



APRENDIZAJE DE MÁQUINA

REGRESIÓN LINEAL

AGENDA

01 Aprendizaje supervisado

Datos de entrenamiento, hipótesis, un ejemplo.

02 Algoritmo LMS

Modelo lineal, función de costo, descenso por gradiente.

03 Ecuaciones normales

Modelo lineal matricial, derivando las ecuaciones normales

04 Interpretación probabilística

Suposiciones para errores, función de verosimilitud

05 Funciones de base

Derivación, funciones, máxima verosimilitud





APRENDIZAJE SUPERVISADO

DEFINICIONES

NOMENCLATURA

APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA



Comenzamos con un conjunto de datos que representan las superficies habitacionales del primer piso y los precios de venta 1,460 casas en Ames, Iowa.

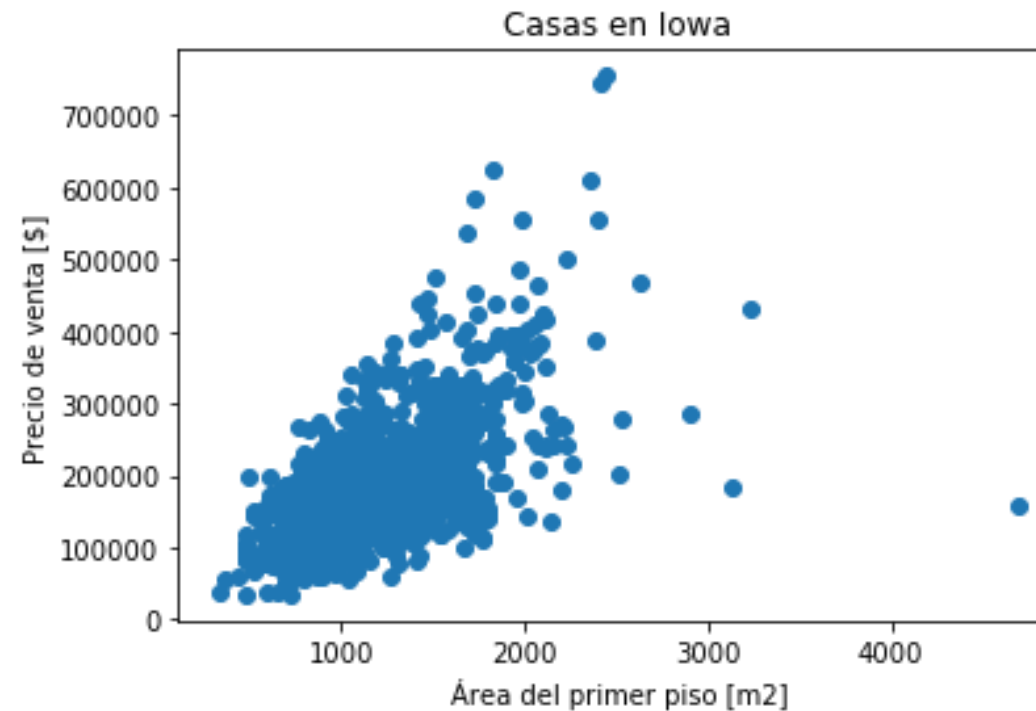
1 st floor square feet	Sale Price
856	208500
1262	181500
920	223500
961	140000
1145	250000
⋮	⋮

APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA

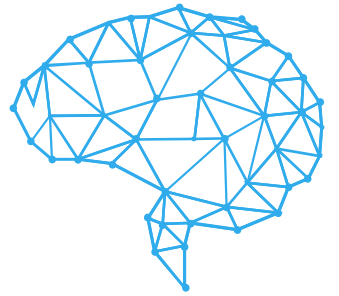


Graficamos los 1,460 datos con Python:



APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA



La pregunta que surge naturalmente sería:

**¿CÓMO PREDECIMOS LOS PRECIOS DE OTRAS
CASAS EN AMES, IOWA?**

APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA



Para responder la pregunta primero se establecerá la nomenclatura que se estará usando:

Variable / Símbolo	Descripción	Ejemplo
$x^{(i)}$	Variable de entrada o características	El área del primer piso de la casa.
$y^{(i)}$	Variable de salida, objetivo o de respuesta	El precio de venta de la casa.
$(x^{(i)}, y^{(i)})$	Ejemplo de entrenamiento	(área, precio)
m	Número de ejemplos de entrenamiento	1460 casas con su área del 1° piso y el precio.
$\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$	Conjunto de datos de entrenamiento	NA
\mathcal{X}	Espacio de los valores de entrada	NA
\mathcal{Y}	Espacio de los valores de salida	

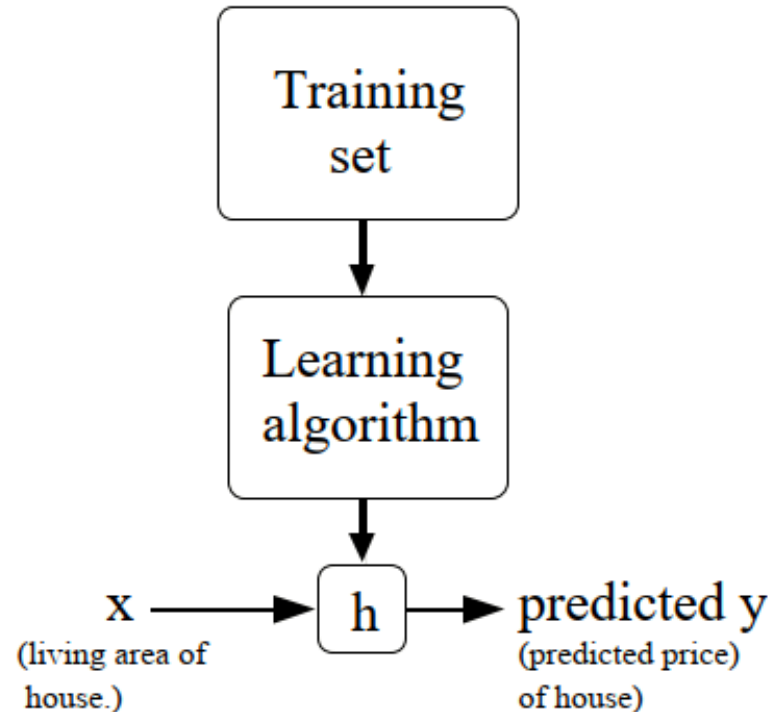
APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA



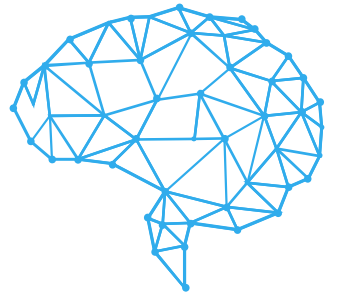
El **objetivo principal** de cualquier **algoritmo supervisado** es **aprender** una **función** $h: \mathcal{X} \rightarrow \mathcal{Y}$, de tal manera que $h(x)$ pueda **predecir** el correspondiente valor de y .

Donde h se define como la **hipótesis**.



APRENDIZAJE SUPERVISADO

DEFINICIONES Y NOMENCLATURA



Existen **2 tipos** de **problemas** de **aprendizaje supervisado**:

Tipo de problema	Descripción
Problema de regresión	y toma valores continuos.
Problema de clasificación	y toma valores discretos.



ALGORITMO LMS

MODELO LINEAL

FUNCIÓN DE COSTO

DESCENSO POR GRADIENTE

ALGORITMO LMS

MODELO LINEAL



Ahora veamos el mismo modelo pero con otra variable más: la superficie habitacional del segundo piso.

1 st floor square feet	2 nd floor square feet	Sale Price
856	854	208500
1262	0	181500
920	866	223500
961	756	140000
1145	1053	250000
⋮	⋮	⋮

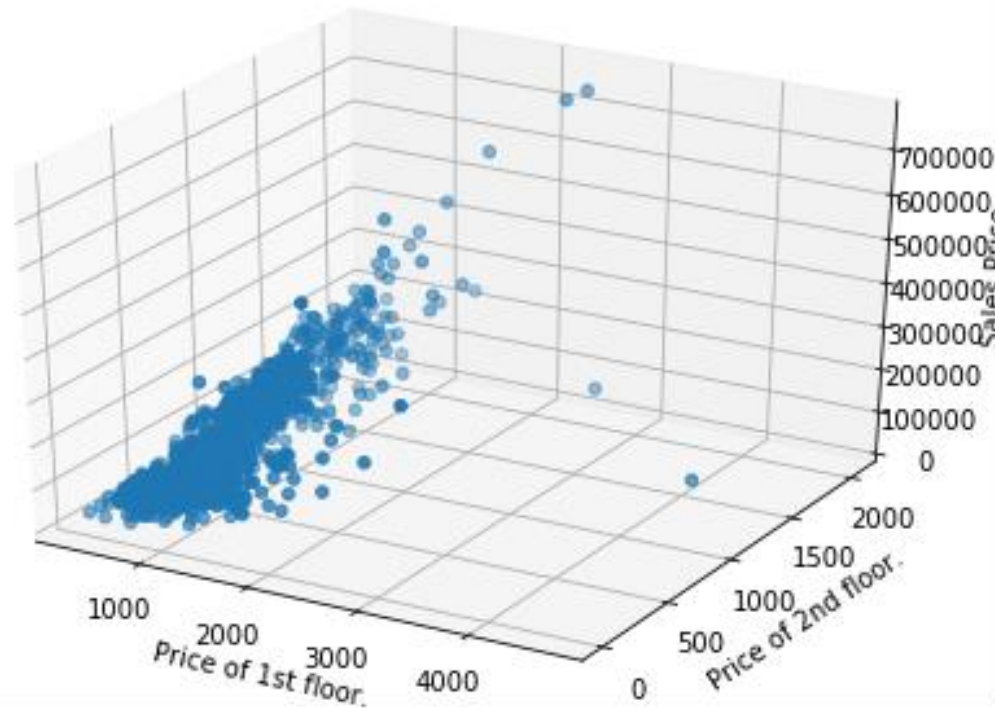
En este caso cada $x^{(i)} \in \mathbb{R}^2$, por lo que se define como un vector $x^{(i)} = [x^{(i)}_1 \ x^{(i)}_2]$.

ALGORITMO LMS

MODELO LINEAL



Graficamos de nuevo los 1,460 datos con Python:



ALGORITMO LMS

MODELO LINEAL



Se plantea la **hipótesis** de que los datos se distribuyen de **manera lineal**, por lo tanto:

$$h_w(x) = w_0 + w_1x_1 + w_2x_2$$

En este caso las **variables** w_i se definen como los **parámetros** o **pesos** que **parametrizan** el **espacio** de todas las **funciones lineales** que mapean de \mathcal{X} a \mathcal{Y} .

Se **simplifica** la expresión en **forma vectorial**, donde n representa el **número** de **variables** de **entrada** (sin contar a x_0), :

$$h(x) = \sum_{i=1}^n w_i x_i = w^T x$$

ALGORITMO LMS

FUNCIÓN DE COSTO

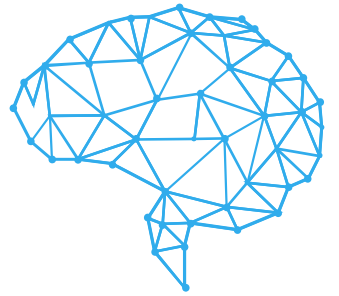


Ahora la pregunta que surge sería:

**¿CUÁL ES LA MEJOR COMBINACIÓN DE PARÁMETROS w
QUE RESULTE EN LA MEJOR HIPÓTESIS h ?**

ALGORITMO LMS

FUNCIÓN DE COSTO



Para **resolver** la **pregunta**, necesitamos **medir** el **error** que producen nuestras **estimaciones** realizadas por h .

Sabemos que podemos **medir** el **error** de una sola **estimación** i como una **diferencia** **entre** la estimación $h_w(x^{(i)})$ y la variable de respuesta $y^{(i)}$.

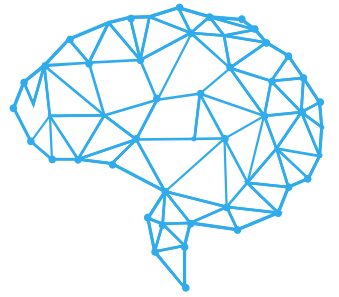
$$e = h_w(x^{(i)}) - y^{(i)}$$

Para tener siempre cantidades positivas de error, se eleva al cuadrado.

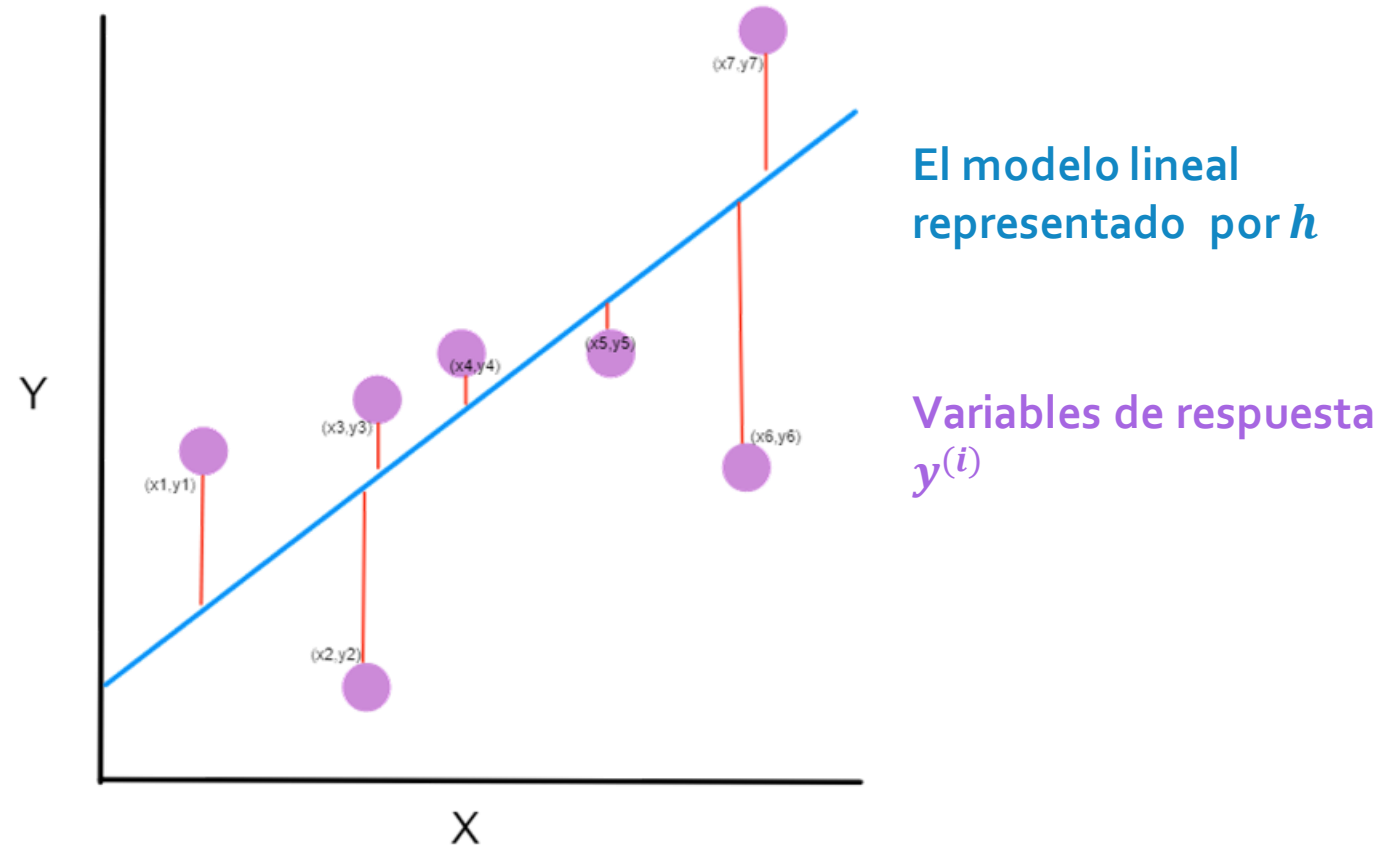
$$e^2 = (h_w(x^{(i)}) - y^{(i)})^2$$

ALGORITMO LMS

FUNCIÓN DE COSTO

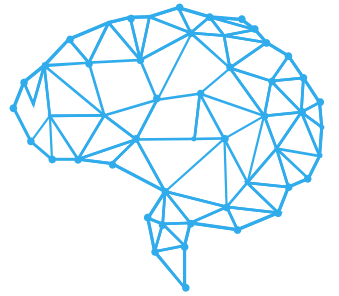


Se **interpreta** el **error** que estamos calculando de **manera gráfica** con un **ejemplo**:



A L G O R I T M O L M S

F U N C I Ó N D E C O S T O



Hasta ahora hemos calculado el **error cuadrático** de un solo **dato** de **entrenamiento**. Se calcula el **promedio** de **error cuadrático medio** MSE para **todos** los **datos** de **entrenamiento**.

$$MSE = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

$$MSE = J(w)$$

Donde $J(w)$ es la **función de costo**.

ALGORITMO LMS

DESCENSO POR GRADIENTE



Por lo tanto la **mejor combinación** de pesos w , es aquella que **minimice** la **función de costo** $J(w)$, que **mide** nuestro **error cuadrático medio**.

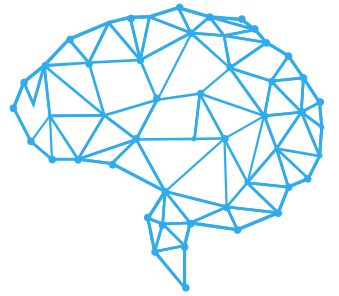
Para **buscar** la **mejor combinación**, diseñemos un **algoritmo** de **búsqueda** que **comience** un **valor inicial** aleatorio de w , y que vaya **actualizando** los **valores** de w hasta que **converja** a un **valor mínimo** de $J(w)$. La **constante** α se define como la **tasa de aprendizaje**.

$$w_j := w_j - \alpha \frac{\partial}{\partial w_j} J(w)$$

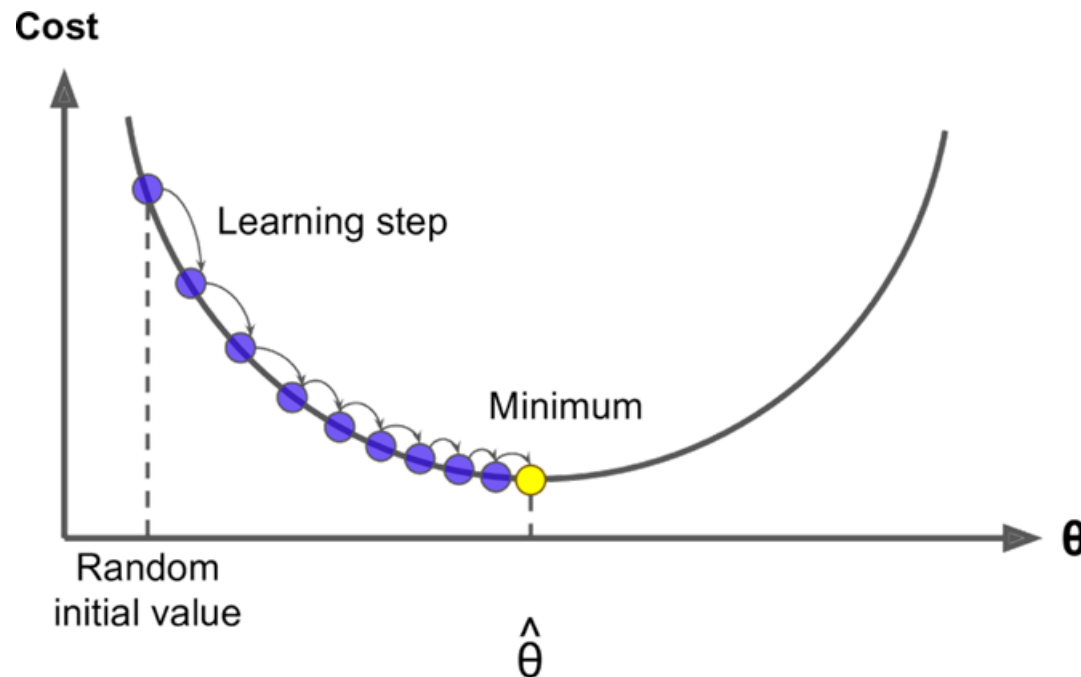
NOTA: la ecuación de arriba **solo actualiza** el **valor** de un solo **peso** w_j de los n pesos que **parametrizan** el **modelo lineal**. En la realidad se **actualizan todos** los **pesos** w_j simultáneamente.

ALGORITMO LMS

DESCENSO POR GRADIENTE



Lo que se hace en el **algoritmo** de **optimización** por **descenso** de **gradiente** es **actualizar** los **pesos** en la **dirección** con **mayor decremento** de $J(w)$.



ALGORITMO LMS

DESCENSO POR GRADIENTE



Se calcula la **derivada** de la **función** de **costo** $J(w)$ con respecto a un **peso específico** w_j .

Derivar el resultado:

$$\frac{\partial}{\partial w_j} J(w) = \frac{1}{m} (h(x) - y) x_j$$

ALGORITMO LMS

DESCENSO POR GRADIENTE



Por lo tanto, la **actualización** de **pesos** por **descenso** de **gradiente** para **un solo dato** de **entrenamiento** quedaría así:

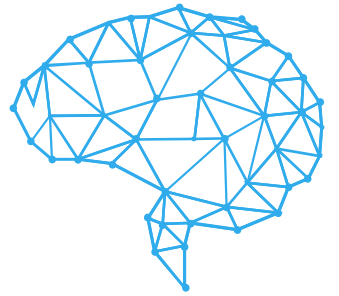
$$w_j := w_j + \frac{\alpha}{m} (y^{(i)} - h(x^{(i)})) x_j$$

A esta **ecuación** se le denomina la **regla** de **actualización LMS** (“Least Mean Squares”) o también la **regla** de **aprendizaje Widrow-Hoff**.

Este **método** también se llama **descenso** de **gradiente** por **lotes**, en donde se **analizan todos** los **datos** de **entrenamiento simultáneamente**.

ALGORITMO LMS

DESCENSO POR GRADIENTE



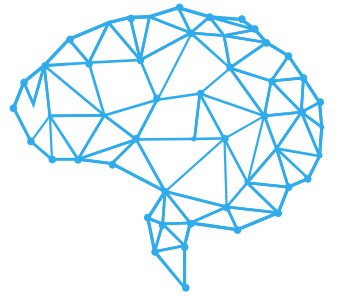
Para m **datos de entrenamiento**, y definiendo como **vectores** a $x^{(i)}$ y a w , la regla de aprendizaje quedaría así:

$$w := w + \frac{\alpha}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x^{(i)}$$

A esta **ecuación** se le denomina la **regla de actualización LMS** (“Least Mean Squares”) o también la **regla de aprendizaje Widrow-Hoff**.

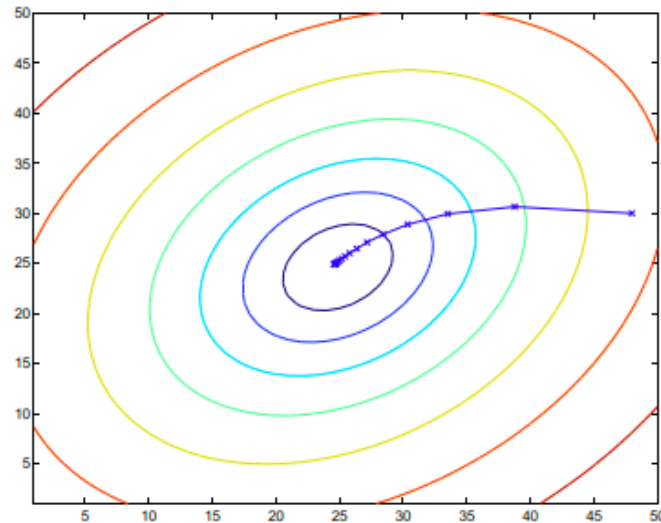
ALGORITMO LMS

DESCENSO POR GRADIENTE



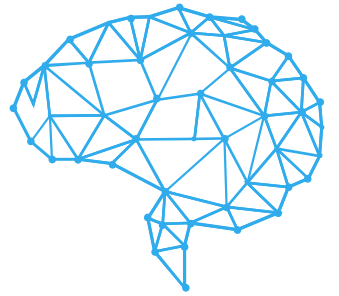
Por lo general, el **método** de **descenso** por **gradiente** puede sufrir de **varios mínimos locales**. En este caso, **para modelos** de regresión **lineal solo existe un único mínimo global**.

En consecuencia, el **algoritmo siempre convergerá** asumiendo que la **tasa de aprendizaje no sea muy alta**.



ALGORITMO LMS

DESCENSO POR GRADIENTE



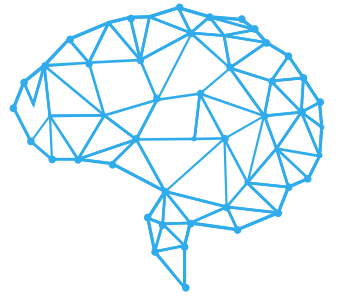
Cuando el **método** solo **observa un dato** de **entrenamiento** a la vez, y **actualiza** los **pesos con un solo dato**, se le denomina como **descenso por gradiente estocástico** o **incremental**.

$$w := w + \alpha(y^{(i)} - h(x^{(i)}))x^{(i)}$$

Esta **variante** del algoritmo **se utiliza** cuando es muy **costosos evaluar** la **actualización** para **conjunto de datos muy grandes** (cuando m es muy grande), pero tiene el problema de **divergencia** en el mínimo.

ALGORITMO LMS

DESCENSO POR GRADIENTE

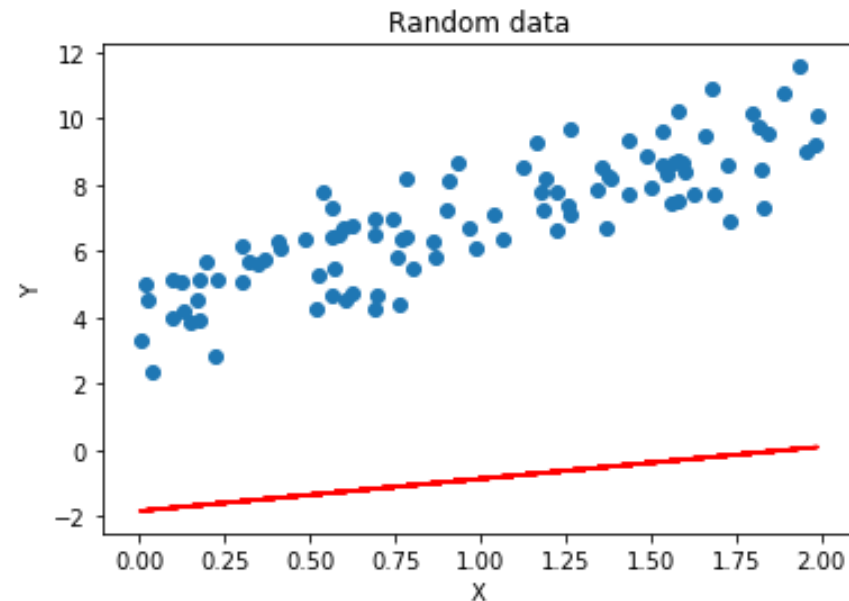


Ejemplo:

Se generaron datos de entrenamiento aleatorios con cierto grado de error e :

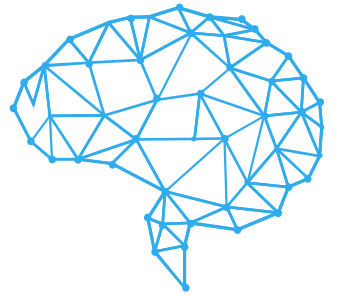
$$y = 3x + 4 + e$$

y se inicializaron aleatoriamente ambos pesos: w_0 y w_1



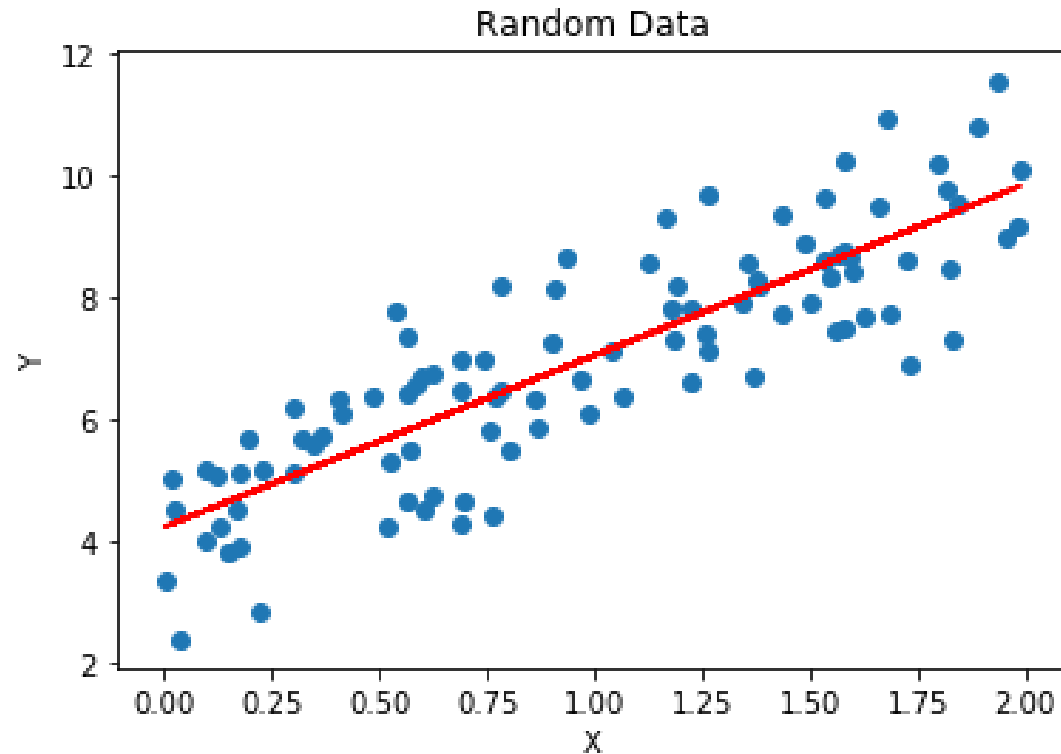
ALGORITMO LMS

DESCENSO POR GRADIENTE



Ejemplo:

Después de 1000 iteraciones se obtuvo el siguiente modelo:



$$w_0 = 4.2174$$
$$w_1 = 2.816$$

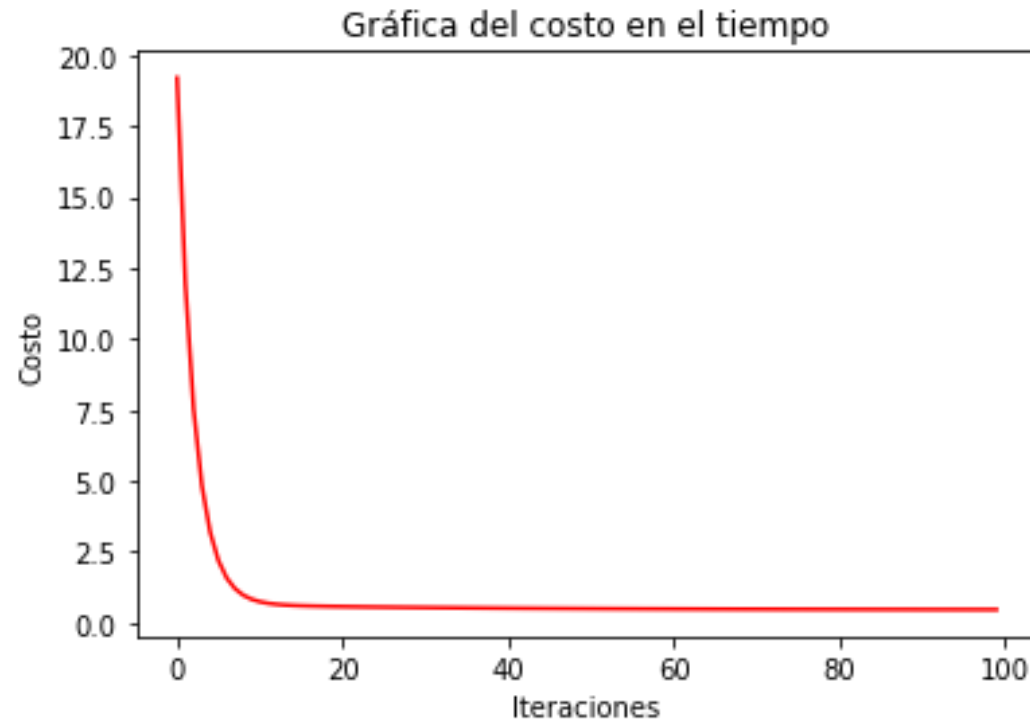
ALGORITMO LMS

DESCENSO POR GRADIENTE



Ejemplo:

Error en función de 100 iteraciones.





ECUACIONES NORMALES

MODELO LINEAL MATRICIAL
DERIVACIÓN

ECUACIONES NORMALES

NOTACIÓN MATRICIAL



Ahora se quiere **calcular** el **mínimo** de la **función** de **costo** de manera **analítica**. Para esto, se introduce la **notación matricial**, lo que **permite expresar** de manera **elegante** las **derivadas** de la **función** de **costo** sin entrar en tanta **verborrea**.

Se representan los **m** datos de **entrenamiento** $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ como una matriz $X \in \mathbb{R}^{m \times n}$, al conjunto de **salidas** como un **vector** $\vec{y} \in \mathbb{R}^m$ y al los **n** **pesos** w_j como un vector $\vec{w} \in \mathbb{R}^n$.

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

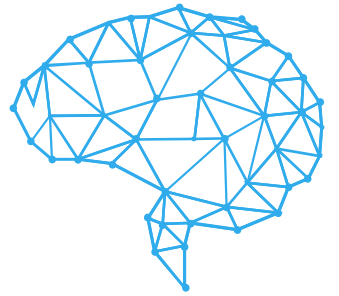
$$X = \begin{bmatrix} -(x^{(1)})^T & - \\ -(x^{(2)})^T & - \\ \vdots & \\ -(x^{(m)})^T & - \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

ECUACIONES NORMALES

NOTACIÓN MATRICIAL



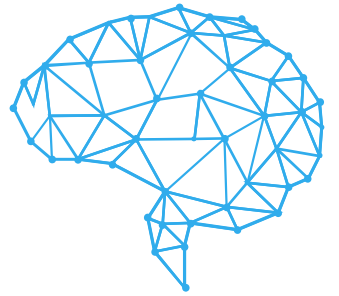
Se representa la **función de error absoluto** en **notación matricial** donde $h_w(x^{(i)}) = (x^{(i)})^T \vec{w}$:

$$X\vec{w} - \vec{y} = \begin{bmatrix} -(x^{(1)})^T \vec{w} - \\ -(x^{(2)})^T \vec{w} - \\ \vdots \\ -(x^{(m)})^T \vec{w} - \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$X\vec{w} - \vec{y} = \begin{bmatrix} h_w(x^{(1)}) - y^{(1)} \\ h_w(x^{(2)}) - y^{(2)} \\ \vdots \\ h_w(x^{(m)}) - y^{(m)} \end{bmatrix}$$

ECUACIONES NORMALES

NOTACIÓN MATRICIAL



El **término del error cuadrático medio** (función de costo):

$$J(\mathbf{w}) = \frac{1}{2m} \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$$

puede ser **representado** en forma **vectorial** usando la **propiedad** $\mathbf{z}^T \mathbf{z} = \sum_i z_i^2$

$$J(\mathbf{w}) = \frac{1}{2m} (\mathbf{X}\vec{\mathbf{w}} - \vec{\mathbf{y}})^T (\mathbf{X}\vec{\mathbf{w}} - \vec{\mathbf{y}})$$

ECUACIONES NORMALES

DERIVACIÓN DE LAS ECUACIONES



Se tienen que **obtener** las **derivadas** de $J(w)$ con respecto al vector w .

$$\nabla_w J(w) = \nabla_w \frac{1}{2m} (X\vec{w} - \vec{y})^T (X\vec{w} - \vec{y})$$

El **gradiente** nos **calcula** todas las **derivadas** con respecto a **todos** los pesos w_j de manera **simultánea**.

Además, gracias a que **representamos** al **conjunto** de **datos** de **entrenamiento** en forma **matricial**, podemos **calcular** el **gradiente** para **todos** los **datos** de **entrenamiento**.

ECUACIONES NORMALES

DERIVACIÓN DE LAS ECUACIONES



Se tiene entonces que:

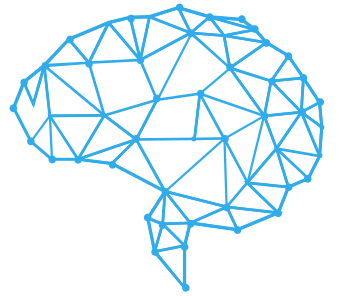
$$\nabla_w J(w) = \nabla_w \frac{1}{2m} (X\vec{w} - \vec{y})^T (X\vec{w} - \vec{y})$$

DERIVAR LA IGUALDAD

$$\nabla_w J(w) = \frac{1}{m} (X^T X\vec{w} - X^T \vec{y})$$

ECUACIONES NORMALES

DERIVACIÓN DE LAS ECUACIONES



Igualando a cero la derivada para encontrar el **mínimo** se obtiene el **vector de pesos** que da ese **mínimo**:

$$\nabla_w J(w) = \frac{1}{m} (X^T X \vec{w} - X^T \vec{y}) = 0$$

$$\vec{w} = (X^T X)^{-1} X^T \vec{y}$$



INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIONES PARA ERRORES
FUNCIÓN DE VEROSIMILITUD

INTERPRETACIÓN PROBABILÍSTICA

MOTIVACIÓN DEL DESARROLLO



¿PORQUÉ MINIMIZAMOS LA FUNCIÓN DE ERROR CUADRÁTICO MEDIO Y NO OTRA FUNCIÓN?

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Se pueden definir a las **variables** de **salida** $y^{(i)}$ como la **hipótesis planteada** h_w más un **error** $\varepsilon^{(i)}$ de **estimación** que captura los **efectos** que **no consideramos** en nuestra **hipótesis** o contempla los **errores aleatorios**.

$$y^{(i)} = h_w(x^{(i)}) + \varepsilon^{(i)}$$

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Se hacen las **siguientes suposiciones** para los **errores** $\varepsilon^{(i)}$ (**muestreo aleatorio**): IID

Errores independientes: la probabilidad de que suscite un error no afecta a la probabilidad de los demás errores.

Errores idénticamente distribuidos: el muestreo de los errores se realiza de la misma distribución de probabilidad.

Distribuidos por una distribución de probabilidad Gaussiana con media $\mu = 0$ y varianza σ^2 .

$$\varepsilon^{(i)} \sim N(0, \sigma^2)$$

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Por lo tanto, la **función de densidad de probabilidad** de $\varepsilon^{(i)}$ está **dada** por:

$$p(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)}$$

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



**¿PORQUÉ PROPONER UN MODELO QUE ESTÁ
DISTRIBUIDO NORMALMENTE?**

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Además, se establece a las **salidas** $y^{(i)}$ y a las **entradas** $x^{(i)}$ como **variables aleatorias**.

Por lo tanto, se realiza la siguiente pregunta:

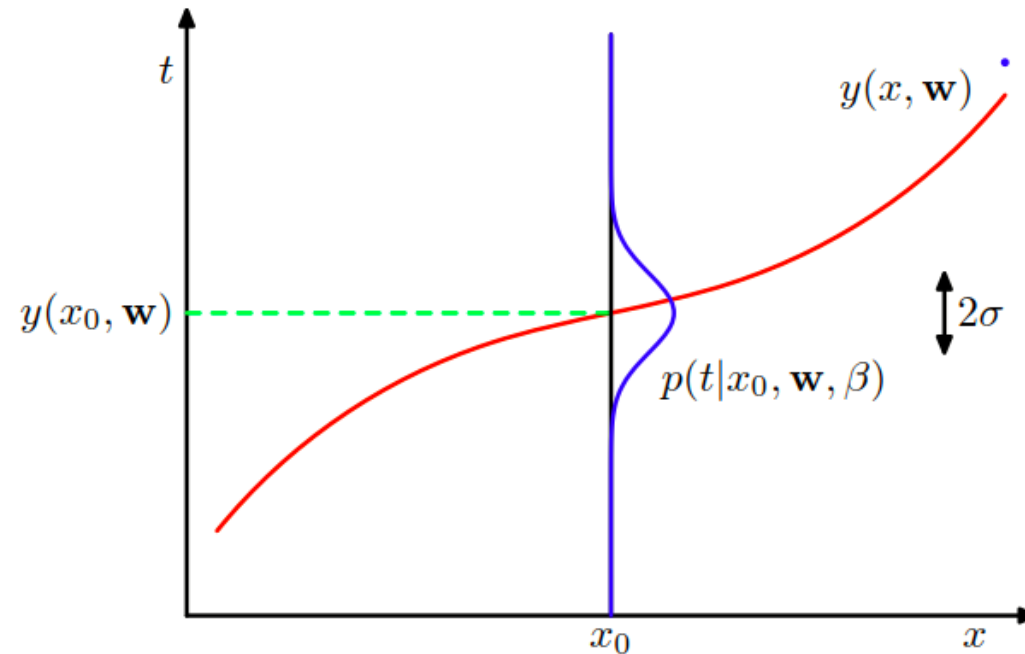
Dado mi conjunto de datos de entrada X y mi vector de pesos \vec{w} ¿cuál es la distribución de probabilidad que modela a mis salidas \vec{y} ?

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Siguiendo las suposiciones de la distribución normal de errores, surge naturalmente que **la distribución de las salidas $y^{(i)}$ siga una forma normal** con media $\mathbf{h}_{\mathbf{w}}(\mathbf{x}^{(i)})$ y **varianza σ^2** .



INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Formalmente quedaría así:

$$p(y^{(i)} / x^{(i)}; w) = \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - h_w(x^{(i)}))^2}{2\sigma^2}\right)}$$

$$p(y^{(i)} / x^{(i)}; w) = \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)}$$

INTERPRETACIÓN PROBABILÍSTICA

SUPOSICIÓN DE ERRORES



Definimos la **probabilidad condicionada** para **todos** los **datos** de **entrenamiento** (suponiendo que las **muestras** se **recolectaron** de manera **independiente**):

$$p(\vec{y}/X; \mathbf{w}) = \prod_{i=1}^m p(y^{(i)} / x^{(i)}; \mathbf{w}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - \mathbf{w}^T x^{(i)})^2}{2\sigma^2}\right)}$$

INTERPRETACIÓN PROBABILÍSTICA

FUNCIÓN DE VEROSIMILITUD



A esta ecuación se le llama función de **verosimilitud**, porque describe la **probabilidad** de que nuestras **creencias** sobre la **realidad** de las **observaciones** (datos) sean **ciertas**. Es decir que la **hipótesis** que planteamos sea la **correcta**.

$$p(\vec{y}/X; \vec{w}) = \prod_{i=1}^m p(y^{(i)} / x^{(i)}; \vec{w}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)}$$

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

INTERPRETACIÓN PROBABILÍSTICA

FUNCIÓN DE VEROSIMILITUD



Por lo tanto, nosotros queremos **maximizar** la **probabilidad** de que **nuestras creencias** sobre como se **distribuyen** los **datos** (de manera **normal**).

Es decir queremos **maximizar** la **función** de **verosimilitud**. Desde una perspectiva **frecuentista** (el valor de \vec{w} no es **aleatorio**), queremos **encontrar** el **valor** de \vec{w} que **maximice** la **probabilidad** de que las **observaciones** que se hacen de la **realidad** sean **ciertas**.

$$\arg \max_{\vec{w}} L(\vec{w}; X, \vec{y}) = \arg \max_{\vec{w}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)}$$

INTERPRETACIÓN PROBABILÍSTICA

VEROSIMILITUD LOGARÍTMICA



Es mucho más fácil **maximizar** una **sumatoria** que una **multiplicación**

$$\arg \max_{\vec{w}} \log(L(\vec{w})) = \arg \max_{\vec{w}} \log\left(\prod_{i=1}^m p(y^{(i)} / x^{(i)}; \vec{w})\right)$$

$$\arg \max_{\vec{w}} \log(L(\vec{w})) = \arg \max_{\vec{w}} \sum_{i=1}^m \log(p(y^{(i)} / x^{(i)}; \vec{w}))$$

INTERPRETACIÓN PROBABILÍSTICA

PÉRDIDA LOGARÍTMICA



Se **convierte** el **problema de maximización** en uno de **minimización** al **escalar** la función por un **signo menos**.

A esta **función** $-\log(L(\vec{w}))$ se le llama **pérdida logarítmica**.

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} - \sum_{i=1}^m \log(p(y^{(i)} / x^{(i)}; \vec{w}))$$

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} - \sum_{i=1}^m \log\left(\frac{1}{\sqrt{2\pi}\sigma} e^{\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)}\right)$$

INTERPRETACIÓN PROBABILÍSTICA

PÉRDIDA LOGARÍTMICA



Desarrollamos la ecuación:

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} - \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \log \left(e^{\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right)} \right)$$

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} - m \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \sum_{i=1}^m -\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}$$

INTERPRETACIÓN PROBABILÍSTICA

PÉRDIDA LOGARÍTMICA



Desarrollamos la ecuación:

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} c + \sum_{i=1}^m \frac{(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2}{2\sigma^2}$$

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

INTERPRETACIÓN PROBABILÍSTICA

ERROR CUADRÁTICO MEDIO



Minimizar la pérdida logarítmica es lo mismo que minimizar el error cuadrático medio $MSE = J(\vec{w})$:

$$\arg \min_{\vec{w}} MSE = \arg \min_{\vec{w}} \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - w^T x^{(i)})^2$$

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - w^T x^{(i)})^2$$



AI

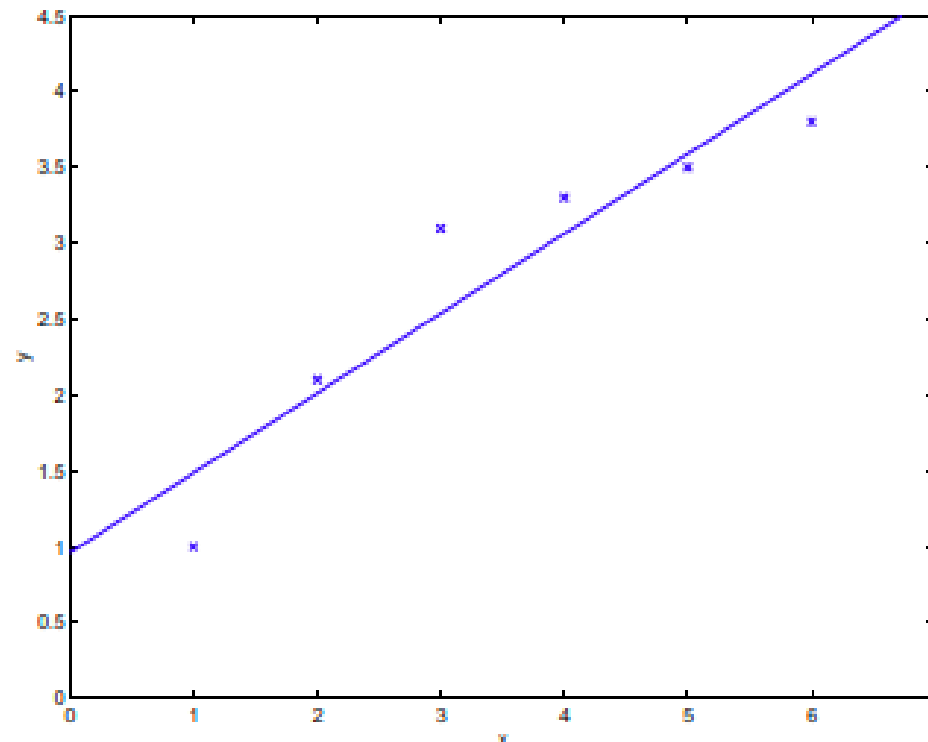
FUNCIONES DE BASE

F U N C I O N E S D E B A S E

P R O B L E M A D E L I N E A L I D A D

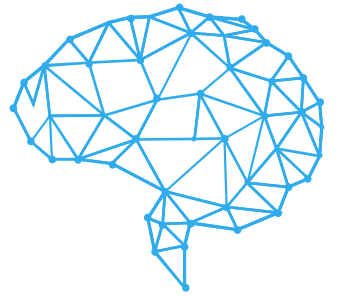


El **mundo real** tiene comportamientos **no lineales**, por lo que los **modelos de regresión lineal** vistos hasta ahora tienen sus **limitantes**.

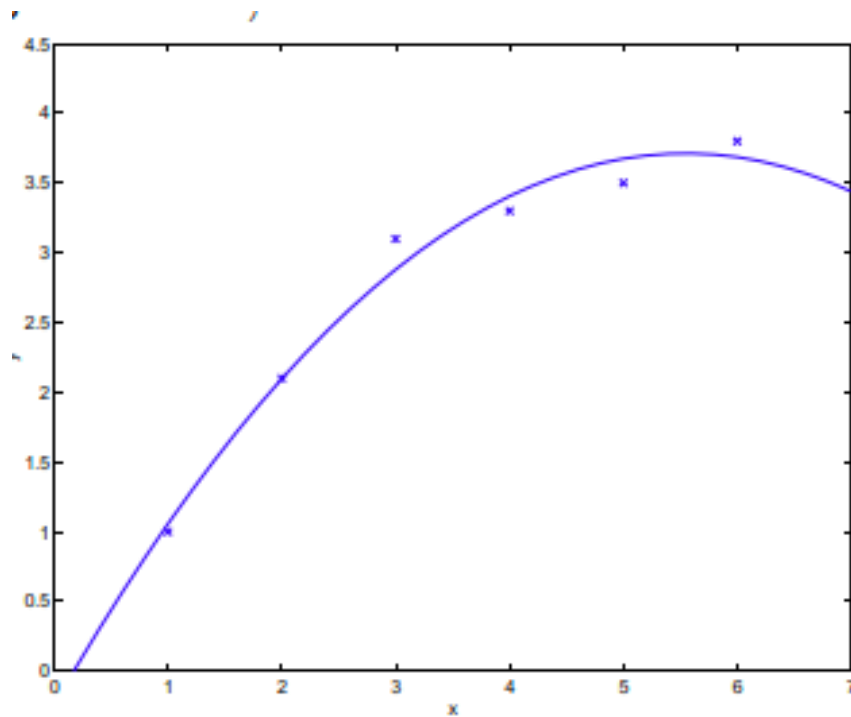


FUNCIONES DE BASE

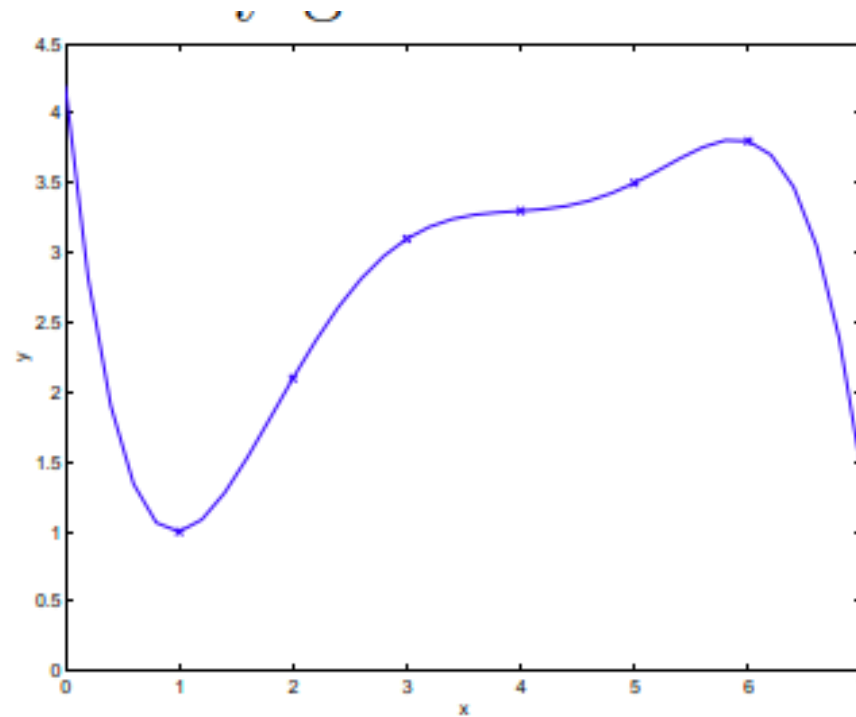
PROBLEMA DE LINEALIDAD



Podemos **establecer** nuestra **hipótesis** como un **modelo polinómico**.



$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$y = \sum_{j=0}^5 \theta_j x^j$$

F U N C I O N E S D E B A S E

INTRODUCCIÓN A FUNCIONES DE BASE



Si llevamos este tipo de pensamiento más allá, podemos **construir** modelos como **combinaciones lineales** de **funciones no lineales** $\phi_j(x)$, llamadas **funciones de base**. Donde $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^k$

$$h_w(x^{(i)}) = \sum_{j=1}^k w_j \phi_j(x^{(i)})$$

$$h_w(x^{(i)}) = w^T \phi(x^{(i)})$$

NOTA

$$w \in \mathbb{R}^k$$

$$x^{(i)} \in \mathbb{R}^n$$

Así logramos **construir** un **modelo no lineal**, pero que sigue estando **parametrizado** por **pesos lineales** w .

F U N C I O N E S D E B A S E

INTRODUCCIÓN A FUNCIONES DE BASE



De manera más detallada, se tiene que:

$$\phi(x^{(i)}) = \begin{bmatrix} \phi_1(x^{(i)}) \\ \phi_j(x^{(i)}) \\ \vdots \\ \phi_k(x^{(i)}) \end{bmatrix}$$

FUNCIONES DE BASE

REGRESIÓN LINEAL CLÁSICA

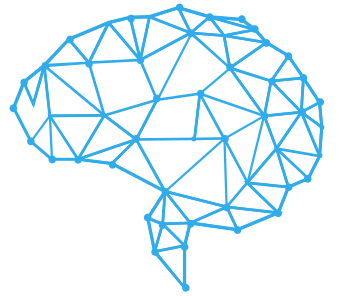


En el caso de **regresión lineal clásica** se tiene **para un solo** dato de entrenamiento $\phi(x) \in \mathbb{R}^k$ en este caso $k = n$:

$$\phi_j(x^{(i)}) = x^{(i)} \quad x^{(i)} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_k \end{bmatrix}$$

F U N C I O N E S D E B A S E

REGRESIÓN LINEAL CLÁSICA



En el caso de **regresión lineal clásica** se tiene **para un solo** dato de entrenamiento $\phi(x) \in \mathbb{R}^k$ en este caso $k = n$:

$$\phi(x) = \begin{bmatrix} \phi_0(x^{(1)}) \\ \phi_1(x^{(i)}) \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ x^{(i)} \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$h_w(x^{(i)}) = w^T \phi(x^{(i)}) = w_0 + w_1(x_1) + \cdots + w_k(x_n)$$

F U N C I O N E S D E B A S E

REGRESIÓN LINEAL CLÁSICA



En el caso de **regresión lineal clásica** se tiene **para m datos** de entrenamiento:

$$\phi(x) = \begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) \\ \vdots & \vdots \\ \phi_0(x^{(m)}) & \phi_1(x^{(m)}) \end{bmatrix} = \begin{bmatrix} \mathbf{1} & x^{(1)T} \\ \vdots & \vdots \\ \mathbf{1} & x^{(m)T} \end{bmatrix}$$

$$\mathbf{w} = [\mathbf{w}_0 \quad \dots \quad \mathbf{w}_k]$$

FUNCIONES DE BASE

REGRESIÓN LINEAL CLÁSICA



En el caso de **regresión lineal clásica** se tiene **para m datos** de entrenamiento:

$$h_w(x) = w^T \phi(x) = \begin{bmatrix} w_0 + \dots + w_k x_n^{(1)} \\ \vdots \\ w_0 + \dots + w_k x_n^{(m)} \end{bmatrix}$$

$$h_w(x) = \begin{bmatrix} h_w(x^{(1)}) \\ \vdots \\ h_w(x^{(m)}) \end{bmatrix}$$

F U N C I O N E S D E B A S E

M A T R I Z D E D I S E Ñ O



Por lo tanto, la **matriz** de **diseño** para cualquier $\phi(x)$ estaría dada por:

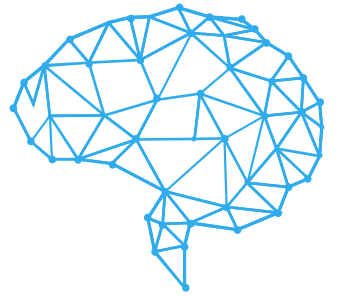
$$\phi(x) = \begin{bmatrix} \phi_0(x^{(1)}) & \cdots & \phi_{k-1}(x^{(1)}) \\ \vdots & \ddots & \vdots \\ \phi_0(x^{(m)}) & \cdots & \phi_{k-1}(x^{(m)}) \end{bmatrix}$$

$$x^{(i)} \in \mathbb{R}^n$$

Donde cada fila de la matriz $\phi(x)$ está dada por $\phi_j = \phi(x^{(i)})^T$

F U N C I O N E S D E B A S E

F U N C I O N E S D E B A S E L I N E A L E S



Funciones lineales clásicas:

$$\phi_j(x^{(i)}) = x^{(i)}; x^{(i)} \in \mathbb{R}^{n+1}$$

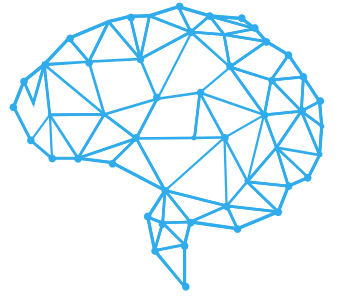
$$\mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix}; k = n$$

$$\phi(x) = \begin{bmatrix} \mathbf{1} & x^{(1)T} \\ \vdots & \vdots \\ \mathbf{1} & x^{(m)T} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \dots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ \mathbf{1} & \dots & x_n^{(m)} \end{bmatrix}$$

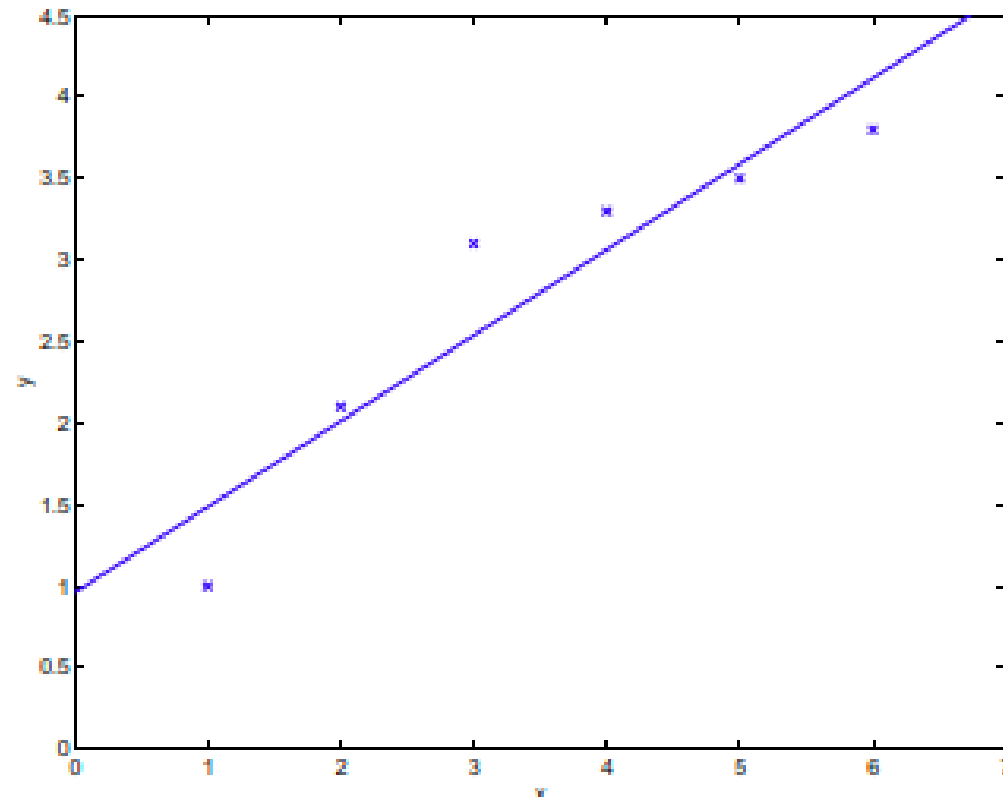
$$h_w(x) = \mathbf{w}^T \phi(x) = \begin{bmatrix} w_0 + w_1 x_1^{(1)} + \dots + w_k x_n^{(1)} \\ \vdots \\ w_0 + w_1 x_1^{(m)} + \dots + w_k x_n^{(m)} \end{bmatrix}$$

F U N C I O N E S D E B A S E

F U N C I O N E S D E B A S E L I N E A L E S

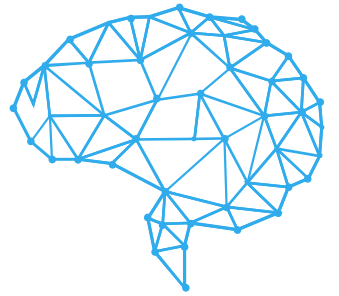


Funciones lineales clásicas:



F U N C I O N E S D E B A S E

F U N C I O N E S D E B A S E P O L I N Ó M I C A S



Funciones de base polinómicas:

$$\phi_j(x) = (x^{(i)})^j; x^{(i)} \in \mathbb{R}^n$$

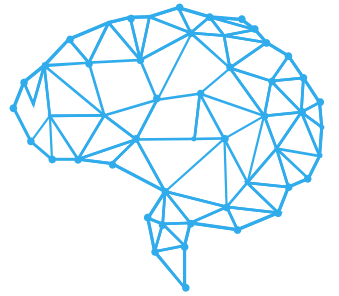
$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{k*n+1} \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} \mathbf{1} & x^{(1)T} & \dots & (x^{(1)T})^{(k-1)} \\ \vdots & \vdots & & \vdots \\ \mathbf{1} & x^{(m)T} & \dots & (x^{(m)T})^{(k-1)} \end{bmatrix}$$

$$h_w(x) = \mathbf{w}^T \phi(x) = \begin{bmatrix} \mathbf{w}_0 + \mathbf{w}_1 x_1^{(1)} + \dots + \mathbf{w}_{k*n} (x_n^{(1)})^{k-1} \\ \vdots \\ \mathbf{w}_0 + \mathbf{w}_1 x_1^{(m)} + \dots + \mathbf{w}_{k*n} (x_n^{(m)})^{k-1} \end{bmatrix}$$

F U N C I O N E S D E B A S E

F U N C I O N E S D E B A S E P O L I N Ó M I C A S



Funciones de base polinómicas: ejemplo concreto polinomio grado 2 y dos datos de entrenamiento $m = 2$

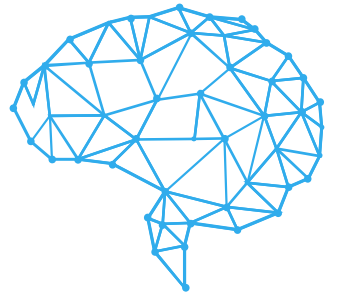
$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & (x_1^{(1)})^2 & (x_2^{(1)})^2 \\ 1 & x_1^{(2)} & x_2^{(2)} & (x_1^{(2)})^2 & (x_2^{(2)})^2 \end{bmatrix}$$

F U N C I O N E S D E B A S E

FUNCIONES DE BASE POLINÓMICAS



Funciones de base polinómicas: ejemplo concreto **polinomio grado 2** y **dos** **datos** de entrenamiento $m = 2$ con **2** **características**.

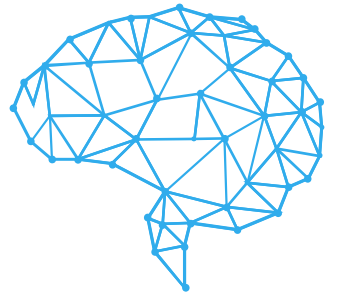
$$h_w(x) = w^T \phi(x)$$

$$h_w(x) = \begin{bmatrix} \mathbf{1} & x_1^{(1)} & x_2^{(1)} & (x_1^{(1)})^2 & (x_2^{(1)})^2 \\ \mathbf{1} & x_1^{(2)} & x_2^{(2)} & (x_1^{(2)})^2 & (x_2^{(2)})^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

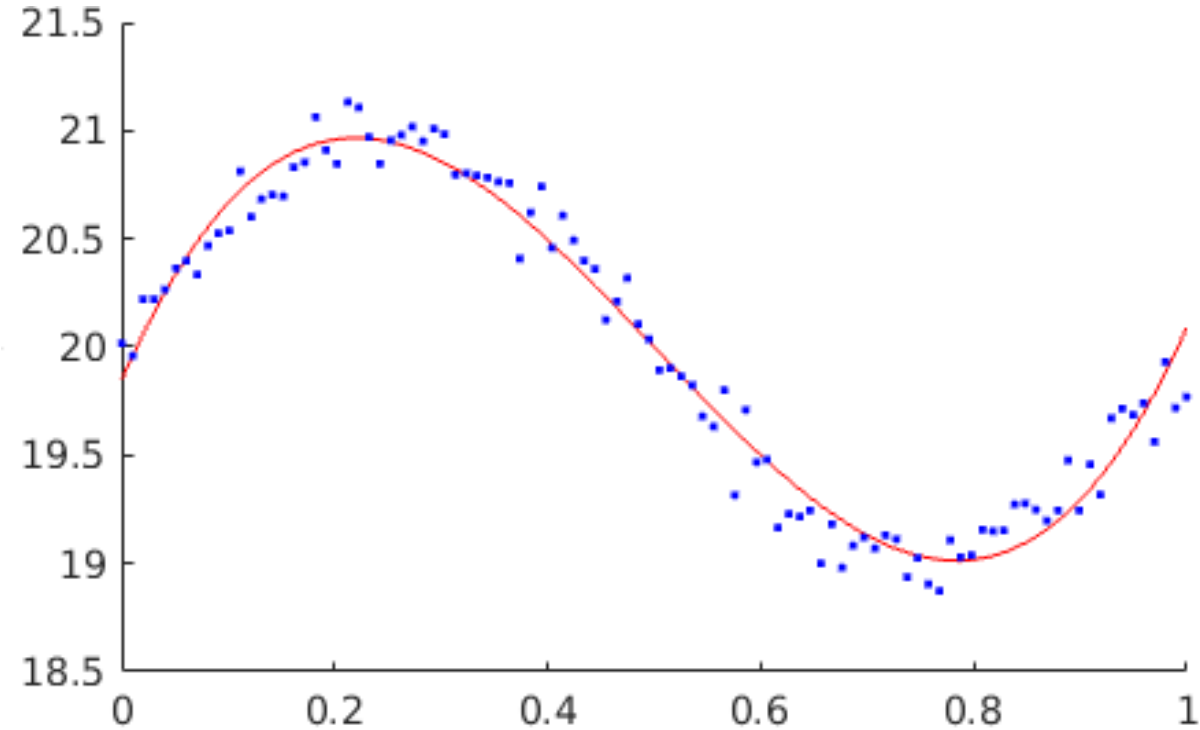
$$h_w(x) = \begin{bmatrix} w_0 + w_1 x_1^{(1)} + w_2 x_2^{(1)} + w_3 (x_1^{(1)})^2 + w_4 (x_2^{(1)})^2 \\ w_0 + w_1 x_1^{(2)} + w_2 x_2^{(2)} + w_3 (x_1^{(2)})^2 + w_4 (x_2^{(2)})^2 \end{bmatrix}$$

F U N C I O N E S D E B A S E

F U N C I O N E S D E B A S E P O L I N Ó M I C A S

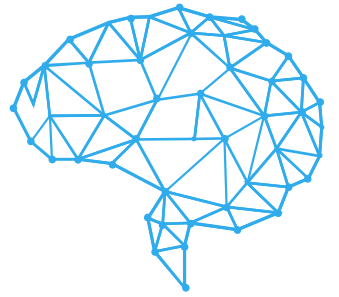


Funciones de base polinómicas:



F U N C I O N E S D E B A S E

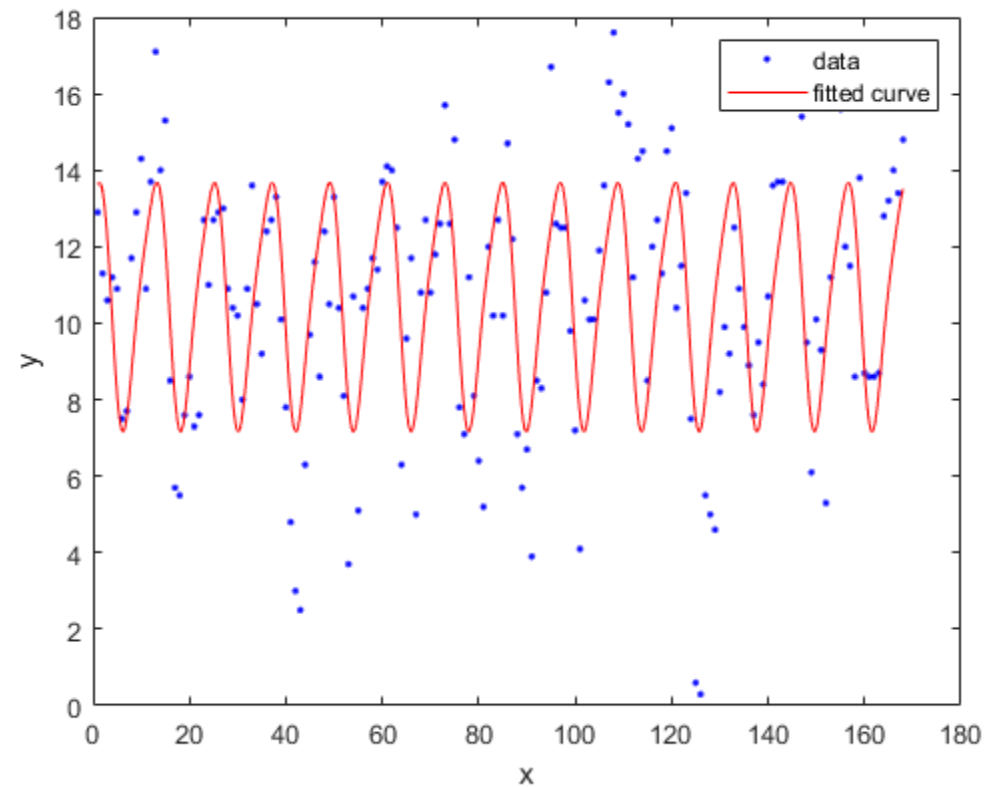
S E R I E S D E F O U R I E R



Series de Fourier

$$\phi_0(x) = 1$$

$$\phi_j(x) = \cos(\varpi_j x^{(i)} + \psi_j) ; j > 0$$



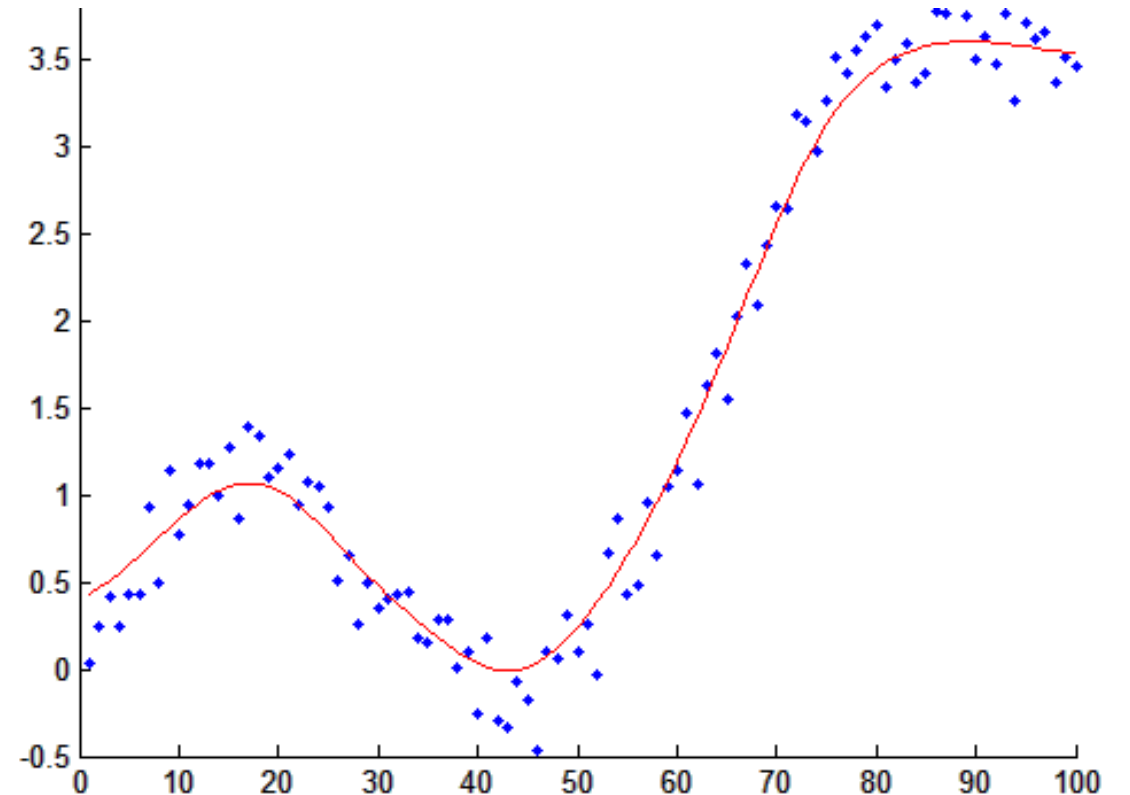
F U N C I O N E S D E B A S E

B A S E R A D I A L



Funciones de base radial

$$\phi_j(x) = e^{-\frac{1}{2l} \sum_{r=1}^n (x_r^{(i)} - \mu_{j,r})^2}$$



F U N C I O N E S D E B A S E

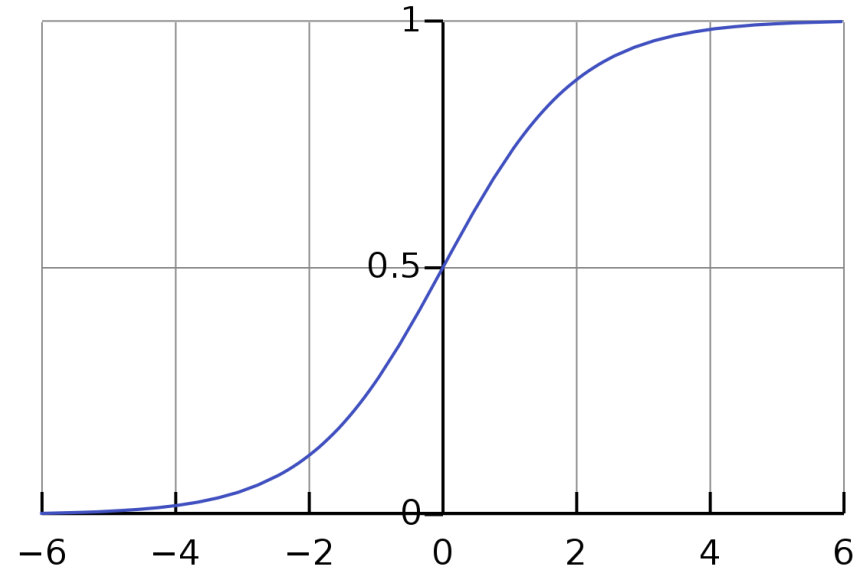
S I G M O I D A L E S



Funciones sigmoidales

$$\phi_j(x^{(i)}) = \sigma\left(\frac{x^{(i)} - \mu_j}{s}\right)$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$



F U N C I O N E S D E B A S E

M Á X I M A V E R O S I M I L T U D



Se propone el **mismo** modelo de **estimación** más un **error**:

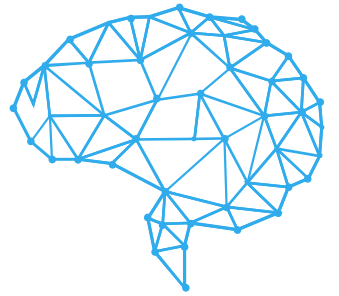
$$\mathbf{y}^{(i)} = \mathbf{h}_w(\mathbf{x}^{(i)}) + \boldsymbol{\varepsilon}^{(i)}$$

donde:

$$\mathbf{h}_w(\mathbf{x}^{(i)}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

F U N C I O N E S D E B A S E

M Á X I M A V E R O S I M I L T U D

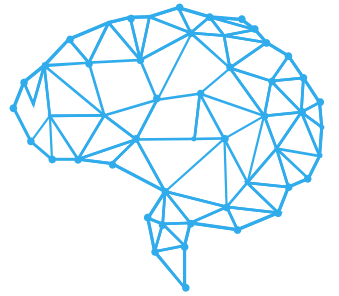


Suponiendo otra vez que el **error** se **distribuye** de manera **normal** con **media 0** y **varianza β^{-1}** . Se define para **m** datos de **entrenamiento**:

$$p(\vec{y}/X, w, \beta) = \prod_{i=1}^m p(y^{(i)}/x^{(i)}; w) = \prod_{i=1}^m \sqrt{\frac{\beta}{2\pi}} e^{\left(-\frac{\beta}{2}(y^{(i)} - w^T \phi(x^{(i)}))^2\right)}$$

F U N C I O N E S D E B A S E

M Á X I M A V E R O S I M I L T U D



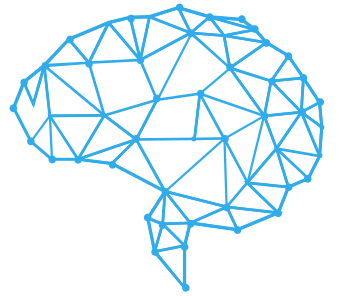
Como la **función** de **verosimilitud** está **parametrizada** en w y β la verosimilitud logarítmica se escribe así:

$$\log p(\vec{y}/w, \beta) = \log \prod_{i=1}^m \sqrt{\frac{\beta}{2\pi}} e^{\left(-\frac{\beta}{2}(y^{(i)} - w^T \phi(x^{(i)}))^2\right)}$$

$$\log p(\vec{y}/w, \beta) = \sum_{i=1}^m \log \sqrt{\frac{\beta}{2\pi}} e^{\left(-\frac{\beta}{2}(y^{(i)} - w^T \phi(x^{(i)}))^2\right)}$$

F U N C I O N E S D E B A S E

M Á X I M A V E R O S I M I L T U D



Desarrollando:

$$\log p(\vec{y}/w, \beta) = \sum_{i=1}^m \log \sqrt{\frac{\beta}{2\pi}} + \log e^{\left(-\frac{\beta}{2}(y^{(i)} - w^T \phi(x^{(i)}))^2\right)}$$

$$\log p(\vec{y}/w, \beta) = m \log \sqrt{\frac{\beta}{2\pi}} - \sum_{i=1}^m \frac{\beta}{2} (y^{(i)} - w^T \phi(x^{(i)}))^2$$

F U N C I O N E S D E B A S E

M Á X I M A V E R O S I M I L T U D



Desarrollando:

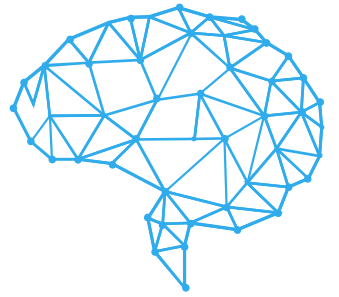
$$\log p(\vec{y}/w, \beta) = \frac{m}{2} \log \beta - \frac{m}{2} \log 2\pi - \sum_{i=1}^m \frac{\beta}{2} \left(y^{(i)} - w^T \phi(x^{(i)}) \right)^2$$

$$\log p(\vec{y}/w, \beta) = \frac{m}{2} \log \beta - \frac{m}{2} \log 2\pi - \beta E_D(w)$$

$$E_D(w) = \frac{1}{2} \left(y^{(i)} - w^T \phi(x^{(i)}) \right)^2$$

FUNCIONES DE BASE

MÁXIMA VEROSIMILITUD



Desarrollando:

$$-\log p(\vec{y}/w, \beta) = -\frac{m}{2} \log \beta + \frac{m}{2} \log 2\pi + \beta E_D(w)$$

$$\arg \min_{\vec{w}} -\log(L(\vec{w})) = \arg \min_{\vec{w}} \frac{\beta}{2} \sum_{i=1}^m \left(y^{(i)} - w^T \phi(x^{(i)}) \right)^2$$

F U N C I O N E S D E B A S E

E C U A C I O N E S N O R M A L E S



Calculando el **gradiente**, **igualando** a **cero** y resolviendo para el **vector w** :

$$w = (\phi^T \phi)^{-1} \phi^T \vec{y}$$