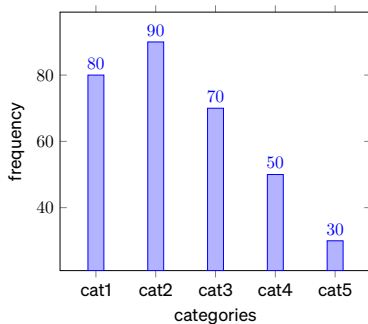POLIMI GRADUATE SCHOOL OF MANAGEMENT

# EXPLORATORY DATA ANALYSIS

PERCORSO EXECUTIVE DATA SCIENCE AND BUSINESS ANALYTICS

Mauricio Soto - mauricioabel.soto@polimi.it

# GRAPHICAL ANALYSIS CATEGORICAL ATTRIBUTE
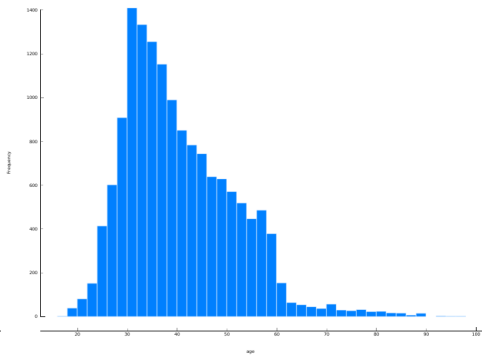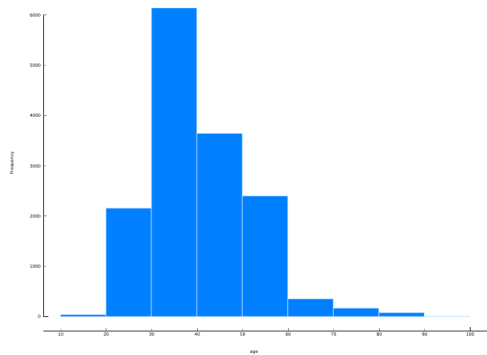


$$f_h = \frac{e_h}{m} = \frac{card\{i \in \mathcal{M} : x_i = cat_h\}}{m}$$

for large samples

$$f_h \approx P(x = cat_h)$$

# GRAPHICAL ANALYSIS NUMERICAL ATTRIBUTE

# CENTRAL TENDENCY

► Mean:

$$\bar{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$$

► Median:

$$x^{\mathsf{med}} = x_{(m+1)/2}, \qquad x^{\mathsf{med}} = (x_{m/2} + x_{m/2+1})/2$$

► Mode

► Midrange:

$$x^{\mathsf{midr}} = (x_{\mathsf{max}} + x_{\mathsf{min}})/2$$

► Geometric mean:

$$\bar{\mu}_{\mathsf{geom}} = \sqrt[m]{\prod_{i}^{m} x_i}$$

# MEASURE OF DISPERSION - NUMERICAL

► Range:

$$x_{\max} - x_{\min}$$

► Mean absolute deviation :

$$MAD = \frac{1}{m} \sum_{i=1}^{m} |x_i - \bar{\mu}|$$

► Sample variance:

$$\bar{\sigma}^2 = \frac{1}{m-1} \sum_{i=1}^{m} (x_i - \bar{\mu})^2$$

► Sample standard deviation :

$$\bar{\sigma} = \sqrt{\bar{\sigma}^2}$$

► Coefficient of Variation:

$$CV = 100 \frac{\bar{\sigma}}{\mu}$$
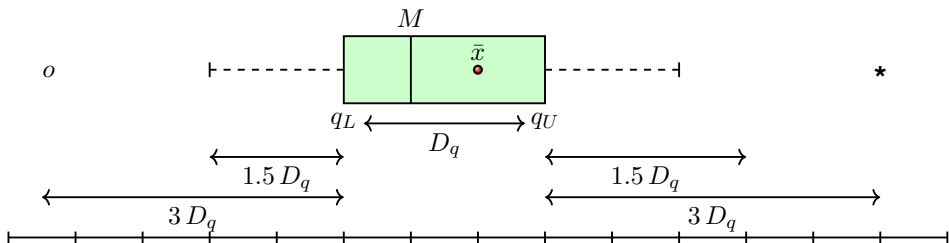
# MEASURE OF DISPERSION - CATEGORICAL

▶ **Gini index:**

$$Gini = 1 - \sum_{h=1}^{H} f_h^2 \qquad \in [0, (H-1)/H]$$

▶ **Entropy index:**

$$Entropy = - \sum_{h=1}^{H} f_h \log_2 f_h \qquad \in [0, \log_2 H]$$

# BOX-PLOT



- Interquartile range $D_q = q_U - q_L = q_{0.75} - q_{0.25}$
- internal lower edge= $q_L - 1.5\, D_q$
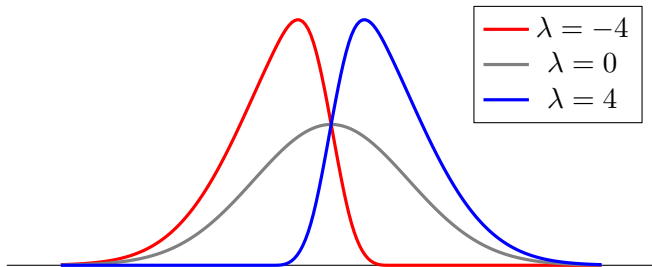- external lower edge= $q_L - 3\, D_q$

# MEASURE RELATIVE LOCATION

- **Mead-Mean**: Mean of values between $q_L$ and $q_U$

- **Trimmed-Mean**: Mean of values between $q_p$ and $q_{(1-p)}$

- **Winsorized-Mean**: Map values smaller (bigger) than $q_p(q_{(1-p)})$ to $q_p(q_{(1-p)})$ and then compute the mean

# ASYMMETRY

$$\bar{\mu}_3 = \frac{1}{m}\sum_{i=1}^{m}(x_i - \bar{\mu})^3, \qquad \text{Skewness} = I_{as} = \frac{\bar{\mu}_3}{\bar{\sigma}^3}$$
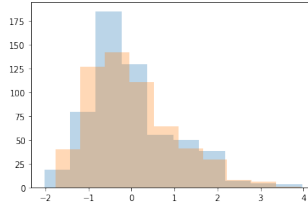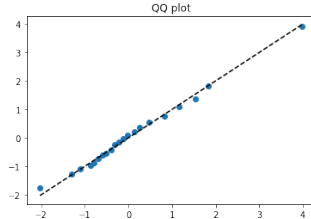
- $I_{as} > 0$ right asymmetry
- $I_{as} < 0$ left asymmetry
- $I_{as} = 0$ symmetric



Legend:
- $\lambda = -4$
- $\lambda = 0$
- $\lambda = 4$

# EMPIRICAL DENSITY

$$\bar{\mu}_4 = \frac{1}{m}\sum_{i=1}^{m}(x_i - \bar{\mu})^4, \qquad \text{Kurtosis} = I_{kurt} = \frac{\bar{\mu}_4}{\bar{\sigma}^4} - 3$$

- ▶ $I_{kurt} > 0$ Hypernormal
- ▶ $I_{kurt} < 0$ Hyponormal
- ▶ $I_{kurt} = 0$ Normal



Leptokurtic Distribution

Normal Distribution

Platykurtic Distribution

# COMPARING DISTRIBUTIONS - QQ-PLOTS

# MEASURE OF CORRELATION

► **Sample covariance:**

$$cov(a_j, a_k) = \frac{1}{m-2} \sum_{i=1}^{m} (x_{ij} - \bar{\mu}_j)(x_{ik} - \bar{\mu}_k)$$

► **Sample LINEAR correlation:**

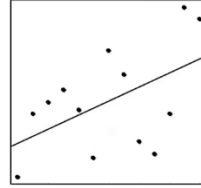$$r_{jk} = \frac{cov(a_j, a_k)}{\bar{\sigma}_j \, \bar{\sigma}_j} \qquad \in [-1, 1]$$
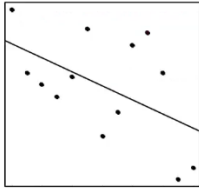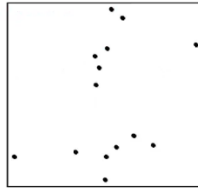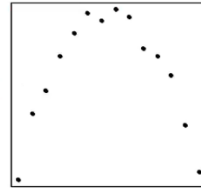
# LINEAR CORRELATION



(a) r=0.96    (b) r=-0.96    (c) r=0.6

(d) r=-0.6    (e) r=0.0    (f) r=0.0

# CORRELATION ON CATEGORICAL ATTRIBUTES

| area | family | | totale |
|---|---|---|---|
| | 0 | 1 | |
| 1 | 2 | 4 | 6 $(f_1)$ |
| 2 | 4 | 2 | 6 |
| 3 | 2 | 5 | 7 |
| 4 | 3 | 3 | 6 |
| totale | 11 $(g_1)$ | 14 $(g_2)$ | 25 |

**Two attributes are independent if**

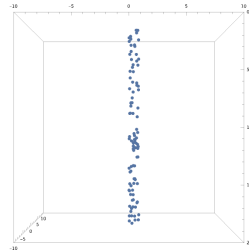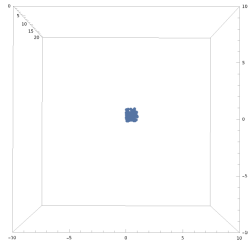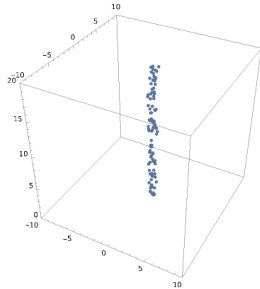$$\frac{t_{r1}}{g_1} = \frac{t_{r2}}{g_2} \qquad r = 1, 2, \ldots, J$$

**Contingency tables**

# PCA: PRINCIPAL COMPONENT ANALYSIS



original data space

PCA

component space

POLIMI GRADUATE SCHOOL OF MANAGEMENT

# PCA-INTUITION



Compute projections that better capture the **variance**

# PCA: PRINCIPAL COMPONENT ANALYSIS

- ► Covariance data matrix $V = X'X$
- ► $\bar{x}_{ij} = x_{ij} - \bar{\mu}_j$
- ► New components $p_j$ obtained as a linear transformation of original data $p_j = X w_j$
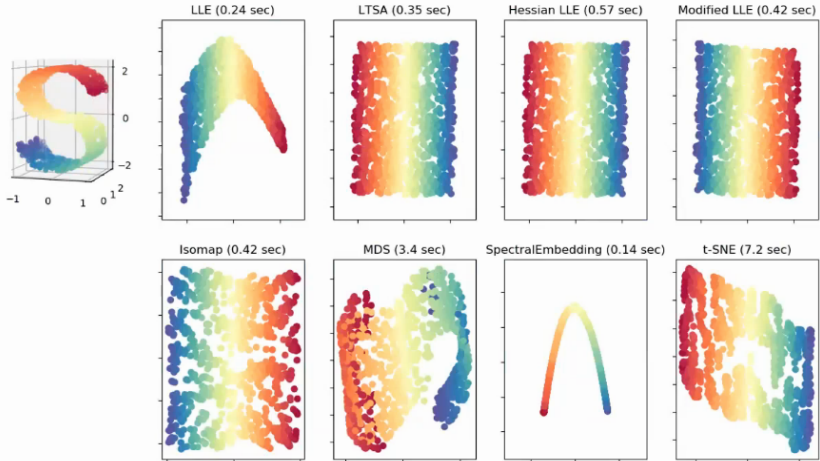- ► Variance of $p_j = w'_j X'X w_j = w'_j V w_j$
- ► Maximizing the variance:

$$\max_{w_1} w'_1 V w_1$$
$$\text{s.t. } w'_1 w_1 = 1$$

- ► $w_j$ is the $j$-th eigenvector of $V$, which explains a variance $\lambda_j$ which is the $j$-th eigenvector
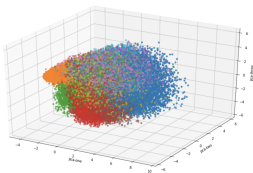
# NONLINEAR REDUCTION



Manifold Learning with 1000 points, 10 neighbors

LLE (0.24 sec) · LTSA (0.35 sec) · Hessian LLE (0.57 sec) · Modified LLE (0.42 sec)

Isomap (0.42 sec) · MDS (3.4 sec) · SpectralEmbedding (0.14 sec) · t-SNE (7.2 sec)

# T-SNE: T-DISTRIBUTED STOCHASTIC NEIGHBOUR EMBEDDING

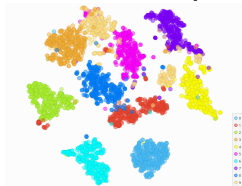We aim to project observations preserving observation distance.

**High-Dimensional Space**



Distance as Normal Distribution

$$p_{ij} = \frac{\exp\left(-\|x_i - x_j\|^2 / 2\sigma_i^2\right)}{\sum_{k \neq l} \exp\left(-\|x_k - x_l\|^2 / 2\sigma_i^2\right)}$$

**Low-Dimensional Space**



Distance as t-Students Distribution

$$q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + \|y_k - y_l\|^2\right)^{-1}}$$

We try to minimize the divergence between the distributions

$$KL(P\|Q) = \sum_i \sum_{j \neq i} p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right) \qquad \text{(Kullback-Leibler)}$$

# EJEMPLO B - MNIST DATASET



https://adamharley.com/nn_vis

POLIMI GRADUATE SCHOOL OF MANAGEMENT

# EJEMPLO C - MNIST FASHION (ZALANDO)

# THANK YOU