



Análisis de conglomerados y clasificación de series de tiempo. Utilizando Kullback-Liebler sobre la densidad espectral. Un ejemplo aplicado a audio de gatos y perros

Facultad de Ciencias Actuariales
Alumno: Campos Martínez Joaquín Nicolás
Asesor: Dr. Cuevas Covarrubias Carlos

Índice general

1. Marco Teórico	1
1.1. Objetivo	1
1.2. Relevancia	1
1.3. Antecedentes	2
2. Como perros y gatos	3
2.1. Distinguiendo su espectro	3
2.2. Densidad espectral	4
2.3. Clasificación cortando las series	4
Referencias	6
2.4. Series de tiempo	8
2.5. Análisis de conglomerados	8
2.5.1. Conglomerados Jerárquicos	9
2.5.2. K-centroides	11
2.6. Análisis espectral	12
2.7. Kullback-Liebler	12
3. Anexo B: Código	14

Introducción

En este practicum queremos mostrar el uso de la densidad espectral para clasificar series de tiempo, algunas de las razones que impulsaron esta idea fueron: primero algunos ejemplos, particularmente la clasificación de temblores basada en esta misma metodología([Shumway, 2003](#))([Yoshihide Kakizawa, 1998](#)), en el caso de Shumway como un medio para monitorear que se respeten las reglas internacionales de pruebas nucleares. También tenemos el caso de la clasificación de flamas según el sonido generado en la combustión.([J. Yuan, 2000](#))

Nuestro ejemplo en particular es la clasificación de audios de perros y gatos que forma parte de una base de datos de 48 clases acústicas adicionales, utilizada en un artículo sobre redes neuronales ([Naoya Takahashi y Gool, 2016](#)).

La base utilizada consta de 164 archivos WAV con sonidos de gatos, maullidos, ronroneos, o similares y 113 sonidos de perros, ladridos principalmente. Decidimos seguir el ejercicio propuesto en el sitio Kaggle([marc moreaux, 2017](#)), donde se propone dividir la muestra en conjunto de entrenamiento con 115 archivos de gatos y 64 de perros, los 49 restantes de cada categoría se asignan al conjunto de prueba.

Los resultados preliminares son bastante satisfactorios, y notamos que el espectro de maullidos es muy distinto de ladridos en términos de duración y armonía.

Al final proponemos como ejercicio el posible uso en la clasificación de series de tiempo económicas, particularmente de precios de acciones o activos, de modo que podemos categorizar por estructuras de covarianza con un análisis de conglomerados, así podríamos decir que activos en un mismo grupo tienen un riesgo similar.

1. Marco Teórico

1.1. Objetivo

Nosotros queremos presentar una metodología de clasificación y conglomeración de series de tiempo basada en la estructura de correlación. Nuestra base es utilizar la medida de discrepancia de Kullback-Liebler sobre la densidad espectral para clasificar en un grupo, según el espectro medio de cada grupo.

1.2. Relevancia

Las series de tiempo pueden ser observadas en diferentes ámbitos, información de ventas, precios de acciones, tasas cambiarias, información sobre el clima, datos biomédicos, etc. ([Saeed Aghabozorgi, 2015](#))

Las aplicaciones de análisis de conglomerados en estos ámbitos son variadas, por ejemplo, en medicina es importante reconocer la diferencia entre señales normales en un ECG, EEG, EMG de aquellas producidas por enfermedades. En sismología, para discriminar las ondas sísmicas, de los movimientos naturales de la tierra ([Corduas Marcella, 2008](#)) o para monitorear violaciones al CTBT (Tratado de prohibición completa de los ensayos nucleares, en español) ([Shumway, 2003](#)).

Otras aplicaciones son, en astronomía las curvas de luz muestran el brillo de una estrella en un periodo de tiempo, en medicina la actividad cerebral, en el medio ambiente y urbanización los niveles de la marea, niveles de contaminantes en el aire ($PM_{2.5}$, PM_{10}) ([Saeed Aghabozorgi, 2015](#)).

Otra razón es la gestión de portafolios de inversión, donde se requiere diversificar el riesgo, y no tener muchas acciones de empresas con giros similares, o que esten altamente correlacionadas.

1.3. Antecedentes

Hay múltiples artículos sobre análisis de conglomerados, y algunos de ellos se enfocan en series de tiempo. *Time-series clustering - A decade review* ([Saeed Aghabozorgi, 2015](#)) tiene como objetivo comparar los enfoques más populares.

«This review will expose four main components of time-series clustering and is aimed to represent an updated investigation on the trend of improvements in efficiency, quality and complexity of clustering time-series approaches during the last decade and enlighten new paths for future works» ([Saeed Aghabozorgi, 2015](#))

Clustering of time series data - a survey busca resumir investigación realizada en el tema y sus aplicaciones en varios campos. Incluye aspectos básicos de los algoritmos más generales que se utilizan en estudios de conglomerados, medidas de similitud y disimilitud, evaluación de conglomerados. ([Liao, 2005](#))

Ambos artículos distinguen 3 principales categorizaciones de los acercamientos a conglomeración de series de tiempo, trabajar con los datos brutos, extraer «características» de los datos, basarse en un modelo al que se ajusten los datos.

«Computational Models of Music Similarity and their Application in Music Information Retrieval» es una tesis que busca mostrar el desarrollo de las técnicas de descubrimiento de música basado en medidas de similitud y disimilitud. Uno de los enfoques consiste en utilizar de función de densidad espectral. Sin embargo también utiliza técnicas aplicables a series de tiempo. ([Pampalk, 2006](#)).

2. Como perros y gatos

2.1. Distinguiendo su espectro

Basandonos en el ejemplo de (Tatman, 2017), decidimos tomar una observación de maullido cuyo sonido sea claro, y un ladrido también nítido. La selección fue el archivo *cat_67.wav* y *dog_barking_105.wav*.

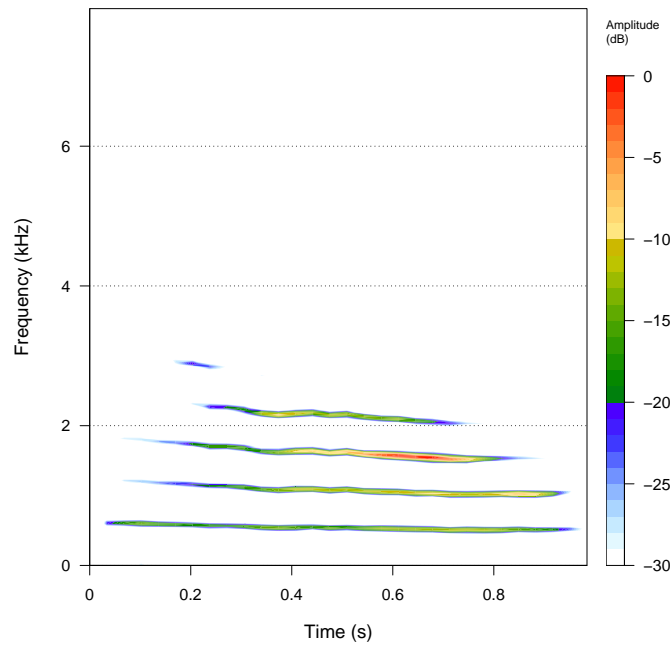


Figura 2.1: Espectrograma maullido

De los espectrogramas 2.1 y 2.2 notamos que los maullidos son en tendencia armónicos y tienen mayor duración que un ladrido de perro. Esta tendencia se sigue para los archivos que tienen menos ruido de fondo y son más nítidos.

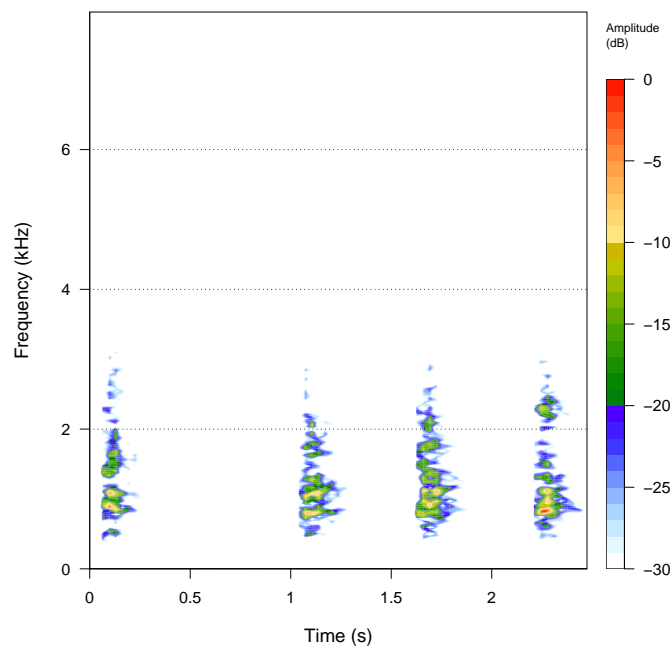


Figura 2.2: Espectrograma ladrido

2.2. Densidad espectral

El estimador que utilizamos de la función de densidad espectral es la corrección al periodograma con la ventana de Hann definida por

$$w(n) = \frac{1}{2} \left(1 - \cos \left(\frac{2\pi n}{N-1} \right) \right)$$

Donde N es el número de observaciones en la serie. De modo que el espectro queda representado por

2.3. Clasificación cortando las series

Los audios son de distinta duración, de modo que los espectros estimados serán de distintas frecuencias, para solucionar este problema cortamos las series con la duración más corta (*cat_55.wav*) de poco menos de un segundo. Partimos del punto medio de las mismas y tomamos aproximadamente medio segundo a la izquierda y medio a la derecha.

Una vez obtenida la densidad espectral, obtuvimos el espectro promedio por clase, la figura 2.3 nos muestra la densidad espectral de los audios de gatos, podemos notar que

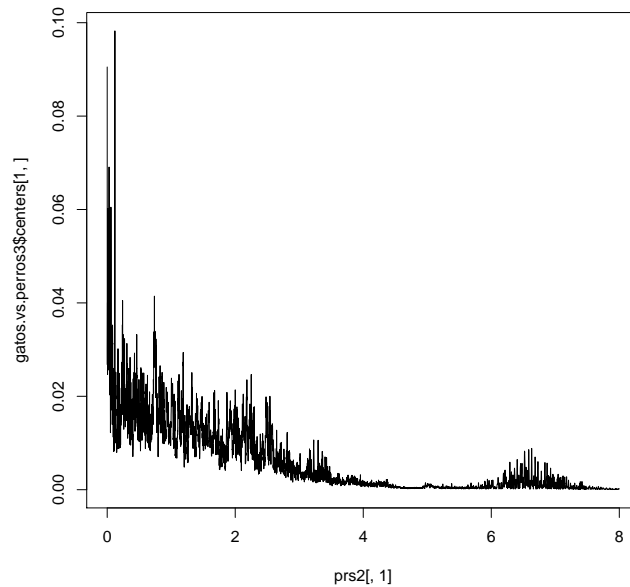


Figura 2.3: Densidad gatos 1

tiene su moda en frecuencias bajas, y luego decae lentamente a 0. En cambio, la densidad espectral para ladridos, 2.4, tienen su moda en frecuencias medias, y decae rápidamente a 0. Ahora que tenemos nuestros espectros promedio, el siguiente paso es la regla de clasificación

$$\text{clasifica en } \left\{ \begin{array}{ll} G & \text{si } I(f_g, f) \leq I(f_p, f) \\ P & \text{si } I(f_g, f) > I(f_p, f) \end{array} \right\}$$

Donde $I(f, g) = \int_x f(x) \log\left(\frac{f(x)}{g(x)}\right) dx$ es la información de Kullback. La interpretación de

I es la información media por observación en favor de f contra g . En este caso nuestras x son las frecuencias con las que estamos trabajando.

Al aplicar esta regla obtenemos la matriz de confusión para el conjunto de entrenamiento 2.3, donde nuestra tasa de error es del 17 %. Para el conjunto de prueba nuestra matriz de confusión nos indica una tasa de error del 26 %

	Gato	Perro
Gatos	77	37
Perro	6	58

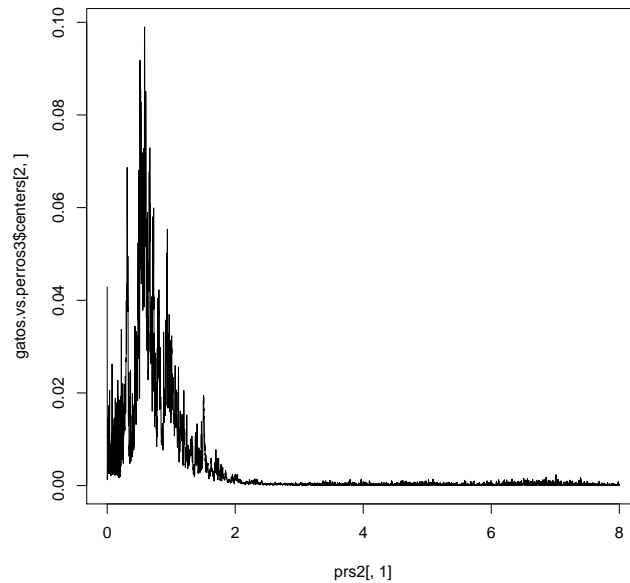


Figura 2.4: Densidad perro 1

	Gato	Perro
Gatos	36	13
Perro	3	46

Referencias

- Chatfield, C. (1989). *The analysis of time series : an introduction*.
- Chen Shihyen, Z. K., Ma Bin. (2009). *On the similarity metric and the distance metric*. <https://www.journals.elsevier.com/theoretical-computer-science>. ([Online; Obtenido 20-November-2017])
- Corduas Marcella, P. D. (2008). *Time series clustering and classification by the autoregressive metric*.
- Hana, R. (2014). *Cluster analysis of economic data*. (University of Economics, Prague, Czech Republic)
- J. Yuan, Y. Z. (2000). *Spectral analysis of combustion noise and flame pattern recognition*.
- Kenneth P. Burnham, D. R. A. (2007). *Model selection and multimodel inference: A practical information-theoretic approach*. ([Online; Obtenido 28-November-2017])
- Kullback, S. (1978). *Information theory and statistics*. ([Online; Obtenido 30-November-2017])

- Liao, T. W. (2005). *Clustering of time series data - a survey*. ([Online; Obtenido 09-Octubre-2017])
- marc moreaux. (2017). *Audio cats and dogs, classify raw sound events*. <https://www.kaggle.com/mmoreaux/audio-cats-and-dogs/home>. ([Online; Obtenido 04-Noviembre-2018])
- Naoya Takahashi, B. P., Michael Gygli, y Gool, L. V. (2016). *Deep convolutional neural networks and data augmentation for acoustic event recognition*.
- Pampalk, E. (2006). *Computational models of music similarity and their application in music information retrieval*. ([Online; Obtenido 09-Octubre-2017])
- Pradeep Rai, S. S. (2010). *A survey of clustering techniques*.
- Saeed Aghabozorgi, T. Y. W., Ali Seyed Shirkhorshidi. (2015). *Time-series clustering – a decade review*.
- Shumway, R. H. (2003). *Time-frequency clustering and discrimination analysis*.
- Tatman, R. (2017). *Visualizing woofs and meows*. <https://www.kaggle.com/rtatman/visualizing-woofs-meows/notebook>. ([Online; Obtenido 04-Noviembre-2018])
- Yoshihide Kakizawa, M. T., Robert H. Shumway. (1998). *Discrimination and clustering for multivariate time series*.

Anexos A: Teoría

2.4. Series de tiempo

Una simplificación de la definición de una serie de tiempo es: un proceso que varía en el tiempo, observaciones tomadas secuencialmente en el tiempo. En la vida real podemos observar varios fenómenos de este tipo, como los mencionados arriba.

Definición. Una serie de tiempo es una colección de de observaciones hechas en una secuencia de tiempo. Se dice que una serie de tiempo es continua si las observaciones son en tiempos continuos. Se dice discreta si cuando las observaciones se toman en tiempos específicos, usualmente separados a intervalos iguales

2.5. Análisis de conglomerados

«Clustering is a division of data into groups of similar objects. Each group, called cluster, consists of objects that are similar between themselves and dissimilar to objects of other groups.» ([Pradeep Rai, 2010](#))

Análisis de conglomerados no es una técnica asociada a un único tipo de problemas. Pero es asociado naturalmente a las ideas de grupos homogéneos, clases de equivalencia, datos multimodales.

Punj & Stewart identifican algunas de las principales aplicaciones en el campo de investigación de mercado como segmentación de mercado, identificación de grupos homogéneos de consumidores, desarrollo de nuevos productos potenciales, selección de un mercado de prueba, como técnica general de reducción.

Las técnicas aplicadas son varias, y algunas de las más tratadas en la literatura son conglomerados jerárquicos, y algunos algoritmos de k-centroides.

Tenemos que reconocer que nuestro objeto de estudio son observaciones multivariadas, medidas en p variables, que pueden ser categóricas (respuestas a una encuesta, sí, no) o

numéricas (valores continuos, como estatura y peso)

2.5.1. Conglomerados Jerárquicos

Se separa en 2 grandes categorías, aglomerativos y divisivos. El primer tipo considera cada observación, elemento de la muestra, como un conglomerado y comienza a agruparlos basado en alguna regla. El segundo comienza con toda la muestra como un gran conglomerado y comienza a separarlo en distintos conglomerados, en su mayoría conforme a una función objetivo.

Un análisis de conglomerados jerárquico es sensible, o depende principalmente de 2 factores.

- La medida de **similitud o distancia** elegida
- El algoritmo para agrupar elementos, la medida de **distancia o similitud** entre conglomerados

Conglomerados jerárquicos aglomerativos suelen ser los típicos ejemplos de libro de texto por su facilidad para explicarse, sin embargo esto no significa que su utilidad sea menor.

Primero debemos abordar el tema de distancia o similitud.

Definición. Dado un conjunto X una medida de distancia es una función $d : X \times X \rightarrow \mathbb{R}$ tal que cumple las siguientes propiedades

- $d(x, y) = 0$ si y sólo si $x = y$
- $d(x, y) = d(y, x)$ para todas $x, y \in X$
- $d(x, y) + d(y, z) \geq d(x, z)$ para todas $x, y, z \in X$

La propiedad de no negatividad $d(x, y) \geq 0$ para todas $x, y \in X$ queda definida por las propiedades anteriores.

La idea de distancia suele ser intuitiva, y se asocia comúnmente con la distancia euclidiana, sin embargo existen muchas otras formas de distancia, como la suma de valores absolutos de las componentes. En cambio, la similitud, aunque la idea es intuitiva, no suele ser tan clara cuando uno maneja valores numéricos o cateogóricos.

Existen algunos intentos sobre la definición de una medida de similitud, entre ellos mencionaremos como ejemplo la definición presentada por Chen, Ma y Zhang ([Chen Shihyen, 2009](#))

Definición. Dado un conjunto X una medida de similitud es una función $s : X \times X \rightarrow \mathbb{R}$ que cumple las siguientes propiedades

- $s(x, y) = s(y, x)$ para todo $x, y \in X$
- $s(x, x) \geq 0$ para todo $x \in X$
- $s(x, x) \geq s(x, y)$ para todo $x, y \in X$
- $s(x, y) + s(y, z) \leq s(x, z) + s(y, y)$ para todo $x, y, z \in X$
- $s(x, x) = s(y, y) = s(x, y)$ si y sólo si $x = y$

La propiedad de simetría es intuitiva, se espera que un objeto sea tan similar a otro, como el otro al primero. La segunda propiedad no es necesaria, pero su importancia radica en el concepto que tenemos de similitud, un objeto no puede tener similitud negativa con respecto de sí mismo. La tercera propiedad parte de la idea de que un ningún objeto es tan similar a otro como a sí mismo.

La cuarta propiedad sería el equivalente a la desigualdad del triángulo, y aunque a primera vista parece poco comprensible, es más fácil analogarla con la idea de intersección de conjuntos, podemos decir que s es una medida de todo lo que tienen en común 2 objetos x, y . La última propiedad nos dice que sólo cuando 2 objetos son iguales se alcanzan las cuotas superiores de la similitud.

No es importante recordar esta definición de similitud, pues nosotros trabajaremos algo más parecido a una distancia que a una similitud.

Lo siguiente es construir una matriz de distancias o similitudes entre los objetos. Si denotamos $s_{ij} = s(x_i, x_j)$ y $d_{ij} = d(x_i, x_j)$, y tenemos un total de n observaciones, nuestra matriz de similitud es

$$\begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1n} \\ s_{21} & s_{22} & \cdots & s_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ s_{n1} & s_{n2} & \cdots & s_{nn} \end{bmatrix}$$

y nuestra matriz de distancias es

$$\begin{bmatrix} d_{11} & d_{12} & \cdots & d_{1n} \\ d_{21} & d_{22} & \cdots & d_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{n1} & d_{n2} & \cdots & d_{nn} \end{bmatrix}$$

En un algoritmo aglomerativo el siguiente paso es juntar las 2 observaciones, distintas, más similares (menos distantes) para formar un conglomerado, luego se recalcula la distancia del resto de los conglomerados (observaciones) al conglomerado recién formado.

¿Cómo se calcula la distancia entre 2 conglomerados?. Existen distintas reglas para calcular la distancia, algunas de las más conocidas son

- Vecino más cercano (single linkage), es decir, la distancia entre 2 conglomerados es la distancia más corta de sus elementos
- Vecino más lejano (complete linkage), es decir, la distancia entre 2 conglomerados es la distancia más grande entre sus elementos
- Promedio (average linkage), media de las distancias entre los elementos de los conglomerados.
- Centroide, se obtiene el centro de masa de los conglomerados y se calcula la distancia entre los mismo

Este proceso se repite hasta tener toda la muestra en un conglomerado, y cuando uno decida hacer el corte, cuando hay r conglomerados, donde el usuario escoge r .

2.5.2. K-centroides

Los algoritmos de k-centroides están basados en la minimización de una función de suma de cuadrados, es decir la suma de cuadrados de la distancia de los elementos de un conglomerado a su centroide.

En comparación con un análisis jerárquico, k-centroides depende de 2 configuraciones

- La primera es el número de conglomerados, k
- El tipo de centroide, centro de gravedad, medoide.

En el caso de **k-means**, el centroide es el centro de gravedad y

$$E = \sum_{h=1}^k \sum_{i=1}^n u_{i,h} d(x_i, \mu_h)^2$$

Es la función que deseamos optimizar ([Hana, 2014](#)). Donde los $u_{i,h}$ cumplen que $\sum_{h=1}^k u_{i,h} =$

1 para $i = 1, \dots, n$ y $\sum_{i=1}^n u_{i,h} = 1$ para $h = 1, \dots, k$, y μ_h es la media del h -ésimo conglomerado. En el caso de medoides sustituimos μ_h por m_h , donde m_h es el medoide del h -ésimo.

Hasta ahora sólo hemos presentado metodologías en las que una observación no puede ser parte de más de un conglomerado sin embargo existen otros algoritmos que asignan «grados de pertenencia» en lugar de un conglomerado fijo.

2.6. Análisis espectral

Inferencias basadas en la función de densidad espectral se dice que es un análisis in el Dominio de las frecuencias.

Teorema 1 (Wiener-Khintchine). Para cualquier proceso estocástico estacionario con función de autocovarianza $\gamma(k)$ existe una función monotonamente creciente $F(W)$ tal que

$$\gamma(k) = \int_0^{\pi} \cos(wk) dF(w)$$

Llamada la representación espectral de la función de autocovarianza, que involucra un tipo de integral llamada de Stieltjes

$F(w)$ tiene una interpretación física directa, es la contribución a la varianza que se puede atribuir a las frecuencias en el rango $(0, w)$ $F(w)$ es monótona en el intervalo $(0, \pi)$, así se puede descomponer en 2 funciones $F_1(w)$ y $F_2(w)$ talque

$$F(w) = F_1(w) + F_2(w)$$

Donde $F_1(w)$ es una función continua no decreciente y $F_2(w)$ es una función escalonada no decreciente (WOLD descomposition). Además $F_1(w)$ se relaciona con la componente puramente no determinista del proceso y $F_2(w)$ se relaciona con la componente determinista. Nuestro estudio es sobre los procesos puramente indeterminados, donde $F_2(w) = 0$, de tal forma que $F(w)$ es continua en $(0, \pi)$. La potencia está directamente relacionado al cuadrado de la amplitud de oscilación. y la varianza es la potencia total

Definición. La forma normalizada de $F(w)$ está dada por

$$F^* = \frac{F(w)}{\sigma_x^2}$$

De modo que $F^*(w)$ es la proporción de varianza atribuida por las frecuencias en el rango $(0, w)$. Como $F^*(\pi) = 1$ y es además monótona creciente tenemos que $F^*(w)$ tiene propiedades similares a una función de distribución acumulada. ([Chatfield, 1989](#))

2.7. Kullback-Liebler

Definición. Dadas dos funciones de densidad de probabilidad f, g definidas sobre un mismo espacio medible, y que son absolutamente continuas entre sí, la divergencia de Kullback-Liebler se define como

$$d_{KL}(f, g) = \int f(x) \ln \left(\frac{f(x)}{g(x)} \right) dx$$

La divergencia de Kullback-Liebler como una medida de divergencia dirigida. Una manera de entender lo anterior es primero considerar un modelo verdadero f y un modelo de aproximación g , entonces la divergencia de Kullback-Liebler nos habla de cuanta información se pierde utilizando g como aproximación ([Kenneth P. Burnham, 2007](#)). La divergencia de Kullback-Liebler no cumple con las propiedades de distancia mencionadas anteriormente.

La propiedad más cercana a una distancia, es la positividad y la nulidad cuando f y g son iguales exceptuando en puntos de probabilidad 0.

Ahora podemos definir la divergencia de Jeffreys

Definición. Dadas dos funciones de densidad de probabilidad f, g definidas sobre un mismo espacio medible, y que son absolutamente continuas entre sí, la divergencia de Jeffreys se define como

$$d_J(f, g) = d_{KL}(f, g) + d_{KL}(g, f) = \int (f(x) - g(x)) \ln \left(\frac{f(x)}{g(x)} \right) dx$$

En este caso, tenemos dos modelos, y la divergencia de Jeffreys nos indica la dificultad para distinguir entre ambos ([Kullback, 1978](#)). Esta nueva medida de divergencia soluciona el problema de la simetría'

3. Anexo B: Código