

Primeros pasos en Machine Learning y regresión lineal

En breve comenzamos...



NEOLAND



Primeros pasos en Machine Learning y regresión lineal

NEOLAND

PONENTE



Josep Miquel Porcar

Data Scientist, matemático y estadístico con varios años de experiencia en proyectos de investigación de Data Science.

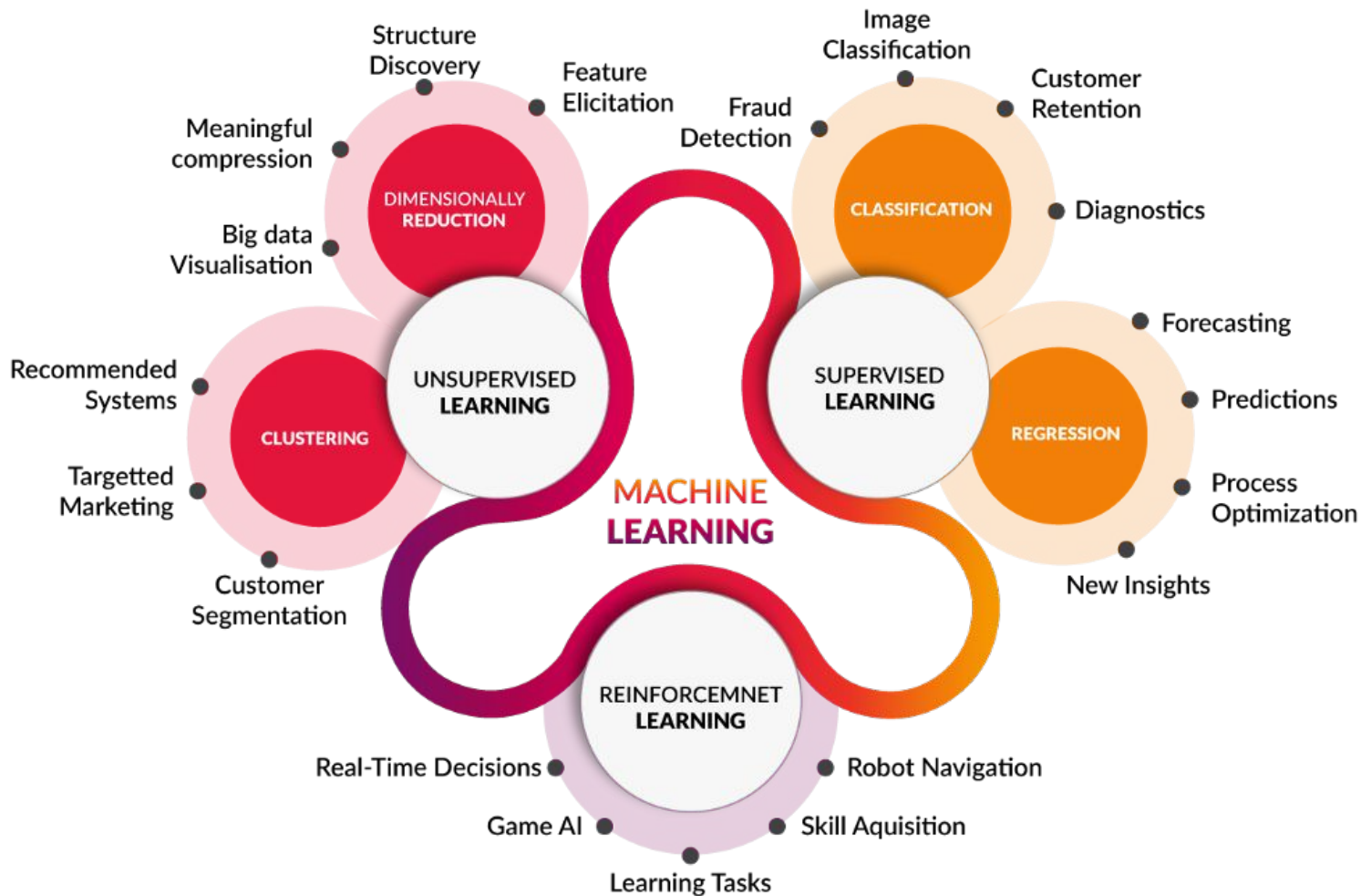
Actualmente es **Head Teacher del Data Science Bootcamp de Barcelona en NEOLAND.**

DATA SCIENCE

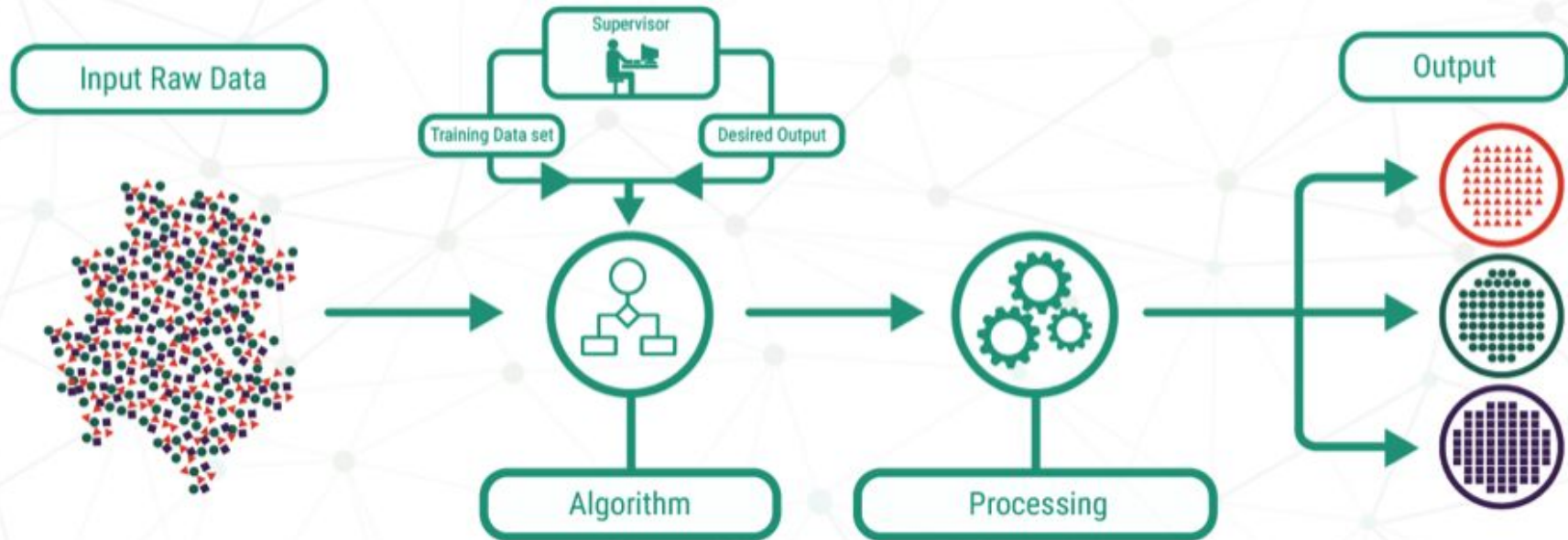
¿Qué veremos?

1. Statistical learning
2. Types of algorithms
3. Supervised learning
4. Linear Regression
5. Features types
6. Metrics
7. Polynomial regression
8. Validation

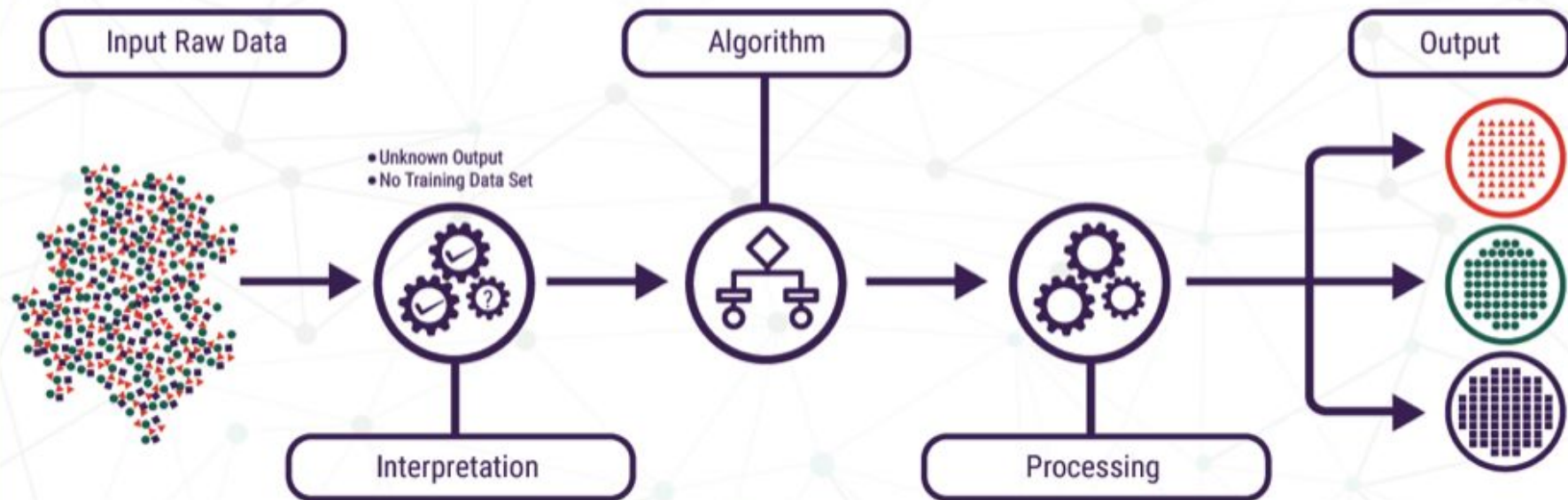




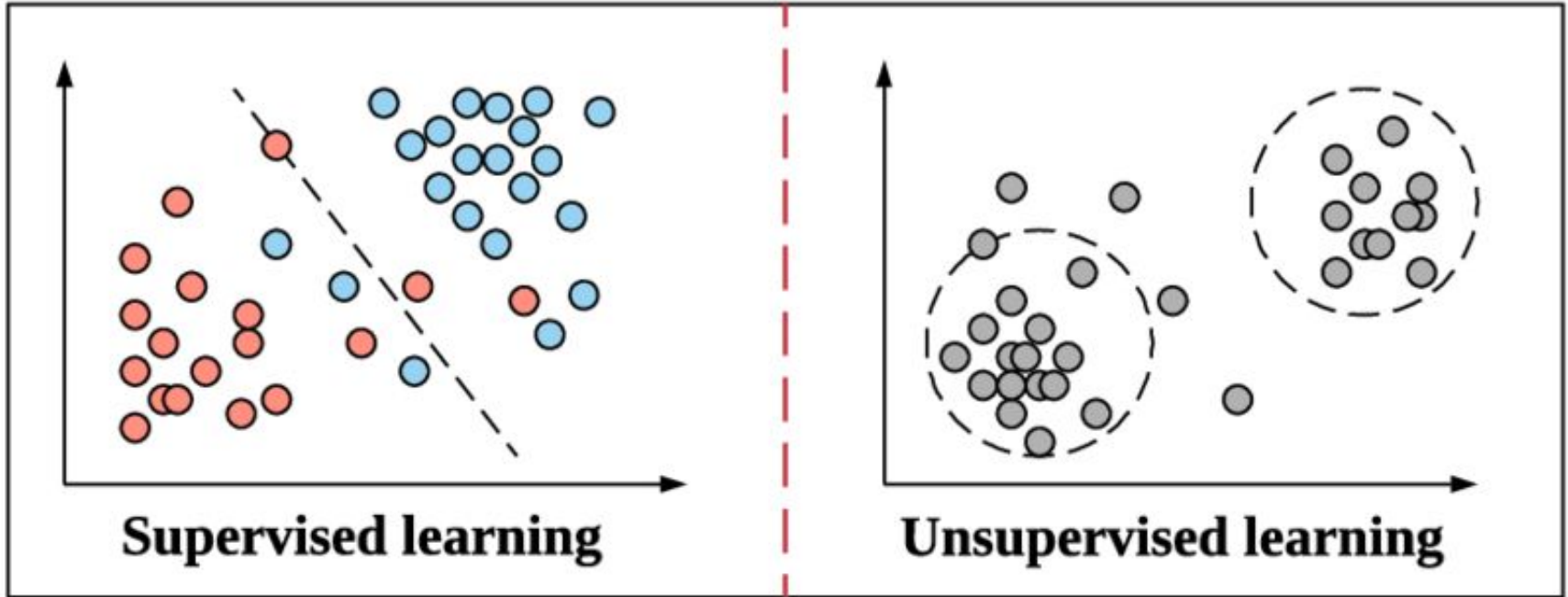
SUPERVISED LEARNING

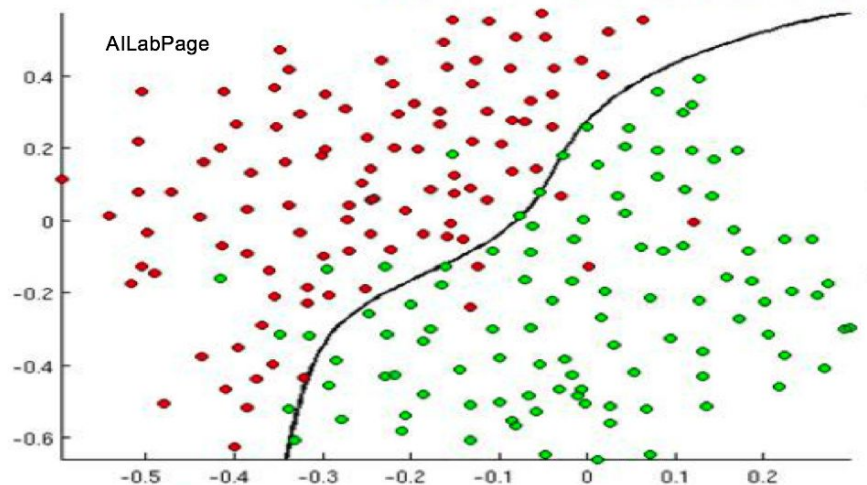
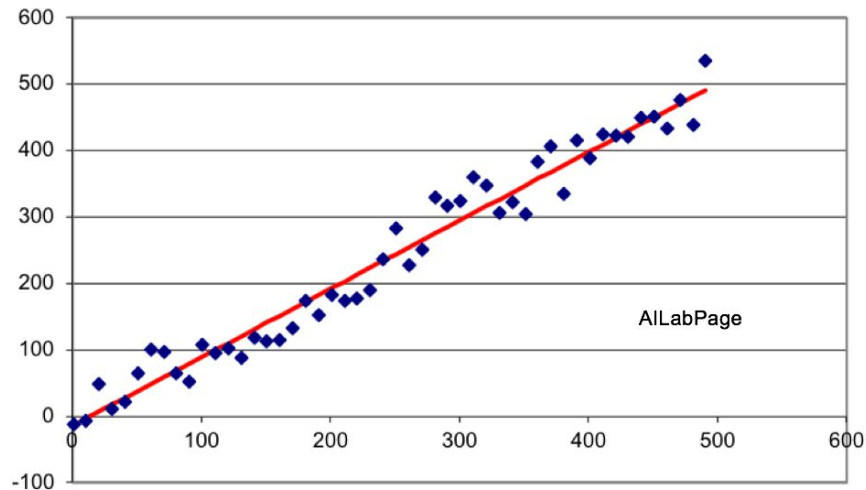


UNSUPERVISED LEARNING



Types of algorithms





Regression

The system attempts to predict a value for an input based on past data.

Example – 1. Temperature for tomorrow



Classification

In classification, predictions are made by classifying them into different categories.

Example – 1. Type of cancer 2. Cancer Y/N

Supervised learning

Supervised learning means learning from data:

- We have a quantitative outcome (regression) or categorical outcome (classification)
- We want to predict the *outcome* based on a set of *features* (supervised)
- We have a *training set*
- We build a prediction model for new unseen objects. The objective is to predict accurately

Vocabulary

Outcome, target, response: Usually denoted by Y

Features, columns, variables: Usually denoted by X (X is a vector of k features)

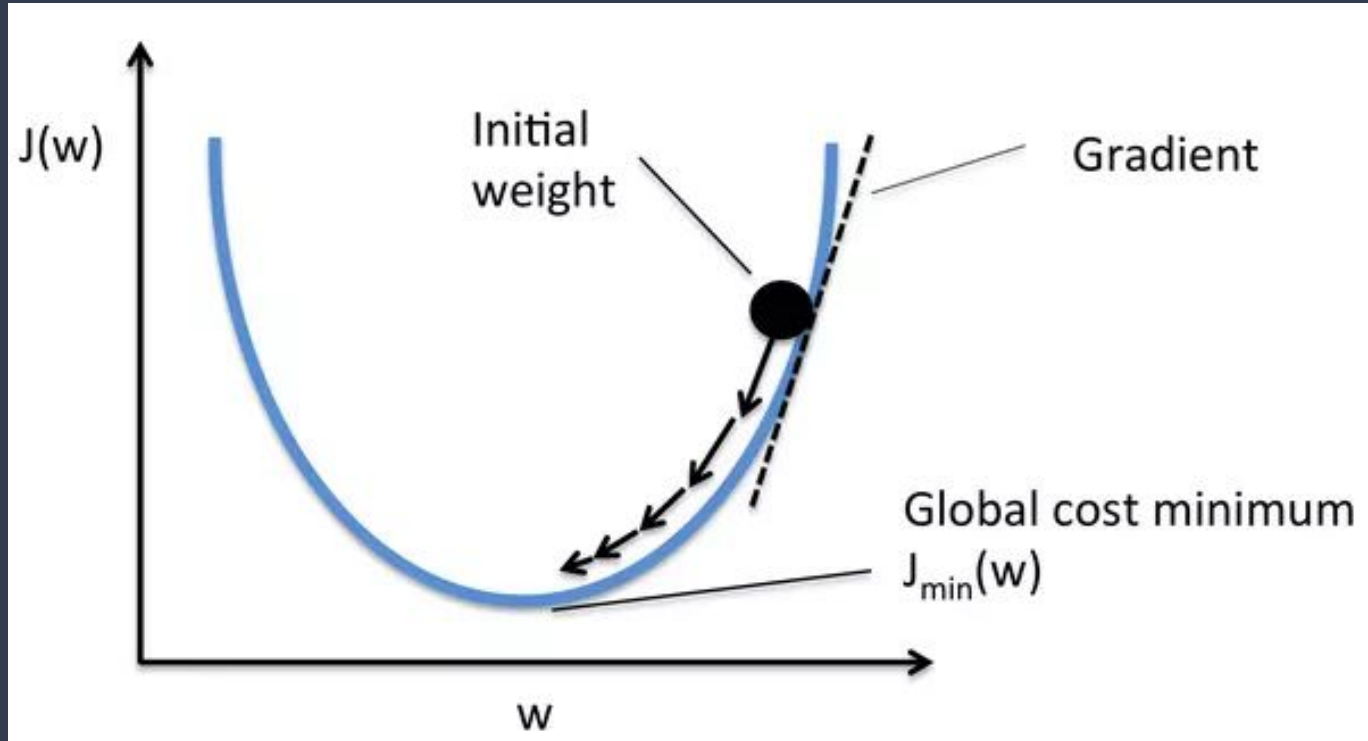
Training set: $(x_1, y_1), \dots, (x_n, y_n)$

Objective: get a good prediction of Y called $\hat{Y} = f(X)$.

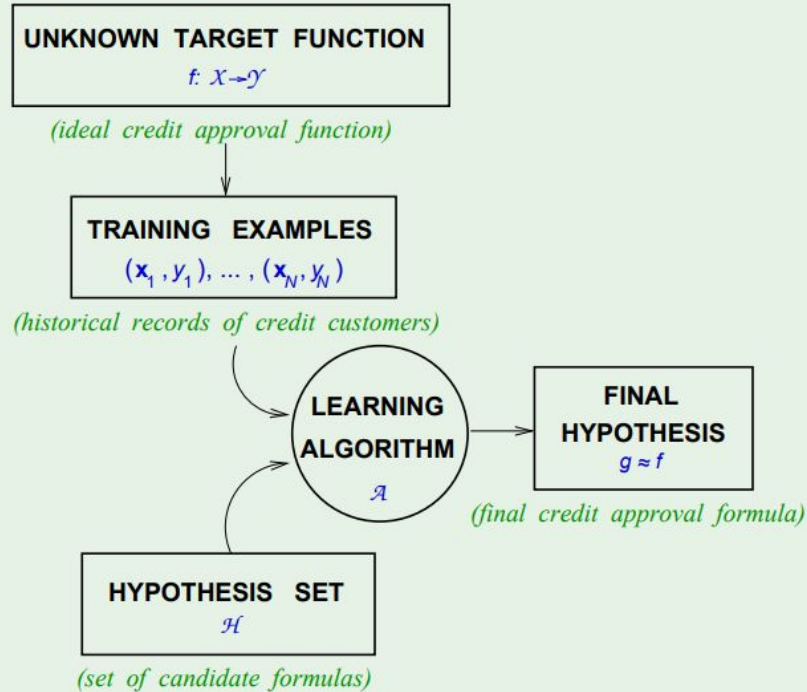
LOSS FUNCTION for penalizing errors (cost function)

Squared loss error $(Y - f(X))^2$

Loss function



Learning process



Linear Regression

A linear regression model assumes that Y is linear in the inputs X :

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$$

$$Y = f(X) + \varepsilon$$

Basic assumptions on errors: $\varepsilon \sim N(0, \sigma^2)$

- Independent
- Mean zero
- Constant variance

Predicting new data

Given a new set of features (X_{nuevo}), we can predict the outcome as:

$$\hat{Y}_{\text{nuevo}} = \hat{\beta} X_{\text{nuevo}} = \hat{\beta}_0 + \hat{\beta}_1 X_{1,\text{nuevo}} + \dots + \hat{\beta}_p X_{p,\text{nuevo}}$$

Features types in Linear regression

Quantitative - Continuous variables:

- Transformations: log, square root...
- Expansions: 2 , 3 , ...
- Interactions: $X_3 = X_1 * X_2$

Qualitative - Categorical variables:

- Dummy coding of the levels. 1 variable with K categories \rightarrow K dummy variables

Metrics

R²:

$$R^2 = 1 - \frac{\Sigma(y - \hat{y})^2}{\Sigma(y - \bar{y})^2}$$

RMSE:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

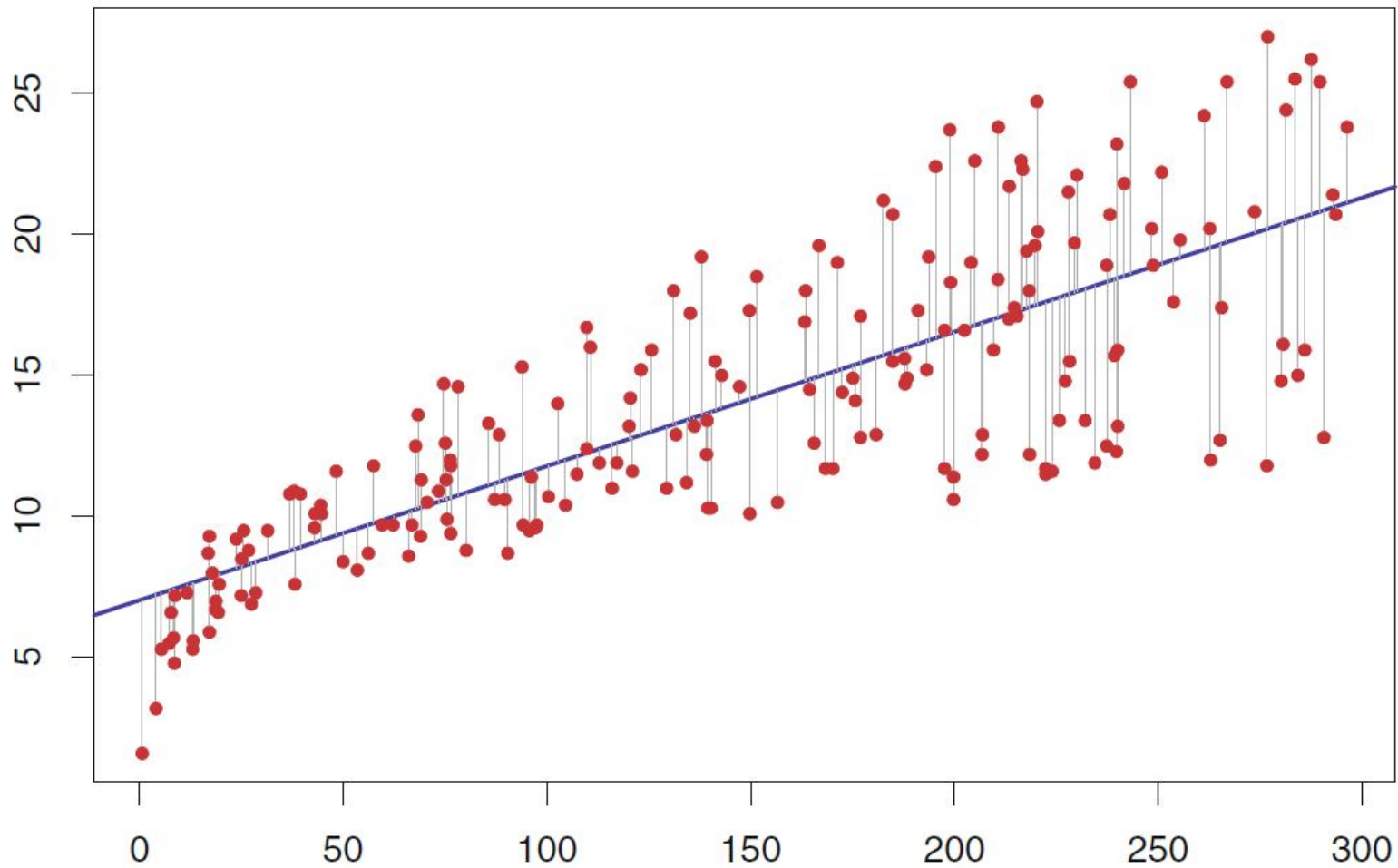
AIC:

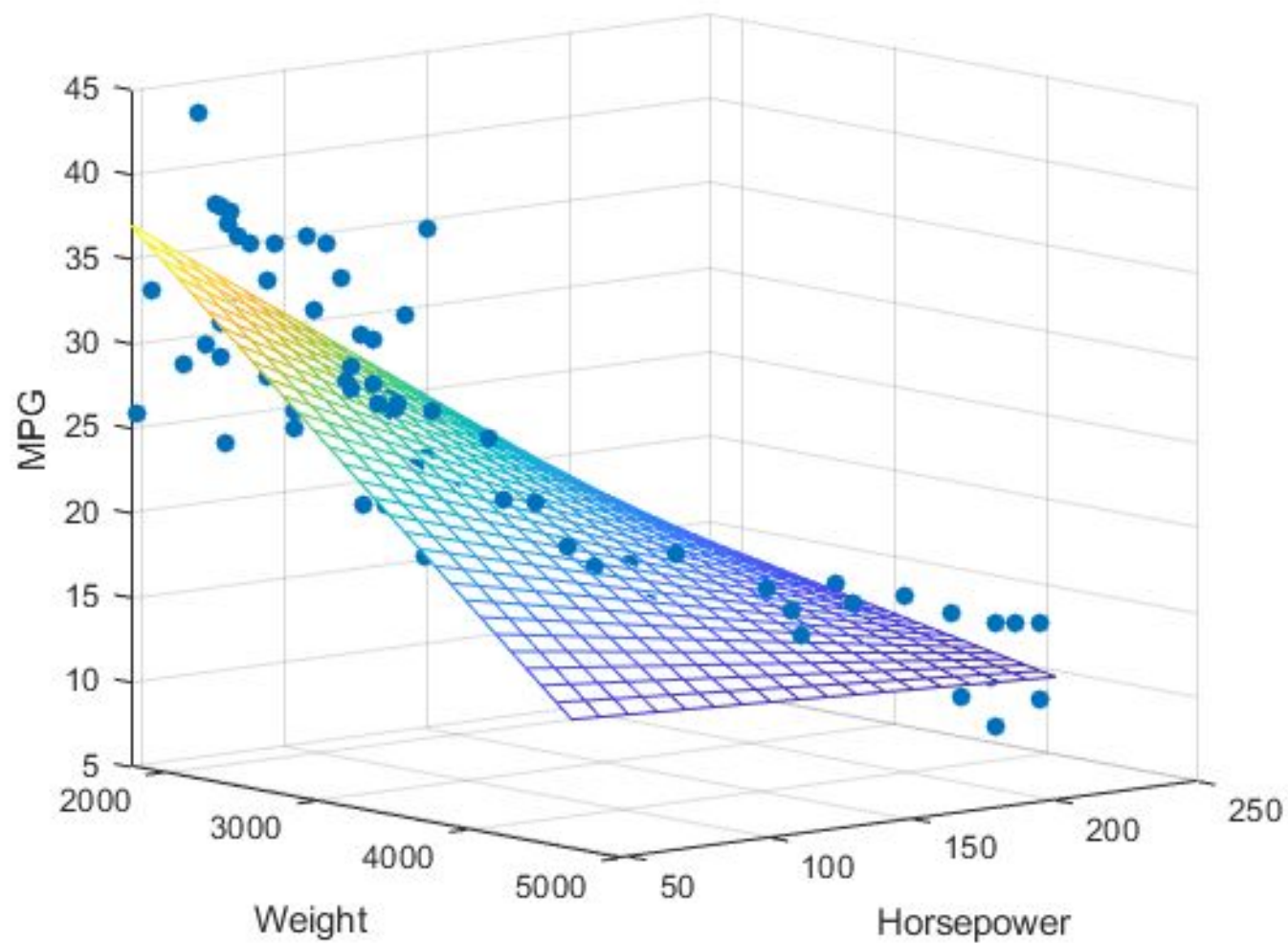
$$AIC = -2\log L + 2q$$

BIC:

$$BIC = 2\log(L) + q\log(N)$$

3. Linear Regression





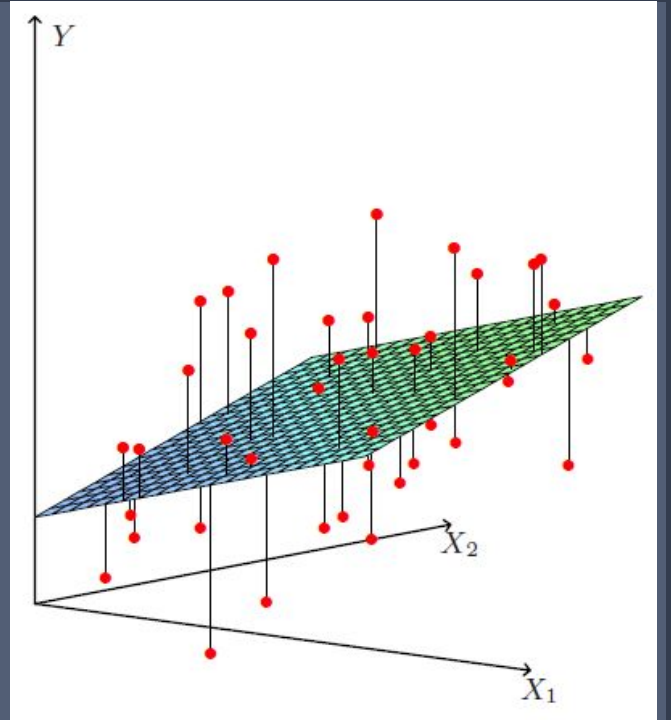
Estimation: Least Squares in LR

$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2.$$

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

$$\frac{\partial \text{RSS}}{\partial \beta} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta)$$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$



Interpreting estimators

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 D_1 + \hat{\beta}_3 D_2$$

X_1 is a continuous variable:

- sign
- size
- marginal effect

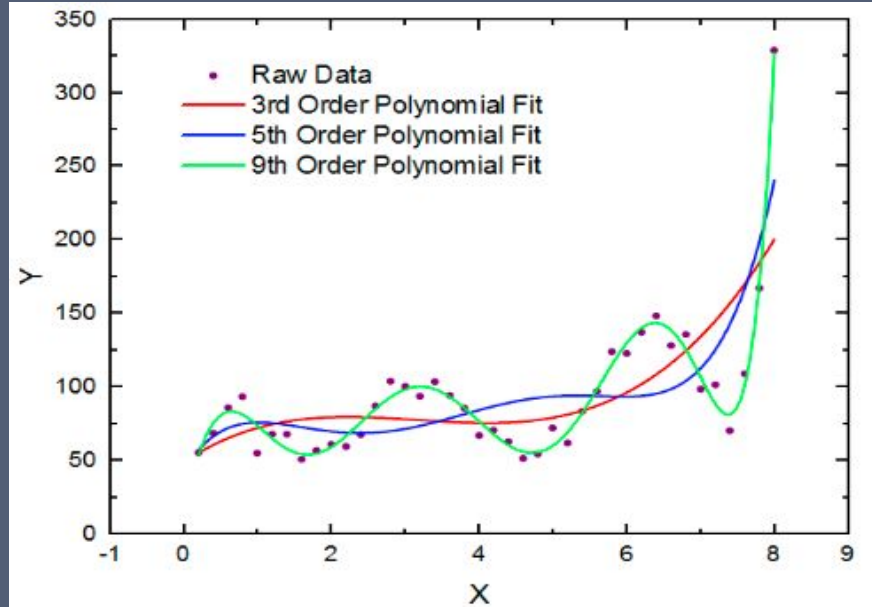
D_1 , D_2 are the dummies of a categorical variable with 3 levels:

- reference category D_3

Polynomial regression

$f(x)$ is a polynomial of order k :

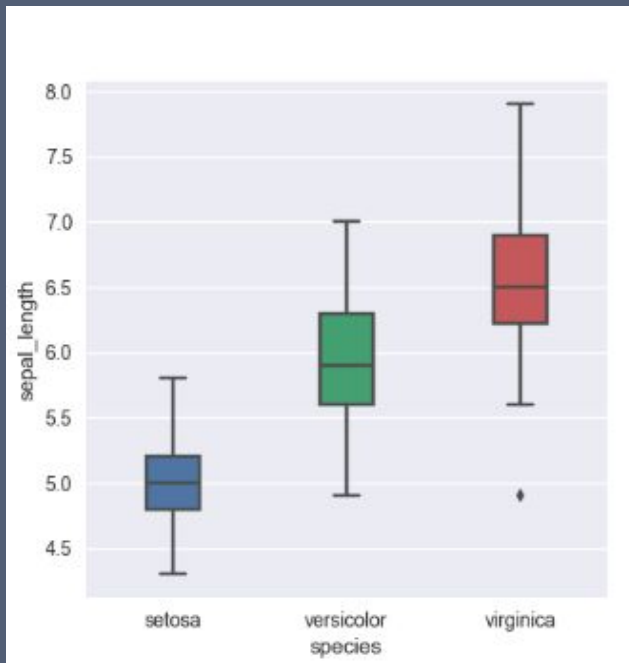
$$f(X) = \sum_{j=0}^k \beta_j X^j$$



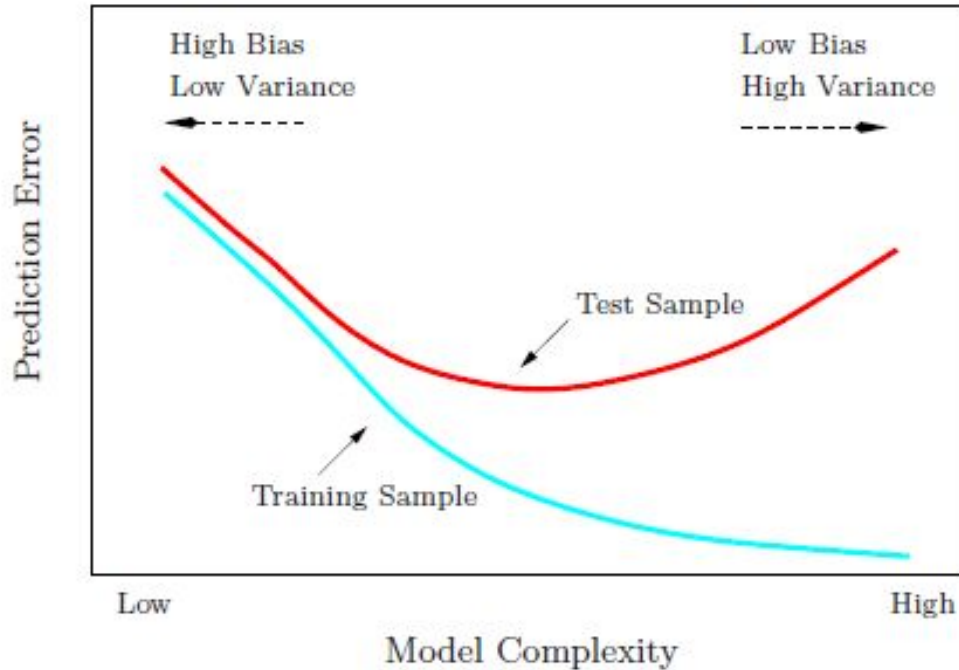
Outliers

In words of Hawkins, 1980:

"An outlier is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism"



Validation





Conectar

Enviar mensaje

Más...

Joseba Elcano Esparza · 2º

Data Scientist | Telecom engineer

Barcelona, Cataluña / Catalunya, España · 159 contactos ·

[Información de contacto](#)



Institute for Bioengineering
of Catalonia (IBEC)



Universitat Pompeu Fabra

Interés por oportunidades

Cargos de Research Assistant, Data Scientist y Data Analyst

[Ver todos los detalles](#)



Javier Fernández · 1er

Data Scientist en inAtlas

Barcelona, Cataluña, España · [Más de 500 contactos](#) ·

[Información de contacto](#)

Enviar mensaje

Más...



inAtlas



University of Havana



Gemma Labraña · 2º

Data Scientist at Martiderm®

Barcelona, Cataluña, España · Más de 500 contactos ·

[Información de contacto](#)

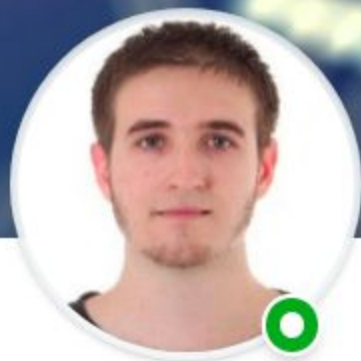
Conectar

 Enviar mensaje

Más...



MartiDerm



Ángel Molina Noguerras · 1er

Data Analyst en CEDEC S.A. Consultoría de Organización Estratégica

Barcelona y alrededores, España · [Más de 500 contactos](#) · [Información de contacto](#)

Enviar mensaje

Más...

cedec

CEDEC S.A. Consultoría de Organización Estratégica



NEOLAND