```
1 !python -V
→ Python 3.11.11
 1 import numpy as np
 2 import pandas as pd
 3 import geopandas as gpd
 4 import plotly.graph_objects as go
 5 import plotly.express as px
 6 import matplotlib.ticker as ticker
 7 from datetime import datetime
 8 from sklearn import tree
 9 from itertools import chain
10
11
12 %matplotlib inline
 1 # data recieved from https://catalog.data.gov/dataset/crash-reporting-drivers-
 2 filename = "/content/maryland_crash_data.csv"
 3 df = pd.read_csv(
 4
       filename,
       dtype={
 6
            "Report Number": str,
           "Local Case Number": str,
 8
           # "Crash Date/Time": datetime,
 9
           "Latitude": str,
           "Longitude": float,
10
           "Vehicle Year": int,
11
12
       },
13
       parse_dates=[
           "Crash Date/Time",
14
15
       ],
```

Project

16)

Data source

As discussed in class, as the professor advised, I am changing my data source to be aligned with the govt data sources. My data source for this mining operation can be found at <u>Data gov</u> site

```
The data has 192183 rows and 39 columns.
  1 features = [
         'Weather',
  3
         'Surface Condition',
         'Light',
         'Driver Substance Abuse',
  6
         'Driver Distracted By',
         'Speed Limit',
  8
         'difference_crash_make_year',
  9]
 10
 11 result = ['Vehicle Damage Extent']
  1 # Sanitize to remove absurd make year values
  2 df = df[df["Vehicle Year"] < 2025]</pre>
  3 df = df[df["Vehicle Year"] > 1970]
  1 df["Vehicle Year"] = pd.to_datetime(df["Vehicle Year"], format='%Y')
  2 df["difference_crash_make_year"] = (df["Crash Date/Time"] - df["Vehicle Year"]
     <ipython-input-6-7da753a0cb32>:1: SettingWithCopyWarning:
     A value is trying to be set on a copy of a slice from a DataFrame.
     Try using .loc[row_indexer,col_indexer] = value instead
     See the caveats in the documentation: <a href="https://pandas.pydata.org/pandas-docs/st">https://pandas.pydata.org/pandas-docs/st</a>
       df["Vehicle Year"] = pd.to_datetime(df["Vehicle Year"], format='%Y')
     <ipython-input-6-7da753a0cb32>:2: SettingWithCopyWarning:
     A value is trying to be set on a copy of a slice from a DataFrame.
     Try using .loc[row_indexer,col_indexer] = value instead
     See the caveats in the documentation: <a href="https://pandas.pydata.org/pandas-docs/st">https://pandas.pydata.org/pandas-docs/st</a>
       df["difference_crash_make_year"] = (df["Crash Date/Time"] - df["Vehicle Year
  1 df = df[df["difference_crash_make_year"] >= 0 ]
Standardize the inputs
This way, futher processing can become easier!
First, let's make the text fields lower-case-d
Then let's filter out records that don't work too well, such as nan or unknown or other values
   1 for coll in chain(features, result):
```

```
df[coll] = df[coll].str.lower()
     except AttributeError as e:
       print(f"Could not convert {coll} to lower case")
   Could not convert Speed Limit to lower case
   Could not convert difference_crash_make_year to lower case
 1 bad_values = {
 2
     'Vehicle Damage Extent': [
         'unknown',
 4
         'other',
         'disabling',
         'vehicle not at scene',
 6
    ],
    'Surface Condition': [
 8
         'unknown',
10
         'other',
11
     ],
     'Light': [
12
13
         'unknown',
14
         'other',
15
16
     'Driver Substance Abuse': [
         'unknown',
17
         'other',
18
19
         'unknown, unknown',
20
    ],
21
     'Driver Distracted By': [
22
         'unknown',
23
         'other',
24
    ],
25 }
26
27 for coll in chain(features, result):
     df = df[df[coll].notna()]
29
30 for coll, values in bad_values.items():
31
32
       df = df[~df[coll].isin(values)]
33
     except Exception as e:
34
       print(f"Could not remove unwanted vals for {coll}")
35
       print(e)
 1 display(df)
                              Local
                                                      ACRS
                                                                Crash
                  Report
                                                                          Route
                               Case
                                      Agency Name
                                                    Report
                                                                Date/
                                                                                 Road
                  Number
                                                                           Type
                             Number
                                                                 Time
                                                      Type
```

	5	MCP3348000Z	230051804	Montgomery County Police	Injury Crash	2023-08-28 11:09:00	เพลryเลทด (State)	NU
	6	MCP302600BD	230046425	Montgomery County Police	Property Damage Crash	2023-07-27 12:30:00	County	GREE
	8	MCP3372001V	230065250	Montgomery County Police	Property Damage Crash	2023-11-10 20:24:00	Maryland (State)	GE
	17	MCP2962008G	230065146	Montgomery County Police	Property Damage Crash	2023-11-08 14:05:00	County	F H
	19	MCP3136006J	230048375	Montgomery County Police	Property Damage Crash	2023-08-08 11:39:00	Maryland (State)	BRAN
	192169	MCP312900BT	250003767	MONTGOMERY	Injury Crash	2025-01-27 07:56:00	County Route	
	192171	MCP3171002L	250004790	MONTGOMERY	Property Damage Crash	2025-02-02 12:30:00	County Route	CEI BLVD
	192172	MCP235300D1	250004139	MONTGOMERY	Property Damage Crash	2025-01-29 06:40:00	Maryland (State) Route	ROC PIKE
	192174	MCP283200B9	250004515	MONTGOMERY	Property Damage Crash	2025-01-31 12:31:00	County Route	MAN
	192176	MCP3404001T	250004298	MONTGOMERY	Injury Crash	2025-01-30 07:00:00	County Route	LOCK

62916 rows × 40 columns

Visualization

Here we can plot all the accidents in our dataset and see how clustered the accidents are and we can also see what year cars are more accident prone and if it has anything to do with how old the cars are.

There might be other interesting ideas that might pop up after checking out the visualizations!

```
1 fig_vehicle_year = px.scatter_geo(df, lat="Latitude", lon="Longitude", color='
2 fig_vehicle_year.show()
```

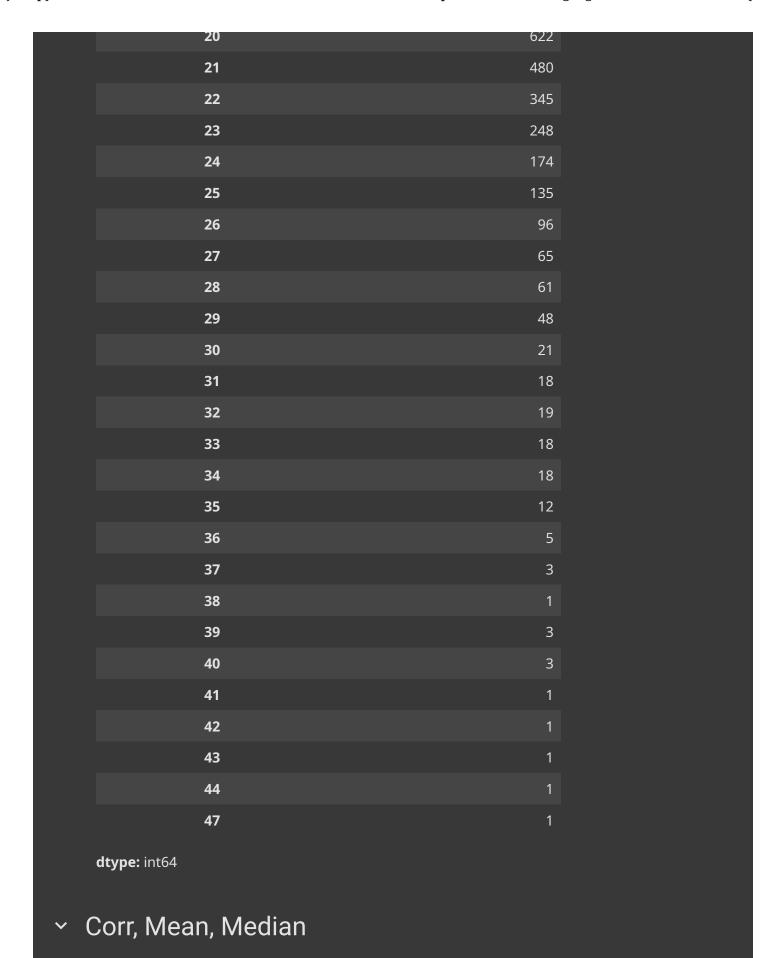
4/4/25, 1:05 AM



1 display(df.groupby(["difference_crash_make_year"])["difference_crash_make_ye						
differe	difference_crash_make_year					
difference_crash_make_year						
0	3177					
1	4646					
2	4615					
3	4608					
4	4580					
5	4352					
6	4147					
7	3789					
8	3708					
9	3406					
10	3257					
11	2892					
12	2594					
13	2369					
14	2003					
15	1716					
16	1506					
17	1281					
18	1031					
19	839					
	100					

4/4/25, 1:05 AM

7 of 15



I'm still understanding what useful knowledge can be extracted to use to make safer cars.

```
I will use Corr, Mean, Median to help determine if cars have gotten safer over time.
  1 df[[
        "Vehicle Year",
        "difference_crash_make_year"
  4 ]].corr(min_periods=3)
                                                                               ᇤ
                                Vehicle Year difference_crash_make_year
            Vehicle Year
                                      1.000000
                                                                   -0.883547
                                                                               11.
                                                                    1.000000
     difference_crash_make_year
                                     -0.883547
  1 df[[
        # "Vehicle Year",
       "difference_crash_make_year"
  4 ]].median(
        axis=0,
  6)
                                  0
     difference_crash_make_year 7.0
    dtype: float64
  1 df[[
        # "Vehicle Year",
        "difference_crash_make_year"
  4 ]].mean(
        axis=0,
  6)
                                       0
     difference_crash_make_year 7.865805
    dtype: float64
  1 df[[
       "Vehicle Make",
  3 ]].value_counts().idxmax()
    ('TOYOTA',)
```

rındıriyə

As we can see, on average, we can expect accidents to start occuring after the car is around 7~8 years of age. This dataset also highlights that the 2012 models are more likely to be the ones in accident. This could also just be a red herring that is throwing off our data since a lot of crashes can go unreported.

Algo selection

I would like to explore a few different ideas. But one of the algorithms we were taught in class was decision tree. I would like to explore more columns with decision tree to determine rates of accidents given the factors, such as:

- Weather
- Driver substance abuse
- Speed limit of road

These are our feature-set.

We will use features to predict how bad the accident turns out to be, based on the damage occurred to the vehicle.

A major reason to use the decision tree is that it is easy to see the outcome. This will produce a very clean algebraic equation that can help customers buy the right car for them!

This is precisely the model I have chosen to implement

After working on it, and enhancing it and playing around with inputs and parameters, this seems to be not a terrible choice!

```
1 \text{ remap} = \{
 2
       'Driver Substance Abuse': {
 3
           'none detected': False,
           'alcohol contributed': True,
           'alcohol present': True,
           'illegal drug present': True,
 6
           'illegal drug contributed': True,
 8
           'medication present': True,
 9
           'combined substance present': True,
           'medication contributed': True,
10
11
           'combination contributed': True,
           'not suspect of alcohol use, not suspect of drug use': False,
12
13
           'suspect of alcohol use, not suspect of drug use': True,
           'suspect of alcohol use, unknown': True,
14
15
           'unknown, not suspect of drug use': False.
```

```
16
           'not suspect of alcohol use, unknown': False,
17
           'suspect of alcohol use, suspect of drug use': True,
           'not suspect of alcohol use, suspect of drug use': True,
18
19
       'Light': {
20
21
           'dark -- unknown lighting': 'dark',
22
           'dark - unknown lighting': 'dark',
23
       },
24 }
25
26 pd.set_option('future.no_silent_downcasting', True)
27 for feature, mapping in remap.items():
28
     try:
29
       df = df.replace({ feature: mapping })
30
    except Exception as e:
31
       print(f"Could not remap {feature}")
32
       print(e)
```

```
1 for feature in chain(features, result):
2 print(f"{feature}: {list(df[feature].unique())}")
```

```
Weather: ['cloudy', 'clear', 'raining', 'snow', 'foggy', 'wintry mix', 'other' Surface Condition: ['dry', 'wet', 'ice', 'slush', 'water(standing/moving)', 's Light: ['daylight', 'dark lights on', 'dawn', 'dusk', 'dark no lights', 'dark' Driver Substance Abuse: [False, True] Driver Distracted By: ['not distracted', 'looked but did not see', 'other dist Speed Limit: [np.int64(30), np.int64(35), np.int64(40), np.int64(5), np.int64(difference_crash_make_year: [np.int64(13), np.int64(1), np.int64(7), np.int64(Vehicle Damage Extent: ['functional', 'superficial', 'no damage', 'destroyed']
```

Narrow down

Here we will narrow down our df to use only some columns, like our feature columns and outcome column(s)

```
1 df = df[chain(features, result)]
```

ONE HOT ENCODING

Other than the great models, I have also learnt that pre-processing is very important! One Hot Encoding is another useful option in preprocessing. This has been absolutely eye-opening for me when it comes to data processing and extracting information and knowledge from it using mining methods.

 $10 { of } 15$

```
1 for coll in chain(features, result):
 2
    if coll in [
         'Driver Substance Abuse',
         'Speed Limit',
         'difference_crash_make_year',
 6
         'Vehicle Damage Extent',
    ]:
 8
       continue
 9
     try:
10
       tmp = pd.get_dummies(df[coll], prefix=coll.lower(), prefix_sep='_')
       df = df.drop(coll, axis = 1)
11
12
       df = df.join(tmp)
13
     except Exception as e:
       print(f"Could not convert {coll} to one hot encoding")
14
15
       print(e)
 1 display(list(df.columns))
   ['Driver Substance Abuse',
    'Speed Limit',
    'difference_crash_make_year',
    'Vehicle Damage Extent',
    'weather_blowing sand, soil, dirt',
    'weather_blowing snow',
    'weather_clear'
    'weather cloudy',
    'weather_fog, smog, smoke',
    'weather_foggy',
    'weather_freezing rain or freezing drizzle',
    'weather_other',
    'weather_rain',
    'weather_raining',
    'weather_severe crosswinds',
    'weather severe winds',
    'weather_sleet',
    'weather_sleet or hail',
    'weather_snow',
    'weather_unknown'
    'weather_wintry mix',
    'surface condition_dry',
    'surface condition_ice',
    'surface condition_ice/frost',
    'surface condition_mud, dirt, gravel',
    'surface condition_oil',
    'surface condition sand'
    'surface condition_slush',
    'surface condition_snow',
    'surface condition_water (standing, moving)',
    'surface condition_water(standing/moving)',
    'surface condition_wet',
    'light_dark',
    'light dark - lighted'
```

```
'light dark - not lighted',
    'light_dark lights on',
    'light_dark no lights',
    'light_dawn',
    'light_daylight',
    'light_dusk',
    'driver distracted by adjusting audio and or climate controls',
    'driver distracted by_by moving object in vehicle',
    'driver distracted by by other occupants',
    'driver distracted by_dialing cellular phone',
    'driver distracted by_distracted by outside person object or event',
    'driver distracted by_eating or drinking',
    'driver distracted by_inattentive or lost in thought',
    'driver distracted by_looked but did not see',
    'driver distracted by_manually operating (dialing, playing game, etc.)',
    'driver distracted by_no driver present',
    'driver distracted by_not distracted',
    'driver distracted by_other action (looking away from task, etc.)',
    'driver distracted by_other cellular phone related',
    'driver distracted by_other distraction',
    'driver distracted by_other electronic device (navigational palm pilot)',
    'driver distracted by_smoking related',
    'driver distracted by_talking or listening to cellular phone',
    'driver distracted by_talking/listening',
    'driver distracted by_texting from a cellular phone',
    'driver distracted by_using device object brought into vehicle',
    'driver distracted by_using other device controls integral to vehicle']
1 results = [
      # 'vehicle damage extent destroyed',
      # 'vehicle damage extent_no damage',
4
      # 'vehicle damage extent_functional',
      # 'vehicle damage extent_superficial',
      'Vehicle Damage Extent',
6
7 ]
8
9 features = [
10
       'Driver Substance Abuse',
       'Speed Limit',
11
       'difference_crash_make_year',
12
13
       'weather_blowing sand, soil, dirt',
       'weather blowing snow',
14
15
       'weather_clear',
       'weather cloudy',
16
17
       'weather_fog, smog, smoke',
18
       'weather_foggy',
19
       'weather freezing rain or freezing drizzle',
20
       'weather_other',
21
       'weather rain',
22
       'weather raining',
23
       'weather_severe crosswinds',
24
       'weather severe winds'.
```

```
25
       'weather_sleet',
       'weather sleet or hail',
26
27
       'weather snow',
28
       'weather unknown',
       'weather wintry mix',
29
30
       'surface condition_dry',
       'surface condition_ice',
31
32
       'surface condition_ice/frost',
33
       'surface condition_mud, dirt, gravel',
       'surface condition oil',
34
       'surface condition_sand',
35
       'surface condition slush',
36
37
       'surface condition_snow',
38
       'surface condition_water (standing, moving)',
       'surface condition_water(standing/moving)',
39
40
       'surface condition_wet',
41
       'light_dark',
42
       'light_dark - lighted',
43
       'light_dark - not lighted',
44
       'light_dark lights on',
45
       'light_dark no lights',
       'light_dawn',
46
47
       'light_daylight',
48
       'light dusk',
       'driver distracted by_adjusting audio and or climate controls',
49
       'driver distracted by_by moving object in vehicle',
50
51
       'driver distracted by by other occupants',
52
       'driver distracted by_dialing cellular phone',
       'driver distracted by_distracted by outside person object or event',
53
       'driver distracted by_eating or drinking',
54
       'driver distracted by_inattentive or lost in thought',
55
       'driver distracted by_looked but did not see',
56
       'driver distracted by_manually operating (dialing, playing game, etc.)',
57
       'driver distracted by no driver present',
58
       'driver distracted by_not distracted',
59
       'driver distracted by_other action (looking away from task, etc.)',
60
       'driver distracted by_other cellular phone related',
61
       'driver distracted by_other distraction',
62
       'driver distracted by_other electronic device (navigational palm pilot)',
63
       'driver distracted by_smoking related',
64
       'driver distracted by_talking or listening to cellular phone',
65
       'driver distracted by_talking/listening',
66
       'driver distracted by_texting from a cellular phone',
67
68
       'driver distracted by_using device object brought into vehicle',
69
       'driver distracted by_using other device controls integral to vehicle',
70 ]
```

Train Test

4/4/25, 1:05 AM

Here we will set up the training data and the test data and use it to generate and score a decision tree.

```
1 msk = np.random.rand(len(df)) < 0.8
2
3 df_train = train = df[msk]
4 df_test = df[~msk]

1 df_features = df_train[features]
2 df_result = df_train[results]
3
4 dtree = tree.DecisionTreeClassifier()
5 dtree = dtree.fit(df_features, df_result)

1 score = dtree.score(df_test[features], df_test[results])
2 print(f"Score: {score}")

Score: 0.45638801261829653</pre>
```

Score

So this is probably my 5th or 7th iteration

Earlier I got even worse scores

But playing around, I was able to see what works.

Some of my learnings:

- One Hot Encoding is a must, or else it can get pretty bad interpretting certain attributes as categories
- Encoding of the resultant category is not as necessary and is better not to encode it
- Geographical data points are ONLY useful if we use geo encodings so that the data points'
 closeness can be interpretted by computers. Seeing it on a map is not enough. A deep dive
 in the near future will lead me to Geographical databases or geo databses. They are super
 cool!
- In terms of findings, I think I am confident to say (after many selections of features and resultants etc.) that these features highlighted in this version of my notebook is really really good! however, further fine tuning is necessary! I wish to expand this project into the future to see its potential!
- I was also reading that there are more types of decision trees even, and this package may
 not have them all! I would love to code one up from scratch!

