

# Regresión y clasificación: especialización ML

miércoles, 16 de agosto de 2023 21:56

$x$ : feature variable, variable de entrada.  $\hat{y}$ : predicción del modelo.

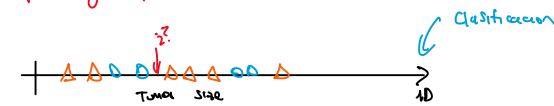
$y$ : target variable, variable de salida.

$m$ : número de ejemplos de entrada.

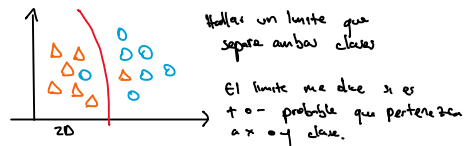
$(x_i, y_i)$ : un solo elemento del conjunto de entrenamiento.

$(x^{(m)}, y^{(m)})$ : nuestro ejemplo de entrenamiento.

## Aprendizaje supervisado:

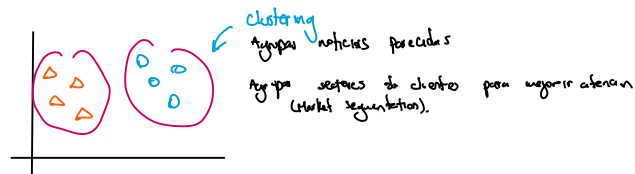


Clase = categoría > salida



## Aprendizaje no supervisado

Datos no etiquetados, el algoritmo debe de hallar un patrón en los datos.



→ Detección de anomalías

→ Reducción de la dimensionalidad.

{ Comprimir un conjunto de datos en otro más pequeño  
preservando la mayor información posible }

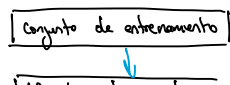
> Predecir un número { 2.2, 1.7, 1.8 } se denomina un **problema de regresión**.

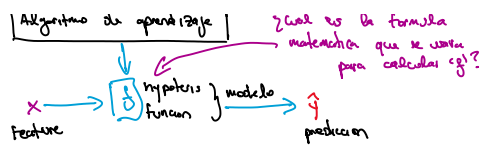
- Número infinito de posibles salidas.

> Predecir una categoría o categoría discreta es un **problema de clasificación**.

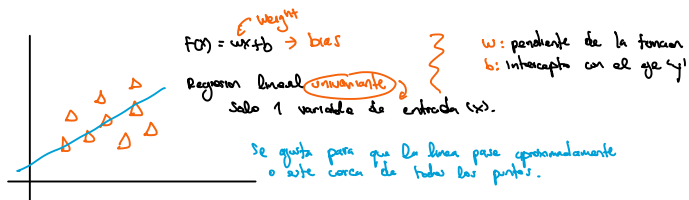
- Número finito, discreto y finito de posibles salidas.

## Funcionamiento





## Regresión lineal (supervisado).



$\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$  error → que tan lejos está la predicción del objeto, del valor real.

(se error cuadrático)

$\frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$  MSE (error cuadrático medio) → Función de coste.

Hay usada en todos los problemas de regresión.

$\hat{y} = g(w, b) = wx + b$  → una predicción es

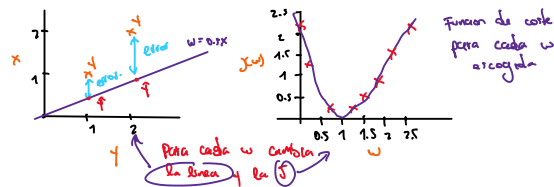
• Hallar valores de  $w$  y  $b$  que hagan más pequeña la función de coste.

Objetivo: minimizar  $J(w, b)$

función d' coste

• Función de coste: No dice que también funciona el modelo.

• Parámetros del modelo son las variables que se pueden ajustar durante el entrenamiento.



Como depende de  $w$  y  $b$  la función de coste se va de mejor manera en un gráfico 2D, sin embargo, también se hace uso de un gráfico de contorno para visualizar las funciones de coste.

Descent de gradiente es un algoritmo que se usa para minimizar cualquier función de entrada.

• Comienza con una hipótesis inicial eligiendo dos parámetros como 0.

• En cada iteración cambia los parámetros tratando de reducir la función de coste.

• Si quisiéramos saber dirección hacia qué punto llegar más rápido. Dirección del descenso más pronunciado.

• En cada iteración el algoritmo también se halla la dirección del descenso más pronunciado.

• Repetir hasta que la función converja, hallar un mínimo local, punto en donde por cada paso los valores de los parámetros no cambian mucho.

→ En cada paso

$w \leftarrow w - \alpha \frac{\partial}{\partial w} J(w, b)$   $b \leftarrow b - \alpha \frac{\partial}{\partial b} J(w, b)$

asignación  $w \leftarrow w - 1$

asignación  $b \leftarrow b - 1$

Valores anteriores de  $w$  ( $w_{i-1}$ ) Valores anteriores de  $b$  ( $b_{i-1}$ )

2: tasa de aprendizaje

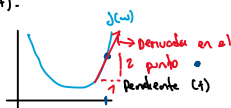
numero positivo entre 0 y 1, controla que tan grande es la distancia para valores a calcular al descenso del gradiente.

$\frac{\partial}{\partial w} J(w, b)$ : Derivada de la función de coste  $J(w, b)$ .

nos dice en que dirección debemos de ir, me ayuda a calcular la pendiente de la función, la dirección de la función.

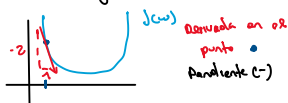
• Se debe de actualizar los parámetros  $w$  y  $b$ .

- Cuando la recta tangente (la derivada) apunta hacia arriba a la derecha la pendiente es positiva (+).



Como  $w$  está positivo se le resta (-)

- Cuando la recta tangente (la derivada) apunta hacia abajo a la derecha (cambia, izquierda) la pendiente es negativa (-).



Como  $w$  es negativa se le suma (+)

- Al elegir un  $\alpha$  pequeño se disminuye el rate de forma muy lenta.

- Si  $\alpha$  es muy grande se puede sobrepasar el mínimo local de la función; puede llegar a no converger e incluso divergir.

- Si ya se está en un mínimo local la pendiente es 0 por lo que,  $w$  quedará sin cambiar por más.

- A medida que se acerca al mínimo local la pendiente (derivada) será cada vez menor, y por ende, el paso también será menor.

- **Batch gradient descent**: En cada paso del descenso de gradiente usamos todas las datos de entrenamiento.

## Múltiples variables

$x_j$  =  $j^{\text{th}}$  Feature

$n$  = # de features

$x^{(i)}$  = features del  $i^{\text{th}}$  ejemplo de entrenamiento.

$x_j^{(i)}$  = valor del feature  $j$  en el  $i^{\text{th}}$  ejemplo de entrenamiento.

**Regresión logística**

$$f(w,b;x) = 0.1x_1 + 0.2x_2 + 0.3x_3 + 0.4x_4 + b$$

$$f(w,b;x) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n] \text{ (vector)}$$

$$\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n] \text{ (vector) features}$$

$$g(\vec{w}, b; \vec{x}) = \vec{w} \cdot \vec{x} + b$$

↑  
producto punto

**Vectorización**: convertir los parámetros a vectores (numpy)

✓ Hace el código más corto

✓ Hace que el código sea más eficiente (realiza operaciones en paralelo).

en un tiempo.

$$\text{top-}w = w - \alpha \frac{\partial}{\partial w} J(w,b)$$

$$\text{top-}b = b - \alpha \frac{\partial}{\partial b} J(w,b)$$

$$w = \text{top-}w$$

$$b = \text{top-}b$$

→ Coge todos los valores de las variables y realiza los cálculos correspondientes en paralelo.

! En el descenso de gradiente con  $n$  parámetros el time que actualizas los  $n$  valores del vector  $\vec{w}$  y luego  $\vec{b}$   $n$  veces para minimizar la función de coste.

! Hay otro algoritmo que sirve para hallar  $w$  y  $b$  solo en regresión lineal y su hipótesis: función normal.

Desventajas:

✓ lento cuando el # de parámetros es alto ( $> 10,000$ )

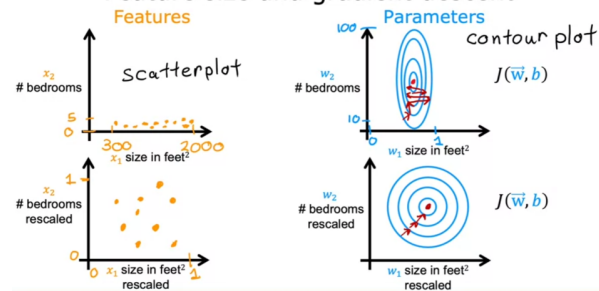
Algunas bibliotecas de ML pueden usarlo para resolver  $w$  y  $b$  (ver the hard)

Feature scaling: técnica.

A los datos de distintas magnitudes al contorno de las funciones asociadas a los mismos puntos se muy grande dado la varianza en su magnitud, en este caso es mejor rescatar un escalado de su dato, para que todos tengan la misma base y magnitud.

El escalado de parámetros permite que el modelo converja rápidamente.

### Feature size and gradient descent



Escalamiento: Forma 1

coger el máximo de  $c_j$  de las características y dividir  $c_j$  de los valores de la característica con su respectivo máximo.

Escalamiento: Forma 2 (Mean normalization).

Se rescatan las características para q' sus valores queden centrados alrededor del 0, sus valores pueden ser  $+ \infty$  o  $- \infty$ , scale estar entre  $-1$  y  $1$ .

✓ Halla la media ( $\mu$ )

✓ Halla el máximo y mínimo del intervalo.

✓  $x_1 = \frac{x_1 - \mu}{\max - \min}$ ; para  $c_j$  de los datos del feature

Escalamiento: Forma 3 (z score)

✓ Calcular la desviación estándar (std), y la media ( $\mu$ ) de la columna.

✓  $x_1 = \frac{x_1 - \mu}{std}$ ; para  $c_j$  de los datos del feature.

! Se debe de tratar de q' luego de escalar el rango del feature oscile lo mas cerca posible de  $-1$  a  $1$ .

Constante de rango muy pequeños ( $-0.001 \leq x \leq 0.001$ )

y en iteraciones muy grandes ( $\sim 10^6$  a  $10^7$ )

Para saber si el descenso del gradiente funciona se puede graficar los valores de  $J$  (función de costo) para cada una de las iteraciones, entre otras se llaman 'learning curve'.

$J$  debería de disminuir en cada una de las iteraciones.

→  $J$  crease, quiere decir que la tasa de aprendizaje es muy alta.

Cuando la curva deja de crecer el descenso del gradiente ha convergido.

Automatic convergence test:

Es como un stopper, si  $J(x_{old}) \leq \epsilon \text{ (stopper)}$   
debe haber convergencia.

Cómo elegir la tasa de aprendizaje ( $\alpha$ )

✓ El  $\alpha$  debería de bajar en cada una de las iteraciones.

✓ Bajar  $\alpha$  a un número muy pequeño y ver si así el  $\alpha$  baja.  
→ degradación.

✓ Verificar con distintos valores de  $\alpha$  y trazar la función de costo.

$\alpha = [0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, 3, \dots]$

Feature engineering: Transformar o combinar las características (datos) originales para facilitar que el algoritmo prediga con más precisión.

Polynomial Regression: Ajusta curvas del orden  $x^0, x^1, x^2$   
puede comenzar con una curva cuadrática y continuar con una función cúbica ( $w_0x + w_1x^2 + w_2x^3 + \dots$ )

! El escalado de parámetros es de suma importancia en estos modelos de  $x^2, x^3$  para no introducir sesgo en el modelo y que los parámetros tengan rangos comparables.

Classification

Classification Binaria: Problemas en los cuales solo hay 2 posibles clases - categorías de salida.

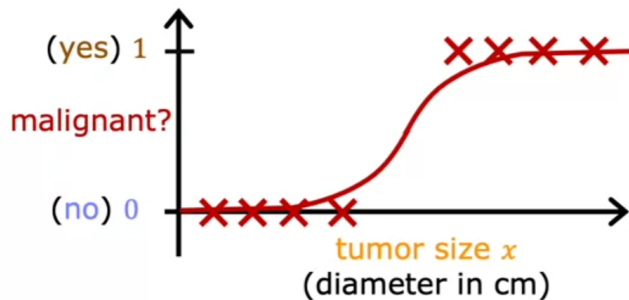
↑ Falso | Verdadero ← clase positiva  
= presencia

clase negativa  
= ausencia

! Regresión lineal no es buena en clasificación  
ya que calcula las líneas y mejor se ajusta a los datos recalculando la frontera de decisión.

Regresión logística (clasificación binaria)

✓ Se ajusta una curva en forma de 'S' al conjunto de datos.



✓ El algoritmo devuelve una probabilidad asociada al dato dado (0-1)

! La probabilidad de que la clase sea igual a 1 dado un determinado  $\vec{x}$ .

Funcion sigmoide o logistica:

$$g(z) = \frac{1}{1+e^{-z}} \quad 1 = \frac{1}{1+e^{-x}}$$

Matematica de la reg. logistica

$$z = \vec{w} \cdot \vec{x} + b \quad \text{key: lineal}$$

$$g(z) = \frac{1}{1+e^{-z}}$$

$$g(\vec{w}, b | \vec{x}) = \text{PCY} = 1 | \vec{x}, \vec{w}, b$$

Probabilidad que  $y$  sea 1, dado un input  $\vec{x}$ , y parametros  $\vec{w}$  y  $b$ .

$$g(\vec{w}, b | \vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x} + b}_z)$$

$$= \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$$

Link de decision: Neutral al predecir si  $y$  es 0 o 1.

$z = \vec{w} \cdot \vec{x} + b = 0$ , este link puede ser lineal

o polinomial.

! El link de decision puede cambiar, puede ser 0.5 así como puede ser 0.9; depende de que tan sensible se quiera al modelo.

Funcion de costo:

Al aplicar MSE (error cuadrático medio) a la

funcion de la regresion logistica me da una

funcion no convexa, con muchos mínimos locales, causando la inestabilidad del descenso de gradiente.

$$L(\vec{w}, b | \vec{x}^{(i)}, y^{(i)}) = \begin{cases} -\log(g(\vec{w}, b | \vec{x}^{(i)})) & \text{si } y^{(i)} = 1 \\ -\log(1 - g(\vec{w}, b | \vec{x}^{(i)})) & \text{si } y^{(i)} = 0 \end{cases}$$

Funcion de perdida

Funcion de perdida: Ver que tan bien o mal es algoritmo

actua para un dato en particular.

Funcion de costo: Que tan bien el algoritmo lo hizo para un conjunto de datos (Promedio de todas las funciones de perdida)

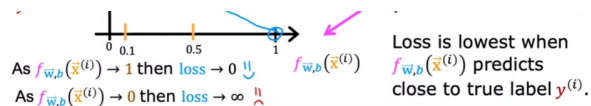
! La perdida es mas baja cuando  $g(\vec{w}, b | \vec{x}^{(i)})$  predice

cercas del valor real  $y^{(i)}$

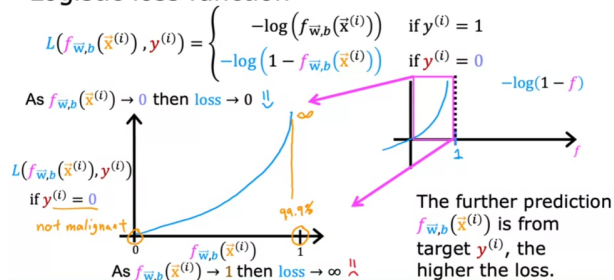
Logistic loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$





## Logistic loss function



! Costo

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(f_{w,b}(\vec{x}^{(i)}), y^{(i)})}_{\text{perdida}}$$

$$\begin{cases} -\log(f_{w,b}(\vec{x}^{(i)})) & \text{si } y^{(i)} = 1 \\ -\log(1 - f_{w,b}(\vec{x}^{(i)})) & \text{si } y^{(i)} = 0 \end{cases}$$

! Este costo genera una función convexa, por lo que el descenso de gradiente puede alcanzar un mínimo global.

$$L(f_{w,b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{w,b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{w,b}(\vec{x}^{(i)}))$$

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{w,b}(\vec{x}^{(i)}), y^{(i)})]$$

## Gradient descent

cost

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{w,b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{w,b}(\vec{x}^{(i)}))]$$

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \quad \rightarrow \quad \frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) \quad \rightarrow \quad \frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)})$$

} simultaneous updates

Underfitting: Es cuando el modelo no se ajusta bien a los datos de entrenamiento.

→ El algoritmo tiene alto sesgo (bias).

Bias { El algoritmo está sesgado al reconocer a una población, etnia, etc.  
El algoritmo no se adapta bien al conjunto de datos.

Generalización: Que el modelo se vea bien produciendo

datos no vistos anteriormente.

**Overfitting:** Se ajusta demasiado bien al conjunto de datos

y no **generaliza** haciendo que sea malo al predecir nuevos datos, el algoritmo tiene una alta **varianza**.

**Abordando el overfitting:**

1. Recopilar mas datos de entrenamiento.

2. Ver si se pueden incluir/excluir características. (feature selection).

Muchas features + datos insuficientes = Overfit.

❌ El algoritmo descarta parte de la información sobre los casos.

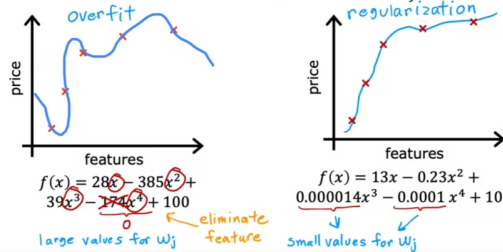
3. Regularización: Reduce el impacto de las features sin llegar al punto de eliminarlas.

✓ Fomentar que el algoritmo baje los valores de los

parámetros ( $w_j$ ) sin exigir que se eliminen (se pongan como 0)

## Regularization

Reduce the size of parameters  $w_j$ .



Solo se regulariza  $w_j$  - un y no b, no hay prácticamente diferencia

si se regulariza b o no.

**Regularización:**

✓ Minimizar  $w_j$  - un al hacerlos  $\approx 0$ .

✓ Entre mas bajos los valores de  $w_j$  mas simple es el modelo, incluso algunos  $w_j$  pueden ser 0 eliminando una característica.

✓ Entre mas simple es el modelo, menos probable que se de sobreajuste.

✓ Se añade a la función de costo  $J(\theta, b)$  todos los valores de  $w_j$  para que se aproximen a 0.

$$J(\theta, b) = \frac{1}{2m} \sum_{i=1}^m (h(\theta, b(x_i)) - y_i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\lambda$ : "lambda" es el **parámetro de regularización**. Representa que tan importante o no es la regularización con respecto a minimizar el error, es por esto que balancea ambos métodos.

✓ Si es 0 no se está regularizando



Si el valor muy grande se pone  
mucha más a la regularización,  
por lo que  $\vec{w} \rightarrow 0$  y  $f(x) \rightarrow 0$   
resultando en una línea horizontal.

Regularized linear regression

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f(\vec{x}^{(i)}) - y^{(i)}) \quad \text{No se regulariza } b$$

✓ En cada iteración se multiplica  $w_j$  por un número  
menor a 1 (por muy poco) haciendo que los valores de  
 $\vec{w}$  bajen en caso de las iteraciones.

$w_j (1 - \alpha \frac{\lambda}{m})$  ———> param  
regularization  
factor  
apartado ———> término del  
conjunto de  
entrenamiento

Regularized logistic regression:

$$J(\vec{w}, b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log(f(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f(\vec{x}^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Min  $\vec{w}, b \rightarrow w_j \downarrow$

