

## Prediction Competition # 5

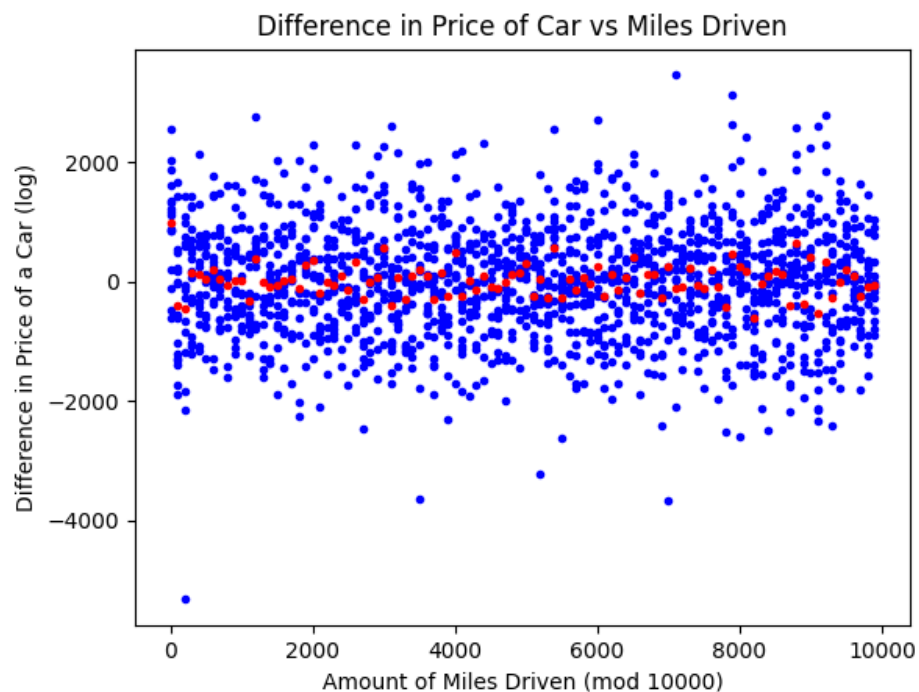
---

**Q1)** To solve this problem, I used a variation of boosting. Due to the fact that across  $k$  feature space the data could still not be accurately linearly estimated (as shown in Q2), what I did was split the data into buckets. Each bucket would have its own boosted model trained on the data in the bucket. Then I would use a boosted model to create a rough estimate for the price, and determine what bucket to use to get a more accurate measure. Much like the previous prediction contest that means we get the following results:

$$MSE = 0.00983$$

$$R^2 = 0.87386$$

**Q2)** Below is a graph showing the difference in price over the distance taken to modulus of 10000. What this does is lets us do is analyse if theres a substantial difference between any of the digits (left most bias) as can be seen on the graph below, the blue shows observations and the red shows average observation value at that modulus value:



As we can see, there are spikes that occur before any of the thousands and after the spikes we see the price grow. This can be seen at 4000 and 6000 very easily. This demonstrates the importance of the thousand digit and how consumers far over weigh this. Thus giving us proof of the left bias.

**Q3)** To solve Q1 chatGPT was given the prompt:

Given the same car prediction model, how can I improve predictions for cars that have chatGPT provided the following response:

If you want to improve predictions for cars that have a cost around a specific threshold

Data Collection and Augmentation:

- a. Gather more data specifically for cars around the threshold price. This can help your
- b. Augment your dataset by creating synthetic data points around the threshold price. Yo

Feature Engineering:

- a. Create new features that capture the unique attributes of cars near the threshold pri
- b. Explore interactions between features and their effects on cars at the specified pric

Model Selection and Tuning:

- a. Choose a model that is well-suited for capturing the nuances of cars at the threshold
- b. Fine-tune hyperparameters, particularly those related to the model's sensitivity to d
- c. Consider using specialized regression techniques, such as Quantile Regression, that f

Ensemble Models:

- a. Use ensemble models like stacking or blending to combine predictions from multiple mo

Feature Importance Analysis:

- a. Analyze the feature importance scores of your model, which can help identify the feat

Custom Threshold Handling:

- a. Implement custom post-processing steps to fine-tune predictions around the threshold

Analyze Residuals:

- a. Examine the residuals (the differences between actual and predicted values) for cars

Cross-Validation:

- a. Use cross-validation techniques that focus on the threshold price range, such as stra

Domain Knowledge:

- a. Consult experts or domain-specific knowledge to understand what features and character

Continuous Monitoring and Refinement:

- a. Continuously monitor the model's performance, and be prepared to update the model as

Remember that improving predictions for a specific price threshold might require a combi

To solve Q2 chatGPT was given the prompt:

Using a sample dataset how can I demonstrate the impact the most significant digit has on the response variable?  
chatGPT provided the following response:

To demonstrate the impact of the most significant digit on the response variable, you can follow these steps:

**Generate or Prepare a Sample Dataset:**

Start by creating or obtaining a sample dataset that includes a numerical feature and a response variable.

Extract the first digit from the numerical feature. You can do this by converting the numerical feature to a string and then extracting the first character.

**Aggregate and Analyze the Data:**

Calculate summary statistics for the response variable based on the first digit of the numerical feature. This can be done using a statistical software package or a programming language like R or Python.

**Visualize the Impact:**

Create visualizations to demonstrate the impact of the first digit. You can create bar charts, line plots, or scatter plots to show the relationship between the first digit and the response variable.

**Interpret the Results:**

Analyze the results to determine if there is a notable impact of the first digit on the response variable. Look for patterns or trends in the data that suggest a relationship.

The analysis demonstrates how the first digit of the numerical feature affects the response variable. For example, if the first digit is 1, the response variable tends to be higher than if the first digit is 2.

## Code for Q1:

---

```
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import xgboost as xgb
import warnings
from pandas.util import hash_pandas_object

warnings.filterwarnings("ignore")
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

# Read in the data
data = pd.read_csv("Econ424_F2023_PC4_training_data_large.csv", low_memory=False)

# ===== General Data Parsing =====
toNumeric = ["bed_length", "back_legroom", "front_legroom", "height", "length", "wheelbase",
             "width", "maximum_seating",
             "fuel_tank_volume"]

data.drop(columns=["major_options", "power", "exterior_color", "city", "interior_color",
                  "listing_color",
                  "transmission", "bed", "bed_height"], inplace=True)

def parseIntoNumeric(x):
    if type(x) == str:
        return float(x.split(" ")[0])
    else:
        return np.nan

for col in data.columns:
    data[col].replace('--', np.nan, inplace=True)
    if col in toNumeric:
        print(col)
        data[col] = data[col].apply(lambda x: parseIntoNumeric(x))

catagories = ["body_type", "engine_type", "fleet", "frame_damaged", "franchise_dealer",
              "franchise_make",
              "fuel_type", "has_accidents", "iscab", "is_certified", "is_cpo", "is_new",
              "is_oemcpo",
              "make_name", "model_name", "salvage", "sp_name", "theft_title",
              "transmission_display",
              "trimid", "trim_name", "vehicle_damage_category", "wheel_system",
              "wheel_system_display"]

data.fillna(0, inplace=True)

for col in data.columns:
    if col in catagories:
        print(data.dtypes[col])
        data[col] = data[col].fillna('None').astype(str)
        encoder = RareLabelEncoder(n_categories=1, max_n_categories=2500, replace_with='Other',
                                   tol=40 / data.shape[0])
        data[col] = encoder.fit_transform(data[[col]])
```

```

# Standardize results
data["price"] = data["price"].map(lambda x: np.log(x))

data = data[data["price"] > np.log(40000)]

catagories = ["body_type", "engine_type", "fleet", "frame_damaged", "franchise_dealer",
              "franchise_make",
              "fuel_type", "has_accidents", "iscab", "is_certified", "is_cpo", "is_new",
              "is_oemcpo",
              "make_name", "model_name", "salvage", "sp_name", "theft_title",
              "transmission_display",
              "trimid", "trim_name", "vehicle_damage_category", "wheel_system",
              "wheel_system_display"]

full_pipeline = ColumnTransformer([('cat', OneHotEncoder(handle_unknown='ignore'), catagories)])

# ===== Data Testing =====

X_Train = [0]*10
Y_Train = [0]*10
X_Test = [0]*10
Y_Test = [0]*10

NUMTESTS = 10
for i in range(0, NUMTESTS):

    train, test = train_test_split(data, test_size=0.1)

    Y_Train[i] = train.iloc[:, 0]
    X_Train[i] = train.iloc[:, 1:]

    Y_Test[i] = test.iloc[:, 0]
    X_Test[i] = test.iloc[:, 1:]

innerModel = xgb.XGBRegressor(n_estimators=500, max_depth=6, eta=0.1, reg_lambda=0.1,
                              colsample_bytree=0.4)
MSE = 0
DIF = 0
for i in range(0, NUMTESTS):
    encoder = full_pipeline.fit(X_Train[i])

    X_train = encoder.transform(X_Train[i])
    innerModel.fit(X_train, Y_Train[i])

    X_test = encoder.transform(X_Test[i])
    predicted = innerModel.predict(X_test)

    localMSE = 0
    localDIF = 0
    for x in range(0, len(predicted)-1):
        localMSE += (predicted[x] - Y_Test[i].iloc[x])**2
        localDIF += (Y_Test[i].iloc[x] - np.mean(Y_Test[i]))**2

```

```

    localMSE /= (len(predicted)-1)
    localDIF /= (len(predicted)-1)
    MSE += localMSE
    DIF += localDIF

MSE /= NUMTESTS
DIF /= NUMTESTS

print(MSE)
print(1 - MSE/DIF)

full_pipeline = ColumnTransformer([('cat', OneHotEncoder(handle_unknown='ignore'), catagories)])
innerModel = xgb.XGBRegressor(n_estimators=500, max_depth=6, eta=0.1, reg_lambda=0.1,
                               colsample_bytree=0.4)

Y_Train = data.iloc[:, 0]
X_Train = data.iloc[:, 1:]
encoder = full_pipeline.fit(X_Train)
X_train = encoder.transform(X_Train)
innerModel.fit(X_train, Y_Train)

test = pd.read_csv("Econ424_F2023_PC5_test_data_without_response_var.csv", low_memory=False)

X_test = encoder.transform(test)

val = innerModel.predict(X_test)

f = open('predictions.csv', 'w')
for estimate in val:
    f.writelines(str(estimate) + ",\n")

# ===== Export Model =====

test["predicted_spec_price"] = test.apply(lambda x: pred[x.loc['New_ID']], axis=1)
val = test["predicted_spec_price"]

f = open('predictions.csv', 'w')
for estimate in val:
    f.writelines(str(estimate) + ",\n")

```

---