



California Traffic Collision: Predictive Analysis

Team 4: 1. Chan Chakrya Menh 2. Priyanka Kumari
3. Riyan Rattan

I. Introduction

Background:

What is Traffic Collision?

A traffic collision or crash occurs when a vehicle collides with another vehicle, pedestrian, road barrier, or a stationary obstacle such as a tree or a utility pole. It may result in injury, death, vehicle damage, possession damage which causes death and disability, and financial burden.

The National Highway Traffic Safety Administration (NHTSA) disclosed its early estimation of traffic fatalities for 2021. NHTSA projects an estimated 42,915 individuals died in motor vehicle traffic crashes last year, a 10.5% expansion from the 38,824 fatalities in 2020. The projection is the highest fatalities since 2005 and the most significant annual percentage increase in the Fatality Analysis Reporting System history.

Motive/Goals:

- This project is to study California traffic collision by using dataset in 2019 to do predictive analysis. This prediction will be achieved by utilizing Orange for machine learning to develop a model with training dataset and eventually with the test dataset.

II. Research Question

- How can we use machine learning to detect the type of collision?
- What type of collision that has more fatalities and injury?
- Which month people was kill and get injure the most from traffic collision?
- What kind of weather has more fatalities and injury?

III. Data Preparation

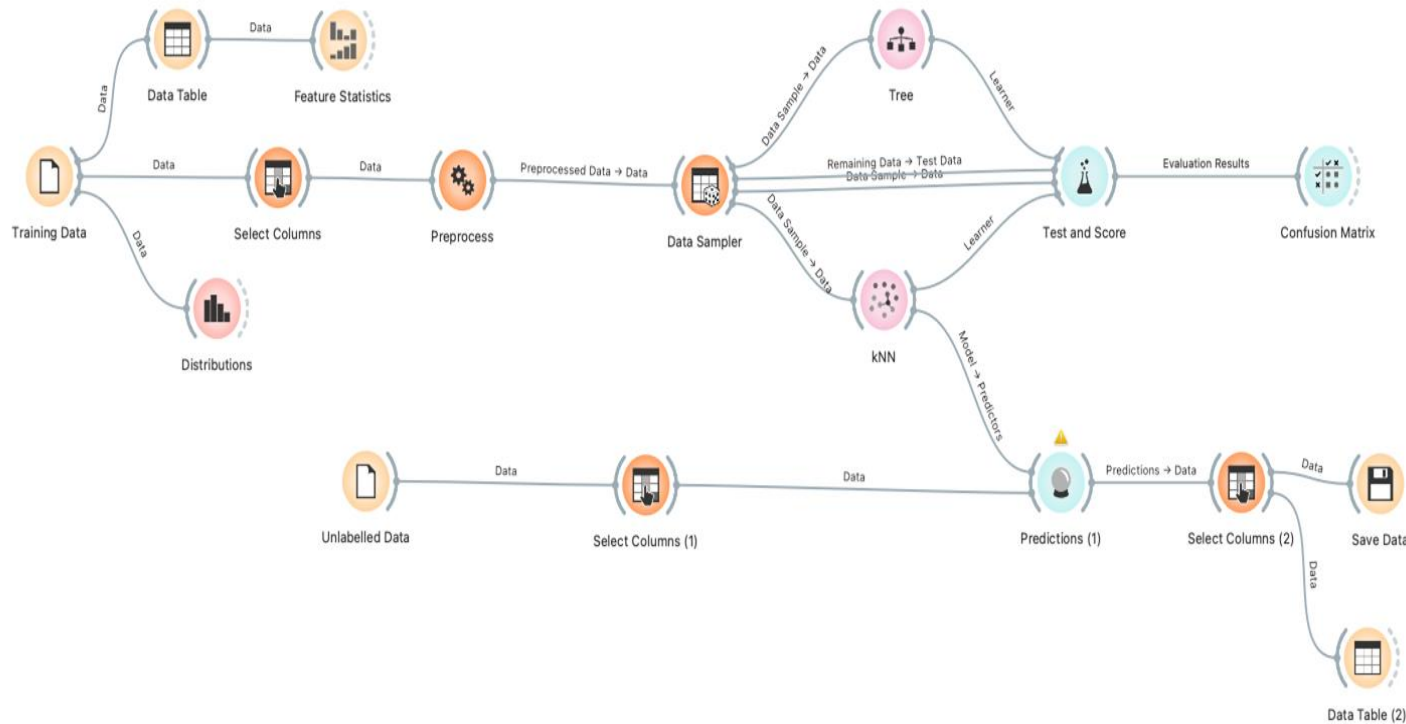
- Cleaning the dataset in csv format
 - Remove irrelevant columns, which eliminated missing values
 - Limit the sample size up to 100K, for faster processing for sake of learning
- Import Data to Orange to check data attributes, statistics, and distribution
- Set Type of Collision as the target in the train data
- Run test score, prediction, and confusion metric
- Use Tableau for data visualization
- Set up second pipeline for test data by removing weather and type of collision to explore possibilities as gender, type of vehicles, etc. in the finalized data model

IV. Data Analysis

- Create training and test data set
- Target response to traffic collision is quantitative method
- Better data models for training the data set:
 - kNN
 - Decision Tree

Orange:

Create two pipeline for
train and test data.



Info

99999 instances (no missing data)
8 features
Target with 9 values
No meta attributes

Variables

- ☐ Show variable labels (if present)
- ☐ Visualize numeric values
- ☒ Color by instance classes

Selection

- ☒ Select full rows

V. Data visualization : Tableau

Fig1:Number Killed Vs Type of Collision.

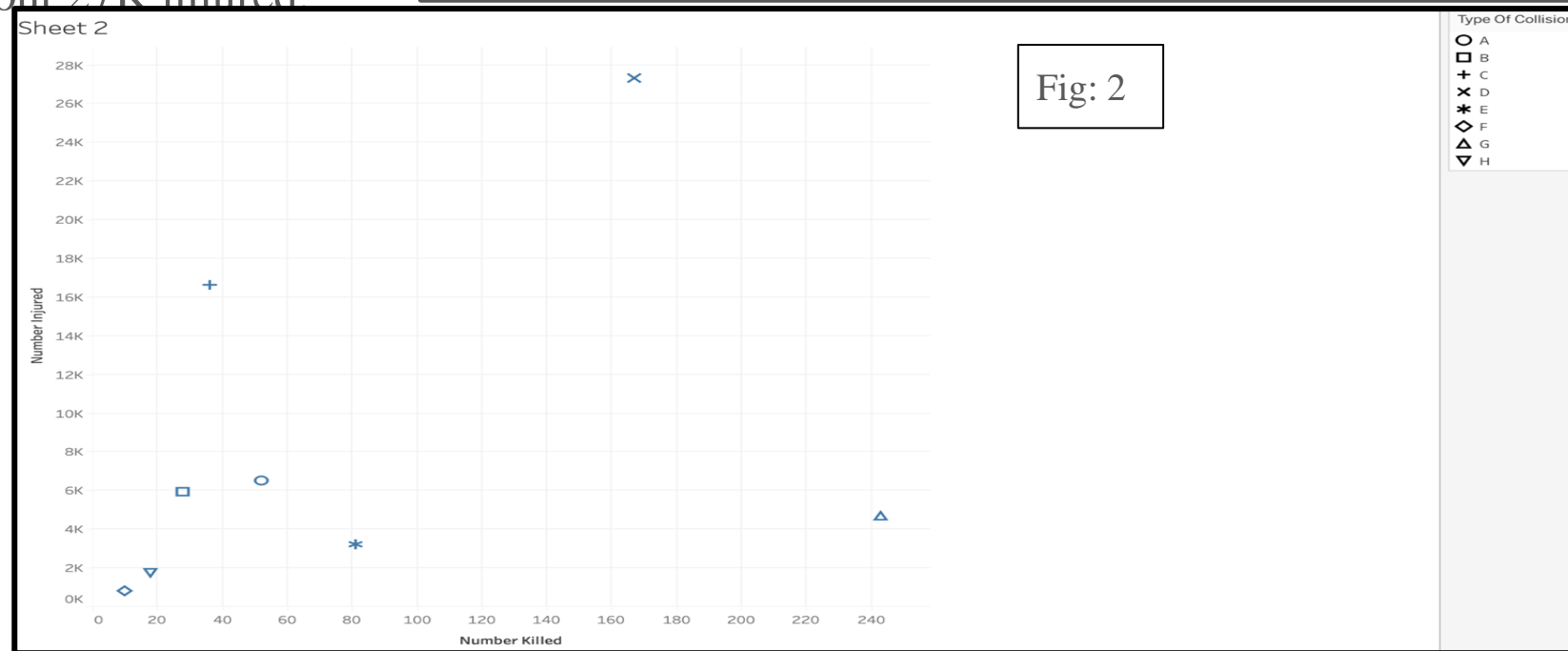
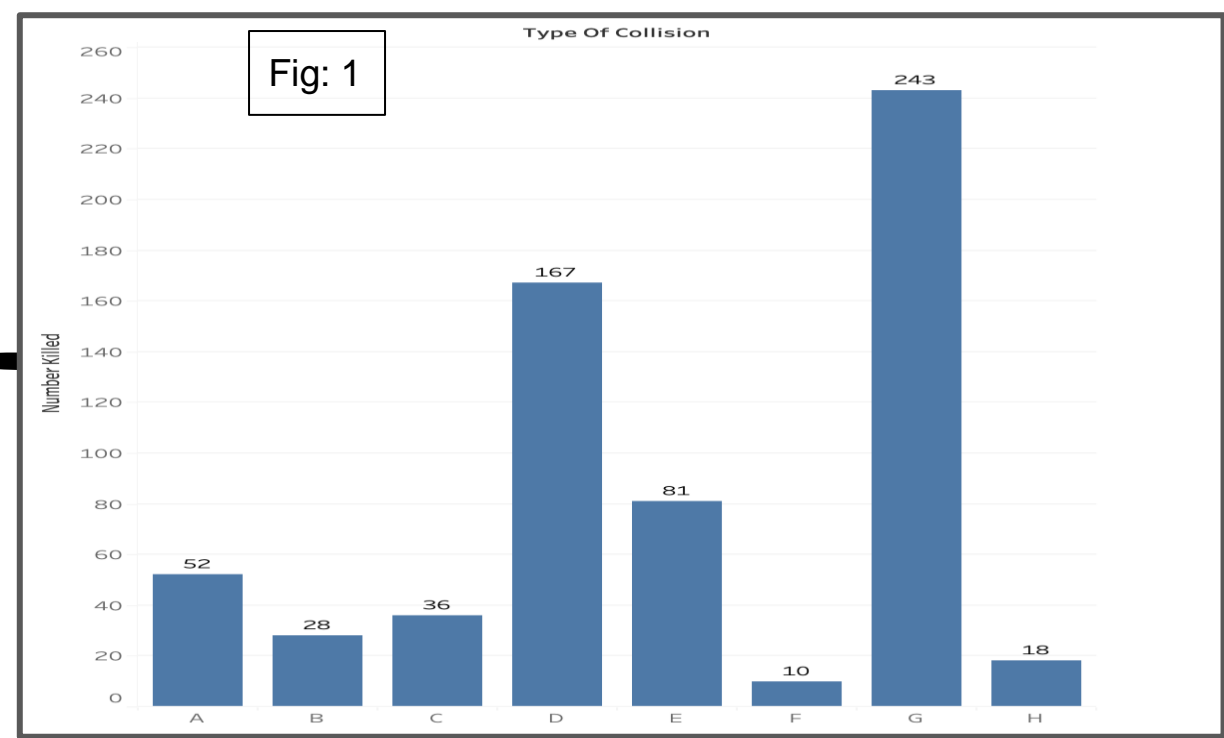
G type has most fatalities, and second is D type.

Fig2: Number Killed Vs Number of Injured,

D type has about 167 fatalities and about 27K injured

Type of Collision

- A - Head-On
- B - Sideswipe
- C - Rear End
- D - Broadside
- E - Hit Object
- F - Overturned
- G - Vehicle/Pedestrian
- H - Other



Weather Vs Number Killed

This graph shows that most people died when the Weather is clear.

Weather 1:

A - Clear

B - Cloudy

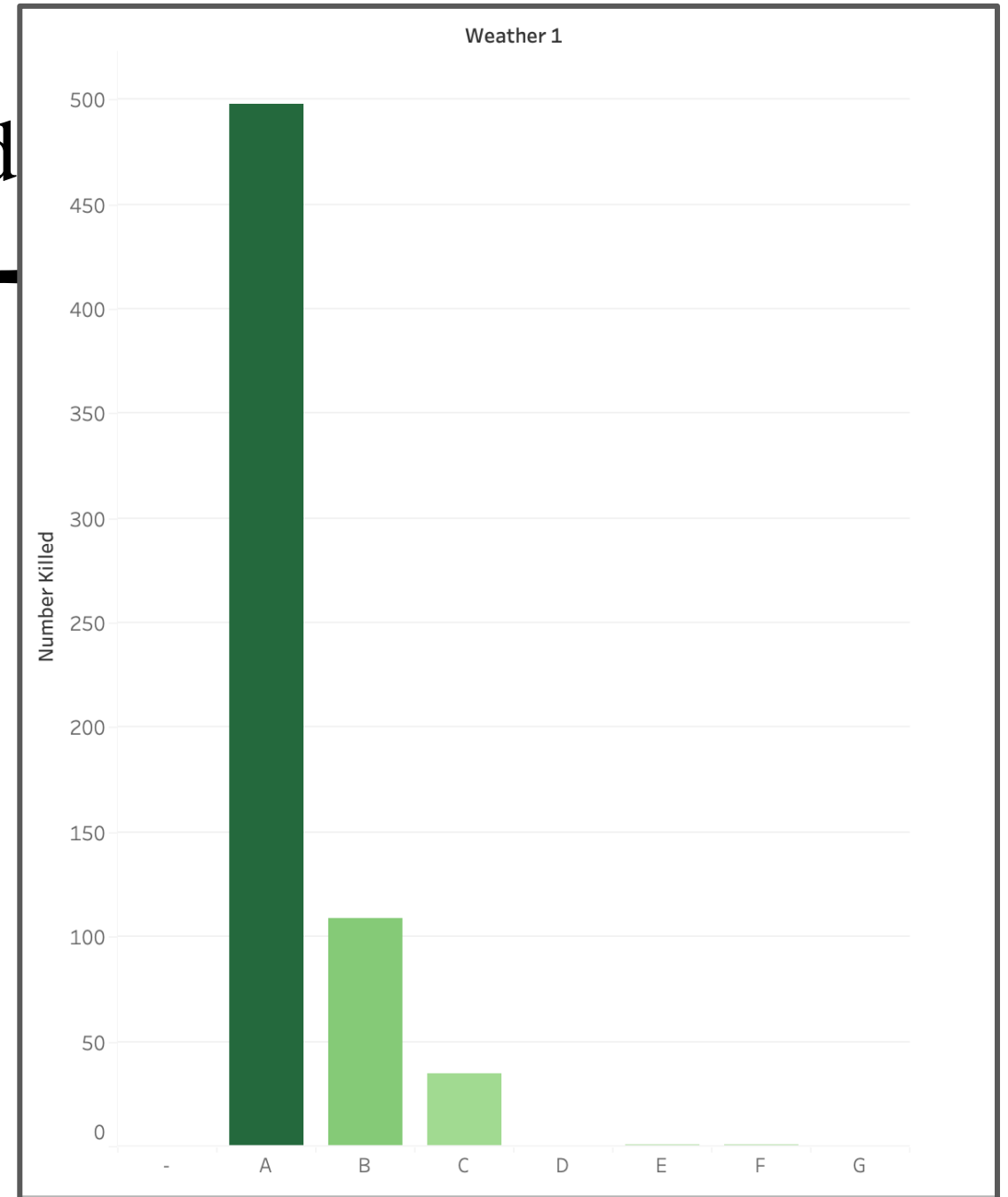
C - Raining

D - Snowing

E - Fog

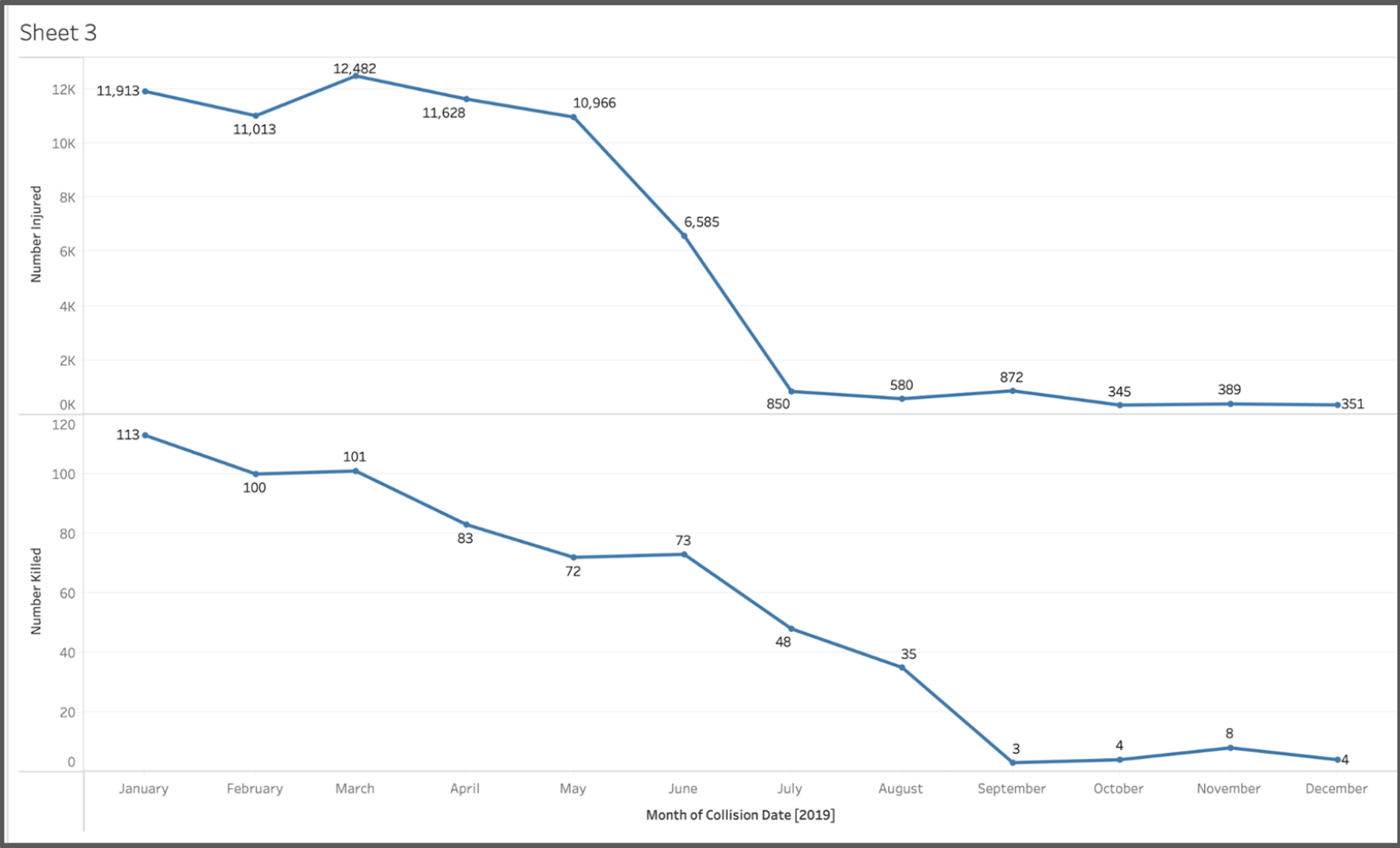
F - Other

G - Wind



Number Killed Vs Number Injured by Month

In comparison, number of injured on monthly bases is more than number of killed.



VI. Model Evaluation (Train Dataset)

Confusion Matrix

Clicking on cells or in headers outputs the corresponding data instances Ok, got it

Show: Number of instances

	Predicted									
	-	A	B	C	D	E	F	G	H	Σ
-	0	0	194	228	542	0	0	29	0	993
A	0	13	1072	1065	2637	0	0	146	0	4933
B	0	7	5326	4905	3063	0	0	125	0	13426
C	0	14	4538	5248	7371	2	0	163	0	17336
D	0	45	3325	4303	11601	5	1	400	0	19680
E	0	4	2518	3014	1666	0	0	212	1	7415
F	0	1	77	99	393	0	0	46	0	616
G	0	3	89	115	2658	0	0	596	0	3461
H	0	1	616	411	989	1	0	122	0	2140
Σ	0	88	17755	19388	30920	8	1	1839	1	70000

Actual

☐ Predictions ☒ Probabilities

☒ Apply Automatically

Select Correct Select Misclassified Clear Selection

Test and Score

Cross validation

Number of folds: 5

☒ Stratified

Cross validation by feature

Random sampling

Repeat train/test: 10

Training set size: 70 %

☒ Stratified

☐ Leave one out

☐ Test on train data

☐ Test on test data

Evaluation results for target (None, show average over classes)

Model	Train time [s]	Test time [s]	AUC	CA	F1	Precision	Recall	Specificity
Tree	0.552	0.009	0.641	0.330	0.280	0.251	0.330	0.778
kNN	0.217	9.542	0.630	0.325	0.277	0.256	0.325	0.782

Compare models by: Area under ROC curve

Negligible diff.: 0.1

	Tree	kNN
Tree		0.999
kNN	0.001	

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

VI. Project Action Plan

No.	Description	Start Date	End Date	Status
1.	Data Cleaning	July 17	July 24	15%
2.	Data Analyzing	July 25	July 31	45%
3.	Data Modeling	Aug 1	Aug 7	65%
4.	Data Visualization	Aug 8	Aug 14	75%
5.	Data Report Submission	Aug 15	Aug 24	100%

VII. Future Plan

- Explore more on data and features from data set, like gender and age to the test data
- Add more visualization
- Working on Modeling and improve accuracy
- Make final slide presentation with all findings