



Treball final de grau

GRAU DE MATEMÀTIQUES

Facultat de Matemàtiques i Informàtica
Universitat de Barcelona

Autor: Juan Tornero Lucas

Director: Dr. Nom tutor
Realitzat a: Departament....
(nom del departament)

Barcelona, May 17, 2017

1. Marco Teórico de las Cadenas de Markov

En este primer capítulo desarrollaremos toda la teoría que involucra a las Cadenas de Markov. Empezaremos viendo qué son las secuencias y, a continuación, los tipos de espacio de probabilidad conocidos como *Procesos*. Además, mediante una versión adaptada del *Teorema de Extensión de Kolmogorov* veremos cómo estos procesos están definidos de forma única. A continuación definiremos los estados de un proceso, para lo cual necesitaremos un par de resultados relacionados con los procesos que cumplen la propiedad de martingala.

Finalmente, veremos la definición de una Cadena de Markov y demostraremos que se trata de un *proceso*.

Secuencias

Comenzamos definiendo el conjunto canónico $X = \{0, 1, \dots, m-1\}$ con $m \in \mathbb{N}$ llamado *alfabeto*, cuyos elementos se denominan *símbolos*, los cuales no tienen por qué ser necesariamente números naturales. Si $x \in X^{\mathbb{Z}}$, entonces $x = \dots x_{-1}x_0x_1\dots$ es una secuencia numerable bi-infinita, donde los índices $i < 0$ denotan el pasado de la secuencia, y los $i \geq 0$ el futuro, particularmente el índice $i = 0$ es el primer símbolo desconocido de la secuencia.

En estos términos definimos una palabra $w \in X^l$ de longitud l , como una l -tupla de X . \emptyset denotará la palabra vacía de longitud 0. Una subsecuencia s es una estructura $s = (w, a, b)$, donde w es una palabra y $a, b \in \mathbb{Z}$ tq $|w| = b - a + 1$. Así, s puede también escribirse $s = s_a \dots s_b$, y a representaría el tiempo inicial, y b el tiempo final. Una secuencia x contendrá a la subsecuencia s si $\forall t \in [a, b], s_t = x_t$. En muchas ocasiones identificaremos s y w indistintamente, ya que lo único que las diferencia es el contexto temporal que le da a w la subsecuencia s .

El conjunto $A_s = \{x \in X^{\mathbb{Z}} | x_i = s_i \forall i \in [a, b]\}$ es el conjunto de las secuencias de $X^{\mathbb{Z}}$ que contienen a s . Si por ejemplo $s = (\emptyset, (a, a-1))$, entonces $A_s = X^{\mathbb{Z}}$.

Por último, el conjunto X^* denotará el de todas las palabras.

Procesos

Un proceso Q es un espacio de probabilidad estacionario en un espacio de secuencias.

Una medida de probabilidad es una función que asigna probabilidades a conjuntos del espacio medible de probabilidades $(X^{\mathbb{Z}}, \mathbb{X})$, donde \mathbb{X} es el conjunto definido como la menor colección de subconjuntos de $X^{\mathbb{Z}}$ tal que:

1. Para toda secuencia s , $A_s \in \mathbb{X}$.
2. \mathbb{X} es cerrado bajo complementos y uniones contables. Es decir, \mathbb{X} es una σ -álgebra.

En definitiva, \mathbb{X} es la menor σ -álgebra que contiene A_s fijado por s , y el par $(X^{\mathbb{Z}}, \mathbb{X})$ asigna probabilidades a estos conjuntos A_s .

Definiremos P la medida de probabilidad sobre $X^{\mathbb{Z}}$ tal que $P(s) = P(A_s)$, tenemos en particular:

$$P(\emptyset) = P(X^{\mathbb{Z}})$$

En nuestra definición de proceso nos referíamos también a ellos como estacionarios. Sea D la función desplazamiento $D : X^{\mathbb{Z}} \rightarrow X^{\mathbb{Z}}$, que actúa sobre todo $x \in X^{\mathbb{Z}}$ de manera que $D(x_t) = x_{t+1}$, es decir desplaza el tiempo de origen.

Decimos que P es una medida de probabilidad estacionaria si $\forall A \in \mathbb{X}, P(D(A)) = P(A)$, como D es de hecho un automorfismo sobre $X^{\mathbb{Z}}$, entonces $P(D^{-1}(A)) = P(D^{-1}(D(A))) = P(A)$. Esto en particular nos permitirá tener bien definido $P(w)$.

Finalmente podemos definir de forma más formal un proceso Q como el espacio de probabilidades estacionario $(X^{\mathbb{Z}}, \mathbb{X}, P)$.

Sea w una palabra y $s = (w, a, b)$, entonces si P estacionario $P(w) = P(s)$. Además de manera trivial podemos obtener que, si W_l es el conjunto de las palabras de longitud $l > 0$:

$$\sum_{z \in W_l} P(wz) = \sum_{z \in W_l} P(zw) = P(w) \quad (1)$$

Además, para $w = \emptyset$:

$$P(\emptyset) = 1 \quad (2)$$

El recíproco también es cierto, y cualquier función de X^* que satisfaga (1) y (2) define un proceso.

Unicidad de los Procesos

Este apartado estará dedicado a demostrar la unicidad de los procesos mediante el siguiente Teorema.

Teorema. Dado $f : X^* \rightarrow [0, 1]$ que satisface:

$$1. f(\emptyset) = 1$$

$$2. \forall w \in X^*, f(w) = \sum_{z \in X^*} f(zw) = \sum_{z \in X^*} f(wz)$$

existe un único proceso $Q = (X^{\mathbb{Z}}, \mathbb{X}, P)$ tq $\forall w \in X^*, P(w) = f(w)$.

Este resultado deriva del Teorema de Extensión de Kolmogorov (TEK). Usando las notaciones $\mathcal{B}(\mathbb{R}^n)$ y $\mathcal{B}(\mathbb{R}^{\mathbb{N}})$ para referirnos a las σ -álgebras en \mathbb{R}^n y $\mathbb{R}^{\mathbb{N}}$ respectivamente.

Teorema de Extensión de Kolmogorov. Suponemos que nos dan un conjunto de medidas de probabilidad μ_n consistentes en $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$, es decir que cumplen:

$$\mu_{n+1}(\prod_{i=1}^n (a_i, b_i] \times \mathbb{R}) = \mu_n(\prod_{i=1}^n (a_i, b_i]).$$

Entonces existe una única medida de probabilidad P' en $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}(\mathbb{R}^{\mathbb{N}}))$, con

$$P'(x|x \in \prod_{i=1}^n (a_i, b_i]) = \mu_n(\prod_{i=1}^n (a_i, b_i]).$$

No demostraremos el TEK directamente, sino que lo usaremos para demostrar una variante del mismo que se adapte a nuestras necesidades, teniendo en cuenta las diferencias que existen entre el espacio de probabilidad $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}(\mathbb{R}^{\mathbb{N}}), P')$ y el de un proceso estacionario $(X^{\mathbb{Z}}, \mathbb{X}, P)$

1. $\mathbb{R}^{\mathbb{N}}$ es producto de copias de \mathbb{R} , mientras que $X^{\mathbb{Z}}$ es producto de copias de un conjunto finito X . Para resolver esta diferencia utilizaremos una aplicación inyectiva $g : X \rightarrow \mathbb{R}$.
2. Los elementos de $\mathbb{R}^{\mathbb{N}}$ son secuencias infinitas mientras que los de $X^{\mathbb{Z}}$ son, de hecho, secuencias bi-infinitas. Utilizaremos de nuevo un mapeado, esta vez biyectivo $h : \mathbb{N} \rightarrow \mathbb{Z}$. Consideraremos los dígitos en el orden $0, 1, -1, 2, -2, \dots$
3. Por último, P' no necesita ser estacionario, así que usaremos el TEK para probar que P existe, y luego veremos que es estacionario.

Introducimos primero los mapas g y h .

Para el caso $g : X \rightarrow \mathbb{R}$ solo ha de cumplir que sea inyectivo, da igual la imagen de los símbolos $x \in X$.

El mapa biyectivo $h : \mathbb{N} \rightarrow \mathbb{Z}$ será

$$h(n) = \begin{cases} n/2 & n \text{ par} \\ (1-n)/2 & n \text{ impar} \end{cases}$$

Su inversa, por lo tanto:

$$h^{-1}(y) = \begin{cases} 2y & y > 0 \\ 1-2y & y \leq 0. \end{cases}$$

Vemos como h envía $1, 2, 3, 4, 5$ a $0, 1, -1, 2, -2$ respectivamente, y viceversa h^{-1} . Definiremos el conjunto $J(n) = \{y \in \mathbb{Z} | h^{-1}(y) \leq n\}$, y los valores $d(n)$ y $c(n)$ como los mayores y menores de $J(n)$ respectivamente, caracterizadas de la siguiente manera:

$$c(n) = \min(h(n), h(n-1)) = \begin{cases} \frac{2-n}{2} & n \text{ par} \\ \frac{1-n}{2} & n \text{ impar} \end{cases}$$

$$d(n) = \max(h(n), h(n-1)) = \begin{cases} \frac{n}{2} & n \text{ par} \\ \frac{n-1}{2} & n \text{ impar} \end{cases}$$

$$J(n) = \{c(n), \dots, d(n)\}$$

Observamos que c y d cumplen las propiedades:

P1. Si n par, entonces $h(n+1) < 0, c(n+1) = c(n) - 1, d(n+1) = d(n)$.

P2. Si n impar, entonces $h(n+1) > 0, c(n+1) = c(n), d(n+1) = d(n) + 1$.

Demostración

P1. Si n par, entonces $n = 2i$, para algún $i \in \mathbb{N} \setminus \{0\}$, y:

$$\begin{aligned} h(n+1) &\stackrel{n+1 \text{ impar}}{=} \frac{1-(n+1)}{2} = \frac{-2i}{2} = -i < 0 \\ \begin{cases} c(n+1) &\stackrel{n+1 \text{ impar}}{=} \frac{1-(2i+1)}{2} = -i \\ c(n) &\stackrel{n \text{ par}}{=} \frac{2-2i}{2} = -i + 1 \end{cases} \end{aligned}$$

Luego $c(n+1) = c(n) - 1$

$$\begin{cases} d(n) &\stackrel{n \text{ par}}{=} \frac{2i}{2} = i \\ d(n+1) &\stackrel{n+1 \text{ impar}}{=} \frac{(2i+1)-1}{2} = i \end{cases}$$

Luego $d(n+1) = d(n)$.

P2. Si n impar, entonces $n = 2i + 1$, para algún $i \in \mathbb{N} \setminus \{0\}$, y:

$$\begin{aligned} h(n+1) &\stackrel{n+1 \text{ par}}{=} \frac{2i+2}{2} = i+1 > 0 \\ \begin{cases} c(n+1) &\stackrel{n+1 \text{ par}}{=} \frac{2-(2i+2)}{2} = -i \\ c(n) &\stackrel{n \text{ impar}}{=} \frac{1-(2i+1)}{2} = -i \end{cases} \end{aligned}$$

Luego $c(n+1) = c(n)$

$$\begin{cases} d(n) &\stackrel{n \text{ impar}}{=} \frac{(2i+1)-1}{2} = i \\ d(n+1) &\stackrel{n+1 \text{ par}}{=} \frac{2i+2}{2} = i+1 \end{cases}$$

Y finalmente $d(n+1) = d(n) + 1$.

Necesitaremos también definir una serie de conjuntos de subsecuencias, una función que nos aplique esos conjuntos y, finalmente, un conjunto producto.

Empezamos definiendo estos conjuntos de subsecuencias:

1. $X^{[a,b]} \subset X^*$ el conjunto de todas las subsecuencias $s = (w, a, b)$ que empiezan en el tiempo a y terminan en el tiempo b . Vemos cómo incluimos a s en el conjunto X^* de las palabras, por la identificación que hacemos de s y w .

2. De forma similar $X^{J(n)} \subset X^*$ denotará el conjunto de todas las subsecuencias $s = (w, (c(n), d(n)))$ que empiezan en $c(n)$ y terminan en $d(n)$. Como $|J(n)| = n$, $X^{J(n)}$ es un producto de n copias de X . La diferencia entre X^n y $X^{J(n)}$ es cómo están etiquetadas

las coordenadas, en el primer caso el orden es $1, 2, \dots, n$ mientras que el segundo lo hace como $c(n), \dots, d(n)$. Por las propiedades P1 y P2 demostradas anteriormente, tenemos que:

Si n es impar:

$$X^{J(n+1)} = X^{J(n)} \times X$$

Si n es par:

$$X^{J(n+1)} = X \times X^{J(n)}$$

Ahora definimos la función H que envía subsecuencias en $X^{J(n)}$ a subsecuencias en \mathbb{R}^n . En el caso que $n = \infty$, H enviará secuencias bi-infinitas de $X^{\mathbb{Z}}$ a secuencias infinitas de $\mathbb{R}^{\mathbb{N}}$. El objetivo de esta función es reordenar los argumentos de las coordenadas y enviarlas a \mathbb{R} . Si $x \in X^{J(n)}$, entonces hay una subsecuencia $v \in \mathbb{R}^n$ que satisface $v_i = g(x_{h(i)}), \forall i \in \{1, \dots, n\}$. Recordamos que g es una función a la que solo imponemos que sea inyectiva, por lo tanto no es invertible. Definimos $H : X^{J(n)} \rightarrow \mathbb{R}^n$ como $H(x) = v$. Si $x = x_{c(n)}, \dots, x_{d(n)} \in X^{J(n)}$, entonces:

$$H(x_{c(n)}x_{c(n)+1} \dots x_{d(n)}) = g(x_0)g(x_1) \dots g(x_{|c(n)|+d(n)+1})$$

Esta función puede ser igualmente definida para el caso $x \in X^{\mathbb{Z}}$. Y aunque H no es invertible, porque g no es invertible, sí que tiene conjuntos inversos. Por ejemplo, si $A \subset \mathbb{R}^n$, entonces:

$$H^{-1}(A) = \{x \in X^{J(n)} | H(x) \in A\}$$

Finalmente, definimos el conjunto producto S como un subconjunto de $X^{[a,b]}$ de la siguiente forma $S = S_a \times \dots \times S_b$, donde $S_i \subset X$. Si $S' \subset X$, entonces $S \times S'$ es un producto contenido en $X^{[a,b+1]}$. De forma natural definimos el cilindro $C(S)$ de la siguiente manera:

$$C(S) = \{x \in X^{\mathbb{Z}} | x_i \in S_i, i \in [a, b]\}$$

De forma parecida a como definíamos A_s , tenemos que $C(S)$ es el conjunto de todas las secuencias de $X^{\mathbb{Z}}$ que coinciden con la subsecuencia en S .

A partir de $C(S)$ definimos también el cilindro desplazado $D(S)$ como:

$$D(S) = \{x \in X^{[a-1,b-1]} | x_i \in S_{i+1}, i+1 \in [a, b]\}.$$

Lema. Los cilindros C y D verifican:

1. $C(S) = C(S \times X) = C(X \times S)$
2. $D(S \times X) = D(S) \times X = X \times D(S)$

$$3. D(C(S)) = C(D(S))$$

Demostración.

Sea $S = S_a \times \dots \times S_b$.

1. Por definición, $\forall i \in \mathbb{Z}, x_i \in X$:

$$C(S) = \{x \in X^{\mathbb{Z}} | x_i \in S_i, i \in [a, b]\} = \{x \in X^{\mathbb{Z}} | x_{a-1} \in X \text{ y } x_i \in S_i, i \in [a, b]\} = C(X \times S)$$

De la misma manera:

$$C(S) = \{x \in X^{\mathbb{Z}} | x_i \in S_i, i \in [a, b]\} = \{x \in X^{\mathbb{Z}} | x_{b+1} \in X \text{ y } x_i \in S_i, i \in [a, b]\} = C(S \times X)$$

2. Teniendo en cuenta $S \times X \subset X^{[a, b+1]}$:

$$\begin{aligned} D(S \times X) &= \{x \in X^{[a-1, b]} | x_{b+1} \in X \text{ y } x_i \in S_{i+1}, i+1 \in [a, b]\} \\ &= \{x \in X^{[a-1, b-1]} | x_i \in S_{i+1}, i+1 \in [a, b]\} \times X \end{aligned}$$

De forma semejante a lo visto anteriormente también tenemos:

$$\begin{aligned} D(S \times X) &= \{x \in X^{[a-2, b-1]} | x_{a-2} \in X \text{ y } x_i \in S_{i+1}, i+1 \in [a, b]\} \\ &= X \times \{x \in X^{[a-1, b-1]} | x_i \in S_{i+1}, i+1 \in [a, b]\} \end{aligned}$$

3.

$$\begin{aligned} D(C(S)) &= \{x \in X^{\mathbb{Z}-1} | x_i \in S_{i+1}, \forall i+1 \in [a, b] \text{ y } x_j \in X, \forall j \notin [a, b]\} \\ &= \{x \in X^{\mathbb{Z}} | x_i \in S_{i+1}, \forall i+1 \in [a, b]\} = C(D(S)) \end{aligned}$$

cqd

Corolario. Sea $S = S_{c(n)} \times \dots \times S_{d(n)} \subset X^{J(n)}$ un conjunto producto, y sea

$$A = C(S) = \{x \in X^{\mathbb{Z}} | x_i \in S_i, i \in [c(n), d(n)]\}$$

Entonces tenemos

$$D(S) = \{x \in X^{[c(n)-1, d(n)-1]} | x_i \in S_{i+1}, i+1 \in [c(n), d(n)]\}$$

y

$$D(A) = \{x \in X^{\mathbb{Z}} | x_i \in S_{i+1}, i+1 \in [c(n), d(n)]\}$$

Además $D(A) = C(D(A))$ y, por el lema $D(A) = C(D(S) \times X)$.

Finalmente, observamos que para cada $S \in X^{[a, b]}$ existe un conjunto $S' \in X^{J(n)}$, donde $n = \min(-2a, 2b+1)$, tal que $C(S) = C(S')$. Como cada conjunto cilíndrico puede ser escrito como $C(S)$ para algún $S \in X^{[a, b]}$, esto significa que cada conjunto cilíndrico puede ser escrito como $C(S')$ para algún $S' \in X^{J(n)}$.

Finalmente, ya estamos en disposición de demostrar nuestra variante del TEK.

Teorema. Suponemos que tenemos una secuencia de medidas v_n en $X^{J(n)}$ que satisfacen que $\forall n$ y para todo $S \subset X^{J_n}$, si n es impar,

$$v_n(S) = v_{n+1}(S \times X) \quad (\text{T.1})$$

y si n es par,

$$v_n(S) = v_{n+1}(X \times S) \quad (\text{T.2})$$

$$v_n(S) = v_{n+1}(D(S) \times X), \quad (\text{T.3})$$

Entonces existe un único proceso estacionario $Q = (X^{\mathbb{Z}}, \mathbb{X}, P)$ tal que, $\forall n, S \subset X^{J_n}$,

$$P(C(S)) = v_n(S)$$

La demostración la haremos en dos partes: primero demostraremos que P existe y, posteriormente comprobaremos que P es estacionario.

Demostración. Definimos una secuencia de medidas μ_n en $\mathbb{R}^{\mathbb{N}}$ de la siguiente manera: si $A \subset \mathbb{R}^{\mathbb{N}}$, entonces definimos $\mu_n(A) = v_n(H^{-1}(A)) = v_n\{x \in X^{J(n)} | H(x) \in A\}$, por la definición de anticonjuntos que establece H .

La consistencia de estas μ 's se puede comprobar:

$$\mu_{n+1}(A \times \mathbb{R}) = v_{n+1}(H^{-1}(A \times \mathbb{R}))$$

Dado que

$$H^{-1}(A \times \mathbb{R}) = \begin{cases} H^{-1}(A) \times X & n \text{ impar} \\ X \times H^{-1}(A) & n \text{ par} \end{cases}$$

Obtenemos para μ_{n+1}

$$\mu_{n+1}(A \times \mathbb{R}) = \begin{cases} v_{n+1}(H^{-1}(A) \times X) & n \text{ impar} \\ v_{n+1}(X \times H^{-1}(A)) & n \text{ par} \end{cases}$$

Ahora, aplicando T.1 y T.2 a la derecha de las igualdades

$$\mu_{n+1}(A \times \mathbb{R}) = v_n(H^{-1}(A)) = \mu_n(A).$$

Aplicando TEK sobre las μ_n 's, obtenemos una única P' en $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}(\mathbb{R}^{\mathbb{N}}))$ que concuerda con μ_n . Por lo tanto, si $A = A_1 \times \dots \times A_n$, $\forall i \in \{1, \dots, n\}$, entonces existe una única probabilidad P' que concuerda con μ_n

$$P'(x | x_i \in A_i, i \in \{1, \dots, n\}) = \mu_n(A)$$

Ahora pasamos a demostrar la segunda parte del Teorema. Definimos P tq $\forall S \in \mathbb{X}$,

$$P(S) = P'(H(x)|x \in S)$$

Probamos primero que P existe. Sea $S \in X^{J(n)}$, tenemos

$$v_n(S) = \mu_n(H(S)) = P'(a \in \mathbb{R}^{\mathbb{N}} | a_i \dots a_n \in H(S)) = P(x \in X^{\mathbb{Z}} | x_{c(n)} \dots x_{d(n)} \in S) = P(C(S))$$

Por lo tanto, esta P es en realidad una extensión de las μ_n 's y existe.

Para ver que P es estacionaria, necesitamos demostrar que $P(D(A)) = P(A)$ para cualquier $A \in \mathbb{X}$ que contenga al cilindro $C(S)$. Llamaremos a esta colección Θ

$$\Theta = \{C(S) | S \subset X^{J(n)} \text{ para algún } n\}$$

Aunque cada $A \in \Theta$ pueden tener asociado más de un par n, S siempre escogeremos el de menor n -que tiene asociado una única S -.
Si n es par, T.3 y el lema demostrado anteriormente nos da

$$\begin{aligned} P(A) &= v_n(S) = v_{n+1}(D(S) \times X) = P(C(D(S)) \times X) \\ &= P(C(D(S))) = P(D(A)) \end{aligned}$$

Y por último, vemos que ocurre lo mismo con n impar

$$\begin{aligned} P(A) &= v_n(S) = v_{n+1}(S \times X) = v_{n+1}(D(S \times X) \times X) = P(C(D(S \times X) \times X)) \\ &= P(C(D(S \times X) \times X)) = P(C(D(S) \times X \times X)) = P(C(D(S) \times X) \times X) \\ &= P(D(A) \times X) = P(D(A)) \end{aligned}$$

Por lo tanto, $P(A) = P(D(A))$. *cqd*

Finalmente estamos en disposición de demostrar el nuestro teorema de partida.

Teorema Dado $f : X^* \rightarrow [0, 1]$ que satisface:

$$1. f(\emptyset) = 1$$

$$2. \forall w \in X^*, f(w) = \sum_{z \in X^*} f(zw) = \sum_{z \in X^*} f(wz)$$

existe un único proceso $Q = (X^{\mathbb{Z}}, \mathbb{X}, P)$ tq $\forall w \in X^*, P(w) = f(w)$.

Demostración. Primero construiremos v_n 's que satisfagan T.1, T.2 y T.3 para a continuación aplicar el Teorema que acabamos de demostrar.

$\forall w \in X^*$ y $|w| = n$, sea $S = \{w\}$; entonces definimos

$$v_n(S) = f(w).$$

Para un $S \subset X^{J(n)}$, S es una unión disjunta de conjuntos de la forma $S_w = \{w\}$, para cada $w \in S$. Podemos entonces definir v_n en todos los subconjuntos de $X^{J(n)}$

$$v_n(S) = \sum_{w \in S} v_n(S_w) = \sum_{w \in S} f(w).$$

Para que las v_n 's sean realmente una medida de probabilidad han de cumplir $v_n(X^{J(n)}) = 1$. Por inducción, si $n = 0$ entonces $X^{J(n)} = \emptyset$ y, por lo tanto $f(\emptyset) = 1$. Para los próximos pasos de la inducción debemos usar la condición 2 del Teorema, separando los casos par e impar.

Suponemos que (1) se cumplen en n impar, entonces

$$\begin{aligned} v_{n+1}(X^{J(n+2)}) &= \sum_{w \in X^{J(n)}} f(w) \\ &= \sum_{w \in X^{J(n+1)}} \sum_{x \in X} f(wx) = \sum_{w \in X^{J(n+1)}} f(w) \\ &= v_n(X^{J(n)}) = 1. \end{aligned}$$

En el caso de n par se hace de forma semejante

$$\begin{aligned} v_{n+1}(X^{J(n+2)}) &= \sum_{w \in X^{J(n)}} f(w) \\ &= \sum_{w \in X^{J(n+1)}} \sum_{x \in X} f(xw) = \sum_{w \in X^{J(n+1)}} f(w) \\ &= v_n(X^{J(n)}) = 1. \end{aligned}$$

Vemos que se satisface T.1 para n impar

$$\begin{aligned} v_n(S) &= \sum_{w \in S} f(w) = \sum_{w \in S} \sum_{x \in X} f(wx) \\ &= \sum_{z \in (S \times X)} f(z) \\ &= v_{n+1}(S \times X). \end{aligned}$$

Para n par, comprobamos T.2

$$\begin{aligned} v_n(S) &= \sum_{w \in S} f(w) = \sum_{w \in S} \sum_{x \in X} f(xw) \\ &= \sum_{z \in (X \times S)} f(z) = v_{n+1}(X \times S) \end{aligned}$$

Por último, asumiendo que n es par vemos T.3

$$\begin{aligned} v_n(S) &= \sum_{w \in S} f(w) = \sum_{w \in S} \sum_{x \in X} f(wx) \\ &= \sum_{z \in (S \times X)} f(z). \end{aligned}$$

Recordamos que, mientras que $X \times S \subset X^{J(n)}$, no ocurre lo mismo con $S \times X \not\subset X^{J(n)}$, sin embargo sí sabemos $D(S \times X) = D(S) \times X \subset X^{J(n)}$, y f no depende del índice temporal, así que

$$v_n(S) = \sum_{z \in (S \times X)} f(z) = \sum_{z \in (D(S) \times X)} f(z) = v_{n+1}(D(S) \times X).$$

Ya podemos aplicar el Teorema anterior para ver que $\exists!$ proceso $Q = (X^{\mathbb{Z}}, \mathbb{X}, P)$ tal que, $\forall n, S \in X^{J(n)}$, tenemos que $P(C(S)) = v_n(S)$. Por lo tanto, si definimos $S = \{w\}$ para cualquier longitud n y palabra w , se cumple

$$f(w) = v_n(S) = P(C(S)) = P(w)$$

cqd

Ejemplo de Proceso

Para ilustrar el concepto de proceso, pondremos un pequeño ejemplo basado en un dado ‘ideal’ de seis caras, $X = \{1, 2, 3, 4, 5, 6\}$ y cualquier palabra de longitud l tendrá asociada una probabilidad de $P(w) = \left(\frac{1}{6}\right)^l$. En particular, si la palabra es de longitud $l = 0$, entonces $P(w) = \left(\frac{1}{6}\right)^0 = 1$. Si además tenemos otra palabra $z \in X$ entonces, a partir de (1) vemos que $f(wz) = f(zw) = \left(\frac{1}{6}\right)^l \sum_{i=1}^6 \frac{1}{6} = \left(\frac{1}{6}\right)^l = f(w)$. Por lo tanto cumple tanto (1) como (2) y es un proceso. *cqd*

Estados de un Proceso

Supongamos que tenemos un proceso $Q = (X^{\mathbb{Z}}, \mathbb{X}, P)$ del que conocemos algunos símbolos recientes. Si $s = (w, a, b)$ es una subsecuencia conocida, con $b = -1$, s induce una distribución condicional en un espacio futuro (DCF) del tipo $P(\cdot|s) = P(\cdot|A_s)$. Esta DCF condensa todo nuestro conocimiento sobre las posibilidades futuras del proceso Q , por lo que al conocer esta distribución, en cierto sentido, podemos olvidar todo lo acontecido en el proceso anteriormente. Esta es la base con la que decimos que $P(\cdot|s)$ es un estado.

Ahora ampliaremos esta noción. Primero definimos $\mathbb{X}^+ = \{x^+ = x_0, x_1, \dots | \forall i \geq 0, x_i \in X\}$, y respectivamente $\mathbb{X}^- = \{x^- = \dots x_{-2}, x_{-1} | \forall i \leq 0, x_i \in X\}$. Los campos σ generados por \mathbb{X}^+

y \mathbb{X}^- son \mathbb{F} y \mathbb{H} , respectivamente. Las DCF están condicionadas por secuencias $s = (w, a, b)$ con w palabra y $w \in \mathbb{X}^-$, por lo que estas subsecuencias pueden o no acabar en la última observación conocida (también conocidas como *historias*), pero nos centraremos únicamente en este último caso.

Además, para lo que veremos a continuación necesitamos definir una serie de σ -álgebras $\mathbb{H}_l, l \in (0, 1, \dots)$, donde \mathbb{H}_l es la σ -álgebra en $X^{\mathbb{Z}}$ generada por todas las historias de longitud l . En definitiva, un evento en \mathbb{H}_l es el conjunto de series infinitas x condicionadas por x_{-l}, \dots, x_{-1} . Además, de forma trivial $\forall l \leq 0, \mathbb{H}_l \subset \mathbb{H}_{l+1}$. A la unión infinita de todas las \mathbb{H}_l la denominamos \mathbb{H} .

Si dos historias, z e y acaban determinando la misma DCF $P(\cdot|z) = P(\cdot|y)$, entonces decimos que $y \sim z$ son equivalentes, es decir podemos definir la clase de equivalencia \sim y $\bar{z} = \{y \in X^* | y \sim z\}$.

Sea π la función proyección sobre las clases de equivalencia de \mathbb{X}^- , es decir $\pi(z) = \bar{z}$. Para toda palabra futura w se cumplirá entonces $P(w|z) = P(w|\bar{z})$.

Utilizando las clases de equivalencia podemos olvidarnos de la historia concreta que ha generado la DCF. Así, imaginemos la historia w , cuya clase de equivalencia \bar{w} induce una DCF sobre x_0 . Así, si escogemos un símbolo k adecuado para x_0 , tenemos la nueva palabra $z = wk$, y tenemos entonces de nuevo una clase de equivalencia \bar{z} que induirá una DCF sobre x_1 , a las diferentes clases de equivalencia \bar{z} que se obtienen a partir de las DCF inducidas por \bar{w} las llamamos *estados del proceso*.

Usando el ejemplo del dado, en este caso todas las historias son equivalentes, luego hay un solo estado del proceso, el cual es visitado recurrentemente. Daremos a continuación una definición más formal de los estados de un proceso.

Sea $s_l = x_{-l} \dots x_{-1}$, con $l > 0$ y $R_l = \{s_l | P(s_l) > 0\}$. Dada la historia s_l , para cualquier evento futuro A :

$$P(A|s_l) = \frac{P(A \cap A_{s_l})}{P(A_{s_l})},$$

En concreto, si $x^- = \lim_{l \rightarrow \infty} s_l$, entonces tenemos

$$P(A|x^-) = \lim_{l \rightarrow \infty} P(A|s_l) \quad (1)$$

Si definimos:

$$1_A = \begin{cases} 1 & \text{if si } x \in A \\ 0 & \text{en caso contrario,} \end{cases}$$

$P(A|s_l)$ es la esperanza condicional $E(1_A|\mathbb{H})(x)$ para cualquier secuencia x que coincide con w . Por lo tanto podemos escribir (1) como $\lim_{l \rightarrow \infty} E(1_A|\mathbb{H})$. Este límite converge casi seguramente a $E(1_A|\mathbb{H})(x)$, mediante el Teorema de convergencia de Martingalas, el cual demostramos al final del capítulo, para no entorpecer la lectura. De esta manera tendremos que $(1)=P(\cdot|x^-)$.

Si x^- tiene una historia s_l con probabilidad nula, entonces $P(A|x^-) = 0$. Denominaremos al conjunto de historias con probabilidad nula $N = \{x^- \in X^{\mathbb{Z}} | P(x^-) = 0\}$.

Un estado del proceso es un DCF que se da al condicionar en una historia de probabilidad no nula. Si $R = \bigcup_{l=0}^{\infty} R_l$, entonces el conjunto de todos los estados del proceso para un proceso dado es:

$$PS = \{P(\cdot|w)|w \in R\} \cup \{P(\cdot|x^-) \in \mathbb{X}^- \setminus N\}$$

Además, definimos la función $G : \mathbb{X}^- \rightarrow PS$ que envía historias no nulas al DCF que inducen, $G(x^-) = P(\cdot|x^-)$. El conjunto inducido de historias de un estado del proceso A es un conjunto que particuliza G mediante $G^{-1}(A)$, y que resulta en todas las historias que induce A .

Proposición. Para cualquier estado del proceso A , $G^{-1}(A) \times \mathbb{X}^+$ es la intersección de un conjunto medible y de un conjunto de medida completa (complementario de medida nula)

Usaremos l para las longitudes de las historias y k para longitudes futuras. Además diremos que $A(w)$ es la probabilidad que asigna A a la palabra w .

Además, diremos que dos estados del proceso B y C son k -equivalentes si, para todas palabras w tq $|w| \leq k$, $B(w) = C(w)$.

Demostración Para cada natural l y palabra w , y para cada historia x^- , definimos:

$$f_l^w(x^-) = P(w|x_{-l} \dots x_{-1})$$

y

$$f^w(x^-) = P(w|x^-),$$

Por el *Teorema de Convergencia de Martingalas* $\lim_{l \rightarrow \infty} f_l^w(x^-) = f^w(x^-)$ casi seguramente. Como f_l^w es medible, su límite también lo será.

Sea A un estado del proceso y $J_k(A)$ el conjunto de historias x^- que inducen estados del proceso $G(x^-)$ que son k -equivalentes a A , $J_k(A) = \{x^- : \forall w \text{ tq } |w| \leq k, P(w|x^-) = A(w)\}$. Es decir:

$$J_k(A) = \bigcap_{w:|w|=k} \{x^- : P(w|x^-) = A(w)\} = \bigcap_{w:|w|=k} (f^w)^{-1}(A(w)).$$

$\{A(w)\}$ es un conjunto medible y f^w es una función medible, así que su la antiimagen $(f^w)^{-1}(A)$ ha de ser medible. Por lo que $J_k(A)$ es intersección finita de conjuntos medibles y por lo tanto medible. Sea $J(A) = \bigcap_{k=1}^{\infty} J_k(A)$, $J(A)$ es intersección numerable de conjuntos medibles, por lo que es medible. Ahora, cada historia $J(A) \cap (X^{\mathbb{Z}} \setminus N)$ induce un estado A , con $X^{\mathbb{Z}} \setminus N$ un conjunto de medida completa, y no hay ninguna historia que no esté en $J(A)$ que induzca A . Por lo tanto $J(A) = \bigcap_{k=1}^{\infty} J_k(A) = G^{-1}(A)$. *cqd*

La idea central es que un estado es una predicción basada en lo sabido sobre el pasado del proceso. En la siguiente sección clasificaremos los diferentes tipos de estados que podemos encontrar.

Martingalas

En esta sección presentamos los resultados relativos a convergencia de martingalas necesarios para el apartado anterior.

Sea X una variable aleatoria en un espacio de probabilidad (Ω, \mathcal{F}, P) . Sea \mathcal{E} una sub σ -álgebra de \mathcal{F} .

Definición. La esperanza condicional de X respecto a \mathcal{E} es cualquier variable aleatoria Y que satisfaga

1. Y es medible respecto a \mathcal{E}
2. $\forall A \in \mathcal{E}$,

$$\int_A X dP = \int_A Y dP$$

Notar que las esperanzas condicionales que cumplen esta definición son únicas en conjuntos de probabilidad no nula.

Si Z es una variable aleatoria en (Ω, \mathcal{F}, P) , entonces Z induce una σ -álgebra $\sigma(Z)$ en Ω , en concreto la menor σ -álgebra que contiene conjuntos de la forma $Z^{-1}(B)$, para conjuntos de Borel B . Sea A un conjunto en \mathcal{F} con función característica 1_A , tal y como hemos definido previamente. Si fijamos $a \in \mathbb{R}$ y evaluamos en $x \in Z^{-1}(a)$ la esperanza condicional $\mathbb{E}(1_A | \sigma(Z))$, encontramos que

$$\mathbb{E}(1_A | \sigma(Z))(x) = P(x \in A | Z = a)$$

Una filtración es un conjunto incremental de σ -álgebras $\{\mathcal{F}\} = \mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$. Por ejemplo, los conjuntos \mathbb{H}_t en $(X^{\mathbb{Z}}, \mathbb{X}, P)$ son filtraciones.

Definición. Una secuencia $\{X\} = X_1, X_2 \dots$ es una martingala respecto a la filtración $\{\mathcal{F}\}$ si

1. $\forall i, X_i$ es medible respecto a \mathcal{F}_i .
2. $\forall s, t \in \mathbb{N}$ tal que $t > s$, $X_s = \mathbb{E}(X_t | \mathcal{F}_s)$.

iii COMPLETAR APARTADO !!!

Tipos de Estado

La definición del conjunto de estados de un proceso nos indica que estos están inducidos por historias, ya sean finitas o infinitas.

Definimos antes de describir la tipología de estados, y algunas de sus características, una propiedad asociada a los estados. Consideramos por ejemplo un estado del proceso $A \in PS$, y consideramos $G^{-1}(A)$, el conjunto de historias que induce A . Podemos asignar a este conjunto una probabilidad $P(A) = P(G^{-1}(A))$. En particular, ante la historia w , $P(A)$ será la probabilidad de que s induzca el estado del proceso A .

Un estado del proceso está *infinitamente precedido* si la historia w que lo genera es de longitud infinita, $|w| = \infty$ y $P(w) > 0$.

Un estado del proceso es *alcanzable* si al menos existe una historia, de longitud cualquiera, w con $P(w) > 0$ que induce este estado.

Un estado del proceso $A \in PS$ es recurrente si $P(A) > 0$.

Proposición.

1. Cada estado del proceso inalcanzable está infinitamente precedido.
2. Cada estado del proceso recurrente está infinitamente precedido.

Demostración.

1. Por definición de estado de un proceso, ya que cada uno de ellos está inducido por alguna historia.
2. Sea A un estado recurrente. Entonces $G^{-1}(A)$ tiene probabilidad no nula y por lo tanto no es vacío. Como $G^{-1}(A)$ contiene solo historias con probabilidad no nula, entonces contiene alguna historia que induce A de probabilidad no nula. *cqd*

A partir de esta proposición podemos asegurar que tres tipos de estado no ocurren:

- Estados inalcanzables que no están infinitamente precedidos.
- Estado alcanzables recurrentes que no están infinitamente precedidos.
- Estados inalcanzables que no son recurrentes ni infinitamente precedidos.

Nos quedan por definir dos tipos de estado.

Un estado del proceso es *transitivo* si es alcanzable pero no recurrente, y si además no está infinitamente precedido se dice que *estrictamente transitivo*.

Un estado del proceso es *elusivo* si es inalcanzable y no recurrente. Este tipo de estado con frecuencia puede ser ignorado, ya que únicamente pueden ser inducidos por eventos de

probabilidad nula, y por lo tanto un conjunto medible de estados elusivos pueden obviarse sin cambiar la medida de probabilidad del proceso. Sin embargo, hay procesos cuyos estados elusivos vienen inducidos por los estados recurrentes, o procesos en los que estos estados elusivos son no numerables y tienen probabilidad positiva. En estos casos tenemos que tenerlos en cuenta.

Incluir tabla de estados

Ahora imaginemos que hay dos observadores: el observador 1 conoce menos información histórica que el observador 2. Puede ocurrir que esta diferencia de conocimiento permita al observador 2 conocer mejor el futuro del proceso, pero también puede ocurrir que esta información extra sea irrelevante. Si eso ocurre decimos que estamos *sincronizados*.

Una palabra w es una *palabra sincronizante* si toda historia de probabilidad no nula que termina en w induce el mismo estado del proceso.

Proposición. Si w es una palabra sincronizante, entonces induce un estado del proceso que es alcanzable y recurrente.

Demostración. Sea A este estado del proceso inducido por w . Como A es inducido por una historia con probabilidad no nula, entonces es alcanzable. w es en sí misma una historia que termina en w , y A es el estado del proceso inducido por todas las historias que terminan en w . Además, como w es una palabra no nula, y en conjunto de historias que contienen w es un subconjunto de $G^{-1}(A)$, tenemos:

$$P(A) \geq P(\text{historias que terminan en } w) \geq P(w) > 0.$$

Por lo tanto A es recurrente. *cqd*

Cadenas de Markov.

Una *Cadena de Markov (MC)* de n -estados es un triplete (S, A, π) , donde S es un conjunto de orden n , A es una matriz cuadrada de dimensión n , y $\pi = \{\pi_1, \dots, \pi_n\}$ es un vector de longitud n , y cumplen las siguientes propiedades:

1. Cada columna de A tiene suma 1.
2. $\sum_i \pi_i = 1$.
3. $\pi A = \pi$.

Los elementos de S son los estados, A es la matriz de transición y π es la distribución estacionaria sobre los estados S .

Si \mathbb{S} es la σ -álgebra definida por el conjunto de los subconjuntos de $S^{\mathbb{Z}}$, entonces $(S^{\mathbb{Z}}, \mathbb{S})$ será nuestro espacio medible, y definimos nuestra medida de probabilidades P de manera que si $v = v_0 v_1 \dots v_{l-1}$, con $v_i \in S$,

$$P(v) = \pi_{v_0} a_{v_0 v_1} \dots a_{v_{l-2} v_{l-1}},$$

Definiendo $P(\emptyset) = 0$, vemos que una MC es un proceso estacionario.

Recordamos que un proceso es estacionario si cumple:

$$\sum_{z \in S} P(wz) = \sum_{z \in S} P(zw) = P(w) \quad (1)$$

$$P(\emptyset) = 1 \quad (2),$$

La condición (2) se cumple trivialmente. Para ver la (1) sea $z \in S$, entonces

$$P(\emptyset z) = P(z\emptyset) = P(z) = \pi_z.$$

Por tanto, aplicando la propiedad 2 de las MC

$$\sum_{z \in S} P(z\emptyset) = \sum_{z \in S} P(\emptyset z) = \sum_{z \in S} \pi_z \stackrel{P2MC}{=} 1 = P(\emptyset)$$

Ahora sea la probabilidad de la unión de v y z , y teniendo en cuenta la primera propiedad de las MC

$$P(vz) = P(v)a_{v_{l-1}v} \rightarrow \sum_{z \in S} P(vz) = P(v) \sum_{z \in S} a_{v_{l-1}z} \stackrel{P1MC}{=} P(v) \cdot 1 = P(v)$$

En el otro sentido, la unión de z con v nos da

$$P(zv) = \pi_z a_{zv_0} \dots a_{v_{l-2}v_{l-1}}$$

Por lo tanto

$$\sum_{z \in S} P(zv) = \left(\sum_{z \in S} \pi_z a_{zv_0} \right) a_{zv_0} \dots a_{v_{l-2}v_{l-1}}$$

Finalmente, aplicando la tercera propiedad de las MC, si tenemos en cuenta que $\sum_{z \in S} \pi_z a_{zv_0}$ es el elemento v_0 del vector πA

$$\sum_{z \in S} \pi_z a_{zv_0} \stackrel{P3MC}{=} \pi_{v_0}$$

Finalmente vemos que

$$\sum_{z \in S} P(zv) = \pi_{v_0} a_{v_0v_1} \dots a_{v_{l-2}v_{l-1}} = P(v)$$

concluimos que (S, A, π) define un proceso en $(S^{\mathbb{Z}}, \mathbb{S}, P)$, *cqd*

Nota: hasta ahora hemos utilizado una notación particular referida a la teoría de secuencias y procesos, a partir de ahora introduciremos notación estándar de estadística para términos que hasta ahora denominábamos *palabras*, *estados del proceso*...

Las MC's pueden considerarse procesos "sin memoria", es decir, que la distribución de probabilidad futura para $X = (X_1, \dots, X_n)$ secuencia de variables aleatorias, depende únicamente del valor actual, más formalmente:

$$P(X_{n+1} = x_{n+1} | X_n = x_n \dots X_1 = x_1) = P(X_{n+1} = x_{n+1} | X_n = x_n)$$

Esta identidad es la denominada *propiedad de Márkov*.

2. Modelos Ocultos de Markov (HMM)

En los Modelos ocultos de Markov (HMM) la MC subyace tras las observaciones. Los estados del HMM solo pueden ser inferidos a partir de los símbolos registrados. Correlacionando las observaciones y las transiciones de estado lo que se busca es encontrar el estado de secuencias más probable.

Un HMM pertenece a los mecanismos de *Machine Learning*, y puede ser entendido como una especie de doble proceso estocástico:

- El primer proceso estocástico es un conjunto finito de estados, cada uno de ellos generalmente asociado a una distribución de probabilidad multidimensional. Mediante la denominada matriz de transición se controla las transiciones entre estados.
- El segundo proceso estocástico es aquel en que cualquier estado puede ser observado, es decir, analizaremos lo observado sin ver en qué estado está ocurriendo, de aquí el epíteto *oculto* que define a este modelo.

Componentes de un HMM

Para definir completamente un HMM, se necesitan cinco elementos:

1. Los N estados del modelo $S = \{S_1, \dots, S_n\}$
2. Las M observaciones de los diferentes símbolos por estado $V = \{V_1, \dots, V_M\}$. Si las observaciones son continuas, obviamente M es infinito.
3. La matriz de transición $A = \{a_{ij}\}$, donde $a_{ij} = P(q_{t+1} = j | q_t = i)$, siendo q_t el estado actual. Esta matriz es equivalente a la vista en la definición de las MC. Hay que observar que si uno de los a_{ij} es definido cero, permanecerá cero durante todo el proceso de entrenamiento.
4. La probabilidad de distribución de observación de los símbolos en cada estado, $B = \{b_j(k)\}$, donde $b_j(k)$ es la probabilidad de observar el símbolo v_k en el estado S_j .

$$b_j(k) = P(o_t = v_k | q_t = j), \forall j \in \{1, \dots, N\}, \forall k \in \{1, \dots, M\}$$

Donde v_k denota el k -ésimo símbolo observado en el alfabeto, y o_t el actual vector de parámetros.

Se deben satisfacer ciertas restricciones:

$$b_j(k) \geq 0, \forall j \in \{1, \dots, N\}, \forall k \in \{1, \dots, M\} \text{ y además } \sum_{k=1}^M b_j(k) = 1, \forall j \in \{1, \dots, N\}$$

En el caso continuo tendremos una función de densidad continua, en vez de un conjunto de probabilidades discretas. En este caso lo que especificamos es el conjunto de parámetros de la función de densidad, aproximada como una suma ponderada de M distribuciones Gaussianas (GHMM),

$$b_j(o_t) = \sum_{m=1}^M c_{jm} N(\mu_{jm}, \Sigma_{jm}, o_t)$$

donde c_{jm} es son los pesos, μ_{jm} es el vector de medias, y Σ_{jm} las matrices de covarianza. Observar que c_{jm} debe satisfacer las condiciones estocásticas $c_{jm} \geq 0, \forall m \in \{1, \dots, M\}$ y

$$\sum_{m=1}^M c_{jm} = 1, \forall j \in \{1, \dots, N\}$$

5. La distribución inicial de estados $\pi = \{\pi_i\}$, donde π_i es la probabilidad de que el modelo está en el estado S_i en el tiempo inicial $t = 0$, con

$$\pi_i = P(q_1 = i), \forall i \in \{1, \dots, N\}$$

Para denotar los parámetros de un HMM con frecuencia se usa:

$$\lambda = (A, B, \pi)$$

Para denotar distribuciones discretas, o bien:

$$\lambda = (A, c_{jm}, \mu_{jm}, \Sigma_{jm}, \pi)$$

Cuando se trata de distribuciones continuas asociadas a funciones de densidad.

Arquitectura de los HMM

???

Algoritmos Relacionados con los HMM's

El objetivo de un HMM, a grosso modo, es aprender a clasificar los datos proporcionados. Es decir, un HMM tiene que ser capaz de discernir las diferencias de comportamiento que posean los diferentes estados del *stream* de observaciones con el que lo alimentamos.

En la historia de los HMM han destacado los estudios de tres problemas:

1. Problema de Evaluación

Dada la secuencia de observaciones $O = \{o_1, \dots, o_m\}$, ¿Cuál es la probabilidad de que O haya sido generada por el modelo $P(O|\lambda)$? Dado un λ .

2. Problema de Decodificación

Dada la secuencia de observaciones $O = \{o_1, \dots, o_m\}$, ¿Cuál es la secuencia de estados más probable dado el modelo λ ?

3. Problema de Aprendizaje

Dada la secuencia de observaciones $O = \{o_1, \dots, o_m\}$, ¿Cómo debemos ajustar los parámetros de λ para maximizar $P(O|\lambda)$?

El problema de evaluación es la piedra angular en muchos estudios de reconocimiento de voz. El problema de decodificación resulta útil a la hora de segmentar, y el problema de aprendizaje debe ser resuelto si queremos entrenar un HMM para su uso en labores de reconocimiento.

Problema de Evaluación: Forward Algorithm

Dada una secuencia de observaciones $O = \{o_1, \dots, o_T\}$ y un modelo $\lambda = (A, B, \pi)$, se trata de averiguar $P(O|\lambda)$. Lo primero que uno se plantea es que este cálculo se puede realizar directamente mediante la fórmula de Bayes: $P(O|\lambda) = \frac{P(\lambda|O)P(O)}{P(\lambda)}$. Sin embargo, el número de operaciones requeridas es del orden de N^T , lo que no resulta efectivo computacionalmente. Existe sin embargo un método de menor complejidad que involucra la utilización de una variable auxiliar *forward* (α) de la forma

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = i | \lambda)$$

Para el cálculo de esta variable podemos definir una relación recursiva de la forma

$$\begin{cases} \alpha_1(j) = \pi_j b_j(o_1), \forall j \in \{1, \dots, N\} \\ \alpha_{t+1}(j) = b_j(o_{t+1}) \sum_{i=1}^N \alpha_t(i) a_{ij}; \forall j \in \{1, \dots, N\}, \forall t \in \{1, \dots, T-1\} \end{cases}$$

Las α_T 's se pueden calcular utilizando la recursión. Así que la probabilidad requerida es dada por

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

Este método es conocido comúnmente como *forward algorithm*.

De forma similar podemos definir una variable *backward* $\beta_t(i)$

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | q_t = i, \lambda)$$

Como el estado actual es i , $\beta_t(i)$ es la probabilidad de tener en t el estado i sabiendo la historia futura parcial $o_{t+1}, o_{t+2}, \dots, o_T$.

Igual que antes, planteamos la fórmula recursiva que nos dará la variable *backward*

$$\begin{cases} \beta_t(i) = \sum_{j=1}^N \beta_{t+1}(j) a_{ij} b_j(o_{t+1}); \forall i \in \{1, \dots, N\}, \forall t \in \{1, \dots, T-1\} \\ \beta_T(i) = 1, \forall i \in \{1, \dots, N\} \end{cases}$$

Combinando ambas obtenemos

$$\alpha_t(i) \beta_t(i) = P(O, q_t = i | \lambda) \quad \forall i \in \{1, \dots, N\}, \forall t \in \{1, \dots, T\}$$

Finalmente, obtenemos que

$$P(O | \lambda) = \sum_{i=1}^N P(O, q_t = i | \lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i)$$

Más adelante desarrollaremos de forma detallada esta combinación *forward-backward* para resolver el *Problema de Aprendizaje* mediante el algoritmo *Baum-Welch*.

Nota: en cada caso, tanto en las variables backward como forward hay que realizar una normalización en cada tiempo t .

Problema de Decodificación: Algoritmo de Viterbi

El problema de decodificación gira en torno a saber cuál es la secuencia de estados más probable dada una secuencia de observaciones $O = o_1, o_2, \dots, o_T$ y un modelo $\lambda = (A, B, \pi)$.

El problema de buscar esta secuencia más probables es que, muchas veces, quedarnos solo con las secuencias de q_t 's más probables resulta en una secuencia de estados poco significativa.

Por ello introducimos el denominado *Algoritmo de Viterbi*, que nos permite encontrar el estado de secuencias que tiene mayor verosimilitud.

Definimos primero

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2, \dots, q_{t-1}, q_t = i, o_1, o_2, \dots, o_{t-1} | \lambda),$$

Es decir, la probabilidad más alta de la secuencia parcial de estados y observaciones hasta t , dado el estado actual i .

A partir de aquí, siendo $\delta_1(j) = \pi_j b_j(o_1)$, $\forall j \in \{1, \dots, N\}$, definimos la fórmula de recursividad

$$\delta_{t+1}(j) = b_j(o_{t+1}) \left[\max_{1 \leq i \leq N} \delta_t(i) a_{ij} \right] \quad \forall i \in \{1, \dots, N\}, \forall t \in \{1, \dots, T-1\}$$

Esto significa que empezamos nuestro cálculo desde $\delta_T(j), \forall j \in \{1, \dots, N\}$. Mantenemos un puntero al estado “ganador”, que será $j^* = \arg \max_{1 \leq i \leq N} \delta_T(j)$. A partir de aquí realizamos un *back-track* en la secuencia redefiniendo j^* en cada paso, obteniendo de esta manera el conjunto de estados que se pide.

Problema de Aprendizaje: Algoritmo Baum-Welch

El problema de aprendizaje se centra en cómo podemos ajustar los parámetros del HMM para que se ajusten de la mejor forma a las observaciones.

Uno de los criterios más utilizados es el de máxima verosimilitud. Dado el HMM de parámetros λ y las observaciones O . La verosimilitud puede ser expresada como

$$L = P(O|\lambda)$$

Conocer el modelo que maximiza $\lambda = (A, B, \pi)$ es una operación que no puede ser resuelta analíticamente. Existen métodos iterativos, como *Baum-Welch* o métodos basados en gradientes, para localizar localmente máximos apropiados para ser parámetros del modelo.

Algoritmo Baum-Welch

El algoritmo Baum-Welch, también conocido como *Forward-Backward*, ya que utiliza estas variables que definimos previamente al tratar el problema de evaluación, utiliza una cantidad auxiliar $Q(\lambda, \lambda')$ para comparar dos modelos λ y λ' .

$$Q(\lambda, \lambda') = \sum_q P(q|O, \lambda) \log[P(O, q, \lambda')]$$

También definimos unas nuevas variable a partir de las *backward* y *forward*.

La primera, referida habitualmente en la literatura como probabilidades *smoothed*, es la probabilidad de estar en el estado i en el instante t dadas las observaciones O y el modelo λ . Se obtiene combinando α y β de acuerdo al Teorema de Bayes.

$$\gamma_t(i) = P(q_t(i)|O, \gamma) = \frac{\alpha_t(i)\beta_t(i)}{\sum_{i=1}^N \alpha_t(i)\beta_t(i)}$$

Observamos que, si sumamos cada $\gamma_t(i)$ en todos los instantes de tiempo, obtenemos el número de transiciones que se realizan desde i : $\sum_{t=1}^{T-1} \gamma_t(i)$.

La segunda variable nueva, también combina las variables *backward* y *forward* para definir la probabilidad de estar en el estado i en el instante t , y en el estado j en el instante $t + 1$; dadas las observaciones O y el modelo γ .

$$\xi_t(i, j) = \frac{P(q_t = i, q_{t+1} = j, O|\lambda)}{P(O|\lambda)} = \frac{\alpha_i(t)a_{ij}\beta_j(t+1)b_j(y_{t+1})}{\sum_{k=1}^N \sum_{l=1}^N \alpha_k(t)a_{kl}\beta_l(t+1)b_l(y_{t+1})}$$

Si hacemos igual que antes y sumamos cada $\xi_t(i, j)$ para cada instante de tiempo, obtenemos el número esperado de transiciones desde el estado i al estado j en las observaciones O : $\sum_{t=1}^{T-1} \xi_t(i, j)$.

Por lo tanto la relación entre $\gamma_t(i)$ y $\xi_t(i)$ viene dada por

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad \forall i \in \{1, \dots, N\}, \forall t \in \{1, \dots, M\}$$

Una vez calculadas estas variables, para resolver el problema de aprendizaje, se definen los nuevos parámetros $\lambda' = (A, B', \pi')$ del HMM utilizando las *smoothed* de acuerdo a las siguiente fórmulas

Probabilidades iniciales: probabilidades de estar en el estado i en $t = 1$:

$$\pi' = \gamma_1(i) \quad \forall i \in \{1, \dots, N\}$$

Matriz de transición: cada a_{ij} viene representado por el cociente entre el número esperado de transiciones entre i y j , entre el número total de transiciones desde i :

$$a'_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad \forall i \in \{1, \dots, N\}, \forall j \in \{1, \dots, N\}$$

Probabilidades de emisión: cociente entre el número de veces que se tiene el estado j y se observa o_k , y el número esperado de veces que se pasa por el estado j

$$b'_j(o_k) = \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad \forall i \in \{1, \dots, N\}, \forall j \in \{1, \dots, N\}$$

En resumen, el algoritmo Baum-Welch se divide en tres fases:

1. Se determina un modelo inicial, que aunque en teoría se puede seleccionar aleatoriamente, es conveniente hacerlo con una cierta aproximación realista para facilitar la convergencia del algoritmo.
2. Se realiza el cálculo de transiciones y símbolos de emisión más probables del modelo inicial, determinando las $\alpha's$, $\beta's$, $\gamma's$ y $\chi's$.
3. Se redefinen los parámetros a partir de lo calculado en el punto 2. para crear un nuevo modelo que mejore en verosimilitud al modelo inicial.

Modelo Gaussiano Oculto de Markov: GHMM

En este apartado describiremos un HMM cuyas observaciones tienen una distribución Gaussiana multivariante de dimensión L , diferente para cada estado. Esto quiere decir que, si tenemos K estados diferentes, cada uno de estos estará asociado a una distribución de parámetros $(\mu_i, \Sigma_i) \quad \forall i \in \{1, \dots, K\}$, con μ_i el vector de medias y Σ_i la matriz de covariancia del estado i .

Al ser un modelo continuo, existen infinitos símbolos (de hecho, $V = \mathbb{R}$). Usaremos la notación para los modelos continuos de $b_j(o_t) = P(o_t|q_t = i), \quad \forall j \in \{1, \dots, N\} \quad \forall t \in \{1, \dots, M\}$; que nos indicará la verosimilitud de la observación o_t respecto i .

Tenemos en cuenta que la función de densidad asociada al estado i será:

$$N(x, \mu_i, \Sigma_i) = \frac{1}{(2\pi\sigma_i^2)^L} e^{(x-\mu_i)\Sigma_i^{-1}(x-\mu_i)^t}, \text{ y por lo tanto } b_j(o(t)) = \sum_{k=1}^K c_{jk} N(x, \mu_j, \Sigma_j).$$

Aplicación de un GHMM

Vemos a continuación cómo podemos emplear un GHMM para modelizar un sistema estocástico como es el de la Bolsa.

Nuestro punto de partida será suponer que los retornos pueden pertenecer a dos estados $X = \{1, 2\}$, uno cuando el mercado tiende al alza (1), y otro cuando tiende a la baja (2), en este caso como las observaciones son escalares nuestros estados tendrán parámetros (μ_i, σ_i) , con $\mu_1 > \mu_2$.

Nota: nosotros suponemos que las observaciones que el HMM clasifique como estado 1 serán las correspondientes al estado alcista porque corresponden a una μ positiva, sin embargo, por lo general las σ 's de este estado son mayores que las del estado 2, y cabe la posibilidad de que en observaciones muy negativas, aunque sean las menos probables, encajen mejor con la distribución del estado 1 que con la del estado 2 en contra de lo intuitivo.

Nuestro caso de estudio será el correspondiente a los valores diarios de apertura del IBEX35 entre el 19/10/1990 y 12/08/1991.

```
## Warning in formatC(x = structure(c(-712662, -712297, -711201, -710836,
## -710470: class of 'x' was discarded
```

```
## \begin{table}[h]
## \centering
## \caption{IBEX 35: Valores de Apertura}
## \label{}
## \begin{tabular}{rrrl}
## \hline
## & Fecha & Valor & Retornos \\
## \hline
## 1 & -712662 & 2208.03 & 0 \\
```

```

##      2 & -712297 & 2268.76 & 2.8\% \\
##      3 & -711201 & 2309.00 & 1.8\% \\
##      4 & -710836 & 2331.88 & 1\% \\
##      5 & -710470 & 2398.47 & 2.9\% \\
##      \hline
## \end{tabular}
## \end{table}

```

Trabajaremos con los retornos (R), y estableceremos como condiciones iniciales $\mu_1 = \overline{R_+}$, donde $R_+ = \{R|R > 0\}$, y $\mu_2 = \overline{R_-}$, con $R_- = \{R|R < 0\}$ De la misma manera $\sigma_1 = \sigma(R_+)$ y $\sigma_2 = \sigma(R_-)$. En cuanto a la matriz de transición y a la probabilidad inicial, los suponemos de la siguiente manera:

$$A = \{a_{ij}\} = \{P(X_t = j | X_{t-1} = i)\} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix}$$

$$\begin{cases} \pi_1 = 0.5 \\ \pi_2 = 1 - \pi_1 = 0.5 \end{cases}$$

Tras aplicar el algoritmo de BaumWelch, obtenemos la estimación de los parámetros:

$$\mu = \begin{pmatrix} 5.8 \cdot 10^{-3} \\ -3.6 \cdot 10^{-3} \end{pmatrix}$$

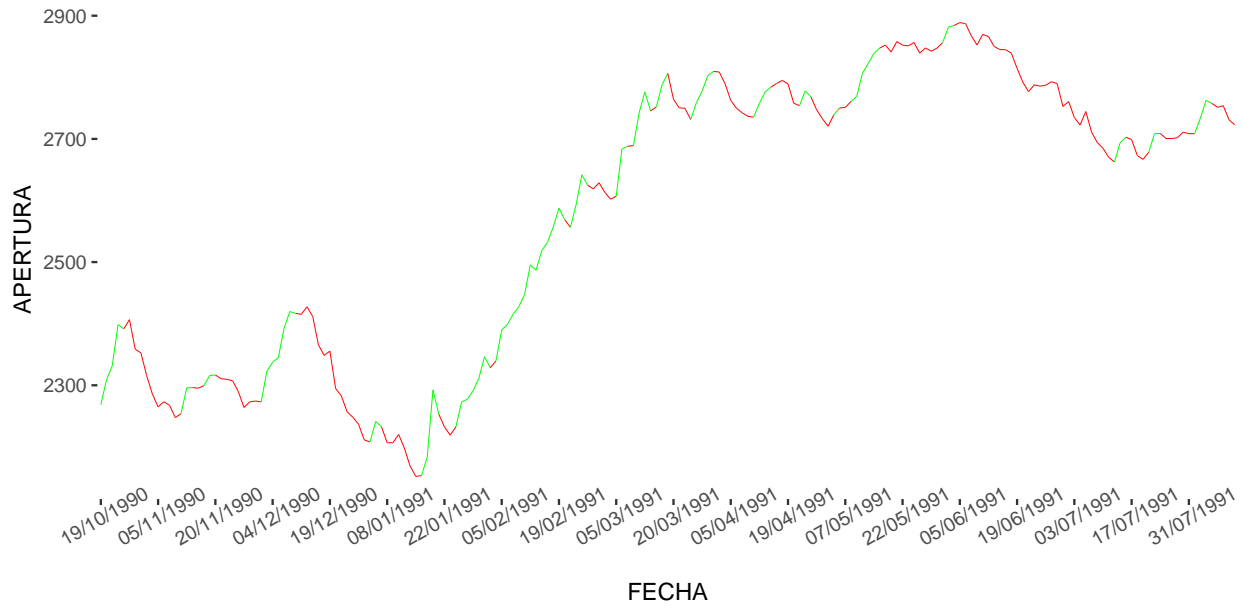
$$\sigma = \begin{pmatrix} 9.43 \cdot 10^{-5} \\ 4.17 \cdot 10^{-5} \end{pmatrix}$$

La nueva matriz de transición y probabilidad inicial son:

$$A = \begin{pmatrix} 0.802 & 0.198 \\ 0.196 & 0.804 \end{pmatrix}$$

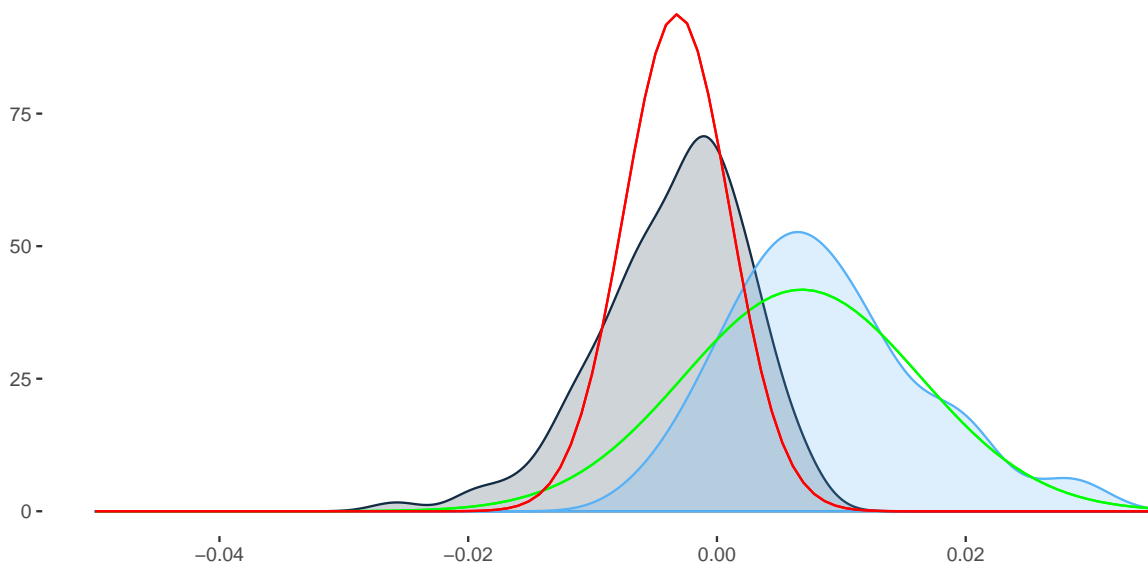
$$\begin{cases} \pi_1 = 1 \\ \pi_2 = 0 \end{cases}$$

Para ilustrar la capacidad clasificadora del HMM, mostramos un gráfico que de la evolución de los valores del IBEX35 en el periodo estudiado, junto a la clasificación que ha hecho el algoritmo de Viterbi de la secuencia de estados más probables.



Gráfica 1. *Gráfico del IBEX 35 entre el 19/10/1990 y 12/08/1991, en verde se han indicado los retornos clasificados como alcistas, y en rojo los bajistas.*

Vemos cómo nuestra clasificación de los retornos se aproxima de manera muy efectiva al comportamiento que muestra la gráfica. Aunque hay retornos positivos y negativos en ambos estados, la distribución a la que se atribuye cada retorno se hace en función del régimen alcista/bajista en que se encuentre la serie.



Gráfica 2. *Distribución de los retornos pertenecientes a cada estado, junto con el gráfico de la Gaussiana asociada.*

3. Redes Neuronales Artificiales

Introducción

En este capítulo estudiaremos otro de los métodos de Machine Learning más utilizados: redes neuronales artificiales (ANN por sus siglas en inglés).

Comentaremos brevemente cuál es la base biológica que incorporan, y a continuación explicaremos los elementos básicos que emplean estas implementaciones.

Luego nos meteremos de lleno a analizar el principal algoritmo que se utiliza a la hora de entrenar redes neuronales artificiales para que sean capaz de realizar tareas concretas de manera efectiva: el algoritmo de *Back-propagation*. Destacaremos el principal problema que implica su implementación: la dificultad de su convergencia; y daremos dos variantes de este algoritmo que tienen como objetivo minimizar este efecto: el Método de α Adaptable y el Algoritmo de *Resilient Back-propagation*.

Base Biológica

El cerebro humano es un complicado mecanismo -o casi siempre lo es- capaz de resolver problemas de gran complejidad. Estamos aún lejos de entender completamente cómo funciona, mucho menos de poder replicarlos artificialmente; sin embargo, sí que tenemos un buen conocimiento para explicar las operaciones básicas que realiza un cerebro.

Para entender cuál es el fundamento de una Red Neuronal Artificial (ANN), primero tenemos que tener un conocimiento básico de los mecanismos internos del cerebro. El cerebro es la parte central del sistema nervioso y está formado por una extensa red neuronal, consistente en aproximadamente 10^{11} neuronas interconectadas.

El centro de cada neurona se llama núcleo. El núcleo está conectado a otros núcleos mediante las dendritas y los axones. Esta conexión es la denominada conexión sináptica.

Las neuronas pueden disparar pulsos eléctricos a través de las conexiones sinápticas, los cuales son recibidos a través de las dendritas de otras neuronas. Cuando una neurona recibe suficiente impulso eléctrico, envía asimismo un pulso eléctrico a través de sus axones, y así sucesivamente. Este proceso permite propagar información a través de la red neuronal.

La cuestión es que esta conexión sináptica cambia durante la vida de las neuronas, fortaleciéndose o debilitándose dependiendo del uso. Esta plasticidad de la sinapsis es lo que permite al cerebro aprender, y es el principio que se aplica en las ANN's para poder entrenar una red neuronal artificial para tareas específicas a través de grandes muestras de ejemplares.

Definición de Red Neuronal Artificial (ANN)

Una neurona artificial puede ser implementada de muchas maneras. La definición matemática general es la siguiente:

$$y(x) = g\left(\sum_{i=0}^n w_i x_i - \theta_i\right)$$

donde x es una neurona con n inputs $x = (x_1, \dots, x_n)$ y un axón de output $y(x)$. Los términos $w = w_0, \dots, w_n$ son los pesos, y determinan cómo deben ser de intensificados los inputs. θ_i es el umbral o sesgo del input i .

g es una función de activación que balancea cómo de potentes debe ser el output (si existe) de la neurona, basado en la suma de los input. En el caso de neuronas reales, g es una función δ de Kronecker que determina una relación binaria de todo o nada para el output. Sin embargo, no se suele implementar de esa manera una neurona artificiales, ya que como veremos luego, los algoritmos para entrenarlas suelen venir determinados por funciones diferenciables.

Entre los tipos de funciones de activación más habituales están:

$$f_1(x) = \alpha x, \alpha \in \mathbb{R} \text{ la función lineal}$$

$$f_2(x) = \begin{cases} S, & \text{si } x > S; \\ x, & \text{si } |x| \leq S; \\ -S, & \text{en cualquier otro caso.} \end{cases}$$

Donde S es el output saturado. Por último están las funciones sigmoidales, que son de las funciones más utilizadas, particularmente la logística (y que será la que nosotros usaremos por defecto).

$$f_3(x) = \frac{m}{1 + e^{-x}}$$

con m el coeficiente de magnitud. Esta función es diferenciable y con un dominio $D(f_3(x)) = (0, m)$.

La función sigmoideal es una función real de variable real diferenciable. Si $y(x) = f_3(x)$,

$$\frac{dy(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = y(x)(1 - y(x))$$

por lo que $y(x)$ tiene una derivada simple

Backpropagation

Entre las ANN's existentes, una de las más empleadas por su alta eficiencia es el Perceptrón Multicapa (MLP) mediante el uso del algoritmo de *Back-Propagation* descrito por Jian-Rong Chen . A diferencia del Perceptrón de una única capa, el MLP puede implementar gran variedad de funciones complejas, incluida la función XOR, que no puede ser realizada por el de una única capa como demostraron Minsky y Papert {ojo} (Perceptrons: an introduction

to computational geometry is a book written by Marvin Minsky and Seymour Papert and published in 1969).

En un MLP una unidad solo puede conectarse a la capa adyacente siguiente, no permitiéndose conexiones recurrentes ni en la misma capa. Sea K el número de capas, M el número de inputs y N el número de outputs. El input en una unidad (siempre que no sea en la capa de entrada) es la suma de los outputs de unidades conectadas en la capa anterior. Sea x_i^j el input de la unidad i en la capa j , w_{ij}^k el peso de la conexión entre la unidad i en la capa k y la unidad j de la capa $k + 1$, y_i^j el output de la unidad i en la capa j y, por último, θ_i^j el umbral (o sesgo) de la unidad i en la capa j . Entonces tenemos que los input de la capa $k + 1$ son:

$$x_j^{k+1} = \sum_i w_{ij}^k y_i^k$$

donde

$$y_i^j = f(x_i^j) \quad (1)$$

siendo f la función de activación. Nosotros utilizaremos como f la función sigmoide

$$f(x_i^j) = \frac{1}{1 + e^{-\frac{x_i^j - \theta_i^j}{T}}}$$

Dependiendo de las aplicaciones un output que tome valores negativos es necesario, y podemos utilizar entonces la función

$$f(x_i^j) = \frac{1 - e^{-\frac{x_i^j - \theta_i^j}{T}}}{1 + e^{-\frac{x_i^j - \theta_i^j}{T}}}$$

Hay dos fases en el algoritmo de *Back-Propagation*, la primera es computar el cálculo a través de las capas utilizando (1). La segunda fase es actualizar los pesos, operación que se realiza computando el error entre el valor esperado y el valor real calculado en la primera fase. Este proceso clasifica el algoritmo de *Back-Propagation* dentro de la categoría de los algoritmos de aprendizaje supervisado. Básicamente el algoritmo de *Back-Propagation* es un algoritmo de gradiente descendente.

A continuación explicaremos cómo se realiza la actualización de los pesos, una vez obtenido el output a través de (1). Primero definiremos una función de error \mathcal{E} que nos dará cuenta de la discrepancia entre el valor calculado y el real. Sea N el número de outputs, d_i el output deseado y y_i el obtenido, con $i \in \{1, \dots, N\}$. Entonces

$$\mathcal{E} = \frac{1}{2} \sum_j (d_j - y_j)^2, \quad (2)$$

El error total será simplemente $\mathcal{E}_{total} = \sum_{k=1}^N \mathcal{E}_k$. Nuestro objetivo será minimizar este error total, para ello utilizamos un algoritmo de gradiente descendente. Este tipo de algoritmo

busca un mínimo en la función (la función error en nuestro caso) dando pasos proporcionales al negativo del gradiente. Es por ello que es tan importante que la función que usamos como función de activación sea derivable.

Así, el ajuste de los pesos es proporcional a la derivada

$$\Delta w_{ij}^k = -\eta \frac{\partial \mathcal{E}}{\partial w_{ij}^k} \quad (3)$$

donde η es el tamaño del paso, importante para asegurar la convergencia. También se puede añadir un término de inercia, mediante una dependencia de recursividad, que nos ayudará a mejorar la convergencia evitando rápidos cambios en Δw_{ij}^k .

$$\Delta w_{ij}^k(t) = -\eta \frac{\partial \mathcal{E}}{\partial w_{ij}^k(t)} + \alpha \Delta w_{ij}^k(t-1) \quad (4)$$

donde α es un positivo real pequeño.

Estrictamente hablando, la \mathcal{E} que se utiliza en (3) debería ser \mathcal{E}_{total} . Sin embargo es mucho más práctico actualizar pesos con cada input de una muestra de entrenamiento, en vez de usar toda la muestra al completo. Estos dos casos se denominan batch y online, respectivamente. En este caso puede utilizarse el \mathcal{E} definido en (2).

Así, para cada capa de un perceptrón simple, tendríamos el output dado por

$$y_j = f\left(\sum_{i=1}^{n-1} w_i I_i + \theta_j\right) = f(s)$$

El sesgo θ_j puede añadirse como input siempre activo ($I_n = 1$) con peso $w_n = \theta_j$

$$y_j = f\left(\sum_{i=1}^n w_i I_i\right) = f(s)$$

Si el output deseado es d_j , entonces de (2) tenemos

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial w_i} &= \frac{\partial \mathcal{E}}{\partial y_j} \frac{\partial y_j}{\partial w_i} \\ &= -(d_j - y_j) f'(s) I_i \end{aligned}$$

Sea

$$\delta_j = (d_j - y_j)$$

Entonces el ajuste de pesos viene dado por

$$\Delta w_i = -\eta \frac{\partial \mathcal{E}}{\partial w_i} = \eta \delta_j f'(s) I_i$$

Este esquema de ajuste de pesos es denominado la regla Delta. El algoritmo de Back-Propagation en realidad es una generalización de la regla Delta.

Sea una red MLP, con N capas y suponemos que nuestra función de activación f es la sigmoide logística, entonces

$$\frac{\partial \mathcal{E}}{\partial y_j^N} = -(d_j - y_j^N)$$

Donde el superíndice N indica que el output es de la capa N y de (1) obtenemos

$$\frac{\partial y_j^N}{\partial x_j^N} = f'(x_j^N) = y_j^N (1 - y_j^N).$$

Como $x_j^{k+1} = \sum_i w_{ij}^k y_i^k$, entonces

$$\frac{\partial x_j^N}{\partial w_{ij}^{N-1}} = y_i^{N-1}$$

Y por lo tanto, podemos escribir la parcial del error respecto de los pesos como

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial w_{ij}^{N-1}} &= \frac{\partial \mathcal{E}}{\partial y_j^N} \frac{\partial y_j^N}{\partial x_j^N} \frac{\partial x_j^N}{\partial w_{ij}^{N-1}} \\ &= -(d_j - y_j^N) y_j^N (1 - y_j^N) y_i^{N-1} \end{aligned}$$

Sea

$$\delta_j^{N-1} = (d_j - y_j^N) y_j^N (1 - y_j^N) \quad (5)$$

Luego

$$\frac{\partial \mathcal{E}}{\partial w_{ij}^{N-1}} = -\delta_j^{N-1} y_i^{N-1} \quad (6)$$

Entonces, teniendo en cuenta (4) y (5) tenemos que en la iteración $t + 1$ los pesos serán

$$w_{ij}^{N-1}(t+1) = w_{ij}^{N-1}(t) + \eta \delta_j^{N-1} y_i^{N-1} + \alpha [w_{ij}^{N-1}(t) - w_{ij}^{N-1}(t-1)] \quad (7)$$

Obviamente, esta es la actualización de los pesos para la capa del output. Pero como la delta la podemos definir así

$$\delta_j^{N-1} = -\frac{\partial \mathcal{E}}{\partial y_j^N} \frac{y_j^N}{\partial x_j^N} = -\frac{\partial \mathcal{E}}{\partial x_j^N}$$

Podemos tener una fórmula general para $w_{ij}^k, \forall k \in \{1, \dots, N-1\}$ de la siguiente manera

$$\begin{aligned}
\frac{\partial \mathcal{E}}{\partial w_{ij}^k} &= \frac{\partial \mathcal{E}}{\partial y_j^{k+1}} \frac{\partial y_j^{k+1}}{\partial x_j^{k+1}} \frac{\partial x_j^{k+1}}{\partial w_{ij}^k} \\
&= \frac{\partial \mathcal{E}}{\partial y_j^{k+1}} y_j^{k+1} (1 - y_j^{k+1}) y_i^k \\
&= y_j^{k+1} (1 - y_j^{k+1}) y_i^k \sum_l \frac{\partial \mathcal{E}}{\partial x_l^{k+2}} \frac{\partial x_l^{k+2}}{\partial y_j^{k+1}} \\
y_i^k y_j^{k+1} (1 - y_j^{k+1}) \frac{\partial \mathcal{E}}{\partial x_l^{k+2}} w_{jl}^{k+1} &= -y_i^k y_j^{k+1} (1 - y_j^{k+1}) \sum_l w_{jl}^{k+1} \delta_l^{k+1}
\end{aligned}$$

Sea

$$\delta_j^k = y_j^{k+1} (1 - y_j^{k+1}) \sum_l w_{jl}^{k+1} \delta_l^{k+1} \quad (8)$$

Entonces tenemos que

$$\frac{\partial \mathcal{E}}{\partial w_{ij}^k} = -\delta_i^k y_i^k,$$

esta ecuación es similar a (6). Finalmente obtenemos el algoritmo de Back-Propagation combinando (7) y (8), llamado así porque las δ 's se propagan hacia atrás, comenzando por (5).

En general, la manera de ejecutar este algoritmo es:

1. Asignar valores pequeños y aleatorios a los pesos y sesgos del MLP.
2. Dependiendo de si el algoritmo se ha programado según el método batch u online, se le pasa un input de la muestra de entrenamiento, o la muestra entera y determinar las δ 's.
3. Actualizar pesos según el algoritmo de Back-propagation.
4. Iterar hasta alcanzar la tolerancia de error deseado.

El problema del algoritmo de Back-propagation, como ocurre en todos los que poseen ciertas características geométricas, es que resulta lento en su convergencia, ya que η depende mucho de la topología local y curvatura de las superficies del error \mathcal{E} . Así, si introducimos una η pequeña el algoritmo convergerá de forma extremadamente lenta, mientras que si η es muy grande, provocaremos que el algoritmo vaya saltando de un lado a otro de los valles de \mathcal{E} , en vez de seguir la curva que marca el camino a su mínimo.

Por esta razón, daremos dos variantes del algoritmo de Back-propagation que se usan frecuentemente para acelerar su convergencia, la *Técnica de α Adaptable* y el *Algoritmo de Resilient Back-propagation* de Martin Riedmiller (1994).

Técnica de α Adaptable

En esta variante del algoritmo de Back-propagation se añade un término de inercia, siguiendo el modelo de (4). En esta técnica introducimos una dependencia de la α (tamaño de salto en el gradiente) que aparece en (4) respecto al número de iteración siguiendo este modelo:

$$\begin{aligned}\alpha(t) &= \alpha(t-1)(1 - h(t)\sqrt{\mathcal{E}(t)}), \quad t \geq 2 \\ h(t) &= a_1 h(t-1) + a_2 \Delta \mathcal{E}(t) \quad t \geq 2 \\ \Delta \mathcal{E}(t) &= \mathcal{E}(t) - \mathcal{E}(t-1) \quad t \geq 2\end{aligned}$$

$\alpha(t)$ es el tamaño de salto en el momento t . $\mathcal{E}(t)$ es la suma cuadrática de errores entre el output deseado y el obtenido en la iteración t

$$\mathcal{E} = \frac{1}{2} \sum_{k=1}^N \sum_{i=1}^P (d_i^k - o_i^k)^2,$$

con N el número de capas del MLP y P el número de outputs. Así, $\Delta \mathcal{E}(t)$, $t \geq 2$ sería el decremento de $\mathcal{E}(t)$, y $h(t)$ será un filtro recursivo de $\Delta \mathcal{E}(t)$: un filtro recursivo reutiliza uno o más de sus outputs como inputs, y permite controlar de esta manera que haya cambios bruscos en $\Delta \mathcal{E}(t)$, creando un versión del algoritmo más estable.

Los parámetros a_1 y a_2 del filtro recursivo $h(t)$ son los encargados de controlar esta adaptación. Para una gran muestra de entrenamiento, un valor pequeño de a_1 y grande en a_2 , combinado con una η de gran tamaño también es lo que más favorece la convergencia **Tesauro, G. and Janssens, B. Scaling Relationships in Back Propagation Learning** *Complex System* vol-2 No.1 pp38-44 1988.

Un análisis del algoritmo nos permite ver que si $h(t)$ es positivo, entonces la tendencia de $\mathcal{E}(t)$ en el pasado inmediato es incrementar, por lo que $1 - h(t)\sqrt{\mathcal{E}(t)} < 1$, y por lo tanto α decrecerá en este paso. De forma similar concluimos que si la tendencia de $\mathcal{E}(t)$ es decrecer, entonces α crecerá. Y finalmente, si $\mathcal{E}(t)$ es muy pequeño, significará que el MLP ya casi ha terminado de aprender, y la adaptación de muy baja magnitud, estabilizando el algoritmo.

Resilient Back-propagation

Este algoritmo, desarrollado por Martin Riedmiller (1994) para MLP con método de aprendizaje batch, tiene como objetivo principal eliminar la influencia del tamaño de la derivada parcial en la actualización de los pesos. Por ello, se considera únicamente el signo de la derivada para indicar la dirección en que tiene que ir esta actualización. Consideramos la conexión de la unidad i , $i \in \{1, \dots, I\}$ de una capa con la unidad j , $j \in \{1, \dots, J\}$ de la siguiente capa en la iteración t :

$$\Delta w_{ij}(t) := \begin{cases} -\Delta_{ij}(t), & \text{si } \frac{\partial \mathcal{E}(t)}{\partial w_{ij}(t)} > 0 \\ +\Delta_{ij}(t), & \text{si } \frac{\partial \mathcal{E}(t)}{\partial w_{ij}(t)} < 0 \\ 0, & \text{si } \frac{\partial \mathcal{E}(t)}{\partial w_{ij}(t)} = 0 \end{cases}$$

Donde $\frac{\partial \mathcal{E}(t)}{\partial w_{ij}(t)}$ es la suma para cada componente de la muestra (batch learning).

Si reemplazamos $\Delta_{ij}(t)$ por una constante se obtiene la regla de actualización de pesos de **Manhattan**.

Para determinar los valores de $\Delta_{ij}(t)$, utilizamos también un filtro de recursivo de la forma

$$\Delta_{ij}(t) := \begin{cases} \eta^+ \Delta_{ij}(t-1), & \text{si } \frac{\partial \mathcal{E}(t-1)}{\partial w_{ij}} \frac{\partial \mathcal{E}(t)}{\partial w_{ij}} > 0 \\ \eta^- \Delta_{ij}(t-1), & \text{si } \frac{\partial \mathcal{E}(t-1)}{\partial w_{ij}} \frac{\partial \mathcal{E}(t)}{\partial w_{ij}} < 0 \\ \Delta_{ij}(t-1), & \text{si } \frac{\partial \mathcal{E}(t-1)}{\partial w_{ij}} \frac{\partial \mathcal{E}(t)}{\partial w_{ij}} = 0 \end{cases}$$

donde $0 < \eta^- < 1 < \eta^+$

Como ya comentamos antes, puede ocurrir que la actualización de pesos haya hecho que nos saltemos un mínimo, lo cual implica que el signo de la derivada parcial del correspondiente w_{ij} cambie, así que disminuimos el valor de actualización $\Delta_{ij}(t)$ un factor η^- . En el caso contrario, para acelerar la convergencia, este valor de actualización es incrementado en $\eta^+ (> 1)$. En caso de que hayamos alcanzado un mínimo, no hay ninguna actualización.

En general, para evitar tener un gran número de parámetros libres, se escogen valores constantes de η^+ y η^- . Dos valores que, de acuerdo a las simulaciones realizadas por Riedmiller, parecen funcionar bien independientemente de la muestra, son $\eta^+ = 1.2$ y $\eta^- = 0.5$.

Hay tres parámetros a considerar: Δ_0 , valor inicial de actualización de pesos, Δ_{max} , que es el límite del tamaño del salto y Δ_{min} , mínimo peso de actualización.

Cuando inicializamos el algoritmo, todos los valores de actualización de pesos son Δ_0 , que debe ser escogido teniendo en cuenta también los valores iniciales de los pesos, si por ejemplo estos pesos están en el rango $[-1,1]$ un buen valor será $\Delta_0 = 0.1$.

El valor de Δ_{max} no es de gran importancia en la mayoría de casos, pero conviene fijar uno para evitar posibles problemas de rebote alrededor de un mínimo. Por último, Δ_{min} sí que es un valor de importancia, ya que fijar un valor adecuado del mismo evitará que el algoritmo se quede atascado si se encuentra con algún mínimo local subóptimo. En su estudio Riedmiller fija un valor de $\Delta_{min} = 1e^{-6}$.

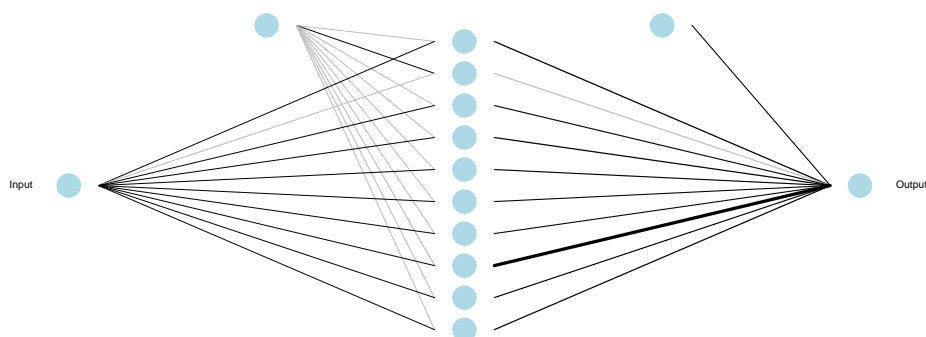
Así, el pseudo-código de implementación del algoritmo de Resilient Back-propagation sería este:

Ejemplo Red Neuronal Artificial

Veremos a continuación un par de aplicaciones de una ANN, usando el algoritmo de *Resilient Backpropagation*.

El primer ejemplo entrenaremos una red neuronal para que sea capaz de calcular el cuadrado de un número.

Usaremos una muestra aleatoria de 50 números entre el 0 y el 10 como inputs, y el cuadrado de estos números como output de una ANN con 10 neuronas ocultas y un error de $\delta = 0.01$.



En el gráfico se muestra un esquema de esta ANN, cuanto mayor es el peso de la conexión más oscuro es el color de la representación de esta conexión.

Veamos la precisión de esta ANN calculando el cuadrado de los naturales del 1 al 10.

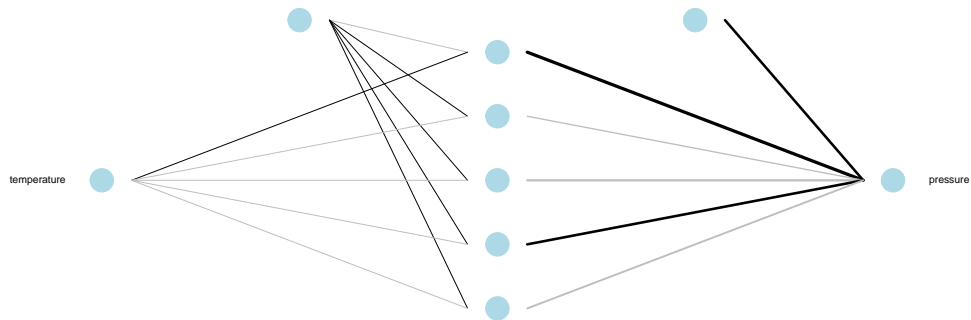
Table 1: Resultados ANN

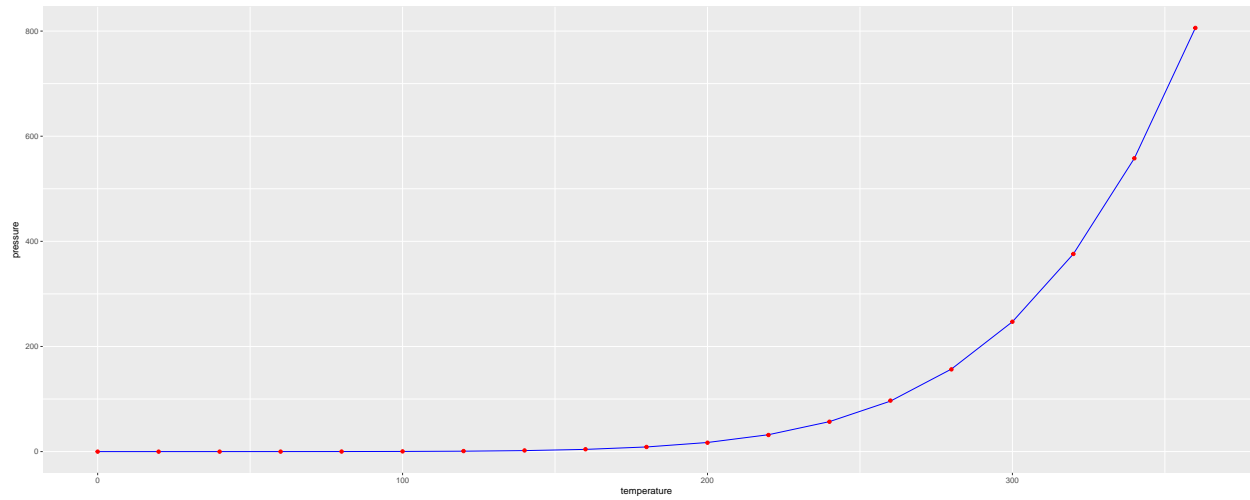
Input	Output	Error
1.00	0.90	0.10
2.00	4.00	0.00
3.00	8.99	0.01
4.00	16.00	0.00
5.00	25.00	0.00
6.00	36.00	0.00
7.00	49.00	0.00
8.00	64.00	0.00
9.00	81.00	0.00
10.00	99.99	0.01

El otro ejemplo que veremos será una ANN entrenada para deducir la relación entre la temperatura ($^{\circ}\text{C}$) y la presión de vapor de mercurio en milímetros (de mercurio), a partir de 19 observaciones.

Table 2: Observaciones

temperature	pressure
0	0.0002
20	0.0012
40	0.0060
60	0.0300
80	0.0900
100	0.2700
120	0.7500
140	1.8500
160	4.2000
180	8.8000
200	17.3000
220	32.1000
240	57.0000
260	96.0000
280	157.0000
300	247.0000
320	376.0000
340	558.0000
360	806.0000





En azul la curva original, y en rojo los puntos determinados por la ANN