

# Ciencia y analítica de datos

Dra. María de la Paz Rico Fernández.



Tecnológico  
de Monterrey

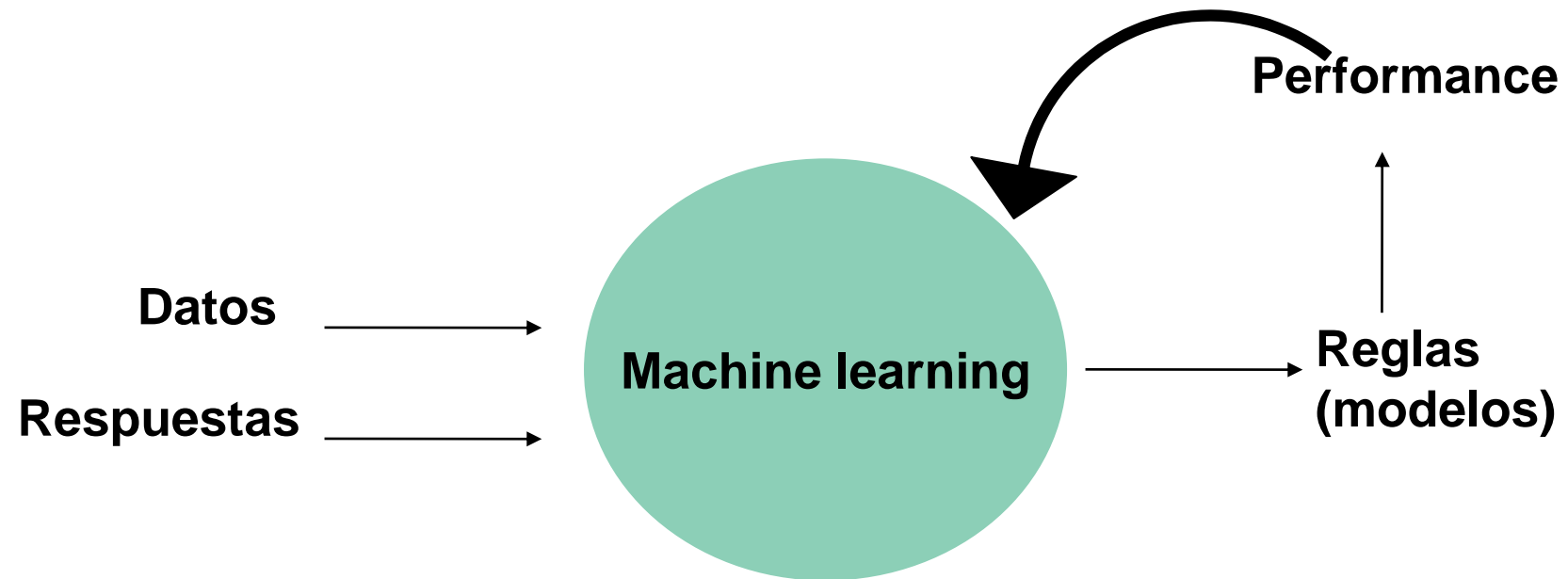


# Presentación

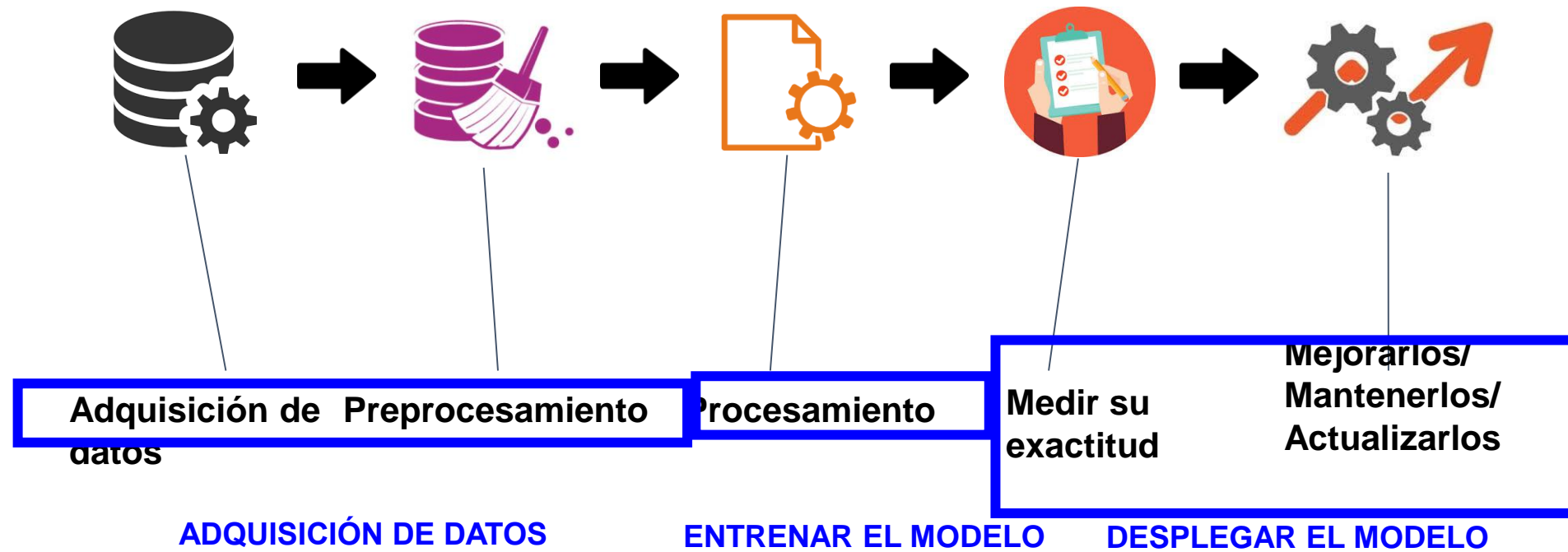
- PhD. María de la Paz Rico Fdz.
- Doctorado en ciencias en Robótica y Manufactura avanzada con especialidad en visión por computadora e inteligencia artificial, Centro de Investigación y Estudios Avanzados del Instituto Politécnico Nacional.
- AI Data Engineer por ANCUD IT, Berlin, Alemania.
- Instructora Certificada de NVIDIA.
- Trabajos:
  - Computer Vision Research Engineer en AIFI Inc, Silicon Valley.
  - Chief Knowledge Officer en Centro de Innovación Industrial en Inteligencia Artificial.
  - City lead Monterrey Women in AI



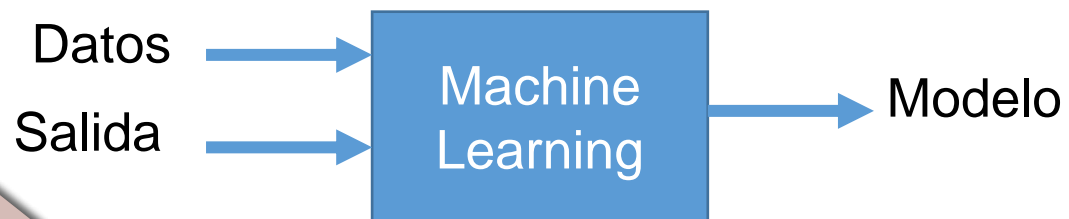
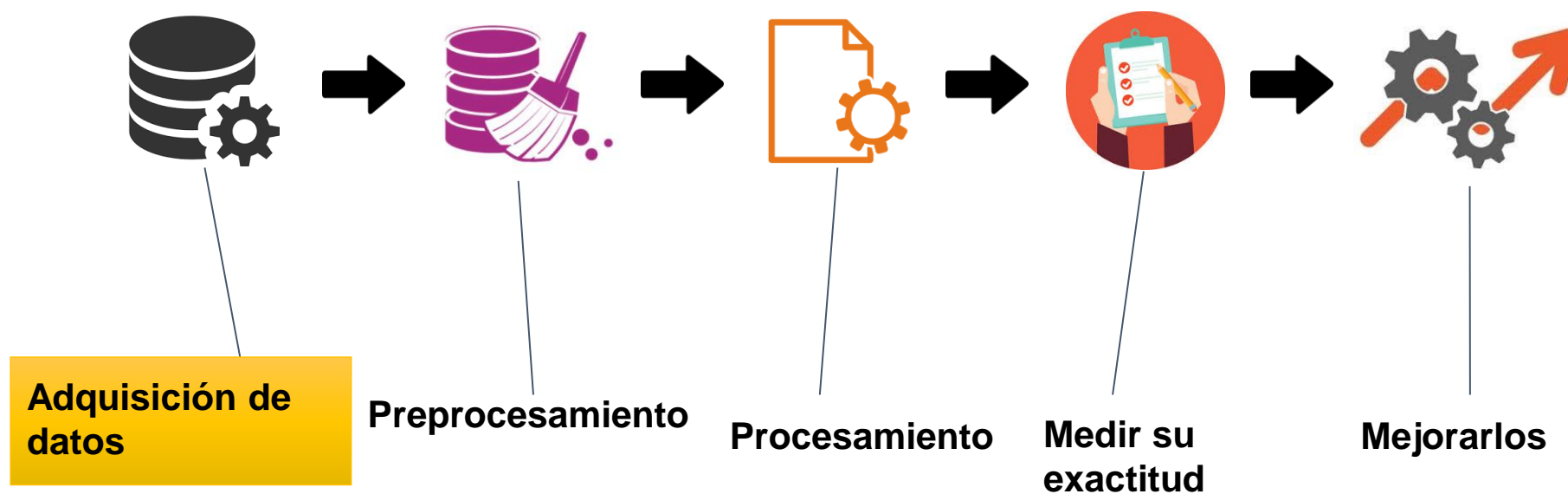
- ¿Qué necesitamos para ml?



## ■ ml pipeline



Entender el problema,  
identificar fuentes de datos  
(etiquetados) y resaltar  
posibles problemas con los  
datos.



# PASOS CLAVE PARA ML PROJECT





## ( ML PIPELINE )

- Adquisición de datos

### Ejemplos de bases de datos

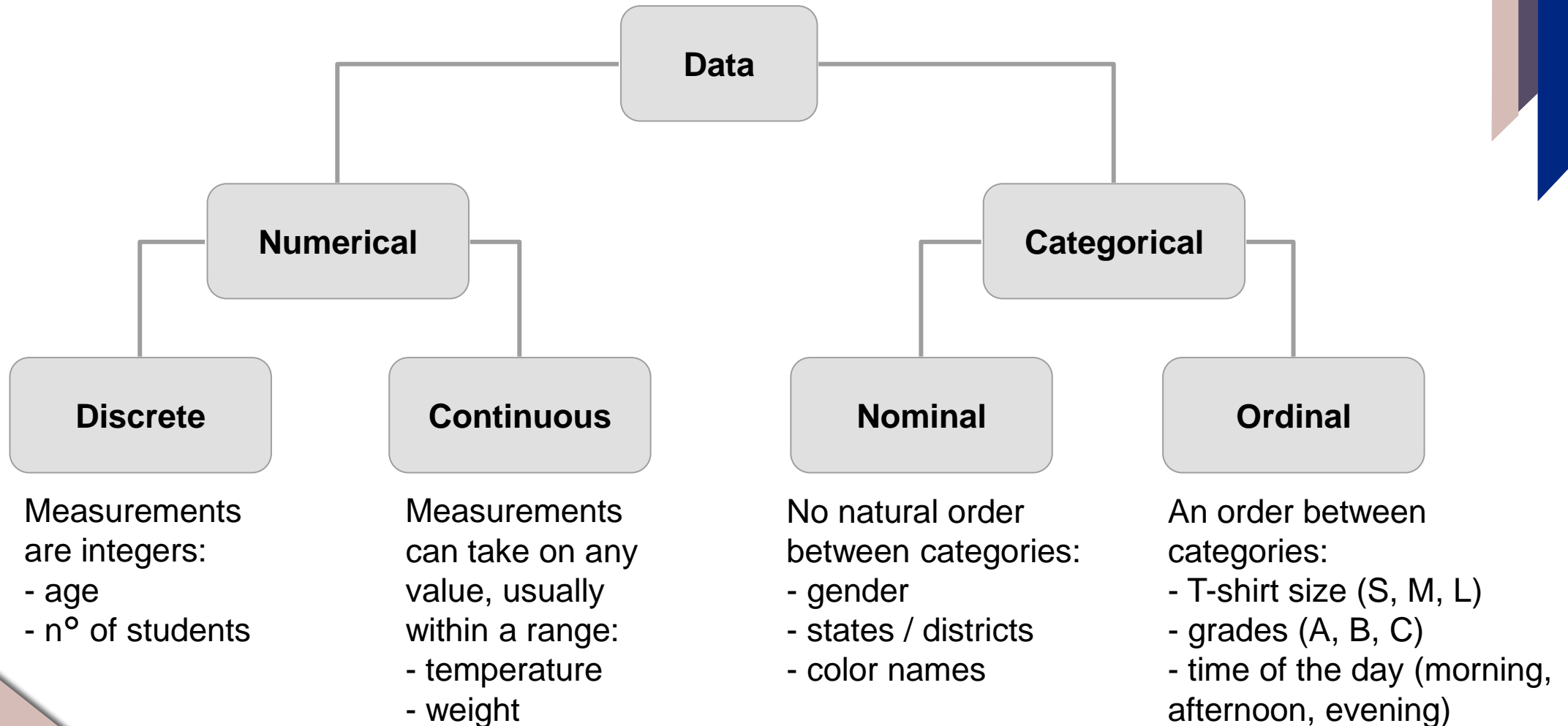
Home prices

size of house (square feet)	# of bedrooms	# of bathrooms	newly renovated	price (1000\$)
523	1	2	N	115
645	1	3	N	150
708	2	1	N	210
1034	3	3	Y	280
2290	4	4	N	355
2545	4	5	Y	440

image	label
	cat
	not cat
	cat
	not cat

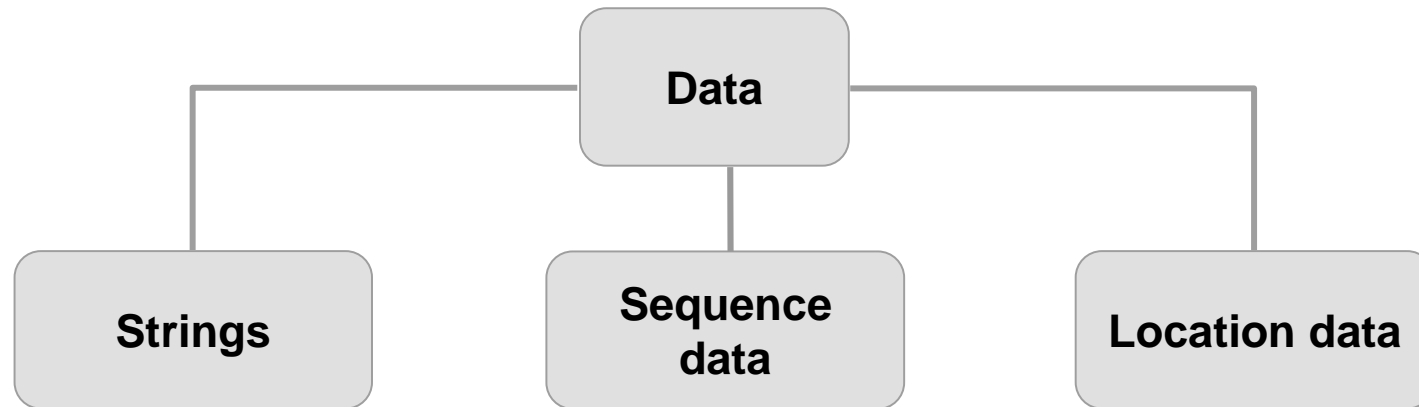
# Machine Learning

## TIPOS DE DATOS



# Machine Learning

## Tipos de datos



- time series (time order)
- sequences of strings (text data)

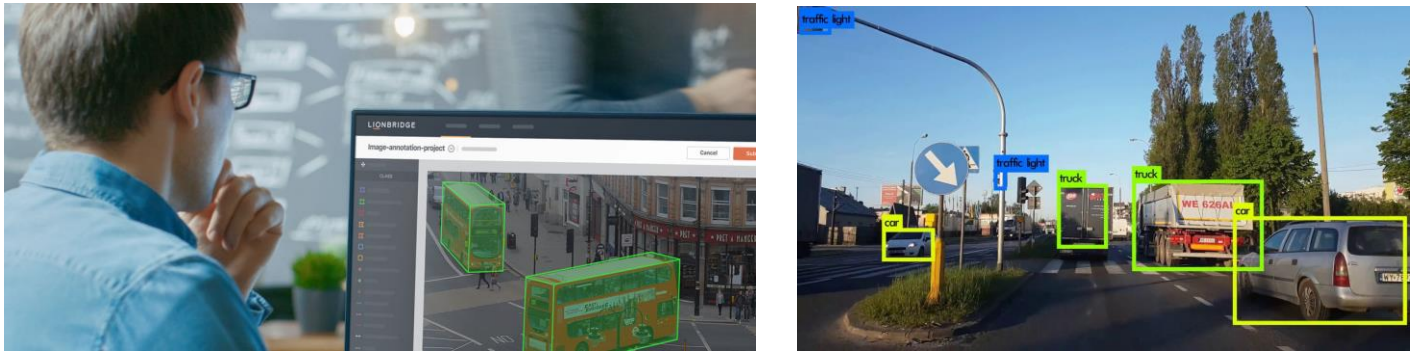


# PASOS CLAVE PARA ML PROJECT ( ML PIPELINE )

- Adquisición de datos

Podemos adquirir las bases de datos por:

1) Etiquetado manual:



2) De datos observados

user ID	time	price (\$)	purchased
4783	Jan 21 08:15.20	7.95	yes
3893	March 3 11:30.15	10.00	yes
8384	June 11 14:15.05	9.50	no
0931	Aug 2 20:30.55	12.90	yes

machine	temperature (°C)	pressure (psi)	machine fault
17987	60	7.65	N
34672	100	25.50	N
08542	140	75.50	Y
98536	165	125.00	Y

3) Descargandola de paginas web o partners.

# PASOS CLAVE PARA ML PROJECT (ML PIPELINE)

- Adquisición de datos y sus problemas

1) Si entran datos ruidosos, estimaciones ruidosas saldr



2) Problemas en los datos

- Etiquetas incorrectas
- Datos faltantes



given: 5  
corrected: 3



given: cat  
corrected: frog



given: lobster  
corrected: crab

3) Varios tipos

- Estructurados: Tablas de datos.
- No estructurados: imágenes, audio, video, texto



Structured  
Data

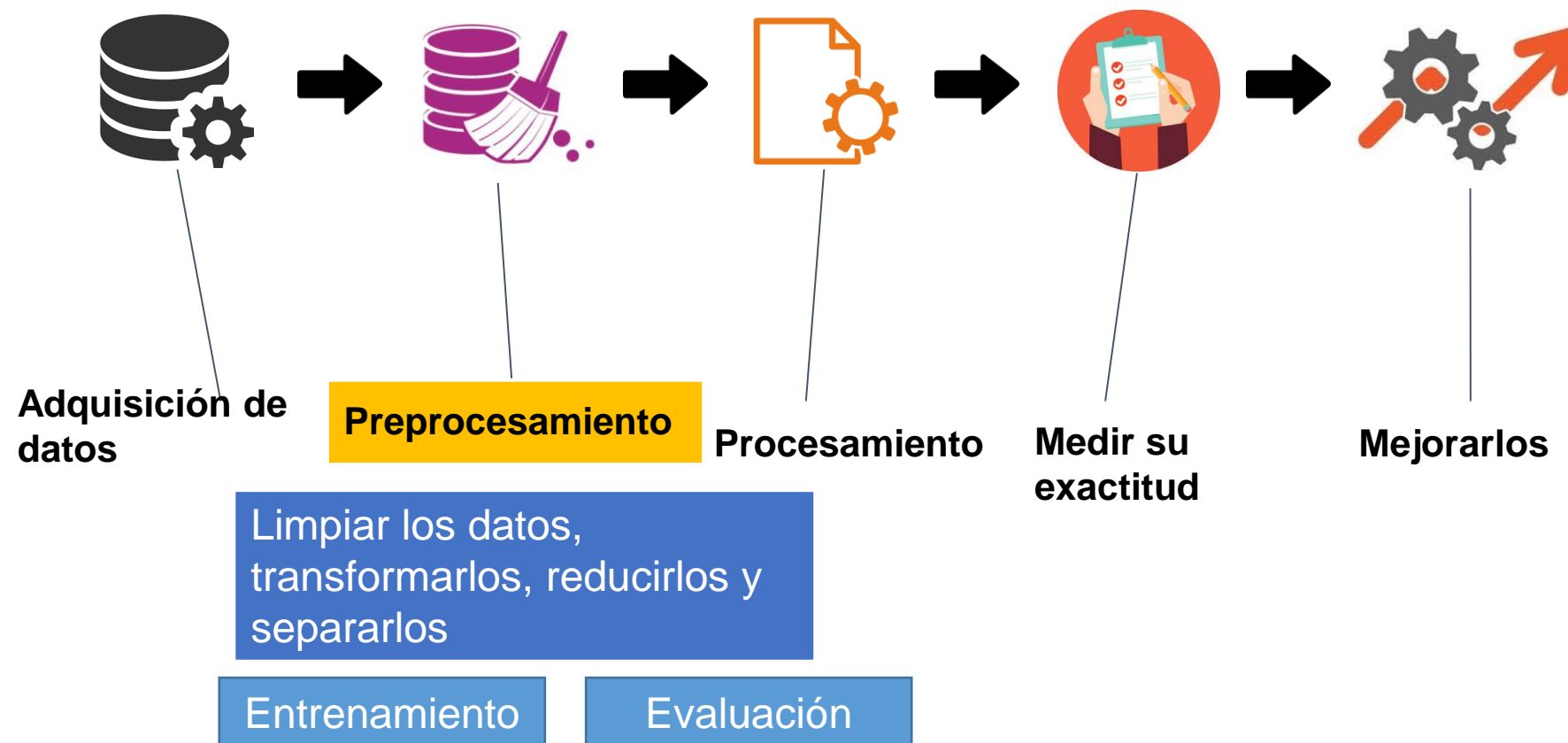


Textual  
Data

Image  
File

Video

Audio

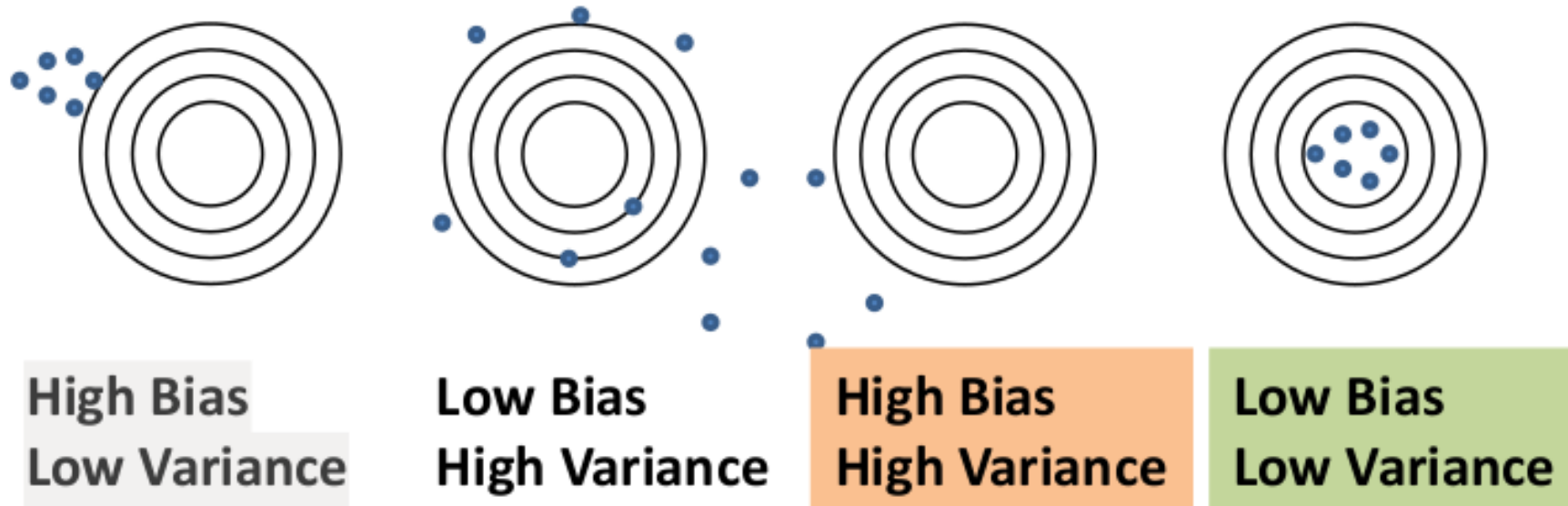


# Machine Learning

## Challenges of Machine Learning

**High-bias models:** consistent but wrong predictions, prone to underfitting.

**High-variance models:** prone to overfitting.

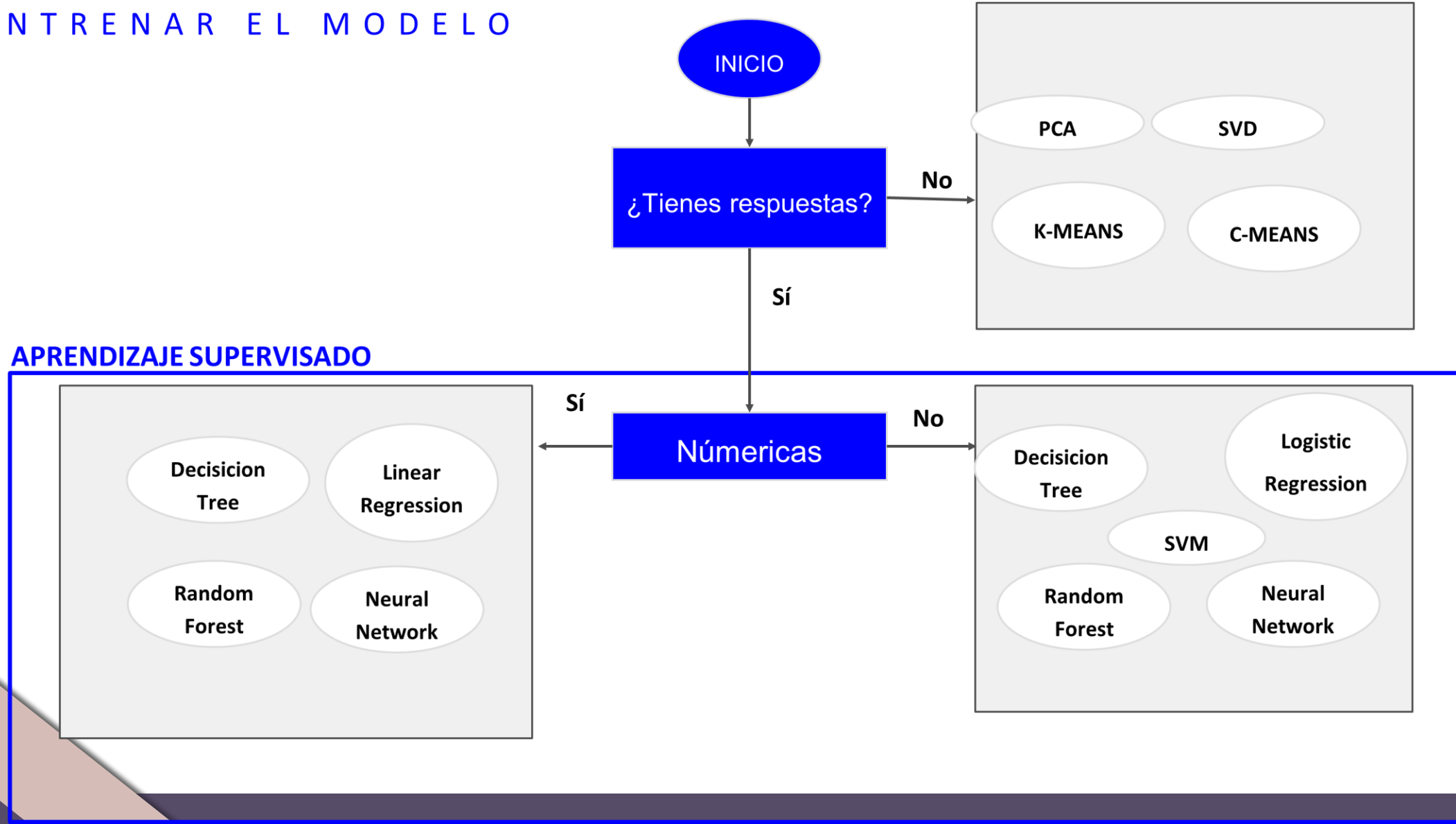


# PASOS CLAVE PARA ML PROJECT (ML PIPELINE)

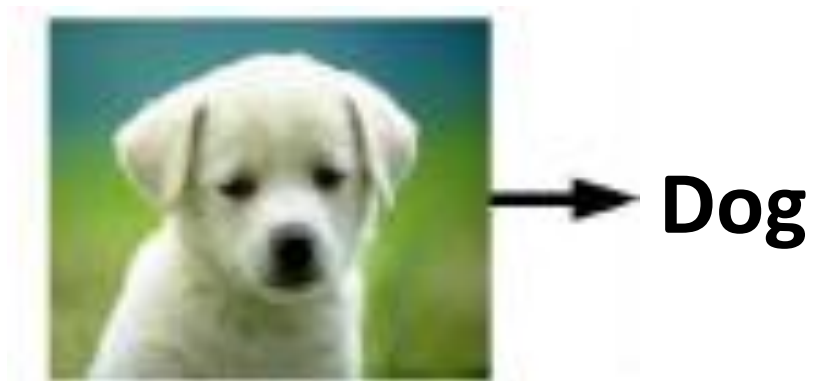
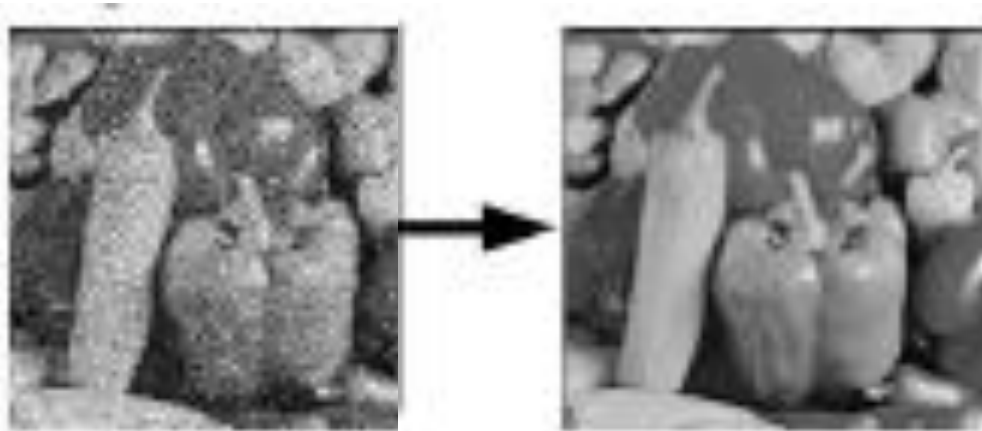
ENTRENAR EL MODELO

APRENDIZAJE NO SUPERVISADO

APRENDIZAJE SUPERVISADO



# Machine Learning



# Machine Learning

## Aprendizaje supervisado

## Aprendizaje no supervisado

Regresión

Clasificación

Reducción de dimensionalidad

Asociación

Métodos:

- Regresión lineal
- Árboles de decisión
- SVM
- Neural network

Métodos:

- K-means
- PCA

# Machine Learning

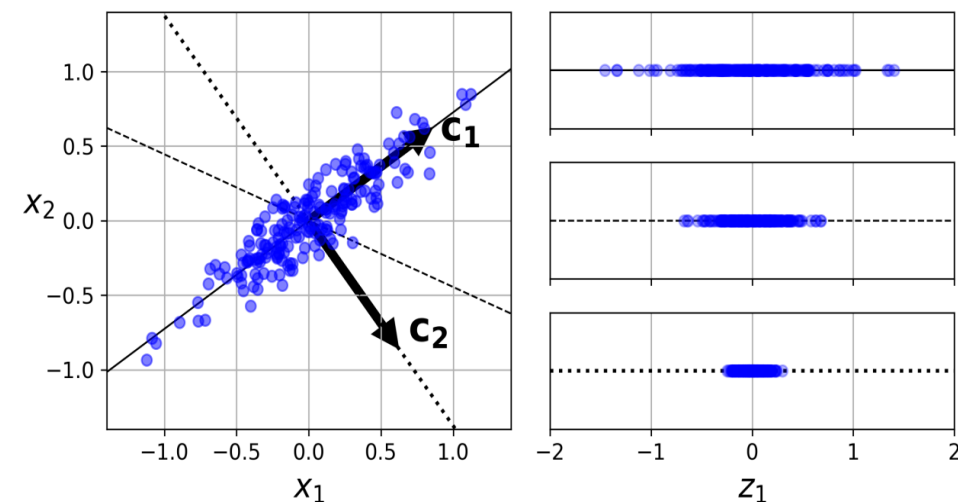
## Unsupervised Learning

### Principal Component Analysis (PCA)

- **Dimensionality reduction:** when there are many features (e.g. thousands or millions) for each training instance it makes training slow and it could be hard to find a good solution.

- **PCA:**

- identifies the hyperplane that lies closest to the data
- projects the data onto the hyperplane
- selects the projection that preserves the maximum amount of variance





# Dimensionality Reduction

- When there are many features (e.g. thousands or millions) for each training instance it makes training slow and it could be hard to find a good solution.
- Reducing dimensionality of the training set before training a model speeds up training.
- Reducing dimensionality does reduce information.
- It is useful for data visualization.

# The Curse of Dimensionality

- High-dimensional datasets are risk of being very **sparse**: most training instances are far away from each other making predictions less reliable since they will be based on larger extrapolations.
- **Sparsity** is a problem for statistical significance, the amount of data needed to support the result often grows exponentially with the dimensionality.

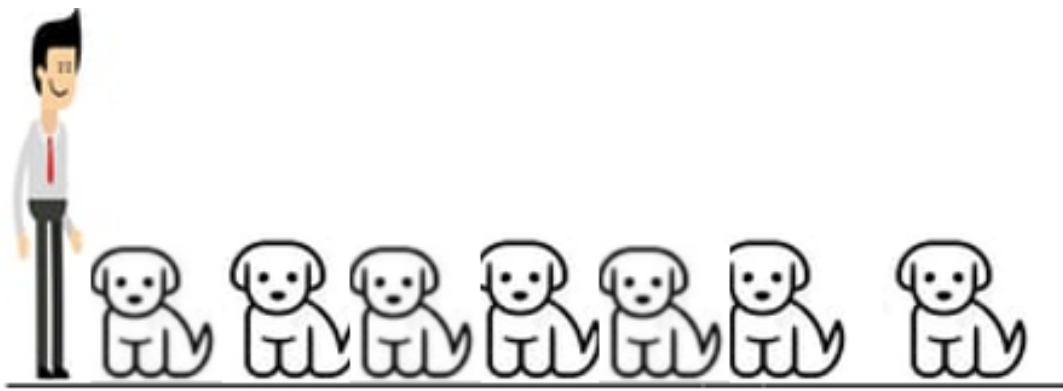


Figure 1 - One-dimension scenario

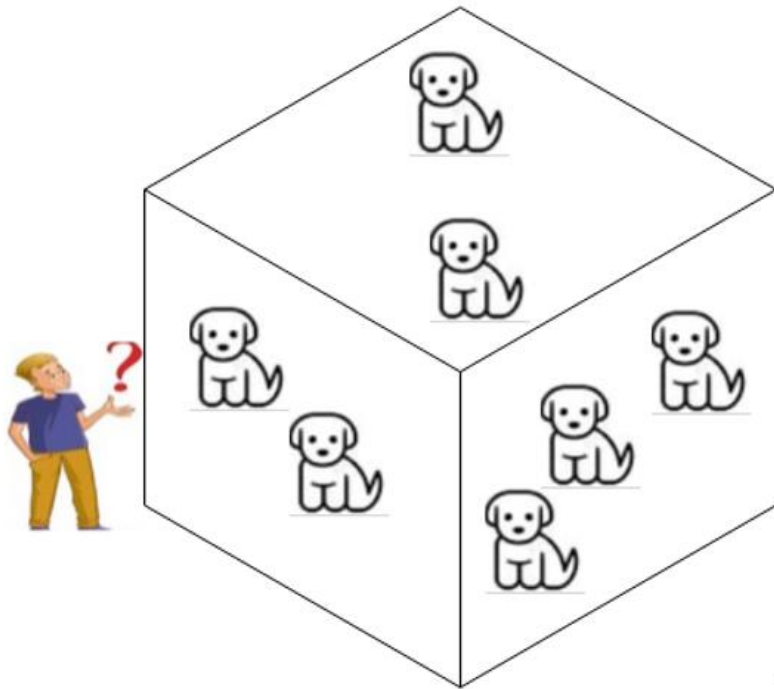


Figure 3 - Three-dimension scenario

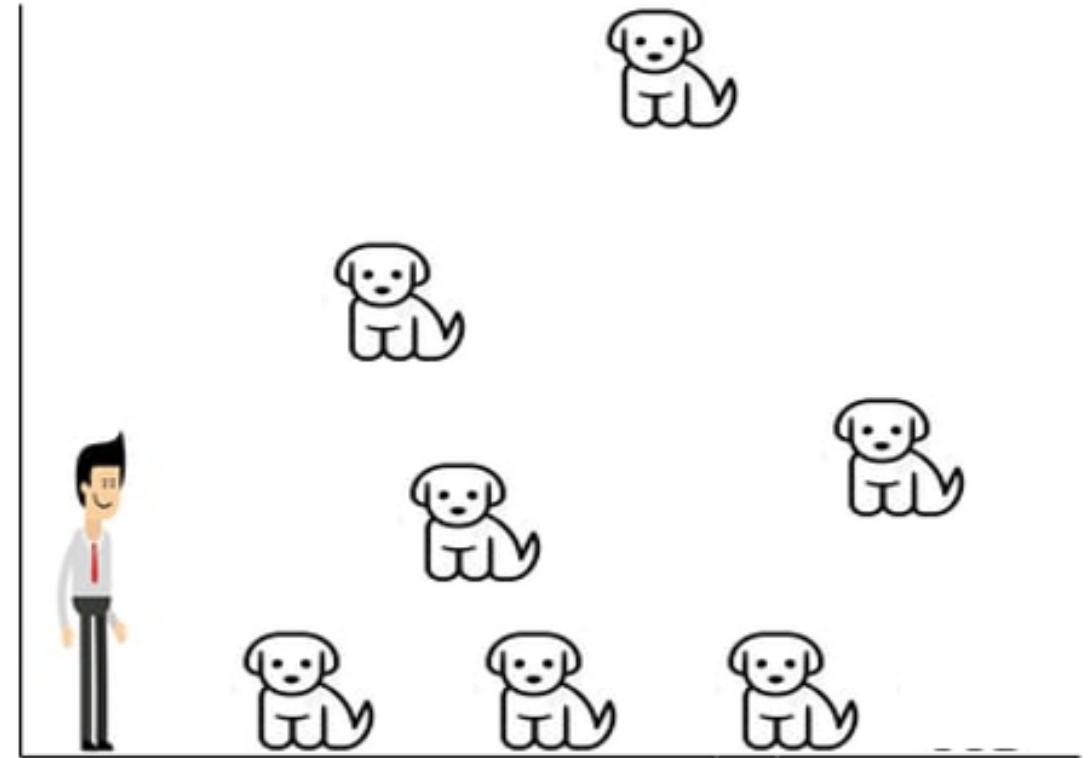
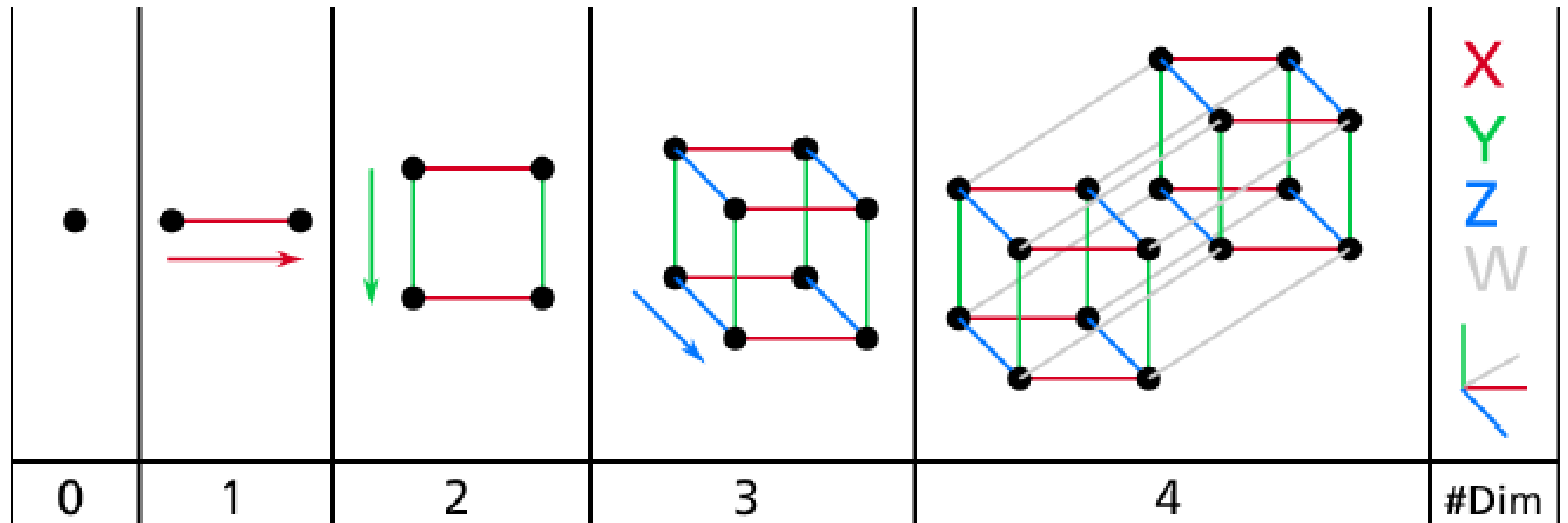


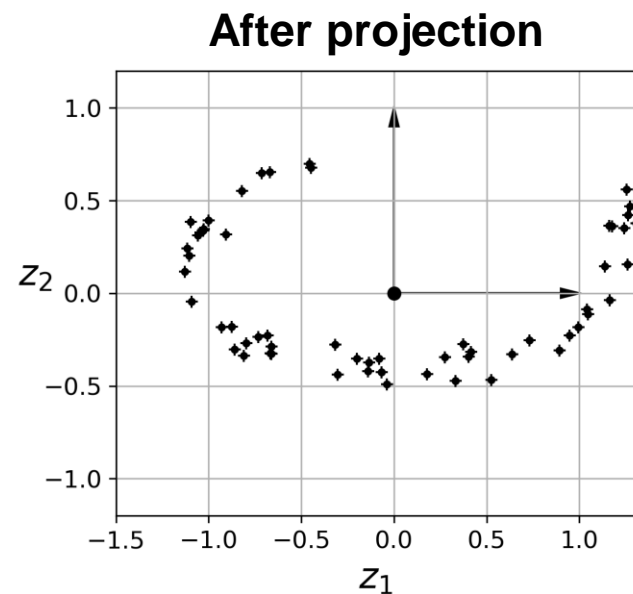
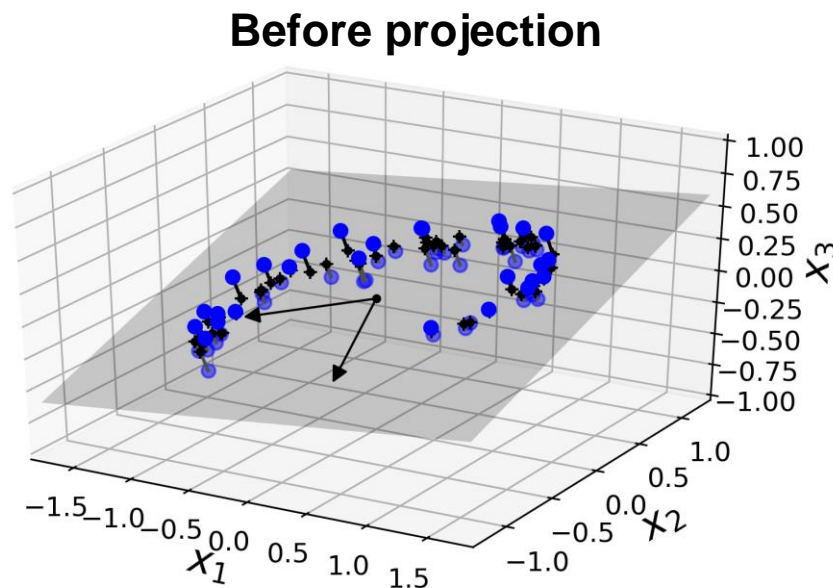
Figure 2- Two-dimension Scenario



# Approaches for Dimensionality Reduction

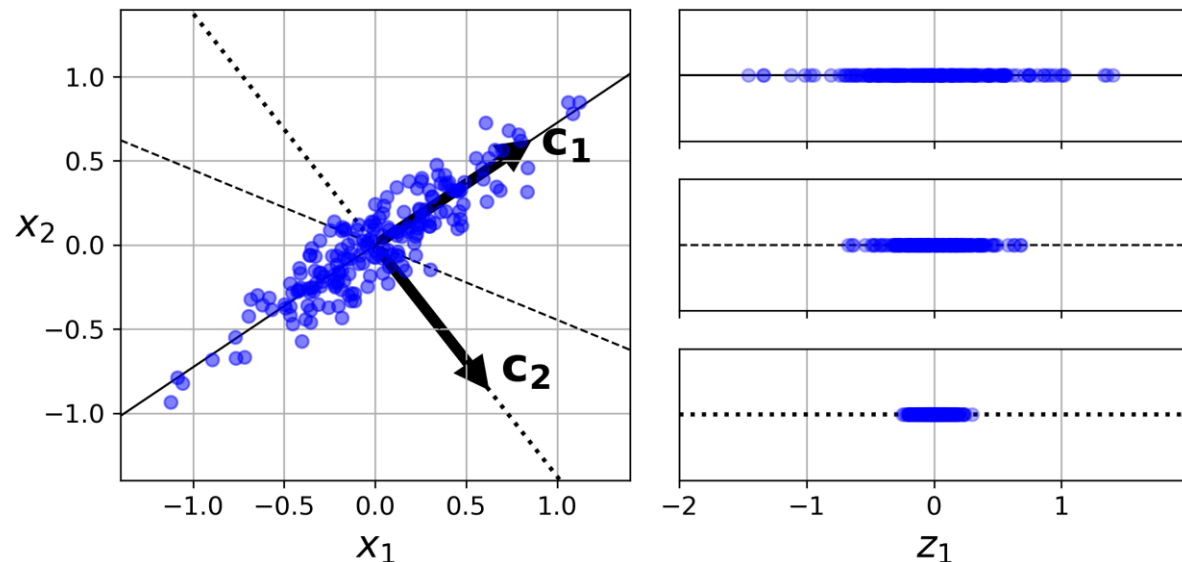
## 1- Projection

- In real world-problems, training instances are not spread out uniformly across all dimensions.
- Then, training instances lie within a lower-dimensional **subspace** of the high-dimensional space.



# Principal Component Analysis(PCA)

- Unsupervised method for dimensionality reduction of the data
- Identifies the hyperplane that lies closest to the data.
- Then, it projects the data onto it.
- Selects the projection that preserves the maximum amount of variance (the axis that minimizes the mean squared distance between the original dataset and its projection onto that axis).



$$X_{dproj} = X \cdot W_d$$

$X$ : matrix training set

$W_d$ : matrix containing the first  $d$  principal components

This projects the training set onto the space defined by the principal components.

# Principal Component Analysis(PCA)

## Principal Components

- PCA identifies the **axis** that accounts for the largest amount of **variance** in the training set.
- PCA finds a second axis, **orthogonal** to the first one, that accounts for the largest amount of remaining variance.
- PCA would also find a third axis, orthogonal to the both previous axis, etc.
- The unit vector that identifies the  $i^{\text{th}}$  axis is called **principal component**.
- PCA assumes that the dataset is centered around the origin. Scikit-Learn's PCA classes centers the data.

# Principal Component Analysis(PCA)

## Scikit-Learn

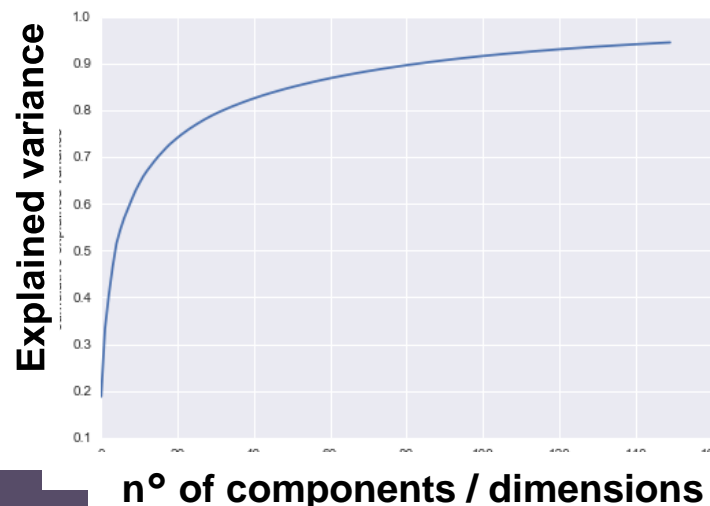
- After fitting PCA transformer to the dataset, the principal components can be accessed using ***components\_*** variable (the first one: `pca.components_T[:, 0]`).
- Explained variance ratio of each principal component: the proportion of the dataset's variance that lies along the axis of each principal component, ***explained\_variance\_ratio\_*** variable.



# PCA (Principal Component Analysis)

## Choosing the Right Number of Dimensions

- It is generally useful to choose the number of dimensions that add up to a large proportion of the variance (i.e. 95%).
- In case of data visualization the dimensionality is usually reduced to 2 or 3 dimensions.
- **Scikit\_Learn**: set `n_components` to a float between 0.0 and 1.0, indicating the ratio of variance to preserve, ***PCA(n\_components=0.95)***.
- Plot the explained variance as a function of the number of dimensions.



# Principal Component Analysis(PCA)

## PCA for Compression

- After dimensionality reduction, the training set takes up less space.
- Speeds up an algorithm like SVM.
- It's possible to decompress the reduced dataset by applying the inverse transformation of the PCA.

# PCA (Principal Component Analysis)

## Disadvantages

- PCA tends to be highly affected by **outliers** in the data.
- PCA assumes that the principle components are a **linear** combination of the original features.
- PCA assumes that the principle components are **orthogonal**.
- PCA uses **variance** as the measure of how important a particular dimension is.
- High variance axes are treated as principle components.
- Low variance axes are treated as noise.