

A Study on Vehicles' accidents in the United Kingdom

Laura C. Carles G.

Universitat Politecnica de Catalunya

June 13, 2018

Structure

Our work this semester

- ① Improvements to the Naive Bayes Classifier
- ② Accidents data preprocessing
- ③ Multivariate Data Visualization techniques

Structure

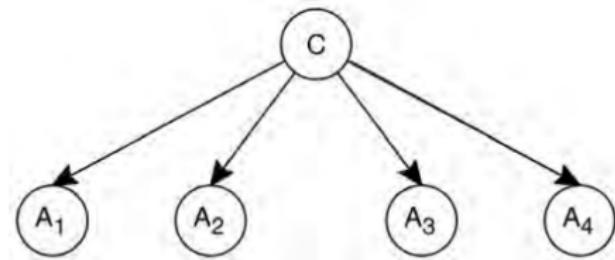
Our work this semester

- ① Improvements to the Naive Bayes Classifier
- ② Accidents data preprocessing
- ③ Multivariate Data Visualization techniques

Naive Bayes

Brief introduction

- C : Class to be estimated
- $A_x, x \in X$, where
 - X the number of attributes.
 - A_x attribute.



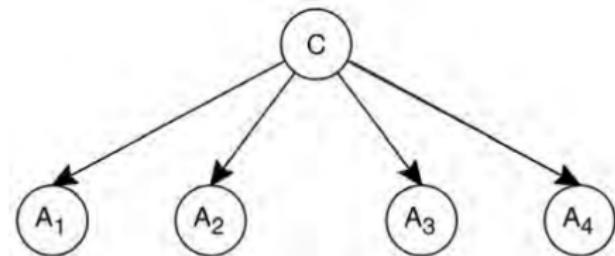
The estimated class can be computed as the class with maximum probability between classes given a_i .

$$c(E) = \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^n P(a_i|c) \quad (1)$$

Naive Bayes

Brief introduction

- C : Class to be estimated
- $A_x, x \in X$, where
 - X the number of attributes.
 - A_x attribute.



The estimated class can be computed as the class with maximum probability between classes given a_i .

$$c(E) = \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^n P(a_i|c) \quad (1)$$

Conditional independence assumption.

Independence assumption

What if data is dependent?

- ① *IRL* data is not independent

Independence assumption

What if data is dependent?

- ① *IRL* data is not independent
- ② Nevertheless, although being a simple approach, **it works quite well**

Independence assumption

What if data is dependent?

- ① IRL data is not independent
- ② Nevertheless, although being a simple approach, **it works quite well**

The big question

Is there a way to improve Naive Bayes to take into account attribute dependencies?

Independence assumption

What if data is dependent?

- ① *IRL* data is not independent
- ② Nevertheless, although being a simple approach, **it works quite well**

The big question

Is there a way to improve Naive Bayes to take into account attribute dependencies?

More importantly...

If so, **would that give us better predictions?**

Improving NB

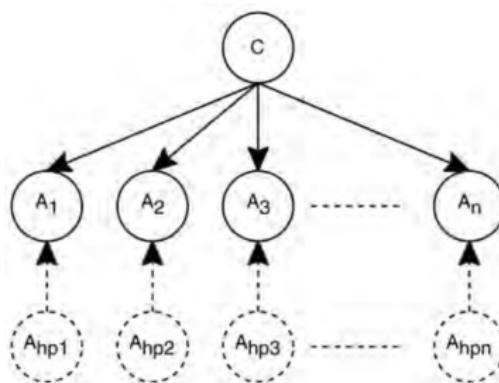
Many different extensions on NB

- Tree Augmented Naive Bayes (TAN)
- NbTree
- Averaged one-dependence estimators (AODE)
- Selective Bayesian Classifiers (SBC)
- **Hidden Naive Bayes (HNB)**

Hidden Naive Bayes

Its inner workings

- C : Class to be estimated
- $A_i, i \in I$, where
 - N the number of attributes.
 - A_i attribute.
 - A_{hp_i} Weight dependency as hidden parent.



Predicted class can be computed as the class with maximum probability between classes given a_i and the relationship with all other attributes a_{hp_i} .

$$c(E) = \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^n P(a_i | a_{hp_i}, c) \quad (2)$$

Hidden Naive Bayes

Its inner workings II

- Taking the influences from other factors as a hidden parent:

$$P(a_i|a_{hp_i}, c) = \sum_{j=1; j \neq i}^n W_{ij} * P(a_i|a_j, c)$$

Hidden Naive Bayes

Its inner workings II

- Taking the influences from other factors as a hidden parent:

$$P(a_i|a_{hp_i}, c) = \sum_{j=1; j \neq i}^n W_{ij} * P(a_i|a_j, c)$$

- Weights are expressed as the conditional mutual information between variables:

$$W_{ij} = \frac{I_p(A_i; A_j | C)}{\sum_{j=1, j \neq i}^n I_p(A_i; A_j | C)}$$

Hidden Naive Bayes

Its inner workings III

- The conditional mutual information can be calculated as:

$$I_p(A_i; A_j | C) = \sum_{a_i, a_j, c} P(a_i, a_j, c) \log \frac{P(a_i, a_j | c)}{P(a_i | c)P(a_j | c)}$$

Hidden Naive Bayes

Its inner workings III

- The conditional mutual information can be calculated as:

$$\text{Ip}(A_i; A_j | C) = \sum_{a_i, a_j, c} P(a_i, a_j, c) \log \frac{P(a_i, a_j | c)}{P(a_i | c)P(a_j | c)}$$

Computing probabilities

- Avoiding zero frequency:

$$P(c) = \frac{F(c) + \frac{1}{\#\text{classes}}}{\#\text{train} + 1}$$

$$P(a_i | a_j, c) = \frac{F(a_i, a_j, c) + \frac{1}{\#\text{val}_i}}{F(a_j, c) + 1}$$

Hidden Naive Bayes

Results

- Titanic dataset



Hidden Naive Bayes

Results

- Titanic dataset



No real independence between
variables!

Hidden Naive Bayes

Results

- Titanic dataset



Attr.	NB (%)	HNB (%)
Class	41	46
Sex	73	78
Age	60	95
Survived	73	76

No real independence between variables!

Structure

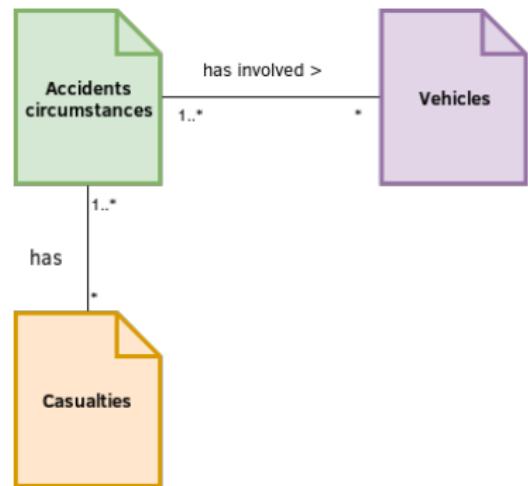
Our work this semester

- ① Improvements to the Naive Bayes Classifier
- ② Accidents data preprocessing
- ③ Multivariate Data Visualization techniques

First glance at data

Files

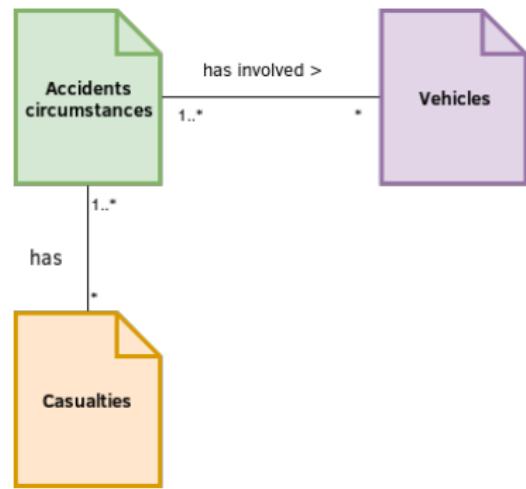
- ① **Accident circumstances**, general circumstances and main information.
- ② **Vehicles involved**.
- ③ **Casualties**, injured and death.



First glance at data

Files

- ① **Accident circumstances**, general circumstances and main information.
- ② **Vehicles involved**.
- ③ **Casualties**, injured and death.



No direct join of data!

First glance at data

Dimensionality

	# Individuals	# Variables
Accident	136.621	32
Vehicles	252.500	24
Casualties	181.384	16
Total	570.505	72

First glance at data

Dimensionality

	# Individuals	# Variables
Accident	136.621	32
Vehicles	252.500	24
Casualties	181.384	16
Total	570.505	72

Pretty large dataset

First glance at data

Missing values

At least 57.000 rows with missing data!

Main issues

to be dealt with

- ① There is no direct join between files.
- ② Dataset is considerably large.
- ③ Many missing values.
- ④ Very complex variables
 - Mainly categorical
 - More than 7 different categories

Preprocessing needed

Preprocessing

Dimensionality reduction

Preprocessing - Dimensionality reduction

Removing variables

- **Objectively:** mostly empty
- **Subjectively:** uninteresting for us or not crucial.

Preprocessing - Dimensionality reduction

Removing variables

- **Objectively:** mostly empty
- **Subjectively:** uninteresting for us or not crucial.
 - Day of week - can be derived from date.
 - Local Authority District & Highway - mostly empty.
 - Road number & class.
 - Pedestrian Crossing Human Control - Only 1% had human control.
 - Police force
 - Urban or rural Area.
 - ...

Preprocessing - Dimensionality reduction

Missing values treatment

Still, many rows with still missing values.

	# rows	# NA	# rows NA	% rows NA
Accidents	136.575	1.939	1.083	0,79
Vehicles	252.500	97.971	81.748	32,38
Casualties	181.384	2.898	2.862	1,5

Preprocessing - Dimensionality reduction

Missing values treatment

Still, many rows with still missing values.

	# rows	# NA	# rows NA	% rows NA
Accidents	136.575	1.939	1.083	0,79
Vehicles	252.500	97.971	81.748	32,38
Casualties	181.384	2.898	2.862	1,5

Decisions taken

- ① Remove rows in *Casualties* and *Vehicles* with NA.
- ② Remove rows with unimputable values (*Latitude*, *Longitude*, *Time*, *Speed Limit*).
- ③ Impute with Knn *Accidents*.

Preprocessing

Feature processing

Preprocessing - Feature processing

Too many possible values that can be reduced to TRUE/FALSE or GOOD/BAD.

Simplifiable variables

- Light Conditions
- Weather Conditions
- Surface Conditions
- Special Conditions at site
- Carriageway Hazards

Preprocessing - Feature processing

Example

Weather

- ① ✓ Fine no high winds
- ② ✗ Raining no high winds
- ③ ✗ Snowing no high winds
- ④ ✗ Fine + high winds
- ⑤ ✗ Raining + high winds
- ⑥ ✗ Snowing + high winds
- ⑦ ✗ Fog or mist
- ⑧ ✗ Other



- ① Good weather (1)
- ② Bad Weather [2 - 8]

Preprocessing - Feature processing

Example

Road Surface Conditions

- | | | |
|--|---|--|
| <ul style="list-style-type: none">① ✓ Dry② ✗ Wet or damp③ ✗ Snow④ ✗ Frost or ice⑤ ✗ Flood over 3cm⑥ ✗ Oil or diesel⑦ ✗ Mud |  | <ul style="list-style-type: none">① Good Road Surface Conditions (1)② Bad Road Surface Conditions [2 - 7] |
|--|---|--|

Preprocessing - Feature processing

Example

Road Surface Conditions

- ① ✓ Dry
 - ② ✗ Wet or damp
 - ③ ✗ Snow
 - ④ ✗ Frost or ice
 - ⑤ ✗ Flood over 3cm
 - ⑥ ✗ Oil or diesel
 - ⑦ ✗ Mud
-
- ① Good Road Surface Conditions (1)
 - ② Bad Road Surface Conditions [2 - 7]

Many more...

Preprocessing - Feature processing

Vehicles type reduction

20 different types of vehicles!

- ① Pedal cycle
- ② Motorcycle 50cc and under
- ③ Motorcycle 125cc and under
- ④ Motorcycle over 125cc and up to 500cc
- ⑤ Motorcycle over 500cc
- ⑥ Taxi/Private hire car
- ⑦ Car
- ⑧ Minibus (8 - 16 passenger seats)
- ⑨ Bus or coach (17 or more passenger seats)
- ⑩ Ridden horse
- ⑪ Agricultural vehicle
- ⑫ Tram
- ⑬ Van / Goods 3.5 tonnes mgw or under
- ⑭ Goods over 3.5t. and under 7.5t
- ⑮ Goods 7.5 tonnes mgw and over
- ⑯ Mobility scooter
- ⑰ Electric motorcycle
- ⑱ Other vehicle
- ⑲ Motorcycle - unknown cc
- ⑳ Goods vehicle - unknown weight

Preprocessing - Feature processing

Vehicles type reduction

Example of Motorcycles grouping

- Motorcycle 50cc and under
 - Motorcycle 125cc and under
 - Motorcycle over 125cc and
up to 500cc
 - Motorcycle over 500cc
 - Mobility scooter
 - Electric motorcycle
- • Motorcycle

Preprocessing - Feature processing

Vehicles type reduction

Example of Cars grouping

- Van / Goods 3.5 tonnes
mgw or under → • Car
- Taxi/Private hire car
- Car

Preprocessing - Feature processing

Vehicles type reduction

Example of Cars grouping

- Van / Goods 3.5 tonnes
mgw or under → • Car
- Taxi/Private hire car
- Car

And more...

Preprocessing

Files join

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident)

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident) → **Count number of vehicles per accident.**

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident) → **Count number of vehicles per accident.**
- ② Aggregation applied to:
 - # Male drivers.
 - # Female drivers.
 - # Male casualties.
 - # Female casualties.
 - ...

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident) → **Count number of vehicles per accident.**
- ② Aggregation applied to:
 - # Male drivers.
 - # Female drivers.
 - # Male casualties.
 - # Female casualties.
 - ...
- ③ Computed the age disparity between
 - Drivers.
 - Vehicles.

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident) → **Count number of vehicles per accident.**
- ② Aggregation applied to:
 - # Male drivers.
 - # Female drivers.
 - # Male casualties.
 - # Female casualties.
 - ...
- ③ Computed the age disparity between
 - Drivers.
 - Vehicles.
- ④ Vehicles type joining.

Preprocessing - Files join

Steps

- ① Correspondence of many to one. (Many Vehicles involved to one accident) → **Count number of vehicles per accident.**
- ② Aggregation applied to:
 - # Male drivers.
 - # Female drivers.
 - # Male casualties.
 - # Female casualties.
 - ...
- ③ Computed the age disparity between
 - Drivers.
 - Vehicles.
- ④ Vehicles type joining.
- ⑤ Worst casualty in accident (*Fatal? Slightly injured?*)
- ⑥ ...

Preprocessing - Files join

Aggregation examples

Max age disparity example

- Driver 1: 60 years old. (MAX)
- Driver 2: 18 years old. (MIN)
- Driver 3: 30 years old.

Preprocessing - Files join

Aggregation examples

Max age disparity example

- Driver 1: 60 years old. (MAX)
- Driver 2: 18 years old. (MIN)
- Driver 3: 30 years old.
- Disparity: $\text{MAX} - \text{MIN} = 42$ years.

Preprocessing - Files join

Aggregation examples

Max age disparity example

- Driver 1: 60 years old. (MAX)
- Driver 2: 18 years old. (MIN)
- Driver 3: 30 years old.
- Disparity: $\text{MAX} - \text{MIN} = 42$ years.

Watch out when interpreting results. When disparity is 0, it may be because there is only one vehicle involved.

Preprocessing - Files join

Aggregation examples

Max age disparity example

- Driver 1: 60 years old. (MAX)
- Driver 2: 18 years old. (MIN)
- Driver 3: 30 years old.
- Disparity: $\text{MAX} - \text{MIN} = 42$ years.

Watch out when interpreting results. When disparity is 0, it may be because there is only one vehicle involved.

Accident type joining decision

- ① Car + Car → Car Accident
- ② Motorcycle + Motorcycle → Motorcycle Accident
- ③ Car + Motorcycle → Mixed
- ④ ...

Preprocessing

Preprocessing Summary

Preprocessing

Main steps summary

① Dimensionality reduction

Step applied

- ① Interesting vars. selection
- ② 57k rows with NA in one variable - **Removed var.**

Preprocessing

Main steps summary

- ① Dimensionality reduction
- ② Feature processing

Step applied

- ① Interesting vars. selection
- ② 57k rows with NA in one variable - **Removed var.**
- ③ Simplification of categorical vars

Preprocessing

Main steps summary

- ① Dimensionality reduction
- ② Feature processing
- ③ Missing Values treatment

Step applied

- ① Interesting vars. selection
- ② 57k rows with NA in one variable - **Removed var.**
- ③ Simplification of categorical vars
- ④ Imputation only on accidents with Knn (1083 rows, 0.79%)

Preprocessing

Main steps summary

- ① Dimensionality reduction
- ② Feature processing
- ③ Missing Values treatment
- ④ **Joining**

Step applied

- ① Interesting vars. selection
- ② 57k rows with NA in one variable - **Removed var.**
- ③ Simplification of categorical vars
- ④ Imputation only on accidents with Knn (1083 rows, 0.79%)
- ⑤ Aggregation of vehicles and casualties.

Preprocessing

Result

Preprocessing result

Separated files dimensionality

	# Ind	# Vars	% Info. ret.	% Info. lost
Accidents	136.575	17	53,125	46,875
Vehicles	252.500	5	20,83	79,17
Casualties	181.384	5	31,25	68,75

Preprocessing result

Separated files dimensionality

	# Ind	# Vars	% Info. ret.	% Info. lost
Accidents	136.575	17	53,125	46,875
Vehicles	252.500	5	20,83	79,17
Casualties	181.384	5	31,25	68,75

Final dimensionality (joined)

	# Ind.	# Vars
Final ds.	136.575	31

Conclusions

- ➊ A good and thorough preprocessing of a dataset this large takes a lot of time.

Conclusions

- ① A good and thorough preprocessing of a dataset this large takes a lot of time.
- ② Must take many decisions.
- ③ Missing data. Not always imputable.

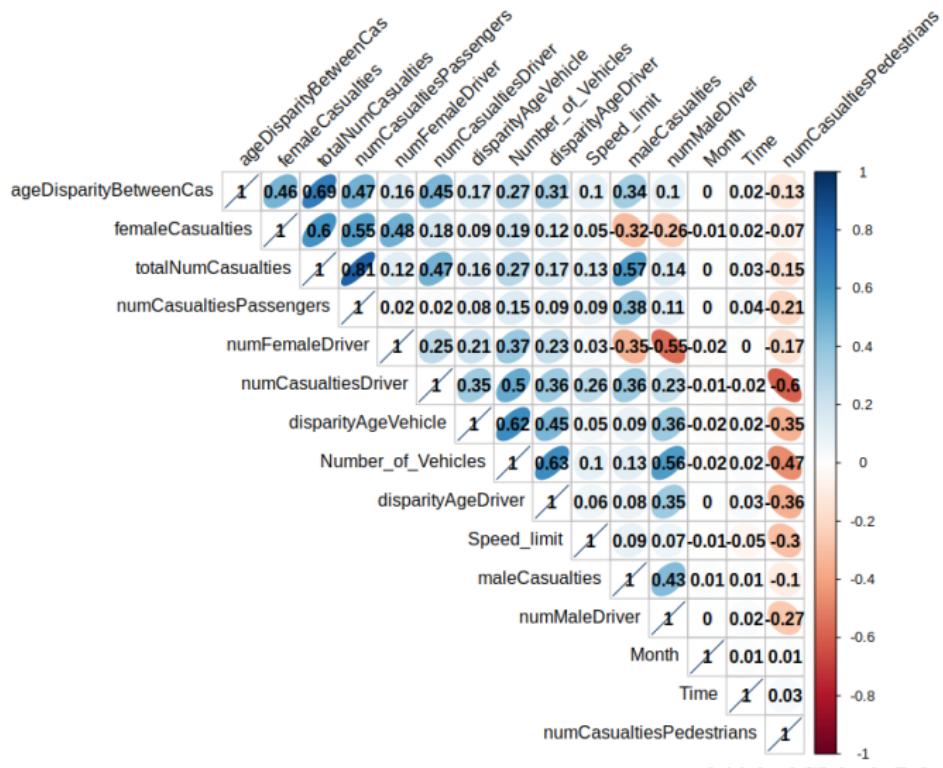
Conclusions

- ① A good and thorough preprocessing of a dataset this large takes a lot of time.
- ② Must take many decisions.
- ③ Missing data. Not always imputable.
- ④ Many source files → High difficulty on managing the data.
- ⑤ Joining is not direct.

The greater the number of variables, the more complex the decisions to take about it.

Multivariate Plotting

Involving variables correlation



Applying the HNB classifier

Attribute as class	% accuracy NB	% accuracy HNB
Accident Severity	67.05	79.04
Road Type	8.15	82.89
Speed Limit	11.64	66.75
Road conditions	78.31	83.46
Weather conditions	82.39	82.29

Structure

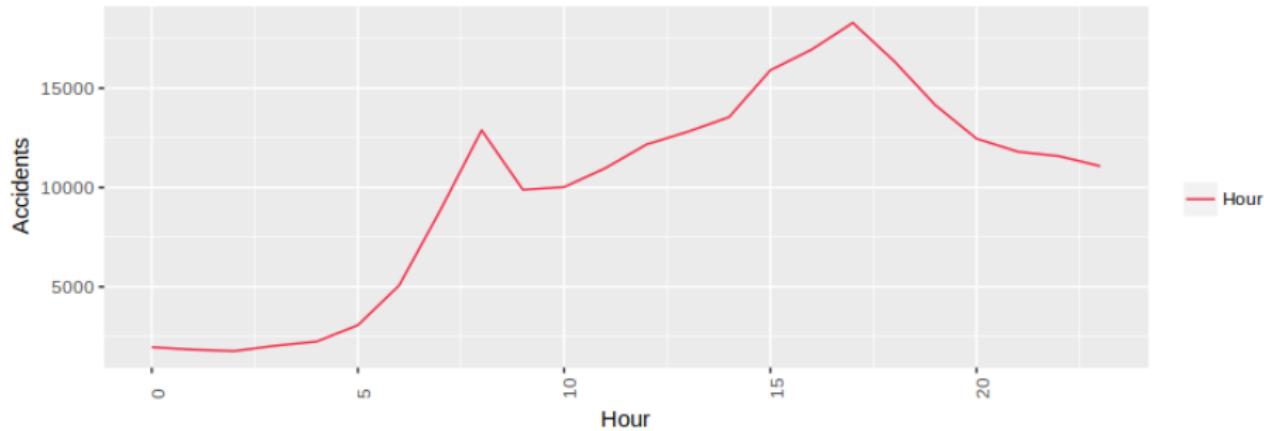
Our work this semester

- ① Improvements to the Naive Bayes Classifier
- ② Accidents data preprocessing
- ③ Multivariate Data Visualization techniques

Simple Plotting

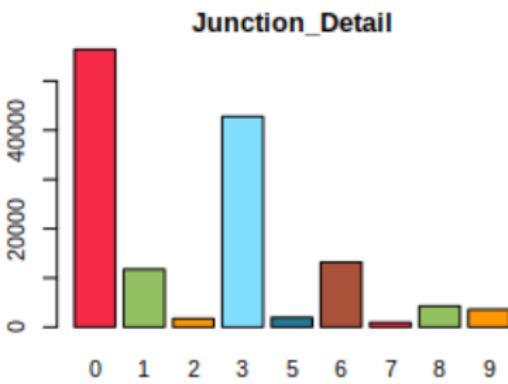
Accidents per hour evolution

Accidents during year by hour



Simple Plotting

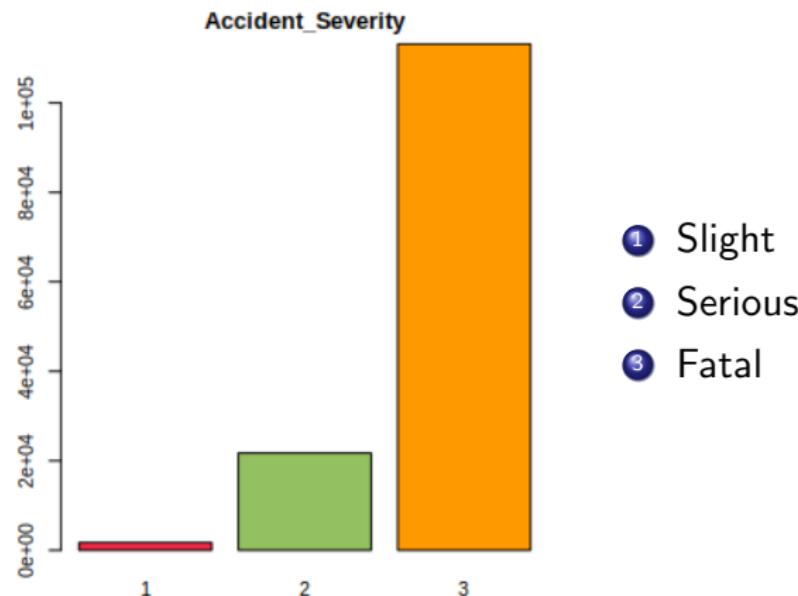
Histograms



- ① Not at junction or within 20 metres
- ② Roundabout
- ③ Mini-roundabout
- ④ T or staggered junction
- ⑤ Slip road
- ⑥ Crossroads
- ⑦ More than 4 arms (not roundabout)
- ⑧ Private drive or entrance
- ⑨ Other junction

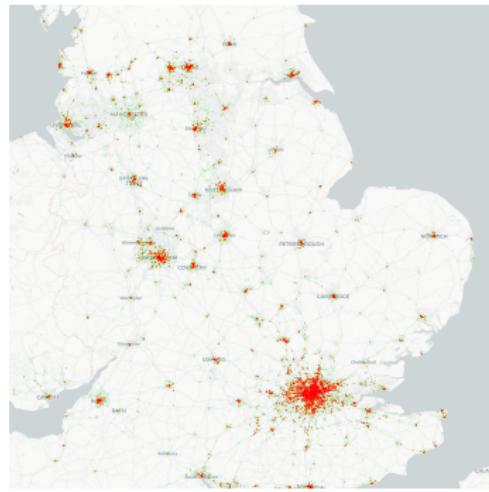
Simple Plotting

Histograms



Multivariate Plotting

Heatmap representation

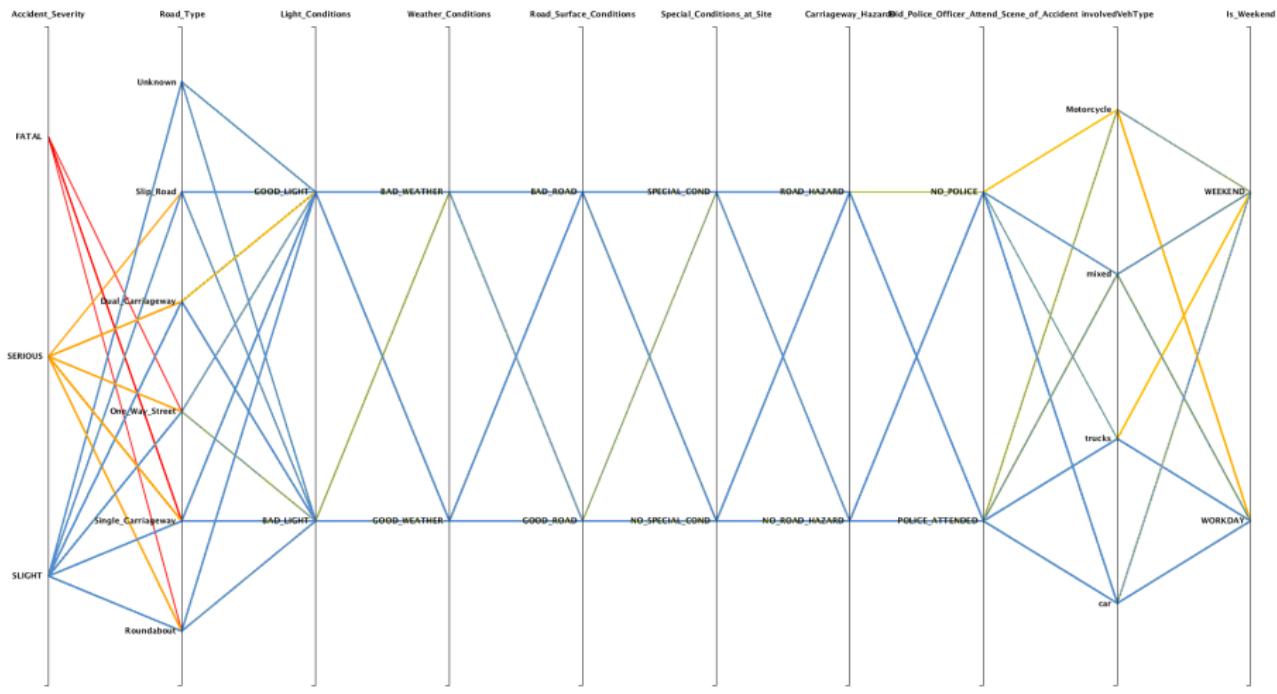


Density of accidents remains in the cities!

Multivariate Plotting

KNIME parallel sets

Parallel Coordinates Plot



Beyond Plotting

Components Analysis

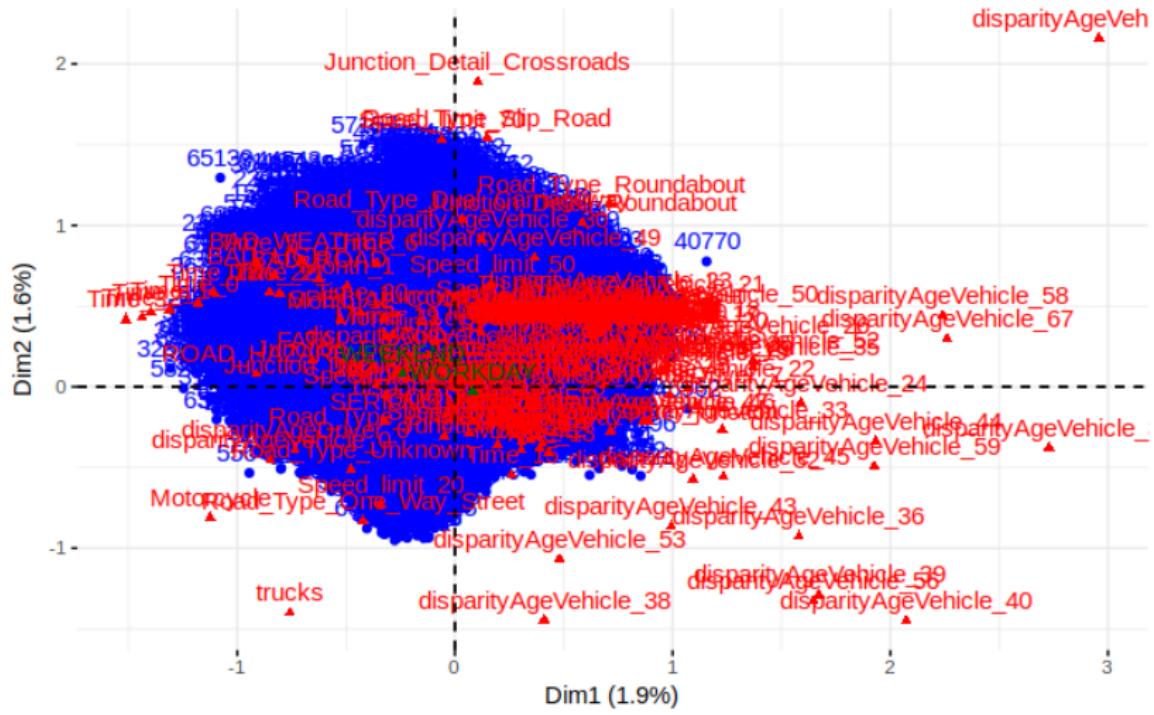
Mixed type attributes

- Principal Components analysis (PCA)
- Multiple Correspondance analysis (MCA)

Beyond Plotting

MCA - Results

MCA - Biplot



Beyond Plotting

MCA - Results

- **Low Inertia** with the factorial plane (3%).
- Dimensionality reduction, retaining up to 90% of inertia.
- **Homogeneous cloud of individuals** in the biplot. No feasible explanation at first. Clustering will be necessary to distinguish clusters.

Beyond Plotting

Clustering

Clustering Approach

- No time based division \implies dealing with the whole dataset.

Beyond Plotting

Clustering

Clustering Approach

- No time based division \implies dealing with the whole dataset.
- **Too much** data for a hierarchical clustering

Beyond Plotting

Clustering

Clustering Approach

- No time based division \implies dealing with the whole dataset.
- **Too much** data for a hierarchical clustering
- Alternative methods: Iterative kmeans

Beyond Plotting

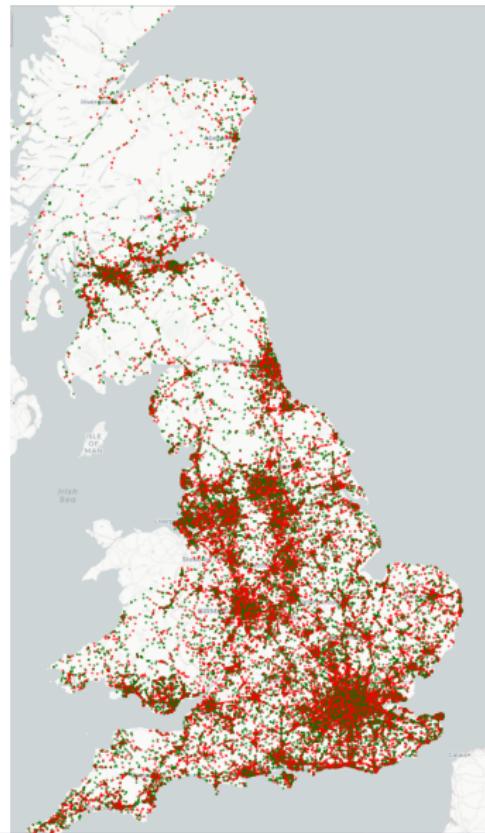
Clustering

Clustering Approach

- No time based division \implies dealing with the whole dataset.
- **Too much** data for a hierarchical clustering
- Alternative methods: Iterative kmeans
- Discovery # num clusters?

Beyond Plotting

Clustering

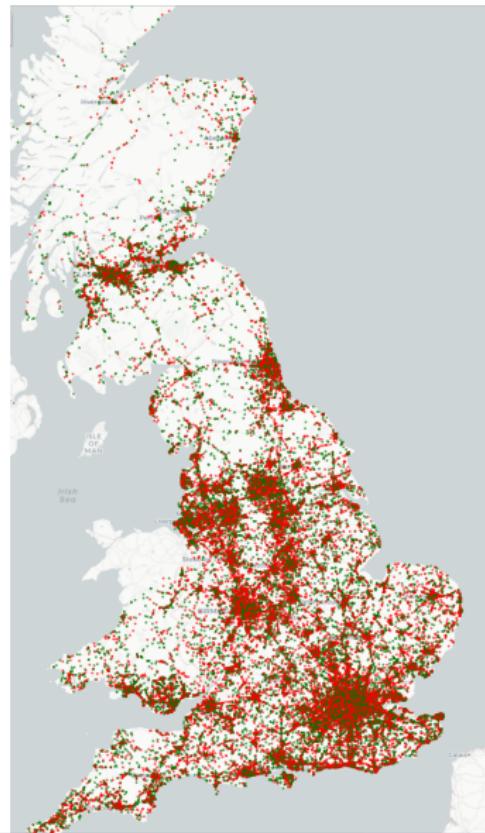


Identified clusters:

- Mixed vehicle type, junction based. No special conditions. Different vehicle age.

Beyond Plotting

Clustering

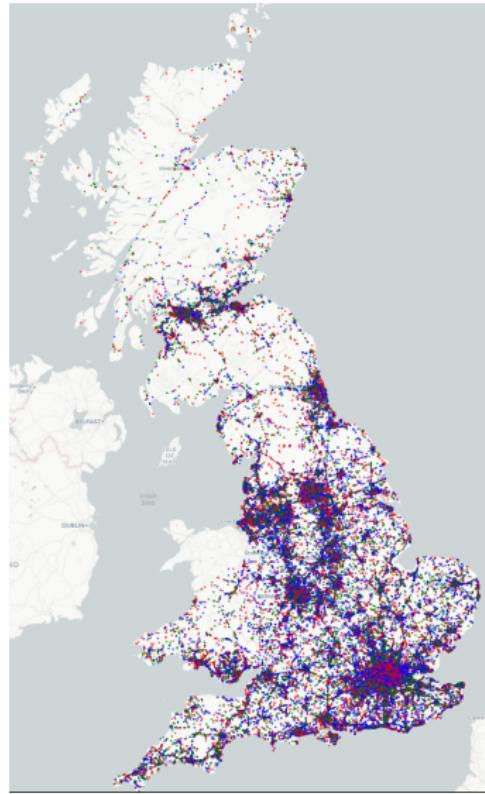


Identified clusters:

- Mixed vehicle type, junction based. No special conditions. Different vehicle age.
- Low speed, truck / motorcycles involved. Accident happened around midnight.

Beyond Plotting

Clustering

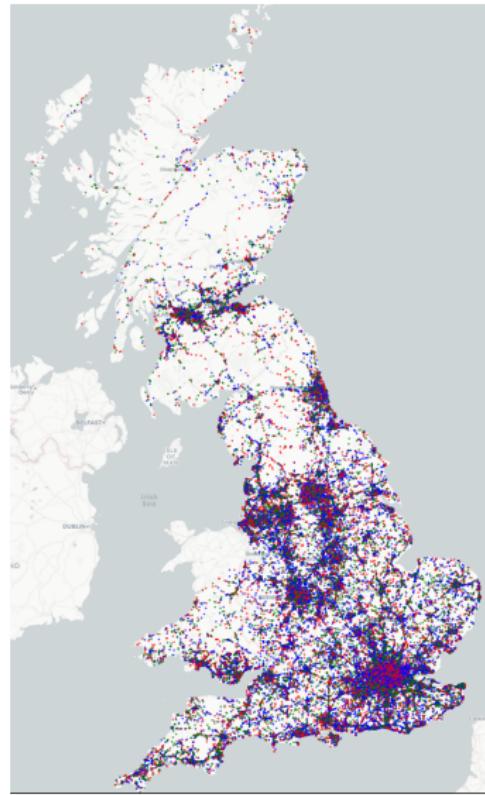


Identified clusters:

- Low speed, truck / motorcycle mixed with good conditions.

Beyond Plotting

Clustering

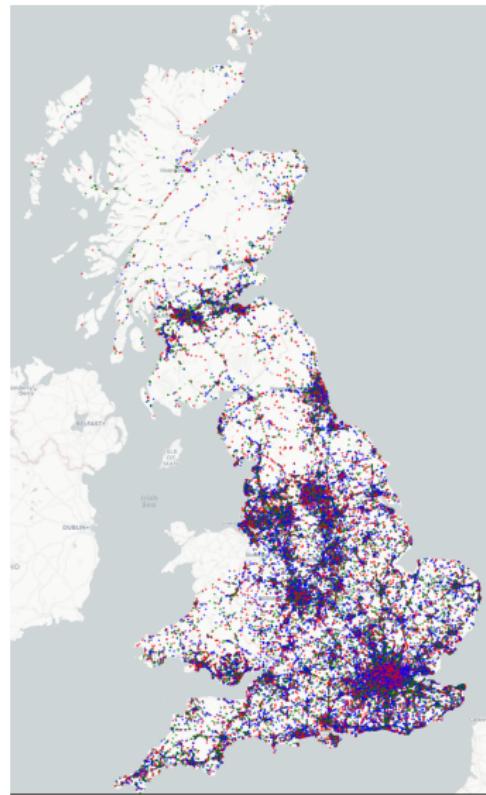


Identified clusters:

- Low speed, truck / motorcycle mixed with good conditions.
- High speed, winter based with bad conditions.

Beyond Plotting

Clustering



Identified clusters:

- Low speed, truck / motorcycle mixed with good conditions.
- High speed, winter based with bad conditions.
- Mixed type, near junction, vehicle age disparity.

Possible future work

Looking back, we would change a couple of things and add others.

To change

- ① Change the age disparity → Has many problems and it dirtied our MCA.
- ② Reselect important variables → i.e. Day of week (recomputed afterwards).

Possible future work

Looking back, we would change a couple of things and add others.

To change

- ① Change the age disparity → Has many problems and it dirtied our MCA.
- ② Reselect important variables → i.e. Day of week (recomputed afterwards).

To add

- ① Time series.
- ② Explore other clustering methods.
- ③ Explore powerful tools and frameworks for Data Viz. (i.e. Tableau)
- ④ Focus on vehicles, not accidents circumstances
- ⑤ Focus on casualties profiles

"If we torture the data long enough, it will confess"

- Darrell Huff

How to Lie With Statistics

1954