



IIC 2433 Minería de Datos

<https://github.com/marcelomendoza/IIC2433>

- OUTLINE -

¿Qué vamos a ver?

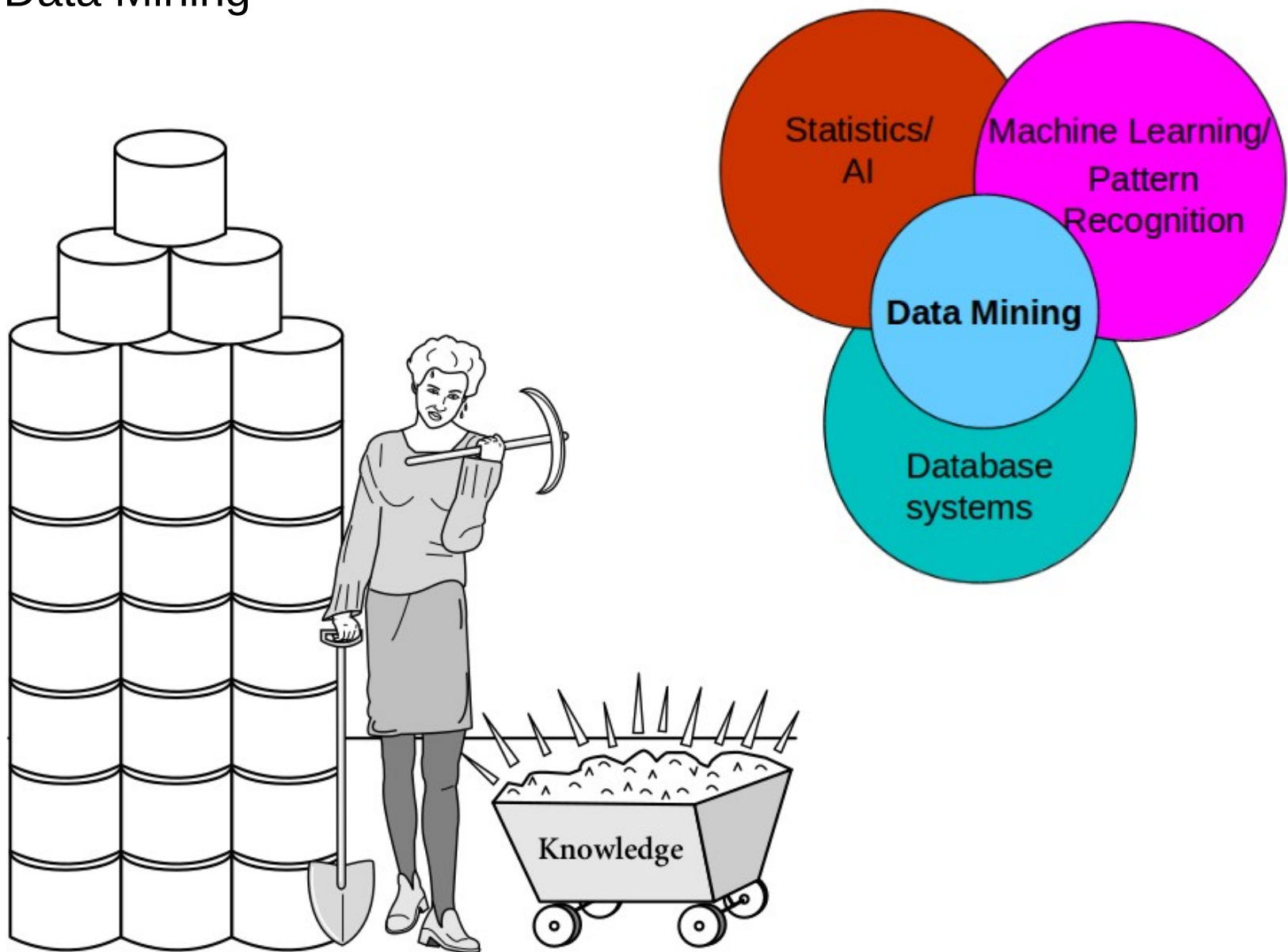
PCA, t-SNE

k-means, HAC, DBSCAN, Louvain

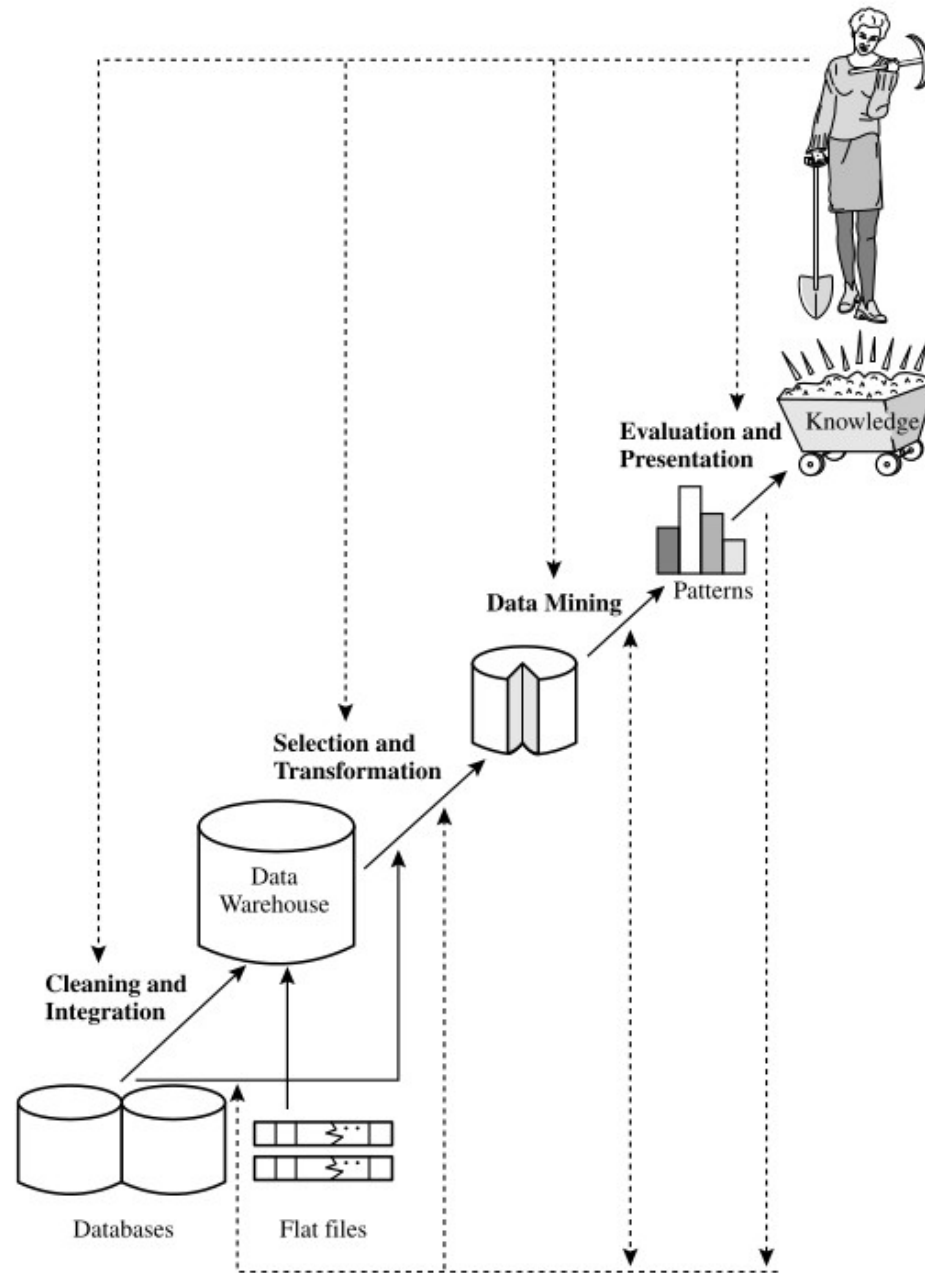
Perceptron, SVM, ensembles

AE, VAE

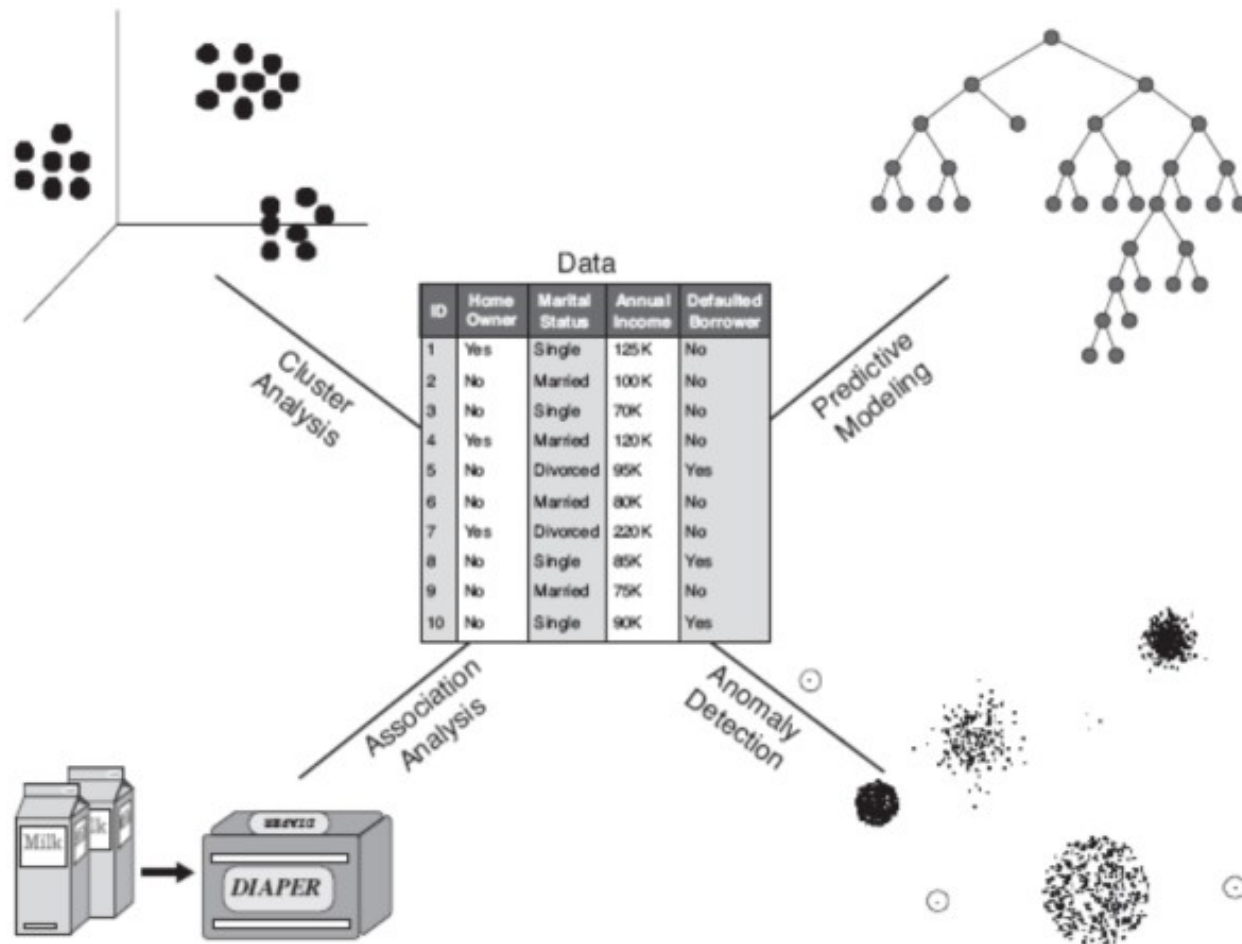
Data Mining



Data Mining



Data Mining



- Preprocesamiento de datos -

Fuentes de Datos

<i>Tid</i>	<i>Refund</i>	<i>Marital Status</i>	<i>Taxable Income</i>	<i>Defaulted Borrower</i>
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

(a) Record data.

<i>TID</i>	<i>ITEMS</i>
1	Bread, Soda, Milk
2	Beer, Bread
3	Beer, Soda, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Soda, Diaper, Milk

(b) Transaction data.

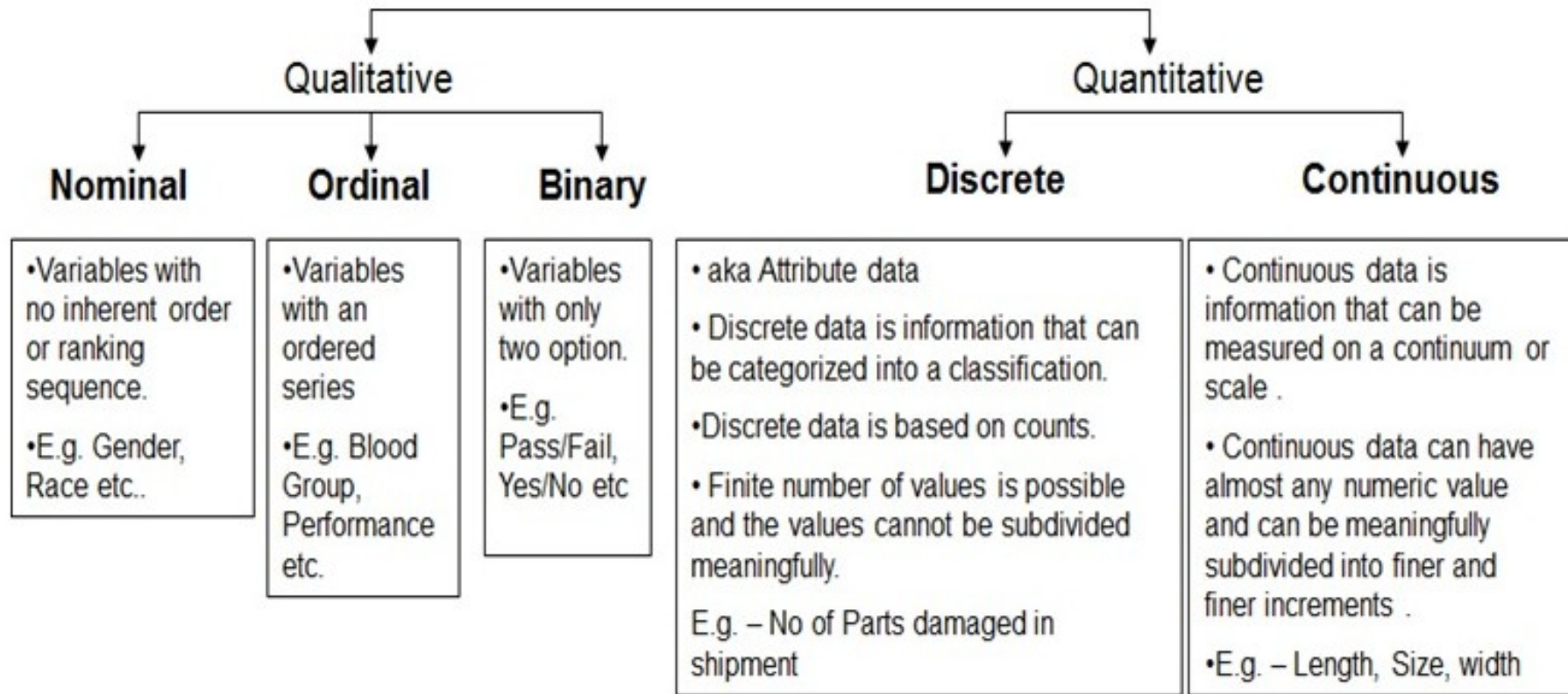
<i>Projection of x Load</i>	<i>Projection of y Load</i>	<i>Distance</i>	<i>Load</i>	<i>Thickness</i>
10.23	5.27	15.22	27	1.2
12.65	6.25	16.22	22	1.1
13.54	7.23	17.34	23	1.2
14.27	8.43	18.45	25	0.9

(c) Data matrix.

	<i>team</i>	<i>coach</i>	<i>play</i>	<i>ball</i>	<i>score</i>	<i>game</i>	<i>win</i>	<i>lost</i>	<i>timeout</i>	<i>season</i>
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

(d) Document-term matrix.

Tipos de características



Características cuantitativas

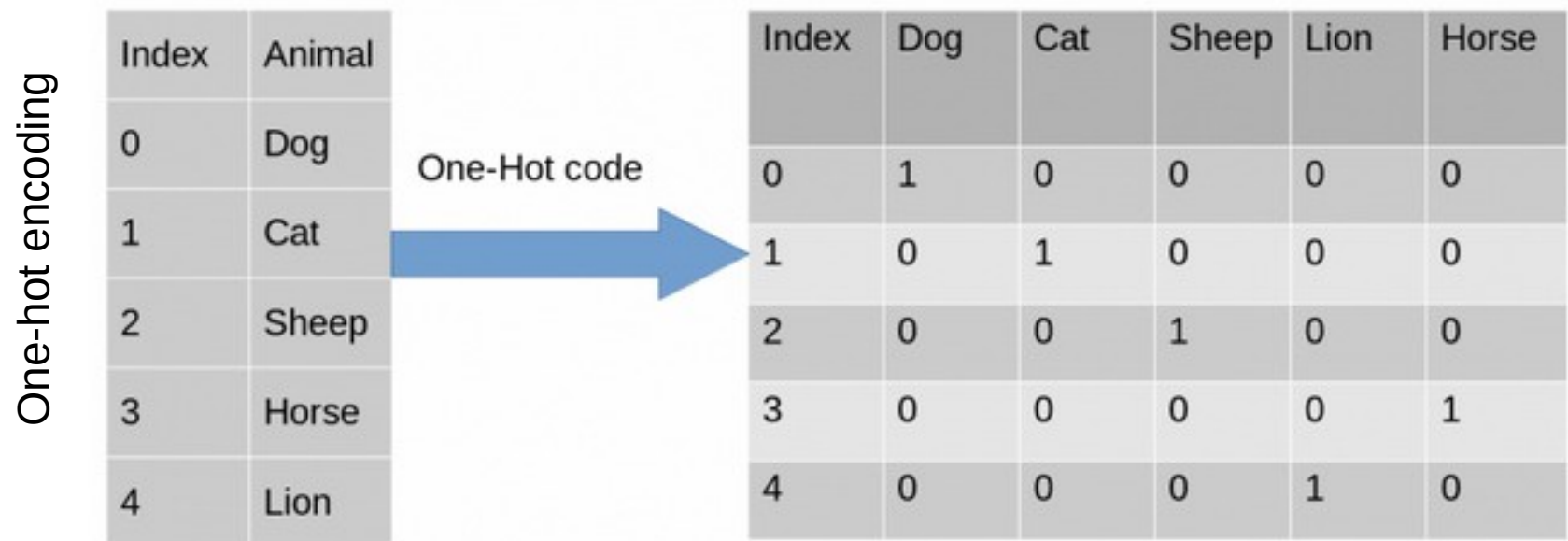
Normalización:

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \longrightarrow [0, 1]$$

$$X_{m-norm} = \frac{X - \mu}{X_{max} - X_{min}} \longrightarrow \text{¿Intervalo?}$$

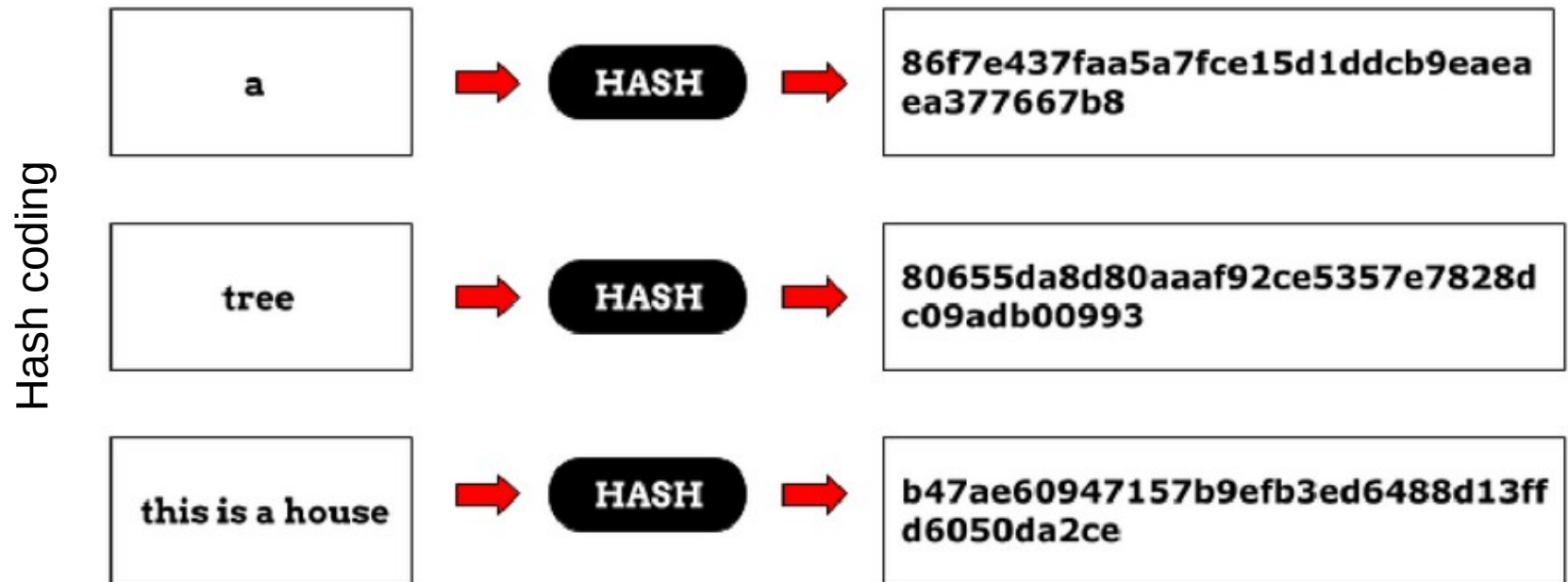
Características cualitativas

Codificación:



Características cualitativas

Codificación:



Vectores y características

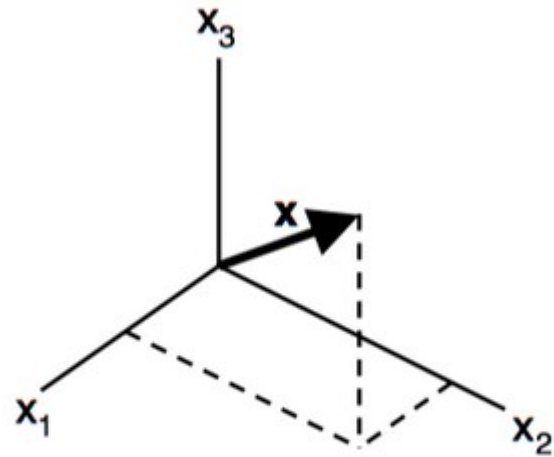
$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_d \end{bmatrix}$$

Feature vector

Vectores y características

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_d \end{bmatrix}$$

Feature vector



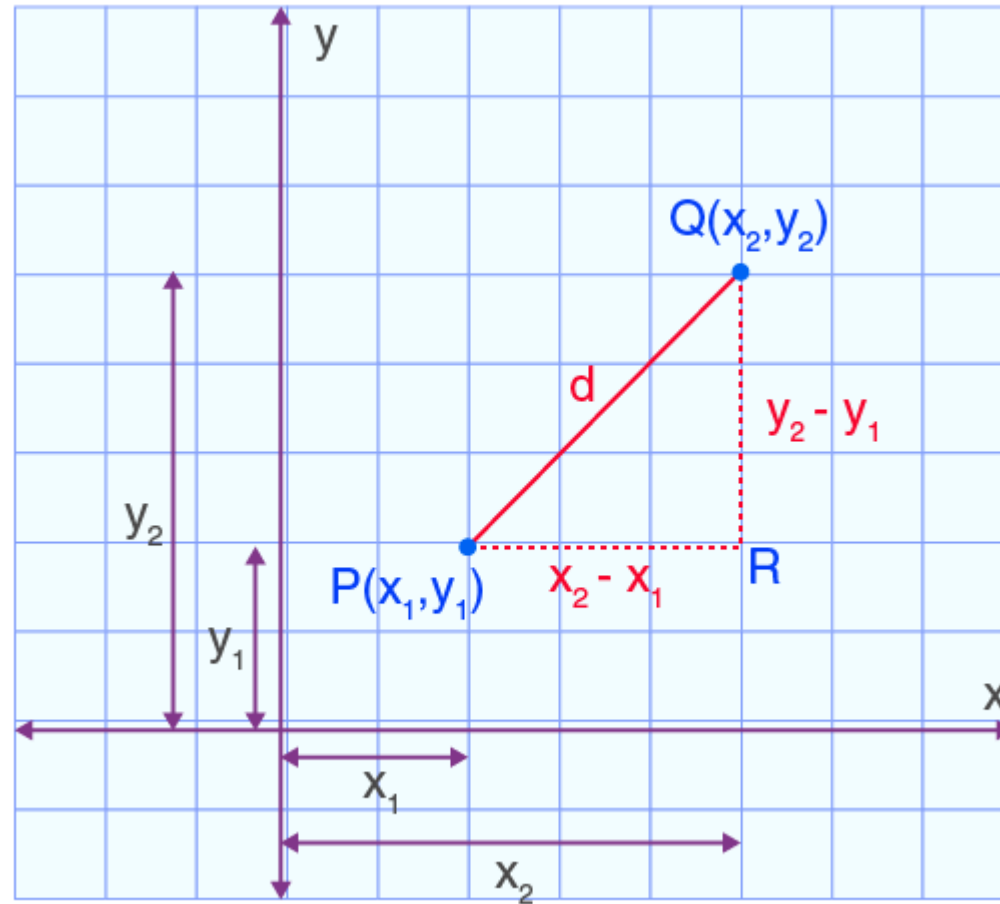
Feature space (3D)

- Distancia y proximidad -



Distancia Euclidea

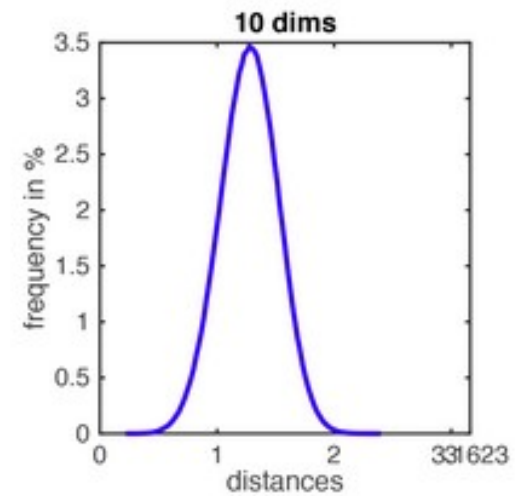
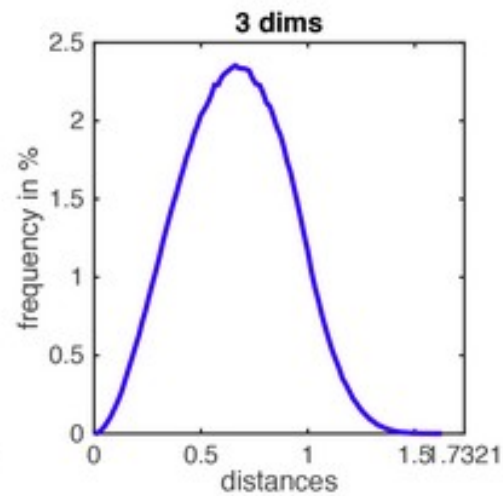
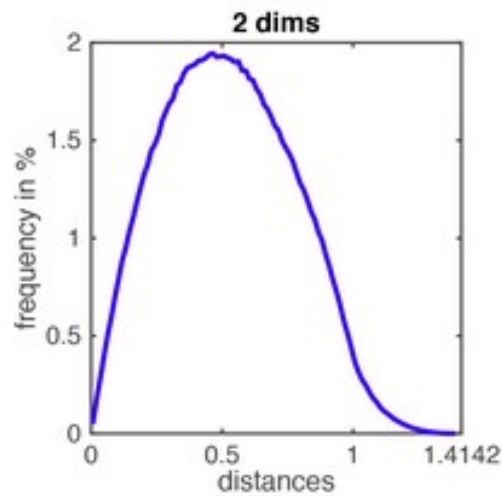
2D:



n-dimensional:
$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$



Distancia Euclideana

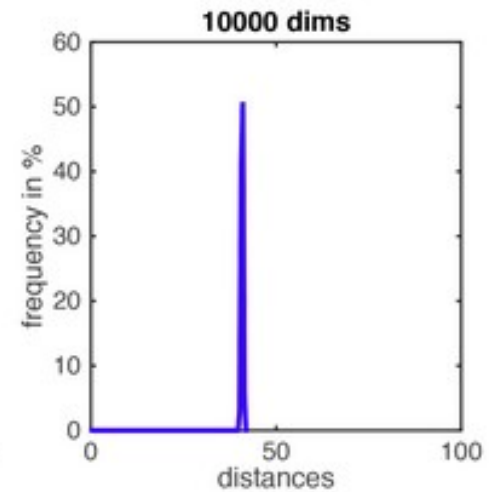
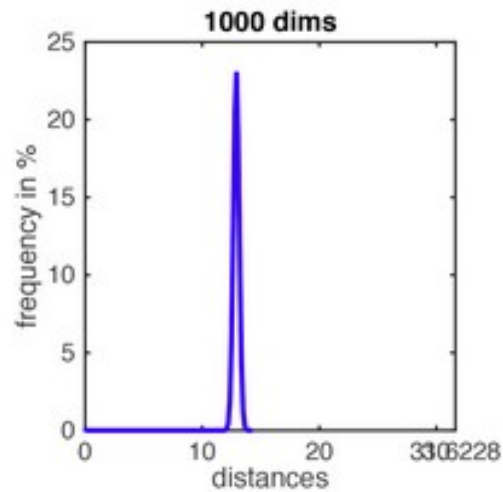
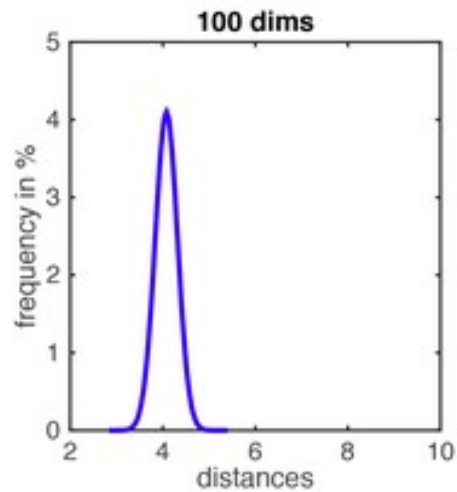


$$\sqrt{2} = 1.41$$

$$\sqrt{3} = 1.73$$

$$\sqrt{10} = 3.16$$

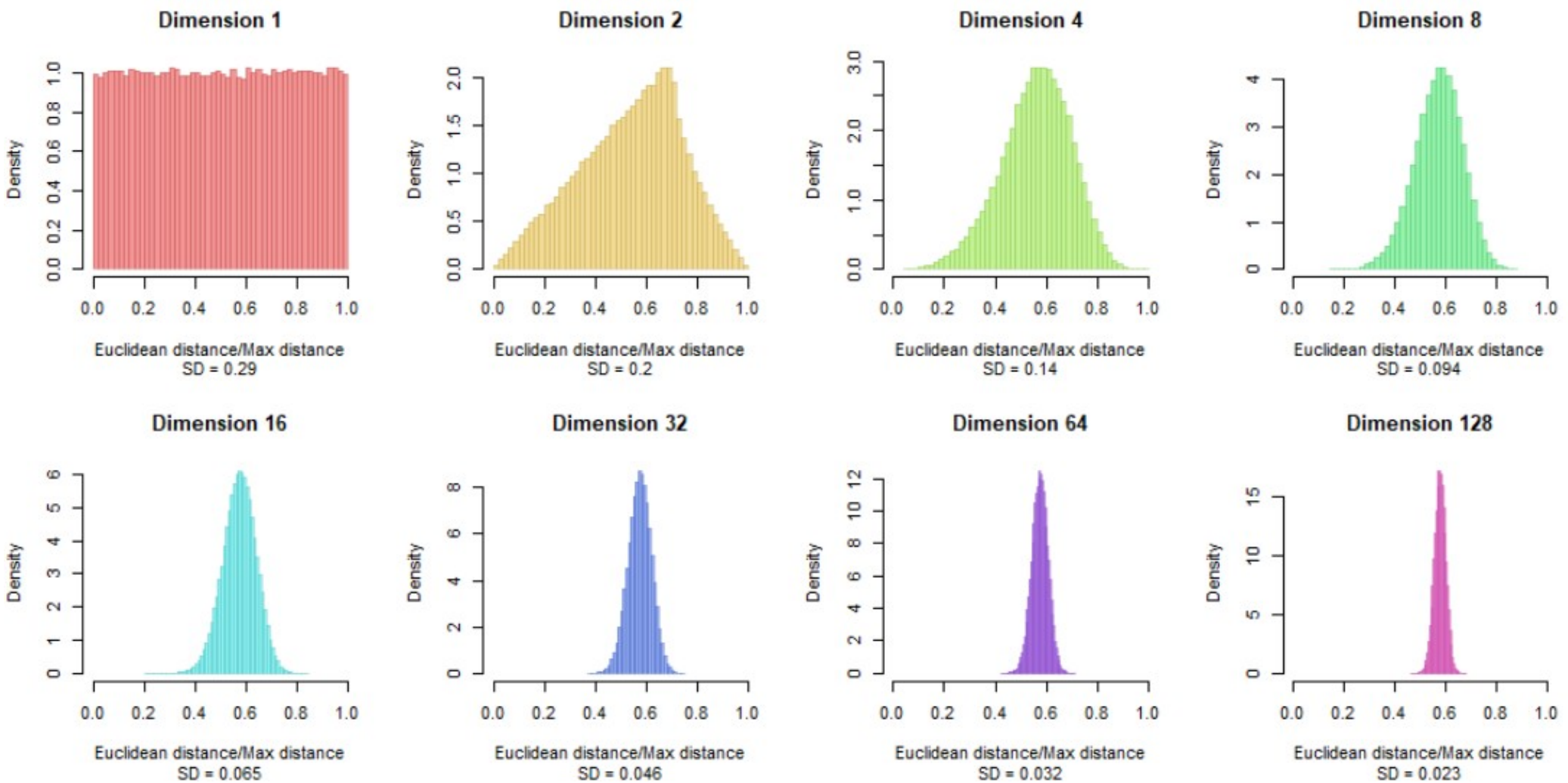
$$\sqrt{1000} = 31.6$$



Maldición de la dimensionalidad

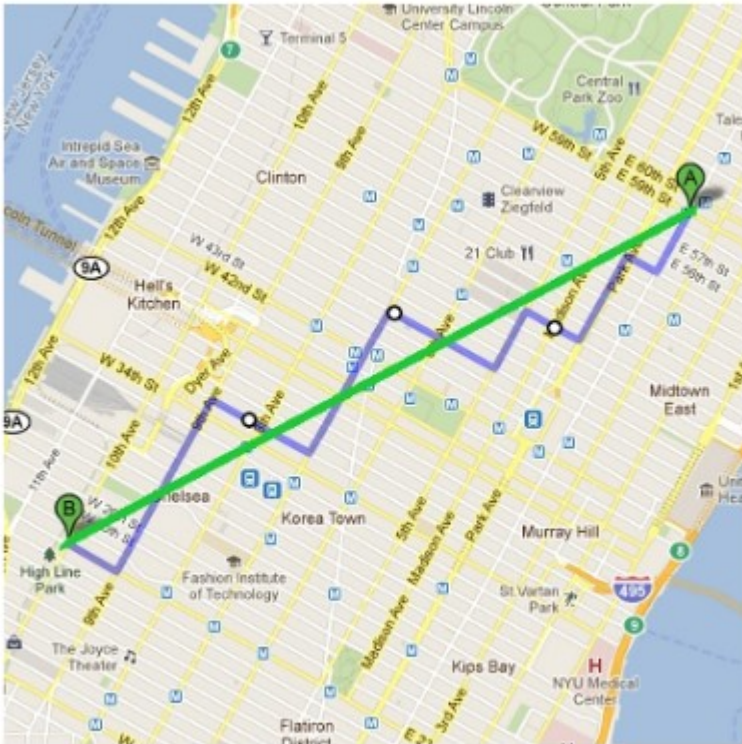


Distancia Euclideana



Maldición de la dimensionalidad (normalizado)

Distancias



Distancia Manhattan

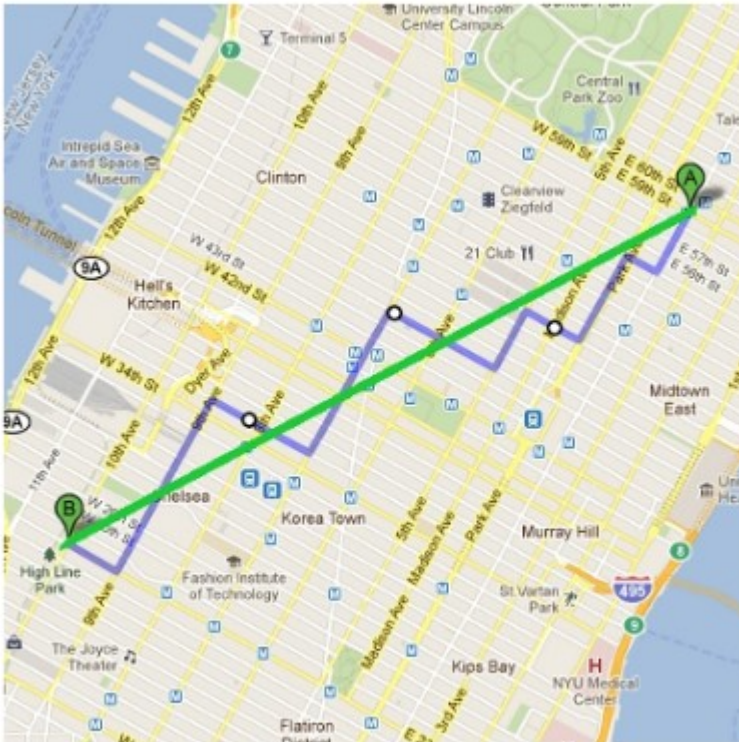
Distancias



Generalización (Minkowski):

$$Dist(\bar{X}, \bar{Y}) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}$$

← *Manhattan* ($p = 1$)



Distancia Manhattan

Distancias

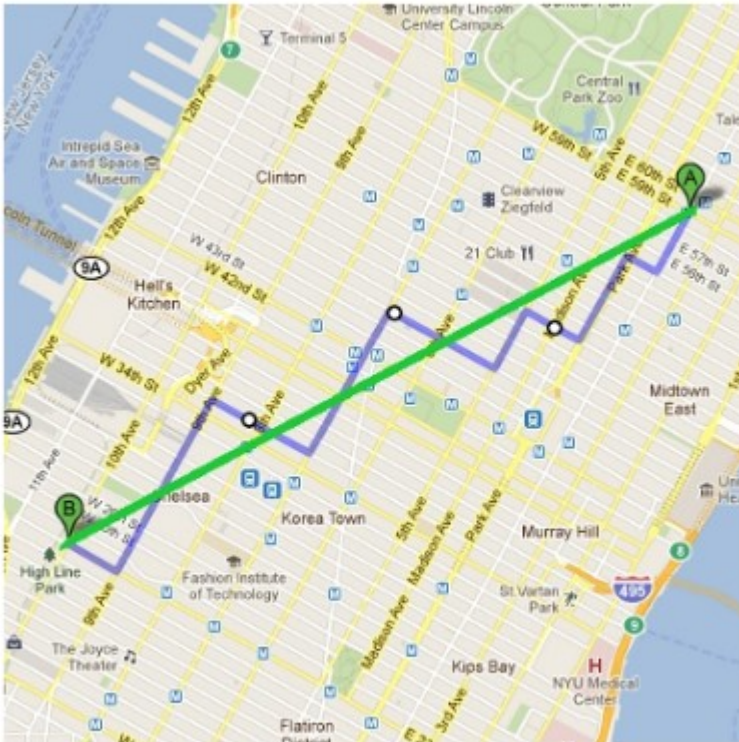


Generalización (Minkowski):

$$Dist(\bar{X}, \bar{Y}) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}$$

← *Manhattan* ($p = 1$)

Euclidean ($p = 2$)

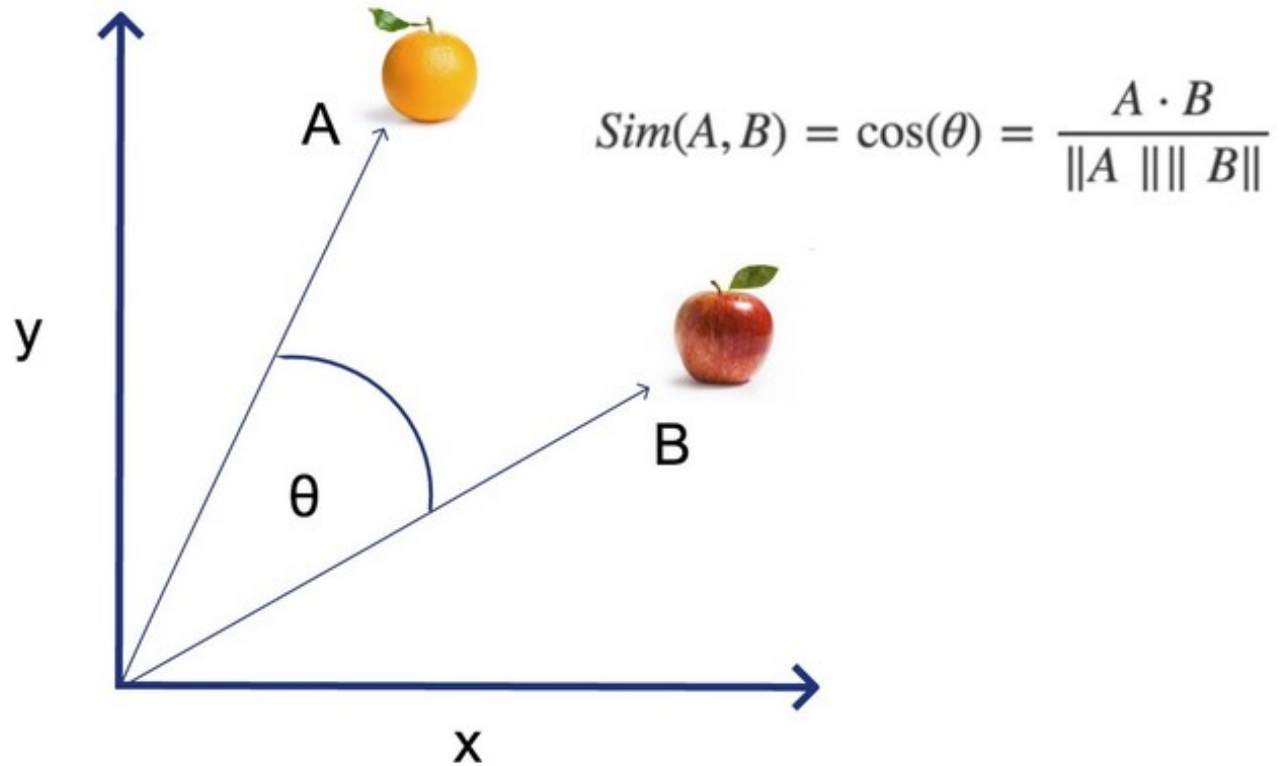


Distancia Manhattan

Proximidades

Proximidad de vectores de alta dimensionalidad:

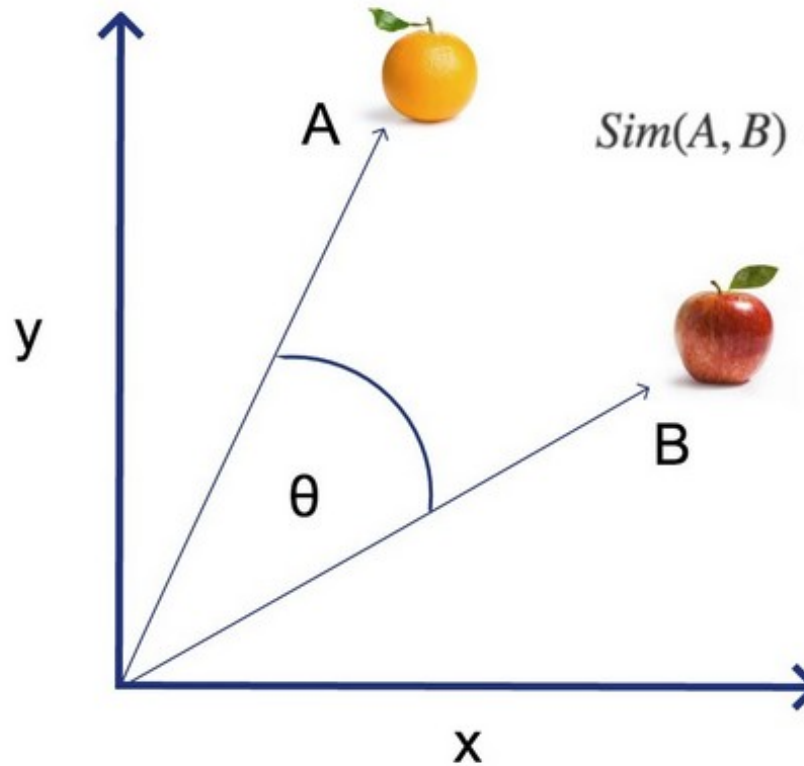
Coseno:



Proximidades

Proximidad de vectores de alta dimensionalidad:

Coseno:



$$Sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

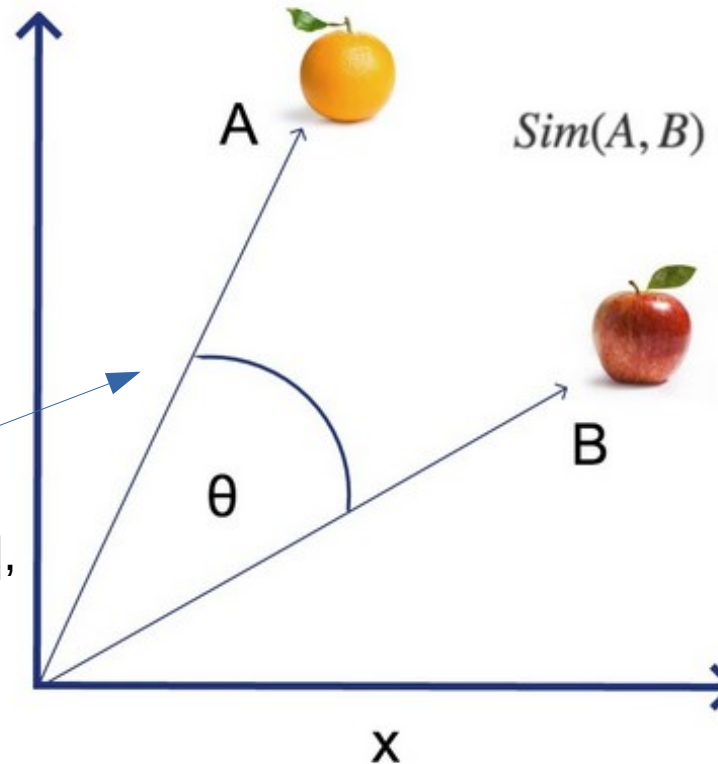
$$\cos(\bar{X}, \bar{Y}) = \frac{\sum_{i=1} x_i \cdot y_i}{\sqrt{\sum_{i=1}^d x_i^2} \cdot \sqrt{\sum_{i=1}^d y_i^2}}$$

Proximidades

Proximidad de vectores de alta dimensionalidad:

Coseno:

Si las características están en $[0,1]$,
los ángulos están en $[0^\circ, 90^\circ]$



$$Sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

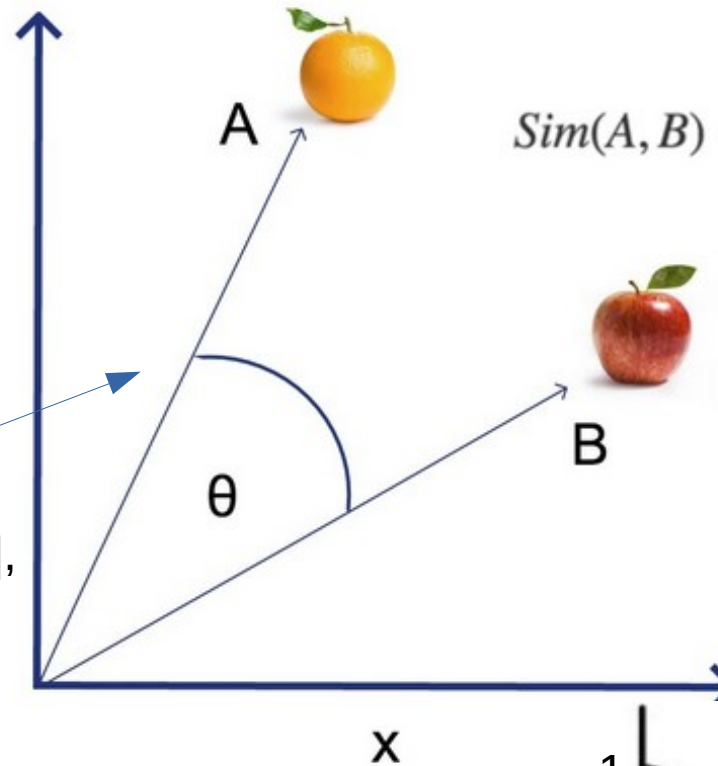
$$\cos(\bar{X}, \bar{Y}) = \frac{\sum_{i=1}^d x_i \cdot y_i}{\sqrt{\sum_{i=1}^d x_i^2} \cdot \sqrt{\sum_{i=1}^d y_i^2}}$$

Proximidades

Proximidad de vectores de alta dimensionalidad:

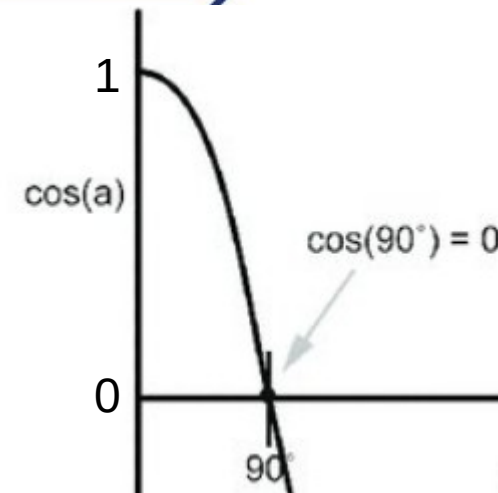
Coseno:

Si las características están en $[0,1]$, los ángulos están en $[0^\circ, 90^\circ]$

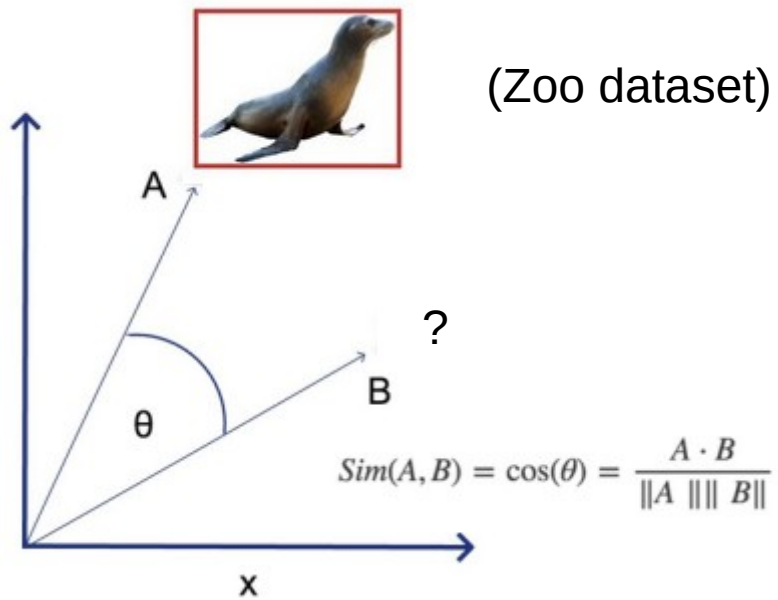


$$Sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

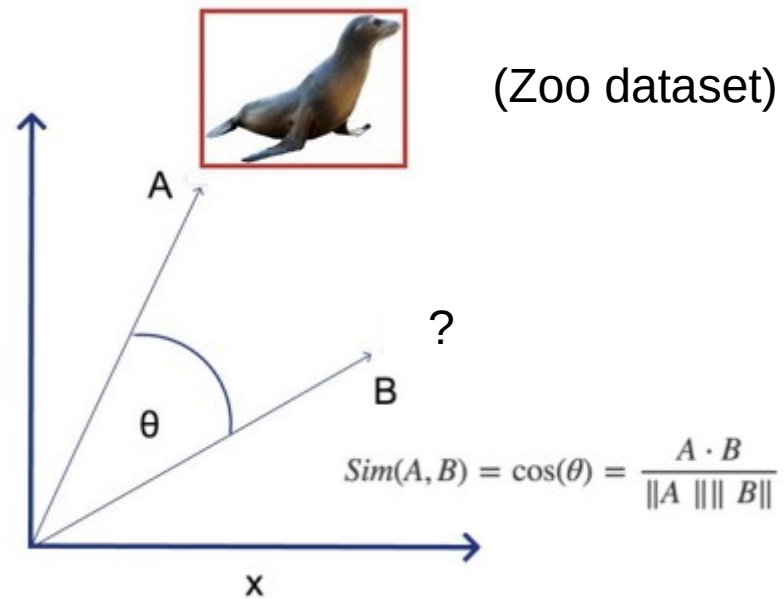
$$\cos(\bar{X}, \bar{Y}) = \frac{\sum_{i=1}^d x_i \cdot y_i}{\sqrt{\sum_{i=1}^d x_i^2} \cdot \sqrt{\sum_{i=1}^d y_i^2}}$$



Proximidades



Proximidades

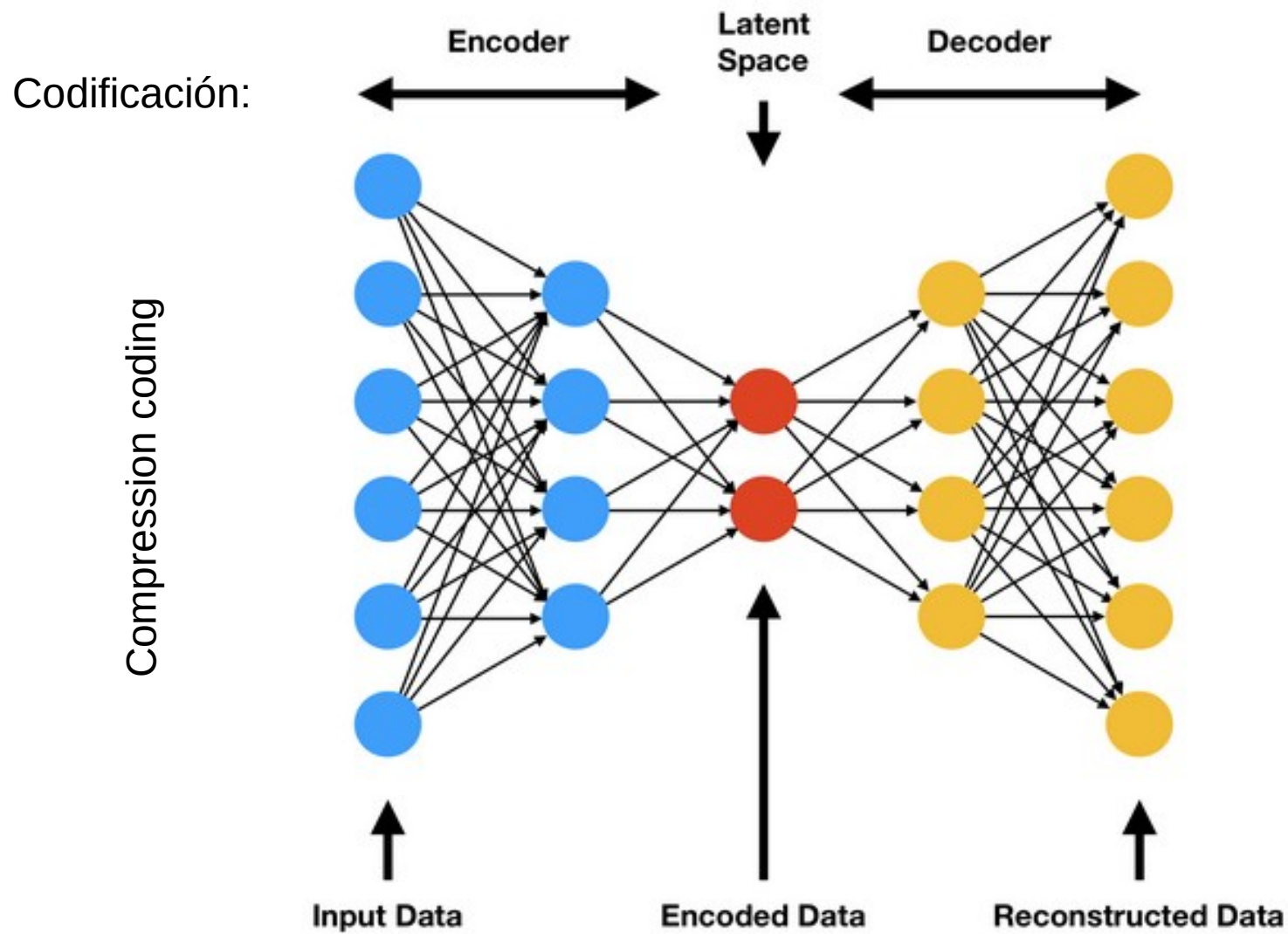


dolphin: 0.875 mink: 0.875 porpoise: 0.875 seal: 0.875 boar: 0.8125 cheetah: 0.8125 leopard: 0.8125 lion: 0.8125



- PCA -

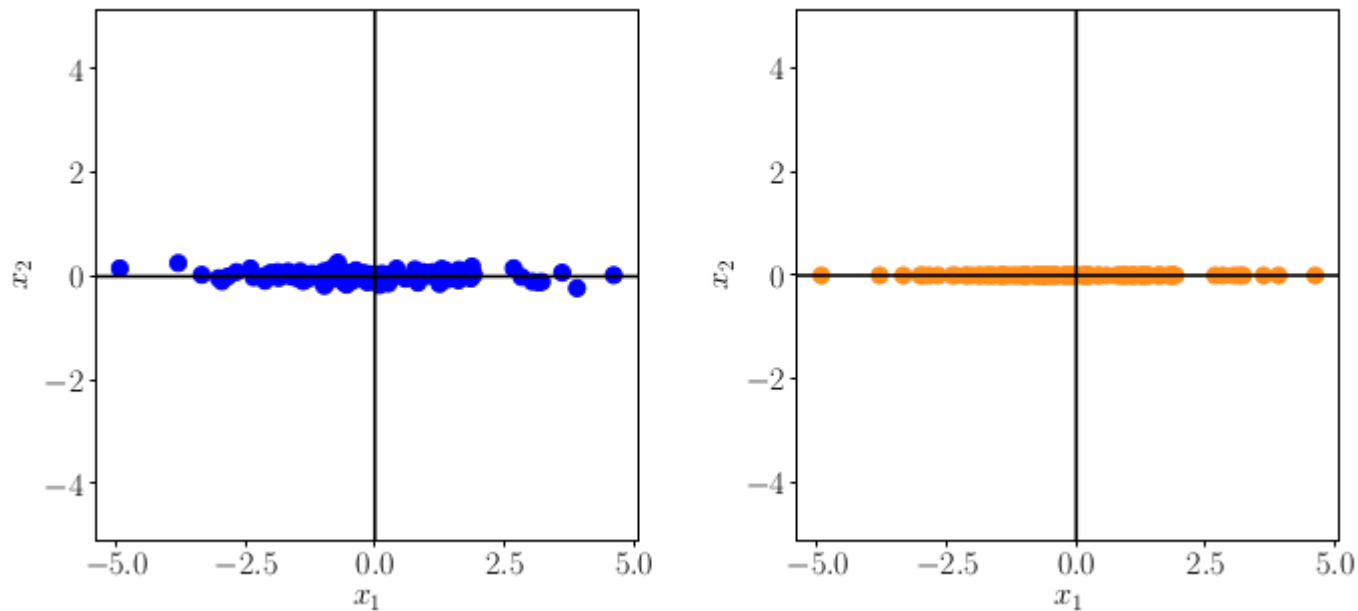
Proyección



Ej.: PCA

- UC - M. Mendoza -

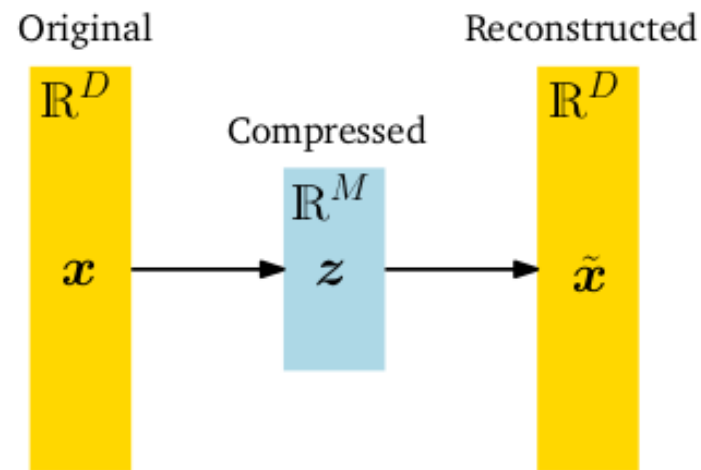
Análisis de Componentes Principales (PCA)



X1 retiene la mayor parte de la varianza por lo que remover x_2 es neutro en términos de compresión.

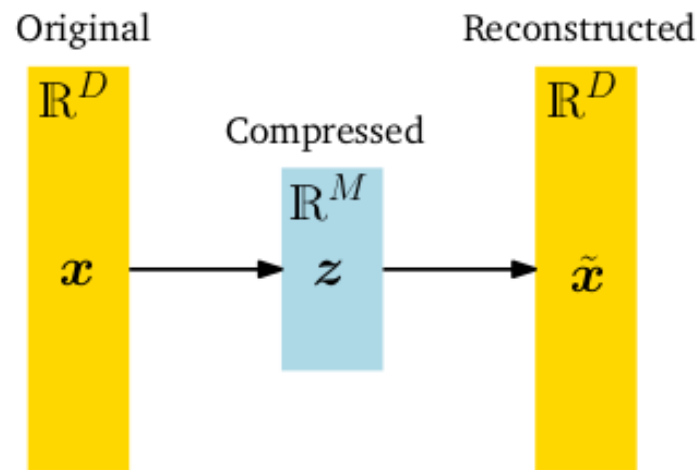
Análisis de Componentes Principales (PCA)

dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$



Análisis de Componentes Principales (PCA)

dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

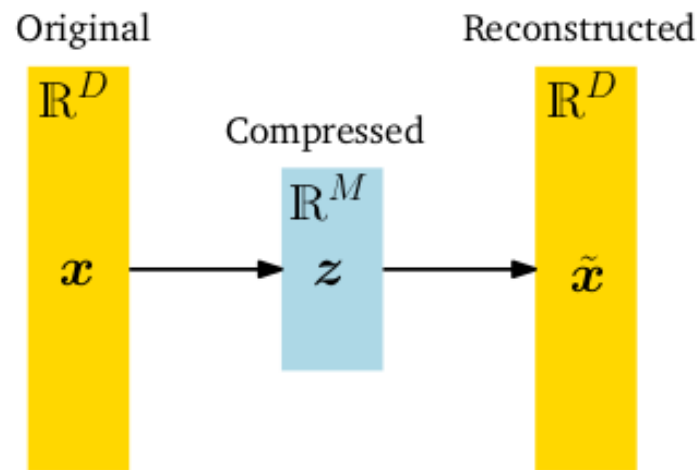


$$\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M \longrightarrow \text{Baja dimensionalidad}$$

└─► Base de la descomposición $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$.

Análisis de Componentes Principales (PCA)

dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$



$$\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M \longrightarrow \text{Baja dimensionalidad}$$

└─► Base de la descomposición $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$.

└─► $\mathbf{b}_i^\top \mathbf{b}_j = 0$ y $\mathbf{b}_i^\top \mathbf{b}_i = 1$.

La proyección se calcula usando la SVD.

Análisis de Componentes Principales (PCA)

Proceso iterativo:

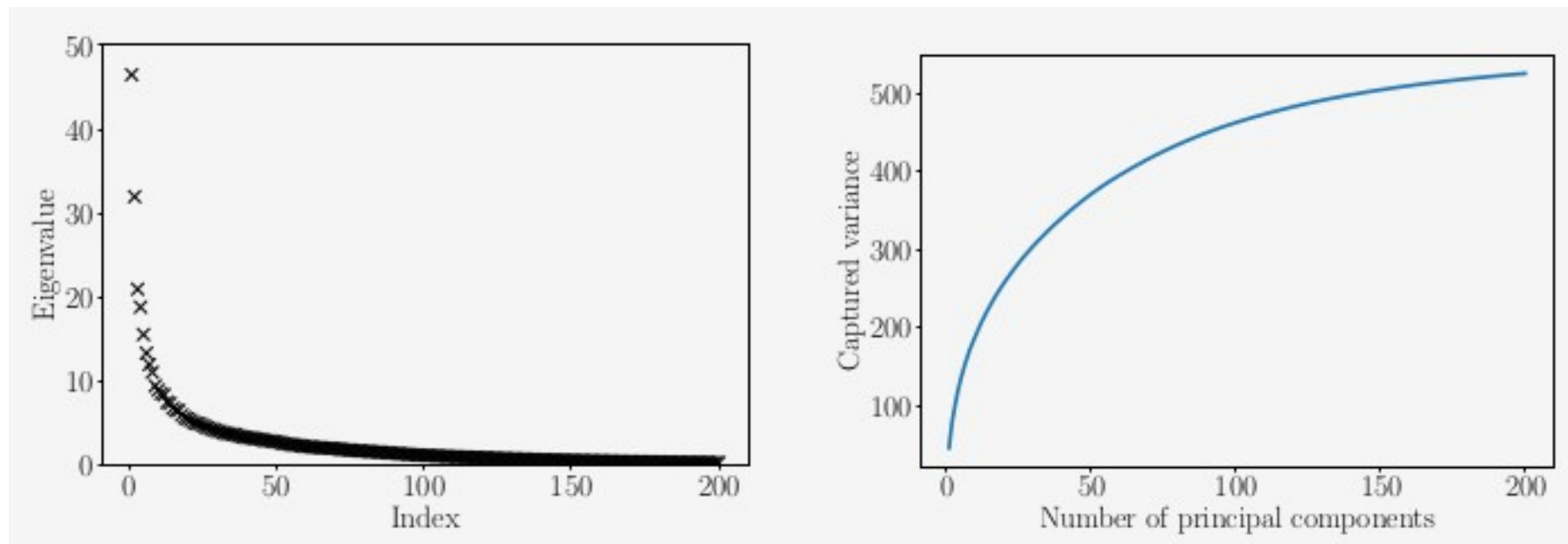
$$\hat{X} := X - \sum_{i=1}^{m-1} b_i b_i^\top X = X - B_{m-1} X, \quad \text{con } X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$
$$\text{y } B_{m-1} := \sum_{i=1}^{m-1} b_i b_i^\top$$

Análisis de Componentes Principales (PCA)

Proceso iterativo:

$$\hat{X} := X - \sum_{i=1}^{m-1} b_i b_i^\top X = X - B_{m-1} X, \quad \text{con } X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$

$$\text{y } B_{m-1} := \sum_{i=1}^{m-1} b_i b_i^\top$$



Análisis de Componentes Principales (PCA)

- Aspectos prácticos:

- Usa la full SVD (LAPACK) para datos densos.
- Usa la SVD truncada (ARPACK) para datos dispersos.

- Implementaciones:

- Python: sklearn

<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>