# MathSoc AutoTrader Hackathon 2024

*Ioan Gwenter, Lourenço Silva, Tom Cassar*
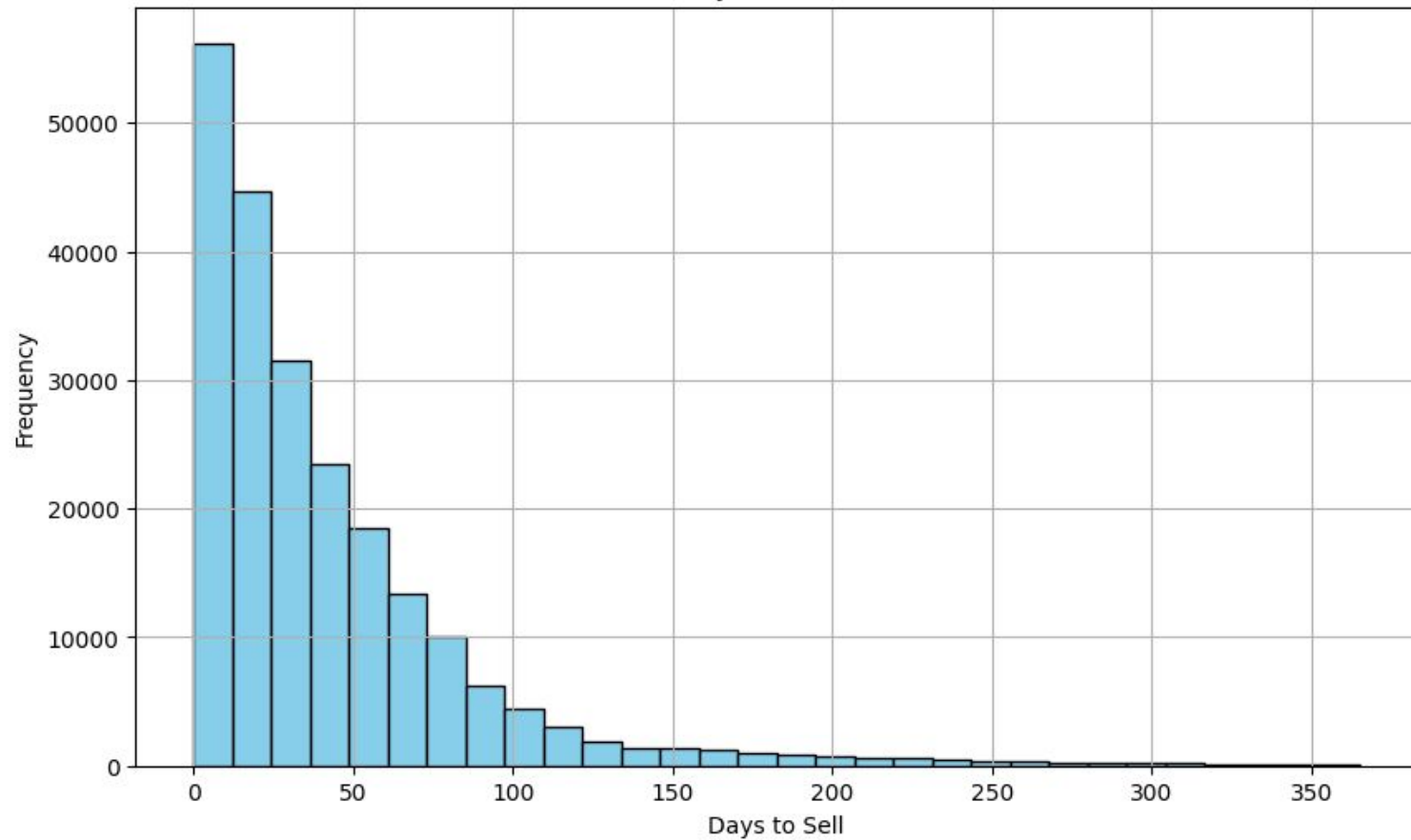
Data Input     Data Cleaning     Pre-processing     Model Training     Deployment

Distribution of Days to Sell

Distribution of Days to Sell (Max Value: 365)

```python
features = ['reg_year',
            'make',
            'model',
            'body_type',
            'fuel_type',
            'transmission_type',
            'drivetrain',
            'colour',
            'price_indicator_rating',
            'postcode_area',
            'first_retailer_asking_price',
            'reviews_per_100_advertised_stock_last_12_months',
            'seats',
            'doors',
            'co2_emission_gpkm',
            'top_speed_mph',
            'zero_to_sixty_mph_seconds',
            'engine_power_bhp',
            'fuel_economy_wltp_combined_mpg',
            'battery_range_miles',
            'battery_usable_capacity_kwh',
            'length_mm',
            'boot_space_seats_up_litres',
            'insurance_group',
            'odometer_reading_miles',
            'adjusted_retail_amount_gbp',
            'predicted_mileage',
            'number_of_images',
            'advert_quality',
            #'percentage_through_year',
            'can_home_deliver',
            'manufacturer_approved',
            'segment'
            ]
```
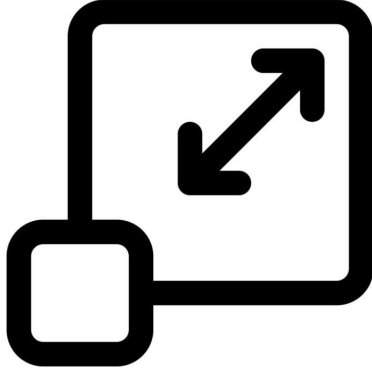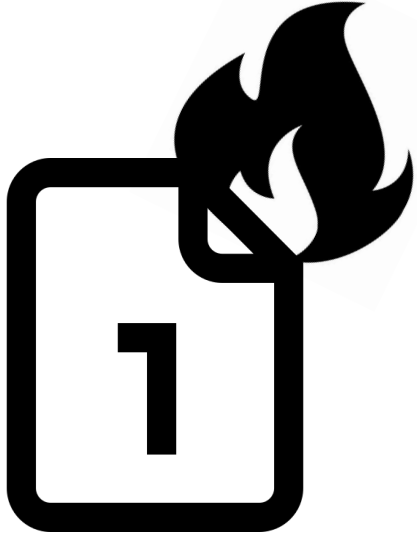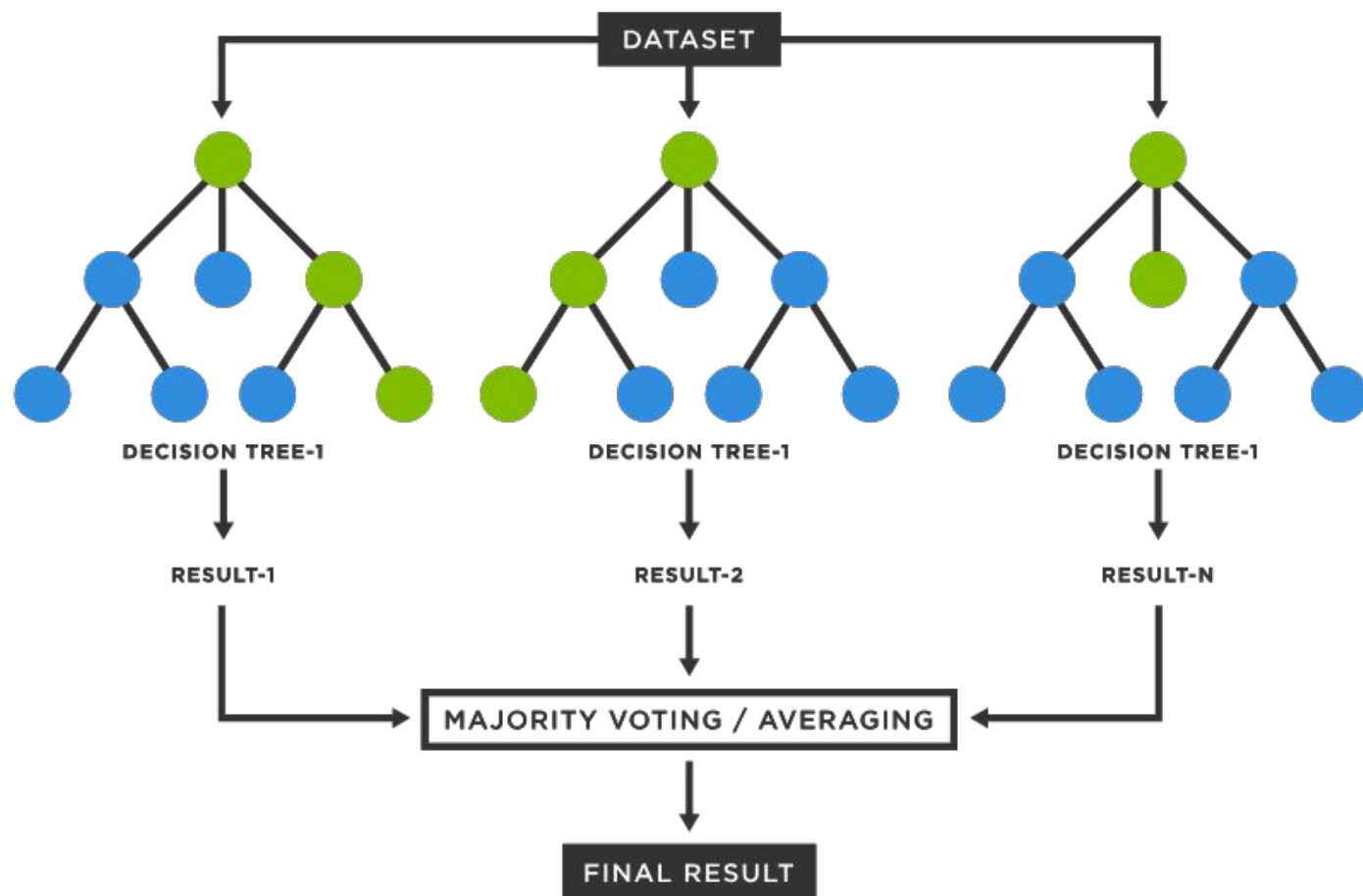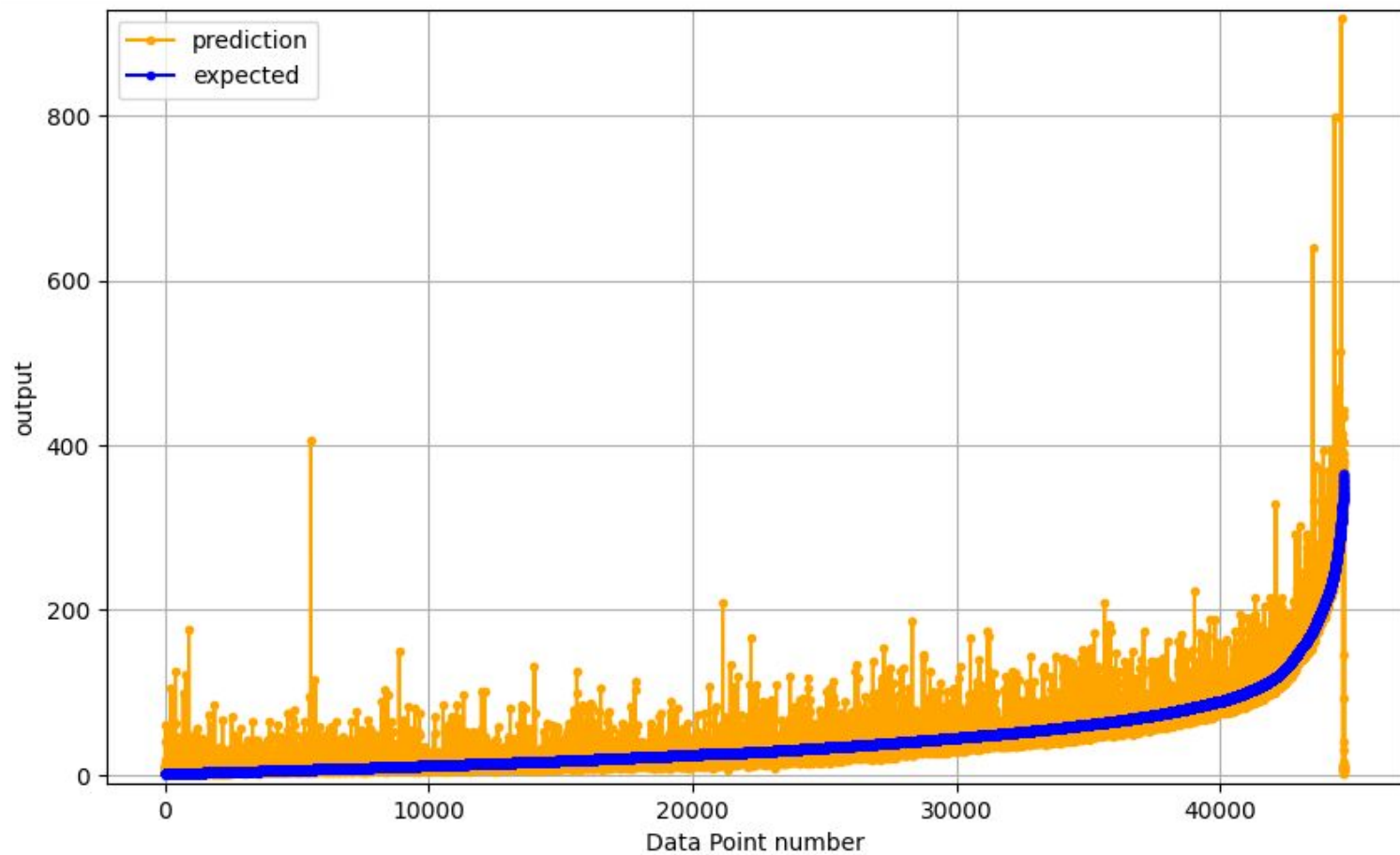
```python
# Split up categorical/nu
categorical_features = ['
```

```python
binary_features = ['can_home_deliver',
```
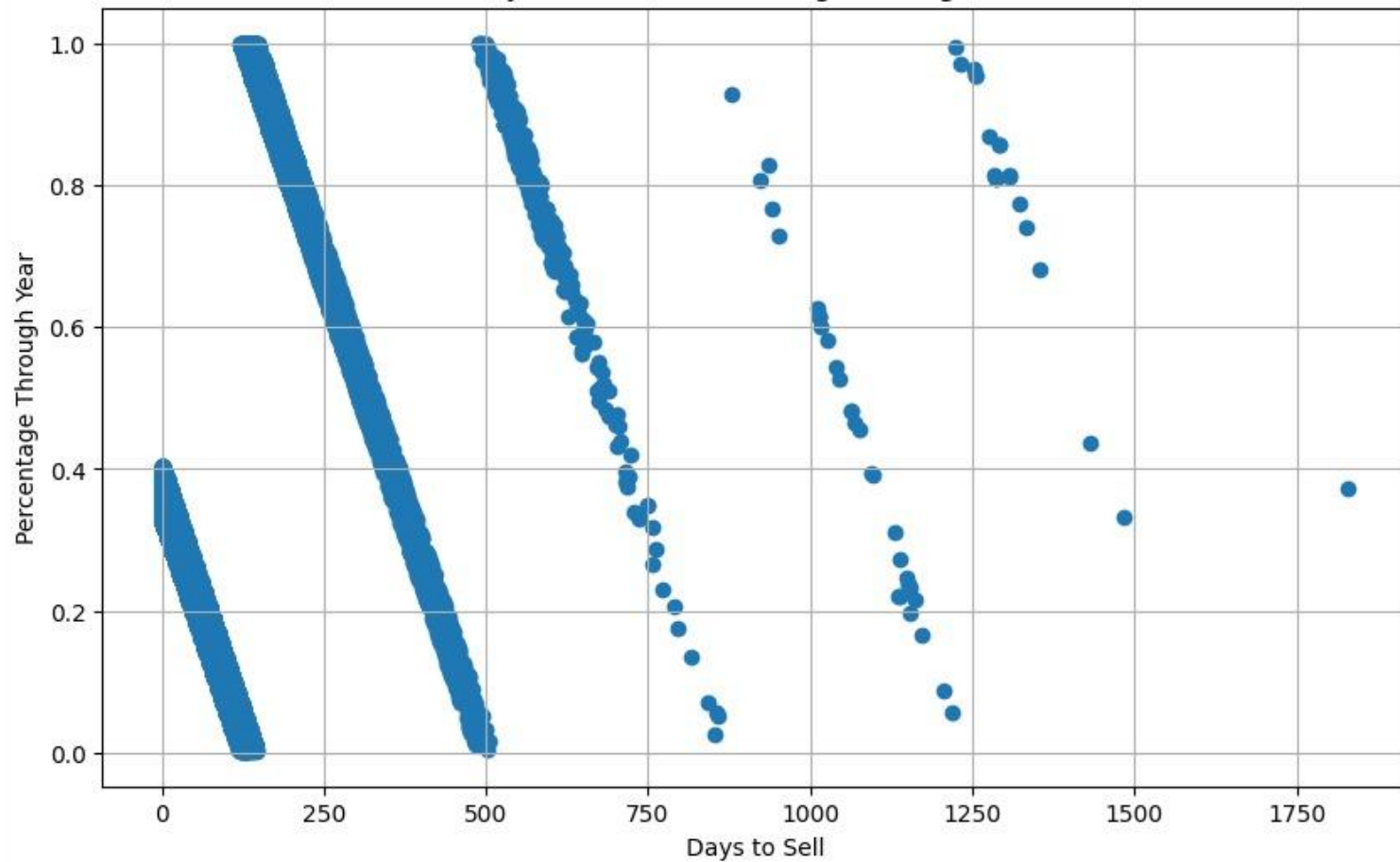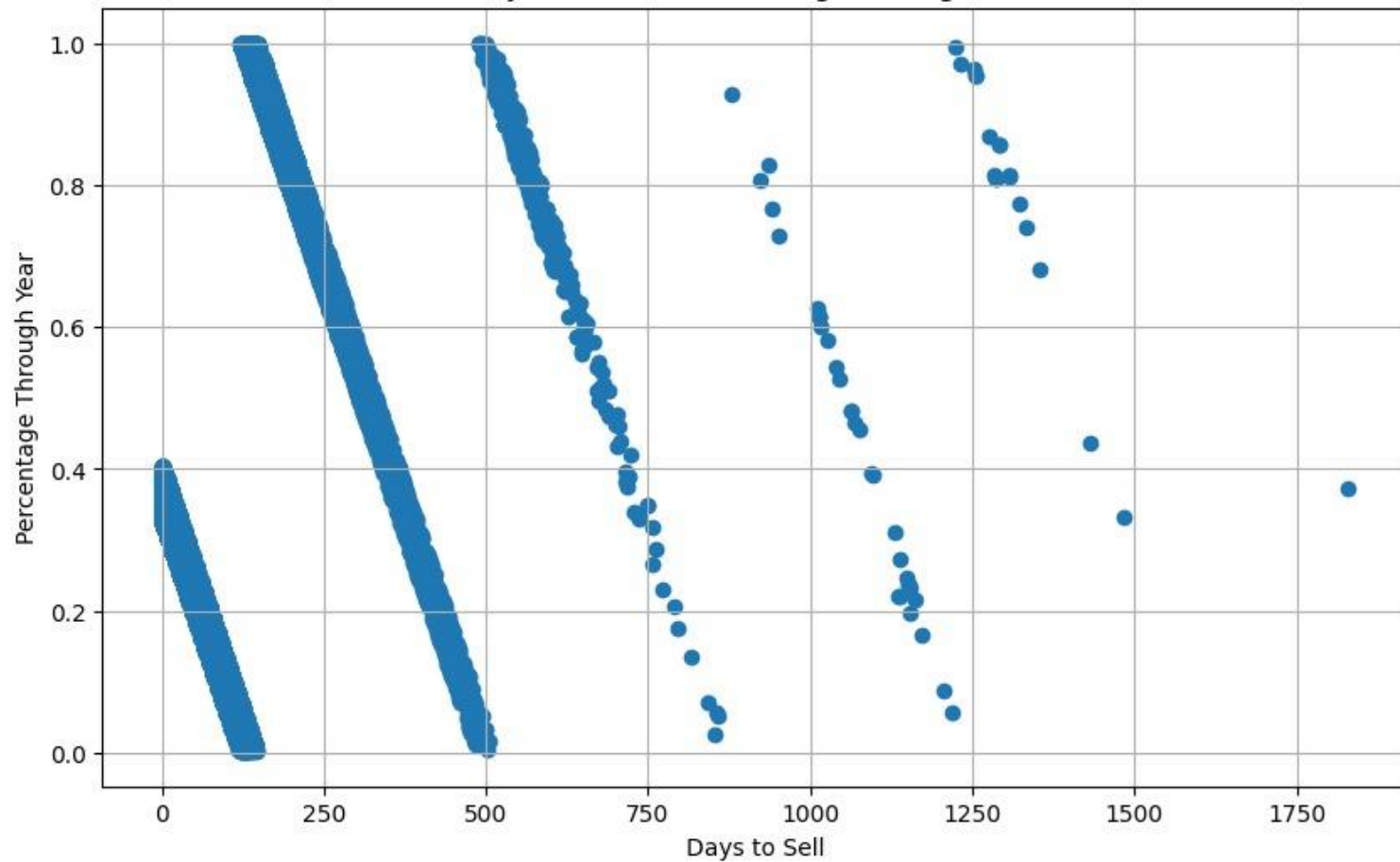
- MAE: 9.00 days
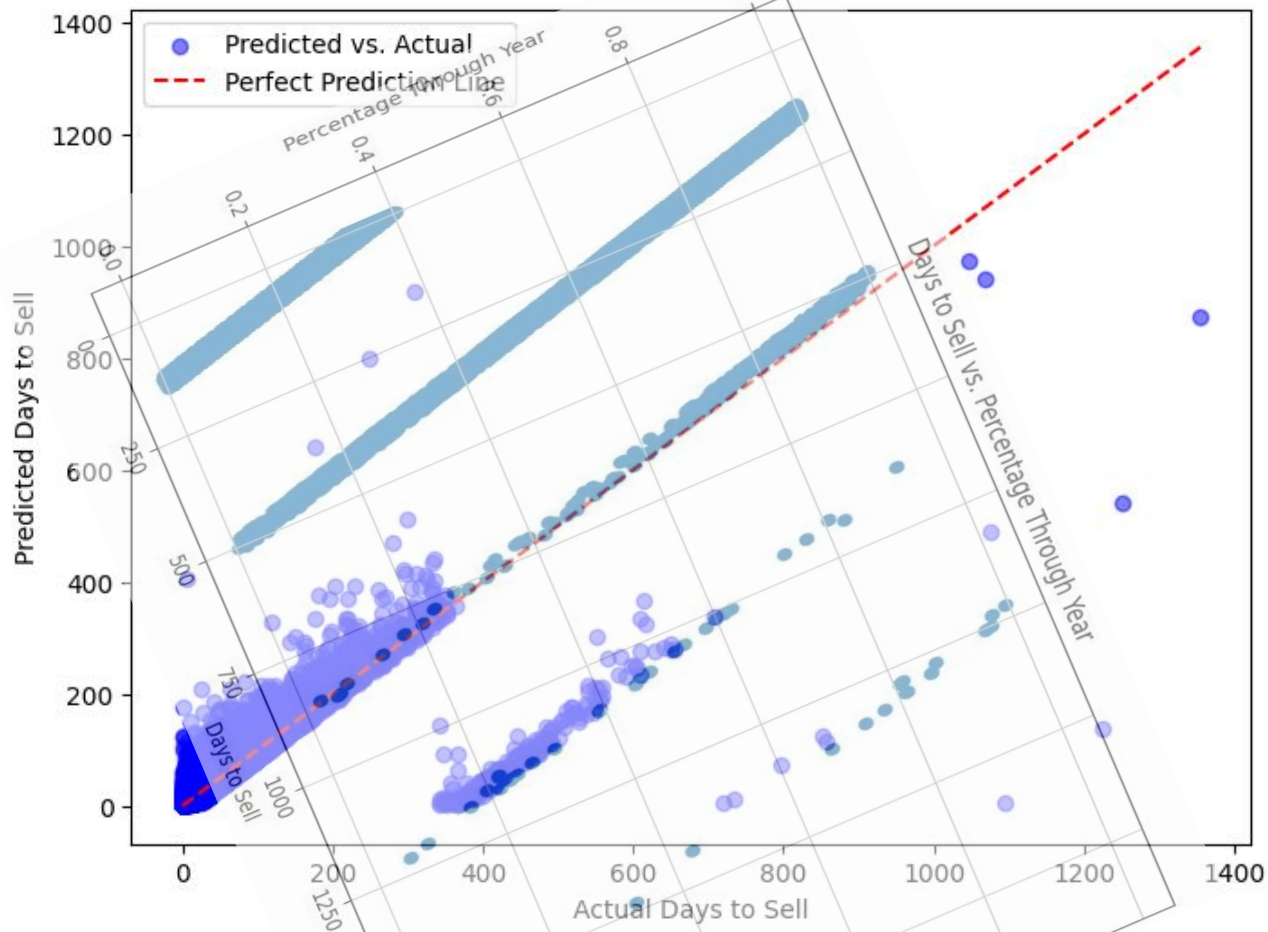- RMSE: ~28.0 days
- R-Squared: 0.711
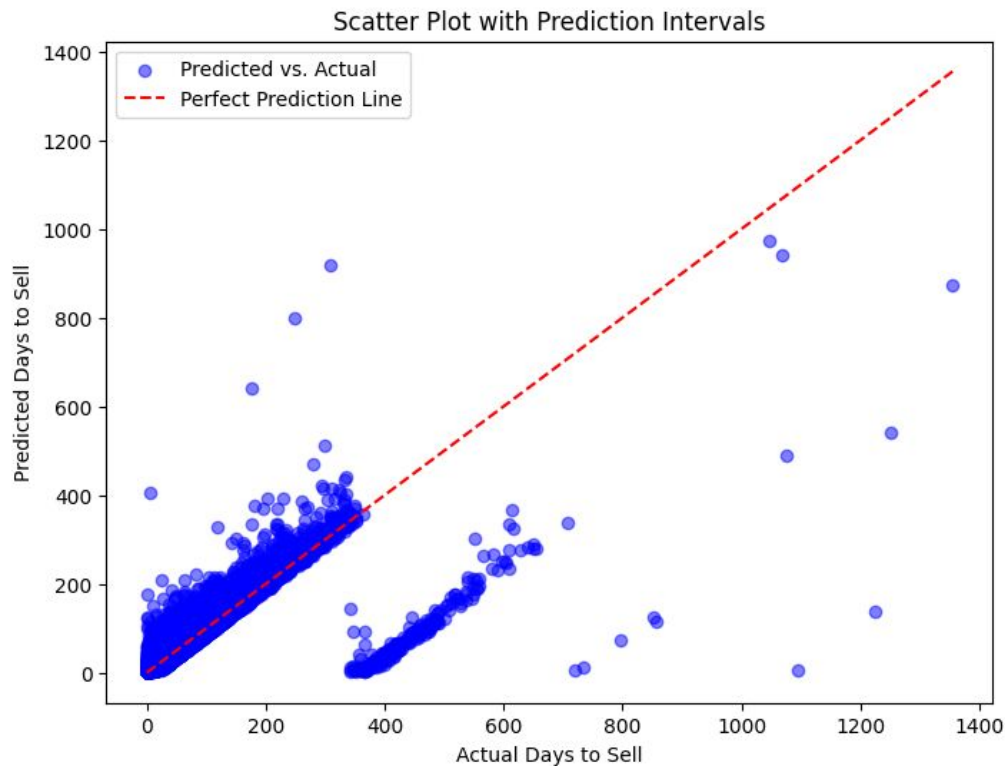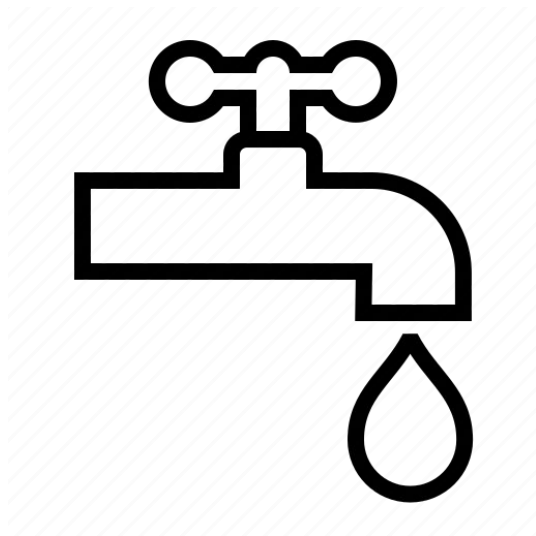
Days to Sell vs. Percentage Through Year

Days to Sell vs. Percentage Through Year

Scatter Plot with Prediction Intervals

Our model was <u>Biased</u>
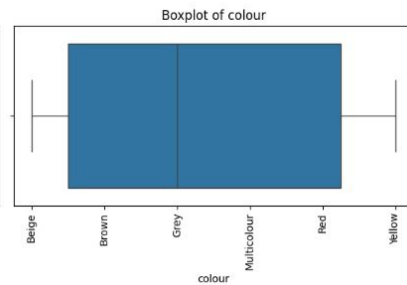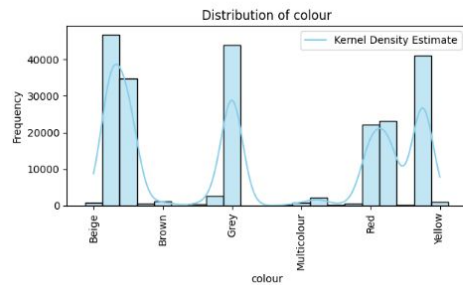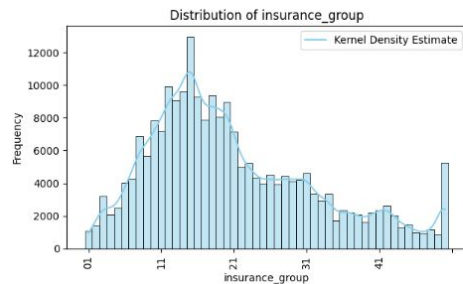
# Data Leakage



Scatter Plot with Prediction Intervals

| | stock_item_id | last_date_seen | first_date_seen | days_to_sell | first_retailer_asking_price | last_retail |
|---|---|---|---|---|---|---|
| 0 | 52ae009b671ab58b3d4ff109a9fbdcf8d847de0fa190e1... | 2023-05-05 | 2021-03-25 | 771 | 6995 | |
| 1 | 32b1bac6934b1f64ff43cffa9df5aa296ead8143c36f9f... | 2023-05-09 | 2021-05-25 | 714 | 13725 | |
| 2 | 21703d22d87eaa95c4dc81a60ba2c8cbe3b90ab659292c... | 2023-05-12 | 2021-11-26 | 532 | 15499 | |
| 3 | 661acafc271373946cea7d30ac7f34257404ab89a1ad33... | 2023-05-16 | 2022-02-17 | 453 | 10995 | |
| 4 | 638216dc92410d965b416fea5b3cec9ca903368795fdde... | 2023-05-04 | 2022-03-21 | 409 | 46000 | |
| 5 | e3c8b08856a8736bb48c38f083d42f43f3e3e8e3466610... | 2023-05-21 | 2022-04-27 | 389 | 1395 | |
| 6 | 82369d8013f2ab13f8f49fb780797298a8dd19974d3b60... | 2023-05-14 | 2022-06-06 | 342 | 8257 | |
| 7 | 1fd13f137d7ed19e993b07dd1708992582537e56efb863... | 2023-05-03 | 2022-06-16 | 321 | 23500 | |
| 8 | c34a29671d55abf60ea1ab1c23ad21a0a7437c8ffea756... | 2023-05-16 | 2022-06-23 | 327 | 96950 | |
| 9 | db6f342f73f5c7819fef4254e6886387eac15e026878ab... | 2023-05-22 | 2022-06-24 | 332 | 15995 | |

# In-Depth Data Exploration

Reviews per 100 advertised stock last 12 months

Advert quality

Top speed mph

Distribution of insurance_group

Boxplot of insurance_group

Distribution of colour

Boxplot of colour

# The Problem with Plots



Segment vs. Days to Sell

Is the `segment` feature worth including, or just **noise**?

# Statistical Tests

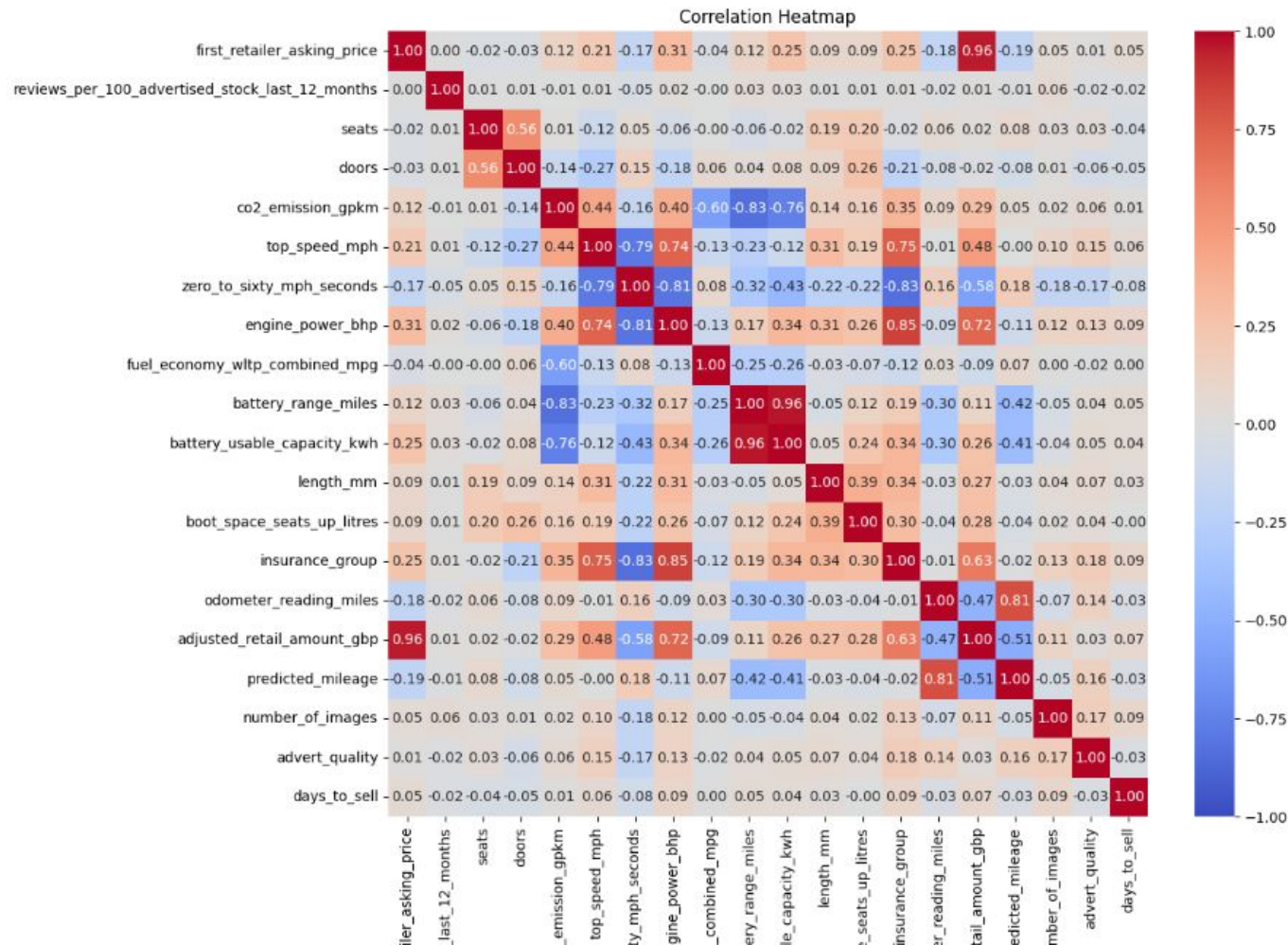*Mann-Whitney U Test* (α=0.05)
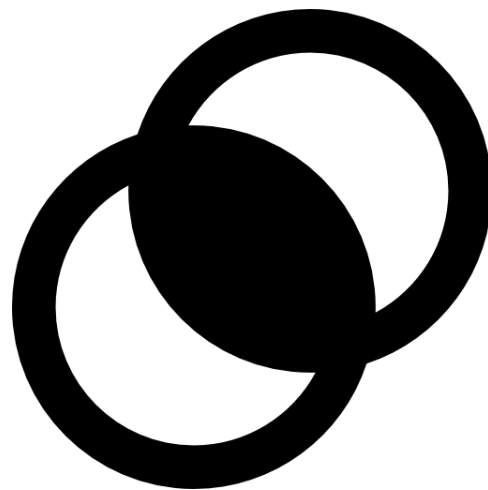
```
Mann—Whitney U test statistic: 153450447.5
P-value: 0.0192223104085375
There is a statistically significant difference in days to sell between the groups.
```
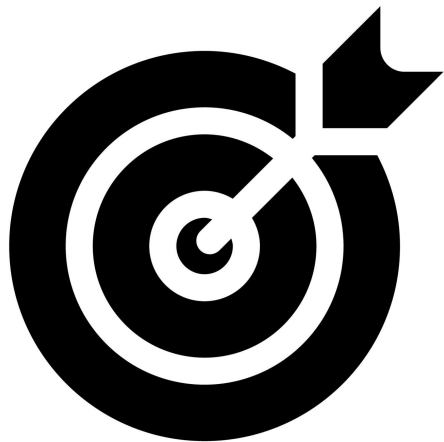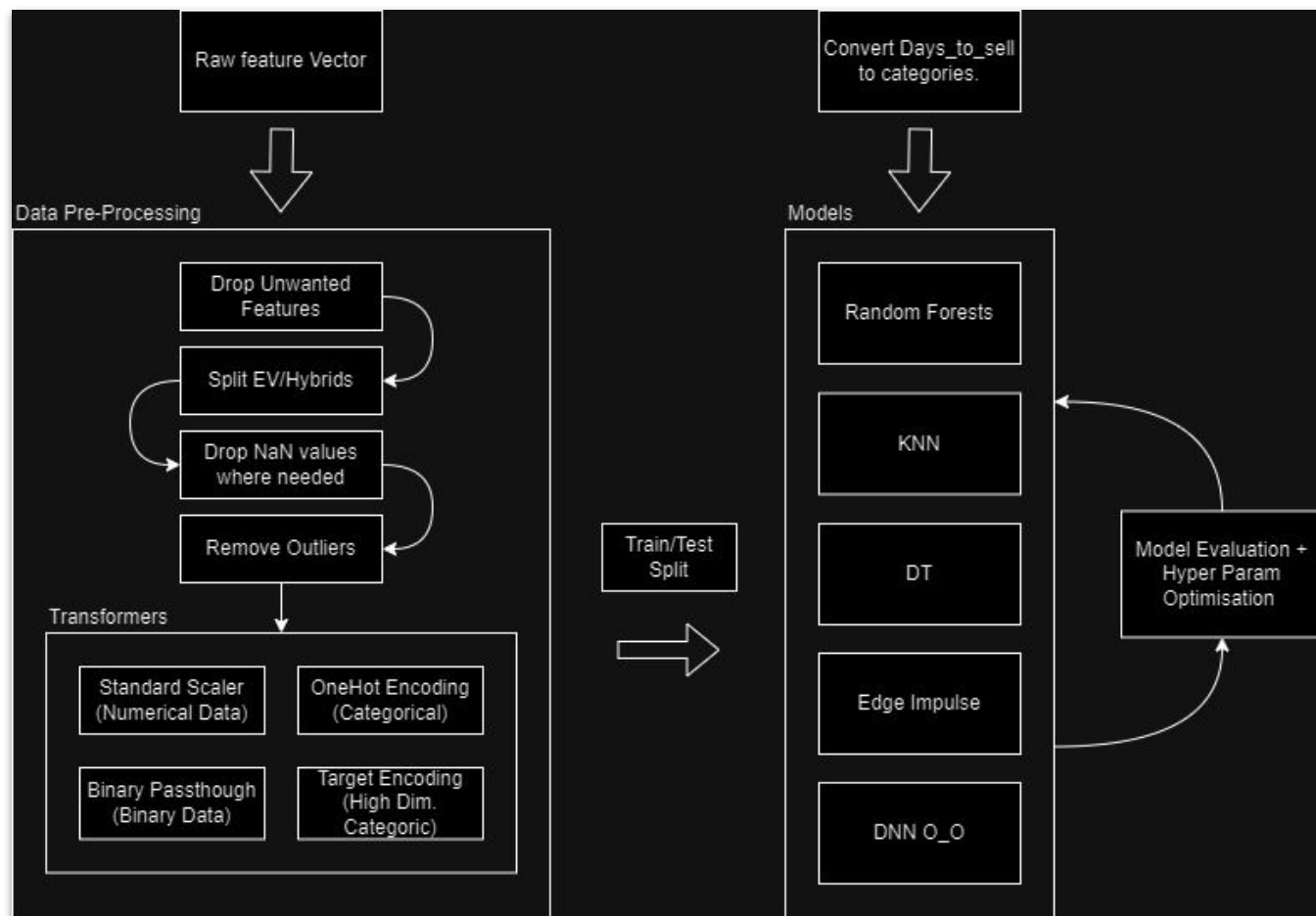
So, **include** the `segment` feature

Correlation Heatmap

*"The Curse of High Dimensionality"*

| first_retailer_asking_price - seats | first_retailer_asking_price - vehicle_age | model_count - number_of_images | model_count - odometer_reading_miles | model_count - reg_year |
|---|---|---|---|---|
| 6990.0 | 6976.0 | 1726.0 | -63224.0 | -228.0 |
| 13720.0 | 13721.0 | 7716.0 | -8287.0 | 5712.0 |
| 15494.0 | 15494.0 | 784.0 | -30287.0 | -1212.0 |
| 10990.0 | 10987.0 | 2797.0 | -76173.0 | 812.0 |
| 45995.0 | 45999.0 | 553.0 | -9639.0 | -1447.0 |
| ... | ... | ... | ... | ... |
| 11040.0 | 11038.0 | 3820.0 | -55171.0 | 1813.0 |
| 8995.0 | 8991.0 | 1154.0 | -36544.0 | -841.0 |
| 11295.0 | 11296.0 | 7727.0 | -43757.0 | 5712.0 |
| 4695.0 | 4687.0 | 26.0 | -38352.0 | -1955.0 |
| 8696.0 | 8689.0 | 2307.0 | -88513.0 | 327.0 |

# Model

Model version: ⑦ [Unoptimized (float32) ▾]

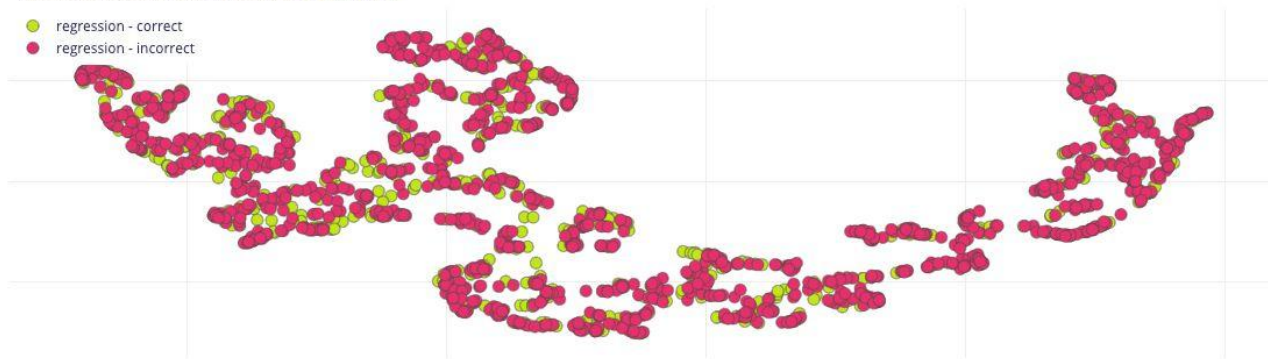## Last training performance (validation set)

📈 **LOSS**
2,759.63

## Feature explorer (full training set) ⑦

Maximum absolute regression error is 20, set thresholds.

- 🟡 regression - correct
- 🔴 regression - incorrect



## On-device performance ⑦

🕐 **INFERENCING TIME**
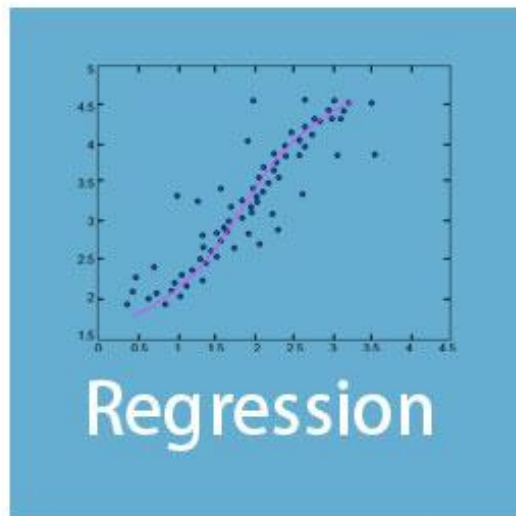1 ms.

🟧 **PEAK RAM USAGE**
1.2K

🟦 **FLASH USAGE**
10.7K

**EDGE IMPULSE**

📊 **DATA COLLECTED**
45h 54m 52s

Regression vs Classification



# Predicting Car Sale Time with Data Analytics and Machine Learning

Hamid Ahaggach, Lylia Abrouk, Sebti Foufou, Eric Lebon

▶ **To cite this version:**

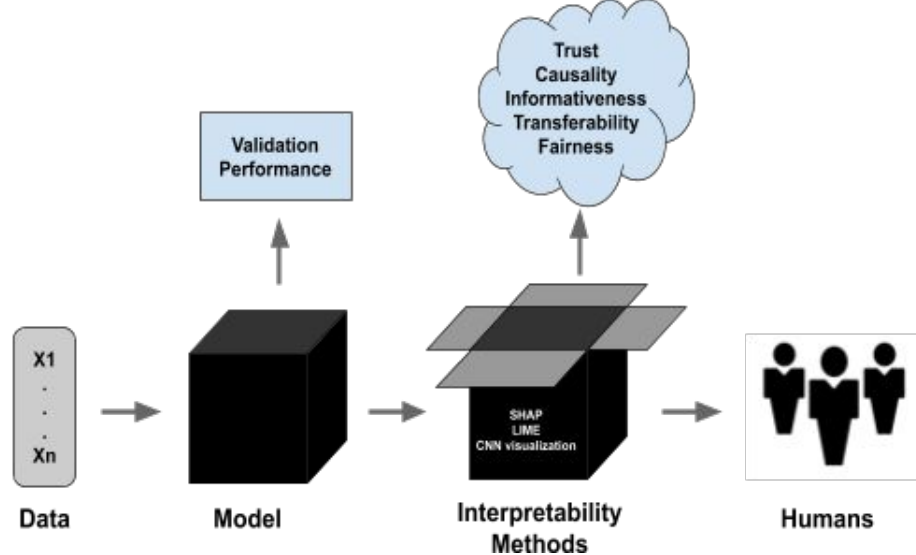**HAL Id: hal-04250878**
**https://hal.science/hal-04250878**

Submitted on 13 Nov 2023

INTERPRETABLE VS EXPLAINABLE MACHINE LEARNING

Data — Model — Interpretability Methods — Humans

Validation Performance

Trust
Causality
Informativeness
Transferability
Fairness

X1 . . . Xn

SHAP
LIME
CNN visualization

Top 10 Most Important Features

| Feature | Importance |
|---|---|
| number_of_images | ~0.34 |
| advert_quality | ~0.21 |
| price_indicator_rating_LOW | ~0.10 |
| price_indicator_rating_GREAT | ~0.09 |
| segment_Independent | ~0.09 |
| first_retailer_asking_price | ~0.08 |
| manufacturer_approved | ~0.04 |
| predicted_mileage | ~0.025 |
| can_home_deliver | ~0.01 |
| advertised_stock_last_12_months | ~0.01 |

Next 10 Most Important Features

- odometer_reading_miles
- top_speed_mph
- make_Mercedes-Benz
- boot_space_seats_up_litres
- length_mm
- co2_emission_gpkm
- price_indicator_rating_GOOD
- fuel_economy_wltp_combined_mpg
- fuel_type_Diesel
- model_Q8

Final Set of Important Features

- model_Caliber
- model_CX-7
- model_CX-5
- model_CX-30
- model_CX-3
- model_CT 200h
- model_CR-Z
- price_indicator_rating_NOANALYSIS
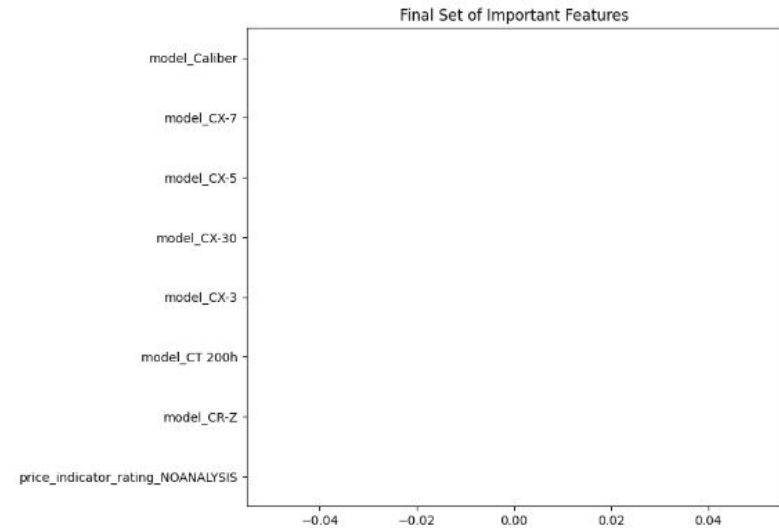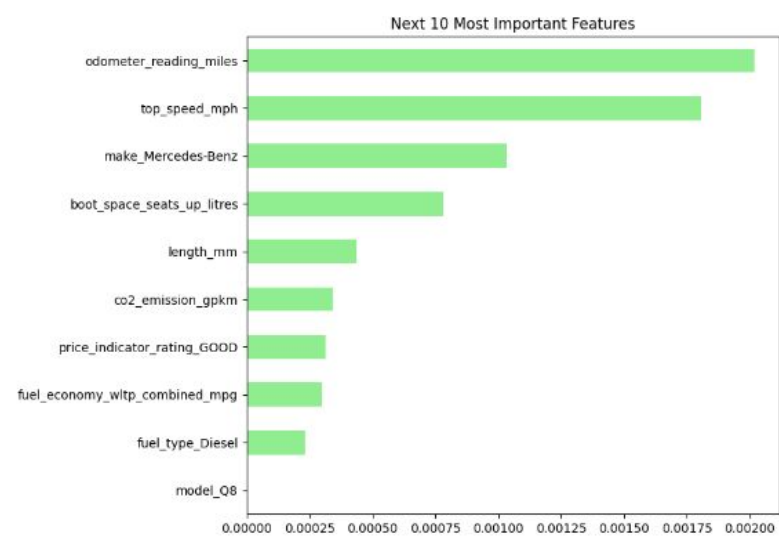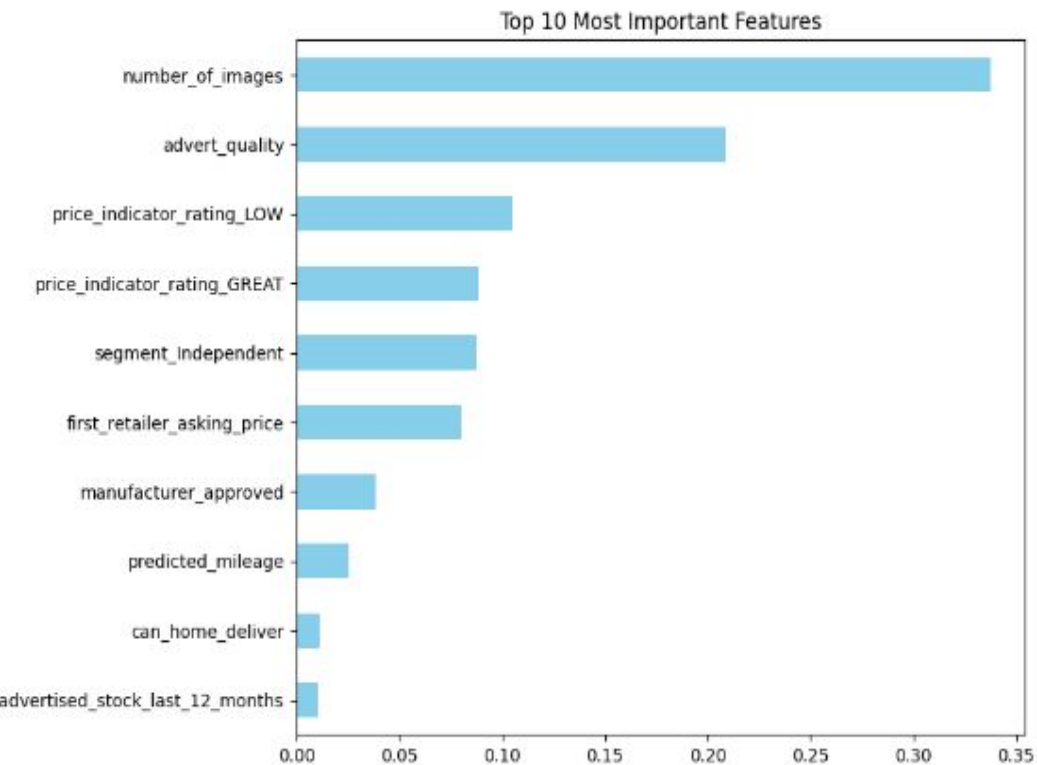
```
Decision Tree (No Optimisation)

Test RMSE:    52.77026090824235  days
 Test MAE:    30.901622579765476 days
```
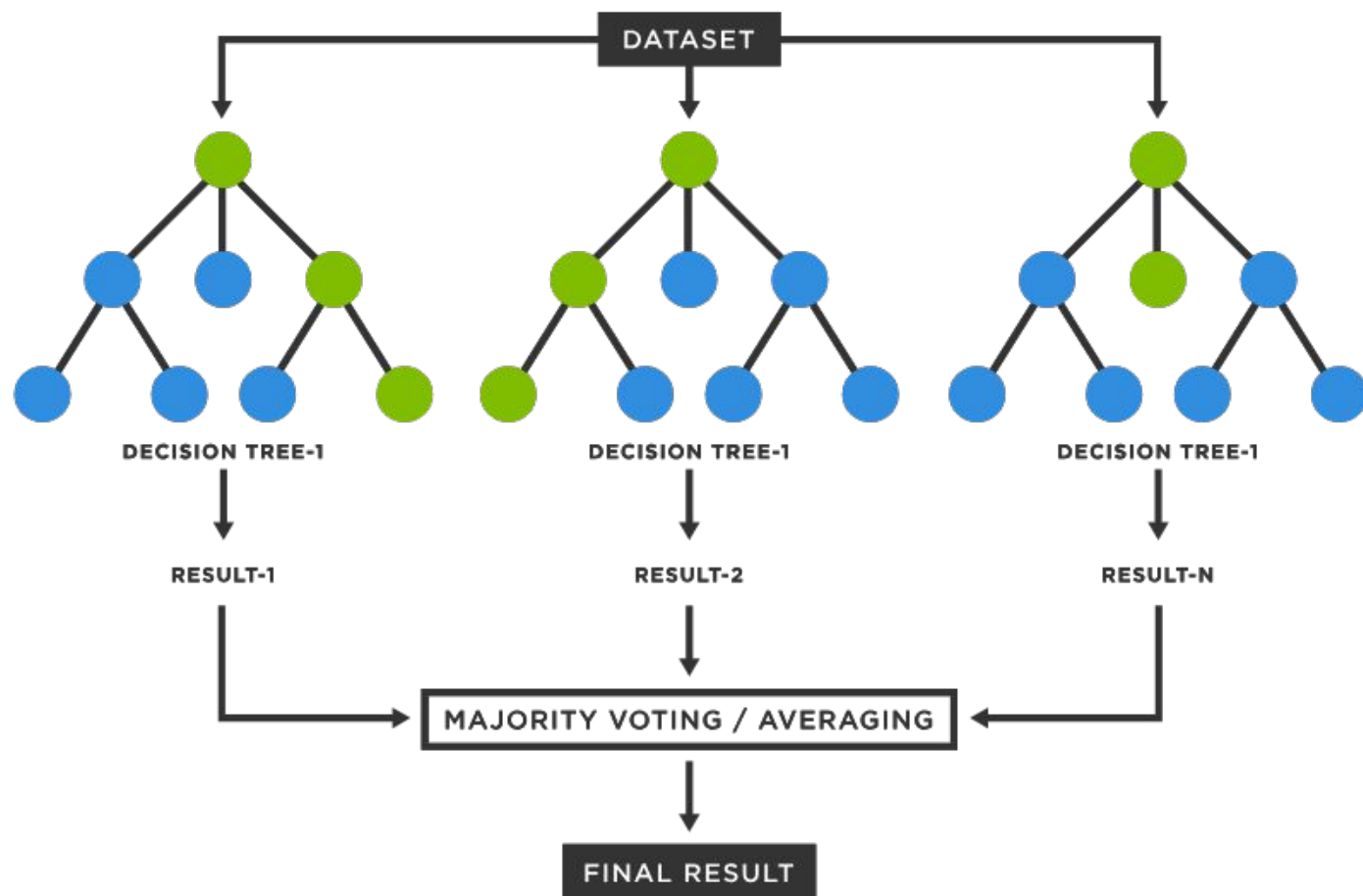
## Decision Tree (with RandomSearchCV Optimiser)

Average MSE across 5 folds: 526.2875447694037 days
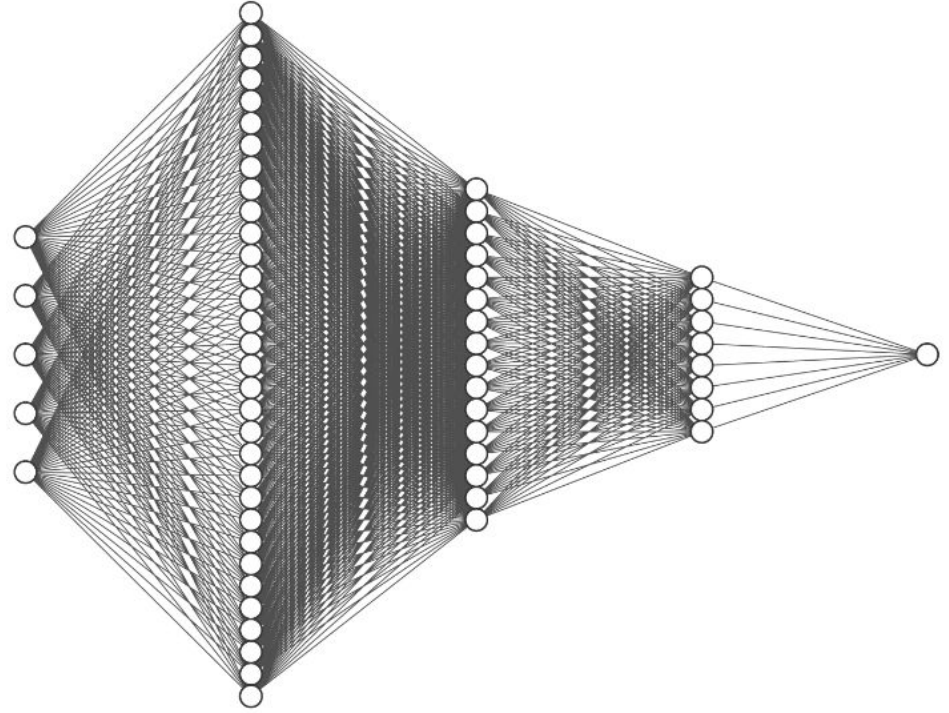Average RMSE across 5 folds: 22.940957799738957 days

## Random Forest Results

```
        Mean Absolute Error:    28.5937290542189    days
         Mean Squared Error:    1810.4967267843872   days
    Root Mean Squared Error:    42.54993215957444    days
                  R-squared:    0.14634715731726555
```
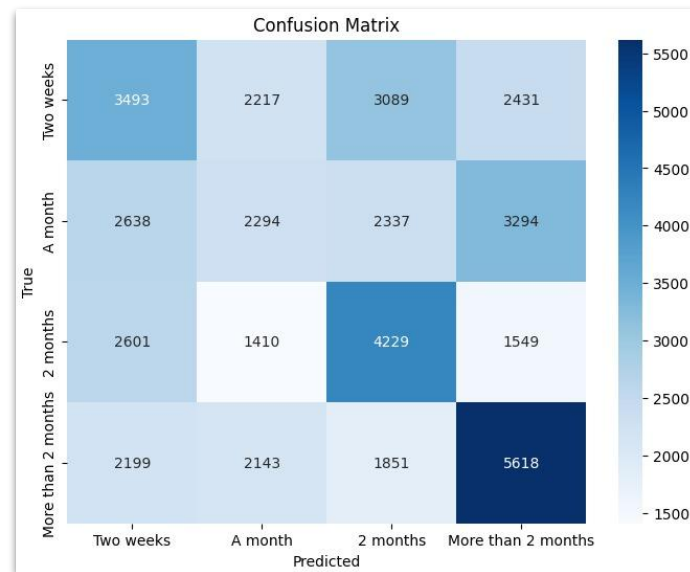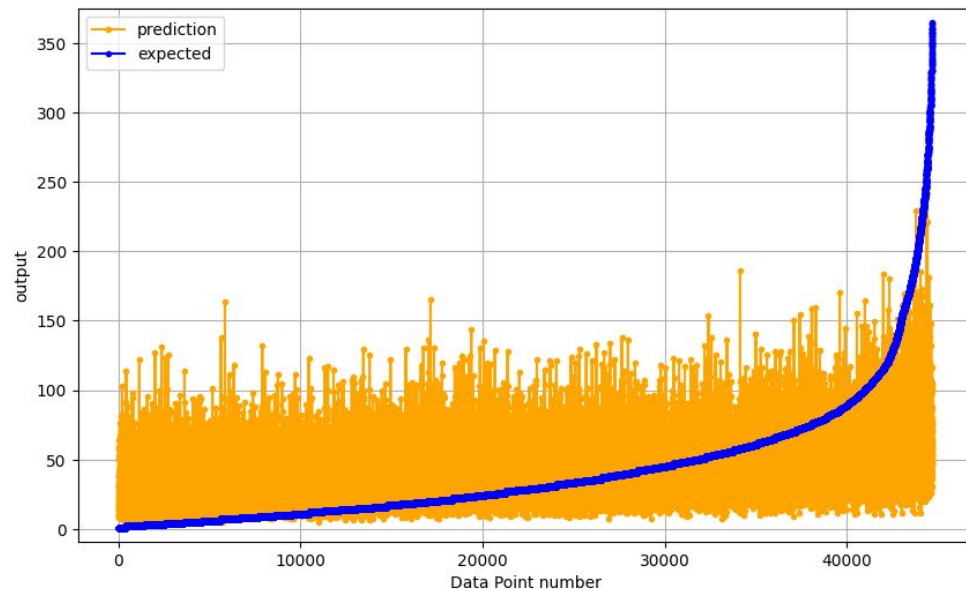
NN Approach

```
Neural Net. Model

Final score (MSE):  3412.760986328125 days
Final score (RMSE): 58.41884231567383 days
```
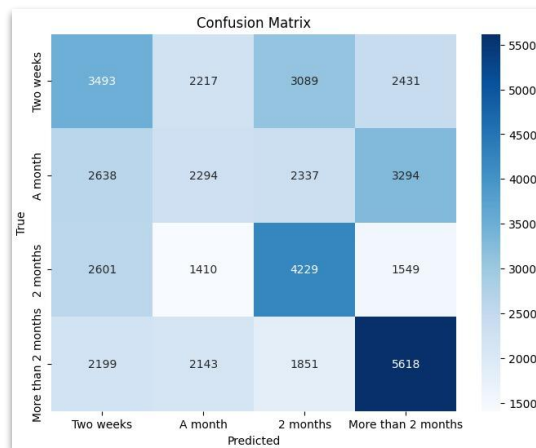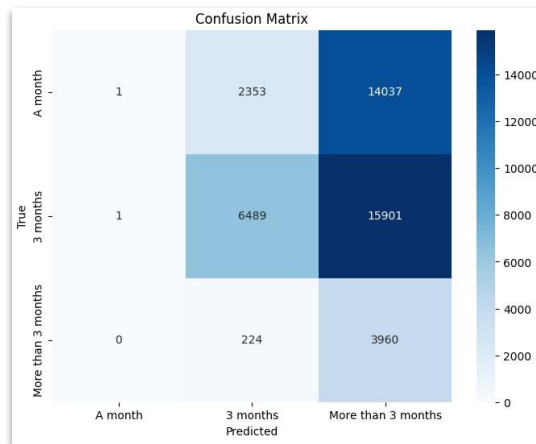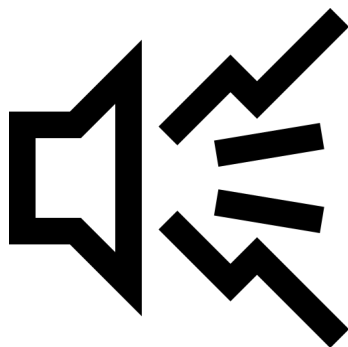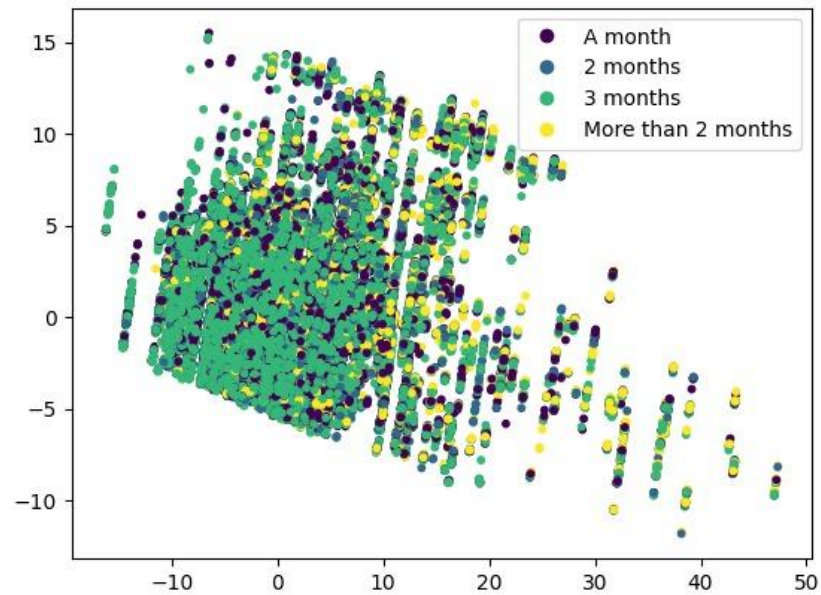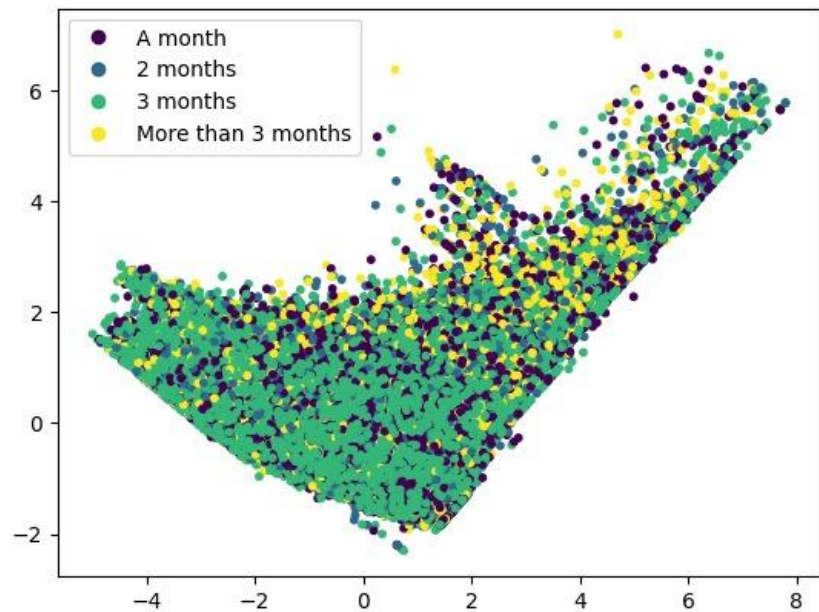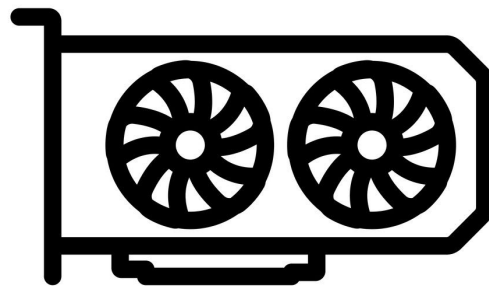
So *what* limited our solution?

*" Simply knowing what a jupyter notebook is*

*DOES NOT*

*make you an ML engineer! "*

*~ Tom Cassar, 2024*

Correlation Heatmap

## Confusion Matrix

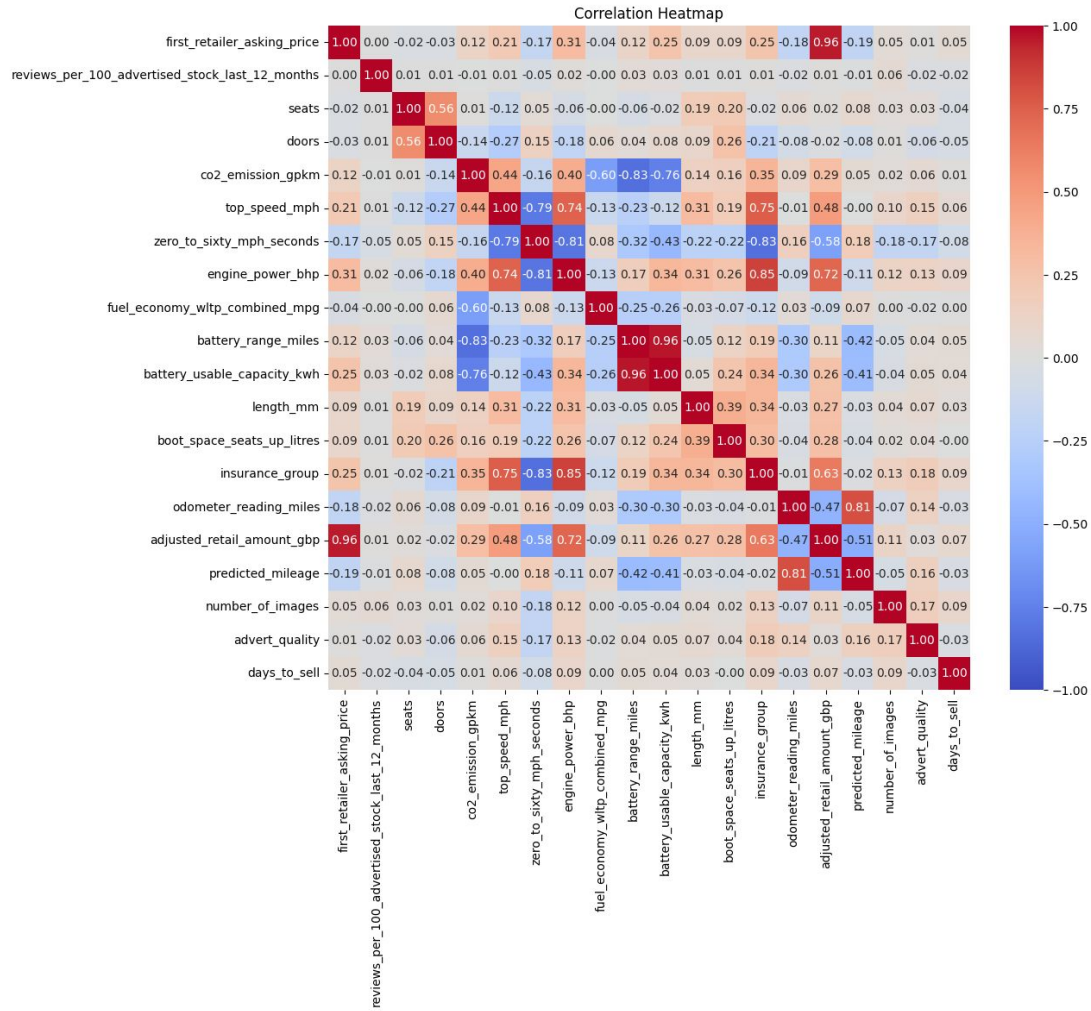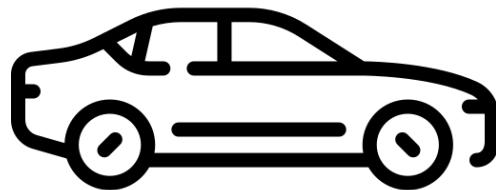|  | Less than 3 months | More than 3 months |
|---|---|---|
| Less than 3 months | 14356 | 8035 |
| More than 3 months | 7989 | 12586 |

## Random Forest:

n_estimators: 1200

min_samples_leaf: 2

Min_samples_split: 2

max_depth: 100

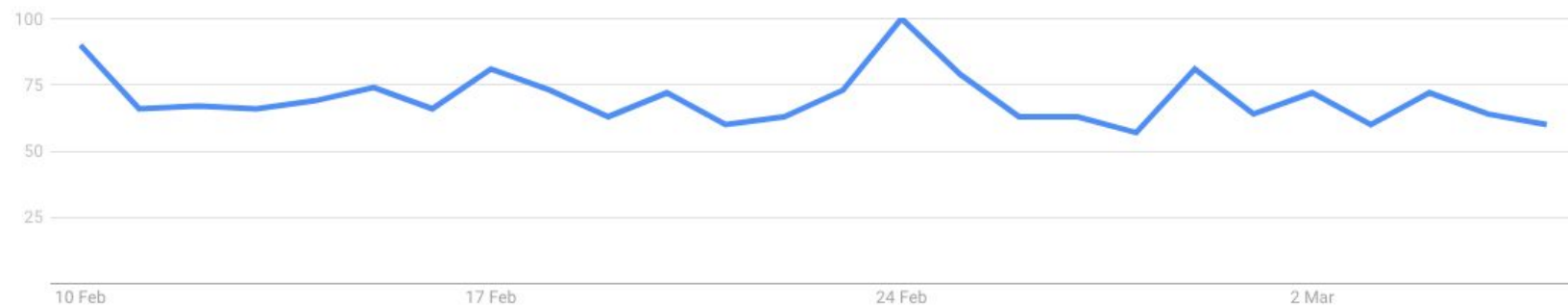bootstrap: true

Class Weights:
< 3 months: 0.1
> 3 months 0.1029

vs

# Improvements
&
Takeaways

## Interest over time ⓘ

# MathSoc AutoTrader Hackathon 2024

*Ioan Gwenter, Lourenço Silva, Tom Cassar*