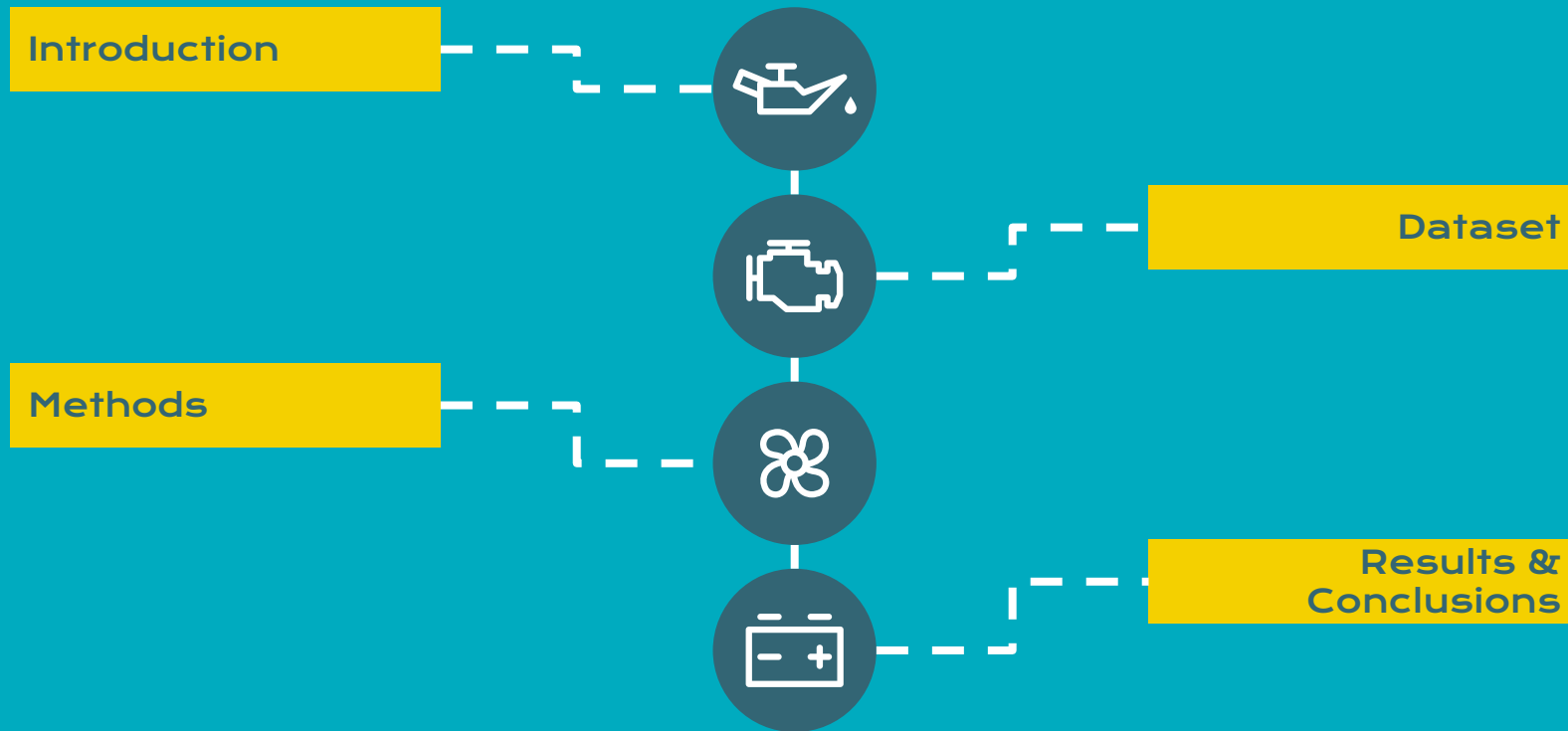


UK Car Accident Classifier

Jason Han
Dana You
Jason Dai
Abdumijit Dolkun



Table of Contents



Project Overview

Goal: use of supervised learning to predict and prevent traffic accidents based on certain patterns identified in the environment, driver, and vehicle

Why Machine Learning?

At a glance, our guiding process:

- Search for dataset
- Clean each subset of data and select columns
- Create a combined dataframe
- Feature engineering
- Experiment with different classifier methods
- Optimization



All About the Dataset

Background

- Kaggle:
<https://www.kaggle.com/datasets/benoit72/uk-accidents-10-years-history-with-many-variables>
- "UK Accidents over a 10 year period with multiple variables"
 - From 2005-2014
- Stats relate only to personal injury accidents that are:
 - 1. On public roads
 - 2. Reported to the British police
 - 3. Recorded using the STATS19 accident reporting form



Structure

- DF #1: Accidents
 - Details about circumstances of accidents from
- DF #2: Vehicles
 - Details about the vehicle(s) involved in recorded accidents
- DF #3: Casualties
 - Details about those injured or harmed in the collision
- xls File: Lookup Tables
 - Information about each variable in the 3 dataframes

Lookup Tables

Feature Examples



Sex of Driver



**Weather
Conditions**



**Casualty
Severity**

code	label	
1	Male	
2	Female	
3	Not known	
-1	Data missing or out of range	



code	label	
1	Fatal	
2	Serious	
3	Slight	

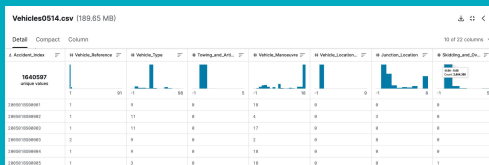
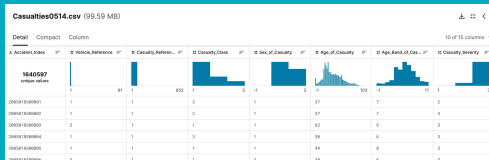
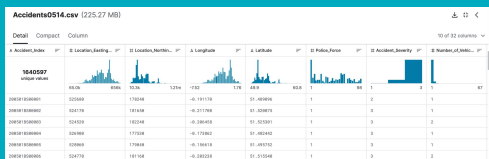


code	label	
1	Fine no high winds	
2	Raining no high winds	
3	Snowing no high winds	
4	Fine + high winds	
5	Raining + high winds	
6	Snowing + high winds	
7	Fog or mist	
8	Other	
9	Unknown	
-1	Data missing or out of range	



Process of Merging

STEP 1



Three Separate
Dataframes

Accident, Casualty, Vehicle

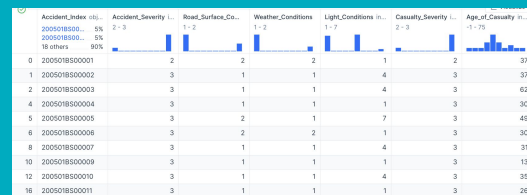
STEP 2

```
merged_df = pd.merge(df_accident_reduced, df_casualty_reduced, on='Accident_Index', how='inner')  
merged_df = pd.merge(merged_df, df_vehicle_reduced, on='Accident_Index', how='inner')
```



Merge Function

STEP 3



One Combined
Dataframe

Duplicate indices removed

Key Stakeholders

01

LAW ENFORCEMENT

Identify accident-prone areas
for preventative action

02

INSURANCE COMPANIES

Assess risk more accurately
to incentivize safer driving

03

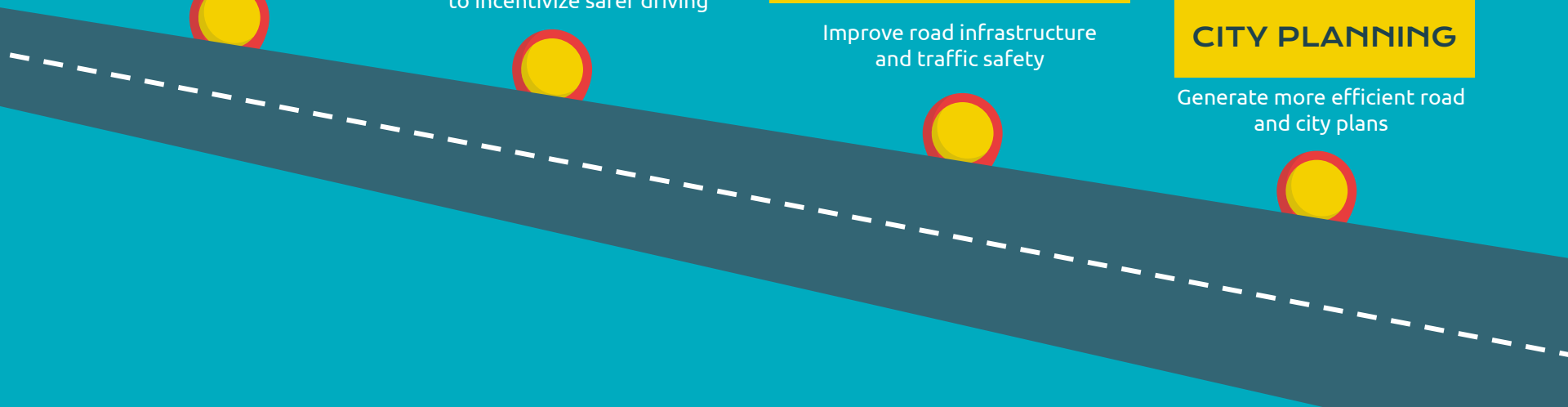
TRANSPORTATION SAFETY

Improve road infrastructure
and traffic safety

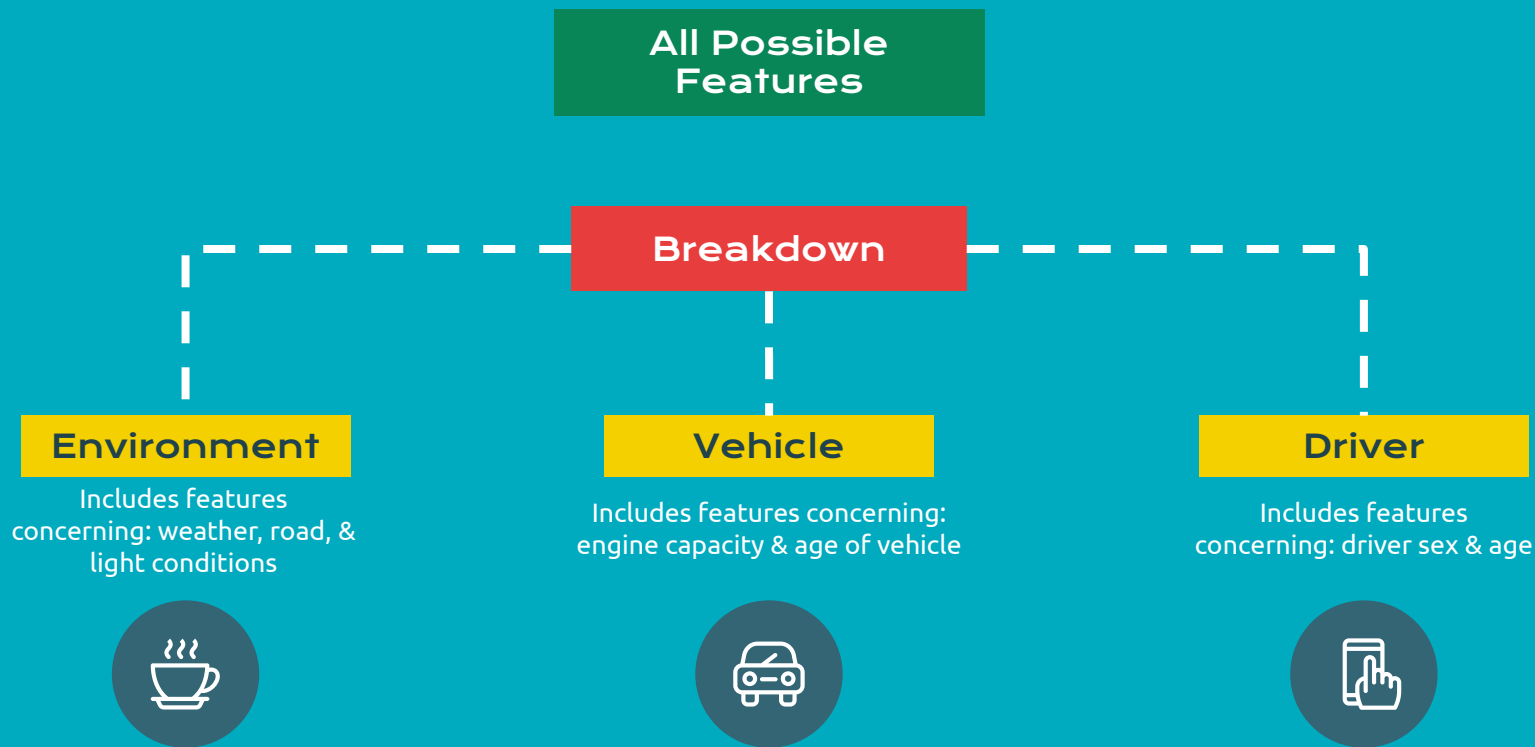
04

CITY PLANNING

Generate more efficient road
and city plans

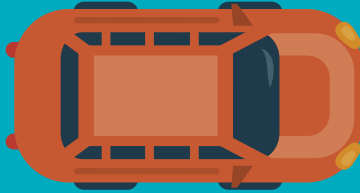


Identifying Categories and Features

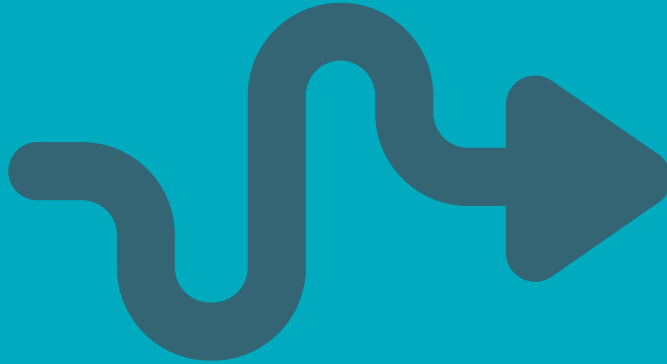


Feature Engineering in Road Accident Analysis

STEP 1



STEP 2



STEP 3



Objective

Enhancing data quality
and predictive power.

Process

Cleaning and transforming
raw data into meaningful
features

Significance

Improving model accuracy
and interpretability.

Cleaning and Encoding Categorical Variables



01

Focused on three key features:

Road_Surface_Conditions,
Weather_Conditions,
Light_Conditions.



02

Invalid entries removal

(-1 and 9) considered as data
entry errors.



03

One-hot encoding

Applied one-hot encoding to
transform these categorical variables
into a format suitable for modeling.



03

Summary

Preserves the informational value of
categorical variables while making them
usable for the ML models.

Code

```
1 # Road surface conditions cleaning
2 rd_null = merged_df['Road_Surface_Conditions'].isnull().sum()
3 print(rd_null)
4
5 cleaned_df = merged_df.copy()
6 n = -1
7 cleaned_df = cleaned_df[cleaned_df['Road_Surface_Conditions'] != n]
8
9 # Weather conditions cleaning (removing -1 and 9)
10 w_null = merged_df['Weather_Conditions'].isnull().sum()
11 print(w_null)
12 m = 9
13 cleaned_df = cleaned_df[cleaned_df['Weather_Conditions'] != m]
14 cleaned_df = cleaned_df[cleaned_df['Weather_Conditions'] != n]
15
16 # Light conditions cleaning
17 l_null = merged_df['Light_Conditions'].isnull().sum()
18 print(l_null)
19 cleaned_df = cleaned_df[cleaned_df['Light_Conditions'] != n]
20
21 # All 3 are categorical and they have no nulls or -1
22 # One hot encoding for all
23 cleaned_df = pd.get_dummies(cleaned_df, columns=['Road_Surface_Conditions', 'Weather_Conditions', 'Light_Conditions'])
24 cleaned_df
```

Feature Transformation - Age and Sex

'Age_of_Driver'

- Converted into age groups to capture generational effects.
- Removal of -1 values to maintain data integrity.



'Sex_of_Driver'

- Encoded as binary variables to simplify the model's understanding.
- Removal of (-1 and 3) values to maintain data integrity.



- **Grouping Ages:** Categorizing ages reveals how different age groups exhibit distinct driving behaviors and accident risks.
- **Encoding Gender:** Binary gender encoding aids in analyzing the impact of gender on driving habits and accident rates.



Code

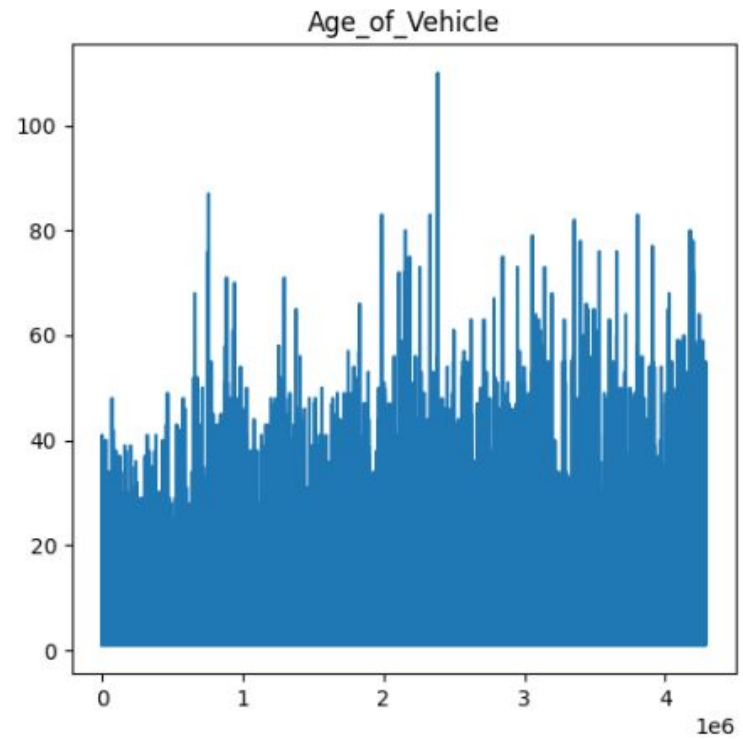
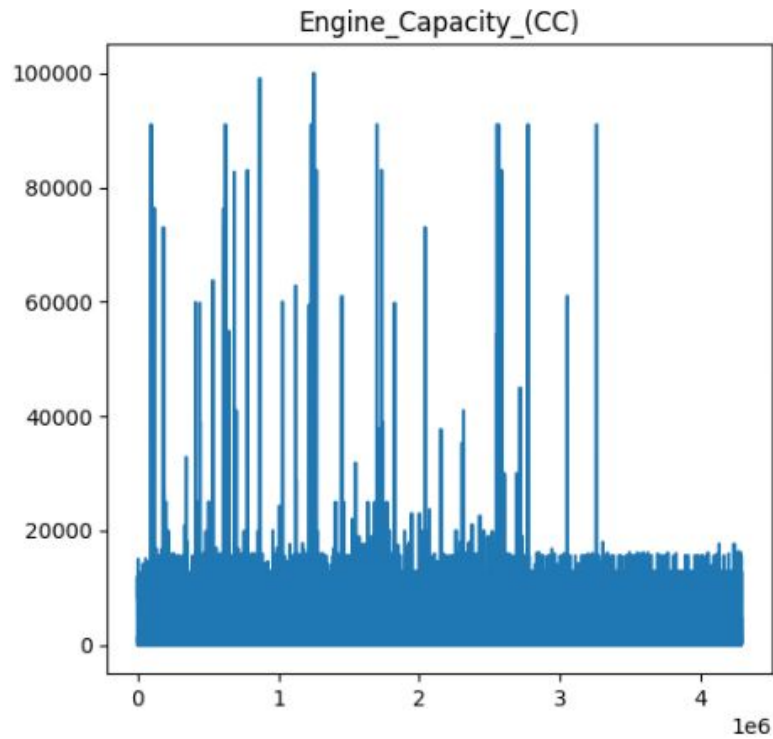
```
1 # Age_of_Driver Cleaning
2 # Check for -1 values and then remove them
3 age_drv_null = cleaned_df['Age_of_Driver'].isnull().sum()
4 age_drv_minus_one = (cleaned_df['Age_of_Driver'] == -1).sum()
5 print('Age_of_Driver - Null:', age_drv_null)
6 print('Age_of_Driver - -1:', age_drv_minus_one)
7
8 cleaned_df = cleaned_df[cleaned_df['Age_of_Driver'] != -1]
9
10 #one-hot encoding age_of_driver
11 bins=[0, 12, 18, 35, 50, float('inf')]
12 labels = ['1', '2', '3', '4', '5']
13 cleaned_df['Driver_Age_Group'] = pd.cut(cleaned_df['Age_of_Driver'], bins=bins, labels=labels, include_lowest=True, r
14
15 # Sex_of_Driver Cleaning
16 # Check for -1 values and then remove them
17 sex_drv_null = cleaned_df['Sex_of_Driver'].isnull().sum()
18 sex_drv_minus_one = (cleaned_df['Sex_of_Driver'] == -1).sum()
19 cleaned_df = cleaned_df[cleaned_df['Sex_of_Driver'] != 3]
20 print('Sex_of_Driver - Null:', sex_drv_null)
21 print('Sex_of_Driver - -1:', sex_drv_minus_one)
22
23 cleaned_df = cleaned_df[cleaned_df['Sex_of_Driver'] != -1]
24
25 # Assuming 'Sex_of_Driver' is a categorical variable that should be one-hot encoded
26 cleaned_df = pd.get_dummies(cleaned_df, columns=['Sex_of_Driver'])
27 cleaned_df = cleaned_df.drop(columns=['Age_of_Driver'])
28 cleaned_df
```

Handling Vehicle and Casualty Features



Binning and Encoding Vehicle & Casualty Data		
01	'Engine_Capacity_(CC)' and 'Age_of_Vehicle'	Binned into categories based on capacity and age for better model interpretation.
02	'Age_of_Casualty'	Grouped into age brackets, similar to 'Age_of_Driver'
03	'Sex_of_Casualty'	One-hot encoded to align with other sex-related features.
04	Summary	Binning continuous variables like engine capacity and vehicle age simplifies analysis by grouping data into interpretable categories, enhancing insights into trends and patterns.

Visual

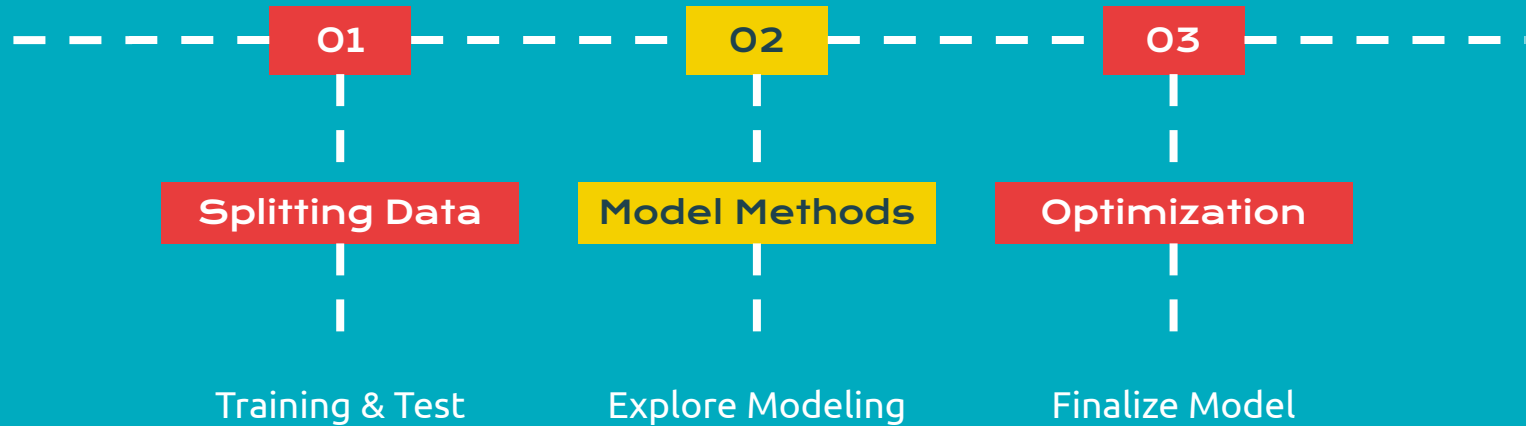


Correlation Graph



Data Modeling

Steps of Data Modeling



Split & Start

```
1 # Splitting dataset into train and test split
2 X = cleaned_df.drop(['Accident_Severity', 'Casualty_Severity'], axis=1)
3 y = cleaned_df['Accident_Severity']
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=42 )
```



K-Nearest-Neighbors

- Instance base
- Handles numerical & categorical
- Classifies based on similarity
- CV: 82.84%
- Test: 83.16%

```
1 k_fold = KFold(n_splits=10, shuffle=True, random_state=42)
2
3 knn_pipeline = Pipeline([
4     ('scaler', StandardScaler()),
5     ('knn', KNeighborsClassifier(n_neighbors=5))
6 ])
7 knn_pipeline.fit(X_train, y_train)
8
9 knn_cv_scores = cross_val_score(knn_pipeline, X_train, y_train, cv=k_fold, scoring='accuracy')
10 avg_cv_accuracy = np.mean(knn_cv_scores)
11
12 y_pred = knn_pipeline.predict(X_test)
13 test_accuracy = accuracy_score(y_test, y_pred)
14
15 print(f"Average KNN Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")
16 print(f"KNN Test Accuracy: {test_accuracy * 100:.2f}%")
```



Average KNN Cross-Validation Accuracy: 82.84%
KNN Test Accuracy: 83.16%

Decision Tree & Random Forest

```
1 decision_tree = DecisionTreeClassifier(random_state=42)
2 decision_tree.fit(X_train, y_train)
3
4 dt_cv_scores = cross_val_score(decision_tree, X_train, y_train, cv=k_fold, scoring='accuracy')
5 avg_cv_accuracy = np.mean(dt_cv_scores)
6
7 y_pred = decision_tree.predict(X_test)
8 test_accuracy = accuracy_score(y_test, y_pred)
9
10 print(f"Average Decision Tree Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")
11 print(f"Decision Tree Test Accuracy: {test_accuracy * 100:.2f}%")
```



Average Decision Tree Cross-Validation Accuracy: 84.48%

Decision Tree Test Accuracy: 84.42%

- Easy to interpret
- Handles high-dimensional data
- Model complex interactions between features
- CV: 84.48% & 84.53%
- Test: 84.42% & 84.45%



```
1 rforest = RandomForestClassifier(n_estimators=100, random_state=42)
2 rforest.fit(X_train, y_train)
3
4 rf_cv_scores = cross_val_score(rforest, X_train, y_train, cv=k_fold, scoring='accuracy')
5 avg_cv_accuracy = np.mean(rf_cv_scores)
6
7 y_pred = rforest.predict(X_test)
8 test_accuracy = accuracy_score(y_test, y_pred)
9
10 print(f"Average Random Forest Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")
11 print(f"Random Forest Test Accuracy: {test_accuracy * 100:.2f}%")
```



Average Random Forest Cross-Validation Accuracy: 84.53%

Random Forest Test Accuracy: 84.45%

Logistic Regression



- Good for binary & multiclass
- Handles numerical & categorical features
- Easy Interpretation
- CV: 84.56%
- Test: 84.49%

Logistic Regression

```
1 scaler = StandardScaler()
2 X_train_scaled = scaler.fit_transform(X_train)
3 X_test_scaled = scaler.fit_transform(X_test)
4
5 lreg = LogisticRegression(max_iter=1000)
6 lreg.fit(X_train_scaled, y_train)
7
8 lreg_cv_scores = cross_val_score(lreg, X_train_scaled, y_train, cv=k_fold, scoring='accuracy')
9 avg_cv_accuracy = np.mean(lreg_cv_scores)
10
11 y_pred = lreg.predict(X_test_scaled)
12 test_accuracy = accuracy_score(y_test, y_pred)
13
14 print(f"Average Logistic Regression Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")
15 print(f"Logistic Regression Test Accuracy: {test_accuracy * 100:.2f}%")
```

Average Logistic Regression Cross-Validation Accuracy: 84.56%
Logistic Regression Test Accuracy: 84.49%

Neural Network



- Captures Complex Patterns
- Model non-linear relationships
- Needs preprocessing
- Tuning hyperparameters
- CV: 84.56%
- Test: 84.49%



```
1 scaler = StandardScaler()
2 X_train_scaled = scaler.fit_transform(X_train)
3 X_test_scaled = scaler.fit_transform(X_test)
4
5 nn = MLPClassifier(
6     hidden_layer_sizes=(30),
7     activation='logistic',
8     solver='lbfgs',
9     random_state=42
10 )
11 nn.fit(X_train_scaled, y_train)
12
13 nn_cv_scores = cross_val_score(nn, X_train_scaled, y_train, cv=k_fold, scoring='accuracy', error_score='raise')
14 avg_cv_accuracy = np.mean(nn_cv_scores)
15
16 y_pred = nn.predict(X_test_scaled)
17 test_accuracy = accuracy_score(y_test, y_pred)
18
19 print(f"Average Neural Networks Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")
20 print(f"Neural Networks Test Accuracy: {test_accuracy * 100:.2f}%")
```



Average Neural Networks Cross-Validation Accuracy: 84.56%

Neural Networks Test Accuracy: 84.49%

Ensemble - Best of both worlds

Log-Reg
Random Forest



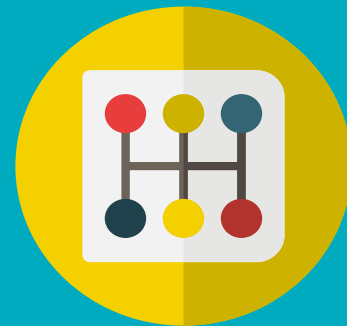
CV: 84.53%
Test: 84.49%

```
1 ensemble = VotingClassifier(  
2     estimators=[  
3         ('lreg', lreg),  
4         ('rf', rforest)  
5     ],  
6     voting='hard')  
7 ensemble.fit(X_train, y_train)  
8  
9 ensemble_cv_scores = cross_val_score(ensemble, X_train, y_train, cv=k_fold, scoring='accuracy')  
10 avg_cv_accuracy = np.mean(ensemble_cv_scores)  
11  
12 y_pred = nn.predict(X_test)  
13 test_accuracy = accuracy_score(y_test, y_pred)  
14  
15 print(f"Average Ensemble Cross-Validation Accuracy: {avg_cv_accuracy * 100:.2f}%")  
16 print(f"Ensemble Test Accuracy: {test_accuracy * 100:.2f}%")
```



Average Ensemble Cross-Validation Accuracy: 84.53%

Ensemble Test Accuracy: 84.49%



Data Results

A

Logistic Regression

Displayed an accuracy of
84.56%

B

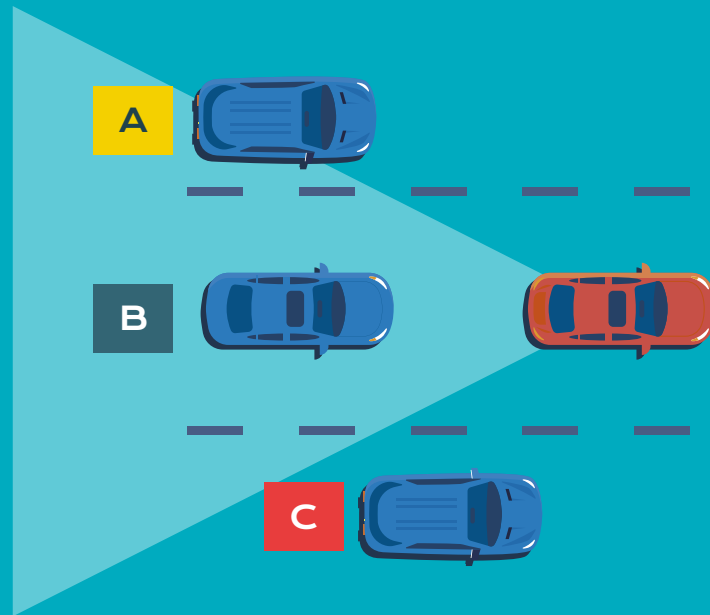
K-Nearest Neighbors

Displayed an accuracy of
83.75%

C

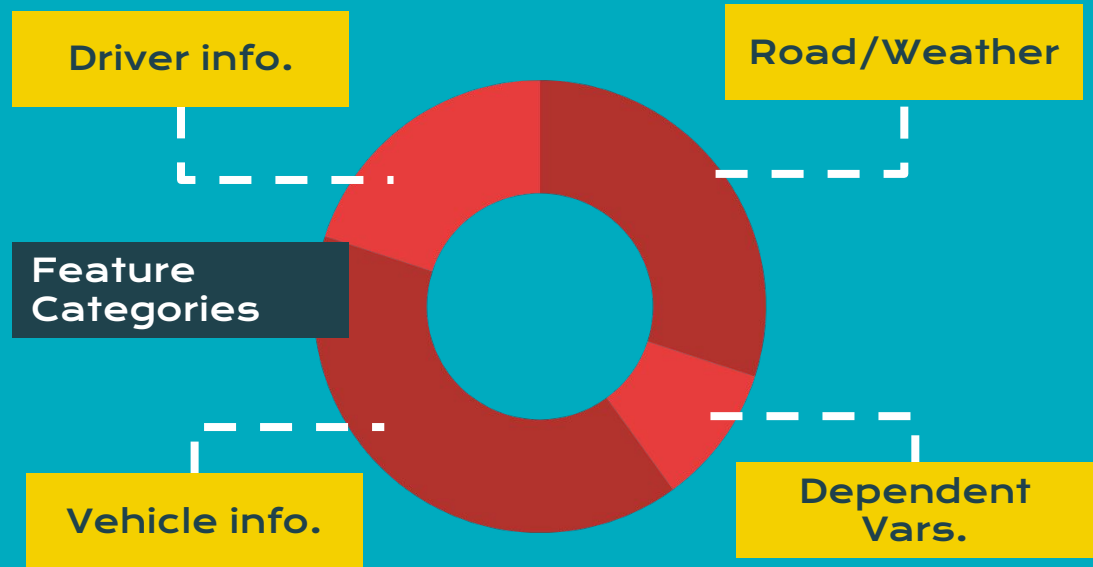
Ensemble Methods

Displayed an accuracy of
84.52%



Data Results (cont.)

- By selecting specific categories of features, we aimed to narrow down main impacting factors
- We noticed that the light condition feature had the highest feature coefficient when examining by significance
- Combining the coefficients in their categories, the Driver info. Category had the highest significance



Conclusions

THEORY

Driver information likely indicates a greater correlation with accidents over other categories

IMPLICATIONS

It's likely that a combination of the weather and driver have the greatest correlation

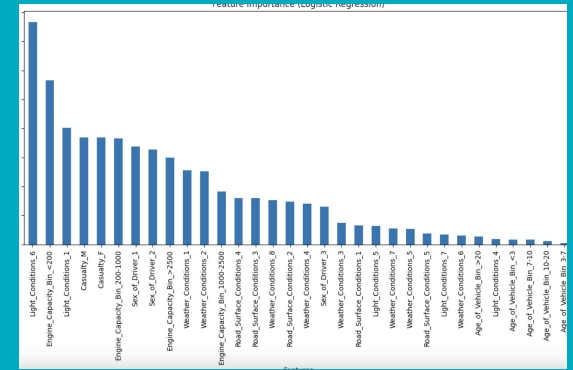


DATA

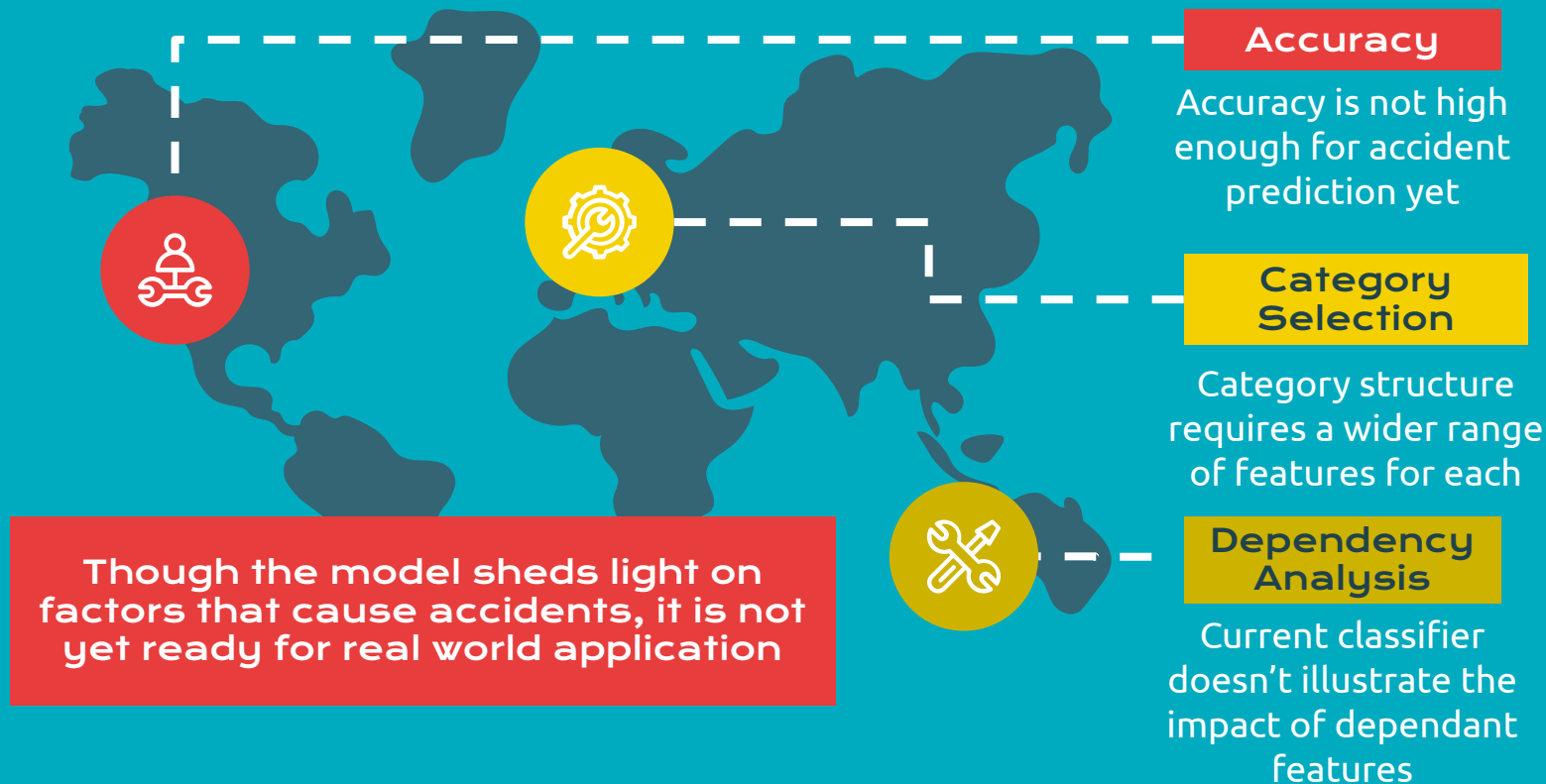
Though light conditions has the highest single influence, driver info has an overall higher combined coefficient

NOTE





There are still many ways to calculate significance and impact



Real World Application and Readiness



Possible improvements and recommendations

			
Data	Classifier	Cleaning	Analysis
Selective but denser feature categories	Use of model and combination	Handle feature specific nulls more accurately	Incorporate further metrics to analyze influence
More intensive dataset	Implementation of further optimizations	Incorporate error analysis	Contextualize dependent features

Thank you!

Any
questions?