

Agenda

- Introduction
- Data cleaning
- Features
- Data exploration
- Feature selections
- Machine learning Models
- Comparisons
- Conclusions
- Future work
- Reference

Introduction

In the research, we want to do some data exploration to answer a question by using machine learning models.

The data source is Kaggle and the format is SQL. This data comes from the California Highway Patrol and covers collisions from January 1st, 2001 until mid October 2020.

Introduction

There are three main tables:

- □ Collisions: Contains information about the collision, where it happened, what vehicles were involved.
- Parties: Contains information about the groups people involved in the collision including age, sex, and sobriety.
- □ Victims: Contains information about the injuries of specific people involved in the collision.

In this project, the three main tables joined where collision date is after 2020-01-01

Data Cleaning

- NAN values: Columns which have more than 50% unknown value have been dropped.
- Outliers: Rows with parties' age of zero have been dropped (32506 out of 351975)

Features

There are 86 columns including categorical and numeric features such Party Number, Victim Sex, Victim Age, Party Race, Vehicle Year, Vehicle Make, Party Age, Party Sex, At Fault

party_sobriety	A	В	C	D	G	Н
party_drug_physical						
E	89	473	14	4	0	0
F	179	7	7	2	0	0
G	0	0	0	0	21596	0
Н	0	0	0	0	0	18630
I	1144	56	39	2	0	0
NA	278607	14686	3087	1957	0	0

Data exploration

After considering the table of parties' use of drugs and the parties' sobriety table, is has been notified that most the related unknown on parties' use of drugs are related to sober parties

FEATURE SELECTIONS

Columns with correlation more than 0.25 has been chosen.

Collision_date Collision		
Mark		100 100 100 100 100 100 100 100 100 100
road_condition_1		1000000
party_count weather_1		
weather_1 0.033		Territoria
primary_collision_factor	· · · · ·	
direction	-	
pcf_violation latitude -		120 (0.20)
latitude 0.064 bicyclist_killed_count 0.068 bicyclist_injured_count 0.074 party_drug_physical 0.075 party_face 0.078 party_face 0.079 chp_vehicle_type_af_fault 0.08 victim_age 0.082 cellphone_use 0.083 chp_vehicle_type_towing 0.085 control_device 0.099 control_device 0.099 control_device 0.099 control_device 0.01 dther_associate_factor_1 0.1 state_highway_indicator 0.11 distance 0.11 distance 0.11 victim_sex 0.11 collision_date 0.11 financial_responsibility 0.12 pedestrian_injured_count 0.12 pedestrian_injured_count 0.14 process_date 0.14 motorcyclist_killed_count 0.14 population 0.18 pedestrian_action 0.19 collision_time 0.19 party_safety_equipment_2 0.19 vehicle_make 0.19 pedestrian_killed_count 0.19 collision_time 0.19 party_safety_equipment_2 0.19 pedestrian_killed_count 0.19 movement_preceding_collision 0.21 statewide_vehicle_type 0.21 beat_type 0.21 beat_type 0.23 bit_and_run 0.23 victim_ejected 0.25 victim_eigected		
bicyclist_injured_count bicyclist_injured_count party_drug_physical party_sex party_sex party_race comparty_race party_sex party_race party_sex party	• =	
bicyclist_injured_count party_drug_physical		CONTROL OF THE CONTRO
party_drug_physical party_sex -		
party_sex - party_sex - party_race - 0.079 chp_vehicle_type_at_fault - 0.08 victim_age - 0.082 cellphone_use - 0.083 chp_vehicle_type_towing - 0.085 control_device - 0.099 party_sobriety - 0.1 state_highway_indicator - 0.11 state_highway_indicator - 0.11 state_responsibility - 0.12 pedestrian_injured_count - 0.12 pedestrian_injured_count - 0.13 ilighting - 0.14 motorcyclist_killed_count - 0.14 process_date - 0.14 motorcyclist_injured_count - 0.14 process_date - 0.14 motorcyclist_injured_count - 0.19 pedestrian_action - 0.19 pedestrian_action - 0.19 pedestrian_action - 0.19 pedestrian_stilled_count - 0.19 party_safety_equipment_2 - 0.19 vehicle_make - 0.19 pedestrian_killed_count - 0.19 statewide_vehicle_type - 0.21 beat_type - 0.21 beat_type - 0.21 beat_type - 0.21 chp_beat_class - 0.22 statewide_vehicle_type_at_fault - 0.23 chp_beat_class - 0.22 statewide_vehicle_type_at_fault - 0.23 statewide_vehicle_type_at_fault - 0.23 victim_ejected - 0.23 victim_aeperted - 0.24 victim_aeperted - 0.23 victim_aeperted - 0.23 victi		
party_race - 0.079		1 TO
Chp_vehicle_type_at_fault victim_age	_	
victim_age	· · · · -	0.09
cellphone use cellphone use control device control device party, sobriety control device party, sobriety control device party, sobriety control device control device party, sobriety control device cont		0.23
Chip vehicle type_towing		
Ontrol device Ontrol device Party sobriety Ontrol device Party sobriety Ontrol device Party sobriety Ontrol device Ontrol de		
party_sobriety other_associate_factor_1 of the content of the cont		3411723
other_associate_factor_1 state_nighway_indicator distance victim_sex collision_date financial_responsibility pedestrian_injured_count party_type lighting motorcyclist_killed_count process_date motorcyclist_injured_count population pedestrian_action collision_time pedestrian_lime party_safety_equipment_2 vehicle_make pedestrian_killed_count movement_preceding_collision statewide_vehicle_type beat_type tow_away dnp_beat_class bit_and_run victim_ejected 0.11 0.12 0.12 0.13 0.14 0.14 0.14 0.16 0.19 0.19 0.19 0.19 0.19 0.19 0.19 0.19	_	2 Mag 2 2 M
state highway indicator		NOTE: The second
distance victim_sex victi		
victim_sex 0.11 -0.00 collision_date 0.11 -0.00 financial_responsibility 0.12 -0.00 pedestrian_injured_count 0.12 -0.12 party_type 0.13 -0.14 lighting 0.14 -0.14 motorcyclist_killed_count 0.14 -0.16 poccess_date 0.14 -0.16 motorcyclist_injured_count 0.18 -0.19 pedestrian_action 0.19 -0.19 collision_time 0.19 -0.19 party_safety_equipment_2 0.19 -0.19 pedestrian_killed_count 0.19 -0.19 movement_preceding_collision 0.21 -0.21 statewide_vehicle_type 0.21 -0.21 beat_type 0.21 -0.22 tow_away 0.22 0.2 statewide_vehicle_type_at_fault 0.23 -0.23 hit_and_run 0.23 -0.23 hit_and_run 0.23 -0.23 reference 0.23<		
Collision_date		011
financial_responsibility pedestrian_injured_count party_type	_	- 0.00
pedestrian_injured_count party_type	_	0.12
part_type		1.000000
lighting		0.13
motorcyclist_killed_count		0.14
process_date		0.14
motorcyclist_injured_ount		0.14
population	_	0.16
pedestrian_action collision_time - 0.19 collision_time - 0.19 party_safety_equipment_2 - 0.19 pedestrian_killed_count - 0.19 pedestrian_killed_count - 0.19 movement_preceding_collision - 0.21 statewide_vehicle_type - 0.21 tow_away - 0.21 chp_beat_class - 0.22 statewide_vehicle_type_at_fault - 0.23 chp_beat_type - 0.23 hit_and_run victim_ejected - 0.23		0.18
collision_time - 0.19 party_safety_equipment_2 - 0.19 vehicle_make - 0.19 pedestrian_killed_count - 0.19 movement_preceding_collision - 0.21 statewide_vehicle_type - 0.21 tow_away - 0.21 chp_beat_class - 0.22 statewide_vehicle_type_at_fault - 0.23 chp_beat_type - 0.23 hit_and_run victim_ejected - 0.23		0.19
party_safety_equipment_2 vehicle_make - 0.19 pedestrian_killed_count - 0.19 movement_preceding_collision - 0.21 statewide_vehicle_type - 0.21 beat_type - 0.21 tow_away - 0.21 chp_beat_class - 0.220. statewide_vehicle_type_at_fault - 0.23 chp_beat_type - 0.23 hit_and_run - 0.23 victim_ejected - 0.23		0.19
vehicle_make 0.19 pedestrian_killed_count 0.19 movement_preceding_collision 0.21 statewide_vehicle_type 0.21 beat_type 0.21 tow_away 0.21 chp_beat_class 0.22 statewide_vehicle_type_af_fault 0.23 chp_beat_type 0.23 hit_and_run 0.23 victim_ejected 0.23	_	0.19
movement_preceding_collision		0.19
movement_preceding_collision	pedestrian killed count -	0.19
statewide_vehicle_type		0.21
tow away - 0.21 chp_beat_class - 0.220. statewide_vehicle_type_at_fault - 0.23 chp_beat_type - 0.23 hit_and_run - 0.23 victim_ejected - 0.23	statewide_vehicle_type -	0.21
chp_beat_class = 0.22 = -0. statewide_vehicle_type_at_fault = 0.23	beat type -	0.21
statewide_vehicle_type_at_fault - 0.23 chp_beat_type - 0.23 hit_and_run - 0.23 victim_ejected - 0.23	tow_away -	0.21
chp_beat_type - 0.23 hit_and_run - 0.23 victim_ejected - 0.23	chp beat class -	
chp_beat_type - 0.23 hit_and_run - 0.23 victim_ejected - 0.23	statewide vehicle type at fault -	0.23
victim_ejected - 0.23		1 TAN 10
victim_ejected - 0.23	hit_and_run =	T-10000
motor vehicle involved with -		
	motor_vehicle_involved_with -	0.23
party_safety_equipment_1 - 0.24	party_safety_equipment_1 -	11.0000000
pcf_violation_category - 0.24	pcf_violation_category -	
victim_safety_equipment_2 - 0.25	victim_safety_equipment_2 -	1.772.6672
victim_seating_position = 0.260.	victim_seating_position -	0.26

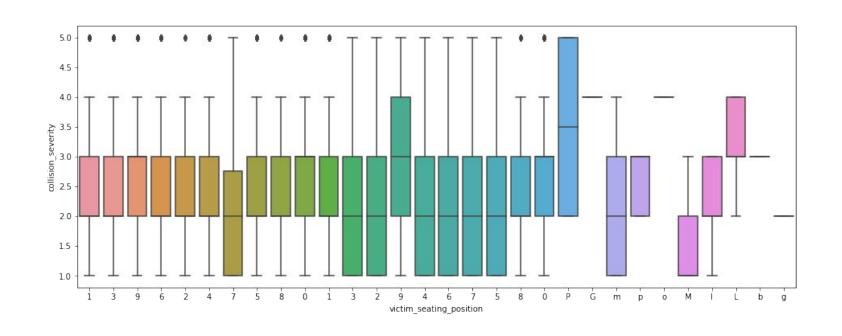
Feature selections

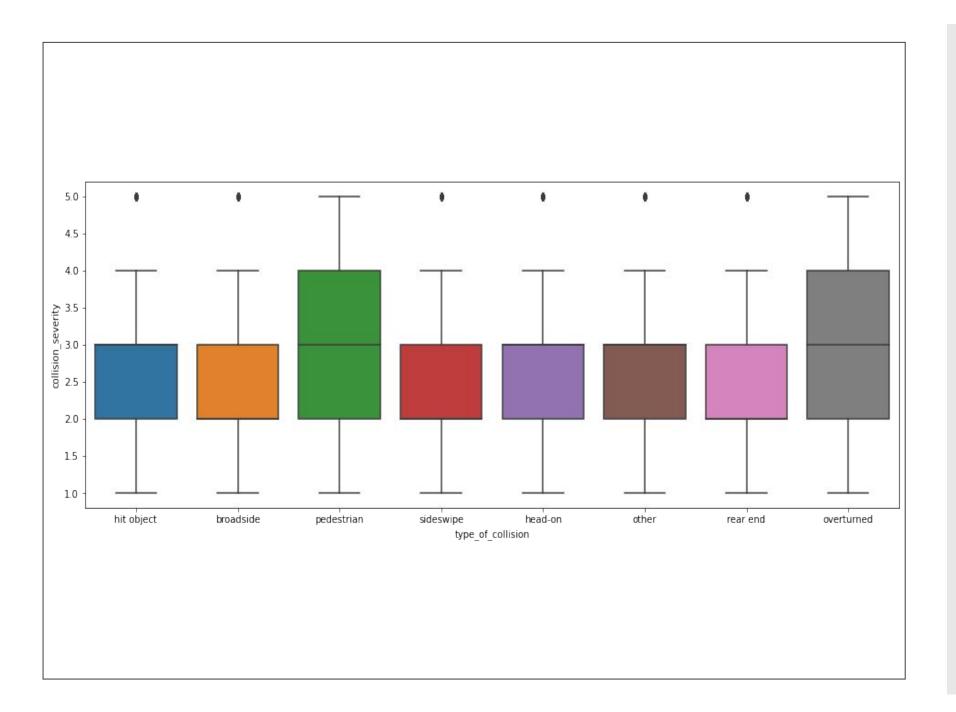
Columns with correlation more than 0.25 has been chosen.

Victim seating position, type of collision, victim role, victim safety equipment 1, primary road, secondary road columns have been chosen to predict their effects on collision severity

Feature selections Victim seating position

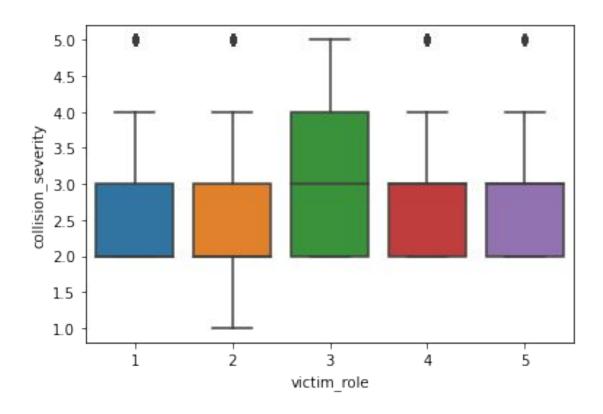
- ∘ 1 Driver
- o 2 thru 6 Passengers
- 7 Station Wagon Rear
- 8 Rear Occupant of Truck or Van
- 9 Position Unknown
- o 0 Other Occupants
- A thru Z BusOccupants
- ∘ - Not Stated





Feature selections type of collision

- A Head-On
- B Sideswipe
- 。 C Rear End
- D Broadside
- E Hit Object
- F Overturned
- G Vehicle/Pedestrian
- H Other
- - Not Stated



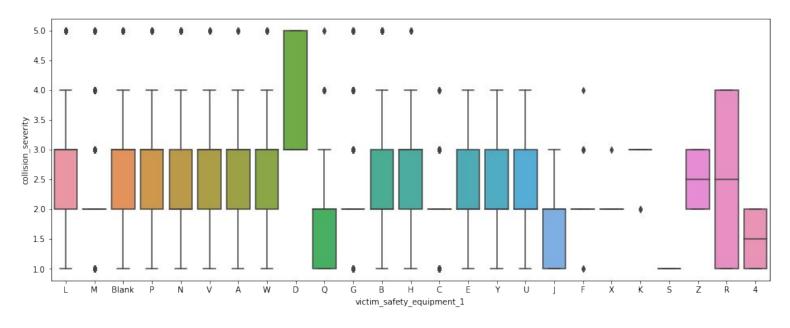
Feature selections victim role

- 1 Driver
- 2 Passenger (includes non-operator on bicycle or any victim on/in parked vehicle or multiple victims on/in non-motor vehicle)
- 3 Pedestrian
- 4 Bicyclist
- 5 Other (single victim on/in non-motor vehicle;
 e.g. ridden animal, horse-drawn carriage, train,
 or building)
- ∘ 6 Non-Injured Party

- A None in Vehicle
- B Unknown
- C Lap Belt Used
- D Lap Belt Not Used
- E Shoulder Harness Used
- F Shoulder Harness Not Used
- G Lap/Shoulder Harness Used
- H Lap/Shoulder Harness Not Used
- J Passive Restraint Used
- K Passive Restraint Not Used
- L Air Bag Deployed
- M Air Bag Not Deployed
- N Other
- P Not Required
- Q Child Restraint in Vehicle Used
- R Child Restraint in Vehicle Not Used
- S Child Restraint in Vehicle, Use Unknown
- T Child Restraint in Vehicle, Improper Use
- U No Child Restraint in Vehicle
- V Driver, Motorcycle Helmet Not Used
- W Driver, Motorcycle Helmet Used
- X Passenger, Motorcycle Helmet Not
- Y Passenger, Motorcycle Helmet Used

· - or blank - Not Stated

Feature selections victim safety equipment1



Feature selections primary road

'NORRIS RD', 'GOSFORD RD', 'FIRST ST',

..., 'LOS ABITOS BL',
'WESTAR DR',

'VENTURA RD FRONTAGE RD'

The number of collisions happened are more than 100 times.

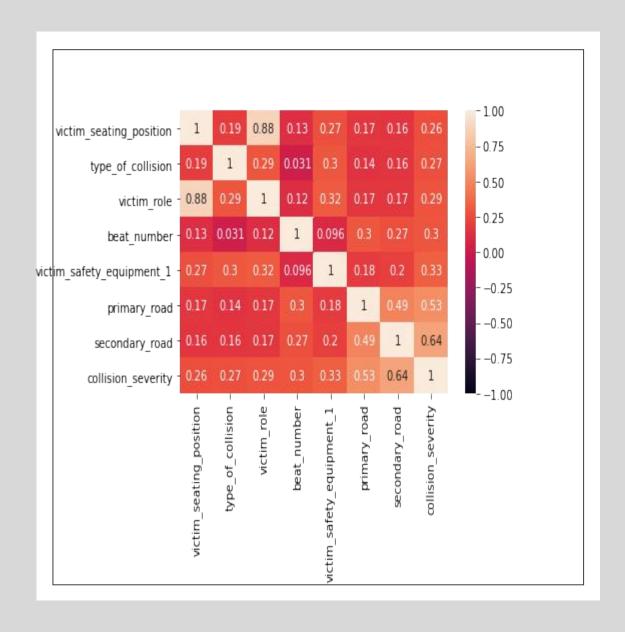
primary_road	avg_severity	count
INTERSTATE 5 SOUTHBOUND	3.421965	173
SR-39 (SAN GABRIEL CANYON RD.)	3.398230	113
US-395	3.382979	141
SR-4	3.329341	167
GENEVA AV	3.304348	161
		
INTERNATIONAL BL	1.780769	260
WALNUT AV	1.731544	149
WILSON WY	1.710744	121
MARCH LN	1.699605	253
SIERRA AV	1.630695	417

Secondary road	Avg severity	count
FRANCISQUITO AVE	4.379747	158
LAKEWOOD BOULEVARD	3.227941	136
ROXFORD ST	2.950820	183
VAN NESS AV	2.926966	178
TYLER ST	2.898734	158
•••		•••
OLIVE AV	1.829268	123
BANCROFT AV	1.771930	114
CITRUS AV	1.767123	146
10TH ST WEST	1.729508	122
SIERRA AV	1.526946	167

Feature selections secondary road

'QUAIL CREEK RD', 'HARRIS RD', 'SIERRA AVE', ..., 'VANGUARD DR', 'EMERSON RD', 'GATESIDE CT'

The number of collisions happened are more than 100 times.



Correlation

Columns with correlation more than 0.25

Machine learning Models

□Random Forest

Logistic Regression

RANDOM FOREST

TUNING

The accuracy of the predictor on the training and test data.

Accuracy= 96%

Train score: 0.9646983585589427

Test score: 0.9646704374271665

RANDOM FOREST IMPORTANCE FOR INTERPRETABILITY

feat	importance
Victim role	0.449560
Victim safety equipment_1	0.241410
Victim seating position	0.143473
Secondary road	0.115767
Primary road	0.034933
Type of collision	0.014857

Pred	property damage only	pain	other injury	severe injury	fatal
Actual					
property damage only	6819	28	29	10	0
pain	0	36178	25	16	0
other injury	0	665	19928	14	0
severe injury	63	352	545	4493	0
fatal	66	74	194	405	462

RANDOM FOREST CONFUSIO N MATRIX

Logistic Regression



MODEL BEST PARAMS: C = 1000



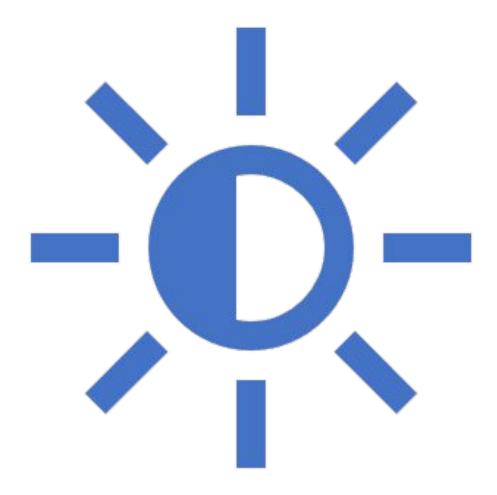
MODEL.SCORE(X_TRAIN, Y_TRAIN)= 0.6853975698145385



MODEL.SCORE(X_TEST, Y_TEST)= 0.6872637353267202

LOGISTIC REGRESSIO N CONFUSIO N MATRIX

Pred	property damage only	pain	other injury	severe injury	fatal
Actual					
property damage only	3701	3013	170	2	0
pain	815	31109	4189	102	4
other injury	56	8668	11549	327	7
severe injury	3	1026	2703	1668	53
fatal	0	106	364	398	333



Comparisons

Random Forest Model can predict the outcome better than Logistic Regression model.

Conclusions

The outcome of the research shows that Victim seating position, type of collision, victim role, victim safety equipment 1, primary road, secondary road are important factors and the most important is victim role on severity of collision.

In the research, a model to predict collision severity based on important features are represented.

The results are important for DMV, Police, or even regular people to manage reducing severity of collisions.

Future work

Our model can be improved by using more features, more data exploration, and more machine learning models, we can answer to some interesting and useful questions such as:

- ☐ When do different makes and models of motorcycles crash?
- ☐ On what days are pedestrians involved in collisions?
- ☐ Are DUIs more likely on certain days?
- ☐ How has COVID19 changed collisions? Are different types of vehicles involved? At different times?
- ☐ What type/color of car/motorcycle is involved in the most crashes?
- ☐ Can we predict injuries based on the other variables in a collision?
- ☐ Are their more accidents around sunrise/sunset?

Reference

California Traffic Collision Data from SWITRS | Kaggle