# US Accidents (2016-2021)

**MVA Assignment D4**

Ander Barrio Campos
Biel Caballero Vergés
Miona Dimic
Gerard Gómez Jorba
Odysseas Kyparissis
Cecilia Pérez Pérez

# Table of contents

# 1. Introduction

- Motivation
  - Tackle the problematic surrounding road accidents using data.

- General description
  - Large dataset
  - Qualitative, quantitative and binary variables are included.
  - Dataset includes information regarding:
    - Where the accident happened?
    - When did it happen?
    - How it affected the road?
    - What severity did it imply?
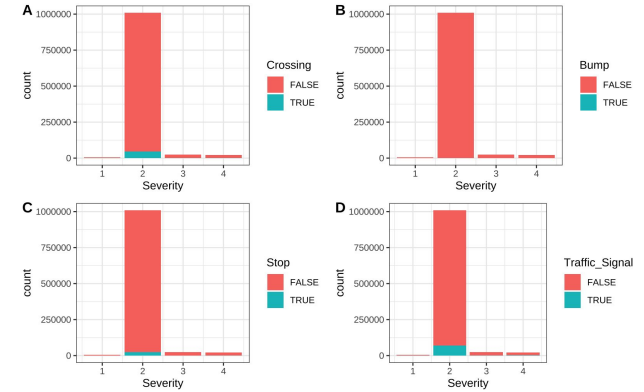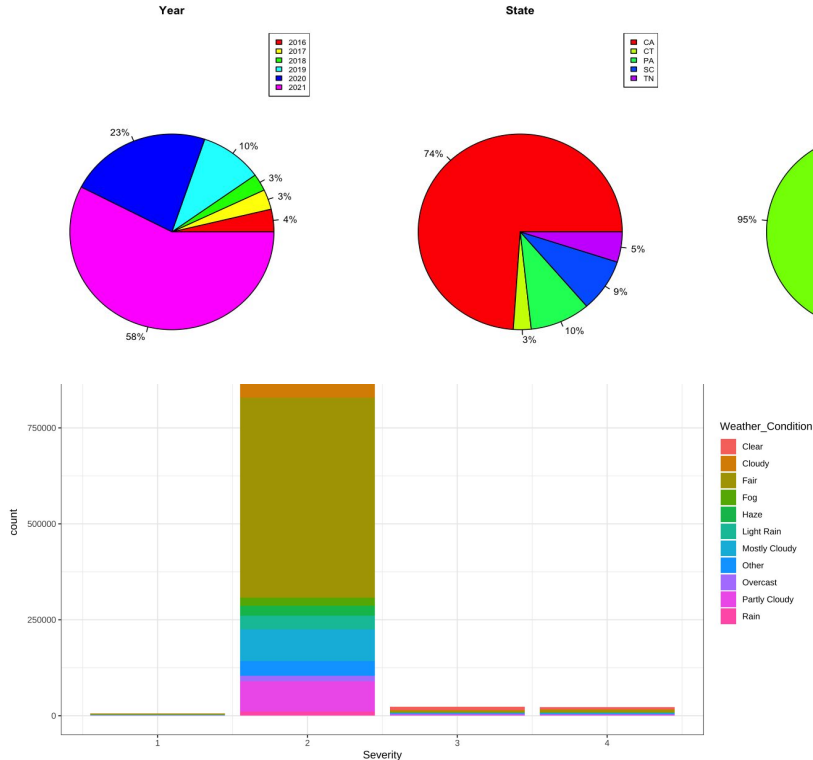    - What where the conditions surrounding it?

# 2. Univariate and Bivariate Descriptive Statistics



Year

| | |
|---|---|
| ■ | 2016 |
| ■ | 2017 |
| ■ | 2018 |
| ■ | 2019 |
| ■ | 2020 |
| ■ | 2021 |

State

| | |
|---|---|
| ■ | CA |
| ■ | CT |
| ■ | PA |
| ■ | SC |
| ■ | TN |

Severity

| | |
|---|---|
| ■ | 1 |
| ■ | 2 |
| ■ | 3 |
| ■ | 4 |

➤ Absence of traffic control indicators

➤ More frequent accident occurrence over years

➤ Majority of accidents with Severity level 2

➤ Minor influence of weather conditions

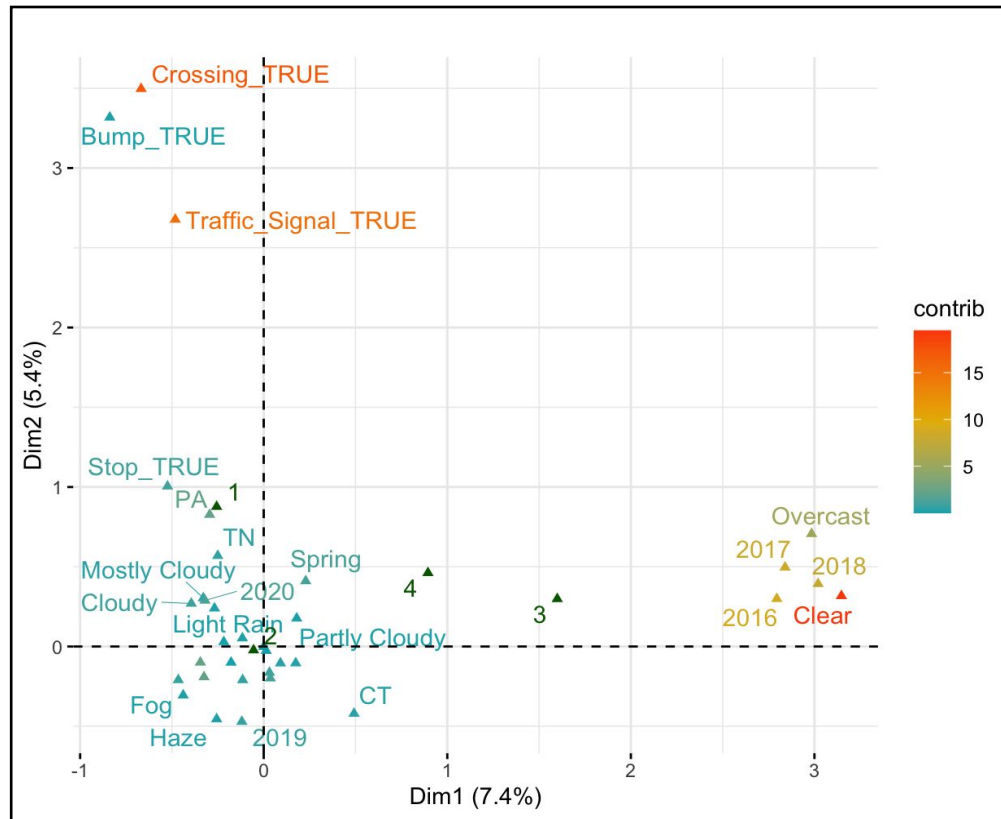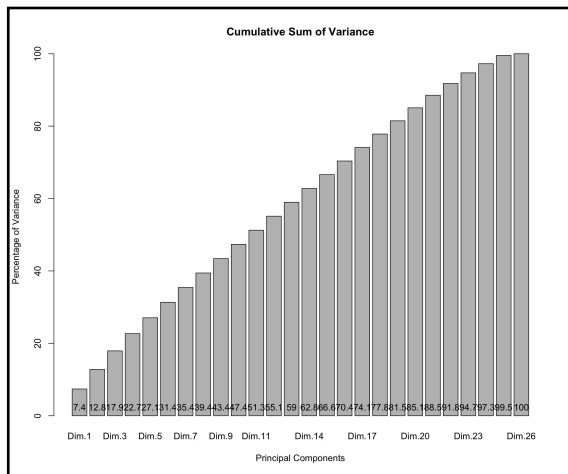➤ 74% of recorded US accidents are in the State of California

# 3. PCA

- Number of selected PCs
- General conclusion ➜ not very valid
- Weather condition associated variables



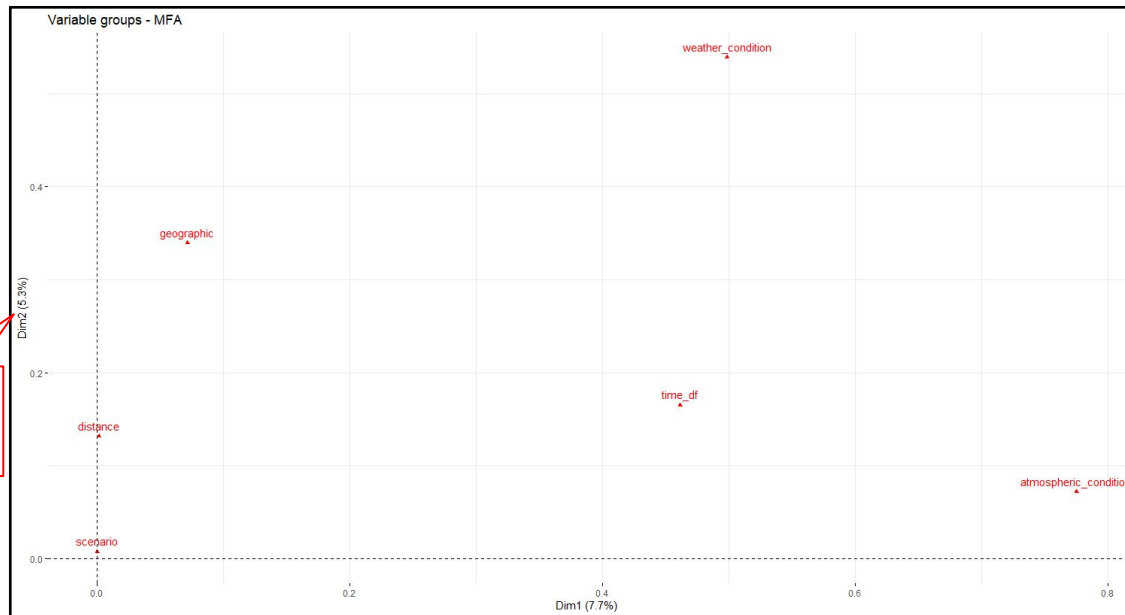Cumulative Sum of Variance



PCA – Biplot

# 4. MCA

- Number of selected dimensions is 12
- General conclusion ➜ two clear groups are formed
- Weather_Condition (without precipitation) and Year explain dimension 1 and Road Elements for dimension 2

# 5. MFA



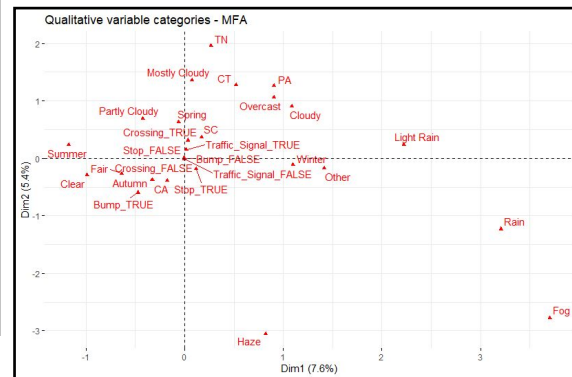Variable groups - MFA

Quantitative variables - MFA
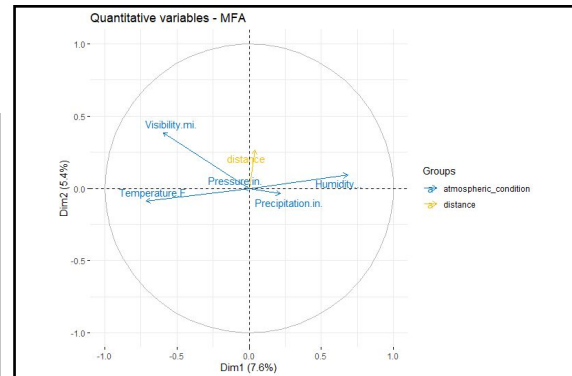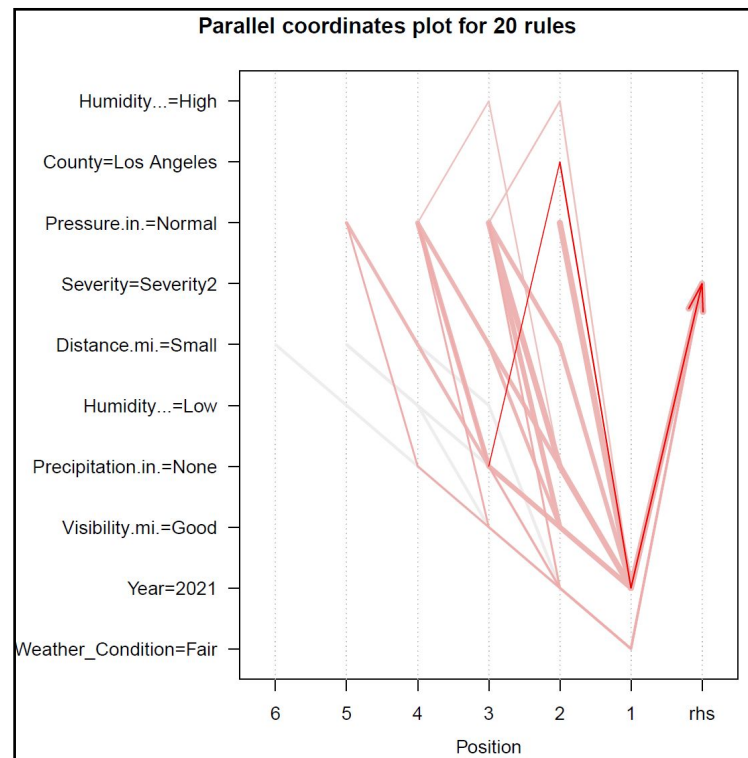
Qualitative variable categories - MFA

Weather condition + location

Weather condition + time

# 6. Association Rules Mining Analysis


Relative Item Frequency Plot


Parallel coordinates plot for 20 rules

# 7. Hierchical Clustering

- Preprocessing
  - Selecting Accidents of California 2020
  - ~20K observations
- Clustering method
  - Agglomerative HC
  - Gower distance
- Metrics
  - Daisy Squared
- Aggregation criteria
  - Ward method

**Cluster Dendrogram**

Height

250

150

50

0

daisy
hclust (*, "ward.D")

# 7. Hierchical Clustering

- Selection of optimal number of clusters
  - NbClust
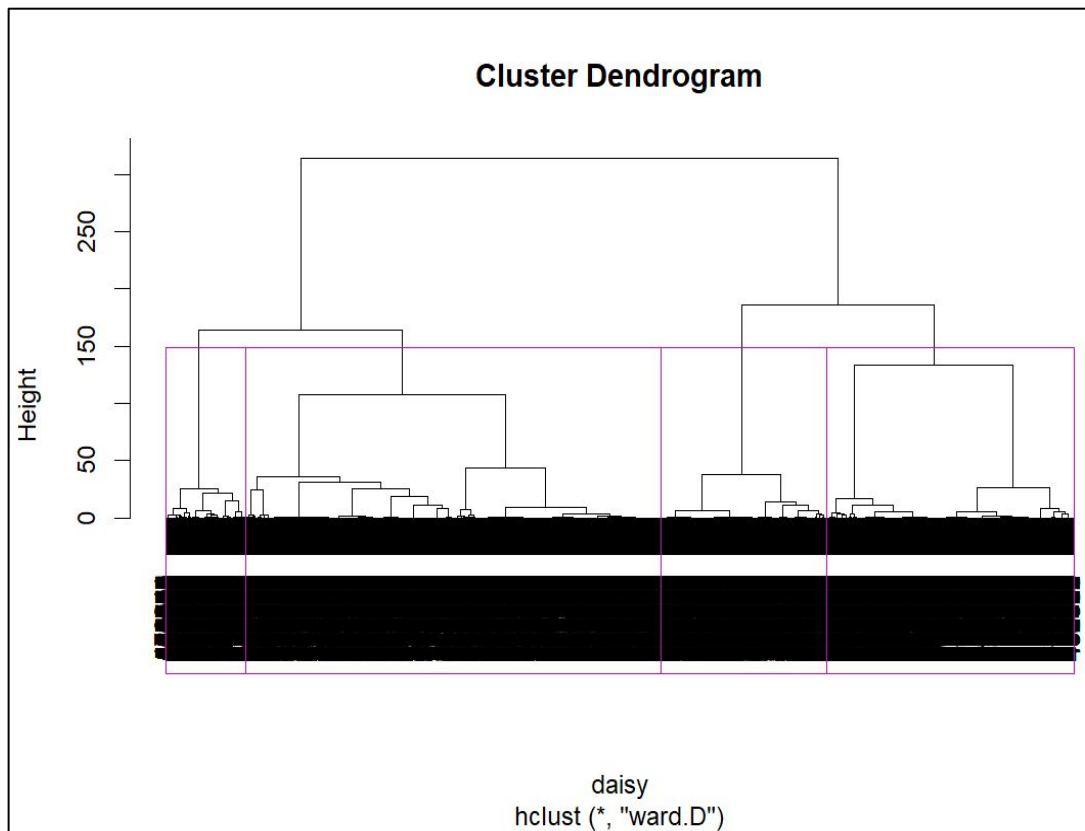    - cindex, dunn, frey, mcclain
  - Silhouette coefficient



Silhouette plot of (x = cut_k, dist = daisy)

n = 20351

4 clusters $C_j$
$j : n_j \mid ave_{i \in C_j} \, s_i$

1 : 3716 | 0.65

2 : 9319 | 0.47

3 : 5532 | 0.40

4 : 1784 | 0.12

Silhouette width $s_i$

# 8. Profiling of Clusters



Temperature distribution for each cluster

Humidity distribution for each cluster

Means of Precipitation.in. per cluster

Means of Distance.mi. per cluster

# 8. Profiling of Clusters

**Prop. of all levels of  Season within each cluster**



**Prop. of all levels of  Severity within each cluster**



**Prop. of all levels of  Weather_Condition within each cluster**



**Prop. of each level of Weather_Condition within all clusters**

# 8. Profiling of Clusters



Prop. of all levels of Crossing within each cluster

Prop. of all levels of Bump within each cluster

Prop. of all levels of Stop within each cluster

Prop. of all levels of Traffic_Signal within each cluster

**Significant features per cluster**

**Cluster 1-** Distance, Precipitation, Crossing and Traffic Signal

**Cluster 2 -** Humidity and Precipitation, Stop, Weather Condition, Crossing and Severity.

**Cluster 3 -** Distance, Crossing, Traffic Signal and Severity

**Cluster 4 -** Season and Severity

# 9. Decision Trees

Tree 4 and 7



| PredClass tree 4 and 7 | 1 | 2 | 3 |
|---|---|---|---|
| Severity 1 | 3029 | 115 | 37 |
| Severity 2 | 705 | 5564 | 474 |
| Severity 3 | 599 | 778 | 1021 |
| Severity 4 | 398 | 466 | 449 |

# 9. Decision Trees

Tree 2



| PredClass tree 2 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Severity 1 | 3144 | 0 | 19 | 18 |
| Severity 2 | 870 | 5683 | 165 | 25 |
| Severity 3 | 638 | 1078 | 547 | 135 |
| Severity 4 | 412 | 565 | 124 | 212 |

# 10. LDA

- <u>Explanatory Continuous Variables Used</u>:
  - Temperature(F), Humidity(%), Pressure(in), Precipitation(in), Distance(mi) and Visibility(mi)
  - Number of variables: 6
- <u>Qualitative Target Variable</u>:
  - Severity
  - 4 Cardinality Levels: Severity1, Severity2, Severity3, Severity4
- <u>Splitting Dataset into Training - Testing Sets</u>:
  - 66% Training Set (13635 observations)
  - 33% Test Set (6716 observations)
- <u>Training LDA on the Training Set</u>
  - Separation explained:

| | Discriminant Functions | | |
|---|---|---|---|
| | **LD1** | **LD2** | **LD3** |
| **Separation Explained** | 62% | 30% | 8% |

# 10. LDA

**LDA1 and LDA2 Projection of Training Data**

# 10. LDA

**Stacked Histogram of the 1st Discriminant Function Values**

**Stacked Histogram of the 2nd Discriminant Function Values**

# 10. LDA

## Training Set - Validation Metrics
- Accuracy: 0.49
- MCR: 0.51
- Precision-Recall-F1 Table

| | Severity1 | Severity2 | Severity3 | Severity4 |
|---|---|---|---|---|
| **Precision** | 0.269 | 0.5046 | 0.340 | 0.407 |
| **Recall** | 0.006 | **0.929** | 0.143 | 0.040 |
| **F1** | 0.012 | **0.6540** | 0.2022 | 0.073 |

# 10. LDA

## Testing Set - Validation Metrics

- Accuracy: 0.48
- MCR: 0.52
- Precision-Recall-F1 Table

| | Severity1 | Severity2 | Severity3 | Severity4 |
|---|---|---|---|---|
| **Precision** | 0.339 | 0.496 | 0.298 | 0.4659 |
| **Recall** | 0.011 | **0.9248** | 0.1265 | 0.0627 |
| **F1** | 0.0216 | **0.6460** | 0.1776 | 0.1106 |

# 11. Gantt Diagram



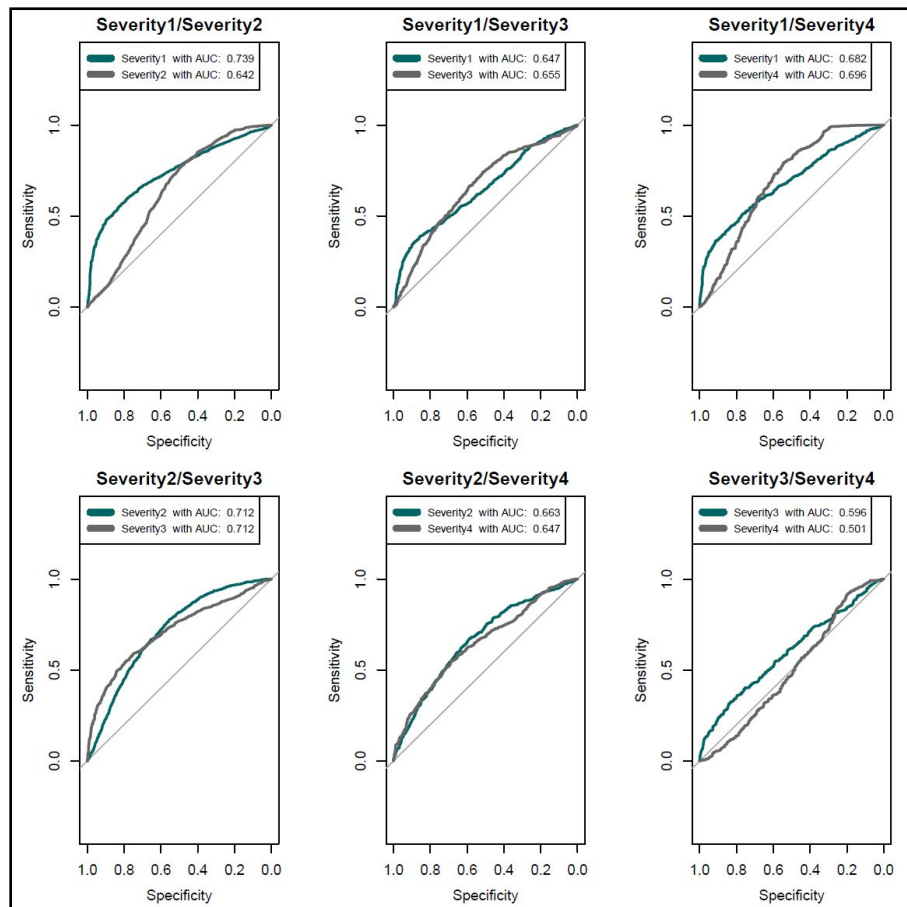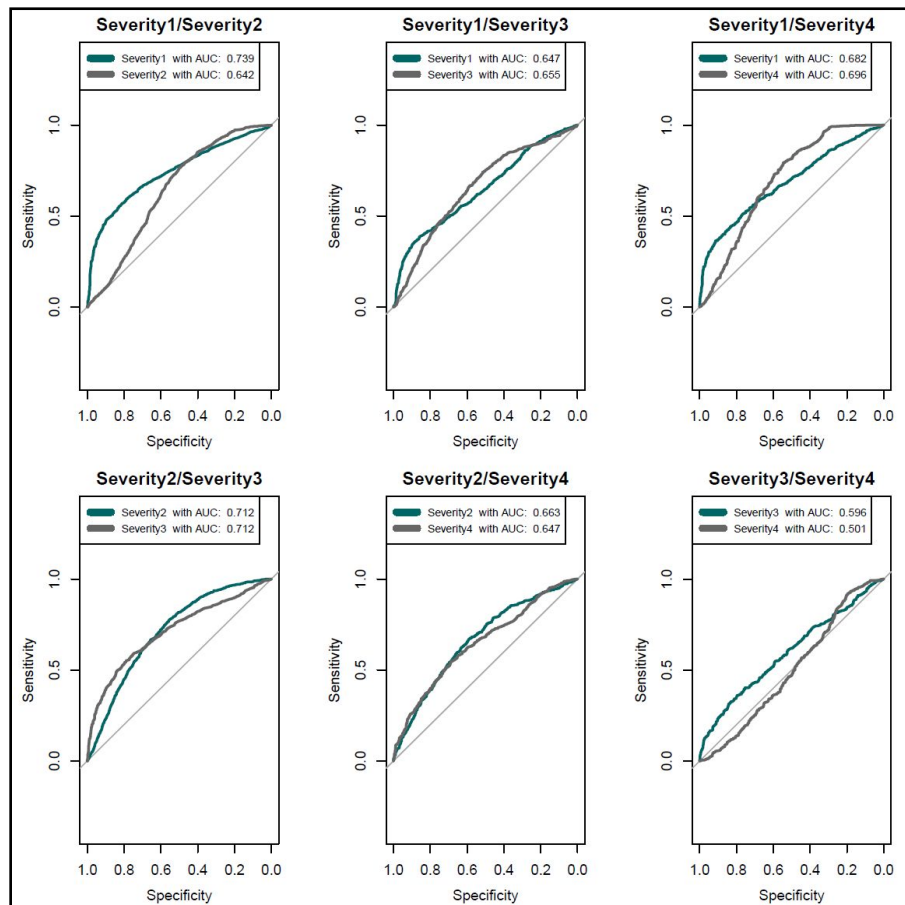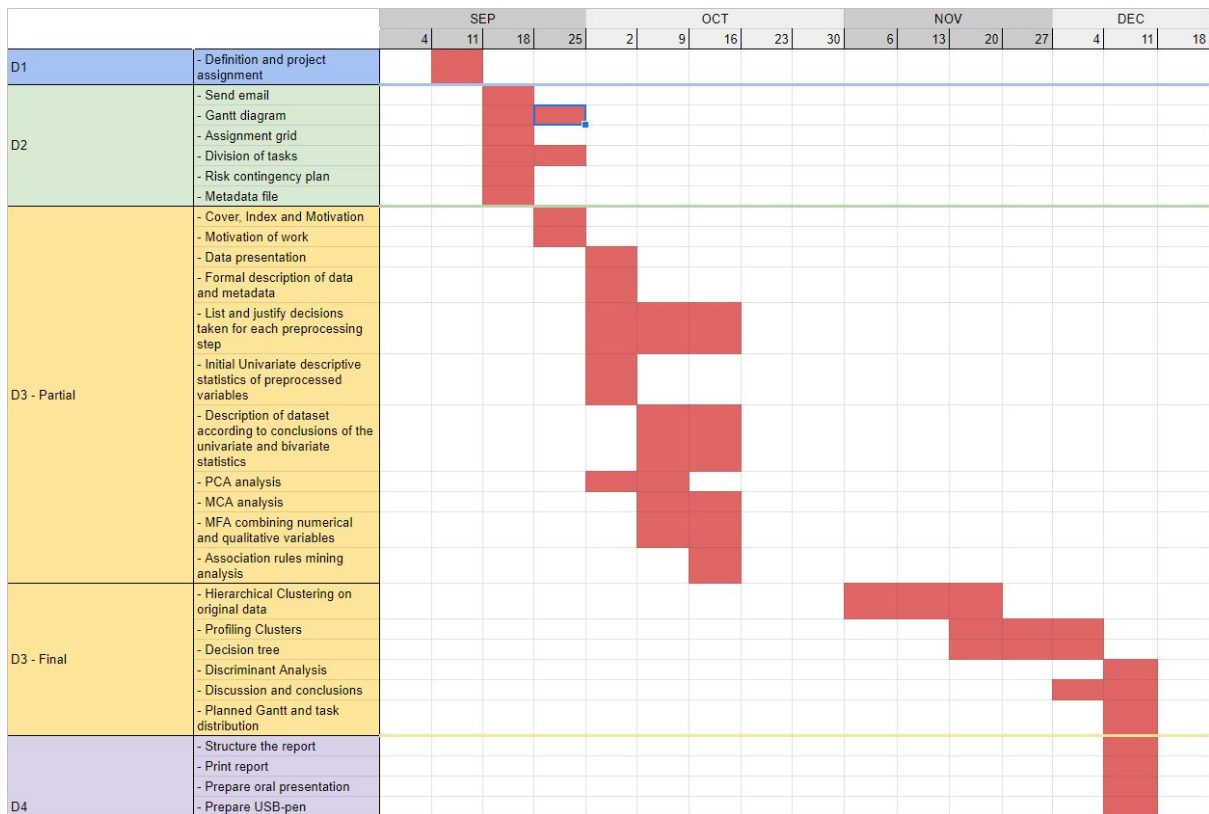| | | SEP | | | | OCT | | | | | NOV | | | | DEC | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 4 | 11 | 18 | 25 | 2 | 9 | 16 | 23 | 30 | 6 | 13 | 20 | 27 | 4 | 11 | 18 |
| D1 | - Definition and project assignment | | ■ | | | | | | | | | | | | | | |
| D2 | - Send email | | | ■ | | | | | | | | | | | | | |
| | - Gantt diagram | | | | ■ | | | | | | | | | | | | |
| | - Assignment grid | | | | ■ | | | | | | | | | | | | |
| | - Division of tasks | | | | ■ | | | | | | | | | | | | |
| | - Risk contingency plan | | | ■ | | | | | | | | | | | | | |
| | - Metadata file | | | ■ | | | | | | | | | | | | | |
| D3 - Partial | - Cover, Index and Motivation | | | | ■ | | | | | | | | | | | | |
| | - Motivation of work | | | | ■ | | | | | | | | | | | | |
| | - Data presentation | | | | | ■ | | | | | | | | | | | |
| | - Formal description of data and metadata | | | | | ■ | | | | | | | | | | | |
| | - List and justify decisions taken for each preprocessing step | | | | | ■ | ■ | | | | | | | | | | |
| | - Initial Univariate descriptive statistics of preprocessed variables | | | | | ■ | | | | | | | | | | | |
| | - Description of dataset according to conclusions of the univariate and bivariate statistics | | | | | | ■ | ■ | | | | | | | | | |
| | - PCA analysis | | | | | ■ | | | | | | | | | | | |
| | - MCA analysis | | | | | | ■ | | | | | | | | | | |
| | - MFA combining numerical and qualitative variables | | | | | | ■ | | | | | | | | | | |
| | - Association rules mining analysis | | | | | | ■ | | | | | | | | | | |
| D3 - Final | - Hierarchical Clustering on original data | | | | | | | | | | ■ | ■ | | | | | |
| | - Profiling Clusters | | | | | | | | | | | ■ | ■ | | | | |
| | - Decision tree | | | | | | | | | | | | ■ | | | | |
| | - Discriminant Analysis | | | | | | | | | | | | ■ | | | | |
| | - Discussion and conclusions | | | | | | | | | | | | | ■ | | | |
| | - Planned Gantt and task distribution | | | | | | | | | | | | | ■ | | | |
| D4 | - Structure the report | | | | | | | | | | | | | | ■ | | |
| | - Print report | | | | | | | | | | | | | | ■ | | |
| | - Prepare oral presentation | | | | | | | | | | | | | | ■ | | |
| | - Prepare USB-pen | | | | | | | | | | | | | | ■ | | |

# 12. Final Conclusions

- Data was challenging to manage and to extract firm conclusions

- Weather Condition and Traffic Signs do not influence severity of the accident

- Location makes a difference in the number of accidents

- Occurrence of accident is strongly correlated with human actions

- To improve the analysis human factors (speed of the vehicle, time of accident, …) should be added

Thank you