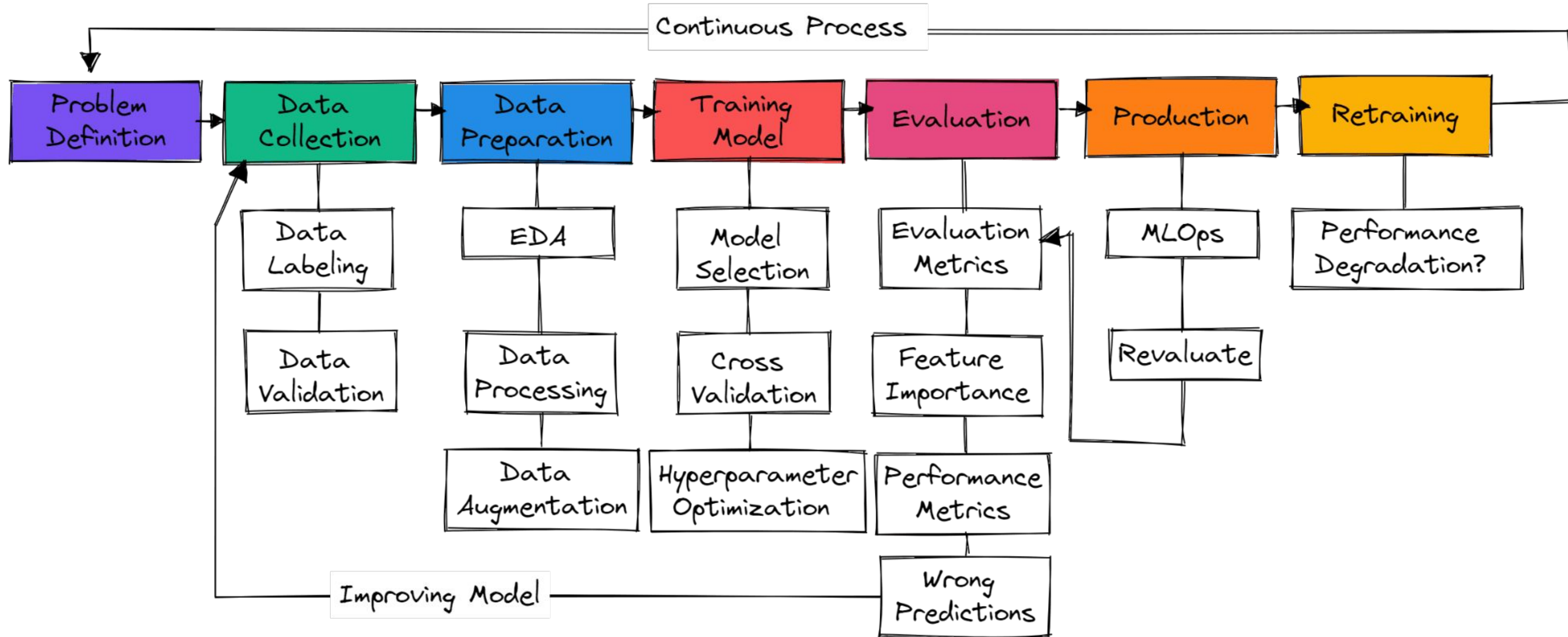


# Resumen ML



# ML Project



# Data Exploration

## IMPORTANCE OF DATA EXPLORATION WHEN PREPARING FOR DATA STORYTELLING





# Types of Machine Learning – At a Glance

## Supervised Learning

- Makes machine Learn explicitly
- Data with clearly defined output is given
- Direct feedback is given
- Predicts outcome/future
- Resolves classification and regression problems



## Unsupervised Learning

- Machine understands the data (Identifies patterns/structures)
- Evaluation is qualitative or indirect
- Does not predict/find anything specific

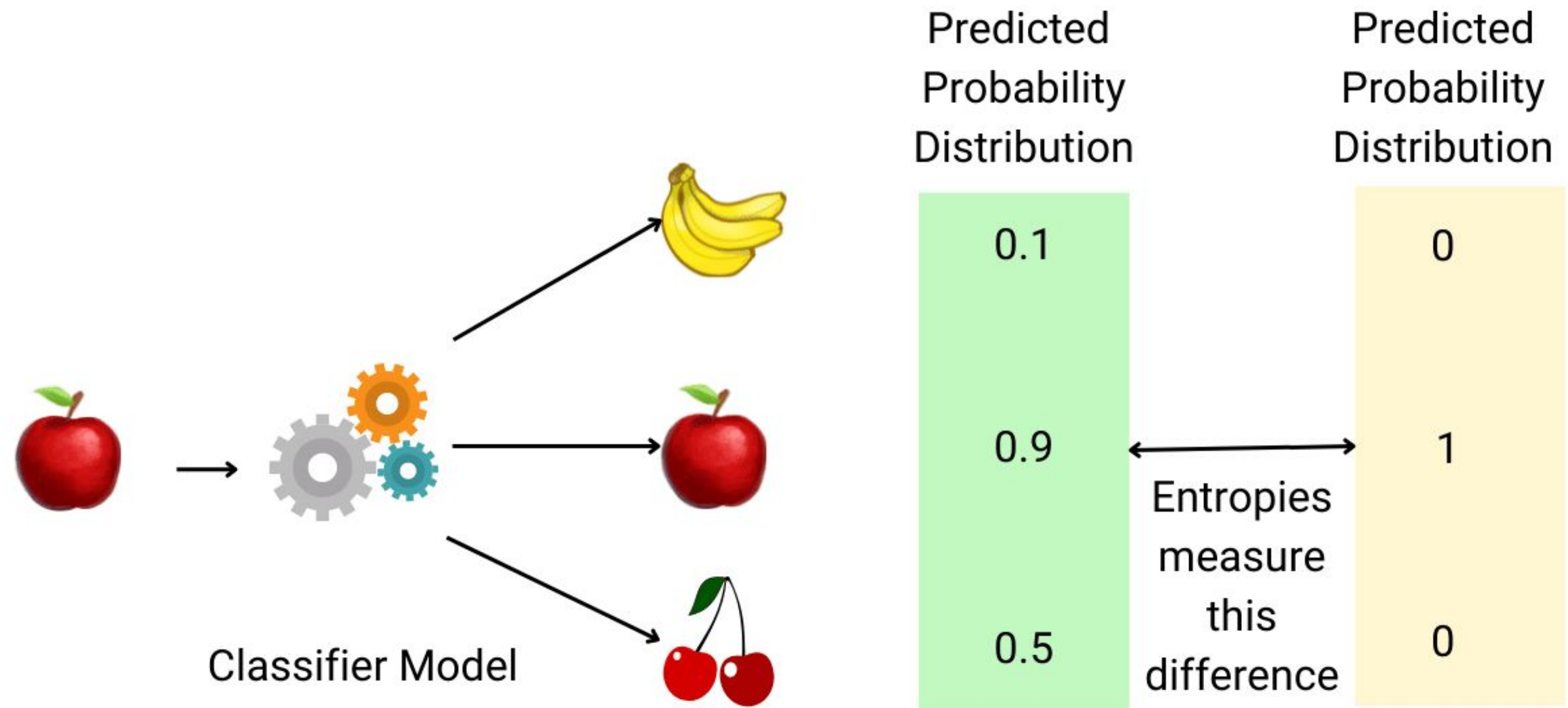


## Reinforcement Learning

- An approach to AI
- Reward based learning
- Learning from +ve & -ve reinforcement
- Machine Learns how to act in a certain environment
- To maximize rewards

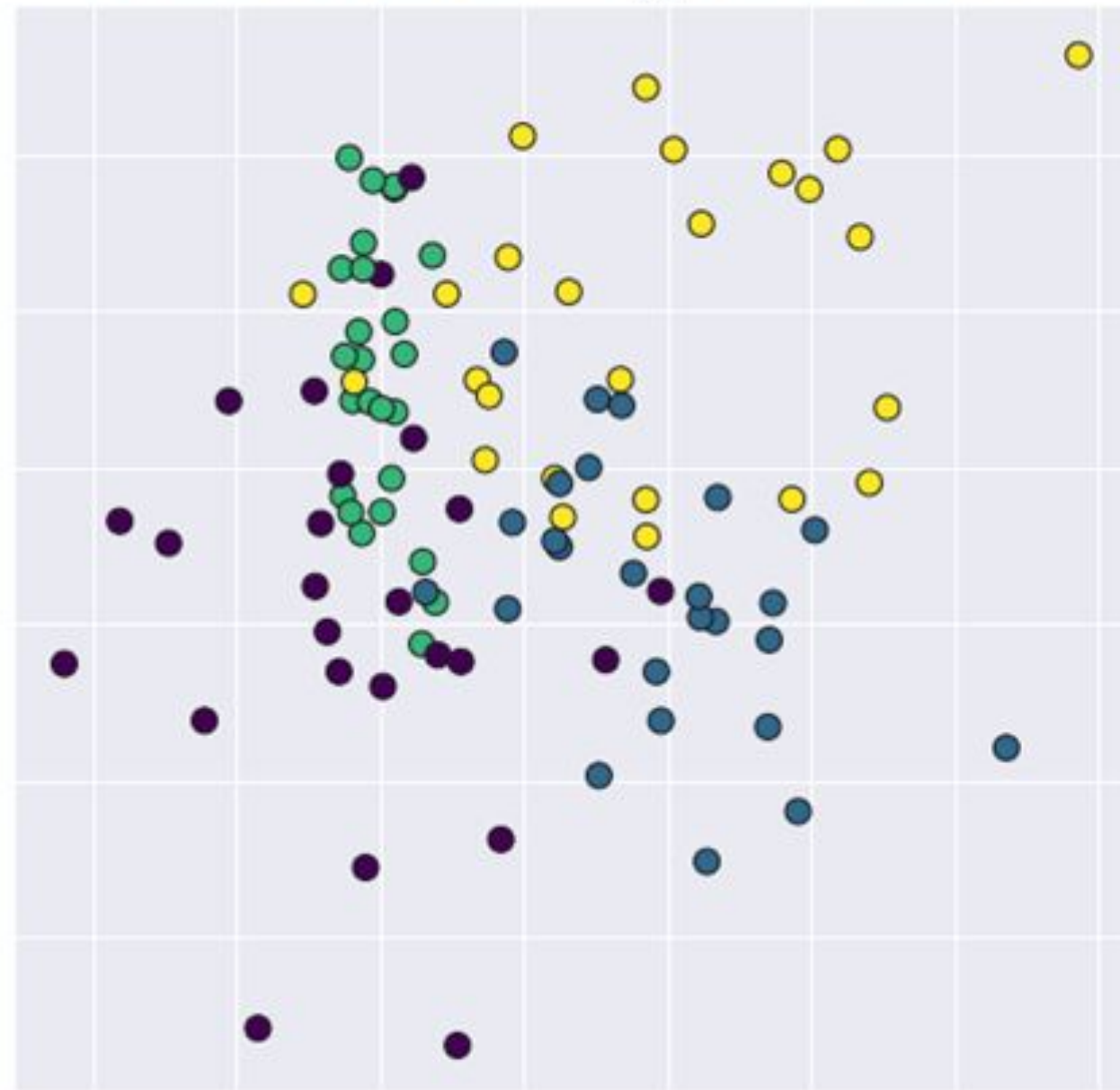


# Supervised-Classification



# Supervised-Multi Classification

Dataset containing four classes



## One-vs-Rest

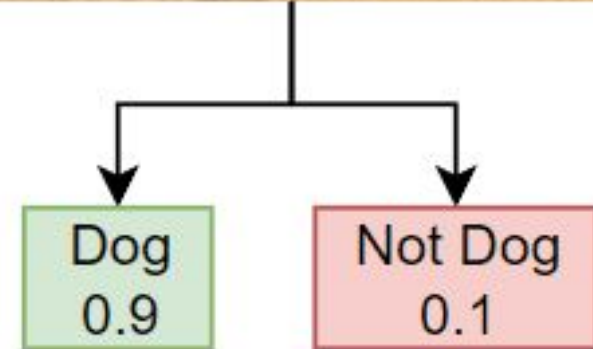


Train one binary model for **each class**, with remaining classes treated as aggregates. Also known as **One-vs-All**. Commonly used with logistic regression.

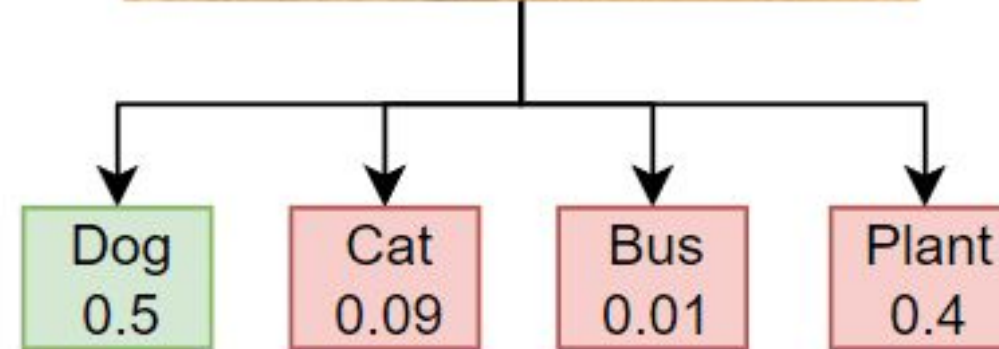


# Supervised-Classification

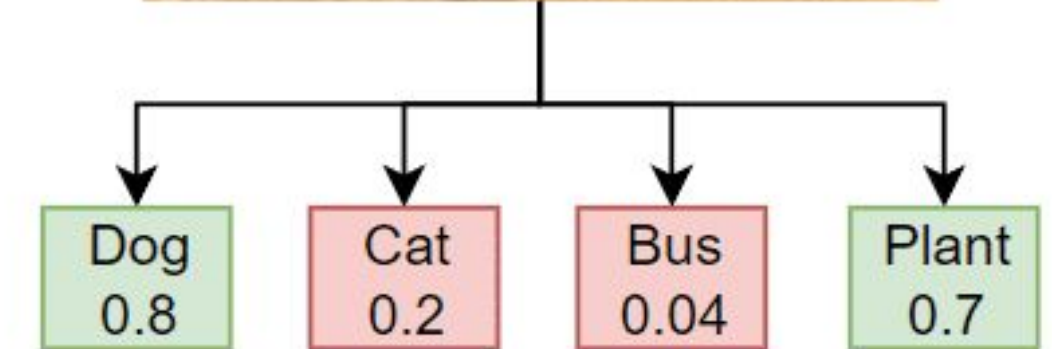
Binary Classification



Multiclass Classification

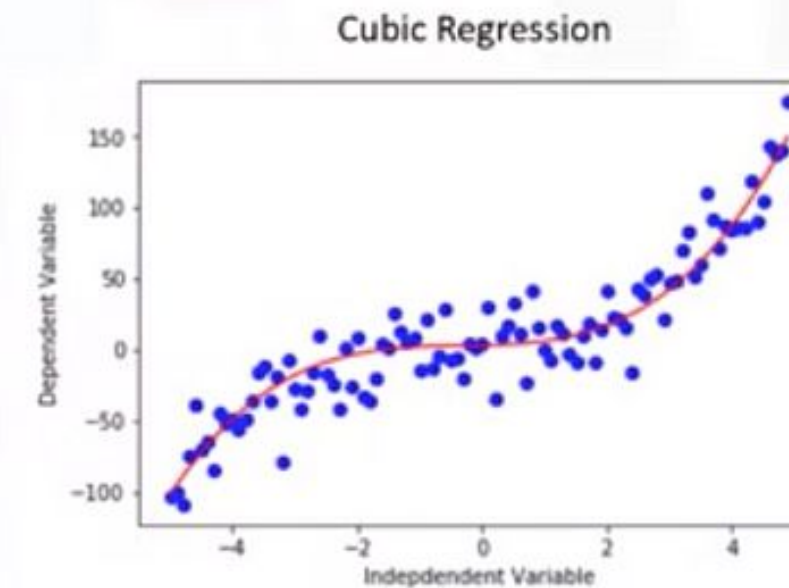
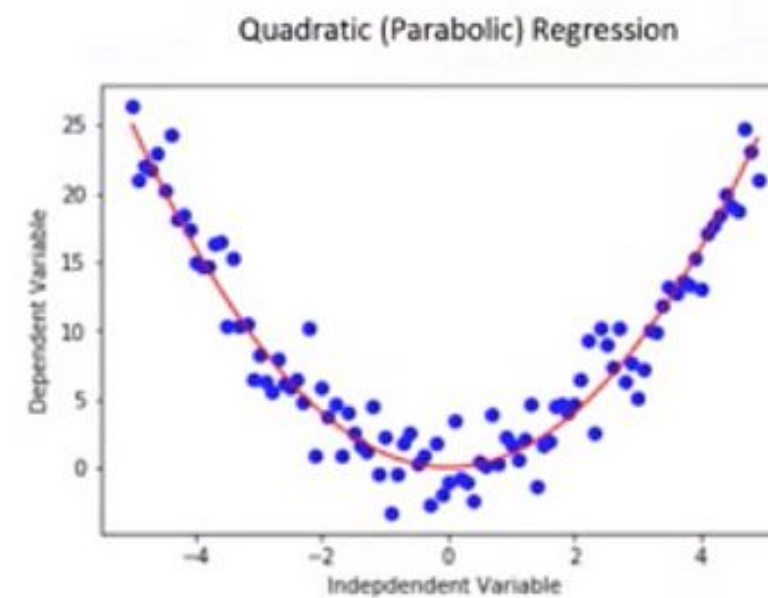
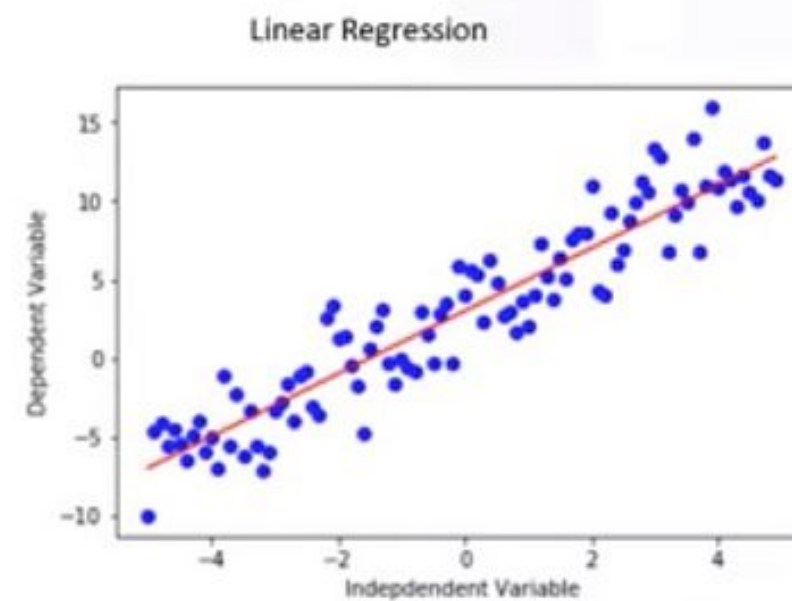


Multilabel Classification

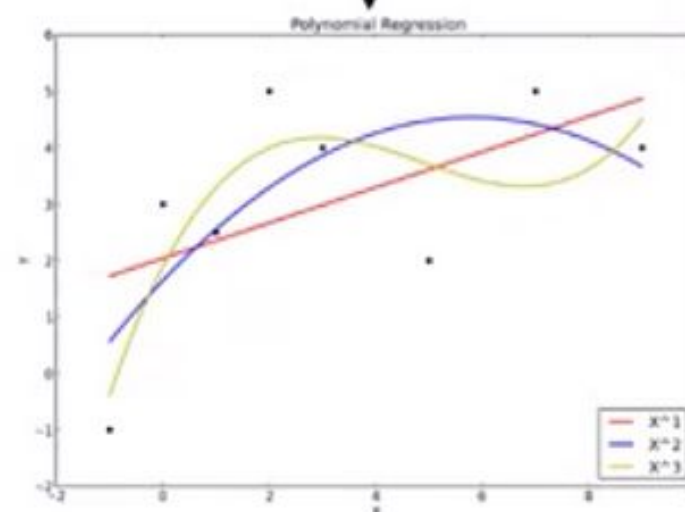


# Supervised-Regression

## Different types of regression

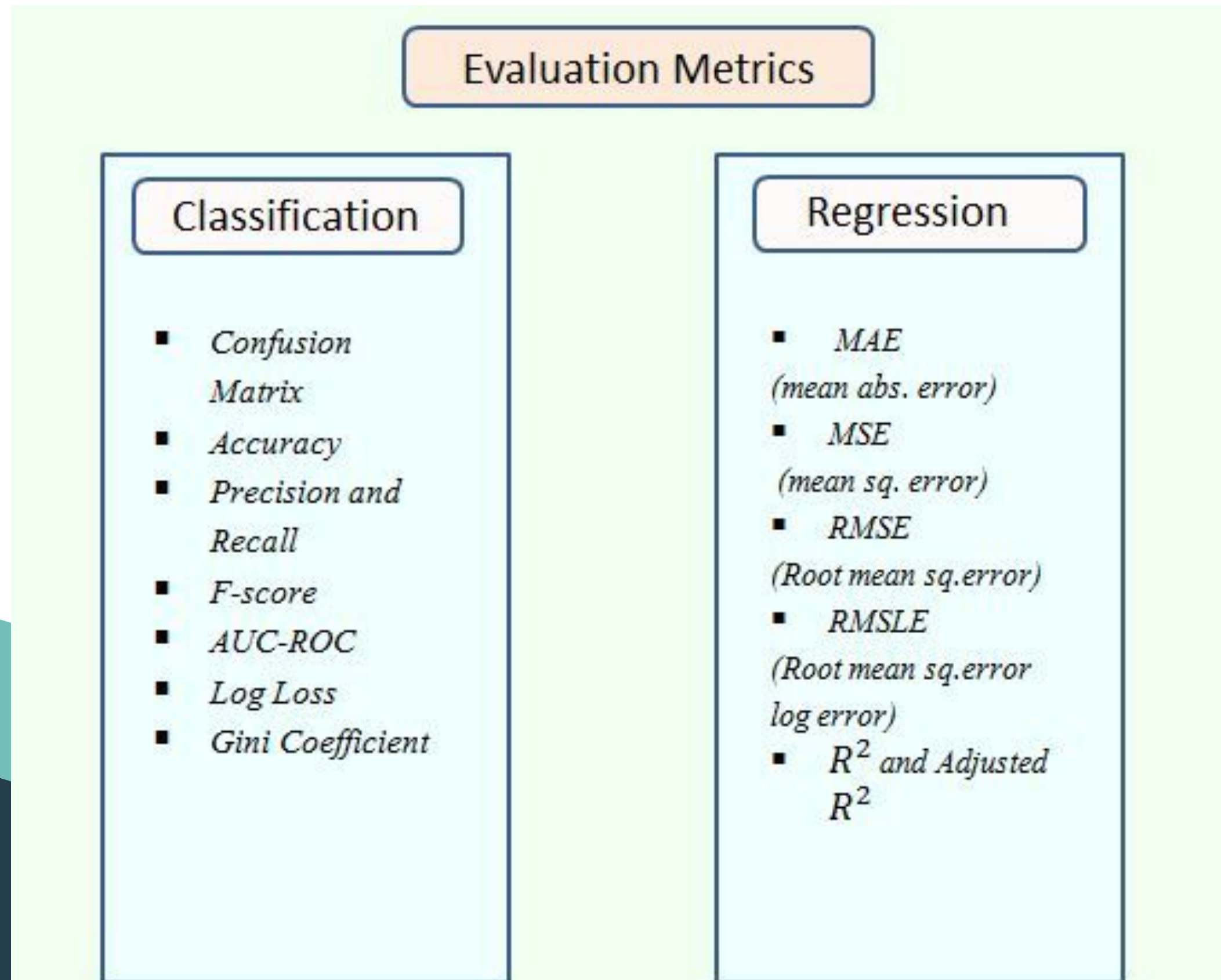


...

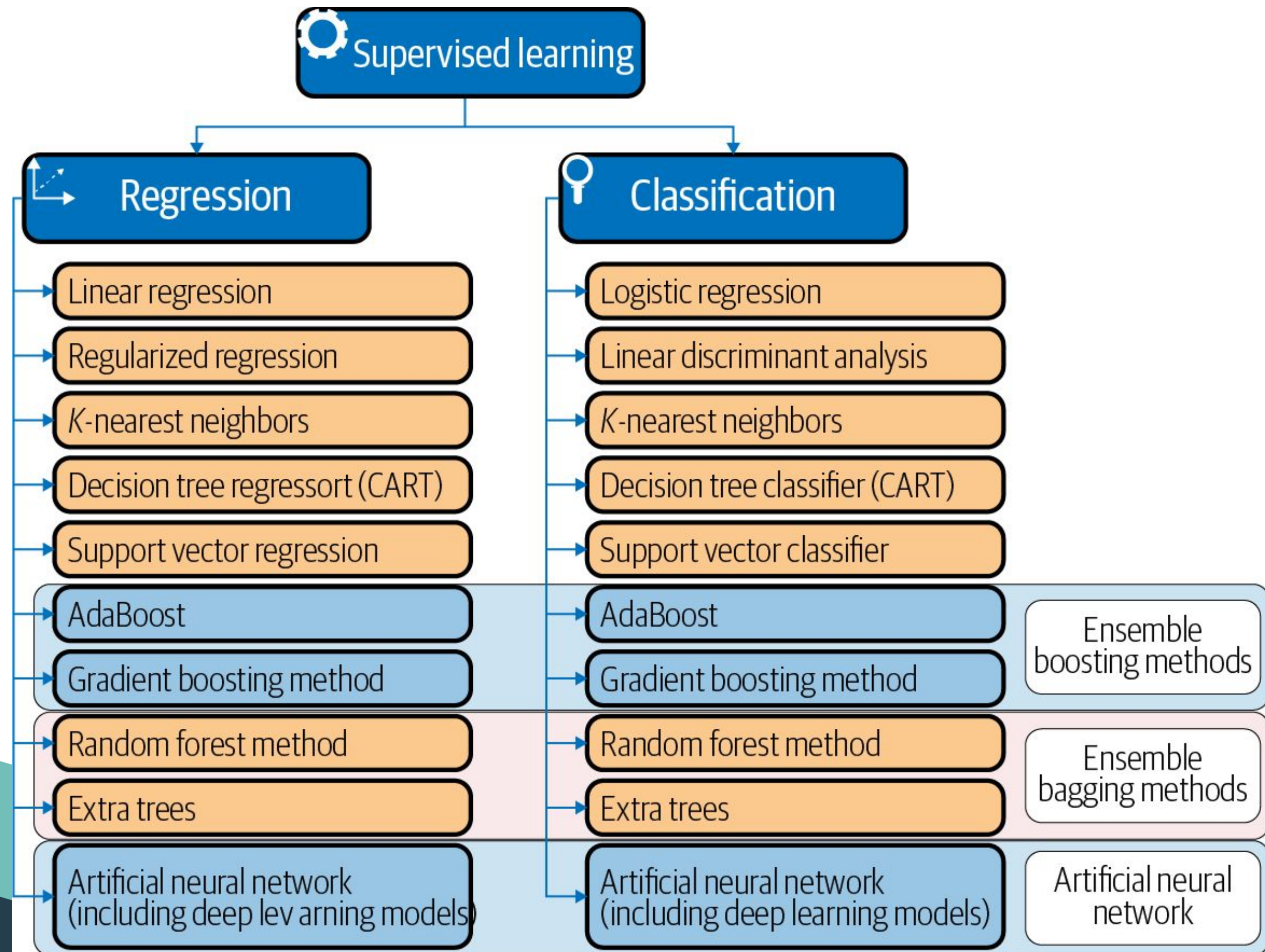




# Supervised-Metrics

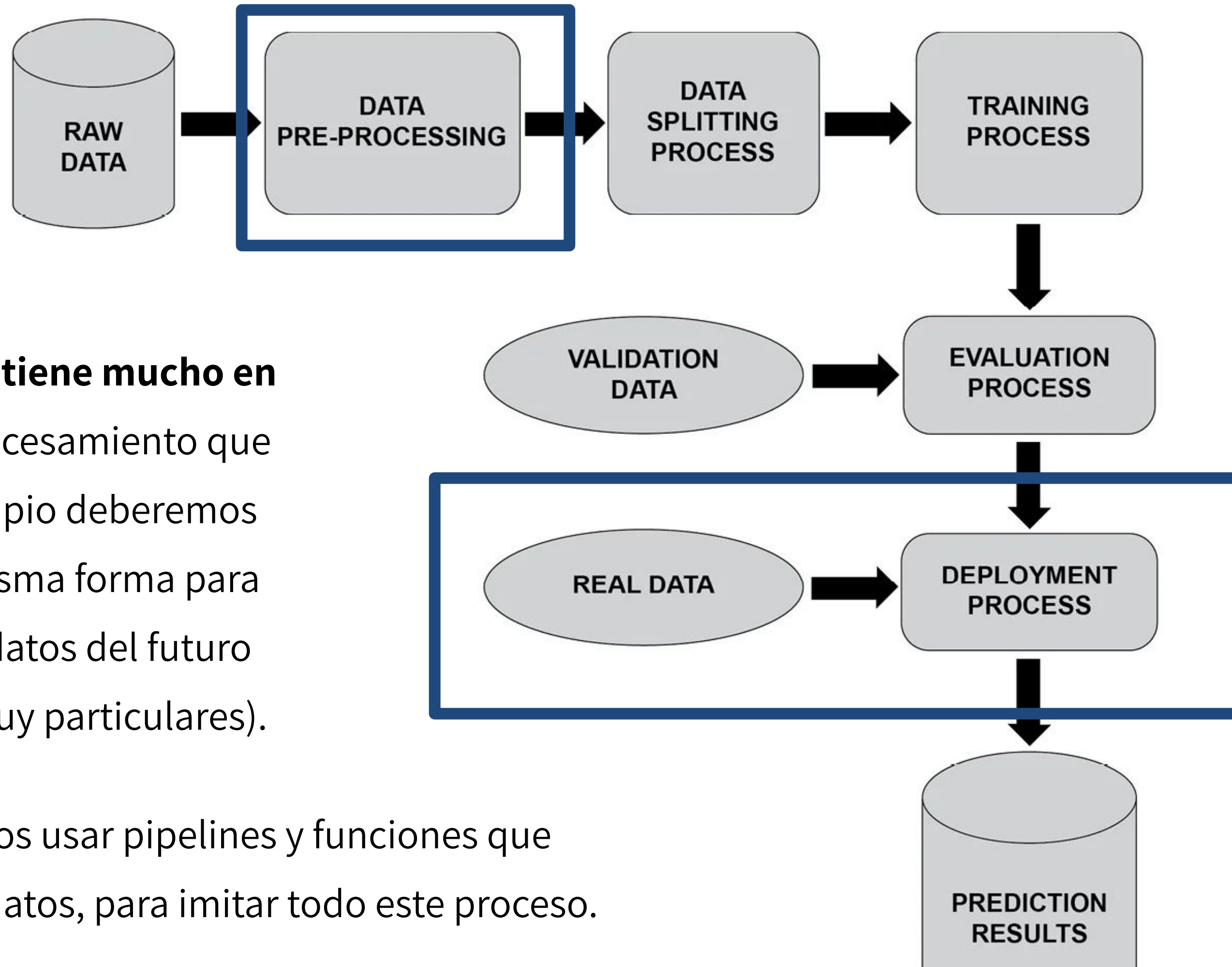


# Supervised Models





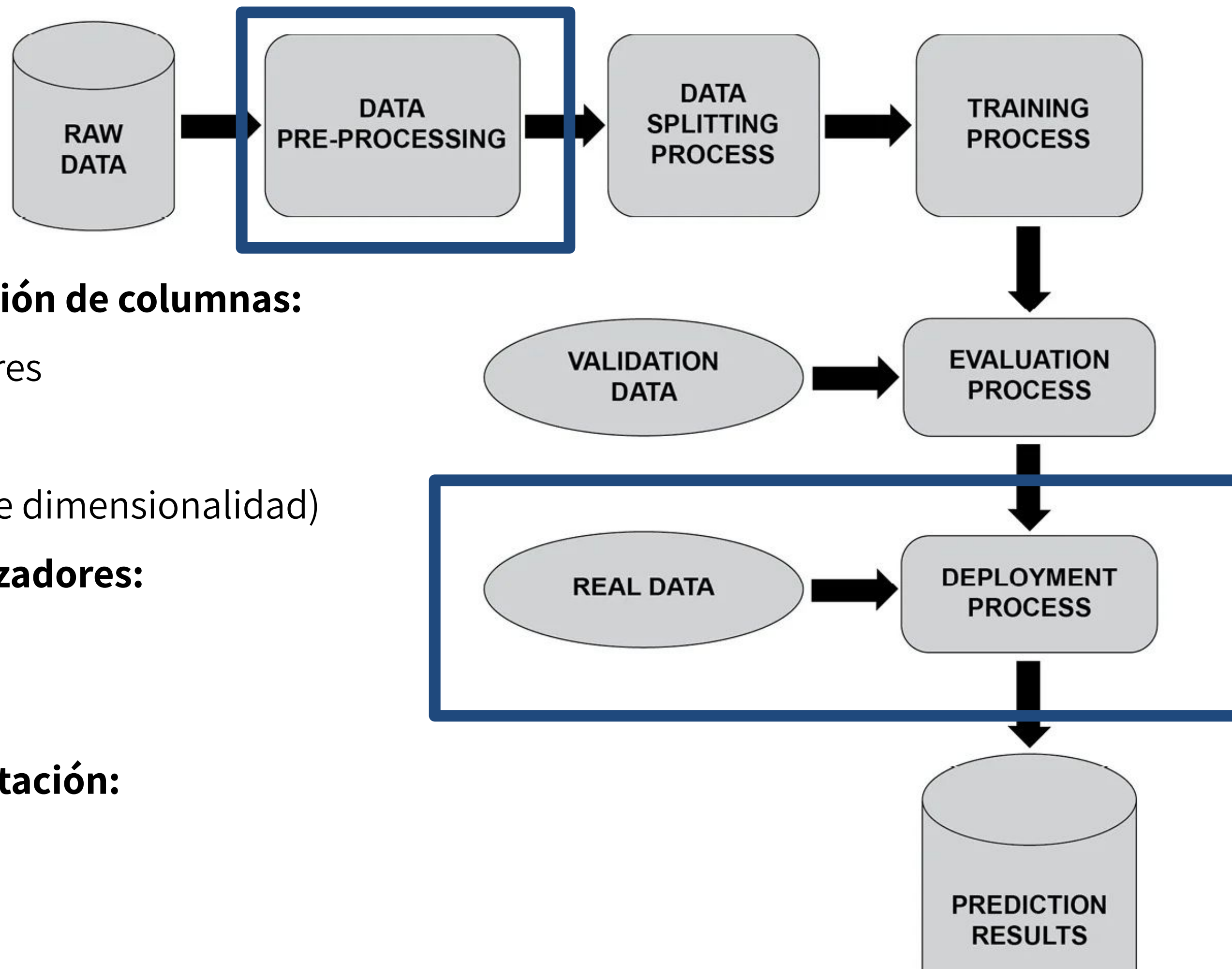
# MODEL BUILDING PROCESS



**Estas dos partes tiene mucho en común.** El preprocesamiento que hacemos al principio deberemos aplicarlo de la misma forma para cualquier set de datos del futuro (excepto casos muy particulares).

Por eso intentamos usar pipelines y funciones que transforman los datos, para imitar todo este proceso.

# MODEL BUILDING PROCESS



## Objetos de creación de columnas:

PolynomialFeatures

OneHotEncoder

PCA (reducción de dimensionalidad)

## Objetos normalizadores:

StandardScaler

MinMaxScaler

## Objetos de imputación:

SimpleImputer

KNNImputer



## Feature Selection

Subsetting the features

Ex: Using correlation with the dependent variable

## Feature Extraction

Creating new features when we could **NOT** have used raw features

Ex: from images to RGB values.  
Automatic methods such as PCA

## Feature Engineering

Creating new features when we could have used raw features

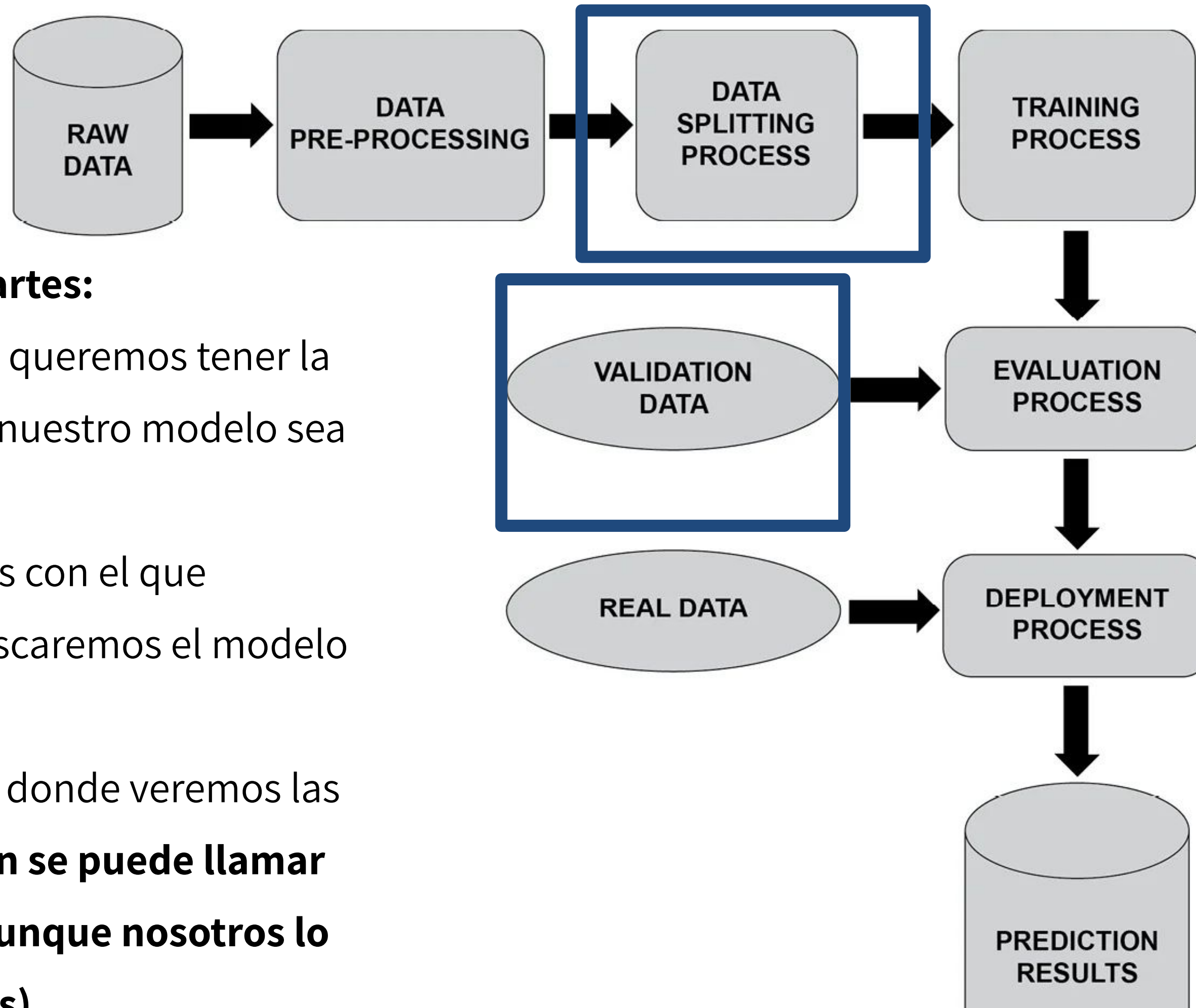
Ex: Creating a new dummy variable for working days

## Feature Learning

Constructing features automatically

Ex: Supervised neural networks, Independent component analysis

# MODEL BUILDING PROCESS



**Dividimos en 2 partes:**

**¿Por que?** Porque queremos tener la seguridad de que nuestro modelo sea estable.

**Train:** Set de datos con el que entrenaremos/buscaremos el modelo

**Test:** Set de datos donde veremos las métricas **(También se puede llamar validation data aunque nosotros lo llamamos al revés)**



# MODEL BUILDING PROCESS

## Training process:

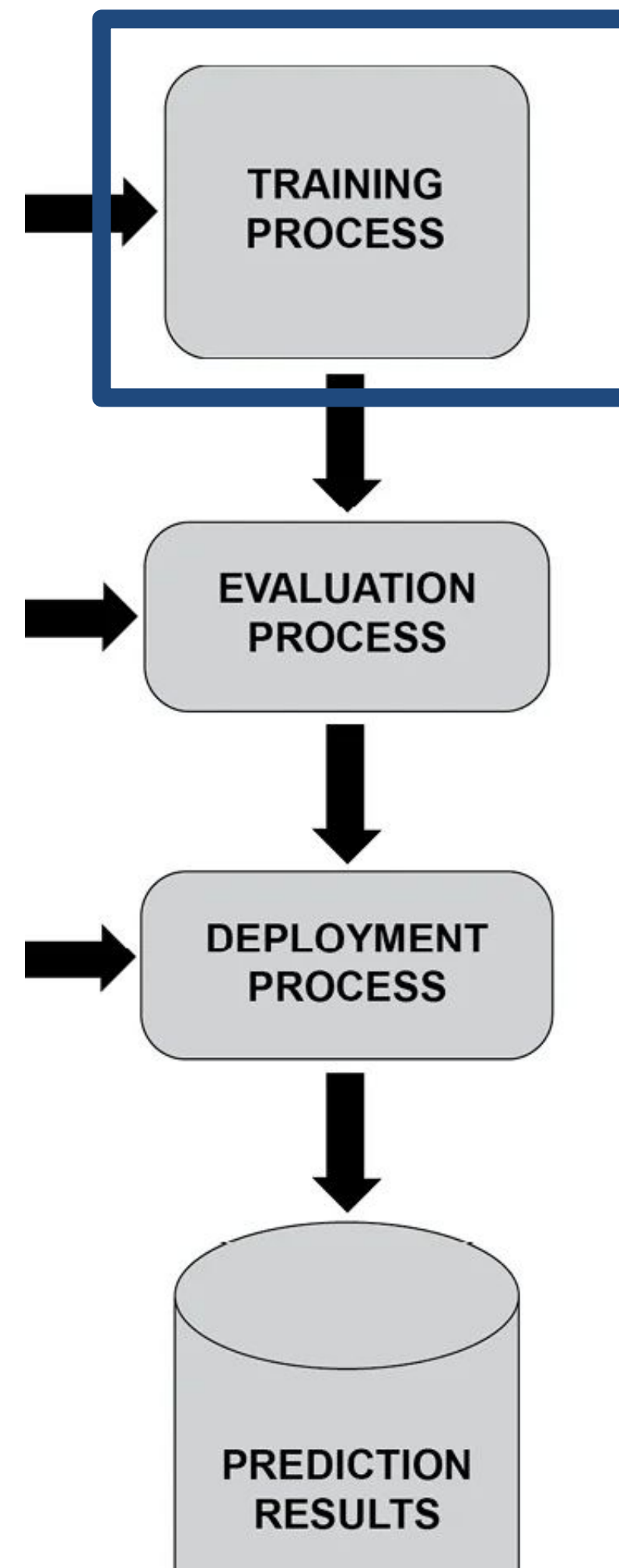
Aquí nosotros hacemos la elección de un tipo de modelo, dependiendo del tipo de problema. Asumimos que elegimos un XGBoost. **Entrenamos un modelo a partir de enseñarles unas X y una Y asociadas.**

## Búsqueda de hiperparametros (get\_params y grid\_searchCV):

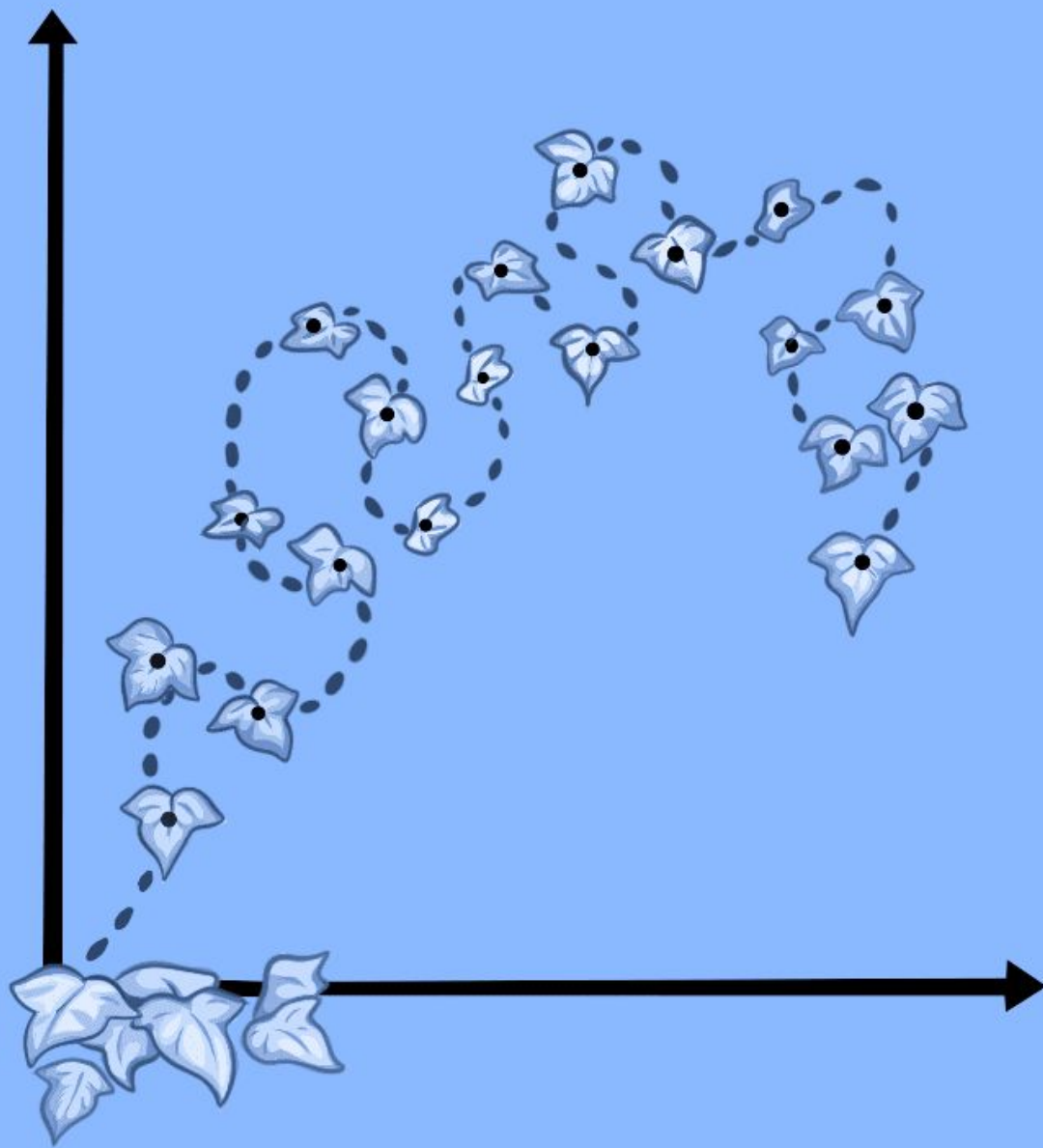
Definimos un conjunto de hiperparametros del modelo, como el max\_depth. Probaremos cada combinación de candidatos y elegiremos la mejor

## Elegimos un método de validación (grid\_searchCV):

Cross-validation: En vez de tener una métrica de cada combinación tendremos K (5 normalmente, hacemos 5 modelos entrenando con 4 partes del train y testeando con 1 parte del train). Con estas K métricas elegimos el mejor modelo.



# Overfitting



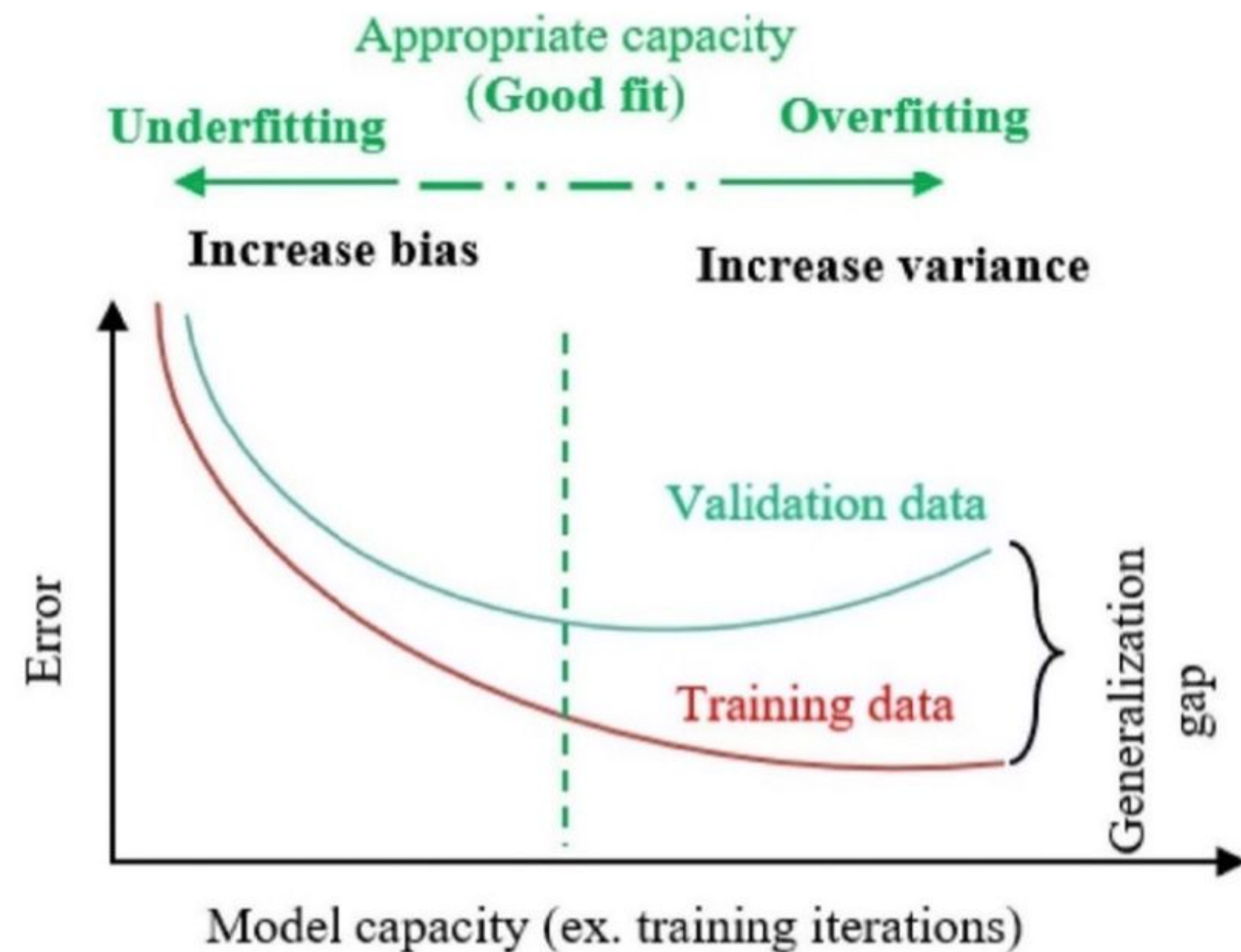
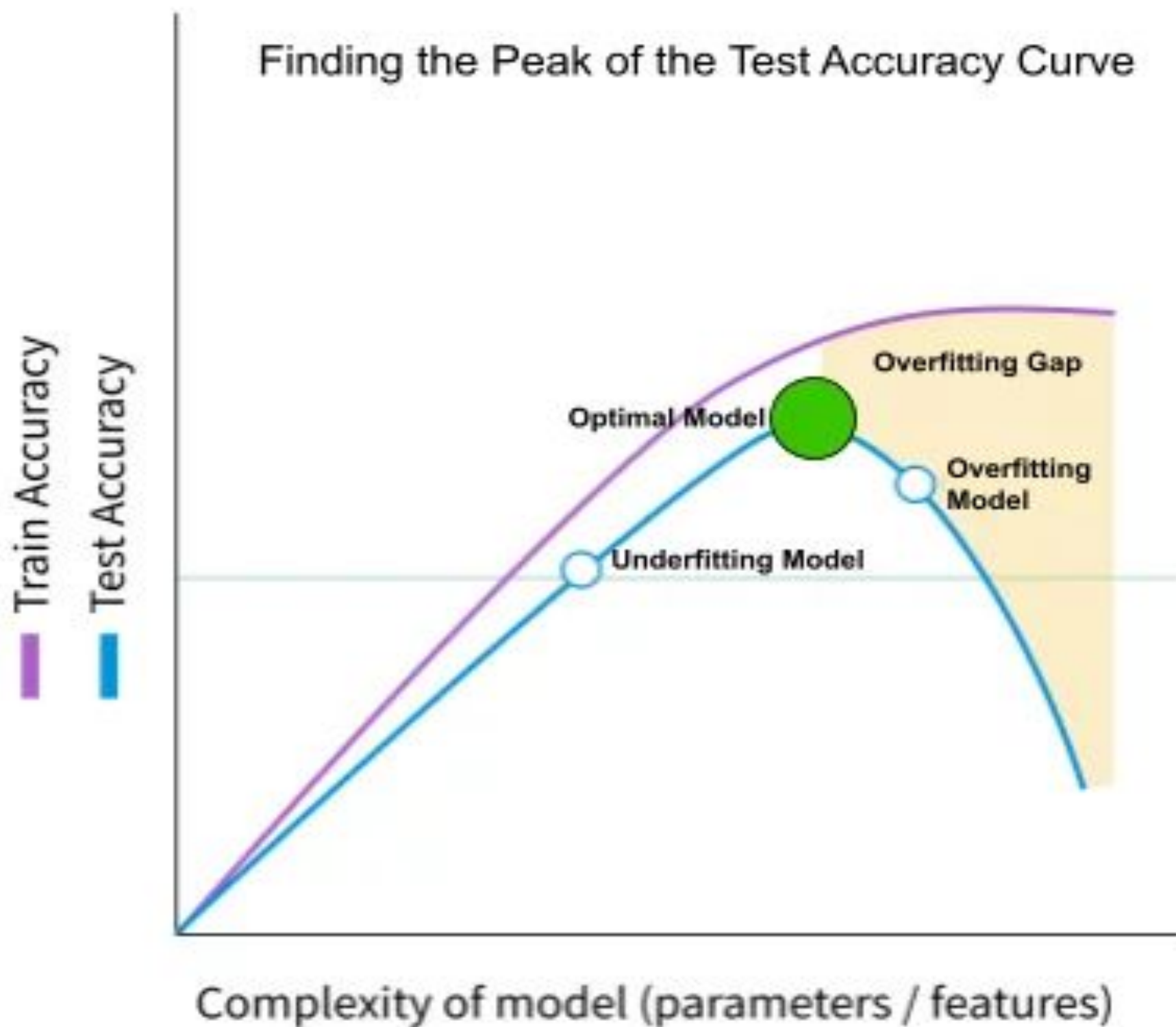
## Overfitting

*['ō-vər-'fi-tiŋ]*

A modeling error in statistics that occurs when a function is too closely aligned to a limited set of data points.



# Overfitting



Training loss	Validation loss	Situation	Solution
High	High	Underfitting	Increase capacity
Low	High	Overfitting	Decrease capacity (shrink or regularization techniques)
Low	Low	Good fit	Run test
High	Low	Unlikely	Debug

# Avoid Overfitting



## How to Tackle Overfitting?



Increase training data



Reduce model complexity



Early stopping during the training phase



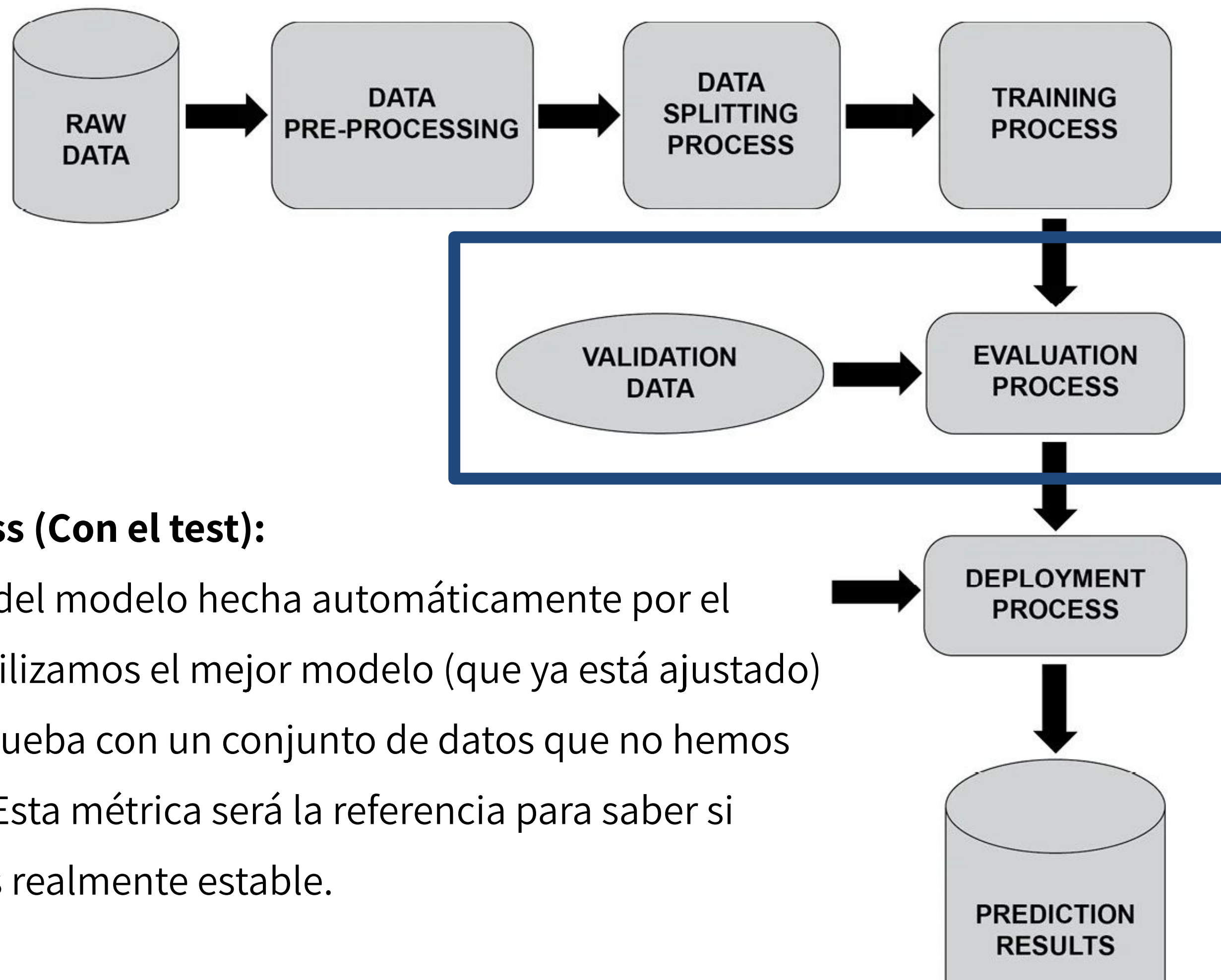
Ridge Regularization and Lasso Regularization



Use dropout for neural networks to tackle overfitting



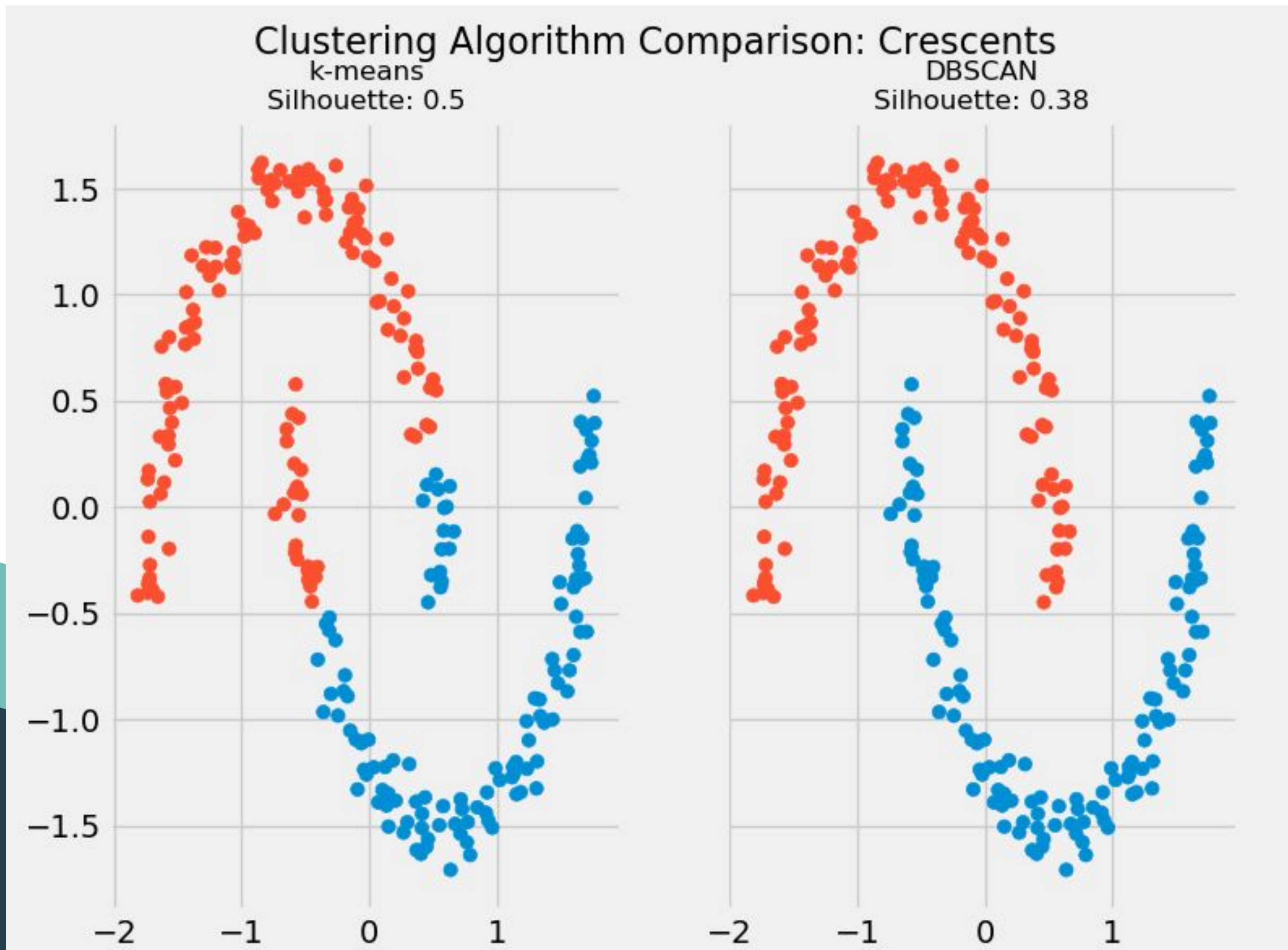
# MODEL BUILDING PROCESS



## Evaluation process (Con el test):

Ya con la elección del modelo hecha automáticamente por el **grid\_searchCV**, utilizamos el mejor modelo (que ya está ajustado) para una última prueba con un conjunto de datos que no hemos utilizado todavía. Esta métrica será la referencia para saber si nuestro modelo es realmente estable.

# Unsupervised-Models





# Unsupervised-Models

