la funcion de similitud signada natural, no verifica la condicion T2 del Teorema 2.2.1, por lo que no quedaria garantizada la unicidad de IC_X .

El optimo se seleccionara mediante:

$$\text{Min } \sum_{i=1}^n \{ [IC(x_i)]^2 (nS_X^2 - (n-1)S_i^2) + [1-IC_X(x_i)]^2 (n-1)S_i^2 \} \quad (3.7)$$

sujeto a

$$(x_j, x_k) C'_X(x_1, x_m) \Leftrightarrow IC_X(x_j) - IC_X(x_k) \leq IC_X(x_1) - IC_X(x_m)$$

con

$$nS_X^2 = \sum_{p=1}^n (x_p - \bar{x})^2 \quad ; \quad (n-1)S_i^2 = \sum_{\substack{p=1 \\ p \neq i}}^n (x_p - \bar{x}_i)^2$$

$$y \quad \bar{x}_1 = \left[\sum_{\substack{p=1 \\ p \neq i}}^n x_p \right] / (n-1)$$

Debido a que la función de C-similitud signada natural es invariante por traslación y escala, y además

$$F(\cdot, a x_1 + b, \dots, IC_X(a x_1 + b), \dots) = a^2 F(\cdot, x_1, \dots, IC_X(x_1), \dots)$$

siendo F el funcional a optimizar en (3.7), a partir del Teorema 2.3.2, se concluye que el óptimo del funcional es invariante por localización y escala.

El óptimo del funcional F , sin restricciones, viene dado por :

$$IC_X^*(x_1) = \frac{(n-1) S_1^2}{n S_X^2}$$

el cual verifica la restricción, ya que:

$$IC_X^*(x_j) - IC_X^*(x_k) \leq IC_X^*(x_1) - IC_X^*(x_m) \Leftrightarrow S_j^2 - S_k^2 \leq S_1^2 - S_m^2$$

$$\Leftrightarrow (x_k - \bar{x})^2 - (x_j - \bar{x})^2 \leq (x_m - \bar{x})^2 - (x_1 - \bar{x})^2 \Leftrightarrow$$

$$\Leftrightarrow (x_j, x_k) C'_X(x_1, x_m)$$

debido a que

$$n S_X^2 = (n-1) S_1^2 + [(n-1)/n] (x_1 - \bar{x}_1)^2$$

y

$$(x_1 - \bar{x})^2 = [(n-1)/n]^2 (x_1 - \bar{x}_1)^2$$

Se ha de notar :

I) $IC_X^*(x) \in [0, 1] \quad \forall x \in X$, lo que se puede interpre-

tar como un conjunto difuso definido sobre X, mediante el criterio C.

$$\text{II) } IC_X^*(x) = 1 \quad \text{si} \quad x = \bar{x}$$

III) En el caso particular de $n=2$

$$S_1^2 = (x_2 - \bar{x}_1)^2 = (x_2 - x_2)^2 = 0 \quad , \quad \text{y análogamente}$$

$$S_2^2 = 0 \quad , \quad \text{por lo que} \quad IC_X^*(x_i) = 0 \quad , \quad i=1,2$$

IV) Dado $\alpha \in (0,1)$

$$IC_X^*(x_j) \leq \alpha \iff S_j^2 \leq [\alpha / (1-\alpha)] [n/(n-1)^2] (x_j - \bar{x})^2$$

y en el caso particular $\alpha = 1/2$

$$IC_X^*(x_j) \leq 1/2 \iff \sum_{x \neq x_j} (x - \bar{x})^2 \leq [(n+1)/(n-1)] (x_j - \bar{x})^2$$

lo que reflejaría que dicha observación tiene una gran influencia sobre la dispersión de la masa de datos.

3.2.1 OUTLIERS.

Grubbs (1950), propone los estadísticos :

$$T_1 = [(n-1) S_{(1)}^2] / [n S_X^2]$$

$$T_2 = [(n-1) S_{(n)}^2] / [n S_X^2]$$

para detectar la presencia de un outlier inferior o superior, respectivamente, en una población normal con media y varianza desconocidas.

Es de destacar que

$$T_1 = IC_X^*(x_{(1)}) \quad \text{y} \quad T_2 = IC_X^*(x_{(n)})$$

Ademas, Tietjen y Moore (1972) proponen el estadistico

$$T_3 = \text{Min} \{ T_1, T_2 \}$$

para detectar la presencia de un outlier superior o inferior en una poblacion normal de media y varianza desconocidas, observandose que

$$T_3 = \text{Min} \{ IC_X^*(x) / x \in X \}$$

3.2.2 CASO DE PARAMETROS POBLACIONALES CONOCIDOS.

En el caso en que se dispusiese de alguna informacion a priori, sobre la poblacion de la que se ha extraido el conjunto de observaciones X , seria conveniente, hacer uso de la misma.

Asi, en el caso en que la media de la poblacion, μ_0 , sea conocida, la relacion dada por (3.5), podria modificarse, considerando:

$$x C_X x' \iff (x - \mu_0)^2 \leq (x' - \mu_0)^2$$

Es de destacar, que en este caso, la relacion binaria C_X es independiente de X .

Los resultados obtenidos en el epigrafe 3.2.1 serian serian facilmente generalizables en esta nueva situacion, obteniendose los estadisticos equivalentes a T_1 , T_2 y T_3 , correspondientes al caso de media conocida.

Es de resaltar que IC_X^* si depende del conjunto de observaciones X .

3.3 CRITERIO DE DISPERSION BASADO EN EL RECORRIDO

Sea $X \subset R / \#(X) = n$, sobre el que se desea analizar el comportamiento de sus elementos respecto del criterio C , que define la relacion binaria :

$x C_X x'$ sii el recorrido de $X - \{x\}$ es mayor ó igual que el de $X - \{x'\}$.

La relacion definida de esta forma, puede expresarse mediante:

$$x_i C_X x_j \iff R_i \geq R_j \quad (3.8)$$

donde $R_i = \text{Max}_{h \neq i} \{x_h\} - \text{Min}_{h \neq i} \{x_h\}$

Si se denota por $R_{(i)}$ al recorrido de $X - \{x_{(i)}\}$ siendo $x_{(i)}$, la i -esima menor de las observaciones, se tiene :

$$\begin{aligned} R_{(1)} &= x_{(n)} - x_{(2)} \\ R_{(i)} &= x_{(n)} - x_{(1)} \quad i=2, \dots, n-1 \\ R_{(n)} &= x_{(n-1)} - x_{(1)} \end{aligned}$$

con lo cual, la relacion definida en (3.8), da lugar a que los conjuntos C_X -dominantes y C_X -dominados, introducidos en la def. 2.1.4, verifiquen :

$$\begin{aligned} x_{(i)} \in DC_X(x_{(j)}) \quad i=2, \dots, n-1 ; j=1, \dots, n \\ dC_X(x_{(i)}) = X \quad i=2, \dots, n-1 \end{aligned}$$

Como quedo reflejado en el epigrafe 3.1, todo subconjunto $X' \subseteq X$, esta propiamente C_X -dominado, y x_* , elemento (C_X d M) de X , verifica :

$$R_{x_*} = \text{Min} \{ R_i, i=1, \dots, n \} = \text{Min} \{ R_{(1)}, R_{(n)} \}$$

por lo que x_* sera $x_{(1)}$ ó $x_{(n)}$, dependiendo de los valores de $R_{(1)}, R_{(n)}$.

El criterio C , es T-evaluable sobre X , y además, elementos $x_{(i)}, i = 2, \dots, n-1$, son C_X -equivalentes.

Por tanto el criterio C , efectua una particion de X , en a lo sumo tres clases de equivalencia.

Para este criterio, las funciones de C-similitud y C-similitud signada, naturales, introducidas mediante las expresiones (3.2) y (3.3), vienen dadas por :

$$SC_X : X \times X \rightarrow R$$

$$SC_X(x_i, x_j) = |R_i - R_j|$$

y

$$SSC_X : X \times X \rightarrow R$$

$$SSC_X(x_i, x_j) = R_i - R_j$$

por lo que :

$$SC_X(x_{(i)}, x_{(j)}) = \begin{cases} 0 & \text{si } 1 < i, j < n \text{ ó } i=j \\ x_{(2)} - x_{(1)} & \text{si } 1 < i < n ; j=1 \\ x_{(n)} - x_{(n-1)} & \text{si } 1 < i < n ; j=n \\ |x_{(n)} - x_{(n-1)} - x_{(2)} + x_{(1)}| & \text{si } i=1 ; j=n \end{cases}$$

Analogamente, se expresa la funcion de C-similitud signada natural.

De forma general, se puede afirmar que la funcion de similitud signada natural, no verifica la condicion T2

del teorema 2.2.1 ,cumplendose en el caso particular en que

$$R_{(1)} = 2 R_{(n)} - R \quad \text{o} \quad R_{(n)} = 2 R_{(1)} - R$$

Debido pues a que no esta garantizada la unicidad de la funcion IC_X , se seleccionara el optimo a partir de:

$$\text{Min} \quad \sum_{i=1}^n \{ [IC_X(x_i)]^2 [R - R_i] + [1-IC_X(x_i)]^2 R_i \} -$$

sujeto a

$$(x_j, x_k) C'_X(x_1, x_m) \Leftrightarrow IC_X(x_j) - IC_X(x_k) \leq IC_X(x_1) - IC_X(x_m)$$

El optimo de este problema ,sera invariante por localizacion y escala, ya que la funcion de C-similitud signada natural lo es , y ademas

$$R_{aX+b} = a R_X \quad a > 0$$

por lo que la invarianza se deduce partir del Teorema 2.3.2.

El optimo del problema anterior , viene dado por :

$$IC_X^*(x_j) = \frac{R_j}{R} = \frac{R_j}{x_{(n)} - x_{(1)}}$$

el cual verifica la restriccion , como se deduce de forma evidente .

Se ha de notar :

$$(I) \quad IC_X^*(x) \in [0,1] \quad \forall x \in X$$

por lo que se puede interpretar como un conjunto difuso definido sobre X , a traves del criterio C.

(II) Para $n > 2$

$$IC_X^*(x_{(i)}) = 1 \quad i = 2, \dots, n-1$$

En el caso particular $n = 2$:

$$IC_X^*(x_i) = 0 \quad i = 1, 2$$

3.3.1 OUTLIERS.

(I)

$$IC_X^*(x_{(1)}) = \frac{x_{(n)} - x_{(2)}}{x_{(n)} - x_{(1)}} \quad \text{por lo que}$$

$1 - IC_X^*(x_{(1)})$ coincide con el estadístico dado por Dixon (1950, 1951), para testar la presencia de un outlier inferior en una muestra extraída de una población exponencial de origen desconocido.

(II)

$$IC_X^*(x_{(n)}) = \frac{x_{(n-1)} - x_{(1)}}{x_{(n)} - x_{(1)}} \quad \text{por lo que}$$

$1 - IC_X^*(x_{(n)})$ coincide con el estadístico dado por Dixon (1950, 1951), para testar la presencia de un outlier superior en la misma situación de (I) y para testar la presencia de un outlier superior en una muestra extraída de una población normal de varianza desconocida.

(III) $\text{Max} \{ 1 - IC_X^*(x_{(1)}), 1 - IC_X^*(x_{(n)}) \}$ coincide con el estadístico propuesto por King (1953), para detectar la presencia de un extremo outlier en una

muestra extraída de una población normal de varianza desconocida.

(IV) Los estadísticos recogidos en (I) , (II) y (III) se utilizan también para testar la presencia de outliers en muestras extraídas de poblaciones uniformes .

3.3.2 CASO DE PARAMETROS CONOCIDOS.

Si el origen de las observaciones fuese un valor \underline{a} conocido , y se considerase :

$$R = x_{(n)} - a$$

$$R_{(i)} = x_{(n)} - a \quad i = 1, \dots, n-1$$

$$R_{(n)} = x_{(n-1)} - a$$

se puede realizar un análisis similar al efectuado anteriormente , llegándose a la conclusión :

$$IC_X^*(x_{(i)}) = 1 \quad i = 1, \dots, n-1$$

$$IC_X^*(x_{(n)}) = \frac{x_{(n-1)} - a}{x_{(n)} - a}$$

y por tanto $1 - IC_X^*(x_{(n)})$, coincide con el estadístico utilizado para testar la presencia de un outlier superior en una muestra extraída de una población exponencial de origen \underline{a} conocido .

3.4 CRITERIO NATURAL (ORIGEN CONOCIDO)

Sea un conjunto $X \subset \mathbb{R}$ / $\#(X) = n$, $x > a$, $\forall x \in X$ sobre el que se desea analizar el comportamiento de sus elementos respecto de un criterio C , que define sobre X la relacion binaria determinada por la ordenacion natural de las observaciones :

$$x C_X x' \iff x \leq x' \quad (x C_X x' \iff x-a \leq x'-a)$$

Es de destacar, que la relacion binaria expresada de esta forma, es independiente de X y da lugar a los conjuntos C_X -dominantes y C_X -dominados :

$$DC_X(x) = \{ x' \in X / x' \leq x \}$$

$$dC_X(x) = \{ x' \in X / x \leq x' \}$$

y ademas, los elementos $(C_X D m)$ y $(C_X d M)$ de X , son $x_{(1)}$ y $x_{(n)}$, respectivamente.

Evidentemente, el criterio C es T-evaluable sobre X y las clases de equivalencia estan formadas por elementos identicos.

En este caso, las funciones de C-similitud y C-similitud signada naturales, son :

$$SC_X(x, x') = x - x'$$

$$SSC_X(x, x') = |x' - x|$$

De forma general, la funcion de C-similitud signada natural no verifica la condicion T2 del teorema 2.2.1, sin embargo, dicha condicion si se verificara para el caso en que :

$$x_{(i)} = x_{(1)} + (i-1)\delta, \quad \delta > 0$$

Si se selecciona el optimo a partir de :

$$\text{Min} \sum_{i=1}^n \{ [IC_X(x_i)]^2 \left[\sum_{\substack{j=1 \\ j \neq i}}^n (x_j - a) \right] + [1 - IC_X(x_i)]^2 (x_i - a) \}$$

sujeto a

$$(x_j, x_k) C'_X(x_l, x_m) \Leftrightarrow IC_X(x_k) - IC_X(x_j) \leq IC_X(x_m) - IC_X(x_l)$$

dicho optimo es invariante por localizacion y escala , segun se deduce del Teorema 2.3.2 , y viene dado por :

$$IC_X^*(x_i) = \frac{x_i - a}{\sum_{j=1}^n (x_j - a)}$$

y como $IC_X^*(x) \in [0,1]$, este puede interpretarse como un conjunto difuso , definido sobre X , a traves del criterio C.

Es de destacar , que aunque la relacion C_X , no depende de X , IC_X^* si que es funcion del conjunto X.

3.4.1 OUTLIERS

(I)

$$IC_X^*(x_{(1)}) = \frac{x_{(1)} - a}{\sum_{j=1}^n (x_j - a)}$$

es el estadístico que estudia Lewis y Fieller (1979) , para testar la presencia de un outlier inferior , en una muestra procedente de una poblacion Gamma , con origen a conocido.

(II)

$$IC_X^*(x_{(n)}) = \frac{x_{(n)} - a}{\sum_{j=1}^n (x_j - a)}$$

es el estadístico que se recoge en los trabajos de Fisher (1929) , Cochran (1941), para testar la presencia de un outlier superior en una población de las mismas características del caso (I) .

3.5 GENERALIZACION AL CASO DE MAS DE UNA OBSERVACION.

En los casos analizados , se ha estudiado el comportamiento de cada elemento del conjunto de observaciones X , respecto de un determinado criterio C . Sin embargo , en ocasiones puede resultar de interes el estudio , del comportamiento de varios elementos de X , respecto de C , de forma simultanea . Para ello se propone el siguiente esquema :

Dado $X = \{ x_1 , \dots , x_n \}$, se considera el conjunto $X^s = X \times \dots \times X$, al cual es posible extender de forma natural los criterios estudiados anteriormente .

3.5.1 GENERALIZACION DEL CRITERIO DE DISPERSION CENTRAL

Dado $X \subset \mathbb{R} / \#(X) = n$, sea $X^s = X \times \dots \times X$, $s < n$, sobre el que se desea analizar el comportamiento de sus elementos respecto del criterio C , que determina la relacion :

$\tilde{x}_i = (x_{i1}, \dots, x_{is}) C_{X^s} (x_{j1}, \dots, x_{js}) = \tilde{x}_j$ si la
 la varianza del conjunto $X - \{ x_{ih} , h=1, \dots, s \}$ es
 mayor o igual que la de $X - \{ x_{jh} , h=1, \dots, s \}$, es
 decir :

$$\tilde{x}_i C_{X^s} \tilde{x}_j \quad \langle == \rangle \quad S_{i1, \dots, is}^2 \geq S_{j1, \dots, js}^2$$

donde $S_{i1, \dots, is}^2 = (1/n_i) \sum_{\substack{p=1 \\ p+i_h, h=1, \dots, s}}^n (x_p - \bar{x}_i)^2$

con $n_i = \#(X - \{ x_{ih} , h = 1, \dots, s \})$

y $\bar{x}_i = (1/n_i) \sum_{\substack{p=1 \\ p+i_h, h=1, \dots, s}}^n x_p$

Es de notar que la relacion definida de esta forma
 generaliza a la introducida en el epigrafe 3.2 , ya que
 dados $x , x' \in X$:

$$x C_X x' \quad \langle == \rangle \quad (x, \dots, x) C_{X^s} (x', \dots, x')$$

incluso , generalizaria a las relaciones que se definan
 de igual forma sobre X^r , $r < s$.

De igual forma a como se hizo en los casos ya anali-
 zados , se determinan los conjuntos C_{X^s} -dominantes y
 C_{X^s} -dominados y todo subconjunto de X^s esta propia-
 mente C_{X^s} -dominado , y la relacion es T-evaluable.

Las funciones de C-similitud y C-similitud signada

naturales se definen de igual forma , y el funcional a optimizar , como generalizacion del dado en el epigrafe 3.2 , viene expresado por :

$$\sum_{\tilde{x}_i \in X^S} \{ [IC_{X^S}(\tilde{x}_i)]^2 (nS_X^2 - n_i S_i^2) + [1 - IC_{X^S}(\tilde{x}_i)]^2 n_i S_i^2 \}$$

sujeto a

$$(\tilde{x}_j, \tilde{x}_k) C'_{X^S}(\tilde{x}_1, \tilde{x}_m) \Leftrightarrow IC_{X^S}(\tilde{x}_j) - IC_{X^S}(\tilde{x}_k) < IC_{X^S}(\tilde{x}_1) - IC_{X^S}(\tilde{x}_m)$$

El optimo de este funcional es invariante por localizacion y escala, segun se dedujo en el epigrafe 3.2 , y viene dado por :

$$IC_{X^S}^*(\tilde{x}_i) = \frac{n_i S_{i1, \dots, is}^2}{n S_X^2}$$

el cual puede interpretarse como un conjunto difuso definido sobre X^S , a traves del criterio C , ya que $IC_{X^S}^*(\tilde{x}_i) \in [0, 1]$, $\forall \tilde{x}_i \in X^S$

OUTLIERS

(I) Debido a que C_{X^S} generaliza a C_X , segun se ha visto , pueden obtenerse como casos particulares los resultados del epigrafe 3.2.1 .

(II) Si se representa mediante $(\dots x_{(i)}, x_{(j)} \dots)$ a cualquier elemento de X^S con unicamente dos componentes distintas : $x_{(i)}$ y $x_{(j)}$, se tiene que

$$IC_{X^S}^* (\dots x_{(1)}, x_{(n)} \dots) = \frac{(n-2) S_{(1),(n)}^2}{n S_X^2}$$

es el estadístico propuesto por Grubbs (1950) , para detectar la presencia de un outlier inferior y uno superior en una muestra extraída de una población normal de parámetros desconocidos.

Analogamente :

$$IC_{X^S}^* (\dots x_{(1)}, x_{(2)} \dots) = \frac{(n-2) S_{(1),(2)}^2}{n S_X^2}$$

$$IC_{X^S}^* (\dots x_{(n-1)}, x_{(n)} \dots) = \frac{(n-2) S_{(n-1),(n)}^2}{n S_X^2}$$

se utilizan para testar la presencia de dos outliers , inferiores y superiores , respectivamente , en la misma situación antes mencionada.

(III) Si se representa por $(\dots x_{(n-k+1)}, \dots, x_{(n)} \dots)$ a cualquier elemento de X^S , con unicamente k componentes distintas : $x_{(n-k+1)}, \dots, x_{(n)}$, se obtiene que

$$IC_{X^s}^* (\dots x_{(n-k+1)}, \dots, x_{(n)} \dots) = \frac{(n-k) S_{(n-k+1), \dots, (n)}^2}{n S_X^2}$$

es el estadístico estudiado por Grubbs (1950, 1969), Dixon (1950), McMillan (1971), Tietjen y Moore (1972) y Fieller (1976), para testar la presencia de k-outliers superiores en las condiciones anteriormente indicadas.

Analogamente se deducirían, como casos particulares los estadísticos utilizados para testar la presencia de k-outliers inferiores ó t-inferiores y k-t superiores ($k > t$).

Para testar k-outliers, sin especificar dirección, se considera el mínimo de $IC_{X^s}^*$, sobre el conjunto de elementos, con k-componentes distintas.

3.5.2 GENERALIZACION DEL CRITERIO DE DISPERSION BASADO EN EL RECORRIDO.

La generalización del criterio dado en epigrafe 3.3, puede realizarse de igual forma en este caso, sin embargo, se prestara atención especial al caso $s = 2$, que es el que presenta mayor interés desde un punto de vista práctico.

Dado $X \subset R / \#(X) = n, n > 2$, sea $X^2 = X \times X$, sobre el que se desea analizar el comportamiento de sus

elementos respecto del criterio C que determina la relación

$(x_1, x_j) C_{X^2} (x_h, x_k)$ sii el recorrido de $X - \{x_1, x_j\}$ es mayor o igual que el de $X - \{x_h, x_k\}$.

Por tanto, la relación puede expresarse mediante:

$$(x_1, x_j) C_{X^2} (x_h, x_k) \Leftrightarrow R_{1,j} \geq R_{h,k}$$

siendo $R_{1,j} = \text{Max}_{h \neq 1,j} \{x_h\} - \text{Min}_{h \neq 1,j} \{x_h\}$

Al igual que el caso estudiado en el epigrafe 3.5.1, esta relación generaliza a la estudiada en el epigrafe 3.3

Además, si se denota por $R_{(i),(j)}$ al recorrido $X - \{x_{(i)}, x_{(j)}\}$ se obtiene que:

$$\begin{aligned} R_{(i),(j)} &= R_{(i)} && \text{si } i = j \\ R_{(i),(j)} &= x_{(n)} - x_{(1)} && \text{si } 1 < i, j < n \\ R_{(1),(j)} &= \begin{cases} x_{(n)} - x_{(3)} & \text{si } j = 2 \\ x_{(n)} - x_{(2)} & \text{si } j = 3, \dots, n-1 \\ x_{(n-1)} - x_{(2)} & \text{si } j = n \end{cases} \\ R_{(i),(n)} &= \begin{cases} x_{(n-1)} - x_{(1)} & \text{si } i = 2, \dots, n-2 \\ x_{(n-2)} - x_{(1)} & \text{si } i = n-1 \end{cases} \end{aligned}$$

Para los conjuntos C_{X^2} -dominantes y C_{X^2} -dominados se verifica:

$$(x_{(i)}, x_{(j)}) \in DC_{X^2}(x_{(h)}, x_{(k)}) \quad 1 < i, j < n; \forall h, k$$

$$dC_{X^2}(x_{(1)}, x_{(j)}) = X \quad 1 < i, j < n$$

y el elemento \tilde{x}_* , $(C_{X^2} \text{ d } M)$ de X^2 , cumple :

$$R_{\tilde{x}_*} = \text{Min}_{i,j} \{ R_{i,j} \} = \text{Min} \{ R_{(1),(2)}, R_{(n),(n-1)}, R_{(1),(n)} \}$$

En esta situación cabe la posibilidad de que el elemento $(C_X \text{ d } M)$ de X , para la relación estudiada en el epigrafe 3.3, no forme parte del par $(C_{X^2} \text{ d } M)$ de X^2 .

Las funciones de C-similitud y C-similitud signada naturales, se definen de la forma usual, y el funcional a optimizar, como generalización del dado en epigrafe 3.3, es :

$$\sum_{(x_i, x_j) \in X^2} \{ [IC_{X^2}(x_i, x_j)]^2 [R - R_{i,j}] + [1 - IC_{X^2}(x_i, x_j)]^2 R_{i,j} \}$$

sujeto a la restricción (3.4).

El óptimo de esta función es invariante por traslación y escala, y viene dado por :

$$IC_{X^2}^*(x_i, x_j) = \frac{R_{i,j}}{R}$$

el cual puede interpretarse, al igual que en los casos anteriores, como un conjunto difuso sobre X^2 .

OUTLIERS

Pueden considerarse como casos particulares de este criterio los resultados obtenidos en el epigrafe 3.3.1, y además :

(I) Los estadísticos propuestos por Dixon (1950)

$$E_1 = \frac{x_{(3)} - x_{(1)}}{x_{(n)} - x_{(1)}}$$

$$E_2 = \frac{x_{(n)} - x_{(n-2)}}{x_{(n)} - x_{(1)}}$$

$$E_3 = \frac{x_{(n-1)} - x_{(2)}}{x_{(n)} - x_{(1)}}$$

para testar la presencia de dos outliers inferiores, dos outliers superiores y uno superior y otro inferior respectivamente, en una muestra extraída de una población exponencial de origen desconocido, se obtienen a partir de las expresiones

$$1 - IC_{X^2}^*(x_{(1)}, x_{(2)})$$

$$1 - IC_{X^2}^*(x_{(n-1)}, x_{(2)})$$

$$1 - IC_{X^2}^*(x_{(1)}, x_{(n)})$$

(IV) En el caso en que se hubiese realizado el estudio de forma generica sobre X^S , se obtiene, como caso particular :

$$\frac{x_{(n-q)} - x_{(p+1)}}{x_{(n)} - x_{(1)}}$$

que es el valor correspondiente a

$$IC_{X^S}^* (\dots, x_{(1)}, \dots, x_{(p)}, x_{(n-q+1)}, \dots, x_{(n)}, \dots)$$

y coincide con el estadístico propuesto por Dixon (1950, 1951), para testar la presencia de q-outliers superiores y p-inferiores en la situación indicada previamente .

3.5.3 GENERALIZACION DEL CRITERIO NATURAL

La generalización de este criterio , para el caso en que se se desee estudiar el comportamiento conjunto de una serie de observaciones ,se realizaria de forma similar a las anteriores . Se han de resaltar los casos particulares :

$$E_1 = \frac{x_{(1)} + \dots + x_{(k)} - k a}{\sum_{i=1}^k (x_i - a)}$$

$$E_2 = \frac{x_{(n-k+1)} + \dots + x_{(n)} - k a}{\sum_{i=1}^k (x_i - a)}$$

que coinciden con los estadísticos tratados en los trabajos de Fieller (1976), Lewis y Fieller (1978), para testar la presencia de k-outliers inferiores ó superiores , respectivamente en una muestra extraída de una población gamma de origen a conocido .

3.6 CRITERIO DE DISPERSION CENTRAL MULTIVARIANTE

Dado $X \in \mathbb{R}^{k \times n}$ / $\#(X) = n$, con $\text{ran}(X) = k < n$,

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1j} & \cdots & x_{1n} \\ \vdots & & \vdots & & \vdots \\ x_{k1} & \cdots & x_{kj} & \cdots & x_{kn} \end{pmatrix} = (x_1 \cdots x_j \cdots x_n)$$

las matrices de suma de cuadrados y suma de productos asociadas a X y $X_{(i)} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$, que se representaran por $S_{(n)}$ y $S_{(n-1)}^i$, respectivamente , vienen dadas por

$$S_{(n)} = X (I_n - (1/n)E_{nn}) X'$$

$$S_{(n-1)}^i = X_{(i)} (I_{n-1} - [1/(n-1)]E_{n-1 \ n-1}) X'_{(i)}$$

donde por I_h se denota a la matriz identidad de orden h y por E_{ab} , a la matriz de dimensiones $a \times b$, con todos sus elementos iguales a la unidad .

En esta situacion , se considera la relacion binaria definida sobre X , por :

$$x_i C_X x_j \Leftrightarrow \det (S_{(n-1)}^i) \geq \det (S_{(n-1)}^j) \quad (3.9)$$

Esta relacion , para el caso particular $k = 1$, coincide con la estudiada en el epigrafe 3.2 por lo que la relacion dada por (3.9) puede considerarse como una generalizacion del Criterio de Dispersion Central analizado en el epigrafe 3.2.

Los conjuntos C_X -dominantes y C_X -dominados, se construyen de la forma usual y evidentemente, el criterio que da lugar a C_X es T-evaluable sobre X.

Las funciones de C-similitud y C-similitud signada naturales son:

$$SC_X(x_i, x_j) = \left| \det (S_{(n-1)}^i) - \det (S_{(n-1)}^j) \right|$$

$$SSC_X(x_i, x_j) = \det (S_{(n-1)}^i) - \det (S_{(n-1)}^j)$$

La selección del IC_X^* , se realizara mediante la minimización de :

$$[IC_X(x_j)]^2 [\det(S_{(n)}) - \det(S_{(n-1)}^i)] + [1 - IC_X(x_j)]^2 \det(S_{(n-1)}^i)$$

bajo la restricción

$$(x_j, x_k) C'_X(x_1, x_m) \Leftrightarrow IC_X(x_j) - IC_X(x_k) \leq IC_X(x_1) - IC_X(x_m)$$

El optimo es invariante por localización y escala, y viene dado por :

$$IC_X^*(x_j) = \frac{\det (S_{(n-1)}^j)}{\det (S_{(n)})}$$

Si se representa por $\bar{x} = (1/n) X E_{n-1}$, se obtiene

$$IC_X^*(x_j) = \frac{\det (S_{(n)} - \frac{n}{n-1} (x_j - \bar{x})(x_j - \bar{x})')}{\det (S_{(n)})} = \dots$$

$$\dots = 1 - \frac{n}{n-1} (x_j - \bar{x})' S_{(n)}^{-1} (x_j - \bar{x})$$

e IC_X^* puede interpretarse como un conjunto difuso sobre X , definido a través del criterio C .

OUTLIERS

(I) Wilks (1963), considera

$$\text{Min}_{j=1, \dots, n} \frac{\det (S_{(n-1)}^j)}{\det (S_{(n)})}$$

(al que denomina "one-outlier scatter ratio"), como estadístico para testar la presencia de un outlier en una muestra extraída de una población normal multivariante de parámetros desconocidos.

(II) La generalización llevada a cabo en el epigrafe 3.5, puede realizarse también en este caso obteniéndose:

$$IC_{X^s}^*(x_{j_1}, \dots, x_{j_s}) = \frac{\det (S_{(n-s)}^{j_1, \dots, j_s})}{\det (S_{(n)})}$$

y

Min $IC_{X^s}^*(x_{j_1}, \dots, x_{j_s})$ coincide con el estadístico propuesto por Wilks (1963), para testar en bloques s observaciones outliers en una muestra extraída en las mismas condiciones de (I).

siguiendo el procedimiento indicado , se podrian obtener algunos otros estadisticos utilizados en distintas situaciones particulares para la identificacion de outliers.

De la definicion dada para el termino outlier en el Capitulo I ,se deduce que el analisis de las observaciones outliers , puede realizarse mediante el estudio del comportamiento de los distintos elementos de la masa de datos respecto del criterio que se desee analizar en los mismos.

El estudio llevado a cabo en el Capitulo II , posibilita que en el Capitulo III , se obtengan algunos de los estadisticos mas utilizados para la identificacion de outliers.

Algunas de las ventajas que presenta el procedimiento que se propone frente a las tecnicas clasicas de identificacion de outliers son las siguientes :

- (I) Los estadisticos se obtienen a traves de un metodo formal.
- (II) No se obtienen bajo hipotesis de modelo poblacional alguno.
- (III) La consideracion de las posibles observaciones outliers carece de cualquier tipo de subjetividad inicial.

Por ultimo , cabe destacar , que el metodo propuesto depende del criterio que se considere y este debera estar en relacion con el analisis estadistico posterior que se desee realizar con el conjunto de observaciones experimentales.

BIBLIOGRAFIA

- Anscombe , F.J. (1960) "Rejection of outliers" .Technometrics,2, 123-147 .
- Barnett , V. y Lewis , T. (1984) "Outliers in Statistical Data".2nd Edition. Ed. John Wiley and Sons.
- Beckman ,R.J. y Cook ,R.D. (1983) "Outlier ... s". Technometrics, 25, 119-163 .
- Boscovich ,R.J. (1757) "De litteraria expeditione per pontificiam ditionem,et synopsis amplioris operis, ac habentur plura ejus ex exemplaria etiam sensorum impressa". Bononiensi Scientiarum et Artum Instuto Atque Academia Commentarii, 4, 353-396.
- Chatfield , C. (1985) "The initial examination of data". Journal Royal Statistical Society . Ser. A, 148, 214-253.
- Cochran , W.G. (1941) "The distribution of the largest of a set of variances as a fraction of their total" . Ann. Eugen.,11,47-52.
- Collet , D., y Lewis , T. (1976) "The subjective nature of outlier rejection procedures".Applied Statistics, 25, 228-237.

- David , H. A. y Paulson , A.S. (1965) "The performance of several tests for outliers" . *Biometrika* ,52, 429-436.
- Dixon , W.J. (1950) "Analysis of extreme values" . *Ann. Math. Statis.* ,21, 488-506.
- Dixon , W.J. (1951) "Ratios involving extreme values" . *Ann. Math. Statis.* ,22, 68-78.
- Dunn , J.C.(1974) "A Fuzzy Relative of the ISODATA Process and its use in detecting Compact Well-Separated Clusters".*Journal of Cybernetics*,3, 32-57.
- Ferguson , T.S. (1961a) "On the rejection of outliers". *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability* ,Vol.1, 253-287.
- Ferguson , T.S. (1961b) "Rules for rejection of outliers". *Rev. Inst. Int. de Statist.*, 29, 29-43.
- Fieller , N.R.J. (1976) "Some Problems related to the Rejection of Outlying Observations".Ph.D.Thesis, University of Sheffield.
- Fisher , R.A. (1929) "Tests of significance in harmonic analysis" . *Proc. Roy. Soc. A*, 125, 54-59.
- Gnanadesikan , R. y Kettenring, J.R. (1972) "Robust estimates, residuals and outlier detection with multi response data" . *Biometrics* ,28, 81-124.

- Grubbs ,F.E. (1950) "Sample criteria for testing outlying observations".Ann. Math. Statis.,21,27-58.
- Grubbs ,F.E. (1969) "Procedures for detecting outlying observations in samples".Technometrics,11, 1-21.
- Hampel F.R. (1968) "Contributions to the Theory of Robust Estimation". Ph.D. dissertation, University of California-Berkeley.
- Hampel,R.R.; Ronchetti,E.M.; Rousseeuw,P.J.; Stahel,W.A. (1986) "Robust Statistics.The approach based on influence functions" . Ed. John Wiley and Sons.
- Hawkins , D.M. (1978) "Analysis of three tests for one or two outliers". Statistica Neerlandica, 32 , 137-148.
- Hawkins , D.M. (1980) "Identificacion of outliers" . Ed. Chapman y Hall.
- Kabe ,D.G. (1970) "Testing outliers from an exponential population" . Metrika ,15, 15-18.
- Kale , B.K. (1976) "Detection of outliers". Sankhya B , 38, 356-363 .
- Kendall ,M.G. y Buckland , W.R. (1957) "A dictionary of Statistical Terms" . Longman.
- King ,E.P. (1953) "On some procedures for the rejection of suspected data" . J. Amer. Statist. Ass. ,48, 531-533.

- Krantz , D.H. ;Luce , R.D. ;Suppes ,P. and Tversky , A.
(1971) "Foundations of Measurement", Vol I ,
Ed. Academic Press.
- Lewis , T. y Fieller , N.R.J. (1979) "A recursive algo-
rithm for null distributions for outliers :
I. Gamma samples" . Technometrics ,21, 371-376..
- McMillan, R.G. (1971) "Tests for one or two outliers in
normal samples with known variance".Technometrics,
13, 75-85.
- Miller, R.G. Jr. (1981) "Simultaneous Statistical In-
ference". 2nd. Edition. Ed. Springer Verlag.
- Pearson, E.S. y Chandra Sekar, C. (1936) "The efficien-
cy of statistical tools and a criterions for the
rejection of outlying observations" .Biometrika ,
28, 308-320.
- Rosado,F.M.F. (1984) "Existencia e Deteccao de outliers.
Uma Abordagem Metodologica" . Ph.D. Thesis. Uni-
versidade de Lisboa.
- Tietjen ,G.L. y Moore , R.H. (1972) "Some Grubbs-type
statistics for the detection of several outliers".
Technometrics ,14, 583-597.
- Tucker,A.W. (1956). "Dual Systems of homogeneous linear
relations". IN H.W. Kuhn and A.W. Tucker (Eds.),
Linear inequalities and relates systems, Annals

of Mathematics Study 38. Princeton University
Press, Princeton, New Jersey.

Wilks, S.S. (1963) "Multivariate statistical outliers".
Sankhya ,A, 25, 407-426.

UNIVERSIDAD DE SEVILLA

Reunido el Tribunal integrado por los abajo firmantes
en el día de la fecha, para juzgar la Tesis Doctoral de
D. Juan L. Moreno Redollo
titulada "Análisis Cualitativo de datos estadísticos"

acordó otorgarle la calificación de APTO CUM LAUDE

Sevilla, 8 de Julio 1987

El Vocál,

El Presidente

El Vocal,

El Secretario,

El Vocal,

El Doctorado,

