



UNIVERSIDAD DE GRANADA

GRADO INGENIERÍA INFORMÁTICA (2017 – 2018)

---

# APRENDIZAJE AUTOMÁTICO

Trabajo 1: Cuestiones de Teoría

Trabajo realizado por Antonio Miguel Morillo Chica.

---

1. Identificar, para cada una de las siguientes tareas, que tipo de aprendizaje es el adecuado (supervisado, no supervisado, por refuerzo) así como los datos de aprendizaje que deberíamos usar en su caso. Si una tarea se ajusta a más de un tipo, explicar como y describir los datos para cada tipo.

- a) **Dada una colección de fotos de caras de personas de distintas razas establecer cuantas razas distintas hay representadas en la colección.**

La detección de caras en fotografías se hace apartir de aprendizaje supervisado para saber que parte de la foto es o no una cara, problema visto en las transparencias de teoría. Tomando esto, supongo que la distinción de las mismas por etnias se ajustaría a un aprendizaje tambien supervisado ya que existen características bien distintivas, tipo de pelo (nivel de rizado), color de piel, anchura de cara, etc.

Las características de los datos podrían ser las ya citadas aunque como son fotos distintas una posible entrada también podría ser colores de los pixeles, distancias entre los elementos de la cara (nariz, ojos etc.) en mi caso:

$$X = \{(x_1, x_2, \dots, x_{10}) \in \mathbb{Z}_2^{10}\}$$

$$Y = \{Raza_1, Raza_2, Raza_3, Raza_4\}$$

Debido al tema controversial de la división por razas o etnias de los seres humanos por parte de biólogos, antropólogos y otros expertos no defino la salida como caucasico, asiatico etc.

- b) **Clasificación automática de cartas por distrito postal.**

Tomando en cuenta que la lectura del digito postal consiste en la digitalización del número ya que este estaría escrito a mano, el problema se transforma en el reconocimiento de números que, como vimos en teoría, era un ejemplo de aprendizaje supervisado.

Los datos para este problema serían unas imágenes que representan a cada dígito individualmente de un tamaño de 16x16 pixeles lo que supone

un vector con 1024 características que definen intensidad o contraste del pixel así pues tendríamos:

$$X = \{(x_1, x_2, \dots, x_{256}) \in \mathbb{Z}_{255}^{256}\}$$

$$Y = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

- c) **Decidir si un determinado índice del mercado de valores subirá o bajará dentro de un periodo de tiempo fijado.**

El aprendizaje supervisado se ajusta perfectamente a este tipo de problemas ya que lo que buscamos es una pronosticación futura, además según un artículo de la Universidad de Medellín, Colombia los modelos no lineales se ajustan mejor a este tipo de problemas, [aquí](#) puede encontrar el artículo como dato interesante. Para mi sería abordado por un aprendizaje por refuerzo, el aprendizaje supervisado necesitaría una cantidad increíble de características ya que muchos factores influyen en la economía, y del mismo modo, para el no supervidado necesitaríamos una ingente cantidad de datos, por ello creo que la mejor linea de aprendizaje sería basarse en el dato numerico que representa el valor del mercado durante un periodo de tiempo de tal forma que:

$$X = \{(x_1, x_2, \dots, x_{90}) \in \mathbb{R}^{90} \times \mathbb{R}\}$$

$$Y = [0, 1]$$

Donde los datos vienen dados por los valores en los 90 días anteriores y una predicción dónde Y indica el valor de éxito de la predicción.

- d) **Aprender un algoritmo que permita a un robot rodear un obstaculo.**

Creo que el aprendizaje por refuerzo puede ser muy util para crear trayectorias “globales”, es decir, para planificadores globales vistos en Tecnologías y Sistemas Inteligentes. Imaginemos que lo que queremos es generar una trayectoria aproximada lo más optima posible, los datos de entrada puede ser la propia trayectoria un vector de posiciones, posteriormente haremos simulaciones para retroalimentar al algoritmo. De forma que si trabajamos en un mapa de  $n \times n$  :

$$X = \{(x_1, x_2, \dots, x_n) \in \mathbb{N}^n \times \mathbb{N}\}$$

$$Y = [0, 1]$$

2. ¿Cuales de los siguientes problemas son más adecuados para una aproximación por aprendizaje y cuales más adecuados para una aproximación por diseño? Justificar la decisión.

a) **Agrupar los animales vertebrados en mamíferos, reptiles, aves, anfibios y peces.**

Aproximación por aprendizaje. Tenemos una serie de características distinguibles que pueden ser usados para diferenciar a los vertebrados, como por ejemplo, tener pico, plumas, pelaje, número de patas, etc. Es un típico problema de aprendizaje supervisado donde:

$$X = \{(x_1, x_2, \dots, x_n) \in \mathbb{N}_2^n\}$$

$$Y = \{Mamifero, Reptiles, Aves, Anfibios, Aves, Peces\}$$

b) **Determinar si se debe aplicar una campaña de vacunación contra una enfermedad.**

Aproximación por aprendizaje por diseño ya que el posible desarrollo de una enfermedad en base a criterios como el contagio puede ser pronosticada, o simulada por procedimientos de algoritmos clásicos.

c) **Determinar si un correo electrónico es de propaganda o no.**

Aproximación por aprendizaje. Es además un tipo clásico de clasificación de spam que se usa mucho como ejemplo en clase o en otras clases de la rama. Poseemos un conjunto finito de respuestas, es o no spam, las características pueden ser la presencia o no de ciertas palabras además de poder dar ejemplos claros de lo que es o no un correo de spam para realizar un conjunto de entrenamiento. De forma que:

$$X = \{(x_1, x_2, \dots, x_n) \in \mathbb{N}_2^n\}$$

$$Y = \{Spam, No Spam\}$$

d) **Determinar el estado de ánimo de una persona a partir de una**

**foto de su cara.**

Aproximación por aprendizaje. Es un problema muy típico al de reconocimiento de caras, una vez reconocidas el proceso sería parecido un aprendizaje supersado de clasificación para distinguir otro tipos de características a las del problema de detectar que es y no una cara.

e) **Determinar el ciclo óptimo para las luces de los semáforos en un cruce con mucho tráfico.**

Aproximación por aprendizaje ya que es inasumible desde el punto de vista del diseño, existen demasiadas variables y demasiados posibles estados finales para la solución.

Además según la Pontífica Universidad Católica de Chile según la tesis sobre “Aplicación de algoritmos de aprendizaje estadístico para predecir velocidades de buses con información en tiempo real” este problema es un problema por aprendizaje.

**3. Construir un problema de aprendizaje desde datos para un problema de clasificación de fruta en una explotación agraria que produce mangos, papayas y guayabas. Identificar los siguientes elementos formales  $X$ ,  $Y$ ,  $D$ ,  $f$  del problema. Dar una descripción de los mismos que pueda ser usada por un computador. ¿Considera que en este problema estamos ante un caso de etiquetas con ruido o sin ruido? Justificar las respuestas.**

La entrada sería una matriz de dos componentes de color y tamaño, donde el color viene definido por la composición de tres colores RGB por ello tiene tres dimensiones, además el tamaño sería siempre de tipo real el cual realmente representaría el radio para determinar el tamaño.

$$\blacksquare \quad X = \{(Colorr, Radio)\} / colorr_n \in \mathbb{R}_{255}^3, radio_n \in \mathbb{R}^{+i}$$

La salida sería la clasificación en las etiquetas de las distintas frutas del problema.

$$\blacksquare \quad Y = \{Mangos, Papayas, Guayabas\}$$

El conjunto de entrenamiento viene dado por los primeros 50 ejemplos de las

distintas 50 primeras características.

$$\blacksquare \quad D = \{(x_i, y_i) \mid i = 1, 2, \dots, 50, x_i \in X, y_i \in Y\}$$

La función  $f$  será la que tendremos que encontrar para poder categorizar los ejemplos y predecir futuros.

$$\blacksquare \quad f: X \rightarrow Y$$

Evidentemente existirá ruido ya que el tamaño por el radio no siempre será perfecto por muy preciso que sea además que de por si ninguna fruta será totalmente esferica.

4. Sea  $X$  una matriz de números reales de dimensiones  $N \times d$ ,  $N > d$ . Sea  $X = UDV^T$  su descomposición en valores singulares (SVD). Calcular la SVD de  $X^T X$  y  $XX^T$  en función de la SVD de  $X$ . Identifique dos propiedades de estas nuevas matrices que no tiene  $X$ . ¿Qué valor representa la suma de la diagonal principal de cada una de las matrices producto?

#### 1. Para $X^T X$

$X^T X = (UDV^T)^T UDV^T$  por propiedad de matrices sabemos que:

$$A = UEV^T = (UEV^T)^T = VE^T U^T, \text{ entonces:}$$

$$X^T X = (UDV^T)^T UDV^T = (VD^T U^T) UDV^T, \text{ también sabemos que } I = U^T U$$

$$X^T X = (UDV^T)^T UDV^T = (VD^T U^T) UDV^T = VD^T D V^T, \text{ como } D \text{ es una matriz con solo diagonal por propiedad de matrices } D^T D = D^2, \text{ entonces:}$$

$$X^T X = (UDV^T)^T UDV^T = (VD^T U^T) UDV^T = VD^T D V^T = VD^2 V^T$$

Propiedades: Siendo  $U$  y  $D$  matrices cuadradas donde  $U$  es una matriz triangularizable y  $D$  la submatriz con los valores propios, entonces  $UD^T$  y  $DU^T$  serían dos matrices simétricas.

#### 2. Para $XX^T$

$$XX^T = (UDV^T)(UDV^T)^T = (UDV^T)(VD^T V^T) = UDD^T U^T$$

Propiedades: Siendo  $U$  y  $D$  matrices cuadradas donde  $U$  es una matriz

triangularizable y D la submatriz con los valores propios, entonces UD es la traspuesta de DU es  $D^T U^T$ .

5. Sean  $x$  e  $y$  dos vectores de características de dimensión  $M \times 1$ . La

expresión  $cov(x, y) = \frac{1}{M} \sum_{i=1}^M (x_i - \bar{x})(y_i - \bar{y})$  define la covarianza entre dichos vectores, donde  $\bar{z}$  representa el valor medio de los elementos de  $z$ .

Considere ahora una matriz  $X$  cuyas columnas representan vectores de características. La matriz de covarianzas asociada a la matriz

$X = (x_1, x_2, \dots, x_N)$  es el conjunto de covarianzas definidas por cada dos de sus vectores columnas. Es decir,

$$cov(X) = \begin{pmatrix} cov(x_1, x_1) & cov(x_1, x_2) & \dots & cov(x_1, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ cov(x_N, x_1) & cov(x_N, x_2) & \dots & cov(x_N, x_N) \end{pmatrix}$$

Sea  $1_M^T = (1, 1, 1, \dots, 1)$  un vector  $M \times 1$  de unos. Mostrar que representan las siguientes expresiones

a)

$$\begin{aligned} E1 = 11^{(T)} X &= \begin{pmatrix} 1_1 \\ 1_2 \\ \vdots \\ 1_M \end{pmatrix} \times (1_1, 1_2, \dots, 1_{M_0}) \times (w_1, w_2, \dots, w_N) = \begin{pmatrix} 1_1 1_1 & 1_1 1_2 & \dots & 1_1 1_M \\ 1_2 1_1 & 1_2 1_2 & \dots & 1_2 1_M \\ \vdots & \vdots & \ddots & \vdots \\ 1_M 1_1 & 1_M 1_2 & \dots & 1_M 1_M \end{pmatrix} \times (w_1, w_2, \dots, w_N) \\ &= \begin{pmatrix} \sum_{i=0}^C w 1_i \sum_{i=0}^C w 2_i & \sum_{i=0}^C w_{ni} \\ \sum_{i=0}^C w 1_i \sum_{i=0}^C w 2_i & \ddots \sum_{i=0}^C w_{ni} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=0}^C w 1_i \sum_{i=0}^C w 2_i & \sum_{i=0}^C w_{ni} \end{pmatrix} = E1 \end{aligned}$$

Es decir, E1 es una matriz donde en cada columna está en las sumatorias de sus características.

$$b) \quad E2 = \left( X - \frac{1}{M} E 1 \right)^T \left( X - \frac{1}{M} E 1 \right) =$$

$$\left( w_1, w_2, \dots, w_N \right) - \begin{pmatrix} \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \\ \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \end{pmatrix}^T \times$$

$$\begin{pmatrix} \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \\ \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{M} \sum_{i=0}^C w_{1i} & \frac{1}{M} \sum_{i=0}^C w_{2i} & \dots & \frac{1}{M} \sum_{i=0}^C w_{Ni} \end{pmatrix}$$

Si usamos  $\frac{1}{M} \sum_{i=0}^C w_{1i}$  como una constante y la llamamos  $\bar{w}_N$  y tras hacer

las restas de los vectores de características en cada paréntesis:

$$\begin{pmatrix} w_{11} - \bar{w}_1 & w_{12} - \bar{w}_1 & w_{11} - \bar{w}_1 & \dots & w_{1N} - \bar{w}_1 \\ w_{21} - \bar{w}_2 & w_{22} - \bar{w}_2 & w_{21} - \bar{w}_2 & \dots & w_{2N} - \bar{w}_2 \\ \dots & \dots & \dots & \ddots & \dots \\ w_{N1} - \bar{w}_N & w_{N2} - \bar{w}_N & w_{N1} - \bar{w}_N & \dots & w_{NN} - \bar{w}_N \end{pmatrix} \times \begin{pmatrix} w_{11} - \bar{w}_1 & w_{21} - \bar{w}_1 & w_{11} - \bar{w}_1 & \dots & w_{1N} - \bar{w}_1 \\ w_{12} - \bar{w}_2 & w_{21} - \bar{w}_2 & w_{21} - \bar{w}_2 & \dots & w_{2N} - \bar{w}_2 \\ \dots & \dots & \dots & \ddots & \dots \\ w_{1N} - \bar{w}_N & w_{N2} - \bar{w}_N & w_{N1} - \bar{w}_N & \dots & w_{NN} - \bar{w}_N \end{pmatrix} =$$

$$\begin{pmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \dots & \text{cov}(x_1, x_N) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) & \dots & \text{cov}(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(x_N, x_1) & \text{cov}(x_N, x_2) & \dots & \text{cov}(x_N, x_N) \end{pmatrix} = E2$$

Es decir, E2 es la matriz de covarianzas del vector de características,  $E2 = \text{cov}(X)$ .



6. Considerar la matriz hat definida en regresión,  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ , donde  $\mathbf{X}$  es una matriz  $N \times (d+1)$ , y  $\mathbf{X}^T\mathbf{X}$  es invertible.

Las demostraciones del ejercicio han sido obtenidas de un canal de youtube que se dedica a la divulgación del aprendizaje automático [aquí](#).

a) **Mostrar que  $\mathbf{H}$  es simétrica**

Suponemos que  $\mathbf{H}^T = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  ya que si fuesen simétricas  $\mathbf{H}$  y  $\mathbf{H}^T$  serían iguales así:

$$\begin{aligned}\mathbf{H}^T &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ &= \mathbf{X} [ (\mathbf{X}^T\mathbf{X})^{-1}]^T\mathbf{X}^T \\ &= \mathbf{X} ( (\mathbf{X}^T\mathbf{X})^T)^{-1}\mathbf{X}^T \\ &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ &= \mathbf{H}\end{aligned}$$

b) **Mostrar que es idempotente  $\mathbf{H}^2 = \mathbf{H}$**

Si  $\mathbf{H}$  fuese idempotente entonces la matriz al cuadrado sería lo mismo que sin elevarla a nada así pues:

$$\begin{aligned}\mathbf{H}^2 &= (\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T)(\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T) \\ &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X})(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ &= \mathbf{X} (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ &= \mathbf{H}\end{aligned}$$

c) **Que representa la matriz  $\mathbf{H}$  en un modelo de regresión?**

Si  $\mathbf{H}$  es un modelo de regresión entonces todos sus vecinos no son negativos, de esta forma:

$$\mathbf{H}\mathbf{w} = \lambda\mathbf{w} \rightarrow$$

$$\mathbf{H}^2\mathbf{w} = \lambda\mathbf{H}\mathbf{w} = \lambda(\lambda\mathbf{w}) ; \quad \text{pero como } \mathbf{H} \text{ es idempotente: } \rightarrow$$

$$\mathbf{H}^2\mathbf{w} = \mathbf{H}\mathbf{w} = \lambda\mathbf{w} \rightarrow$$

$$\lambda^2\mathbf{w} = \lambda\mathbf{w} \rightarrow$$

$$\lambda(\lambda-1)w = 0 \rightarrow \lambda = 0 \text{ o } 1$$

7. La regla de adaptación de los pesos del Perceptron ( $w_{\text{new}} = w_{\text{old}} + yx$ ) tiene la interesante propiedad de que los mueve en la dirección adecuada para clasificar  $x$  de forma correcta. Suponga el vector de pesos  $w$  de un modelo y un dato  $x(t)$  mal clasificado respecto de dicho modelo. Probar que la regla de adaptación de pesos siempre produce un movimiento en la dirección correcta para clasificar bien  $x(t)$ .

$x(t)$  clasifica según  $w^t x(t)y(t) > 0$  si está bien clasificada y  $w^t x(t)y(t) \leq 0$  si no tal y como vemos en teoría si estuviera mal clasificada entonces

$y(t) \neq \text{sing}(W^{T(t)} + x(t)y(t))$  en cambio  $y(t) = \text{sing}(W^{T(t)} + x(t)y(t))$  si lo estuviese bien. Teniendo esto en cuenta y siguiendo la regla de adaptación vemos que pase en  $t+1$ :  $W^{T(t+1)}x(t)y(t)$  donde  $W^{T(t+1)} = W^{T(t)} + x(t)y(t)$

De esta forma si  $W^{T(t+1)}x(t)y(t) > W^{T(t)}x(t)y(t) \rightarrow (W^{T(t)} + x(t)y(t))x(t)y(t) > W^{T(t)}x(t)y(t)$  lo que supone añadir un peso extra  $x(t)y(t)$  que hará que el valor de clasificación sea mejor en  $x(t+1)$  que en  $x(t)$ . Si no está bien clasificado lo provocará es que pueda ser clasificado más fácilmente en los  $x(t+n)$ .

8. Sea un problema probabilístico de clasificación binaria cuyas etiquetas son  $\{0,1\}$ , es decir  $P(Y=1) = h(x)$  y  $P(Y=0) = 1-h(x)$ .

- a) Dar una expresión para  $P(Y)$  que sea válida tanto para  $Y=1$  como para  $Y=0$ .

Al igual que en las transparencias de teoría tenemos que:

$$P(y|x) = \begin{cases} h(x) & \text{for } y=+1 \\ 1-h(x) & \text{for } y=0 \end{cases}$$

- b) Considere una muestra de  $N$  v.a. independientes. Escribir la función de Máxima Verosimilitud para dicha muestra.
- c) Mostrar que la función  $h$  que maximiza la verosimilitud de la muestra es la misma que minimiza

$$E_{in}(w) = \sum_{n=1}^N \mathbb{I}(y_n=1) \ln\left(\frac{1}{h(x_n)}\right) + \mathbb{I}(y_n=0) \ln\left(\frac{1}{1-h(x_n)}\right)$$

donde  $\mathbb{I}(\cdot)$  vale 1 o 0 según que sea verdad o falso respectivamente la expresión en su interior.

- d) Para el caso  $h(x) = \sigma(w^T x)$  mostrar que minimizar el error de la muestra en el apartado anterior es equivalente a minimizar el error muestral.

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

9. Mostrar que la regresión logística se verifica:

$$\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{-y_n w^T x_n}} = \frac{1}{N} \sum_{n=1}^N -y_n x_n \sigma(-y_n w^T x_n)$$

Argumentar sobre si un ejemplo mal clasificado contribuye al gradiente más que un ejemplo bien clasificado.

10. Definimos el error en un punto  $(x_n, y_n)$  por  $e_n(w) = \max(0, -y_n w^T x_n)$

Argumentar si con esta función de error el algoritmo PLA puede interpretarse como SGD sobre  $e_n$  con tasa de aprendizaje  $\eta = 1$ .

Sabemos que la forma de actualizar los pesos en el SGD es en primer lugar es ver el gradiente en el instante  $t$  con unos pesos  $w$ :

$$g_t = \text{gradiente}(E_{in}(w(t)))$$

Cambiamos la dirección en la que nos movemos y además actualizamos los pesos como en las transparencias de teoría:

$$v = -g_t$$

$$w(t+1) = w(t) + \eta v \rightarrow w(t+1) = w(t) + \eta \text{gradiente}(E_{in}(w(t)))$$

Ahora tenemos en cuenta que el error en un punto en un determinado instante

$t$  es  $e_n(\mathbf{w}) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$ , derivamos respecto de  $\mathbf{w}^T$ , obteniendo el gradiente de la función  $e_n(\mathbf{w})$ :

$$f(\mathbf{w}) = -y_n \mathbf{w}^T \mathbf{x}_n \rightarrow f'(\mathbf{w}) = -y_n \mathbf{x}_n$$

Actualizamos los pesos para  $t+1$ , además hay que tener en cuenta que en el enunciado pone que  $\mu = 1$ :

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \mu (-y_n \mathbf{x}_n) = \mathbf{w}(t) - y_n \mathbf{x}_n$$

Esta es la misma fórmula que la que tiene la actualización de los pesos en el algoritmo PLA.