# Car Price Prediction

By: Adkeme Berhe
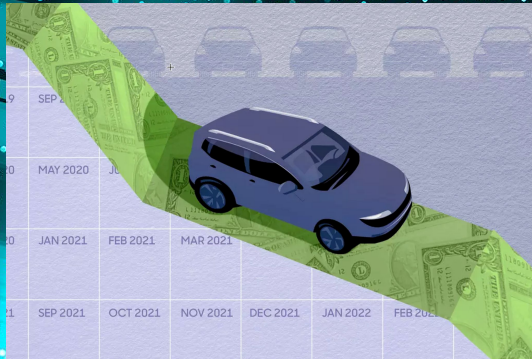
# Today's Agenda

1. Background
2. Business problem
3. Data Sources
4. Data Preparation
5. Modeling
6. Evaluation

# Business Problem

- Due to an ongoing chip shortage and increased cost in raw materials the U.S average price for a new car was up 6%

- Ahead of 2023's first quarter, Group 11 I have been tasked with creating a Price Prediction model to help determine car prices.

# Data Sources

- Obtained from Kaggle at
  https://www.kaggle.com/datasets/deepcontractor/car-price-prediction-challenge/discussion/367511

- **19,237** rows, **18** columns. **(2.2mb)**
  Most notable Car features
  - Price
  - Manufacturer
  - Model
  - Fuel Type
  - Production Year

# Data Preparation



- Before any cleaning or pre processing model ran at a 6% accuracy
- Went through and checked for NULL Values
- After further analyzing the data, 313 duplicate values were fund and dropped.
- With prices ranging from $0 to 2.38 billion we made it more practical, removing <1,000 and >100,000
- Mileage included KM unit which needed to be removed for numerical reasons
- Missing Levy feature resulted in dropping

FROM: 19,237 —> TO: 8,027 rows

(75%-25% split) of our remaining 8,027 rows

FROM: 18 —>17 columns



```
In [262]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=5)
```

# Model

```python
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.neighbors import KNeighborsRegressor
```

```python
In [261]: x = cars_f.drop(['Price', 'Color'], axis=1).values
          y = cars_f['Price'].values
          x = MinMaxScaler().fit_transform(x)
```

```python
In [262]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=5)
```

```python
In [263]: #lr  = LinearRegression()
          knn = KNeighborsRegressor(n_neighbors=8)
          dt  = DecisionTreeRegressor(max_depth = 5)
          rf  = RandomForestRegressor(n_estimators=100, max_features= 7)

          regressors = [ ('K Nearest Neighbours', knn),('Decision Tree', dt), ('Random Forest', rf)]
```

```python
In [264]: for regressor_name, regressor in regressors:

              regressor.fit(x_train, y_train)

              y_pred = regressor.predict(x_test)
              accuracy = round(r2_score(y_test,y_pred),3)*100

              print('{:s} : {:.0f} %'.format(regressor_name, accuracy))
              plt.rcParams["figure.figsize"] = (20,8)
              plt.bar(regressor_name,accuracy)
```
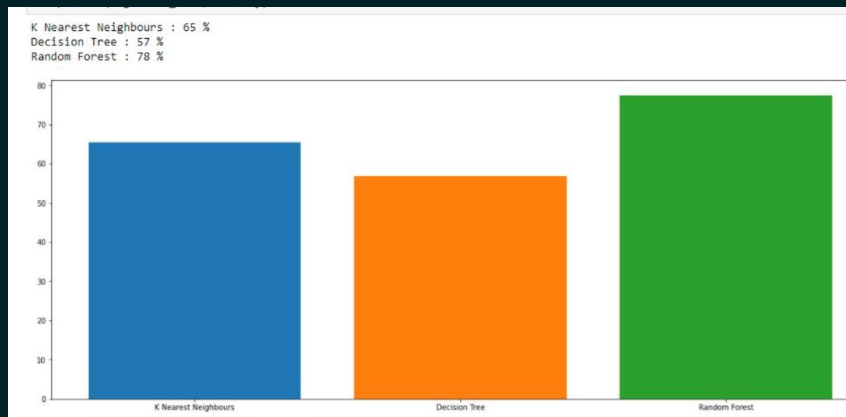
knn=8 neighbors

dt= depth of 5

rf= 7 features, 100 estimators

Stored in regressors, fitted our testing

Printed our prediction & accuracy

Created a bar graph for depiction

# Evaluation

| Model Type | KNN | Decision Tree | Random Forest |
|---|---|---|---|
| Parameters | N neighbors value | Max Depth | Estimator & max features |
| Values | 8 | 5 | 100 and 7 |
| Accuracy % | 65% | 48% | 77% |



K Nearest Neighbours : 65 %
Decision Tree : 57 %
Random Forest : 78 %

# Recommendations & Lessons Learned

While we are satisfied with our Random Forest Model, more time may have allowed us the chance to further increase our model accuracy

Buyer and seller discretion is advised due to accuracy level being moderate.

---------------------------------------------------------------------

Data preparation is the most tedious (and annoying) part of everything. Its importance cant be understated.

Dont forget to consider business value and ROI for companies

# Future work



- Currently no plans to extend project any further.

- I will utilize my project as a resume portfolio piece .

*Adkeme Berhe*
https://www.linkedin.com/in/adkeme-berhe-2bb456240/

# Thank you for your time !