



MACHINE LEARNING

LINEAR REGRESSION

AGENDA

01 Supervised Learning

Training data, hypotheses, an example.

02 LMS Algorithm

Linear model, cost function, gradient descent.

03 Normal Equations

Matrix linear model, deriving the normal equations

04 Probabilistic Interpretation

Assumptions for errors, likelihood function

05 Base Function

Derivation, functions, maximum likelihood





SUPERVISED LEARNING

DEFINITIONS

NOMENCLATURE

SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



We start with a data set that represents first-floor living space and sales prices for 1,460 homes in Ames, Iowa.

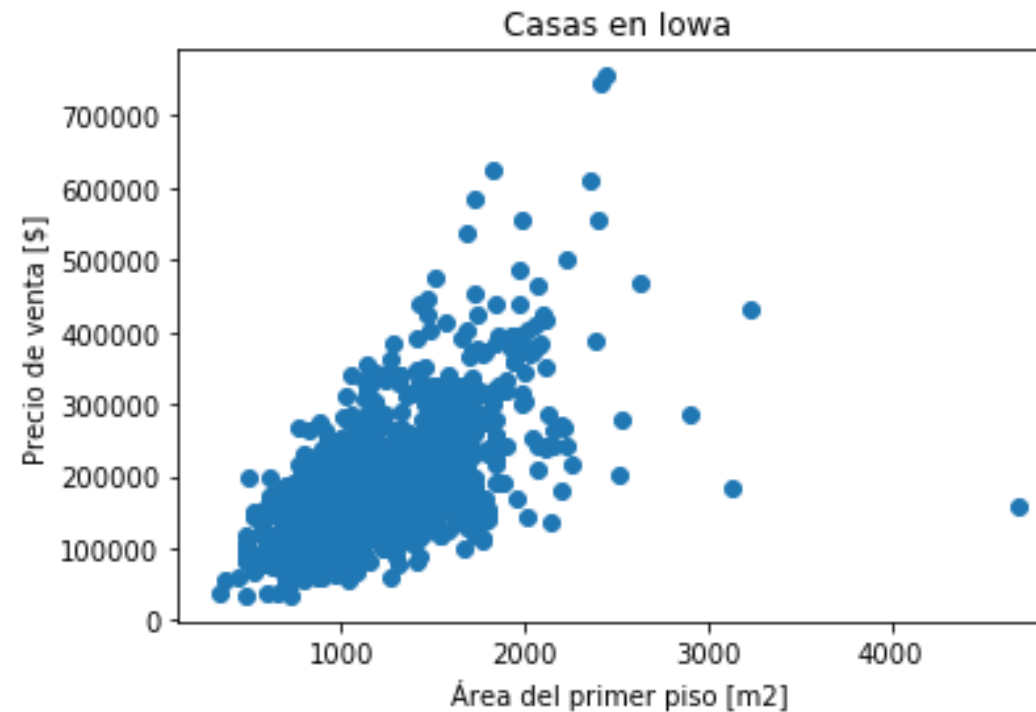
1 st floor square feet	Sale Price
856	208500
1262	181500
920	223500
961	140000
1145	250000
⋮	⋮

SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



We graph the 1,460 elements using Python:



SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



The question that naturally arises would be:

**HOW DO WE PREDICT THE PRICES FOR OTHER
HOMES IN AMES, IOWA?**

SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



To answer the question, we must first define the nomenclature that will be used:

Variable / Symbol	Description	Example
$x^{(i)}$	Input variable or characteristics	The house's first floor area
$y^{(i)}$	Output, target, or response variable	The house's sale price
$(x^{(i)}, y^{(i)})$	Training example	(area, price)
m	Number of training examples	1460 houses with their first floor area and their price.
$\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$	Training data set	NA
\mathcal{X}	Input value space	NA
\mathcal{Y}	Output value space	

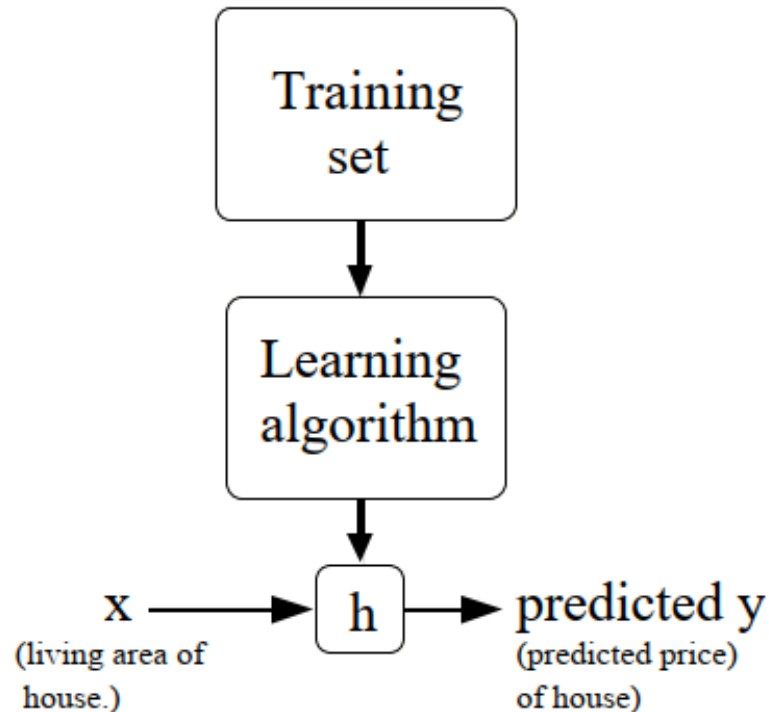
SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



The **main goal** of any **supervised algorithm** is to **learn** a function $h: \mathcal{X} \rightarrow \mathcal{Y}$, such that $h(x)$ can **predict** the corresponding value of y .

Where h is defined as the **hypothesis**.



SUPERVISED LEARNING

DEFINITIONS AND NOMENCLATURE



There are **2 types** of **supervised learning problems**:

Type of problem	Description
Regression problem	y takes continuous values.
Classification problem	y takes discrete values



LMS ALGORITHM

LINEAR MODEL

COST FUNCTION

GRADIENT DESCENT

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

L I N E A R M O D E L



Now let's look at the same model but with another variable: the living space on the second floor.

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
⋮	⋮	⋮

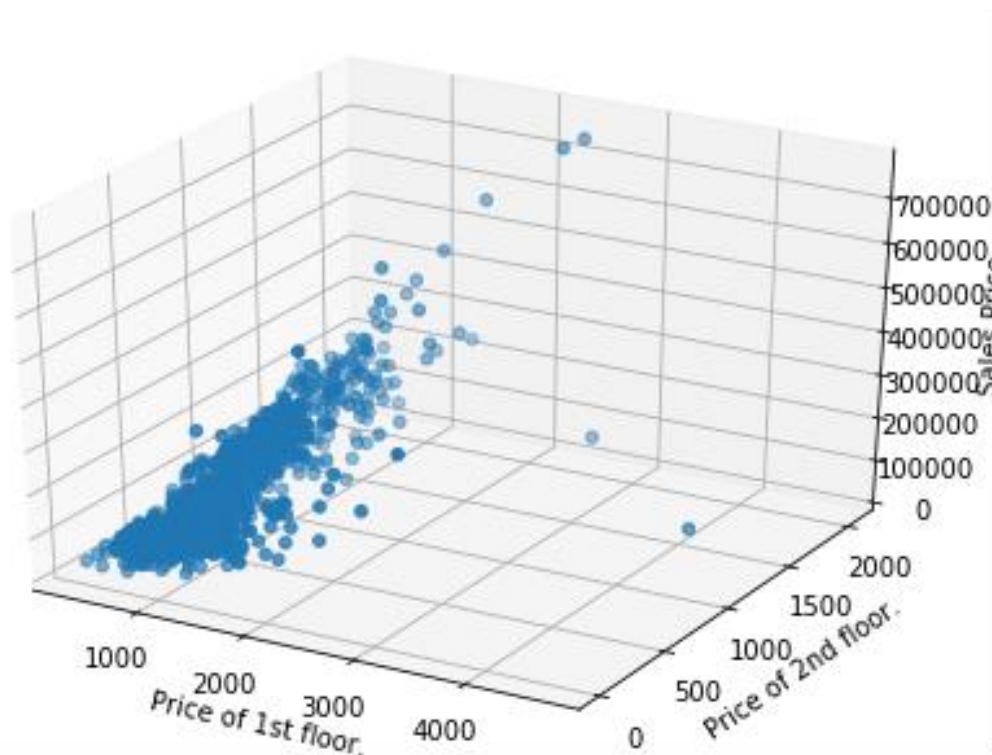
In this case each $x^{(i)} \in \mathbb{R}^2$, so it is defined as a vector $x^{(i)} = [x^{(i)}_1 \ x^{(i)}_2]$.

LS ALGORITHM

LINEAR MODEL

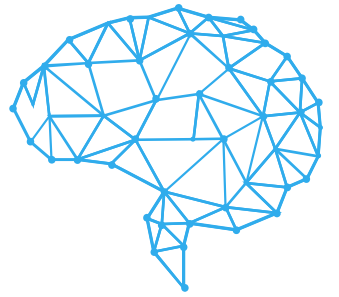


We graph the 1,460 elements again using Python:



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

L I N E A R M O D E L



It is **hypothesized** that the data is distributed in a **linear fashion**, therefore:

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

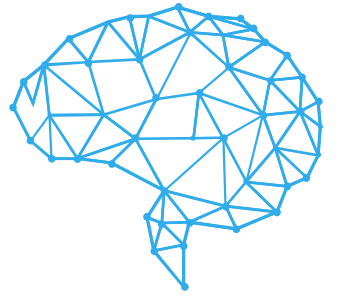
In this case the **variables** w_i are defined as the **parameters** or **weights** that **parameterize** the **space** of all **linear functions** that map from \mathcal{X} to \mathcal{Y} .

The expression is **simplified** into **vector form**, where n represents the **number of input variables** (omitting x_0),:

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

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

C O S T F U N C T I O N



Then, the question would be:

**WHAT IS THE BEST COMBINATION OF PARAMETERS w
THAT RESULTS IN THE BEST HYPOTHESIS h ?**

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

C O S T F U N C T I O N



To **answer** the **question**, we need to **measure** the **error** produced by the **estimates** produced by h .

We know that we can **measure** the **error** of a single **estimate** i as a **difference** **between** the estimate $h_w(x^{(i)})$ and the response variable $y^{(i)}$.

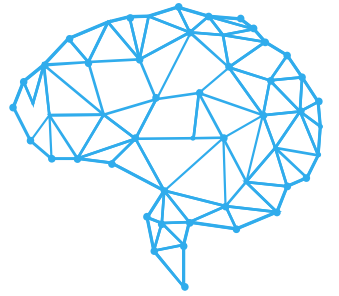
$$e = h_w(x^{(i)}) - y^{(i)}$$

To ensure positive amounts of error, the answer is squared:

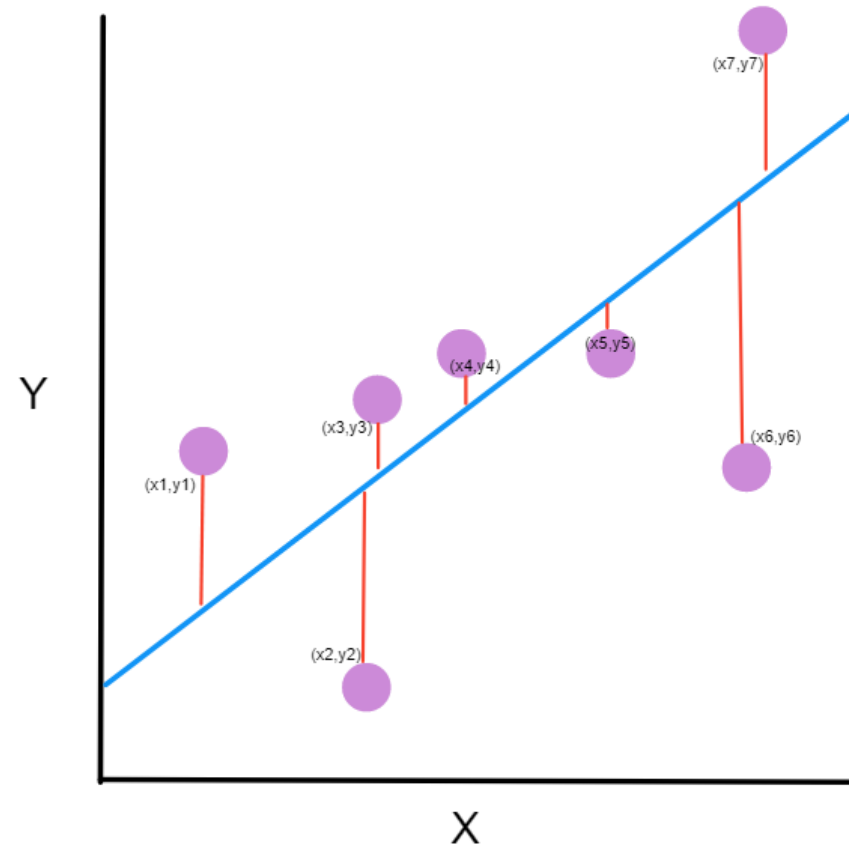
$$e^2 = (h_w(x^{(i)}) - y^{(i)})^2$$

LEAST SQUARES ALGORITHM

COST FUNCTION



The **error** that we are **calculating graphically** is **interpreted** with an **example**:

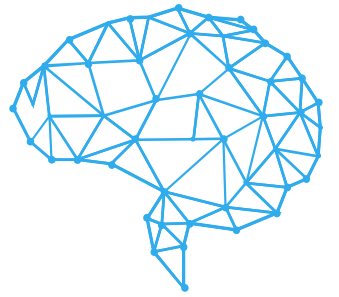


The linear model
represented by h

Response variables $y^{(i)}$

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

C O S T F U N C T I O N



So far, we have calculated the **squared error** of a **single training data**. The **mean square error** *MSE* is calculated for **all training data**.

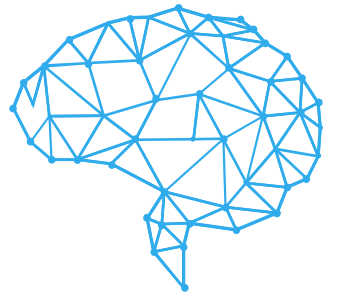
$$MSE = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

$$MSE = J(w)$$

Where $J(w)$ is the **cost function**.

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

G R A D I E N T D E S C E N T



Therefore, the **best combination** of weights w is the one that **minimizes** the **cost function** $J(w)$, which **measures** our **mean square error**.

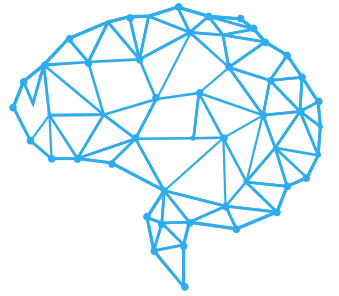
To find the **best combination**, let's design a **search algorithm** that **starts** with a **random initial value** of w , and **updates** the **values** of w until it **converges** to a **minimum value** of $J(w)$. The **constant** α is defined as the **learning rate**.

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

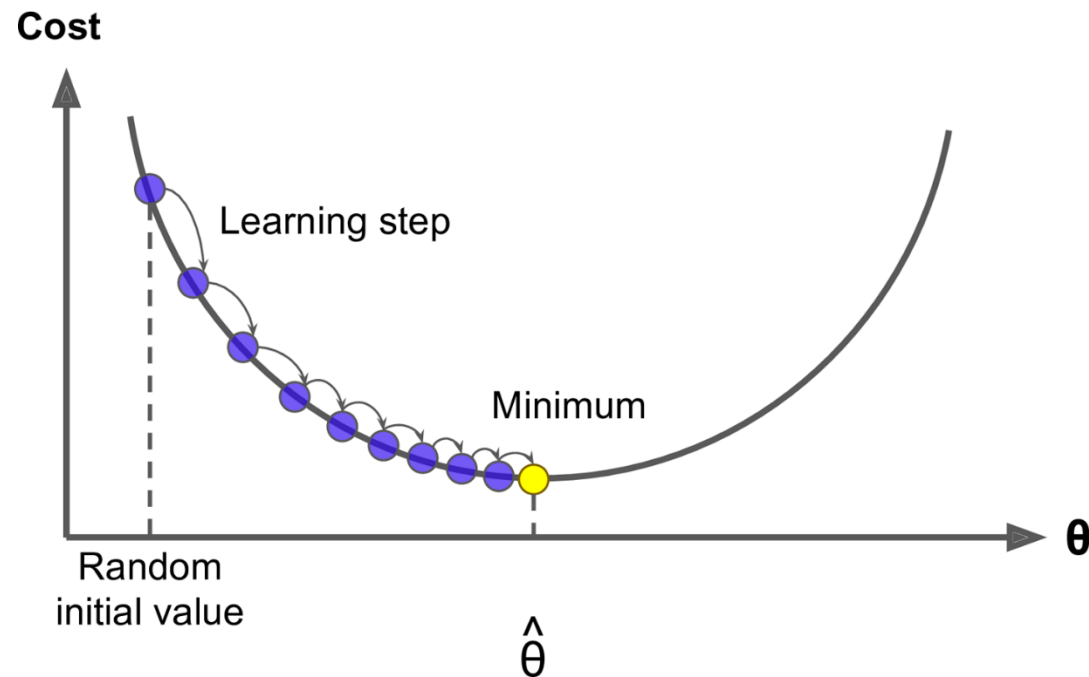
NOTE: the above equation **only updates** the **value** of a **single weight** w_j of the n weights that **parameterize** the **linear model**. In reality, **all weights** w_j are **updated simultaneously**.

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

G R A D I E N T D E S C E N T

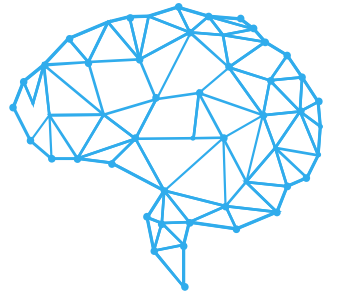


In the **gradient descent optimization algorithm**, we **update the weights** in the **direction** with the **greatest decrement of $J(w)$** .



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

G R A D I E N T D E S C E N T



The **derivative** of the **cost function** $J(\mathbf{w})$ with respect to a **specific weight** w_j is calculated.

Derive the result:

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

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

G R A D I E N T D E S C E N T



Therefore, the **gradient descent weights update** for a **single training data** would look like this:

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

This **equation** is called the **LMS update rule** (“**Least Mean Squares**”) or also the **Widrow-Hoff learning rule**.

This **method** is also called **batch gradient descent**, where **all training data** is **analyzed simultaneously**.

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

G R A D I E N T D E S C E N T



For m training data, and defining $x^{(i)}$ and w as **vectors**, the learning rule would look like this:

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

This **equation** is called the **LMS** update rule (“Least Mean Squares”) or also the **Widrow-Hoff learning rule**.

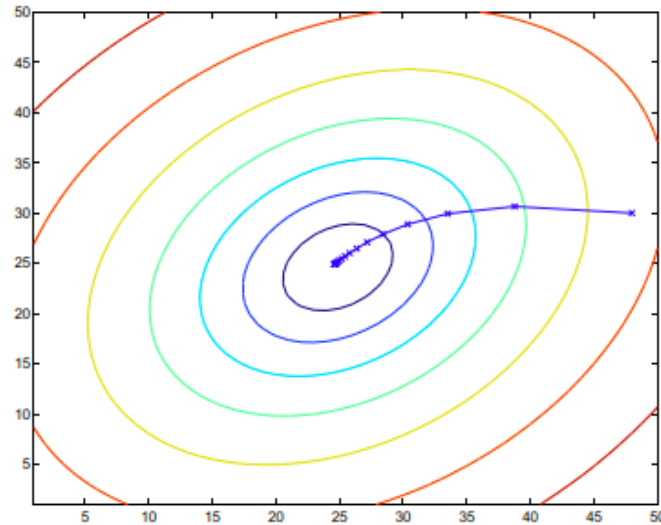
L M S A L G O R I T H M

G R A D I E N T D E S C E N T



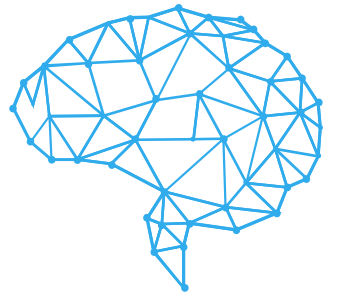
In general, the **gradient descent method** can suffer from **several local minima**. In this case, for **linear regression models** there is only a **single global minimum**.

Consequently, the **algorithm will always converge** assuming that the **learning rate is not too high**.



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

G R A D I E N T D E S C E N T



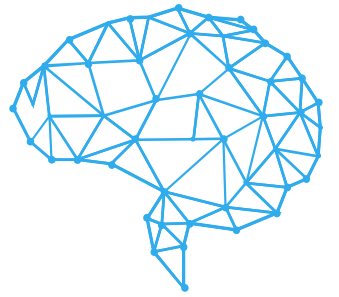
When the **method** only **observes one training data** at a time, and **updates the weights with a single data**, it is called **stochastic or incremental gradient descent**.

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

This **variant** of the algorithm is **used** when it is **very expensive to evaluate** the **update** for **very large data sets** (when m is **very large**), but it has the minimal **divergence** problem.

LEAST SQUARES ALGORITHM

GRADIENT DESCENT

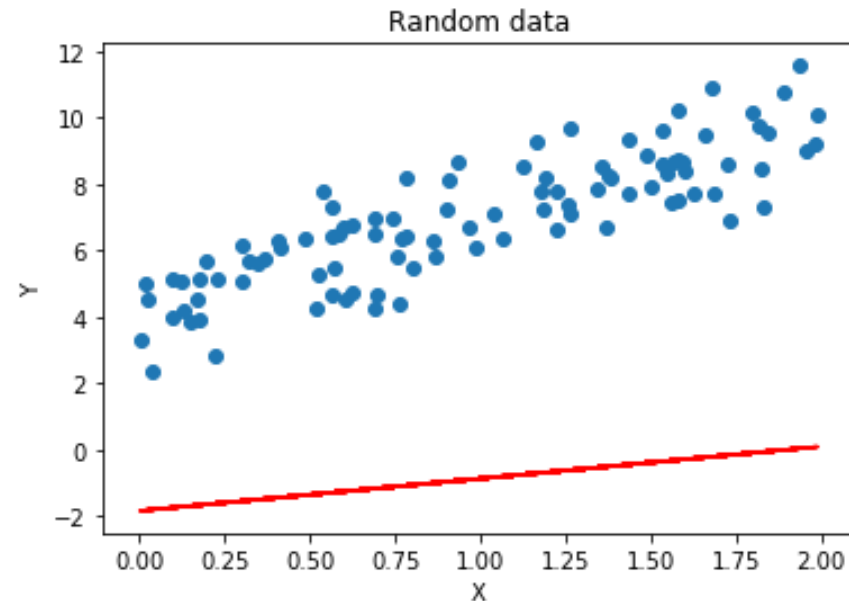


Example:

Random training data was generated with a certain degree of error e :

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

and both weights were randomly initialized: w_0 and w_1



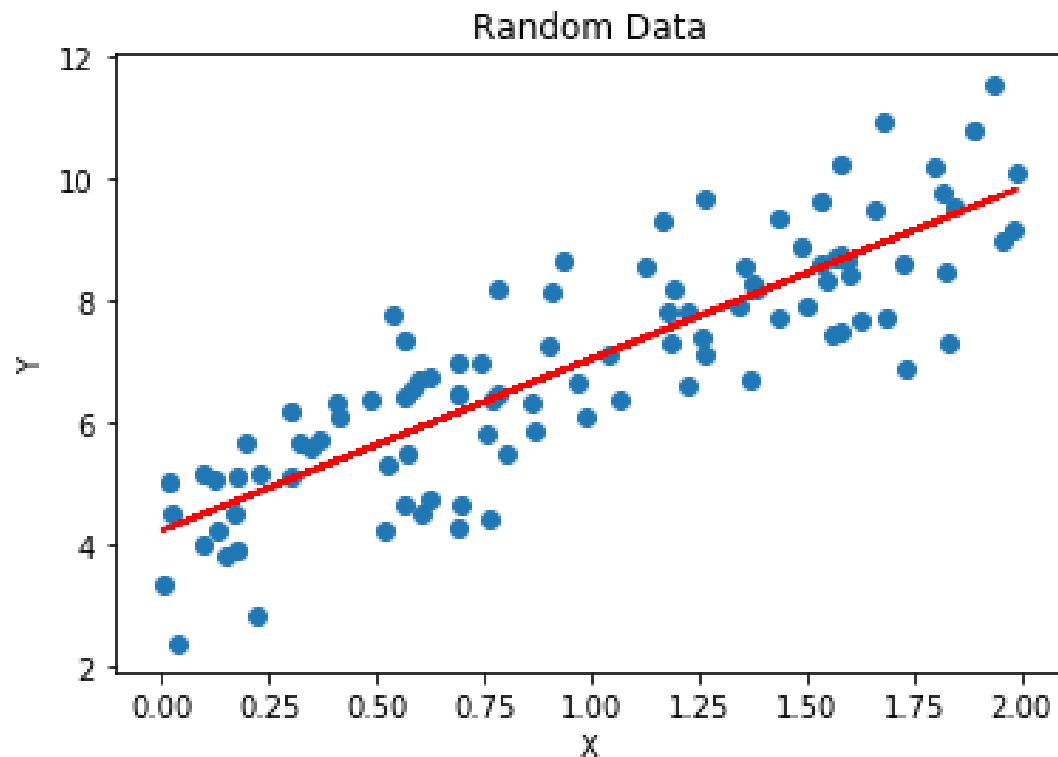
L M S A L G O R I T H M

G R A D I E N T D E S C E N T



Example:

After 1000 iterations the following model was obtained:



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

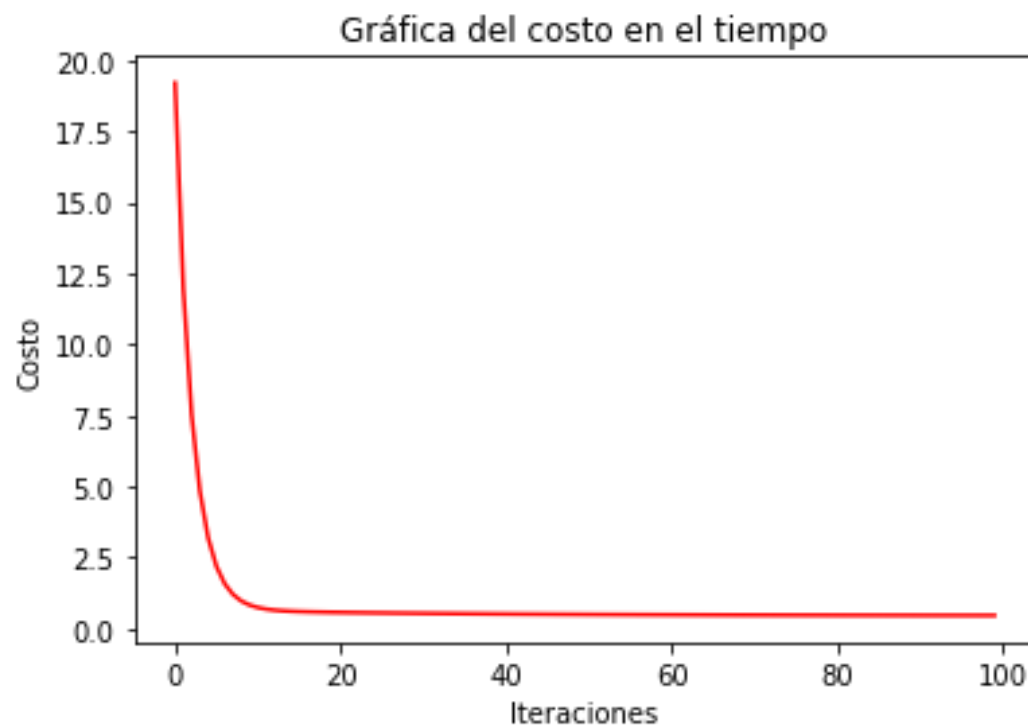
L M S A L G O R I T H M

G R A D I E N T D E S C E N T



Example:

Error based on 100 iterations.





NORMAL EQUATIONS

LINEAR MATRIX MODEL
DERIVATION

NORMAL EQUATIONS

MATRIX NOTATION



Now we want to **calculate** the **minimum** of the **cost function analytically**. For this, **matrix notation** is introduced, which **allows** the **derivatives** of the **cost function** to be **expressed** in an **elegant way** without getting into so much **verbiage**.

The **m training data** $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ is represented as a matrix $X \in \mathbb{R}^{m \times n}$, the set of **outputs** as a **vector** $\vec{y} \in \mathbb{R}^m$ and the **n weights** w_j as a vector $\vec{w} \in \mathbb{R}^n$.

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad X = \begin{bmatrix} -(x^{(1)})^T & - \\ -(x^{(2)})^T & - \\ \vdots & \\ -(x^{(m)})^T & - \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

NORMAL EQUATIONS

MATRIX NOTATION



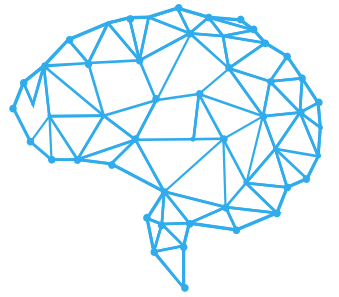
The **absolute error function** is represented in **matrix notation** where $\mathbf{h}_w(\mathbf{x}^{(i)}) = (\mathbf{x}^{(i)})^T \vec{w}$:

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

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

NORMAL EQUATIONS

MATRIX NOTATION



The **mean square error term** (cost function):

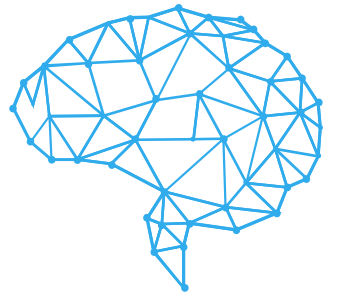
$$J(\mathbf{w}) = \frac{1}{2m} \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$$

can be **represented** in **vector** form using the **property** $\mathbf{z}^T \mathbf{z} = \sum_i z_i^2$

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

NORMAL EQUATIONS

DERIVATION OF EQUATIONS



We have to **obtain** the **derivatives** of $J(w)$ with respect to the vector w .

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

The **gradient** **calculates** all the **derivatives** with **respect** to **all** the weights w_j **simultaneously**.

Also, because we **represent** the **training dataset** in **matrix form**, we can **calculate** the **gradient** for **all training data**.

NORMAL EQUATIONS

DERIVATION OF EQUATIONS



Therefore, we obtain:

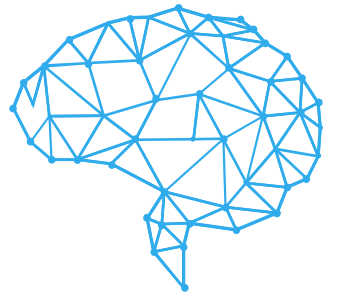
$$\nabla_w J(\mathbf{w}) = \nabla_w \frac{1}{2m} (\mathbf{X}\vec{\mathbf{w}} - \vec{\mathbf{y}})^T (\mathbf{X}\vec{\mathbf{w}} - \vec{\mathbf{y}})$$

DERIVE THE EQUALITY

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

NORMAL EQUATIONS

DERIVATION OF EQUATIONS



Equating the derivative to zero to find the **minimum**, we obtain the **vector of weights** that gives this **minimum**:

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

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



PROBABILISTIC INTERPRETATION

ASSUMPTIONS FOR ERRORS
LIKELIHOOD FUNCTION

PROBABILISTIC INTERPRETATION

DEVELOPMENT MOTIVATION



**WHY DO WE MINIMIZE THE MEAN QUADRATIC
ERROR FUNCTION AND NOT ANOTHER FUNCTION?**

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



The **output variables** $y^{(i)}$ can be defined as the **hypothesis** h_w **plus** an **estimation error** $\varepsilon^{(i)}$ that captures the **effects** that we **did not consider** in our **hypothesis** or contemplates the **random errors**.

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

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



The **following assumptions** are made for the **errors $\varepsilon^{(i)}$ (random sampling): IID**

Independent Errors - The probability of one error does not affect the probability of other errors.

Identically distributed errors: errors are sampled from the same probability distribution.

Distributed by a Gaussian probability distribution with mean $\mu = 0$ and variance σ^2 .

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

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



Therefore, the **probability density function** of $\varepsilon^{(i)}$ is given by:

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

INTERPRETACIÓN PROBABILÍSTICA

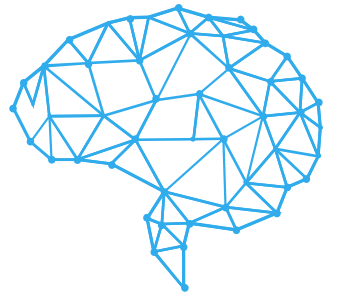
ASSUMPTION OF ERRORS



**WHY PROPOSE A MODEL THAT IS NORMALLY
DISTRIBUTED?**

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



Furthermore, the **outputs** $y^{(i)}$ and the **inputs** $x^{(i)}$ are set as **random variables**.

Consequently, the following question is posed:

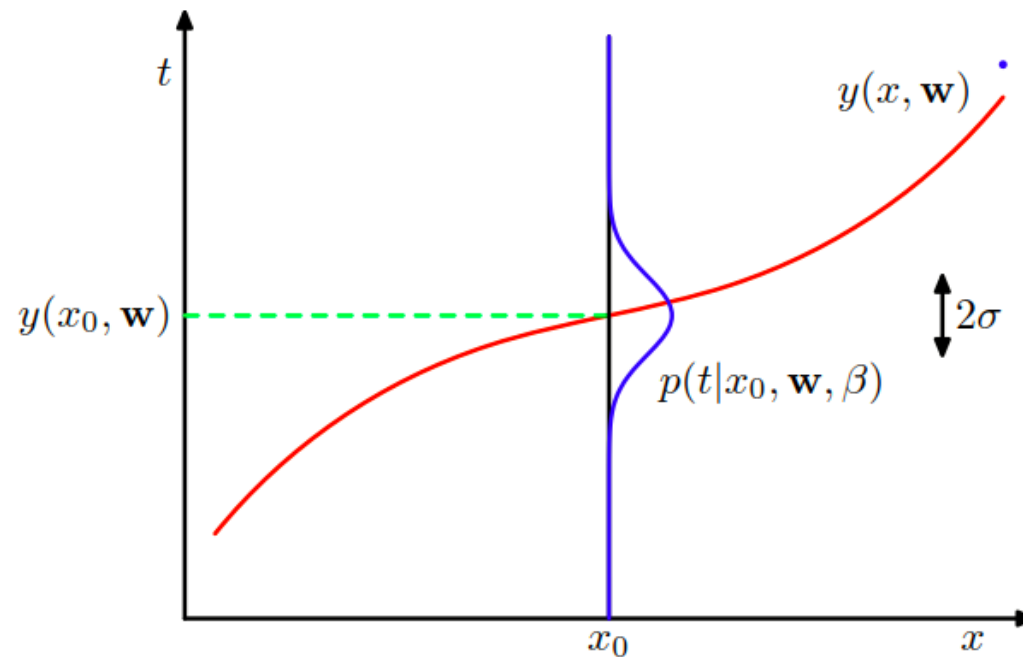
**Given my input data set X and my weight vector \vec{w} ,
what is the probability distribution that models my
outputs \vec{y} ?**

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS

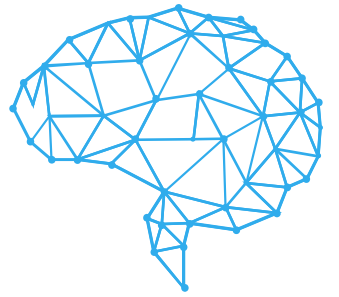


Following the **assumptions** from the **normal distribution of errors**, it naturally occurs that the **distribution** of the outputs $y^{(i)}$ follows a **normal form** with mean $h_w(x^{(i)})$ and variance σ^2 .



INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



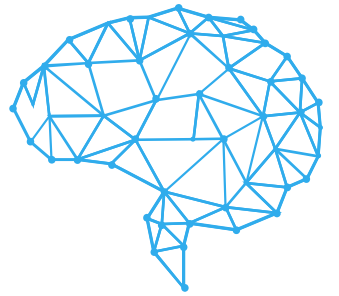
Formally it would look like this:

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

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

INTERPRETACIÓN PROBABILÍSTICA

ASSUMPTION OF ERRORS



We define the **conditional probability** for **all training data** (assuming the **samples** were **collected independently**):

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

INTERPRETACIÓN PROBABILÍSTICA

LIKELIHOOD FUNCTION



This equation is called the **likelihood function**, because it describes the **probability** that our **beliefs** about the **reality** of the **observations** (data) are **true**. In other words, that the **hypothesis** we proposed is **correct**.

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

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

INTERPRETACIÓN PROBABILÍSTICA

LIKELIHOOD FUNCTION



Therefore, we want to **maximize** the **probability** that **our beliefs** about how the **data is distributed** (in a **normal** way).

That is, we want to **maximize** the **likelihood function**. From a **frequentist** perspective (the value of \vec{w} is **not random**), we want to **find** the **value** of \vec{w} that **maximizes** the **probability** that the **observations** made of **reality** are **true**.

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

INTERPRETACIÓN PROBABILÍSTICA

LOGARITHMIC LIKELIHOOD



Maximizing a summation is much easier than doing so with a **multiplication**

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

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

INTERPRETACIÓN PROBABILÍSTICA

LOGARITHMIC LOSS



The **maximization problem** is **converted** into a **minimization** one by **scaling** the function with a **minus sign**.

This **function** $-\log(L(\vec{w}))$ is called **logarithmic loss**.

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

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

INTERPRETACIÓN PROBABILÍSTICA

LOGARITHMIC LOSS



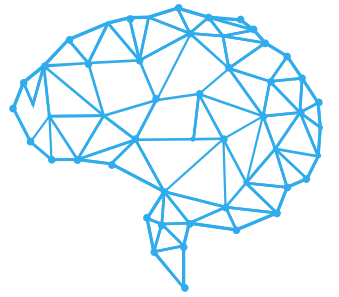
We develop the equation:

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

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

INTERPRETACIÓN PROBABILÍSTICA

LOGARITHMIC LOSS



We develop the equation:

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

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

INTERPRETACIÓN PROBABILÍSTICA

MEAN SQUARED ERROR



Minimizing the logarithmic loss is the same as minimizing the mean square error $MSE = J(\vec{w})$:

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

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



AI

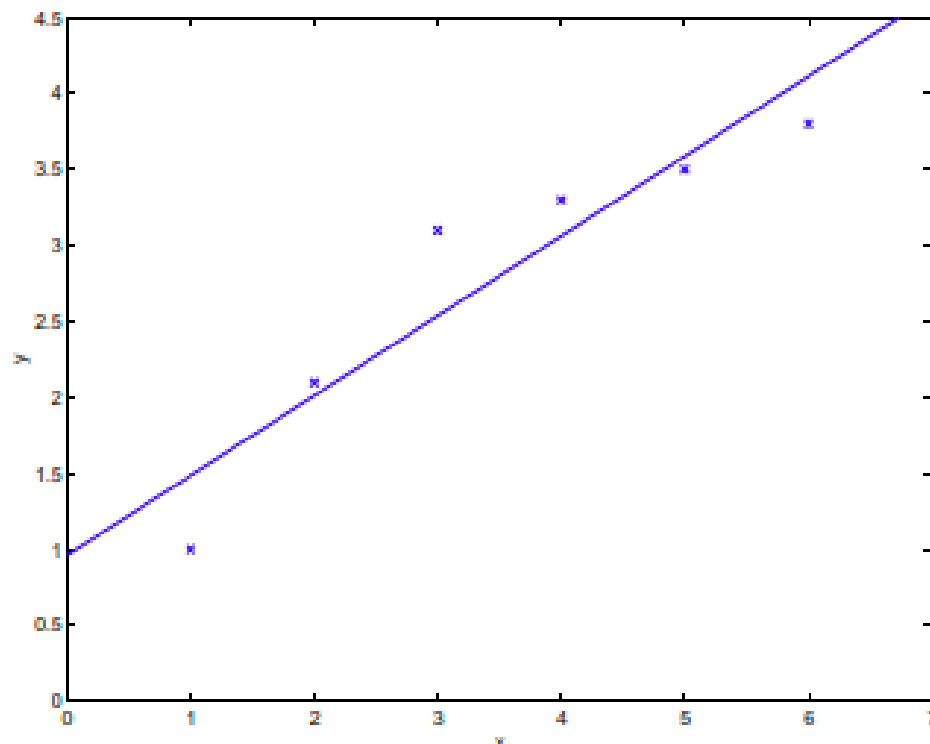
BASIS FUNCTIONS

B A S I S F U N C T I O N S

L I N E A R I T Y P R O B L E M

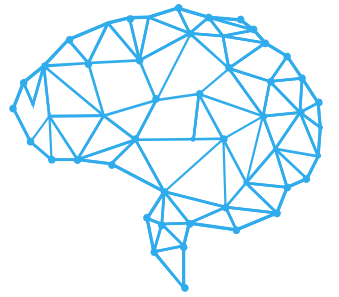


The **real world** has **non-linear** behaviors, so the **linear regression models** seen so far have their **limitations**.

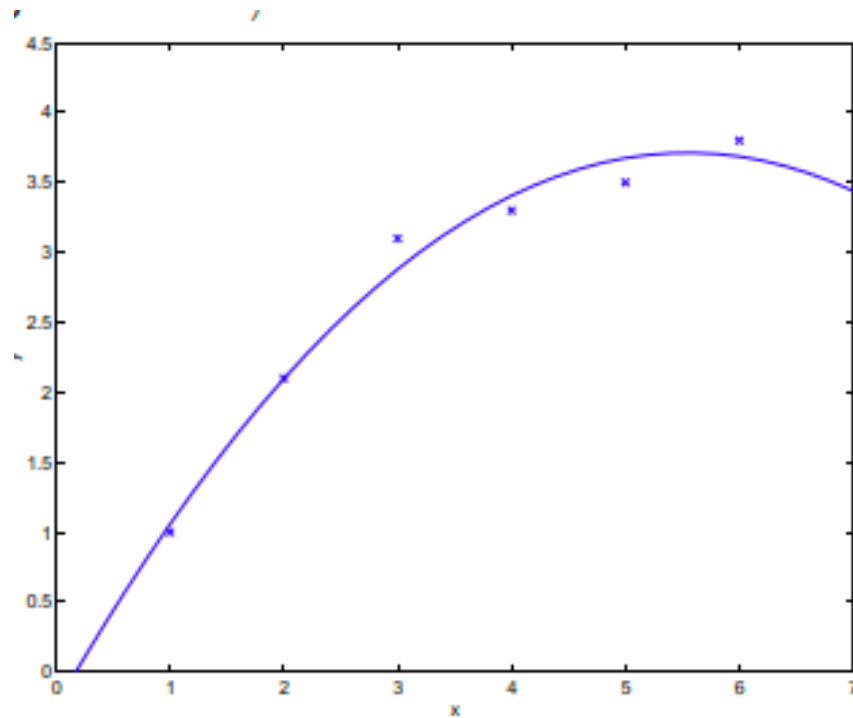


BASIS FUNCTIONS

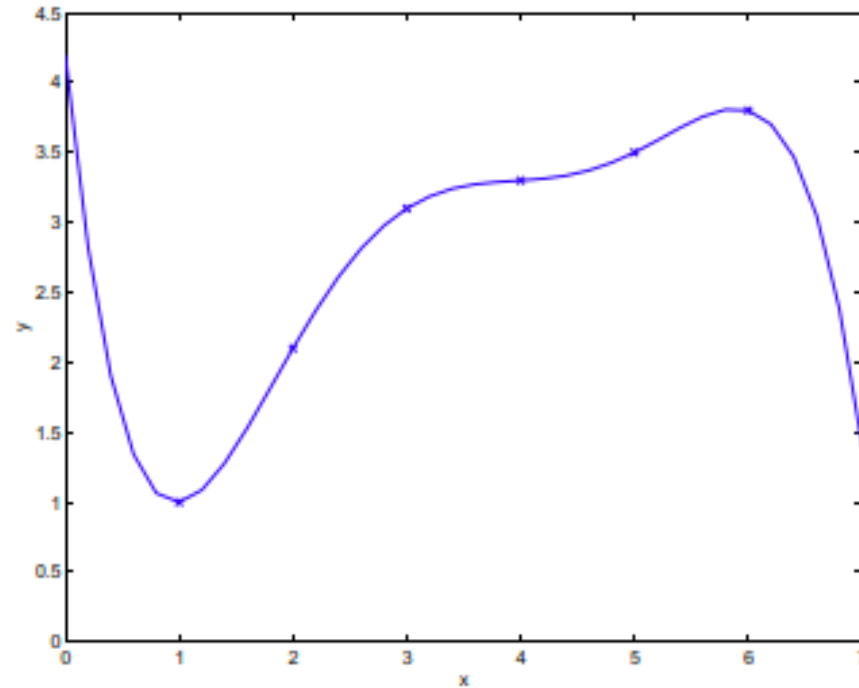
LINEARITY PROBLEM



We can **establish** our **hypothesis** as a **polynomial model**.



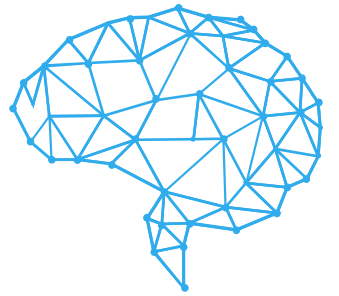
$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$



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

B A S I S F U N C T I O N S

AN INTRODUCTION TO BASIS FUNCTIONS



If we take this kind of reasoning further, we can **build** models as **linear combinations** of **nonlinear functions** $\phi_j(x)$, called **base functions**. Where $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^k$

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

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

NOTA

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

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

Thus we can **build** a **non-linear model**, which is still **parametrized** by **linear weights** w .

B A S I S F U N C T I O N S

AN INTRODUCTION TO BASIS FUNCTIONS



In a more detailed fashion, we have that:

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

B A S I S F U N C T I O N S

CLASSICAL LINEAR REGRESSION



In the case of **classical linear regression** we have for a **single** training data $\phi(x) \in \mathbb{R}^k$ in this case $k = n$:

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

B A S I S F U N C T I O N S

CLASSICAL LINEAR REGRESSION



In the case of **classical linear regression** we have for a **single** training data $\phi(x) \in \mathbb{R}^k$ in this case $k = n$:

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

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

B A S I S F U N C T I O N S

CLASSICAL LINEAR REGRESSION



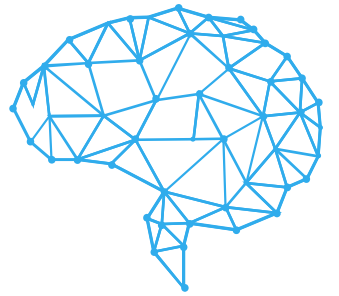
In the case of **classical linear regression**, we have for m training **data**:

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

$$w = [w_0 \quad \dots \quad w_k]$$

B A S I S F U N C T I O N S

CLASSICAL LINEAR REGRESSION



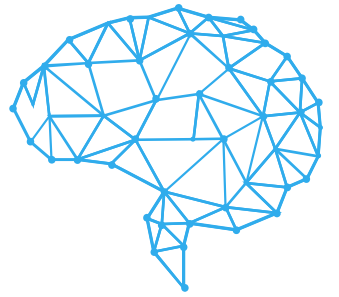
In the case of **classical linear regression**, we have for m training **data**:

$$\mathbf{h}_w(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \mathbf{w}_0 + \cdots + \mathbf{w}_k \mathbf{x}_n^{(1)} \\ \vdots \\ \mathbf{w}_0 + \cdots + \mathbf{w}_k \mathbf{x}_n^{(m)} \end{bmatrix}$$

$$\mathbf{h}_w(\mathbf{x}) = \begin{bmatrix} \mathbf{h}_w(\mathbf{x}^{(1)}) \\ \vdots \\ \mathbf{h}_w(\mathbf{x}^{(m)}) \end{bmatrix}$$

B A S I S F U N C T I O N S

D E S I G N M A T R I X



Therefore, the **design matrix** for any $\phi(x)$ would be given by:

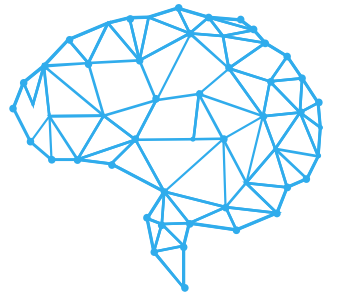
$$\phi(x) = \begin{bmatrix} \phi_0(x^{(1)}) & \cdots & \phi_{k-1}(x^{(1)}) \\ \vdots & \ddots & \vdots \\ \phi_0(x^{(m)}) & \cdots & \phi_{k-1}(x^{(m)}) \end{bmatrix}$$

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

Where each row of the matrix $\phi(x)$ is given by $\phi_j = \phi(x^{(i)})^T$

BASIS FUNCTIONS

LINEAR BASIS FUNCTIONS



Classic linear functions:

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

$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_k \end{bmatrix}; k = n$$

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

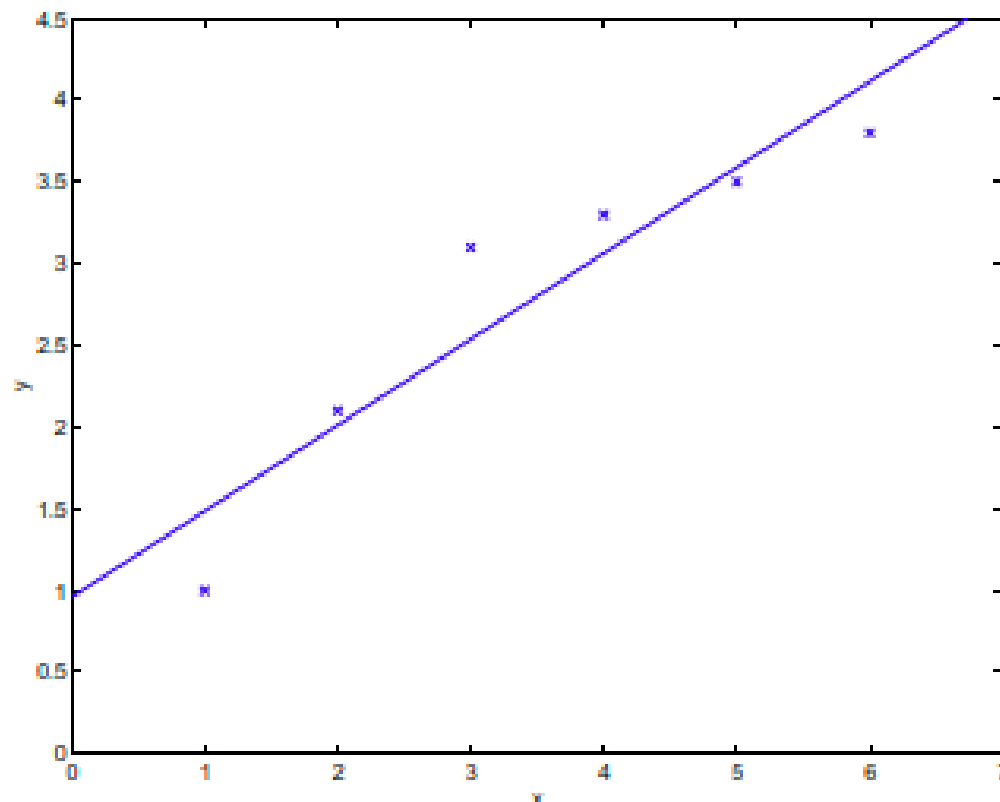
$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \begin{bmatrix} \mathbf{w}_0 + \mathbf{w}_1 \mathbf{x}_1^{(1)} + \cdots + \mathbf{w}_k \mathbf{x}_n^{(1)} \\ \vdots \\ \mathbf{w}_0 + \mathbf{w}_1 \mathbf{x}_1^{(m)} + \cdots + \mathbf{w}_k \mathbf{x}_n^{(m)} \end{bmatrix}$$

B A S I S F U N C T I O N S

LINEAR BASIS FUNCTIONS



Classic linear functions:



B A S I S F U N C T I O N S

POLYNOMIAL BASIS FUNCTIONS



Polynomial basis functions:

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

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

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

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

BASIS FUNCTIONS

POLYNOMIAL BASIS FUNCTIONS



Polynomial basis functions: concrete example polynomial degree 2 and two training data $m = 2$

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

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

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

B A S I S F U N C T I O N S

POLYNOMIAL BASIS FUNCTIONS



Polynomial basis functions: concrete example **polynomial degree 2** and two training data $m = 2$ with **2 characteristics**.

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

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

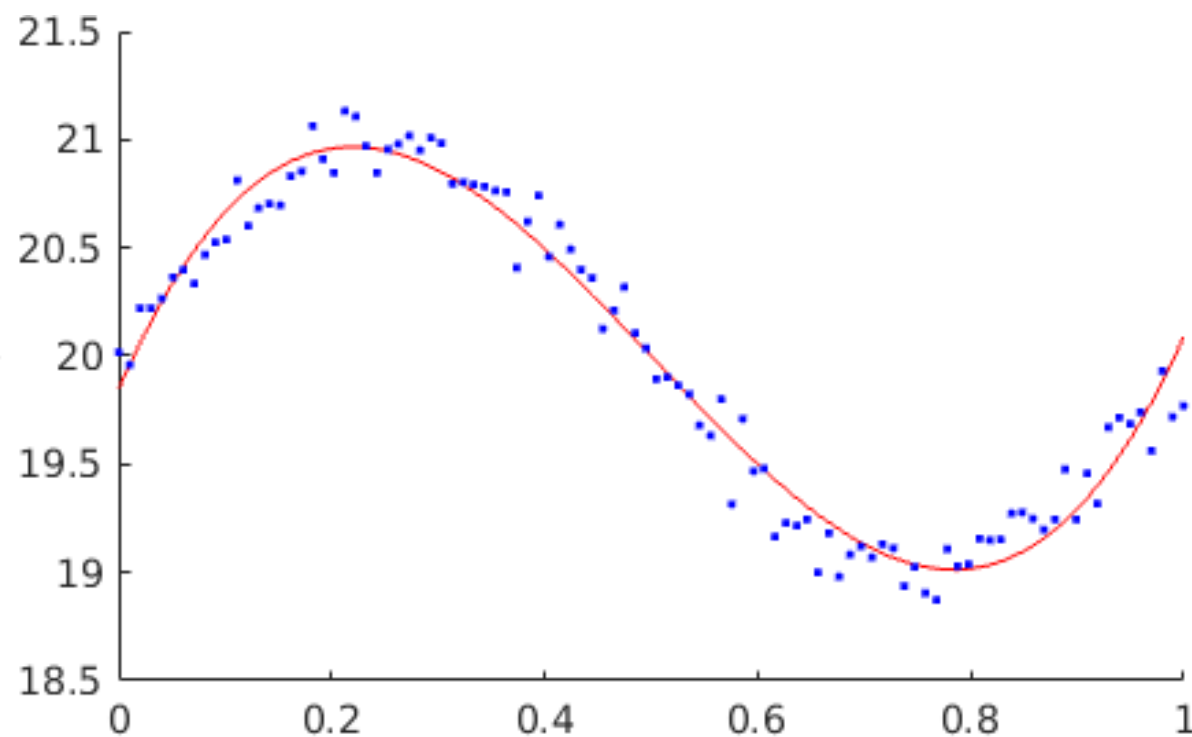
$$h_w(x) = \begin{bmatrix} w_0 + w_1 x_1^{(1)} + w_2 x_2^{(1)} + w_3 (x_1^{(1)})^2 + w_4 (x_2^{(1)})^2 \\ w_0 + w_1 x_1^{(2)} + w_2 x_2^{(2)} + w_3 (x_1^{(2)})^2 + w_4 (x_2^{(2)})^2 \end{bmatrix}$$

B A S I S F U N C T I O N S

POLYNOMIAL BASIS FUNCTIONS



Polynomial basis functions:



B A S I S F U N C T I O N S

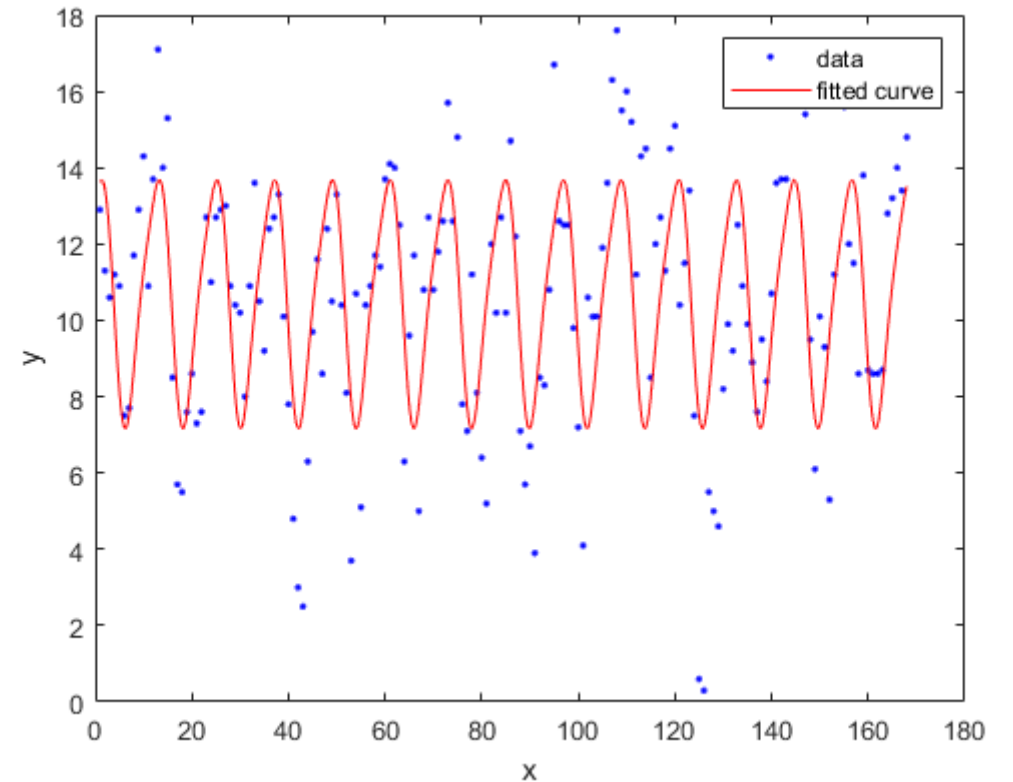
F O U R I E R S E R I E S



Fourier series

$$\phi_0(x) = 1$$

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



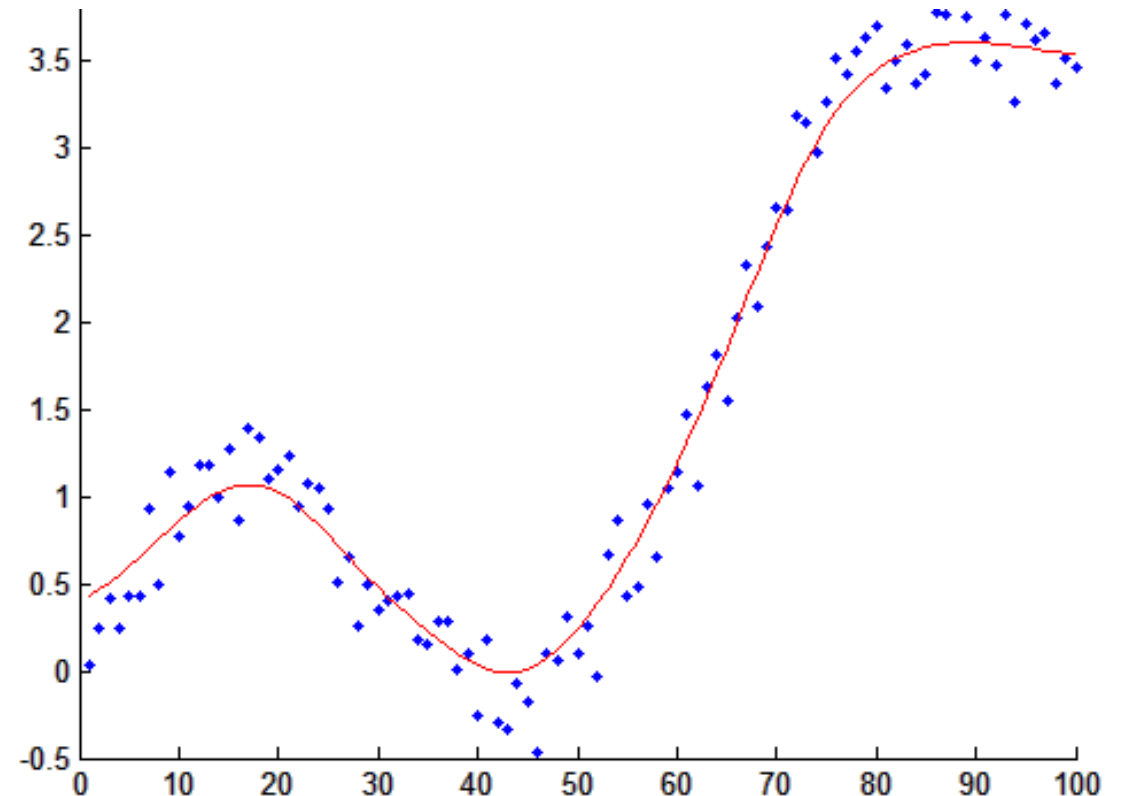
B A S I S F U N C T I O N S

R A D I A L B A S I S



Radial basis function:

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



B A S I S F U N C T I O N S

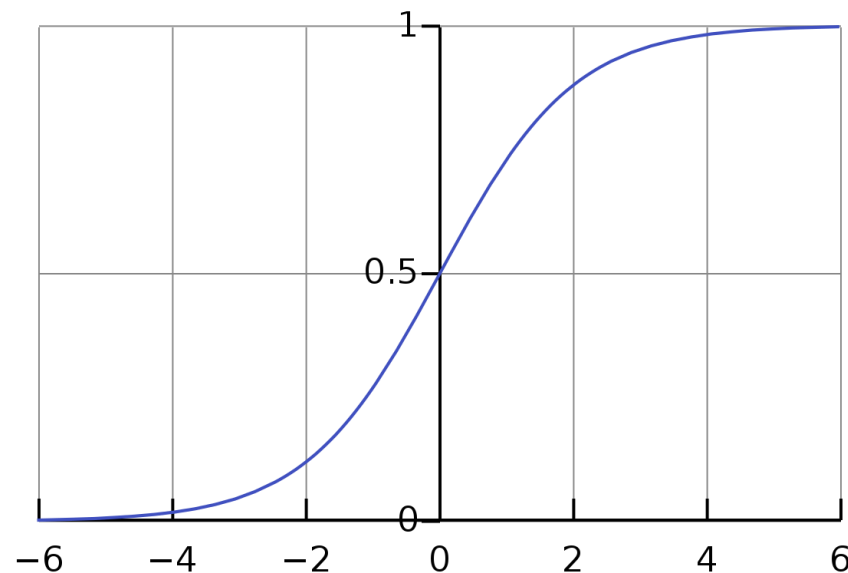
S I G M O I D A L S



Sigmoid function:

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

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



B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D



The **same** estimation **model** is proposed **plus** an **error**:

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

Where:

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

B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D

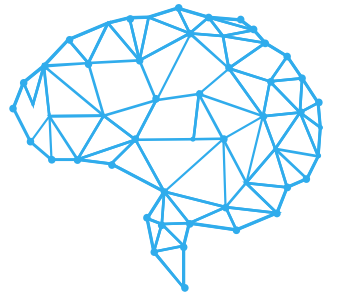


Assuming again that the **error** is **normally distributed** with **mean 0** and **variance β^{-1}** . It is defined for **m training data**:

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

B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D



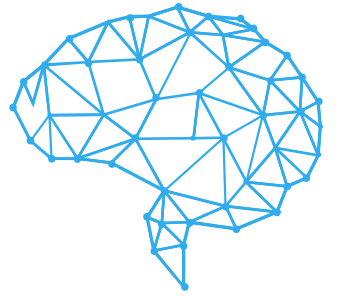
Since the **likelihood function** is **parameterized** in w and β the logarithmic likelihood is written like this:

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

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

B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D



Developing:

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

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

B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D



Developing:

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

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

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

B A S I S F U N C T I O N S

M A X I M U M L I K E L I H O O D



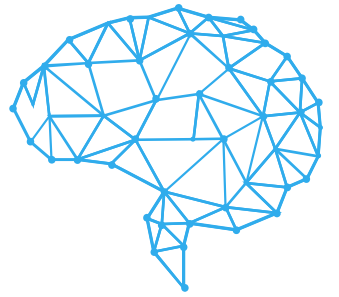
Developing:

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

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

B A S I S F U N C T I O N S

N O R M A L E Q U A T I O N S



Calculating the **gradient**, **setting** it equal to **zero** and solving for the **vector w** :

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