

Quality of Life

Clustering of Italian Provinces



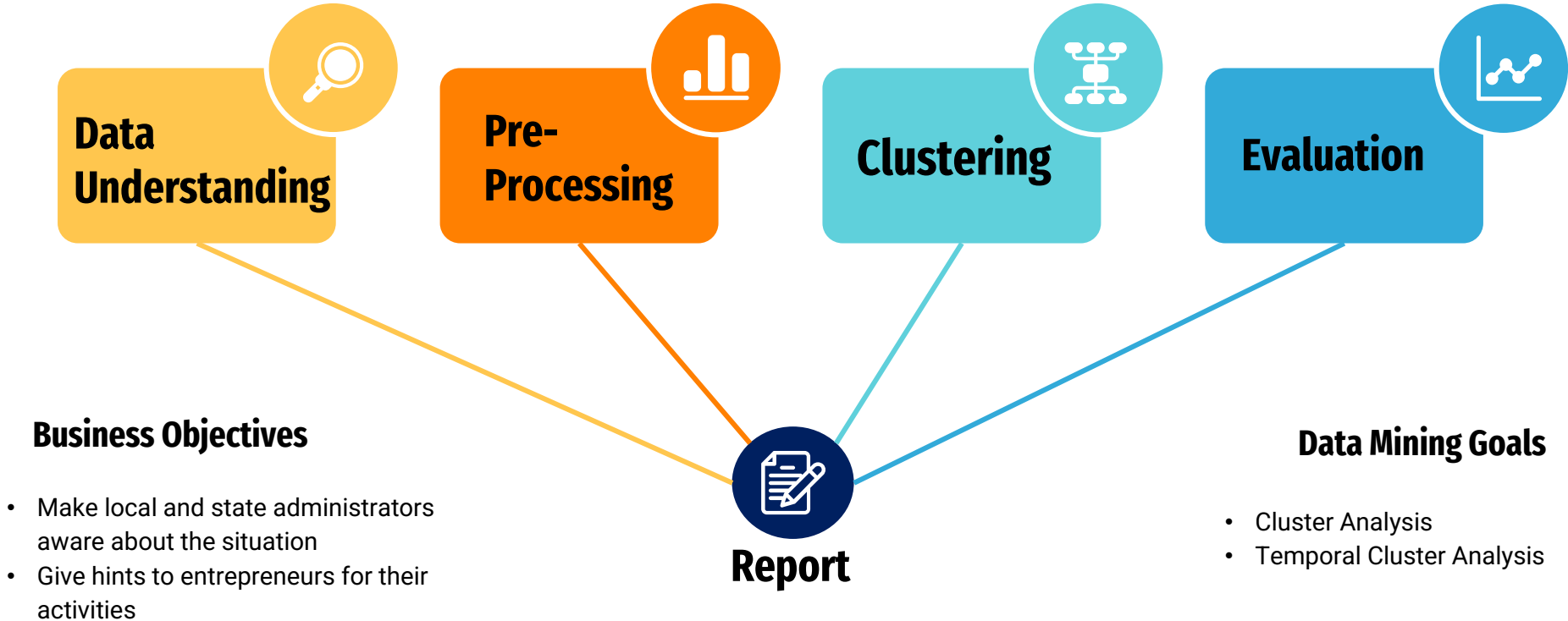
UNIVERSITÀ DI PISA

Data Mining and Machine Learning



By Iacopo Bicchierini

Highlights



Dataset

**Wealth and
Consumption**



**Culture and
Leisure**

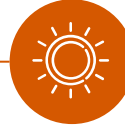


15 indicators each group

**Business
and Work**



**Environment
and Services**



**Justice and
Security**



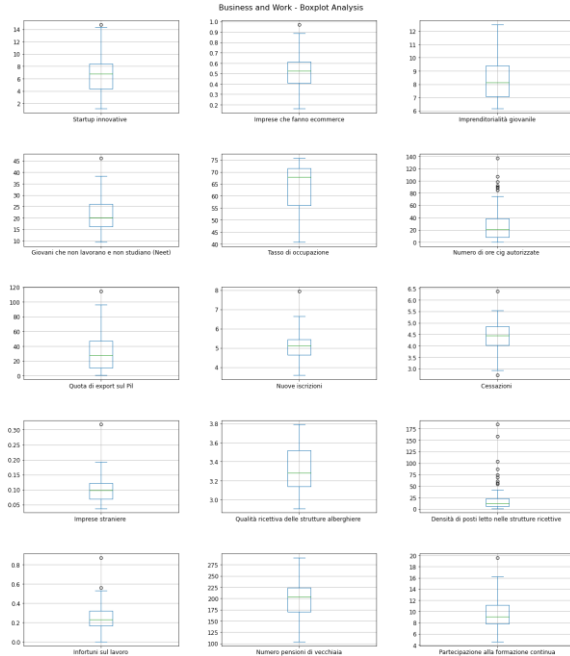
**Demography
and Society**



Source: //Sole240re year 2020, 2021 and 2022

Exploratory Data Analysis

Box Plot Analysis



Correlation Analysis

	Cessioni	Densità di posti letto nelle strutture ricettive	Giovani che non lavorano e non studiano (Neet)	Imprenditorialità giovanile	Imprese che fanno ecommerce	Imprese straniere	Infortuni sul lavoro	Numero di ore cig autorizzate	Numero pensioni di vecchiaia
Cessioni									
Densità di posti letto nelle strutture ricettive									
Giovani che non lavorano e non studiano (Neet)									-0.73334
Imprenditorialità giovanile									
Imprese che fanno ecommerce									
Imprese straniere									
Infortuni sul lavoro									
Numero di ore cig autorizzate									
Numero pensioni di vecchiaia									-0.73334
Nuove iscrizioni									
Partecipazione alla formazione continua									
Qualità ricettiva delle strutture alberghiere									
Quota di export sul PIL									
Startup innovative									
Tasso di occupazione	0.701139		-0.893498	-0.77364					0.832873

Preprocessing



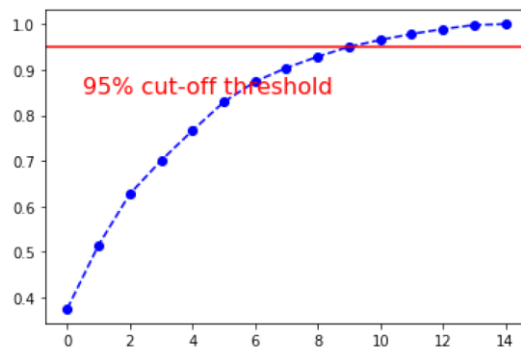
Normalization

$$z = \frac{x - \mu}{\sigma}$$

Z-Score



Feature Reduction



PCA with 95% variance explained

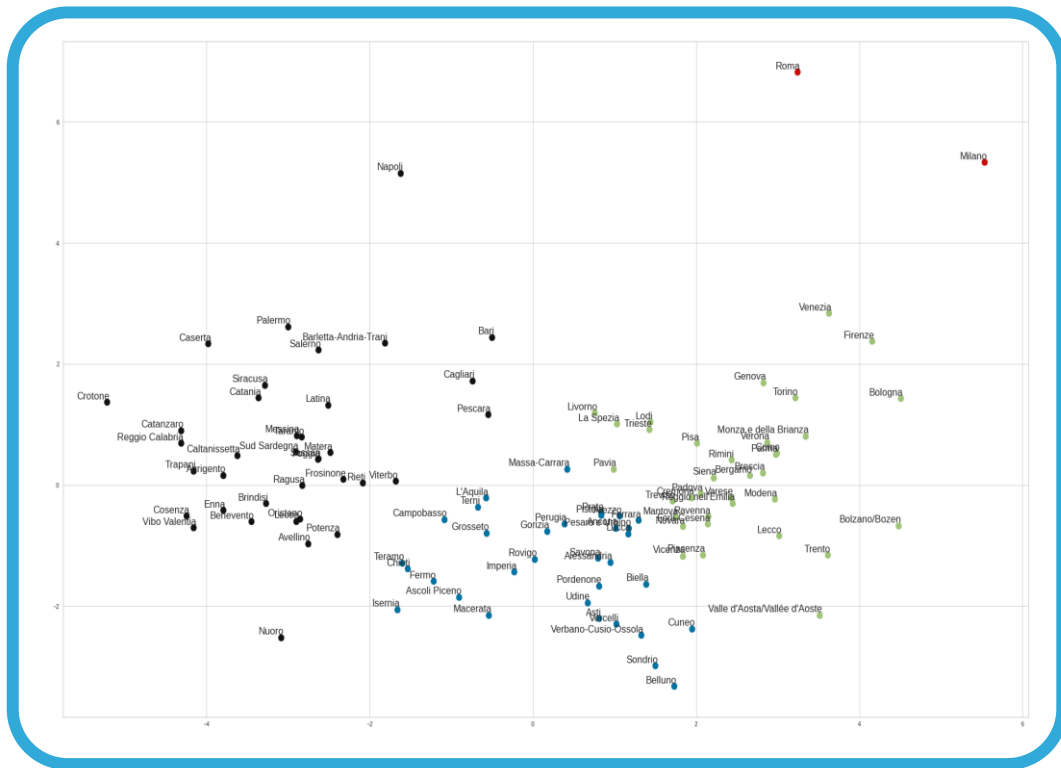


Cluster Tendency

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

Hopkins Statistic
of around 0.7 for all groups

Partitioning Clustering



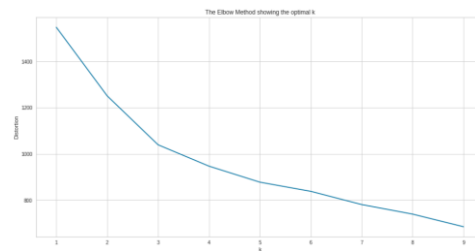
Wealth and Consumption 2022

Methods

- K-Means
- CLARANS

Choosing Number of Clusters

- Elbow Method
- Silhouette Analysis



Hierarchical Clustering

-

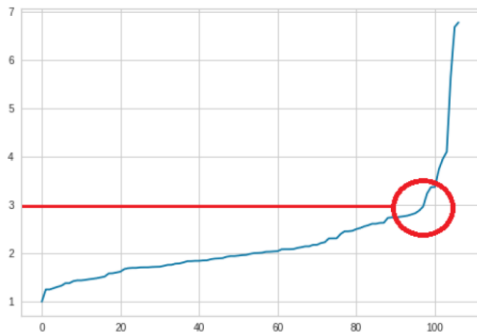
Density-based Clustering

Method

DBSCAN

Determine Parameters Eps and MinPts

Using k-dist graph



Total Fail

No global density values
for this dataset

Temporal Cluster Analysis

Matching Matrix

How many object from cluster i of window $t-1$ goes to cluster j of window t by the total number of objects

Clustering Method

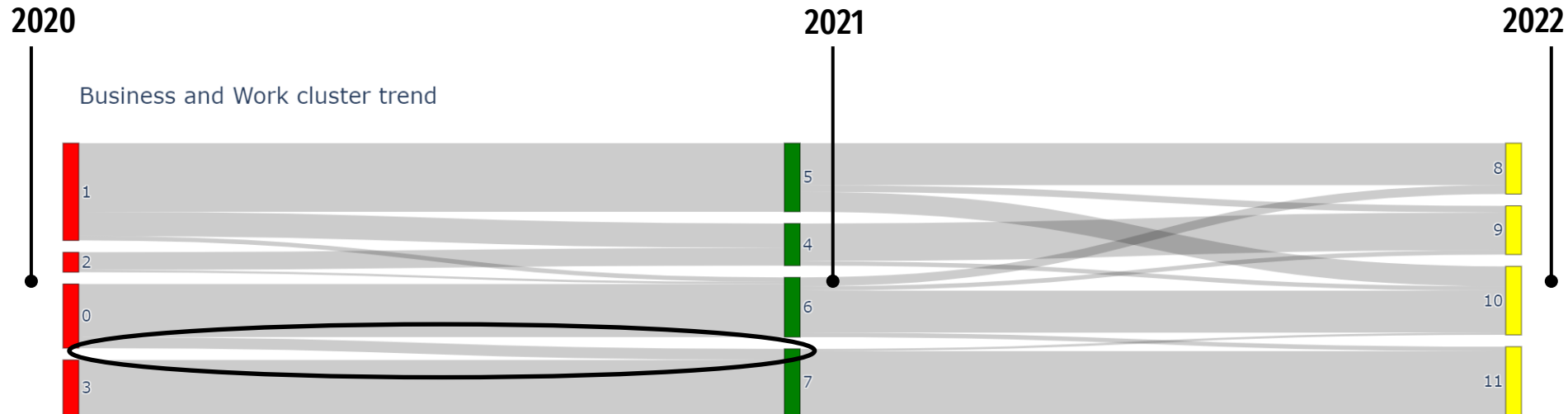
Same Clustering method used for the 3 years

Visualization

Using **Sankey Diagram**

Result

Intersection between clusters reveals who goes where



Evaluation

Silhouette Coefficient

Poor values

Best Clustering

K-Means

Best Indicator

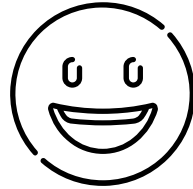
Demography and
Society

Worst Clustering

CLARANS

	K-Means	CLARANS	Agglomerative	BIRCH
<i>Business and Work</i>	0.1567	0,096	0.1582	0.1614
<i>Culture and Leisure</i>	0.1913	0.1753	0.1816	0.1519
<i>Demography and Society</i>	0.199	0.1748	0.184	0.1796
<i>Environment and Services</i>	0.1557	0.085	0.1182	0.1097
<i>Justice and Security</i>	0.17	0.1	0.133	0.136
<i>Wealth and Consumption</i>	0.1965	0.1815	0.1785	0.123
Average	0.1782	0.1354	0.159	0.1436

THANK FOR YOUR ATTENTION



GitHub: <https://github.com/Bicchie/Quality-of-Life-Italian-Provinces-Clustering>