

Non-linear Regression

Cristian López Del Alamo
clopezd@utec.edu.pe
GINIA - Research group

2024

1

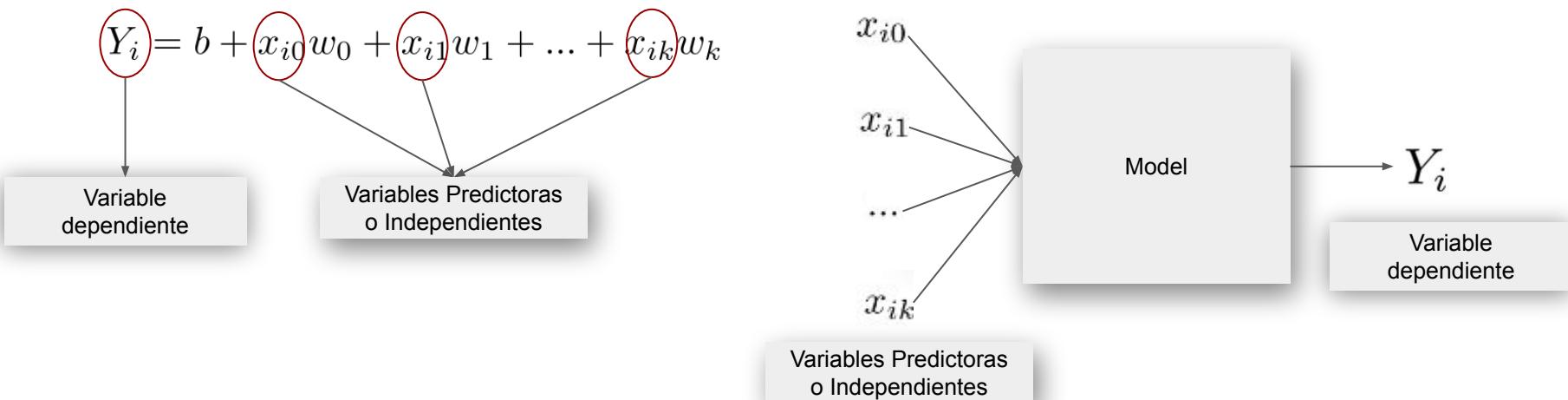
Remembering : Linear Regression



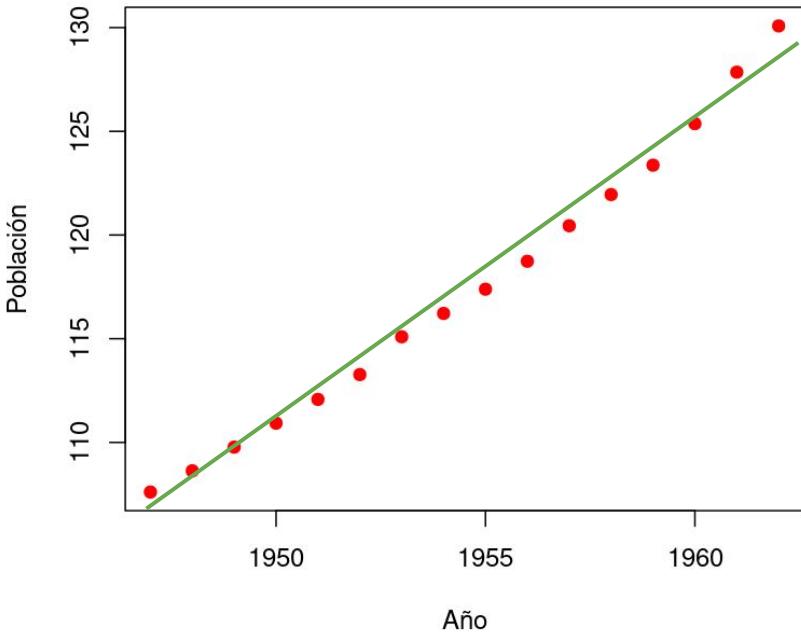
Cristian López Del Alamo

Objective: In this class, the goal is to understand the similarity between linear regression and non-linear regression, their applications, and the related machine learning algorithm.





Simple Linear Regression



1. Hypothesis

$$h(x_i) = b + wx_i$$

2. Loss Function

$$\mathcal{L} = \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n}$$

3. Derivative Calculus

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\sum_{i=0}^n (y_i - h(x_i))(-x_i)}{n}$$

4. Weights Update

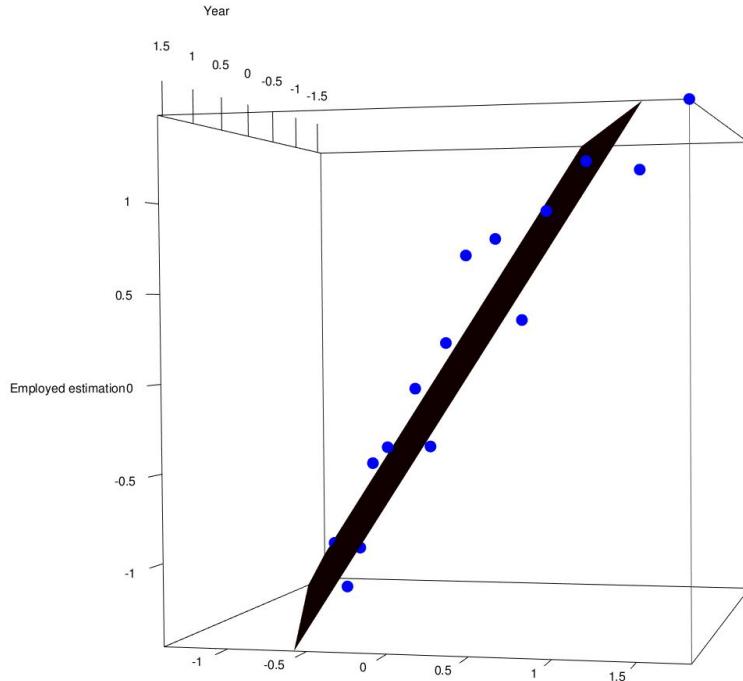
$$w = w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\sum_{i=0}^n (y_i - h(x_i))(-1)}{n}$$

$$b = b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$

Cristian López Del Alamo

Multiple Linear Regression: Case of 2 variables



1. Hypothesis

$$h(x_i) = b + x_{(i,1)}w_1 + x_{(i,2)}w_2$$

2. Loss Function

$$\mathcal{L} = \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n}$$

3. Derivative Calculus

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\sum_{i=0}^n (y_i - h(x_i))(-x_{(i,j)})}{n}$$

4. Weights Update

$$w_j = w_j - \alpha \frac{\partial \mathcal{L}}{\partial w_j}$$

$$b = b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\sum_{i=0}^n (y_i - h(x_i))(-1)}{n}$$

Multiple Linear Regression: Case of k variables

Hypothesis

$$h(x_i) = b + x_{(i,1)}w_1 + x_{(i,2)}w_2 + \dots + x_{(i,k)}w_k$$

$$h(x_i) = w_0 + x_{(i,1)}w_1 + x_{(i,2)}w_2 + \dots + x_{(i,k)}w_k$$

$$h(x_i) = x_{(i,0)}w_0 + x_{(i,1)}w_1 + x_{(i,2)}w_2 + \dots + x_{(i,k)}w_k$$

$$h(x_i) = [x_{(i,0)} \ x_{(i,1)} \ x_{(i,2)} \ \dots \ x_{(i,k)}][w_0 \ w_1 \ w_2 \ \dots \ w_k]^T$$

$$h(x_i) = x_i * w^T$$

Multiple Linear Regression: Case of k variables

Hypothesis

$$X = \begin{bmatrix} 1 & x_{01} & \cdots & x_{0p} \\ 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix} \otimes W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_i \\ \vdots \\ w_p \end{bmatrix} \quad h(X) = \begin{bmatrix} h(0) \\ h(1) \\ \vdots \\ h(i) \\ \vdots \\ h(p) \end{bmatrix}$$

h(x_i) = $[1, x_{i,1}, \dots, x_{i,p}] \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix}$

$$h(X) = XW^t$$

Multiple Linear Regression: Case of k variables

Loss Function

$$Y = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix} = h(X) = \begin{bmatrix} h(x_0) \\ h(x_1) \\ \vdots \\ h(x_n) \end{bmatrix} \rightarrow Loss = \begin{bmatrix} (y_0 - h(x_0))^2 \\ (y_1 - h(x_1))^2 \\ \vdots \\ (y_n - h(x_n))^2 \end{bmatrix} \rightarrow \text{Normal } L_2$$

$$\mathcal{L} = \|y - h(X)\|_2^2$$

Multiple Linear Regression: Case of k variables

Derivative Computation

$$(Y - h(X)) = \begin{bmatrix} y_0 - h(x_0) \\ y_1 - h(x_1) \\ \vdots \\ y_i - h(x_i) \\ \vdots \\ y_n - h(x_n) \end{bmatrix}^T \otimes (-1) X = \begin{bmatrix} 1 & x_{01} & \cdots & x_{0p} \\ 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix} \rightarrow \frac{\partial \mathcal{L}}{\partial w} = \frac{(Y - h(X))^T * (-1 * X)}{n}$$

Multiple Linear Regression: Case of k variables

Updating w

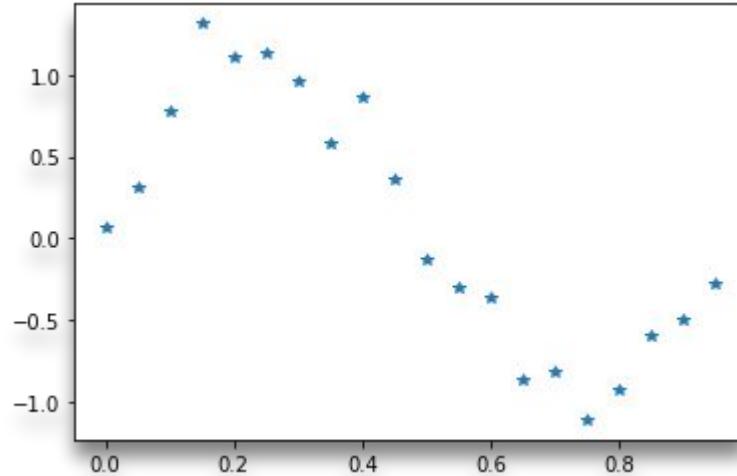
$$W = W - \alpha * \frac{\partial \mathcal{L}}{\partial w}$$

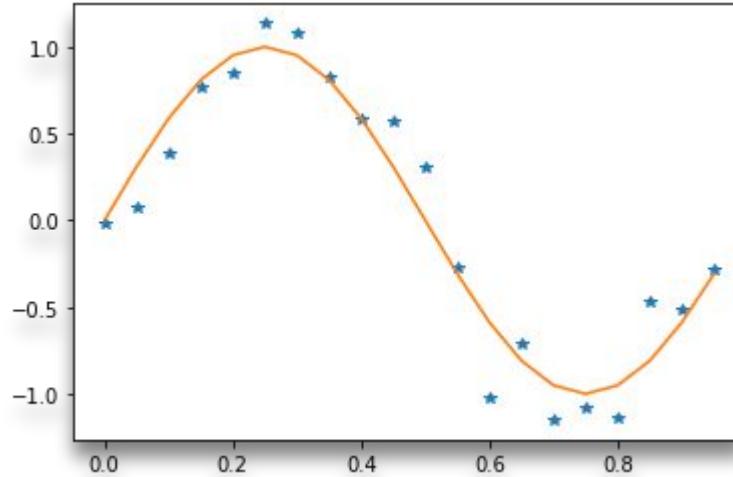
2

Nonlinear Regression



Objective: Understand the mathematics behind multivariable nonlinear regression.





Hypothesis: Polynomial

$$h(x_i) = b + x_i w_1 + x_i^2 w_2 + x_i^3 w_3 + \dots + x_i^p w_p$$

$$h(x_i) = b + \sum_{j=1}^p x_i^j w_j$$

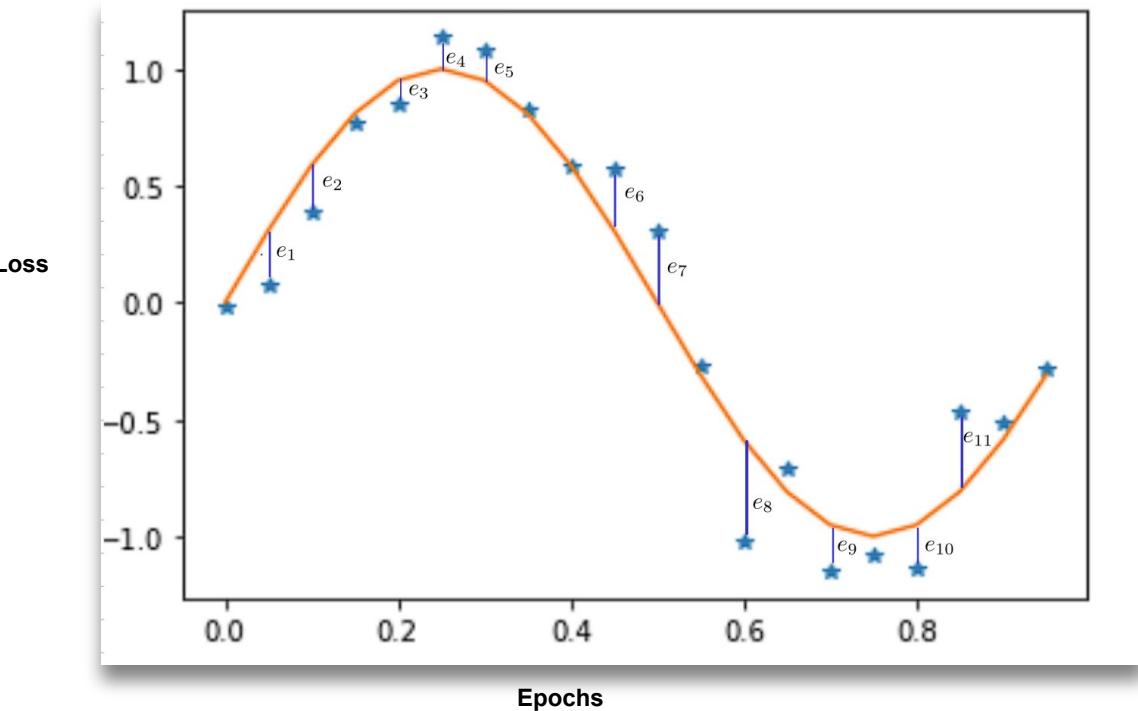
$$h(x_i) = x_i^0 w_0 + \sum_{j=1}^p x_i^j w_j$$

$$h(x_i) = \sum_{j=0}^p x_i^j w_j$$

$$h(X) = \begin{pmatrix} h(x_1) \\ h(x_2) \\ \vdots \\ \boxed{h(x_i)} \\ \vdots \\ h(x_n) \end{pmatrix} \quad \square \quad X = \begin{pmatrix} x_1^0 & x_1^1 & \cdots & x_1^j & \cdots & x_1^p \\ x_2^0 & x_2^1 & \cdots & x_2^j & \cdots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots \\ \boxed{x_i^0 & x_i^1 & \cdots & x_i^j & \cdots & x_i^p} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots \\ x_n^0 & x_n^1 & \cdots & x_n^j & \cdots & x_n^p \end{pmatrix} \quad \times \quad W^T = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_i \\ \vdots \\ w_p \end{pmatrix}$$

$$h(X) = XW^T$$

Loss Function



Loss Function

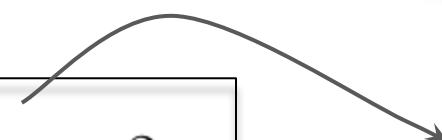
$$\mathcal{L} = \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix} \quad \square \quad h(X) = \begin{pmatrix} h(x_1) \\ h(x_2) \\ \vdots \\ h(x_i) \\ \vdots \\ h(x_n) \end{pmatrix} \quad \square \quad Y - \hat{Y} = \begin{pmatrix} y_1 - h(x_1) \\ y_2 - h(x_2) \\ \vdots \\ y_i - h(x_i) \\ \vdots \\ y_n - h(x_n) \end{pmatrix} \quad \square \quad E = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_i \\ \vdots \\ e_n \end{pmatrix}$$

$$E = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_i \\ \vdots \\ e_n \end{pmatrix} \otimes E^T = (e_1 \ e_2 \ \cdots \ e_i \ \cdots \ e_n) \equiv E \ E^T$$

$$L = \frac{EE^T}{2n}$$

$$\mathcal{L} = \frac{\|Y - h(X)\|_2^2}{2n}$$



L2 norm

```
1 def train(x, y, umbral, alfa, p):
2     np.random.seed(2001)
3     w = [np.random.rand() for i in range(1,p+1)]
4     y_pred = h(x,w)
5     L = Error(y, y_pred)
6     while (L > umbral):
7         dw = derivada(x, y, w)
8         w = update(w,db, alfa)
9         y = h(x,w)
10        L = Error(y, y_pred)
11    return w,L
```

```
1 def train(x, y, umbral, alfa, p):
2     np.random.seed(2001)
3     w = [np.random.rand() for i in range(1,p+1)]
4     y_pred = h(x,w)
5     L = Error(y, y_pred)
6     while (L > umbral):
7         dw = derivada(x, y, w)
8         w = update(w,db, alfa)
9         y = h(x,w)
10        L = Error(y, y_pred)
11    return w,L
```

```
1 def train(x, y, umbral, alfa, p):
2     np.random.seed(2001)
3     w = [np.random.rand() for i in range(1,p+1)]
4     y_pred = h(x,w)
5     L = Error(y, y_pred)
6     while (L > umbral):
7         dw = derivada(x, y, w)
8         w = update(w,db, alfa)
9         y = h(x,w)
10        L = Error(y, y_pred)
11    return w,L
```

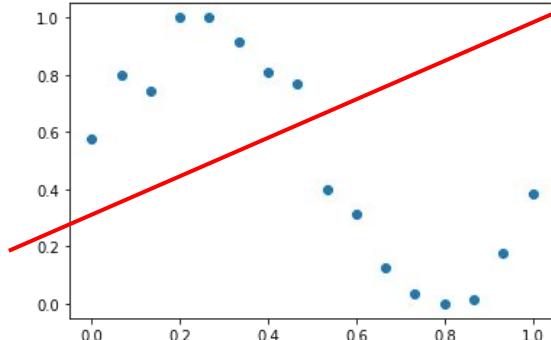
```
1 def train(x, y, umbral, alfa, p):
2     np.random.seed(2001)
3     w = [np.random.rand() for i in range(1,p+1)]
4     y_pred = h(x,w)
5     L = Error(y, y_pred)
6     while (L > umbral):
7         dw = derivada(x, y, w)
8         w = update(w,dw, alfa)
9         y = h(x,w)
10        L = Error(y, y_pred)
11    return w,L
```

```
1 def train(x, y, umbral, alfa, p):
2     np.random.seed(2001)
3     w = [np.random.rand() for i in range(1,p+1)]
4     y_pred = h(x,w)
5     L = Error(y, y_pred)
6     while (L > umbral):
7         dw = derivada(x, y, w)
8         w = update(w,db, alfa)
9         y = h(x,w)
10        L = Error(y, y_pred)
11    return w,L
```

3

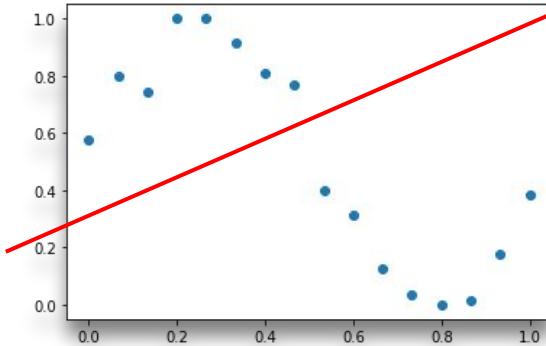
Regularization

Underfitting



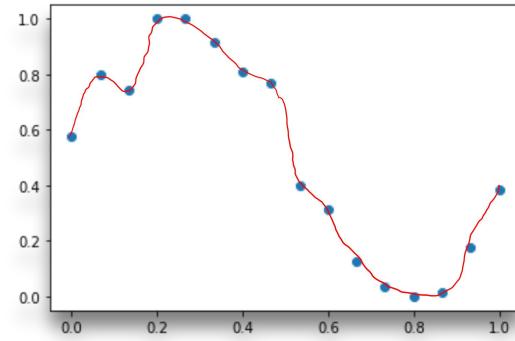
$$h(x_i) = x_i w + b$$

Underfitting



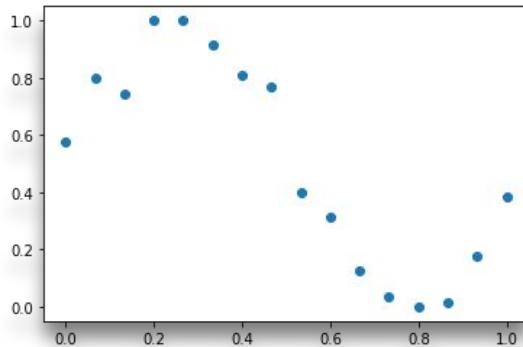
$$h(x_i) = x_i w + b$$

Overfitting



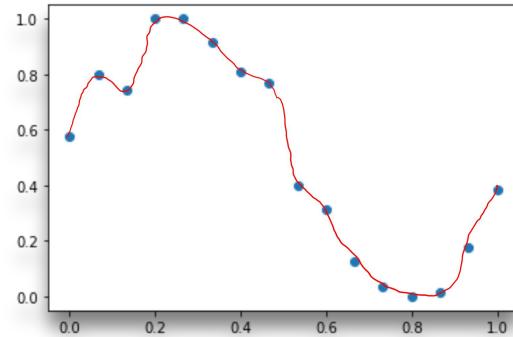
$$h(x_i) = x_i^0 w_0 + x_i^1 w_1 + \dots + x_i^{20} w_{20}$$

Underfitting



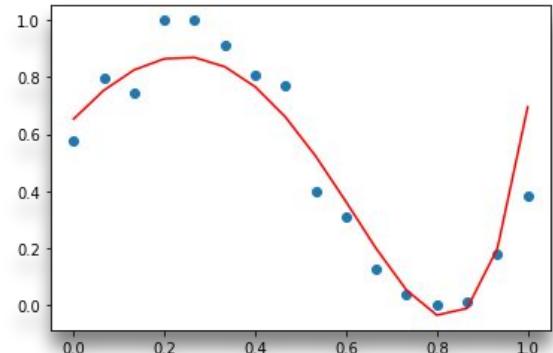
$$h(x_i) = x_i w + b$$

Overfitting

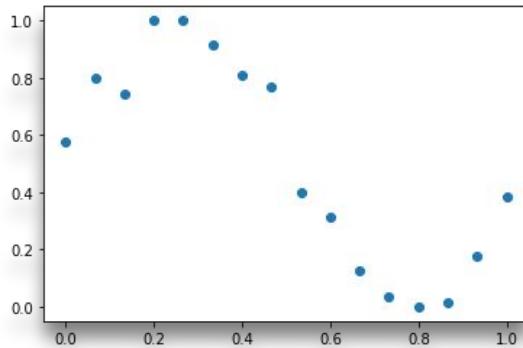


$$h(x_i) = x_i^0 w_0 + x_i^1 w_1 + \dots + x_i^{20} w_{20}$$

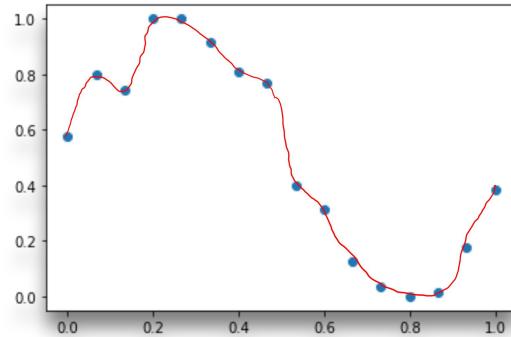
good



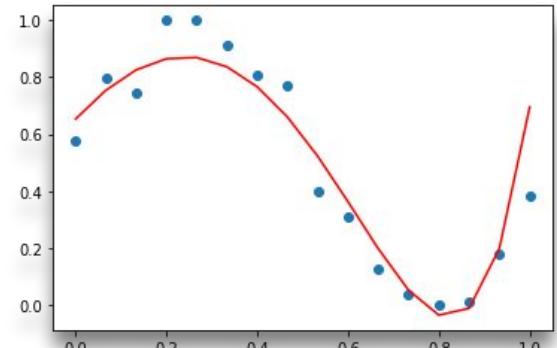
$$h(x_i) = x_i^0 w_0 + x_i^1 w_1 + \dots + x_i^3 w_3$$

Underfitting

$$h(x_i) = x_i w + b$$

Overfitting

$$h(x_i) = x_i^0 w_0 + x_i^1 w_1 + \dots + x_i^{20} w_{20}$$

good

$$h(x_i) = x_i^0 w_0 + x_i^1 w_1 + \dots + x_i^3 w_3$$

- Simple Model
- Low Capacity Model

- Highly Complex Model
- High Capacity Model

- Data-Fitting Model
- Model with Adequate Capacity



How do we decrease the complexity of our model?

Regularization

Regularization term

$$\mathcal{L} = \frac{||Y - h(X)||_2^2}{2n} + \boxed{\mathcal{R}(w)}$$

$$\mathcal{L} = \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n} + \boxed{\frac{\lambda}{n} \sum_{j=1}^p w_j^2}$$

$$\mathcal{L} = \min \left\{ \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n} + \frac{\lambda}{n} \sum_{j=1}^p w_j^2 \right\}$$

$$\mathcal{L} = \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n} + \boxed{\lambda \sum_{j=1}^p w_j^2}$$

$$\mathcal{L} = \min \left\{ \frac{\sum_{i=0}^n (y_i - h(x_i))^2}{2n} + \lambda \sum_{j=1}^p w_j^2 \right\}$$

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \lambda \|W\|_2^2$$

If the parameter λ is very large, what happens?

$$\lambda = 10000$$

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \lambda \|W\|_2^2$$

What happens if the parameter λ is very small?

$$\lambda = 0.01$$

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \lambda \|W\|_2^2$$



Who is affected or what changes if we modify the loss function

The Derivatives

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \left(\frac{\|Y - XW^t\|_2^2}{2n} + \lambda |W|^2 \right)}{\partial w_i}$$

1. Hypothesis

$$h(X) = X * w^T$$

2. Loss Function

$$\mathcal{L} = \frac{||Y - XW^t||_2^2}{2n} + \boxed{\frac{\lambda}{n} ||W||_2^2}$$

3. Derivatives

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\sum_{i=0}^n (y_i - h(x_i)) * (-x_i^j)}{n} + \boxed{\frac{2\lambda w_j}{n}}$$



2.1 Regularization Methods in Regression



Loss Function

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \lambda \|W\|_2^2$$

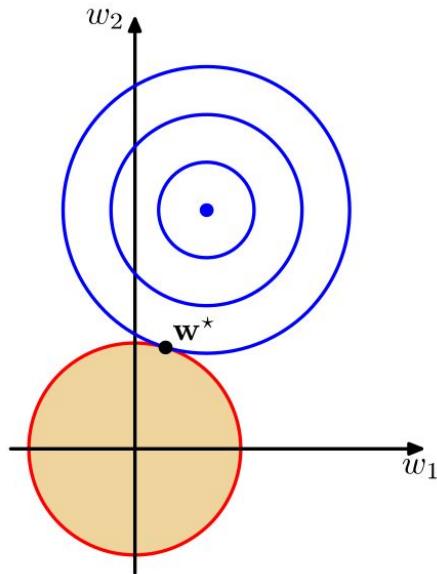
Ridge

Loss Function

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \lambda \|W\|_1$$

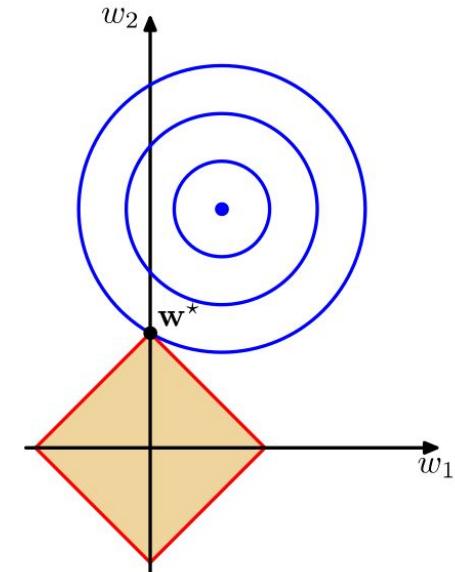
Lasso

Ridge



$$w_1^2 + w_2^2 \leq s$$

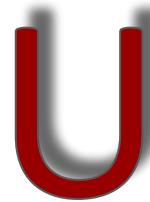
Lasso



$$|w_1| + |w_2| \leq s$$

$$\mathcal{L} = \frac{||Y - XW^t||_2^2}{2n} + \lambda ||W||_2^2$$

Small values of W



$$\mathcal{L} = \frac{||Y - XW^t||_2^2}{2n} + \lambda ||W||_1$$

Vector with multiple zero values



Elastinet

$$\mathcal{L} = \frac{\|Y - XW^t\|_2^2}{2n} + \rho\lambda\|W\|_1 + (1 - \rho)\lambda\|W\|_2$$

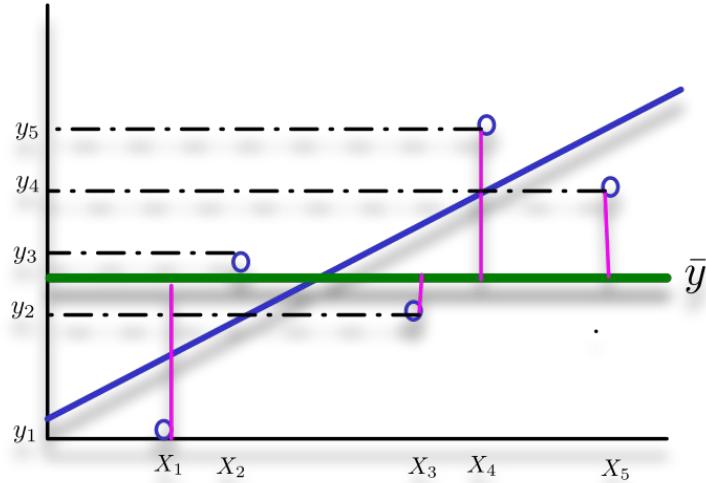


2.3 Measure of Regression Quality

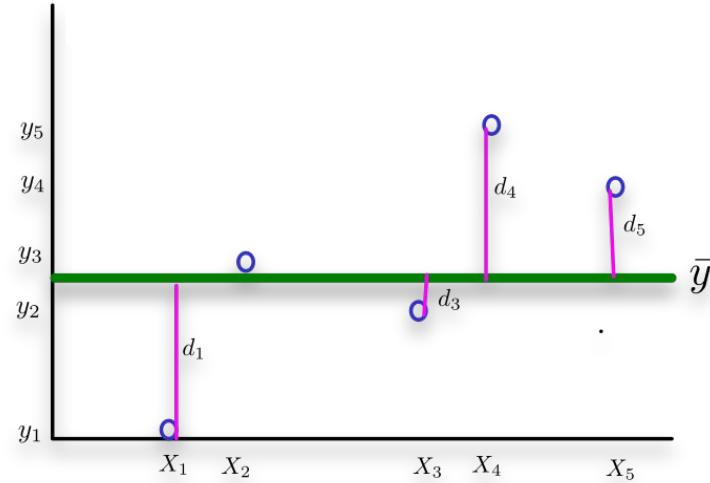
- Coefficient of determination R^2
- Adjusted R^2

Coefficient of determination R^2

$Y : Y \text{ Real}$



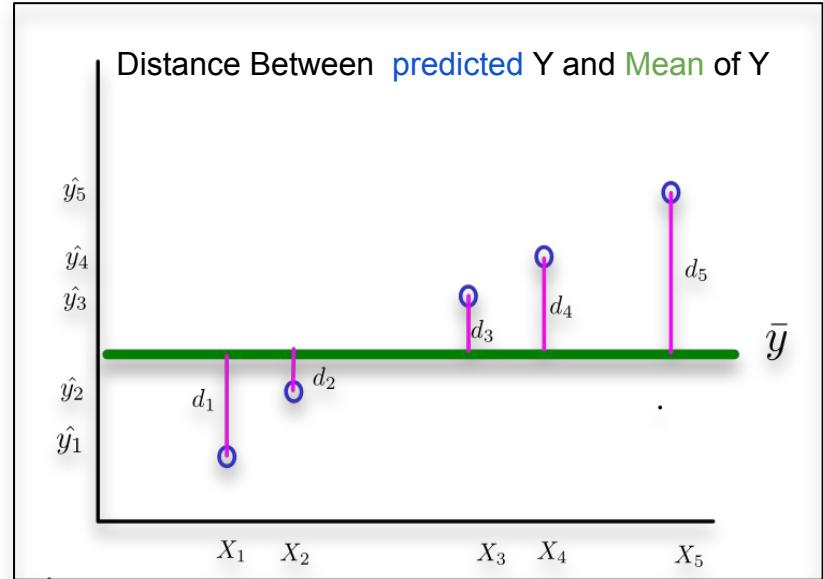
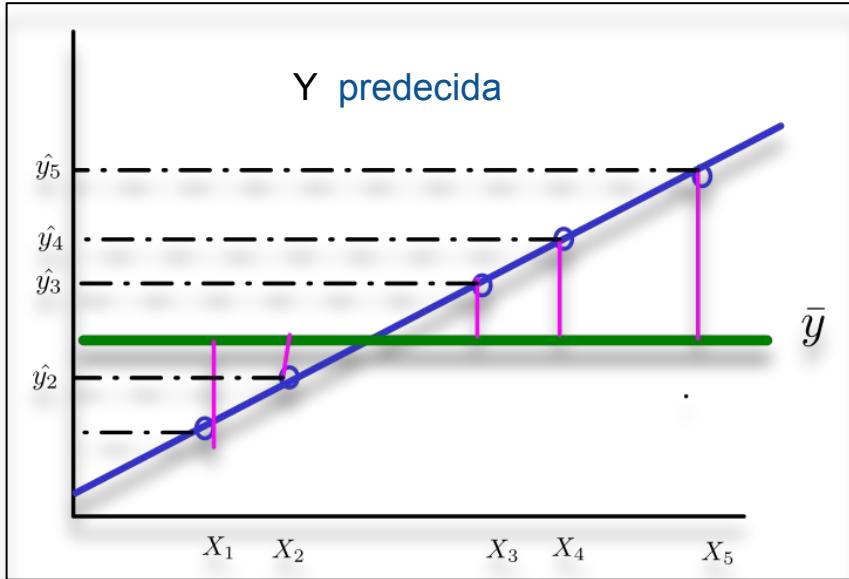
Distance Between real Y and Mean of Y



Distance Between real Y and Mean of Y :

$$(Y - \bar{Y})^2$$

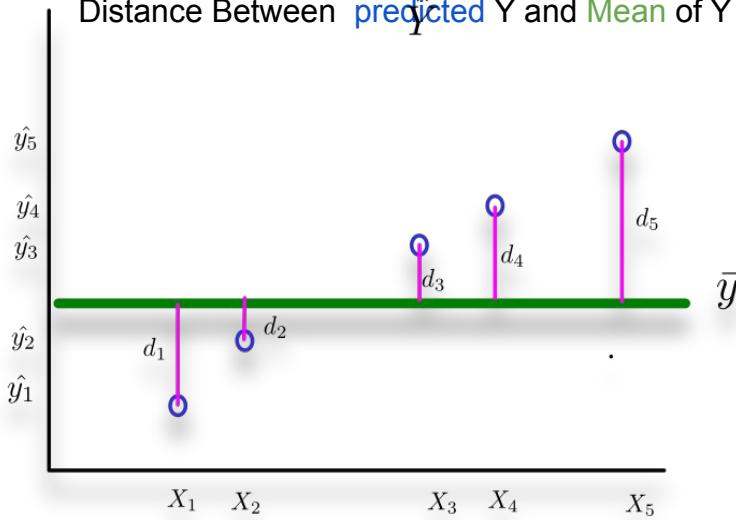
Coefficient of determination R^2



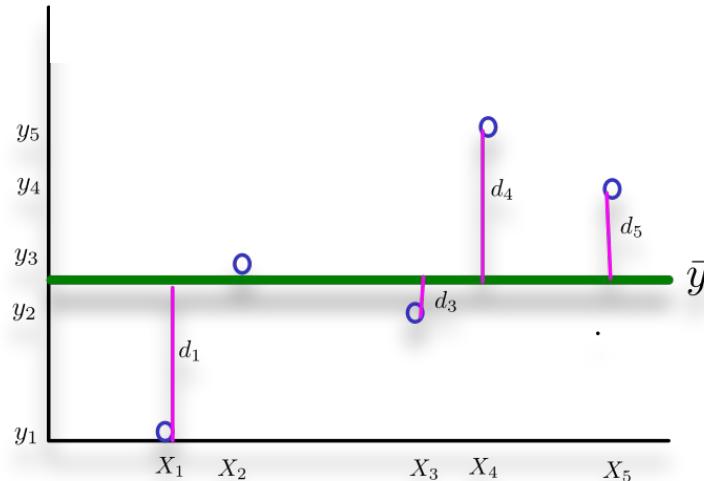
Distance Between predicted Y and Mean of Y : $(\hat{Y} - \bar{Y})^2$

Coefficient of determination R^2

Distance Between \hat{Y} and Mean of Y



Distance Between real Y and Mean of Y



Coefficient of Determination

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{Y})^2}{\sum_{i=1}^n (y_i - \bar{Y})^2}$$

$$\frac{\|\hat{Y} - \bar{Y}\|_2^2}{\|Y - \bar{Y}\|_2^2}$$

Coefficient of determination R^2

Interpretation

- $R^2 = 1$: The model completely explains the variation of the dependent variable.
- $R^2 = 0$: The model does not explain the variation of the dependent variable at all.
- Value between 0 and 1: The closer the R-squared value is to 1, the better the model.

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{Y})^2}{\sum_{i=1}^n (y_i - \bar{Y})^2}$$

Disadvantages

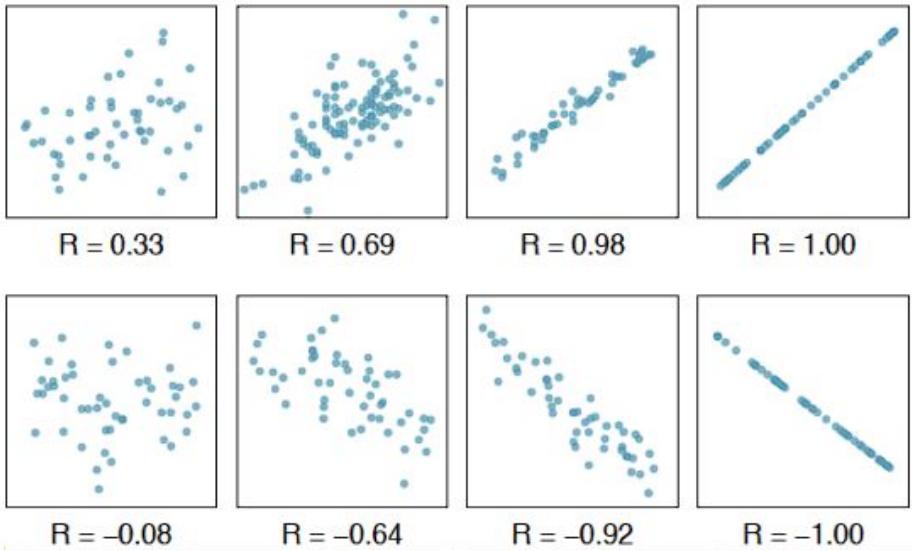
- Does not indicate causality.
- Sensitive to the number of predictors.
- Not suitable for non-linear models
- Not suitable when there are outliers
- No ajusta por tamaño de la muestra

Adjusted R^2

$$R_{ajustado}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

- n es el número total de observaciones
- p es el número de predictores independientes en el modelo,
- R^2 es el coeficiente de determinación no ajustado.
- Puede tomar valores negativos. En este caso las variables predictoras no tiene relación con la variable dependiente Y.

Example of correlation coefficient



Source: [Click](#)



3

Training, Bias, and Variance

Data Set

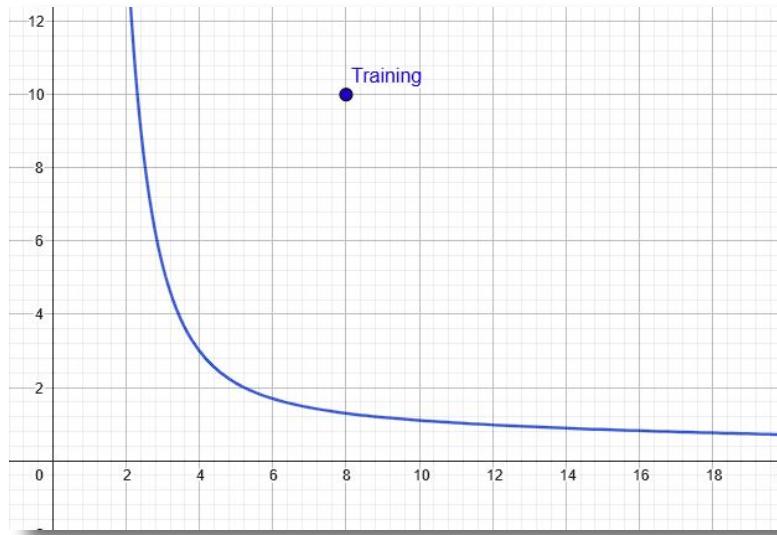
Training data

Testing data

80%

30%

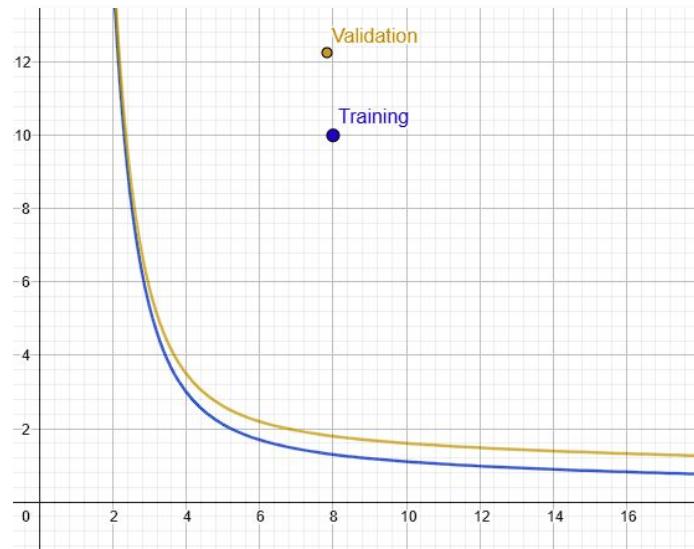
loss, w = training(x_train, y_trraig, alpha, epochs)



Data Set



loss, w = training(x_train, y_traing, alpha, epochs)



```
traing(x_train, y_train, alpha, epochs):
    w = # Initialize parameters
    for i in range(epochs):
        x,y = get_data(x_train, y_train, batch)
        dw = derivatives(x,y,w)
        w = Change_Parameters(w,dw, alpha)
        loss_training = Error(x,y,w,batch)
        loss_val = Error(x_val,y_val,w,batch)
        L_T.append(loss_training)
        L_V.append(loss_val)
    return L_T, L_V, w
```

```
traing(x_train, y_train, alpha, epochs):
    w = # Initialize parameters
    for i in range(epochs):
        x,y = get_data(x_train, y_train, batch)
        dw = derivatives(x,y,w)
        w = Change_Parameters(w,dw, alpha)
        loss_training = Error(x,y,w,batch)
        loss_val = Error(x_val,y_val,w,batch)
        L_T.append(loss_training)
        L_V.append(loss_val)
    return L_T, L_V, w
```

```
traing(x_train, y_train, alpha, epochs):
    w = # Initialize parameters
    for i in range(epochs):
        x,y = get_data(x_train, y_train, batch)
        dw = derivatives(x,y,w)
        w = Change_Parameters(w,dw, alpha)
        loss_training = Error(x,y,w,batch)
        loss_val = Error(x_val,y_val,w,batch)   
        L_T.append(loss_training)
        L_V.append(loss_val)
    return L_T, L_V, w
```

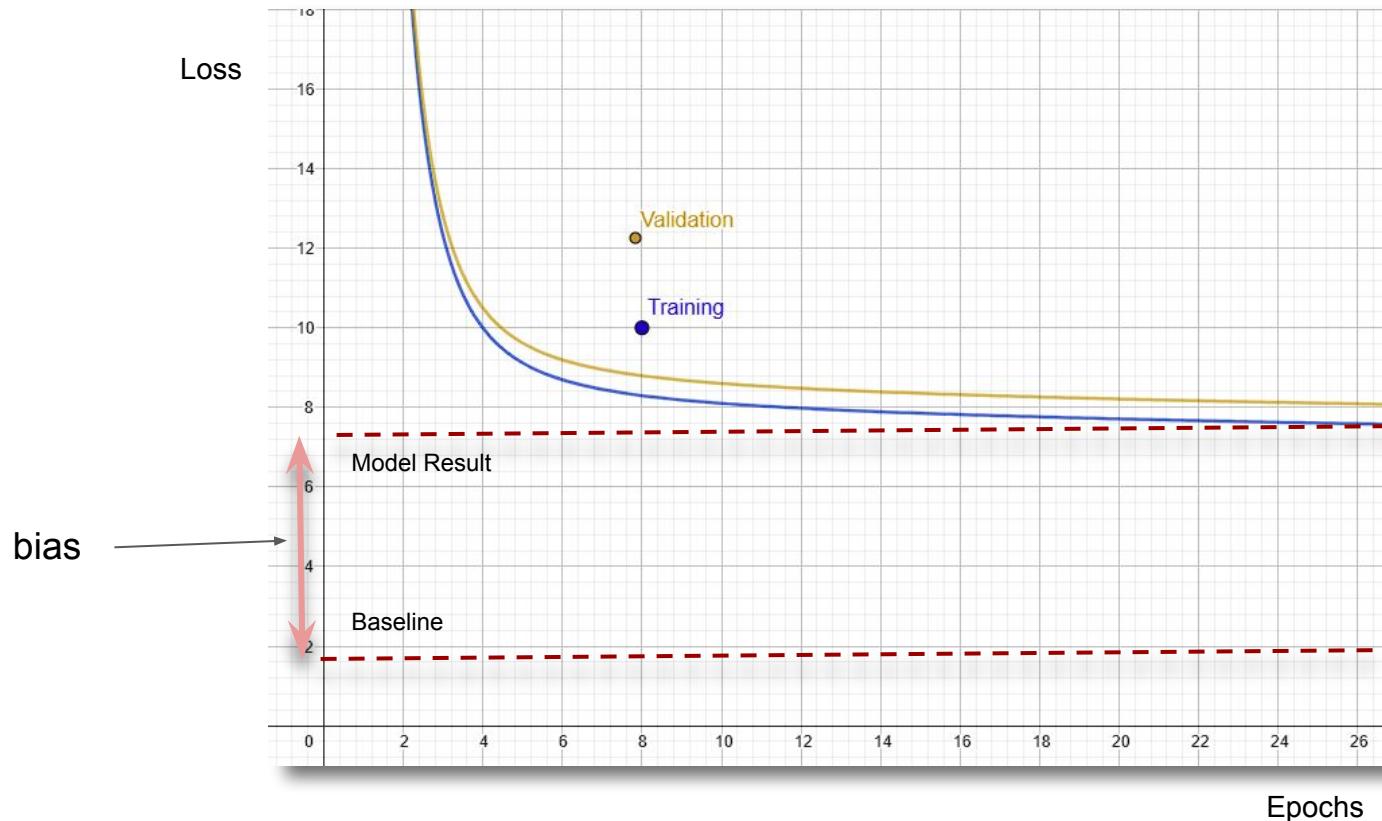
```
traing(x_train, y_train, alpha, epochs):
    w = # Initialize parameters
    for i in range(epochs):
        x,y = get_data(x_train, y_train, batch)
        dw = derivatives(x,y,w)
        w = Change_Parameters(w,dw, alpha)
        loss_training = Error(x,y,w,batch)
        loss_val = Error(x_val,y_val,w,batch)
        L_T.append(loss_training)
        L_V.append(loss_val)
    return L_T, L_V, w

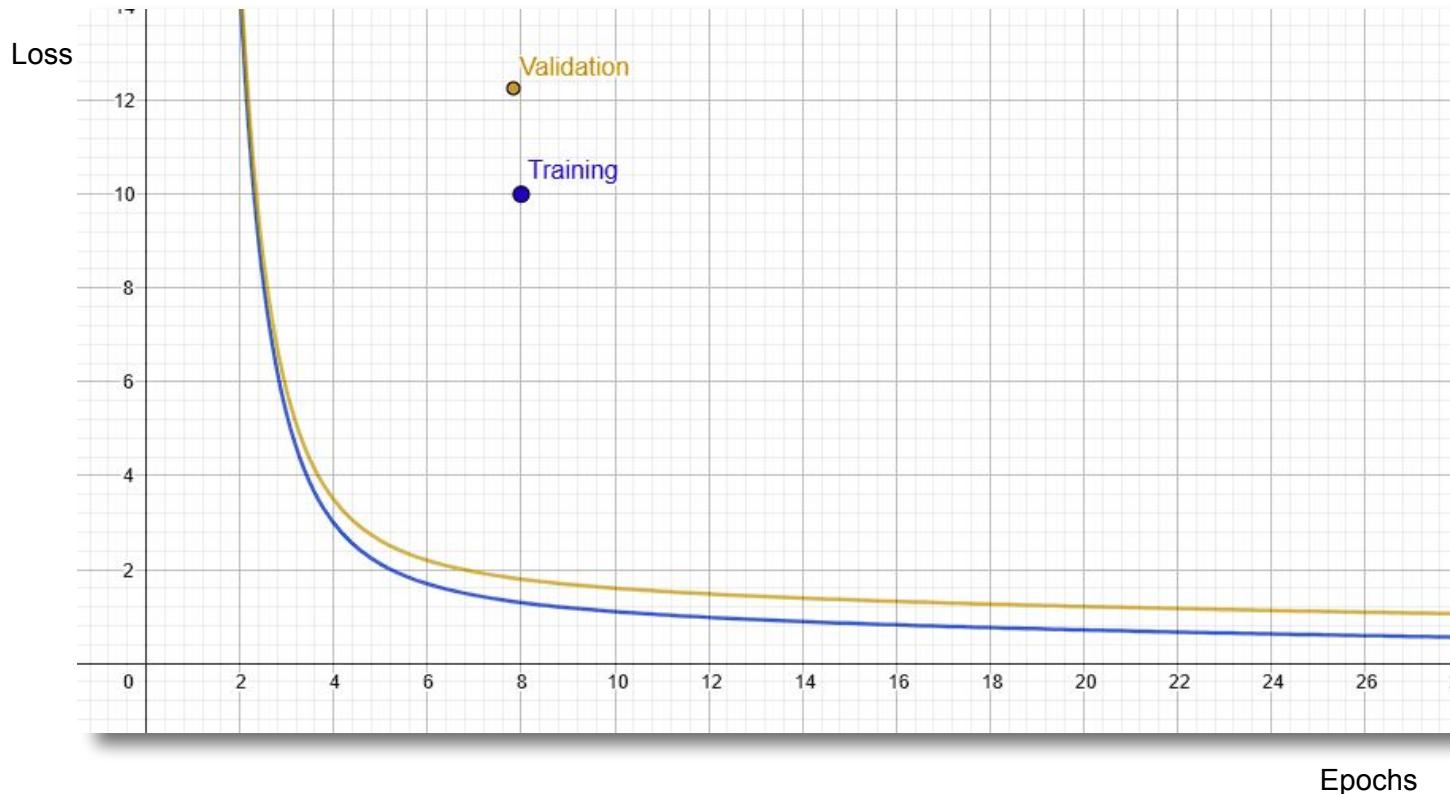
plot(loss_t, loss_v, epochs)
```

```
traing(x_train, y_train, alpha, epochs):
    w = # Initialize parameters
    for i in range(epochs):
        x,y = get_data(x_train, y_train, batch)
        dw = derivatives(x,y,w)
        w = Change_Parameters(w,dw, alpha)
        loss_training = Error(x,y,w,batch)
        loss_val = Error(x_val,y_val,w,batch)
        L_T.append(loss_training)
        L_V.append(loss_val)
    return L_T, L_V, w
```

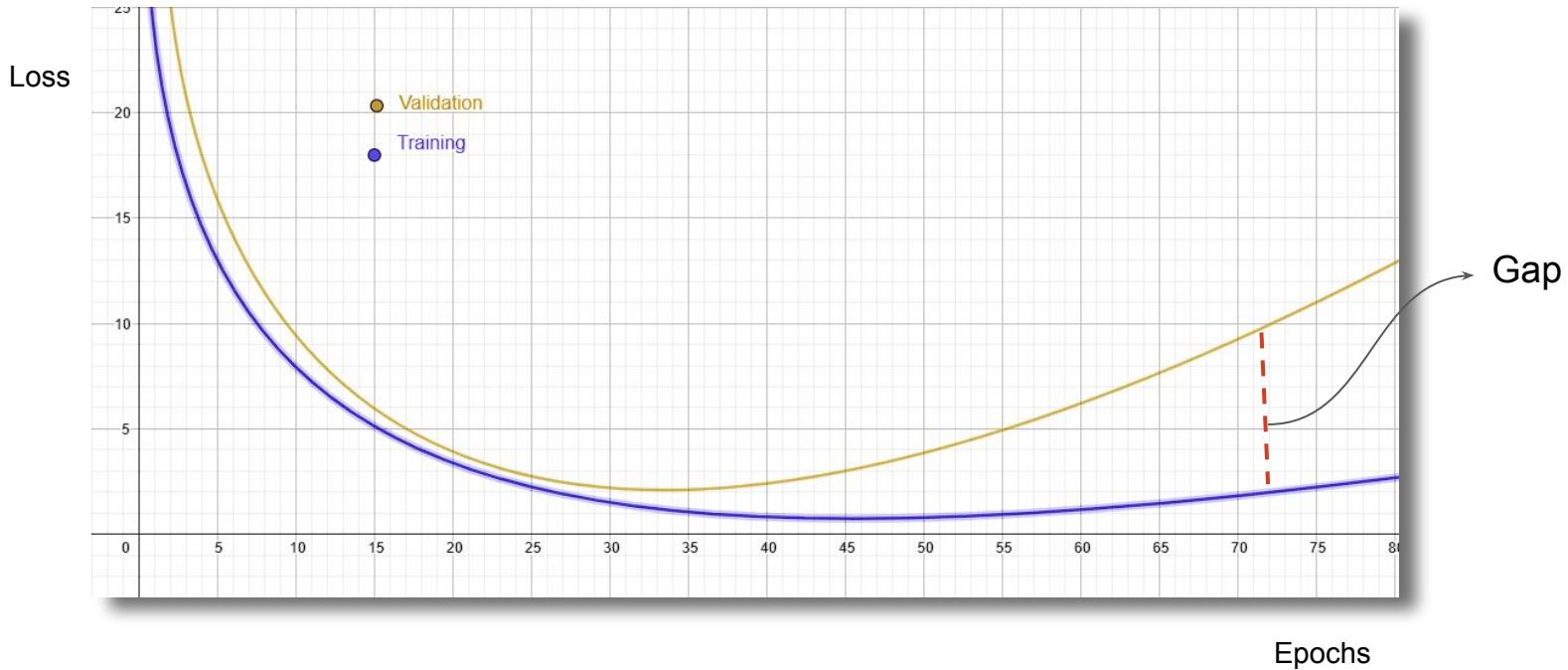
```
plot(loss_t, loss_v, epochs)
```

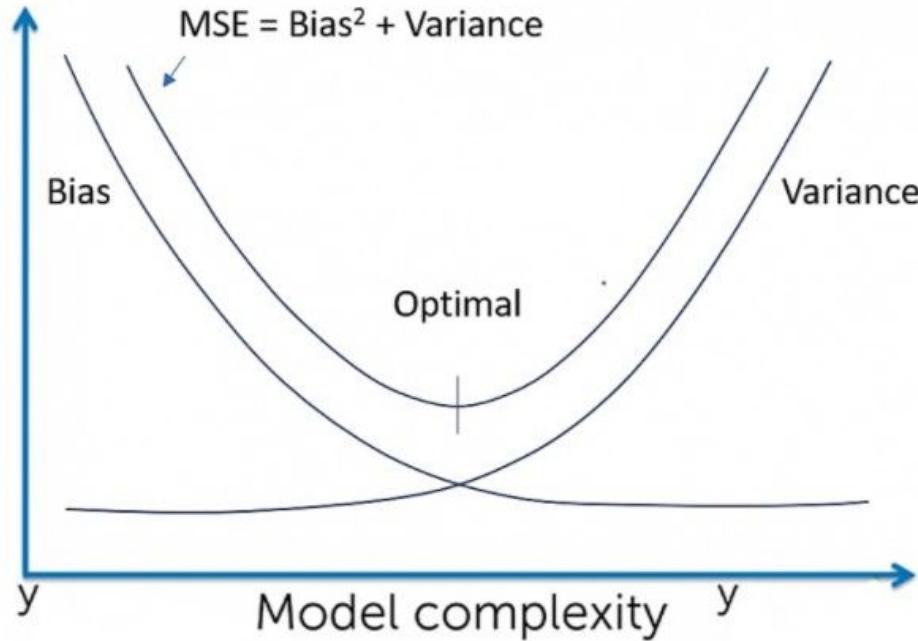
Underfitting



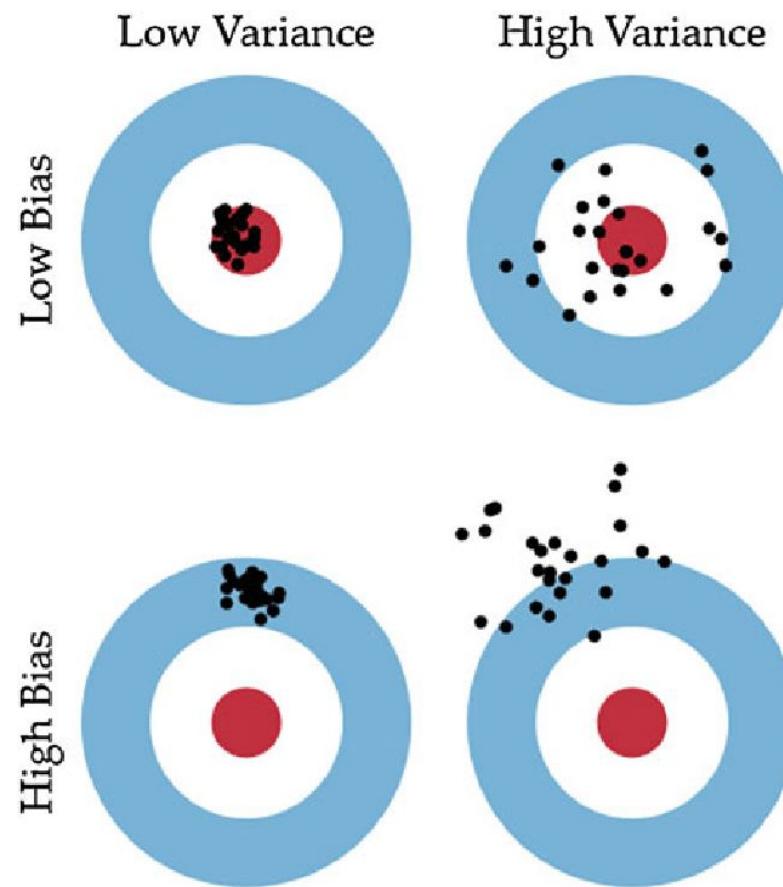


Overfitting





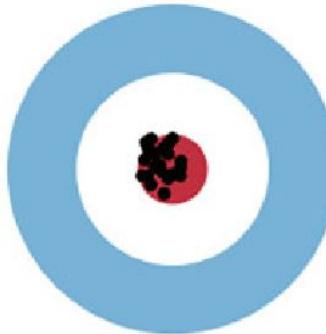
Source: <https://editor.analyticsvidhya.com/uploads/983161.png>



Low Variance

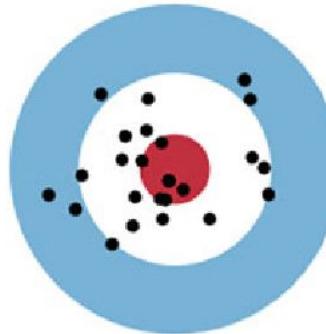
Low Bias

The model is consistent and accurate.



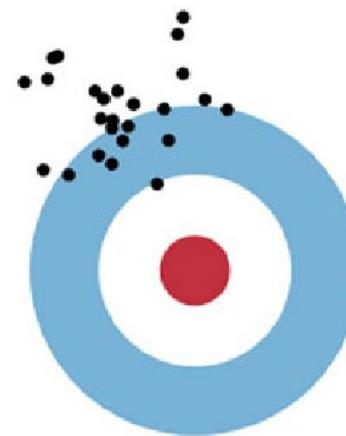
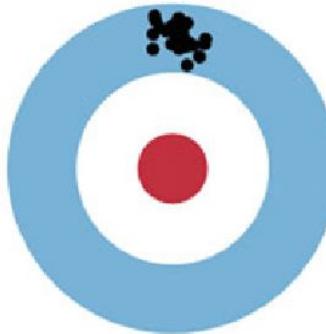
High Variance

The model is uncertain but accurate



High Bias

The model is consistent but inaccurate.

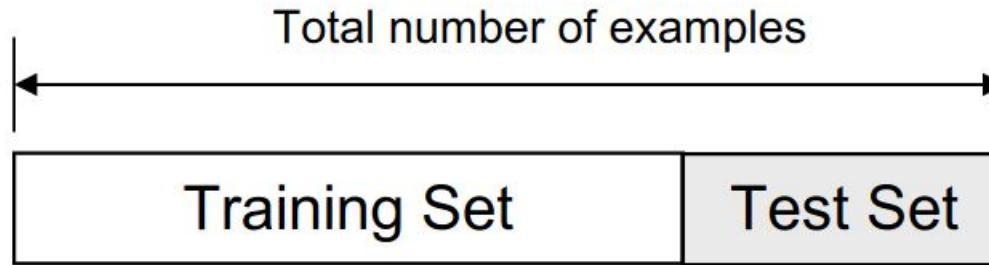


The model is uncertain and on average inaccurate

4

Method for evaluating the quality of generalization

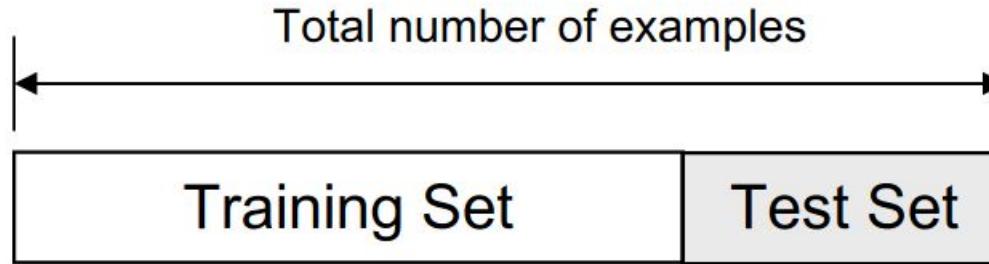
Holdout method



Problems

- What happens if we have a sparse database, or a very small one?
- Since it is only a training and testing experiment, the estimation of the error rate will be misleading in case we have a faulty split.

El método holdout



Problemas

- ¿Qué ocurre si tenemos una base de datos poco densa, o muy pequeña?
- Dado que sólo es un experimento de entrenamiento y prueba, la estimación de la tasa de error será engañosa en caso de que tengamos una división defectuosa.

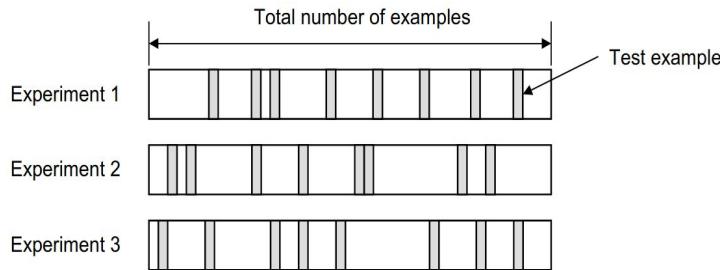
Resampling method

- Cross Validation
 - Random Subsampling
 - K-Fold Cross-Validation
 - Leave-one-out-Validation
- Bootstrap

Problems

- High computational cost

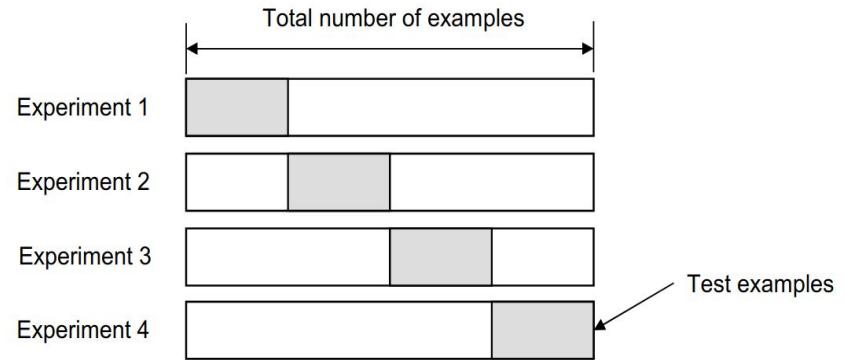
Random Subsampling



The error estimate is obtained:

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

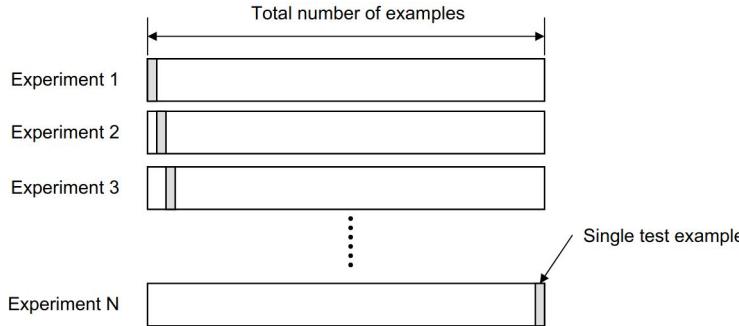
K-Fold Cross-Validation



The error estimate is obtained:

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

Leave-one-out-Validation



- Para un dataset con N ejemplos, se realizan N experimentos.
- Para cada experimento se utiliza $N-1$ ejemplos para entrenamiento y el restante para testing.

The error estimate is obtained:

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

How many folds are necessary?

A high number of folds

- The bias of the estimated error ratio will be small (the estimator will be more precise).
- The variance of the error estimator will be higher.
- The computational cost will be greater as we will have more experiments.

A low number of folds

- The number of experiments will be lower and therefore the computational cost.
- The variance of the estimator will be lower.
- The bias of the estimator will be higher (less precise).

In practice

- For large databases, even 3-fold Cross Validation will be sufficient.
- A common choice in K-fold Cross Validation is K=10.
- For small databases, we should use leave-one-out.

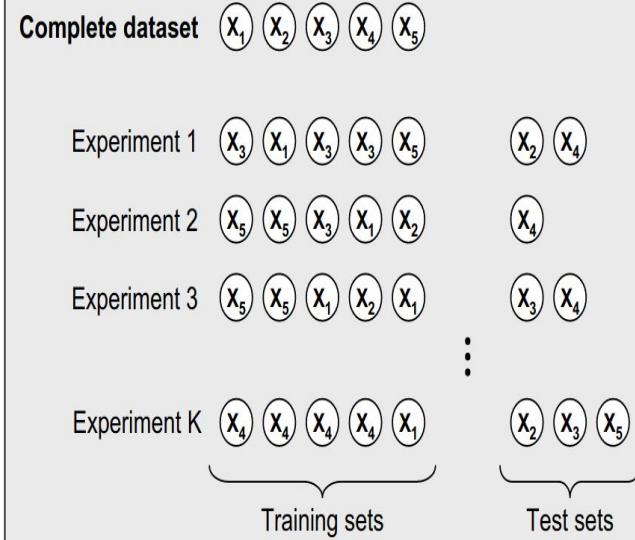


Método Bootstrap

- Select N elements (with replacement) and use them for training.
- The remaining are used for testing.
- Repeat the process for k folds.

The error estimate is obtained:

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$



Chapter summary

Regression models predict a real value from a set of predictor variables $R^k \rightarrow R$.

What should we do in the process?

- Explore the data to intuit a possible hypothesis.
- Clean the data in case of missing data
- Formulate the regression model (hypothesis)
- Normalize the data if necessary.
- Randomly split the data into training, validation, and testing.
- Evaluate the quality of generalization using cross-validation techniques or other methods.
- Use adjusted R^2 to check if the model explains the variability of the data.
- Conduct experiments to fine-tune the hyperparameters.

Teamwork Time



Classification

Cristian López Del Alamo
clopezd@utec.edu.pe
GINIA - Research group

2022

Objective: In this class, the objective is to understand what classification is from the perspective of machine learning, what binary classification is and how it works, and finally to implement logistic regression.

Classification



- Aprendizaje Supervisado.
- Proceso de Entrenamiento
- Datos de Entrenamiento.



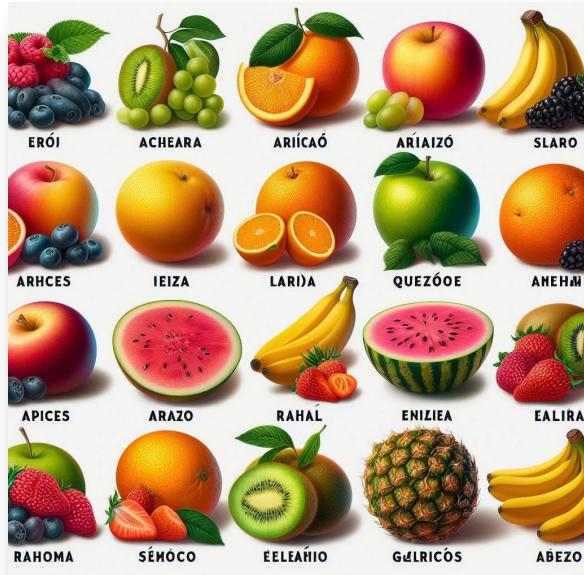
- Proceso de Testing



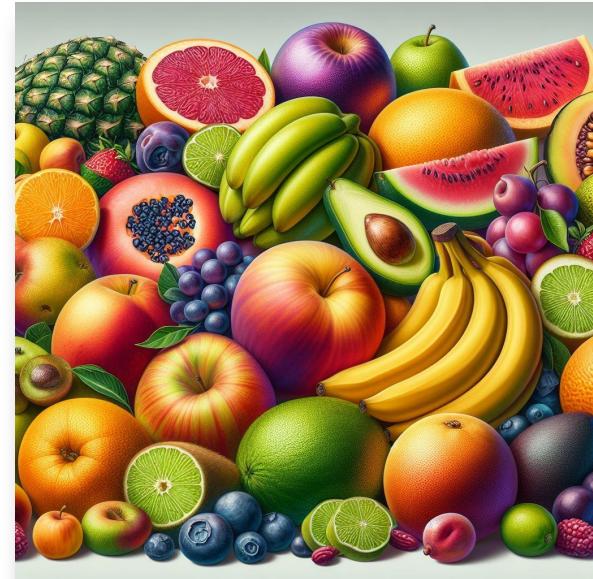
- Proceso de Inferencia

Classification

TrainingData = {(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)}

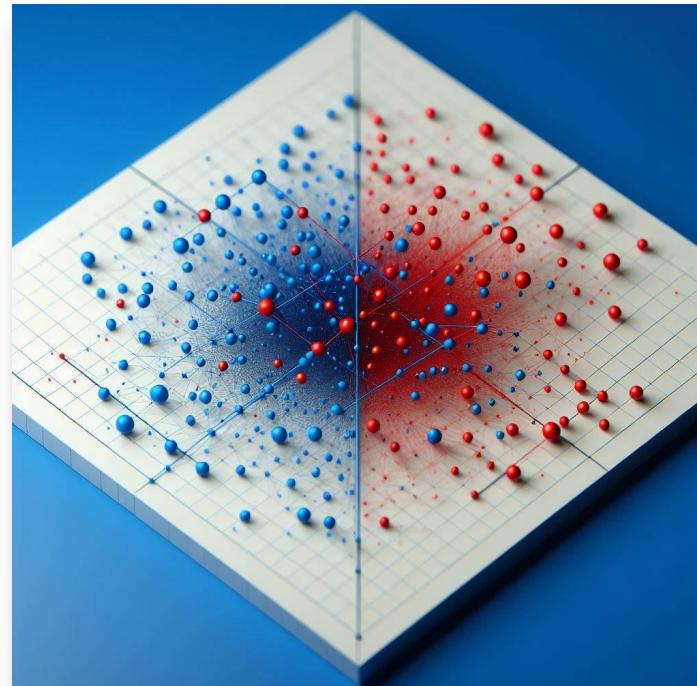


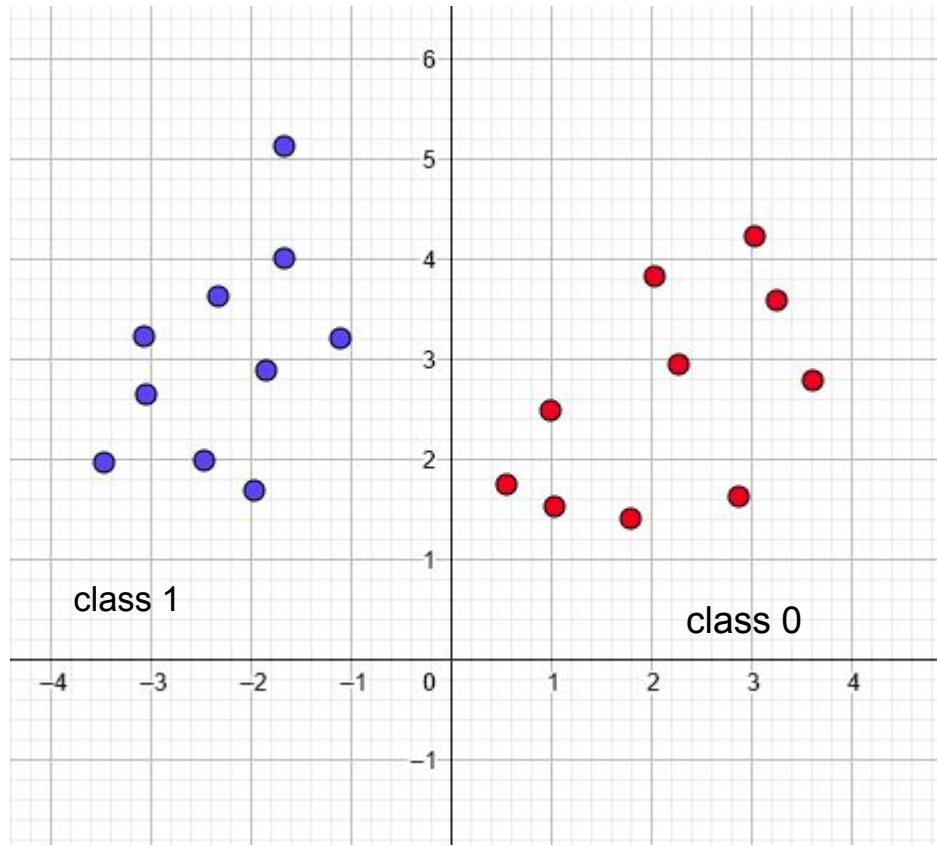
TestData = { x_1, x_2, \dots, x_n }



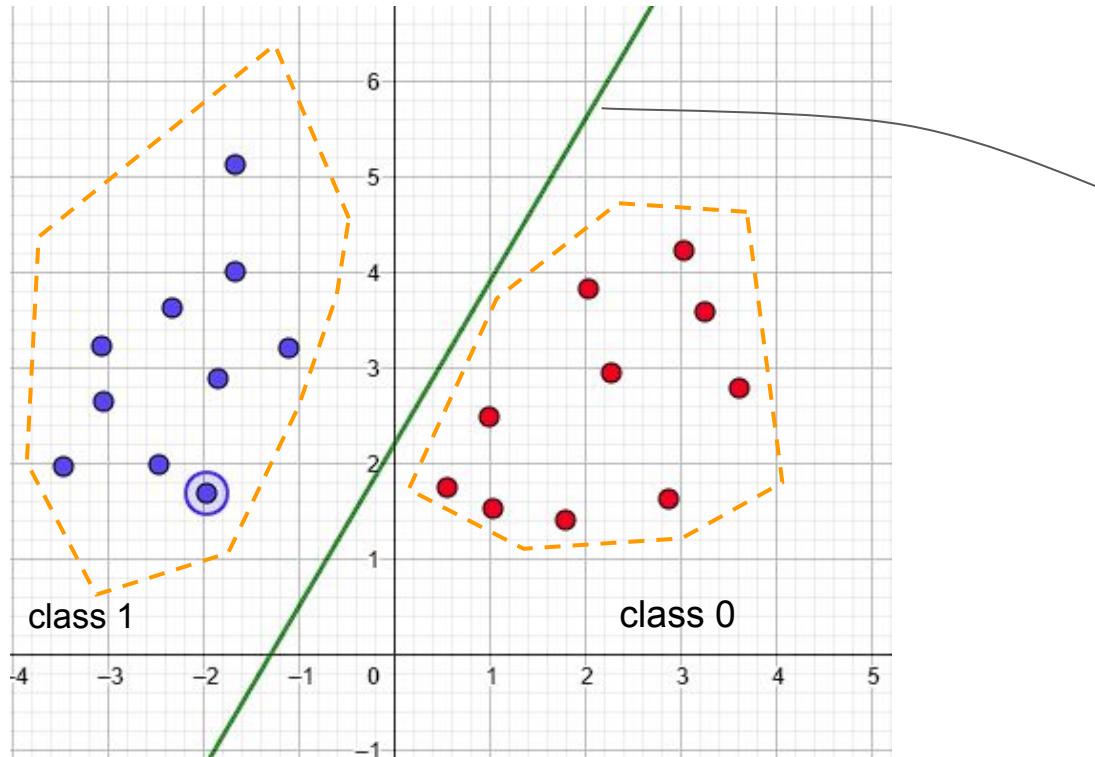
Binary Classification

> Reinventa el mundo <



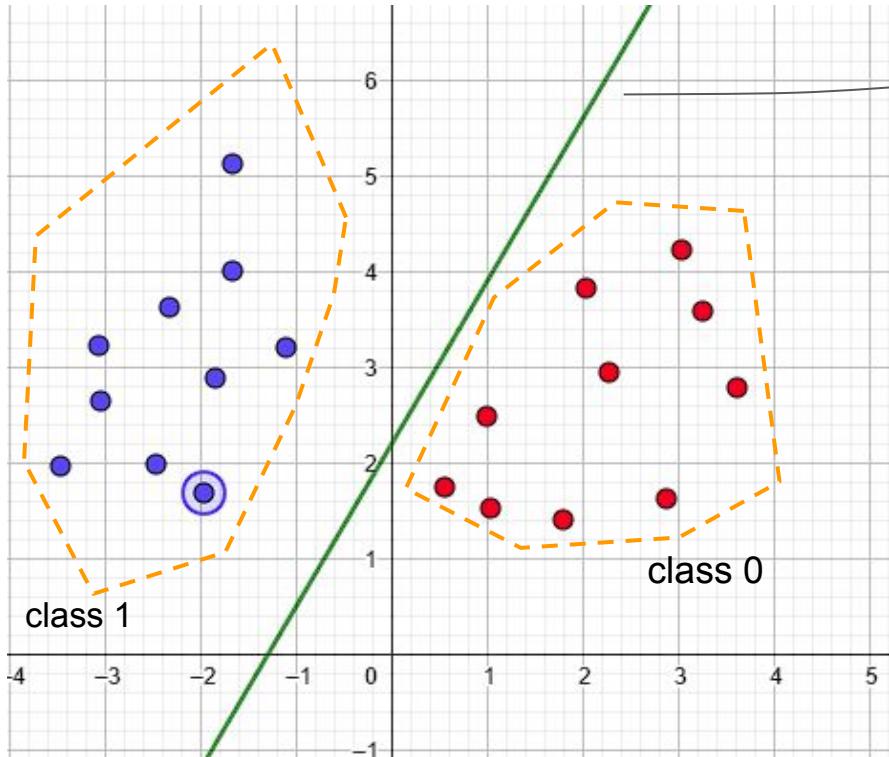


Binary Classification



$$h(x_i) = w_0 + x_1w_1$$

Binary Classification

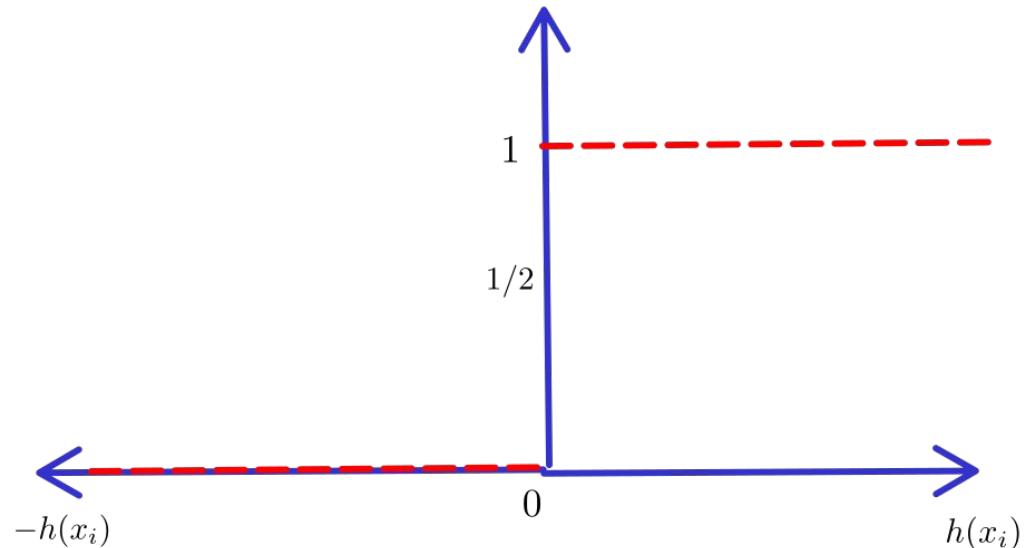


$$h(x_i) = w_0 + x_1 w_1$$

$$s(x_i) = \frac{1}{1 + e^{-h(x_i)}}$$

Logistic Function

$$s(x_i) = \frac{1}{1 + e^{-h(x_i)}}$$

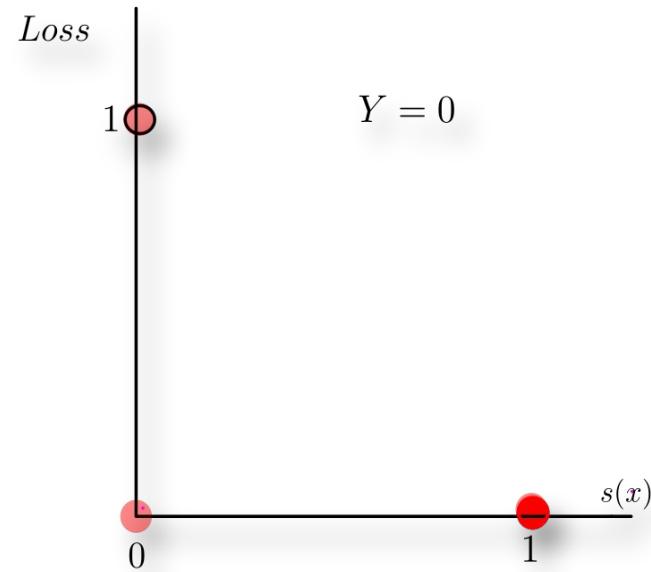
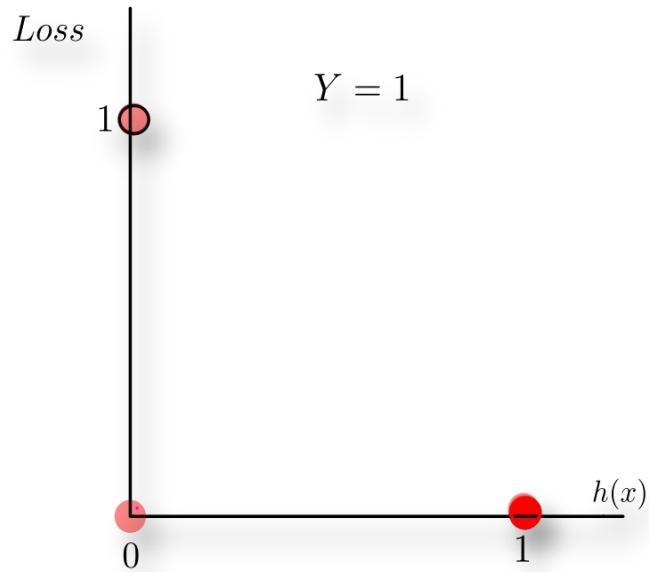


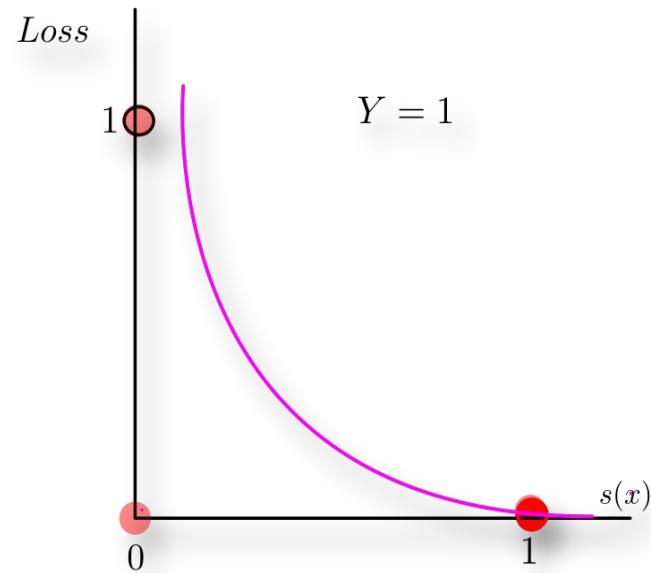
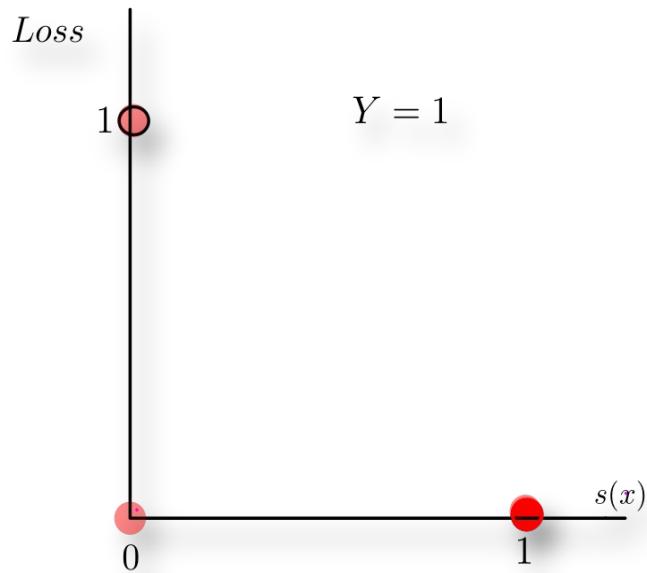
1. Hypothesis : $s(x_i) = \frac{1}{1 + e^{-h(x_i)}}$

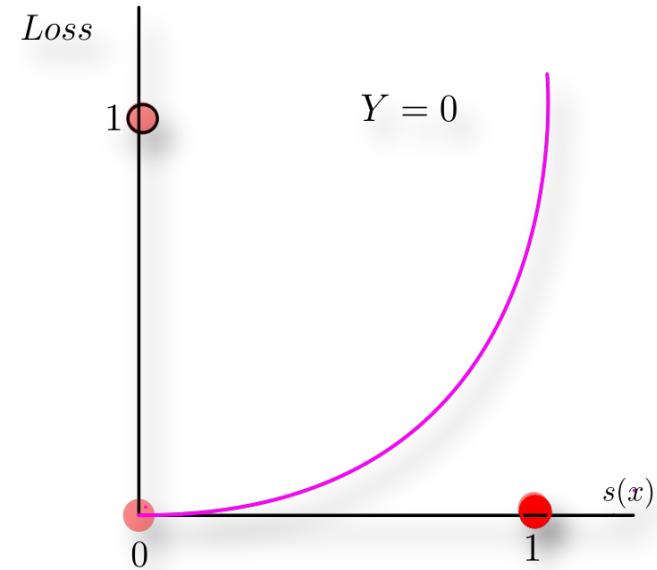
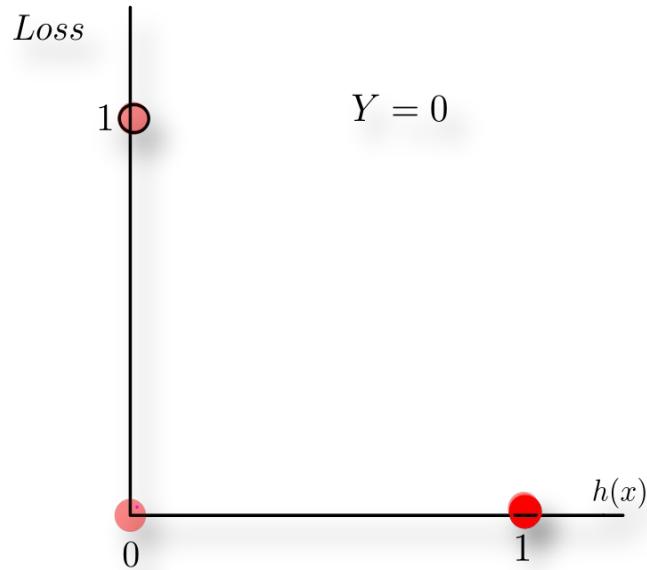
What would be the error function?



Cristian López Del Alamo







Si $y_i = 1 \rightarrow \log(s(x_i))$

Si $y_i = 0 \rightarrow \log(1 - s(x_i))$

2. Loss Function : $\mathcal{L} = - \sum_{i=1}^n (y_i \log(s(x_i)) + (1 - y_i) \log(1 - s(x_i)))$

1. Hypothesis : $s(x_i) = \frac{1}{1 + e^{-h(x_i)}}$

2. Loss Function : $\mathcal{L} = - \sum_{i=1}^n (y_i \log(s(x_i)) + (1 - y_i) \log(1 - s(x_i)))$

3. Derivadas : $\frac{\partial L}{w_j} = \frac{1}{n} \sum_{i=1}^n (y_i - s(x_i))(-x_{ij})$

```
1 def train(X, Y, epochs, alfa, lam):
2     np.random.seed(2001)
3     W = np.array([np.random.rand() for i in range(X.shape[1])])
4     L = Error(X,W,Y, lam)
5     loss = []
6     for i in range(epochs):
7         dW = derivada(X, W, Y, lam)
8         W = update(W, dW, alfa)
9         L = Error(X, W, Y, lam)
10        loss.append(L)
11        if i%10000==0:
12            print(L)
13    return W, loss
```

Chapter summary

Binary classification is a type of supervised learning problem where a target variable with two possible values is predicted.

What should we do in the process?

- Explore the data to intuit a possible hypothesis.
- Clean the data in case of missing data
- We formulate the binary classification model.
- Normalize the data if necessary.
- Randomly split the data into training, validation, and testing.
- Evaluate the quality of generalization using cross-validation techniques or other methods.
- How do we measure the quality of the cluster? ... **Coming soon**



UTEC

UNIVERSIDAD DE INGENIERIA
Y TECNOLOGIA