



Minería de Datos

¿Qué vamos a ver?

Representación (PCA, t-SNE, ...)

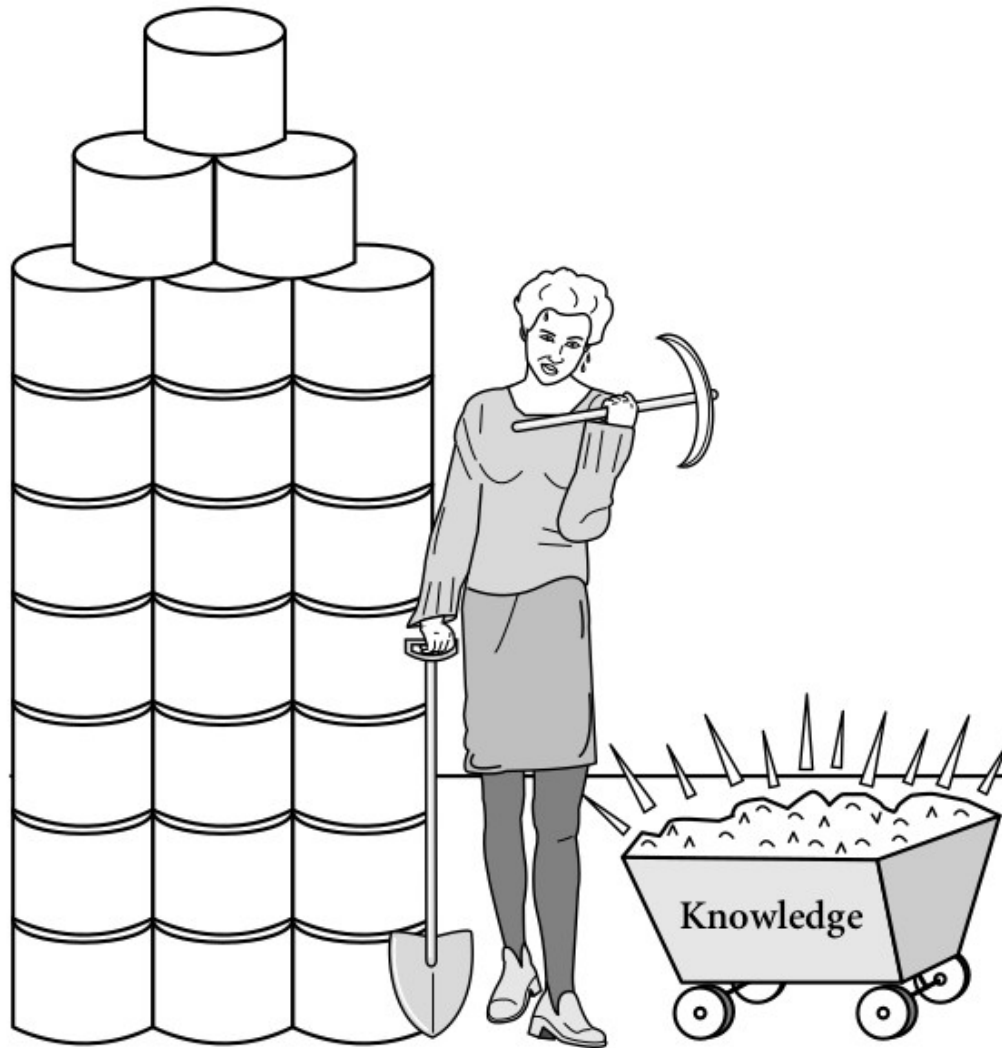
Asociación (apriori)

Clustering (k-means, HAC)

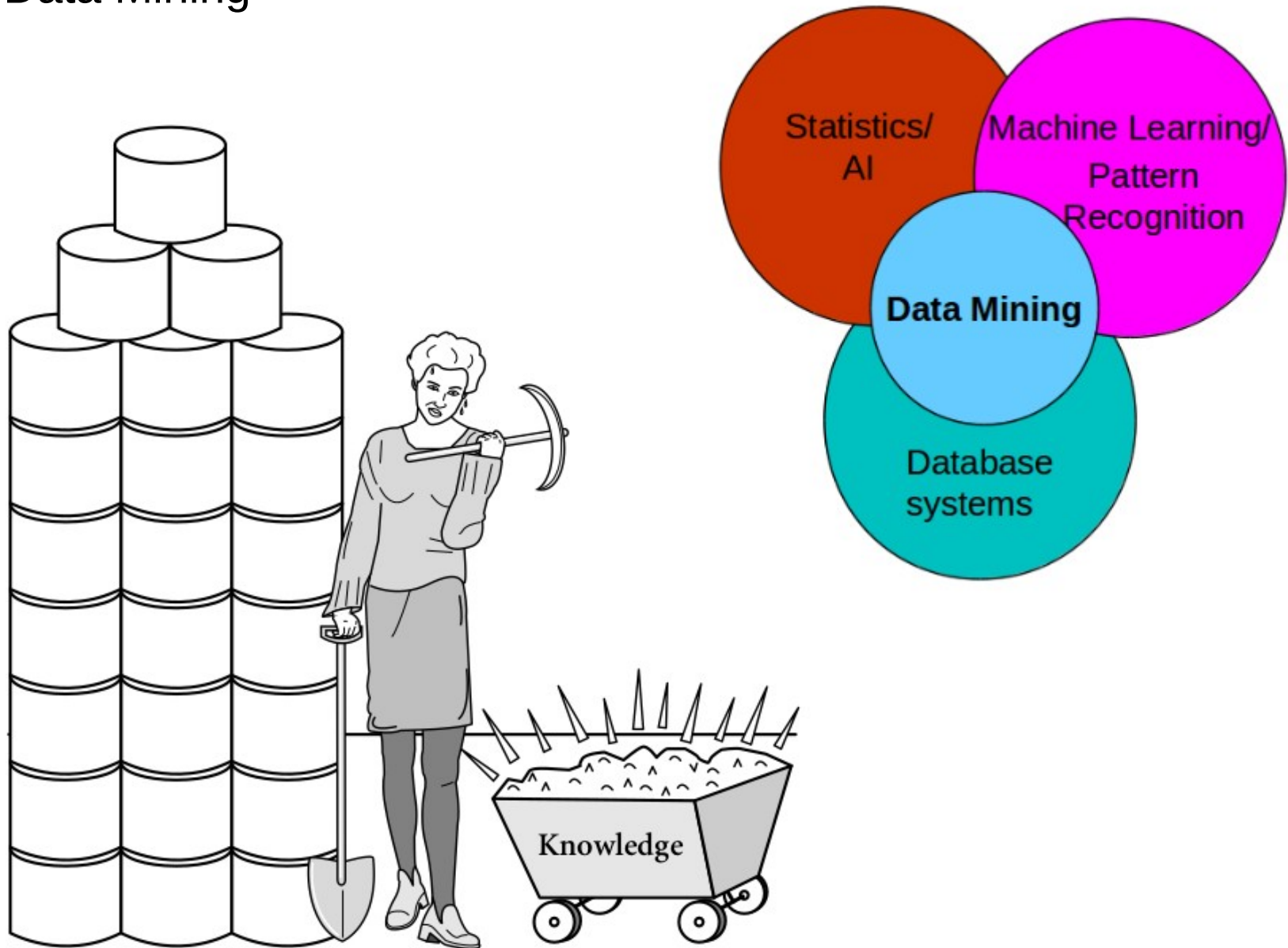
Clasificación

Auto-encoders

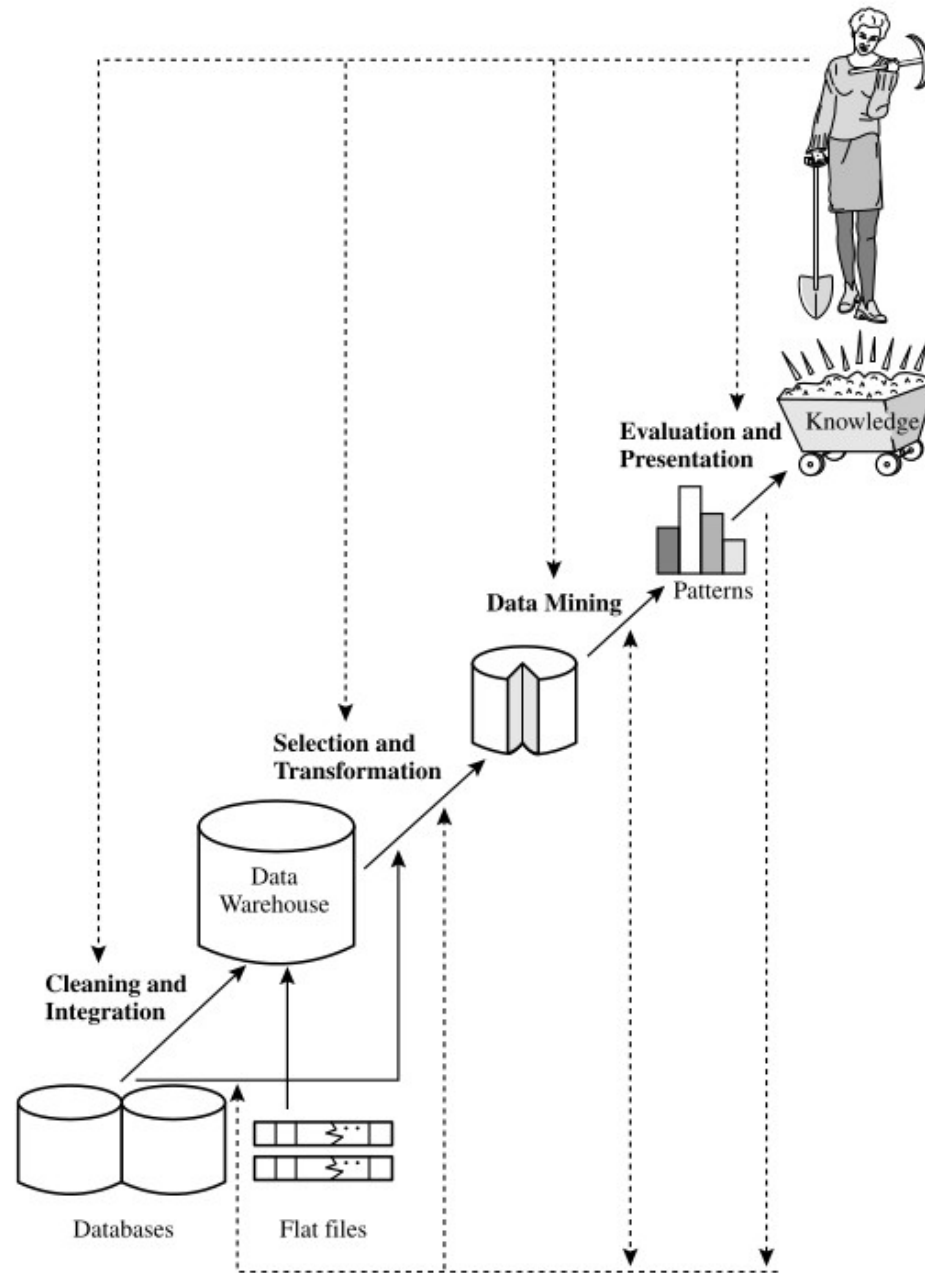
Data Mining



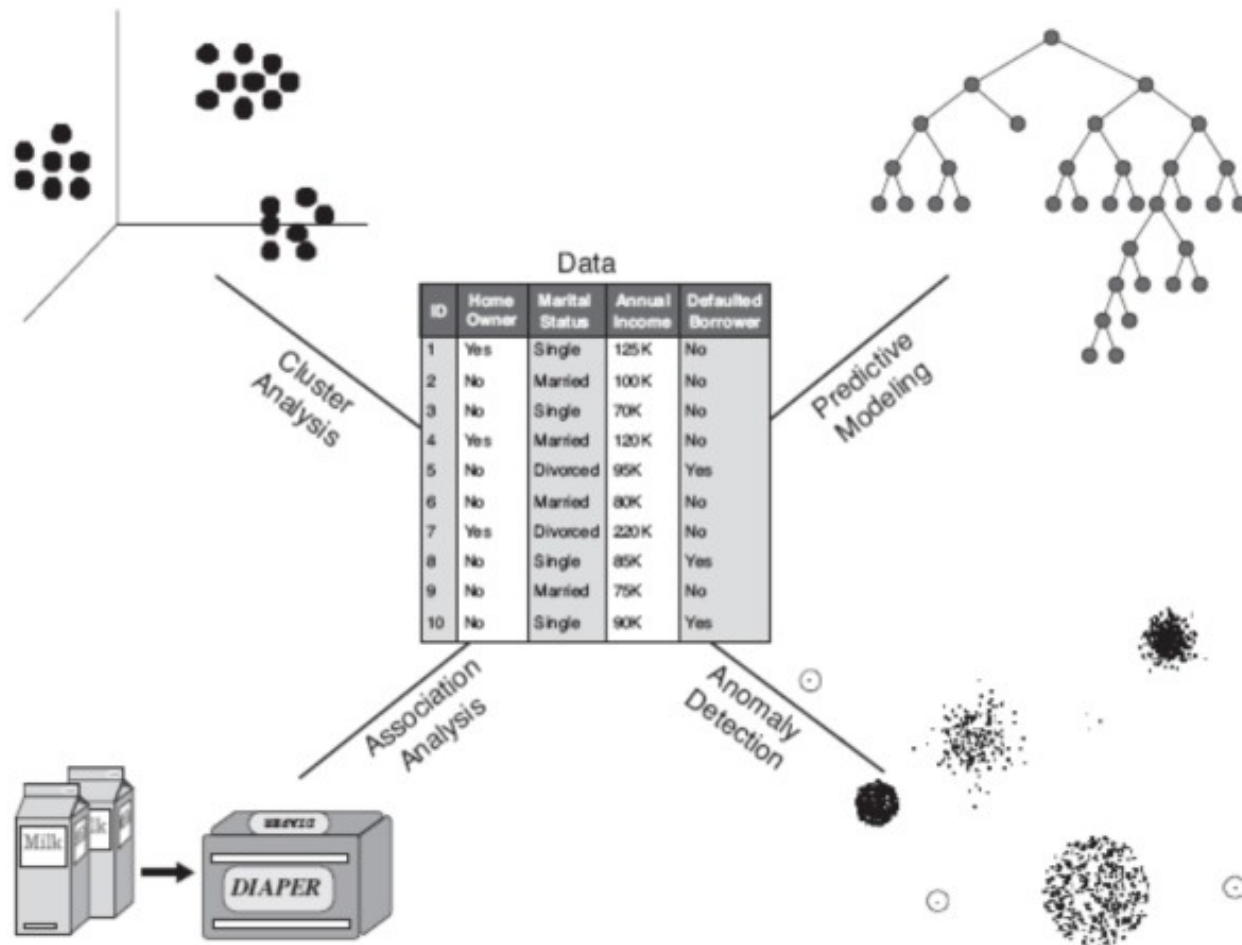
Data Mining



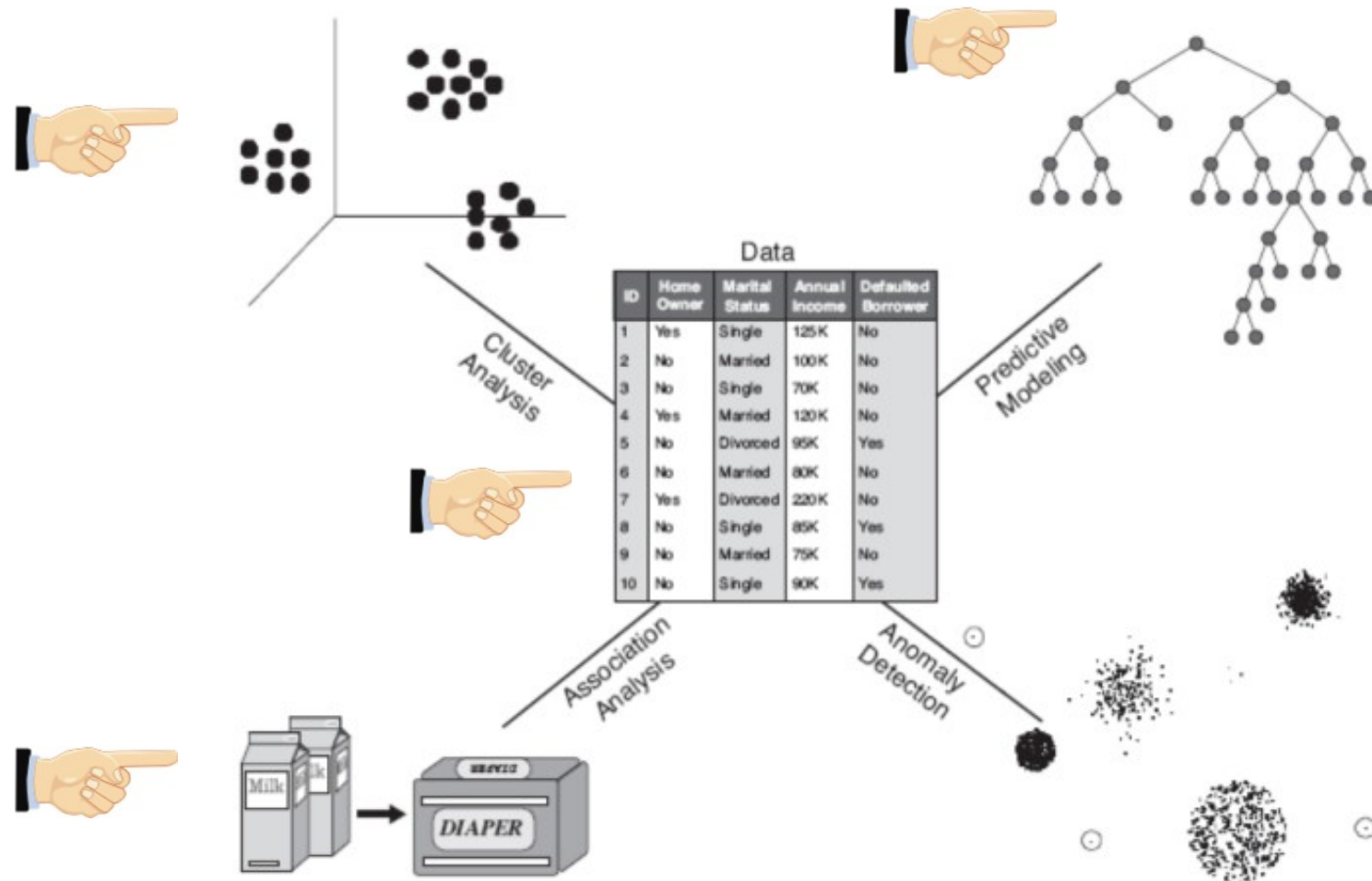
Data Mining



Data Mining



Data Mining



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.

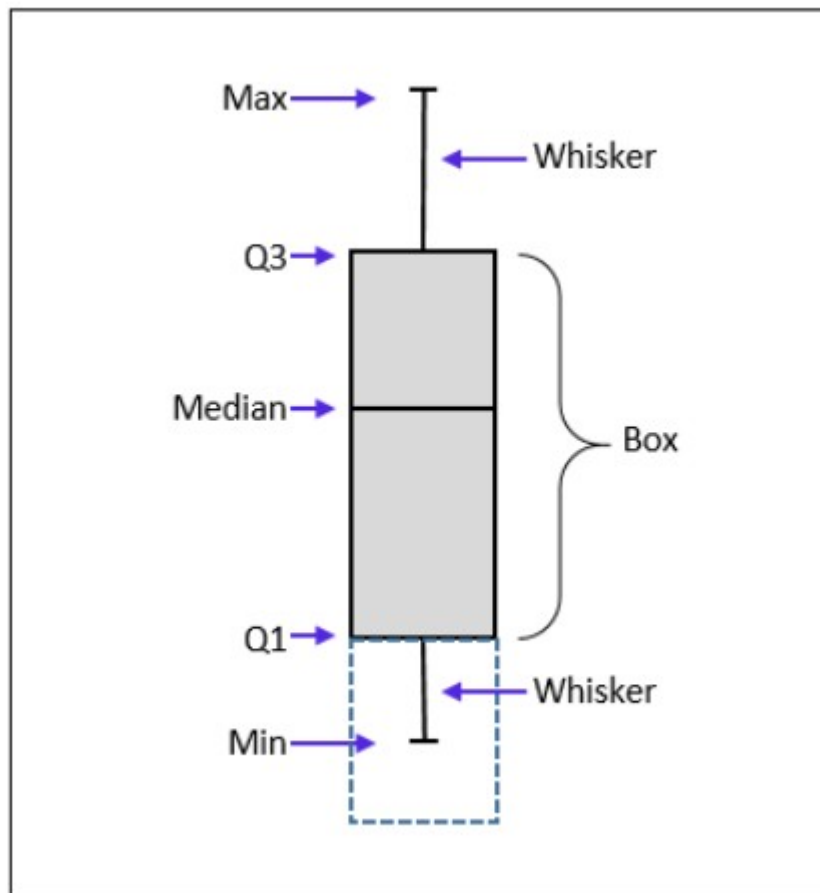
<i>Projection of x Load</i>	<i>Projection of y Load</i>	<i>Distance</i>	<i>Load</i>	<i>Thickness</i>
10.23	5.27	15.22	27	1.2
12.65	6.25	16.22	22	1.1
13.54	7.23	17.34	23	1.2
14.27	8.43	18.45	25	0.9

(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

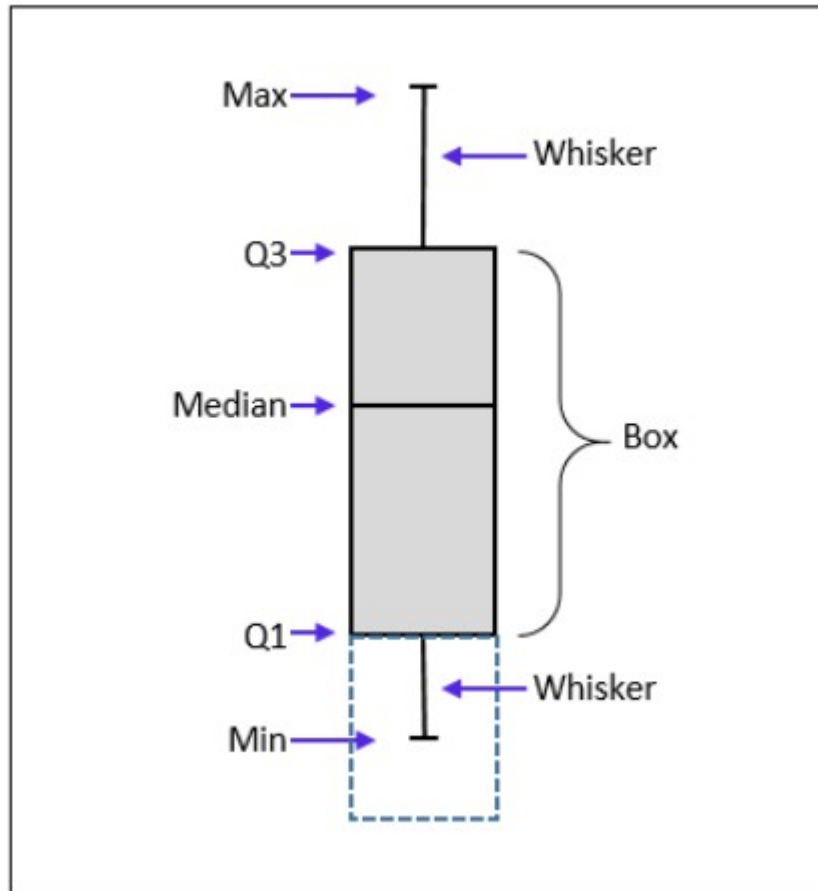
(d) Document-term matrix.

Herramientas de análisis exploratorio

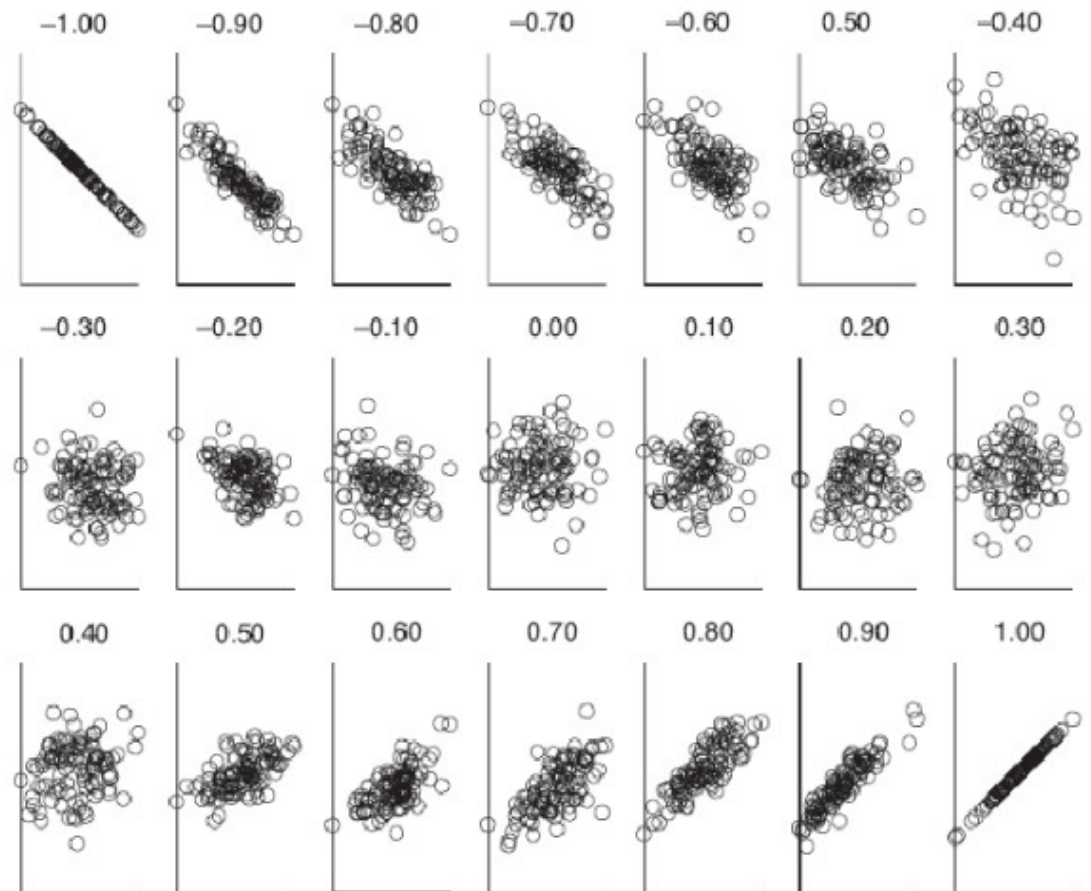


Box Plots

Herramientas de análisis exploratorio



Box Plots

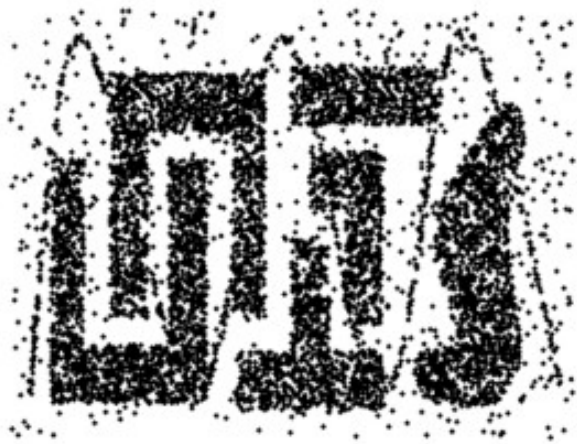


Scatter Plots

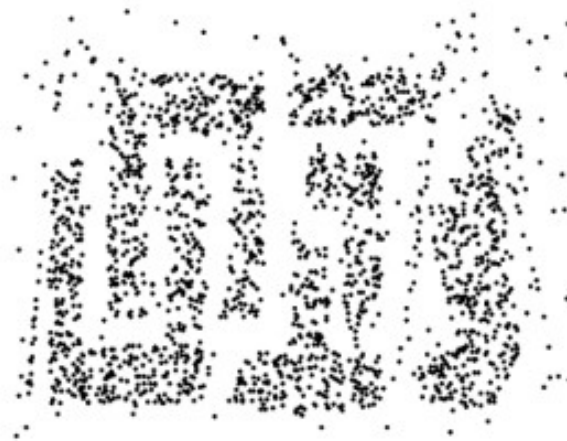
Captura de Datos

- ▶ Muestreo aleatorio: equiprobable en la selección de datos.
- ▶ Muestreo con reemplazo: los datos no son removidos de la colección (puede ser muestreado más de una vez).
- ▶ Muestreo sin reemplazo: los datos son removidos de la colección si son muestreados (puede ser muestreado solo una vez).
- ▶ Muestreo estratificado: los datos son segmentados de acuerdo a algún criterio. Luego, en cada segmento se conduce muestreo aleatorio.

Captura de Datos



8000 points



2000 Points

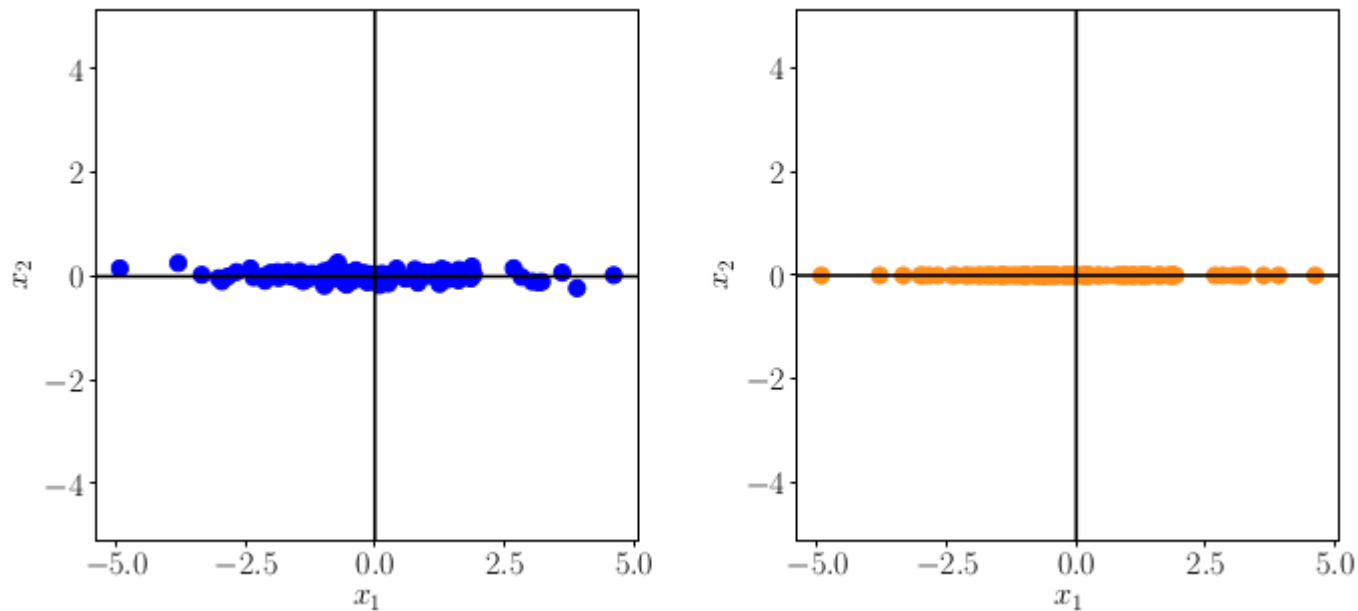


500 Points

Riesgo empírico en muestreo

Representaciones y reducción de dimensionalidad

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)

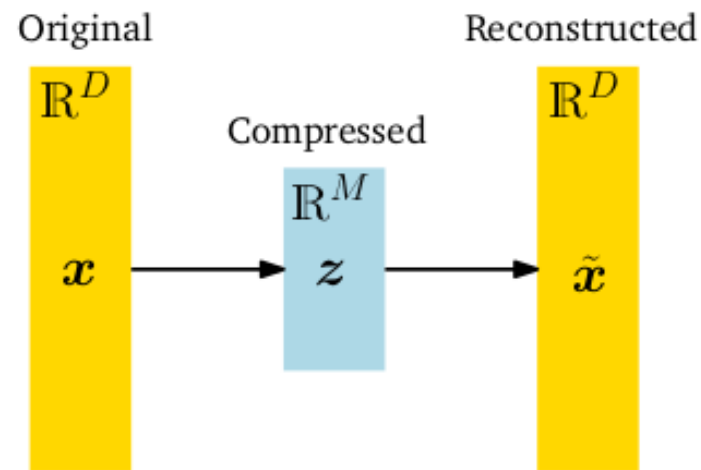
dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

Análisis de Componentes Principales (PCA)

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

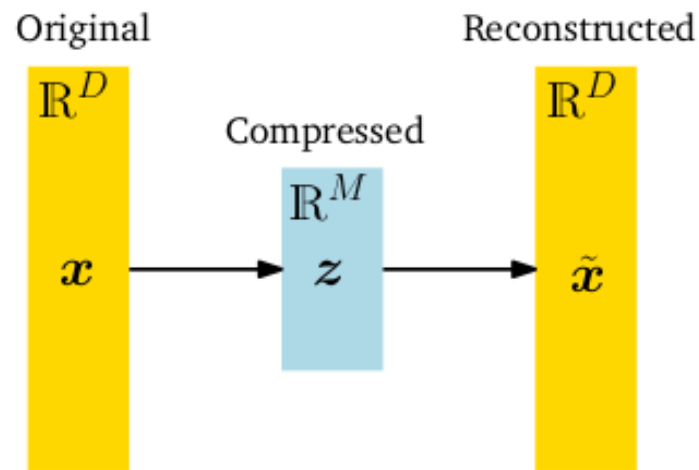
Análisis de Componentes Principales (PCA)

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



Análisis de Componentes Principales (PCA)

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

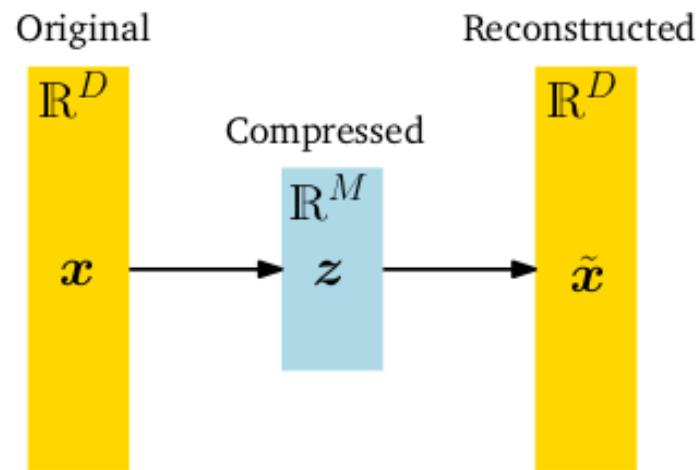


$$\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M \longrightarrow \text{Baja dimensionalidad}$$

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



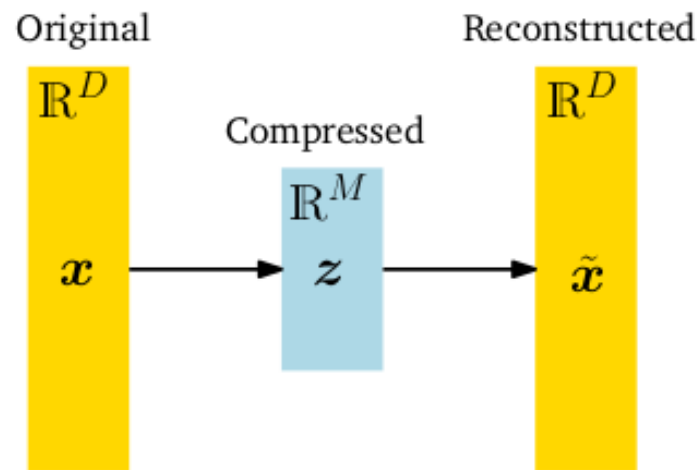
$$\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M \longrightarrow \text{Baja dimensionalidad}$$

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

$$\longrightarrow \mathbf{b}_i^\top \mathbf{b}_j = 0 \quad \text{y} \quad \mathbf{b}_i^\top \mathbf{b}_i = 1.$$

Análisis de Componentes Principales (PCA)

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



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

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

$$\longrightarrow \mathbf{b}_i^\top \mathbf{b}_j = 0 \quad \text{y} \quad \mathbf{b}_i^\top \mathbf{b}_i = 1.$$

Minimizar la pérdida de información implica capturar la mayor cantidad de varianza en la descomposición.

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[\mathbf{z}] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[z] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Por lo tanto, corregimos los datos a media 0.

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[z] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Por lo tanto, corregimos los datos a media 0.

PCA (enfoque secuencial): Comenzamos buscando $\mathbf{b}_1 \in \mathbb{R}^D$:

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[z] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Por lo tanto, corregimos los datos a media 0.

PCA (enfoque secuencial): Comenzamos buscando $\mathbf{b}_1 \in \mathbb{R}^D$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2$$

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[\mathbf{z}] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Por lo tanto, corregimos los datos a media 0.

PCA (enfoque secuencial): Comenzamos buscando $\mathbf{b}_1 \in \mathbb{R}^D$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 . \text{ Dado que } z_{1n} = \mathbf{b}_1^\top \mathbf{x}_n ,$$

$$\begin{aligned} V_1 &= \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_1 \\ &= \mathbf{b}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 , \end{aligned}$$

Análisis de Componentes Principales (PCA)

Importante: la varianza de la descomposición es independiente de la media.

$$\mathbb{V}_z[z] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}]$$

Por lo tanto, corregimos los datos a media 0.

PCA (enfoque secuencial): Comenzamos buscando $\mathbf{b}_1 \in \mathbb{R}^D$:


$$\begin{aligned} V_1 &:= \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 . \text{ Dado que } z_{1n} = \mathbf{b}_1^\top \mathbf{x}_n , \\ V_1 &= \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_1 \\ &= \mathbf{b}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 , \end{aligned}$$

Por lo tanto, maximizar la varianza corresponde a:

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano 

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\hookrightarrow \mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\hookrightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

$$\begin{aligned} = 0 \quad & \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1 , \\ & \mathbf{b}_1^\top \mathbf{b}_1 = 1 . \end{aligned}$$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

$$= 0 \quad \begin{aligned} \mathbf{S} \mathbf{b}_1 &= \lambda_1 \mathbf{b}_1 , \\ \mathbf{b}_1^\top \mathbf{b}_1 &= 1 . \end{aligned}$$

Reescribimos: $V_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

$$= 0 \quad \begin{aligned} \mathbf{S} \mathbf{b}_1 &= \lambda_1 \mathbf{b}_1 , \\ \mathbf{b}_1^\top \mathbf{b}_1 &= 1 . \end{aligned}$$

Se calcula usando la SVD

Reescribimos: $V_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

$$= 0 \quad \begin{aligned} \mathbf{S} \mathbf{b}_1 &= \lambda_1 \mathbf{b}_1 , \\ \mathbf{b}_1^\top \mathbf{b}_1 &= 1 . \end{aligned}$$

Se calcula usando la SVD

Reescribimos: $V_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1$

La reconstrucción es: $\tilde{\mathbf{x}}_n = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^\top \mathbf{x}_n \in \mathbb{R}^D$

Análisis de Componentes Principales (PCA)

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad , \text{ sujeto a } \|\mathbf{b}_1\|^2 = 1 .$$

Lagrangiano

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top , \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

$$= 0 \quad \begin{aligned} \mathbf{S} \mathbf{b}_1 &= \lambda_1 \mathbf{b}_1 , \\ \mathbf{b}_1^\top \mathbf{b}_1 &= 1 . \end{aligned}$$

Se calcula usando la SVD

Reescribimos: $V_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1$

La reconstrucción es: $\tilde{\mathbf{x}}_n = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^\top \mathbf{x}_n \in \mathbb{R}^D$

Análisis de Componentes Principales (PCA)

Proceso iterativo:

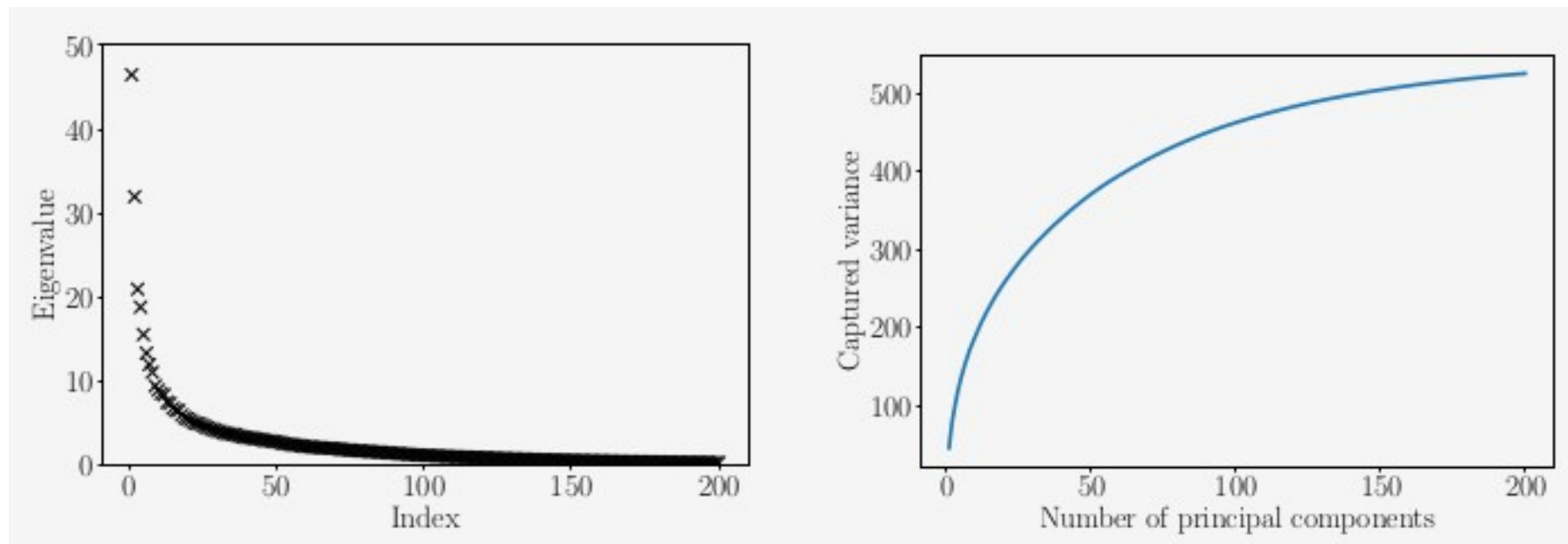
$$\hat{X} := X - \sum_{i=1}^{m-1} b_i b_i^\top X = X - B_{m-1} X, \quad \text{con } X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$
$$\text{y } B_{m-1} := \sum_{i=1}^{m-1} b_i b_i^\top$$

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.
- Se puede usar MLE (reconstrucción – data) para estimar el # de componentes.

- Implementaciones:

- Python: sklearn

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