

# Accidents with cyclists — A guide to avoid them

Author: Patrick Queiroz dos Anjos

## Business Understanding

The use of bicycles as a means of transport, as sports or just for recreation, has increased dramatically during the recent years. The bicycle is a vehicle that has no pollution (direct, in its use) and, therefore, is considered one of the cleanest means of transport, in terms of the environment, to use.

But, unfortunately like any other means of transport, the bicycle and the cyclist are liable to suffer some accidents, from simple impacts to severe actions. To ride a bicycle a cyclist must pay attention to several issues on the road, such as weather conditions on the day.

These conditions can help in the fact that in a rain condition, for example, the tendency that the cyclist has to lose control of his bicycle is greater, due to the aspect that, at a given speed, the control of his vehicle is impaired. Other climates are more severe, such as snow and milder ones, such as its use on sunny days.

Lighting in front of the road where the bicycle will be used is also important. Road conditions must be suitable for use, such as a well-lit environment. Environments in low light or with different colors, such as at dusk, can be prone to possible accidents.

The environment where the cyclist will be used, usually on roads, is also relevant. Conditions where there is a road with sand, or with a road after rain indicate a greater inclination for accidents. The opposite is also true, as conditions where the road is dry denotes greater initial safety for cyclists.

The environment where the bicycle will be used can also influence accidents, such as at intersections where, by some different means, such as rain and a wet road, they can increase the likelihood of accidents.

So, this project aims to identify certain parameters, such as weather, road conditions, light conditions and the collision address type to demonstrate and indicate patterns where accidents occur with cyclists and thus avoid, or be attentive to, the use of bicycles in the indicated parameters, through of a decision tree.

## Data understanding

The data used in this project will use the .csv file (<https://s3.us.cloud-object-storage.appdomain.cloud/cf-courses-data/CognitiveClass/DP0701EN/version-2/Data-Collisions.csv>) made available in the first week of the Applied Data Science Capstone course, in the “Downloading Example Dataset”. It shows a file where can get various information, such as weather conditions, number of pedestrians involved in the specified accident and the severity of the accident, among others.

The dataset has 194673 different rows with a total of 38 different columns, each detailing some variable inherent to the accident. A summary of the file can be seen at the link <https://s3.us.cloud-object-storage.appdomain.cloud/cf-courses-data/CognitiveClass/DP0701EN/version-2/Metadata.pdf>.

The data that will be collected, processed to suit the later machine learning and primary data analysis processes will be the columns ADDRTYPE, ROADCOND, LIGHTCOND, WEATHER, PEDCYLCOUNT and SEVERITYCODE.

The column ADDRTYPE indicates the collision address type, where it shows the environment where the accident occurred and the ROADCOND column shows the conditions of the road where the accident occurred.

The LIGHTCOND column denotes the lighting conditions in the accident and the WEATHER column indicates the weather conditions that the accident occurred.

The PEDCYLCOUNT column and the SEVERITYCODE are the high impact columns in the final table, a table suitable for the machine learning algorithm, in this case a decision tree.

PEDCYLCOUNT is the applicant table for the accidents involveds with cyclists. There are, clearly, accidents that involved pedestrians only, only the driver of the vehicle involved in the accident, or both. But the focus of this project is on accidents involving cyclists, so this column is extremely important.

SEVERITYCODE is the column that establishes the relationship between the accident and a code that returns a value. This value denotes how severe the accident was. So, therefore, this will be the target column for later machine learning methods.

Then, by means of a dataframe equipped with these specified columns, we will have the process of modeling through machine learning, through the decision tree algorithm.

## Data Preparation

The original dataframe was converted to have only 6 columns, the columns ADDRTYPE, ROADCOND, LIGHTCOND, WEATHER, PEDCYLCOUNT and SEVERITYCODE, described earlier. The ‘df’ dataframe has a total of 187525 rows and the 6 columns chosen. From that choice, some methods of choosing variables, transforming and analyzing were specified:

— The df was specified to have data only with accidents involving cyclists; The target values in the decision tree algorithm, SEVERITYCODE, now only have a value of 1 or 2;

— The accident severity code values are 1 for "Property Damage" and 2 for "Injury". But the amount of data with value 1 ("Property Damage") and value 2 ("Injury") is unbalanced, with 4793 data for value 2 and 676 data for value 1.

Then the dataframe was balanced so that there was 50% of each dataframe that will be submitted to the machine learning process. The Under-Sampling technique was used and, subsequently, the amount of data with value 2 was equal to that of value 1.

The dataframe now has 676 data with values 1 and 2, totaling 1352 rows. Image 1, just below, shows some information from the pre-processed dataframe:

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1352 entries, 192552 to 194585
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   ADDRTYPE        1352 non-null   object
1   ROADCOND        1352 non-null   object
2   LIGHTCOND       1352 non-null   object
3   WEATHER         1352 non-null   object
4   PEDCYLCOUNT     1352 non-null   int64
5   SEVERITYCODE    1352 non-null   int64
dtypes: int64(2), object(4)
memory usage: 73.9+ KB
```

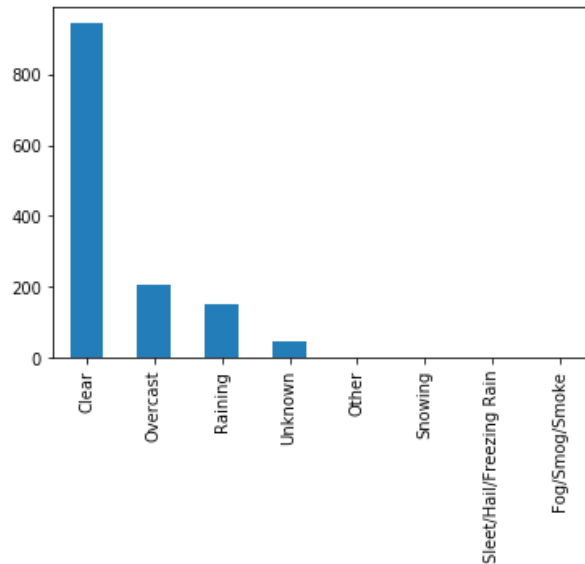
*Figure 1 – Dataset partially pre-processed*

The image shows the number of columns, non-null count, dtype and other parameters for information methods.

After these processes, information was collected about the chosen columns:

— The WEATHER column has 8 values specified in a table and displayed in a bar graph, with the count of each categorical value. They can be seen in image 2 below:

```
(<matplotlib.axes._subplots.AxesSubplot at 0x1c24802e448>,
Clear          944
Overcast       207
Raining        150
Unknown        46
Other          2
Snowing        1
Sleet/Hail/Freezing Rain  1
Fog/Smog/Smoke 1
Name: WEATHER, dtype: int64)
```

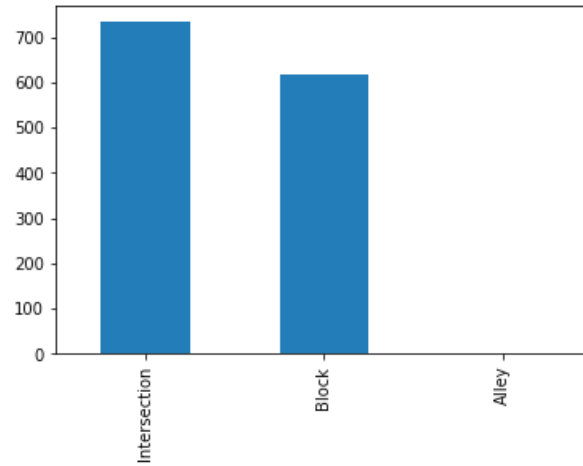


*Figure 2 - WEATHER column description*

The categorical values "Unknown" and "Other" were removed, as they do not provide useful information in the face of an initial situation for the cyclist.

— The ADDRTYPE column has 3 values with the demonstrated count and a descriptive bar graph. Figure 3 below shows these values:

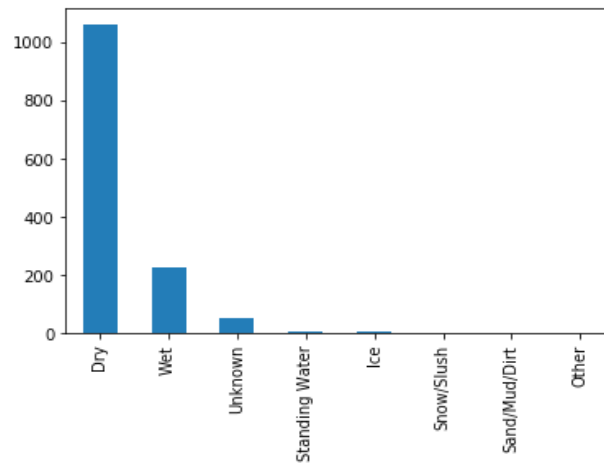
```
(<matplotlib.axes._subplots.AxesSubplot at 0x1c248095d88>,  
Intersection    733  
Block           617  
Alley           2  
Name: ADDRTYPE, dtype: int64)
```



*Figure 3 - ADDRTYPE column description*

— The ROADCOND column has 8 categorical values. The count of each category and therefore the bar graph are shown in image 4 below:

```
(<matplotlib.axes._subplots.AxesSubplot at 0x1c255b8eac8>,  
Dry          1061  
Wet          227  
Unknown      53  
Standing Water  3  
Ice          3  
Snow/Slush   2  
Sand/Mud/Dirt 2  
Other        1  
Name: ROADCOND, dtype: int64)
```

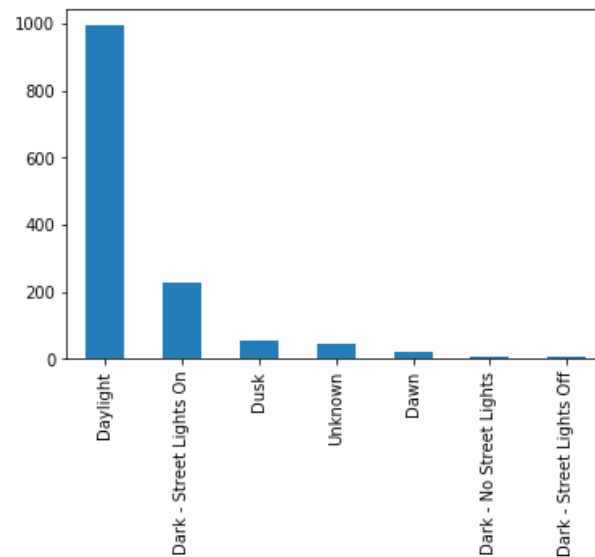


*Figure 4 - ROADCOND column description*

The “Unknown” and “Other” categories have been removed as well.

— The LIGHTCOND column has 7 categories. The count for each category and also the intrinsic bar graph are shown in image 5 below:

```
(<matplotlib.axes._subplots.AxesSubplot at 0x1c248185988>,
Daylight          994
Dark - Street Lights On    225
Dusk                53
Unknown            46
Dawn                19
Dark - No Street Lights    8
Dark - Street Lights Off   7
Name: LIGHTCOND, dtype: int64)
```



*Figure 5 - LIGHTCOND column description*

The “Unknown” column has been removed.

All 4 columns underwent the One Hot Endoding process and were concatenated and converted to values that will be submitted to the machine learning algorithm, the decision tree.

The target values in the SEVERITYCODE column have been modified. The value "1" was reset to "Property Damage" and the value "2" to "Injury".

The values of each predictor column and the target column were submitted to Train-Test-Split, to partition the training and test values, for training the decision tree and for its test, respectively. The result is seen below, in image 6:

```
from sklearn.model_selection import train_test_split

# normalized, value 0 or 1 after One Hot Encoding
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=0)

print ('Train set:\n'
      '- The x_train shape is', x_train.shape, 'and y_train shape is', y_train.shape, '\n',
      'Test set:\n'
      '- The x_test shape is', x_test.shape, 'and y_test shape is', y_test.shape)

Train set:
- The x_train shape is (946, 21) and y_train shape is (946,)
Test set:
- The x_test shape is (406, 21) and y_test shape is (406,)
```

*Figure 6 - Train-Test-Split applied*

Day — 5  
Month — October  
Year — 2020