# Catalonian Road Accidents 2010-21

MARINA ALAPONT VIDAL
DANIEL PULIDO  GÁLVEZ
SIMÓN HELMUTH OLIVA STARK
JOEL CARDONA SAUS
DAVID LATORRE ROMERO

Presentation: 21/10/2022

# Outline

1. Goals and Database overview

2. Data Mining process

3. Descriptive analysis

4. Preprocessing

5. PCA and Clustering

6. Profiling

7. Conclusions

8. Task Scheduling

# Topics and Goals

## Topics

- Factors that can affect accidents

  - Velocity

  - Weather

  - State of the road

- Consequences of the accident

  - Number of injuries

  - Entities involved

## Goals

- What kinds of situations facilitate accidents

- Help the society studying accidents

- Accidents affects everyone

# Database Overview

### ORIGINAL DIMENSIONS

- 21.161 individuals

- 58 features

### AFTER DATA SELECTION

- 5.000 individuals

- 23 features:

  - 7 numerical
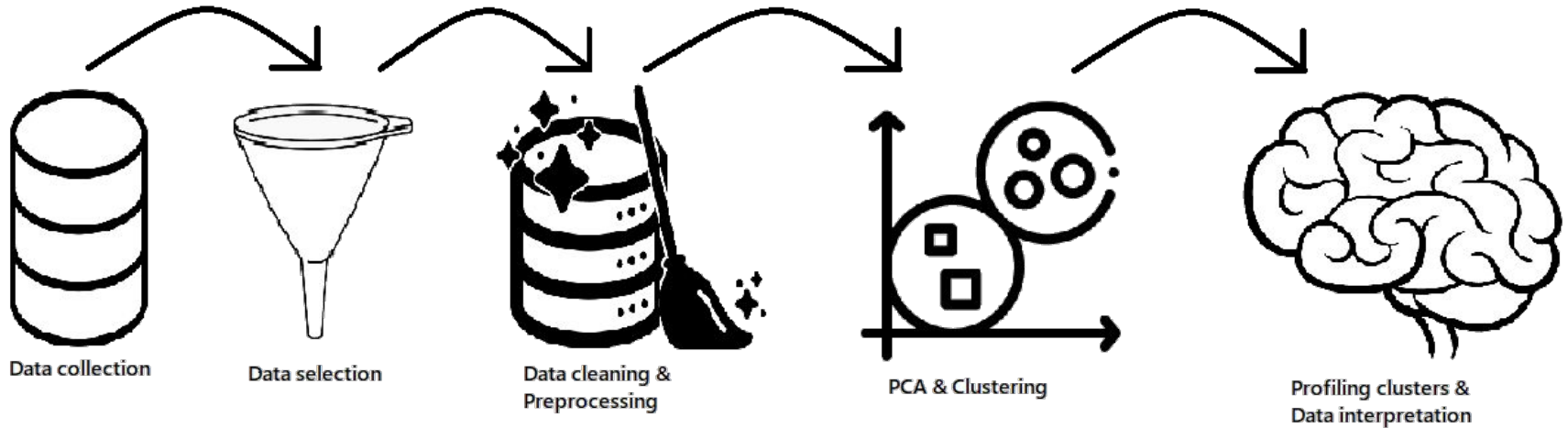
  - 11 categorical

  - 5 boolean

Data Source:

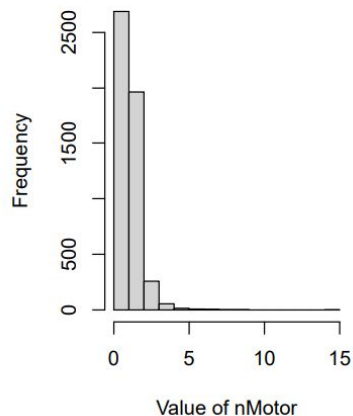https://analisi.transparenciacatalunya.cat/Transport/Accidents-de-tr-nsit-amb-morts-o-ferits-greus-a-Ca/rmgc-ncpb
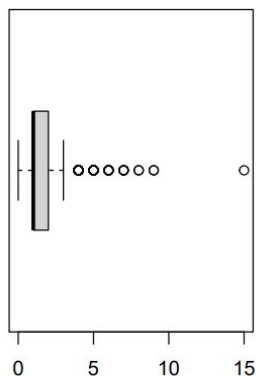
# Data Mining Process



Data collection

Data selection

Data cleaning & Preprocessing

PCA & Clustering

Profiling clusters & Data interpretation

# Descriptive Analysis

**Histogram of nMotor**



**Boxplot of nMotor**



| Min. | 1st Qu. | Median | Mean | 3rd Qu. |
|------|---------|--------|------|---------|
| 0 | 1 | 1 | 1.5214 | 2 |

| Max. | sd | vc. | Missing |
|------|-----|-----|---------|
| 15 | 0.8232117 | 0.5410883 | 0 |

**Pie of Prov**



**Barplot of Prov**



| Prov | Frequency | Proportion |
|------|-----------|------------|
| Barcelona | 2993 | 0.5986 |
| Girona | 783 | 0.1566 |
| Tarragona | 698 | 0.1396 |
| Lleida | 526 | 0.1052 |

# Preprocessing Steps (1)

➔ **Factorization**, levels, sorting

| |
|---|
| ~~Character~~ |
| Factor -> 17 |
| Date -> 1 |
| Integer -> 7 |



➔ **Renaming**: Variables, levels



~~D_INFLUIT_CIRCULACIO~~ -> TrafficInf

~~D_INFLUIT_ESTAT_CLIMA~~ -> WeatherInf
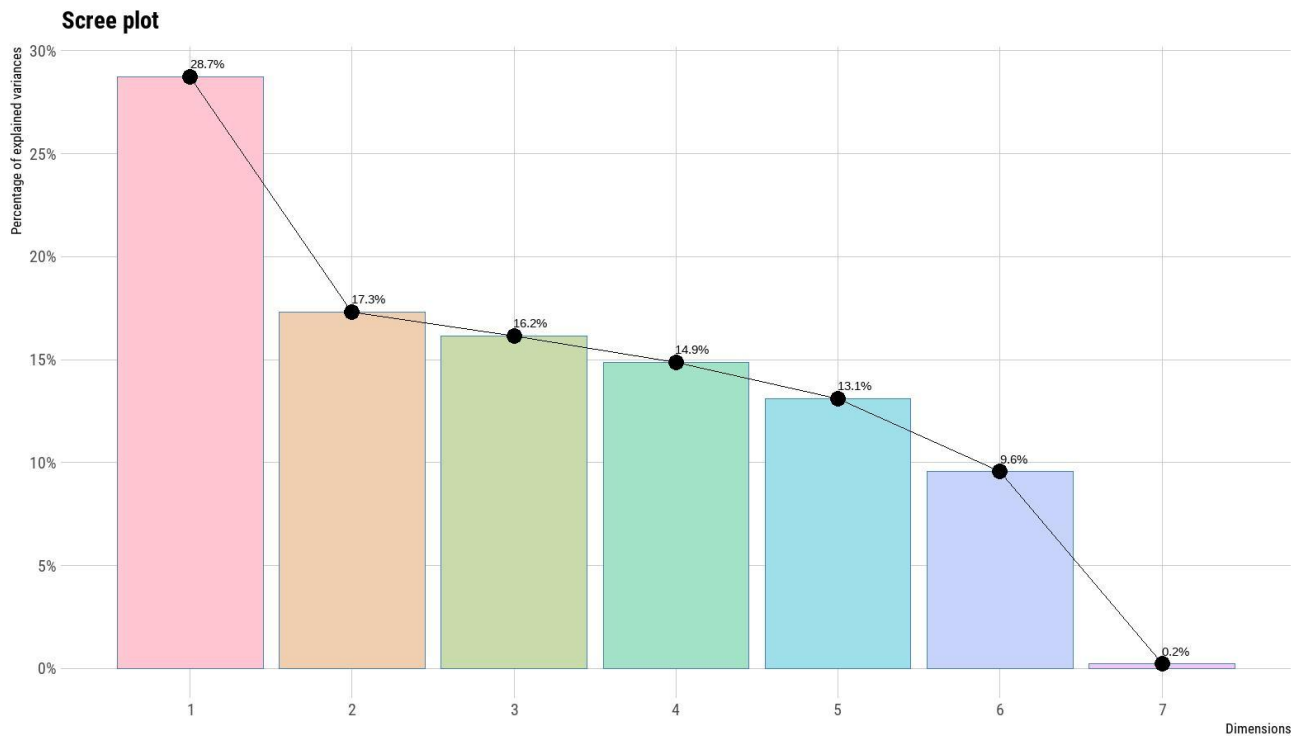
etc...

~~Catalan~~ -> English

# **Preprocessing Steps (2)**

➔ **MISSINGS**

**Random,** no pattern

| Variable name | Proportion of missing values |
|---|---|
| Vel | **18,56%** |
| Escaped | 0,82% |
| Weather | 0,02% |
| TrafficInf | 0,02% |
| WeatherInf | 0,02% |
| LightInf | 0,02% |
| VisionInf | 6,68% |
| Surface | 0,02% |

# PCA: Specifications



Scree plot
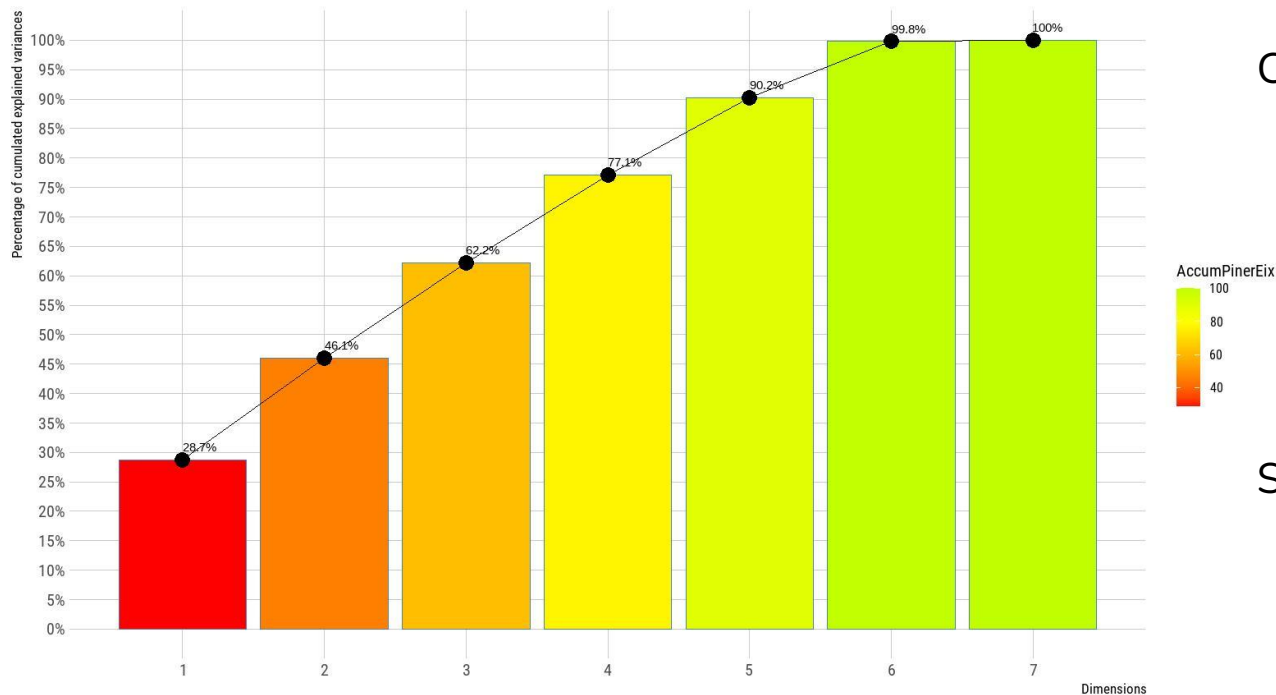
# PCA : Specifications



Cumulative Scree plot

Criteria: keep 80% of the inertia

Selected dimensions: 1 to 4

# PCA: Specifications (2)

QUALITY -> ^2 cosine of **correlation**

[0, 1]



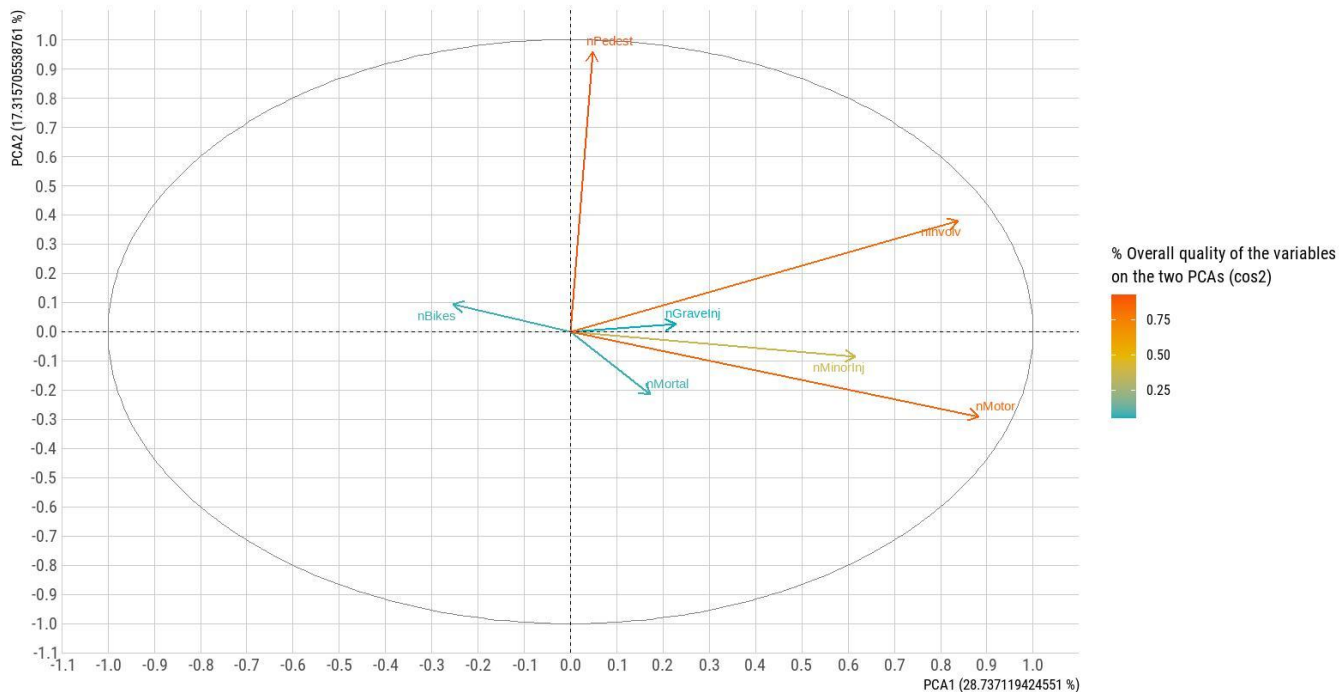Quality of the variables on the PCAs (square cosine of the correlation)

# PCA: First Factorial Plane
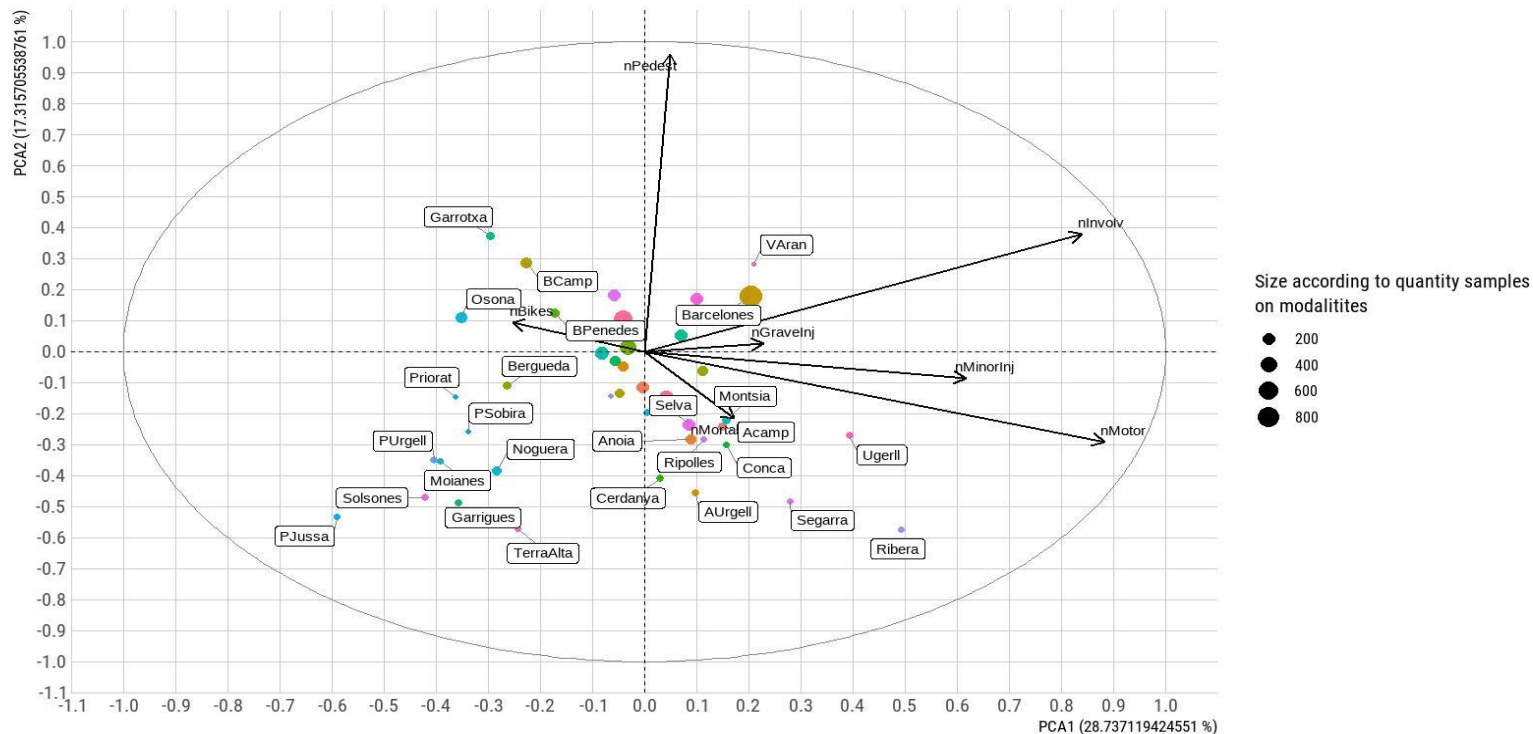


**Correlation circle**

and variables quality (calculated as the sum (for the two PCAs) of the square cosine of the correlation)

# PCA: First Factorial Plane (2)

**Correlation circle, and representation of all modalitites of Region cat. variable**

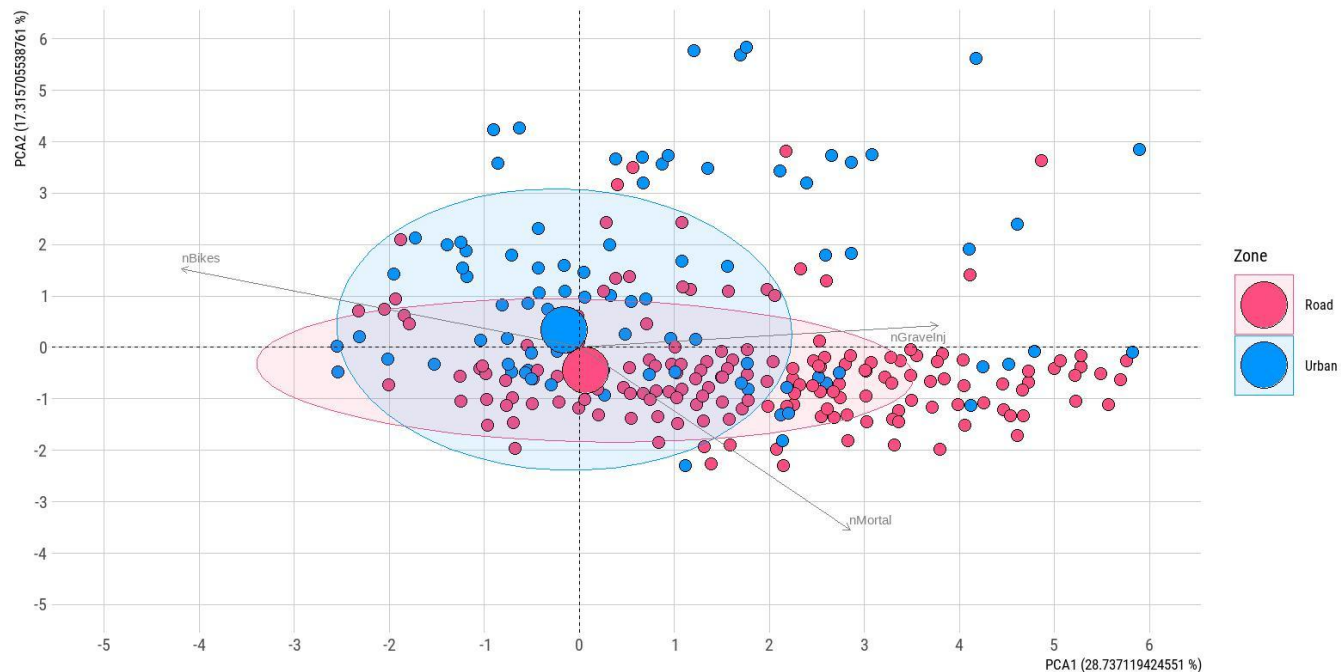and modality quantities as point sizes

# PCA: Conclusions (1)



**Correlation circle, Individuals, and representation of the Zone categorical variable (ZOOMED IN)**

Correlation vectors are scaled for clarity.
Concentration ellipses (using multivariate normal distribution) are drawn. Mean points for the levels are also drawn.

# PCA: Conclusions (2)



Correlation circle, Individuals, and representation of the Surface categorical variable (ZOOMED IN)

Correlation vectors are scaled for clarity.
Concentration ellipses (using multivariate normal distribution) are drawn. Mean points for the levels are also drawn.
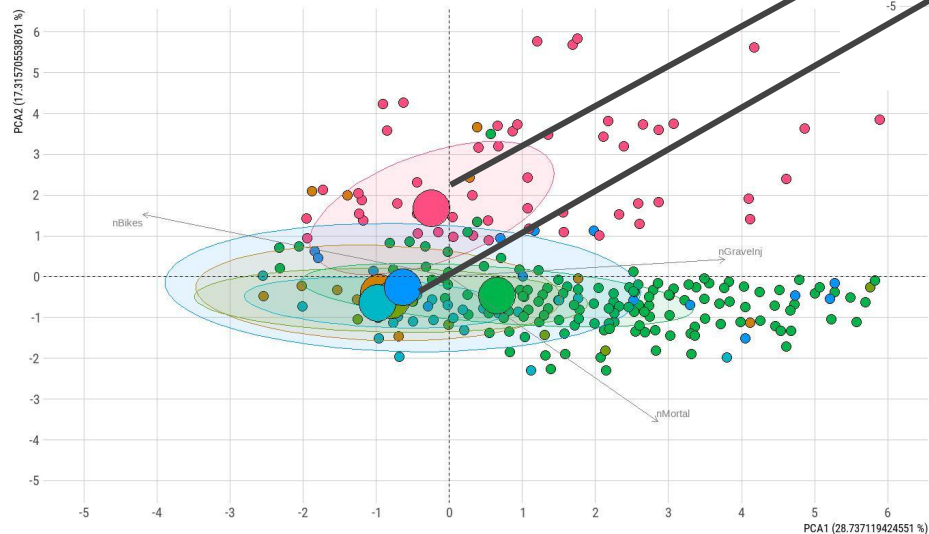
Surface
- Dry&Clean
- Slippery
- Wet
- Flooded
- Icy
- Snowy
- UnkSrfc

Correlation circle, Individuals, and representation of the AccType categorical variable (ZOOMED IN)
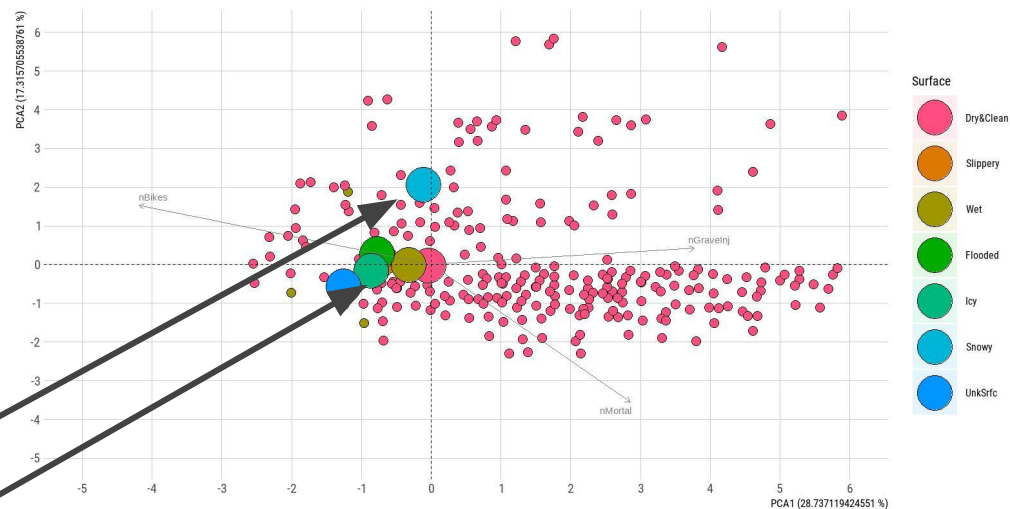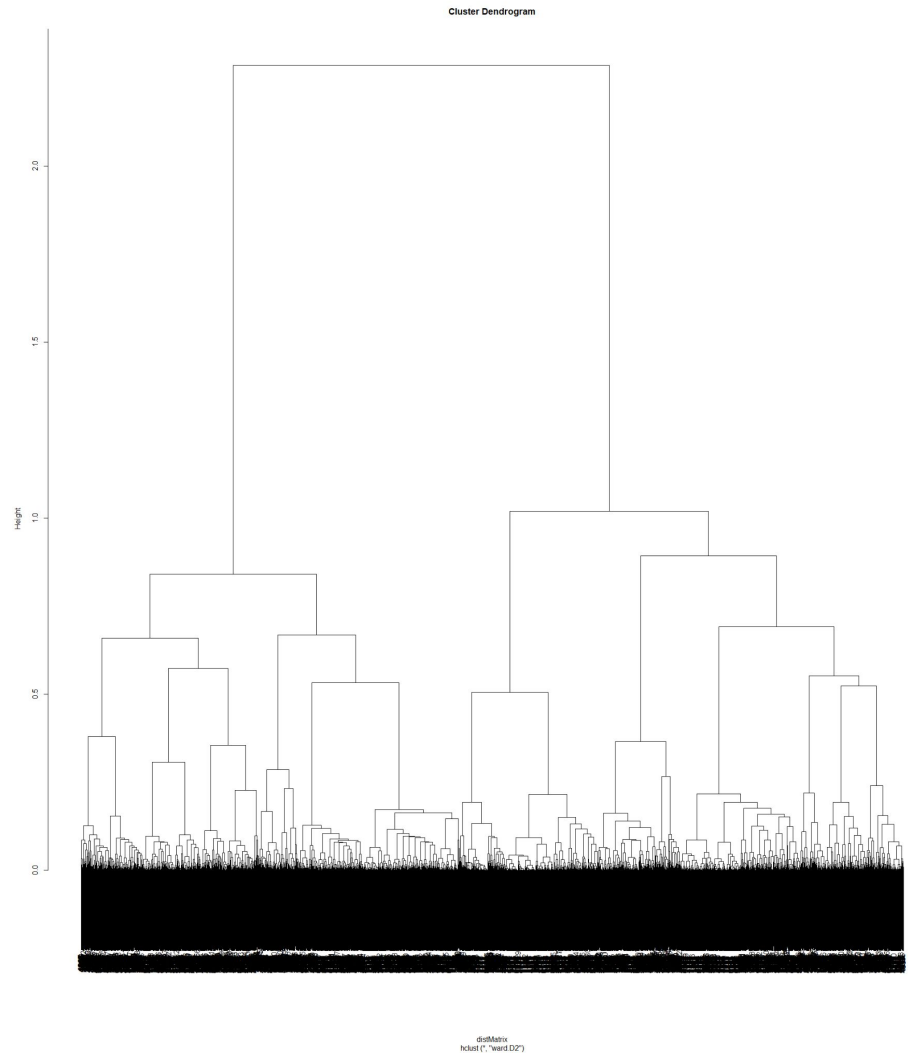
Correlation vectors are scaled for clarity.
Concentration ellipses (using multivariate normal distribution) are drawn. Mean points for the levels are also drawn.

AccType
- RunOver
- Rollover
- HitObstacle
- HitVehicle/s
- NARoadExit
- Other

# **Clustering process**

- Date decomposed into Year and Month

- Ward's D2 method

- Gower mixed distance

- Minimize inter-class inertia loss



Cluster Dendrogram

distMatrix
hclust (*, "ward.D2")

# Cluster Number Discussion

- KPI's

- PCA results

- Expert's opinion

| | k=2 | k=3 | k=4 | k=5 | k=6 | k=7 | k=8 | k=9 | k=10 |
|---|---|---|---|---|---|---|---|---|---|
| Cindex | 0.3788 | 0.3510 | 0.3312 | 0.3232 | 0.3465 | 0.3329 | 0.3285 | 0.3229 | 0.3162 |
| Mcclain | 0.7592 | 1.5408 | 1.9721 | 2.7166 | 3.0248 | 3.4005 | 4.3610 | 5.0876 | 5.3953 |
| Silhouette | 0.1582 | 0.1388 | 0.1684 | 0.1309 | 0.1574 | 0.1662 | 0.1391 | 0.1254 | 0.1407 |
| Dunn | 0.0815 | 0.0815 | 0.0815 | 0.0815 | 0.0915 | 0.0791 | 0.0458 | 0.0458 | 0.0460 |

| | Cluster 1 | Cluster 2 |
|---|---|---|
| Number of individuals | 2696 | 2304 |

**Number of clusters = 2**

# Class Interpretation Tools
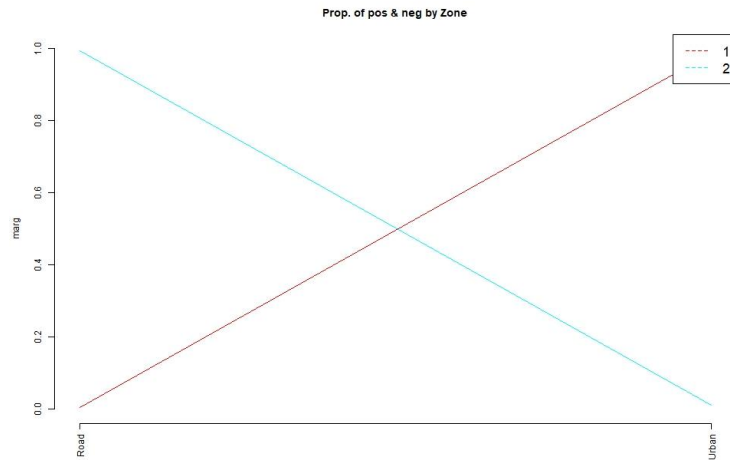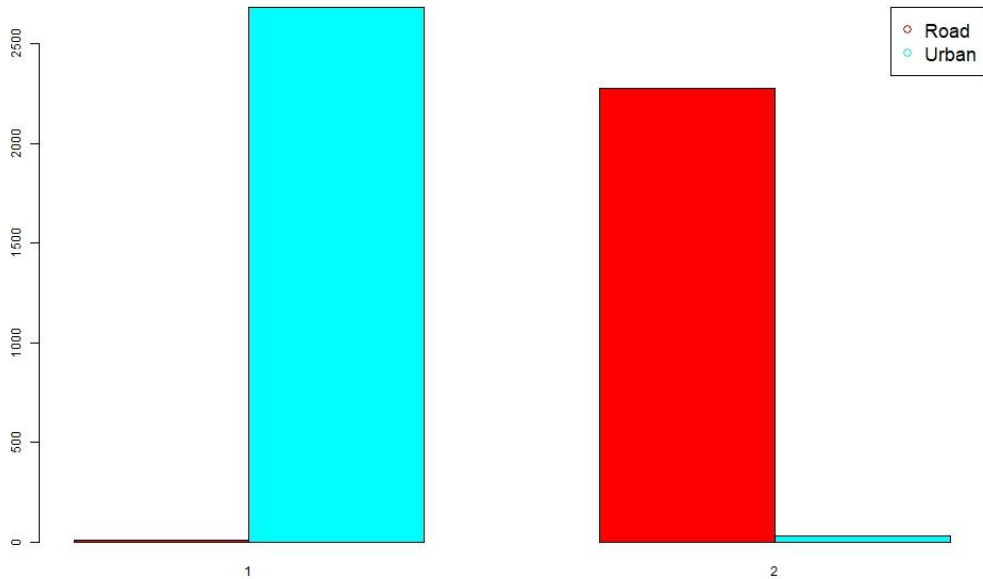


Profiling plots + statistics

# Profiling graphs or numerical information about our clusters to be highlighted
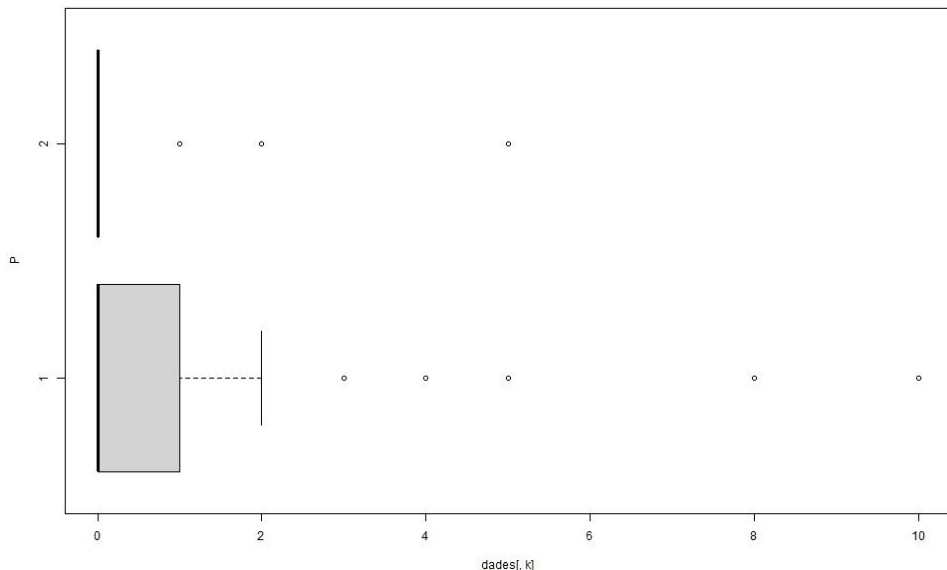


Prop. of pos & neg by Zone

[1] "Distribucions condicionades a columnes:"

```
P           Road          Urban
  1 0.005249344 0.988946205
  2 0.994750656 0.011053795
```
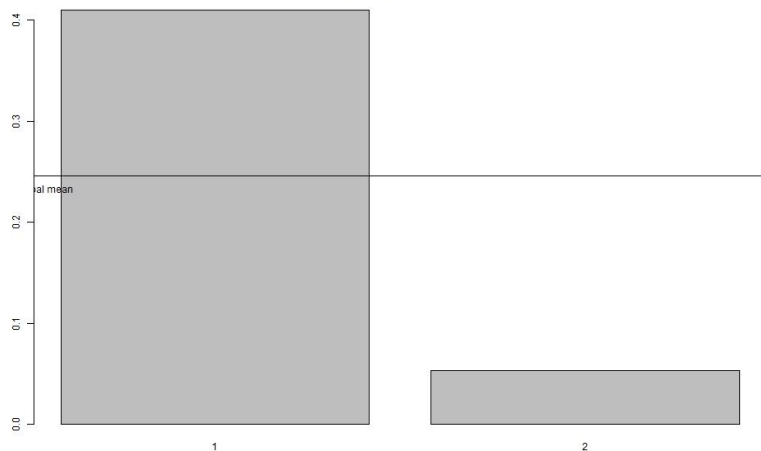
# Profiling graphs or numerical information about our clusters to be highlighted



Boxplot of nPedest vs Class



Means of nPedest by Class

```
[1] "Anàlisi per classes de la variable: nPedest"
[1] "Estadístics per groups:"
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 0.0000  0.0000  0.0000  0.4102  1.0000 10.0000
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.00000 0.05295 0.00000 5.00000
```

# Final Class Profiling

**CLUSTER 1:**

Urban zones

Run overs

↓ mortality

↓ minor injured victims

↑ number of pedestrians

**CLUSTER 2:**

Road zones

Unknown exits from the road

↑ mortality

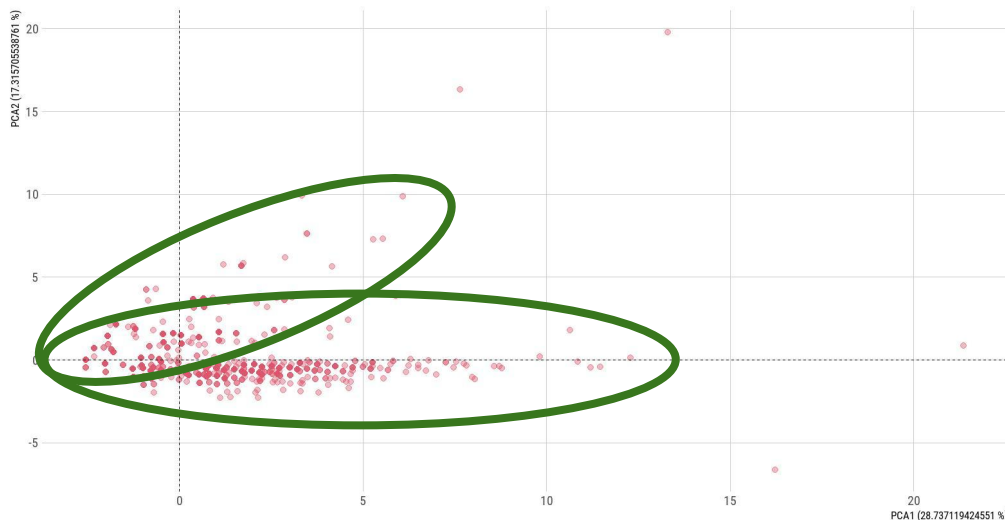↑ minor injured victims

↓ number of pedestrians

# PCA and Clustering Conclusions

➔ **SIMILAR**

**Individuals plot**
Point transparency according to the frequency on that exact point



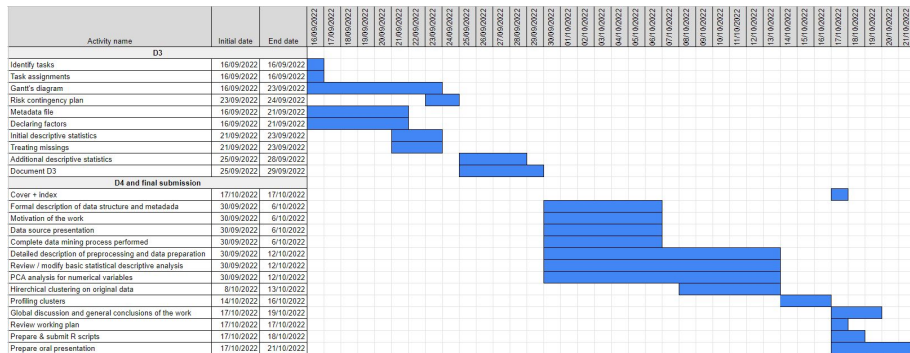| CLUSTER 1 | CLUSTER 2 |
|---|---|
| not so s. accidents | serious accidents |
| low speeds | high speeds |
| urban roads | rural roads |
| populated regions | rural regions |
| not so many un. imp. | many un. implicated |

similar conclusions on PCA!

# General Conclusions

Urban, populated, low speeds $\longrightarrow$ focus **protect pedestrians**

rural roads, high speed $\longrightarrow$ most dangerous, most units implicated.
**More sector radars**

car, bicycle $\longrightarrow$ **Increase distance**

icy surface $\longrightarrow$ Very mortal. Extreme safety measures

snowy surface $\longrightarrow$ Related with Runovers

Wet, slippery, flooded $\longrightarrow$ Related with Rollowers, hitting obstacles, road exits.
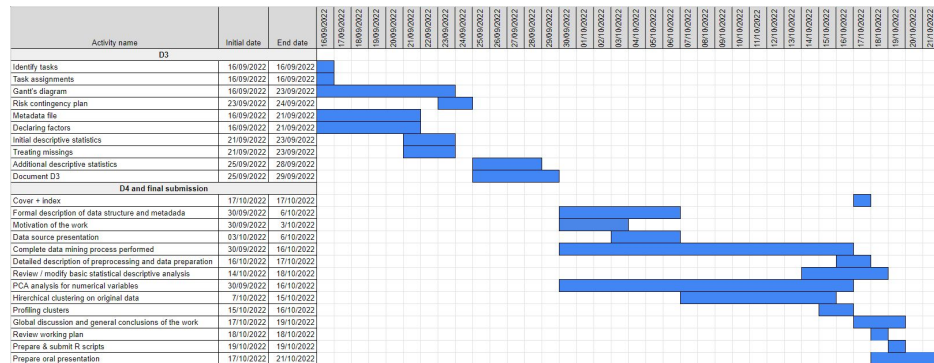
# Task Scheduling

## Original Gantt Chart



## Final Gantt Chart

# THANKS

Any questions?