

Capstone Project - 2

Appliances Energy Prediction

Kartika Sharma

- [illegible]

Problem statement

- This Dataset is speculative models driven by power consumption data. The data used includes measurements of temperature and humidity from the wireless network, weather from the nearest airport and the recorded power consumption of the lighting equipment. Filtering data to remove unpredictable parameters and feature editing is a common function of this database. From a wireless network, data from the kitchen, laundry and living room is considered very important in power forecasting. Predictability models with weather data only, have selected atmospheric pressure (corresponding to wind speed) as the flexibility of the most accurate weather data in the forecast. Therefore, atmospheric pressure may be necessary to incorporate it into power forecasting models and to create a performance model.

Data Summary

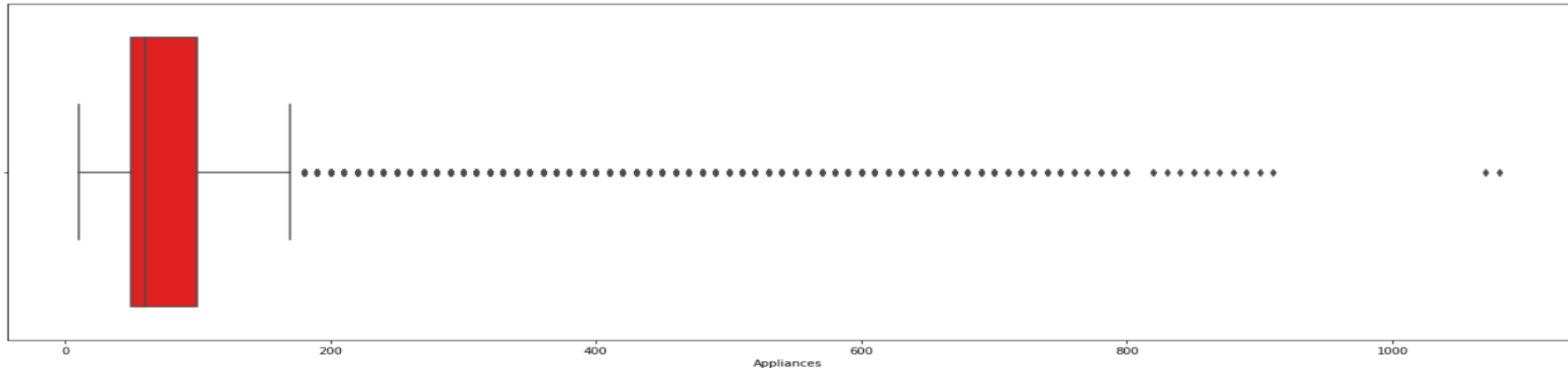
- Date - time year-month- day of energy consumption
- Appliances - energy use in WH (Dependent variable)
- Lights - energy use of light fixtures in the house in Wh (Drop this column)
- T1 - Temperature in kitchen area, in Celsius
- RH1 - Humidity in kitchen area, in %
- T2 - Temperature in living room area, in Celsius
- RH2 - Humidity in living room area, in %
- T3 - Temperature in laundry room area
- RH3 - Humidity in laundry room area, in %
- T4 - Temperature in office room, in Celsius
- RH4 - Humidity in office room, in %

Data Summary

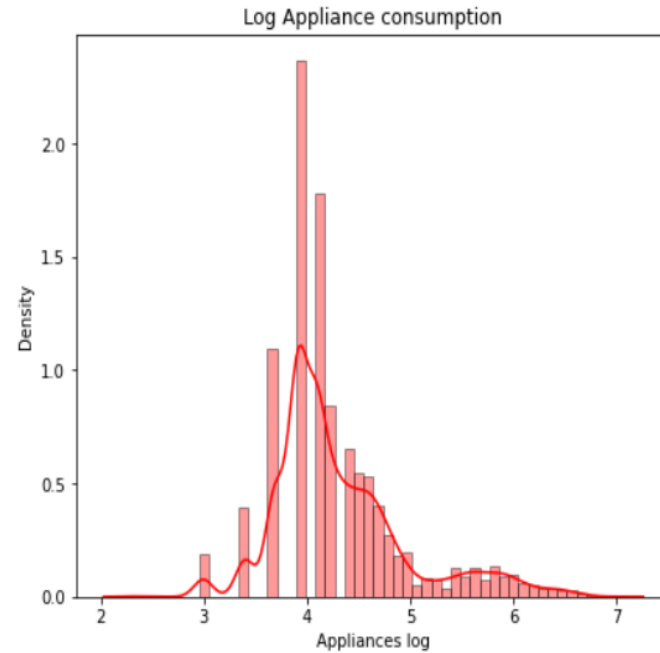
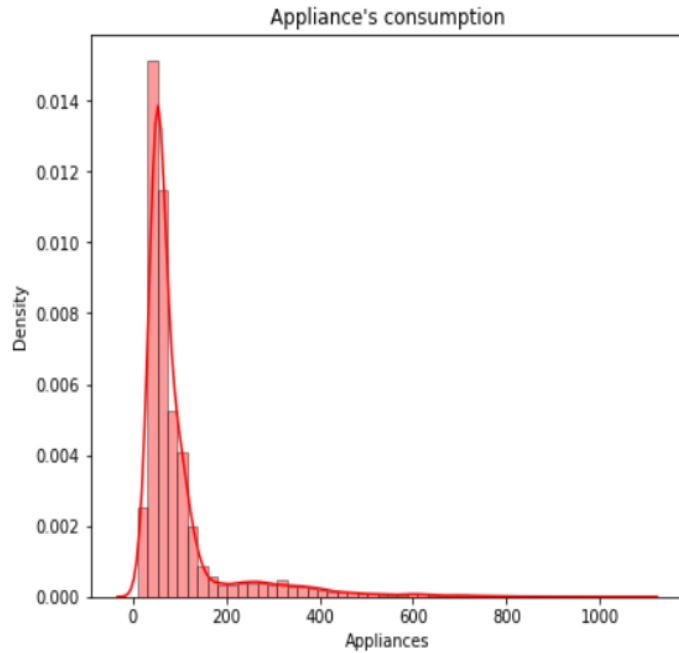
- T9 - Temperature in parent's room, in Celsius
- RH9 - Humidity in parent's room, in %
- T_out - Temperature outside (from Chievres weather station), in Celsius
- Press_mm_hg - Pressure (from Chievres weather station), in mm Hg
- Rhout - Humidity outside (from Chievres weather station), in %
- Wind speed - (from Chievres weather station), in m/s
- Visibility - (from Chievres weather station), in km
- Tdewpoint - (from Chievres weather station), $\hat{A}^{\circ}\text{C}$
- rv1 - Random variable 1, nondimensional
- rv2 - Random variable 2, nondimensional

Target Variable

- 75% of Appliance consumption is less than 100 Wh
- With the maximum consumption of 1080 WH, there will be outliers in this column and there are small number of cases where consumption is very high
- This column is positively skewed, most the values are around mean 100 Wh



Normalising Outliers...

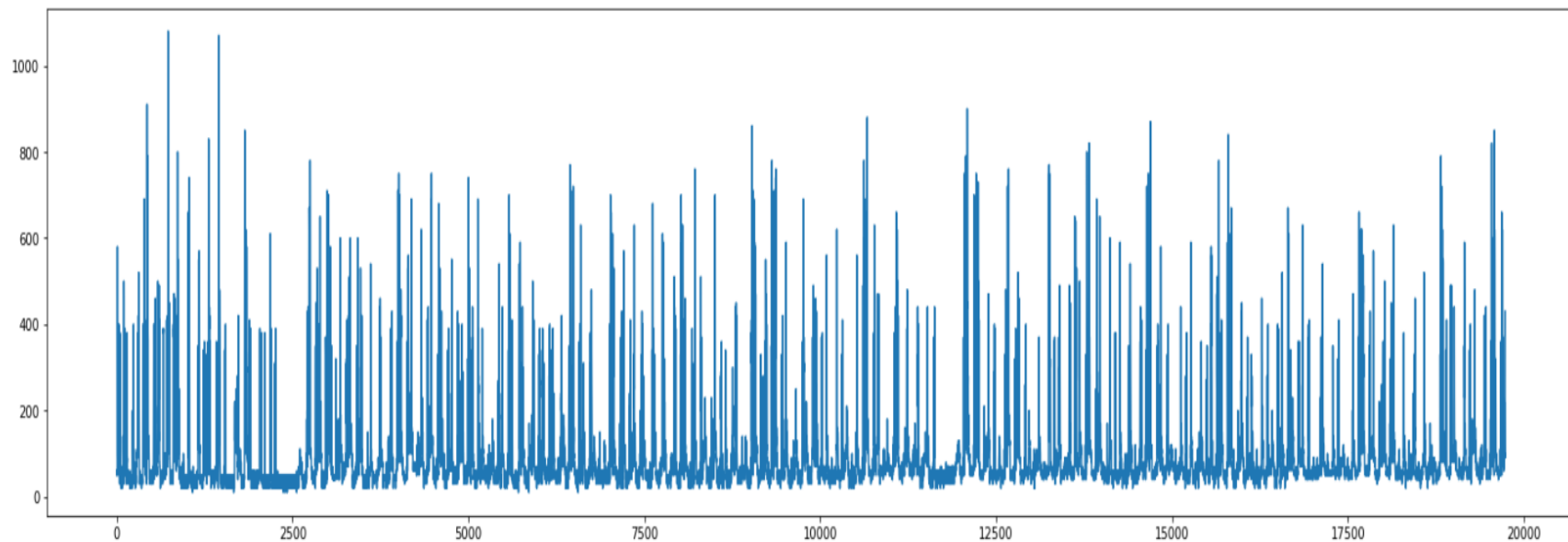


Feature extraction...

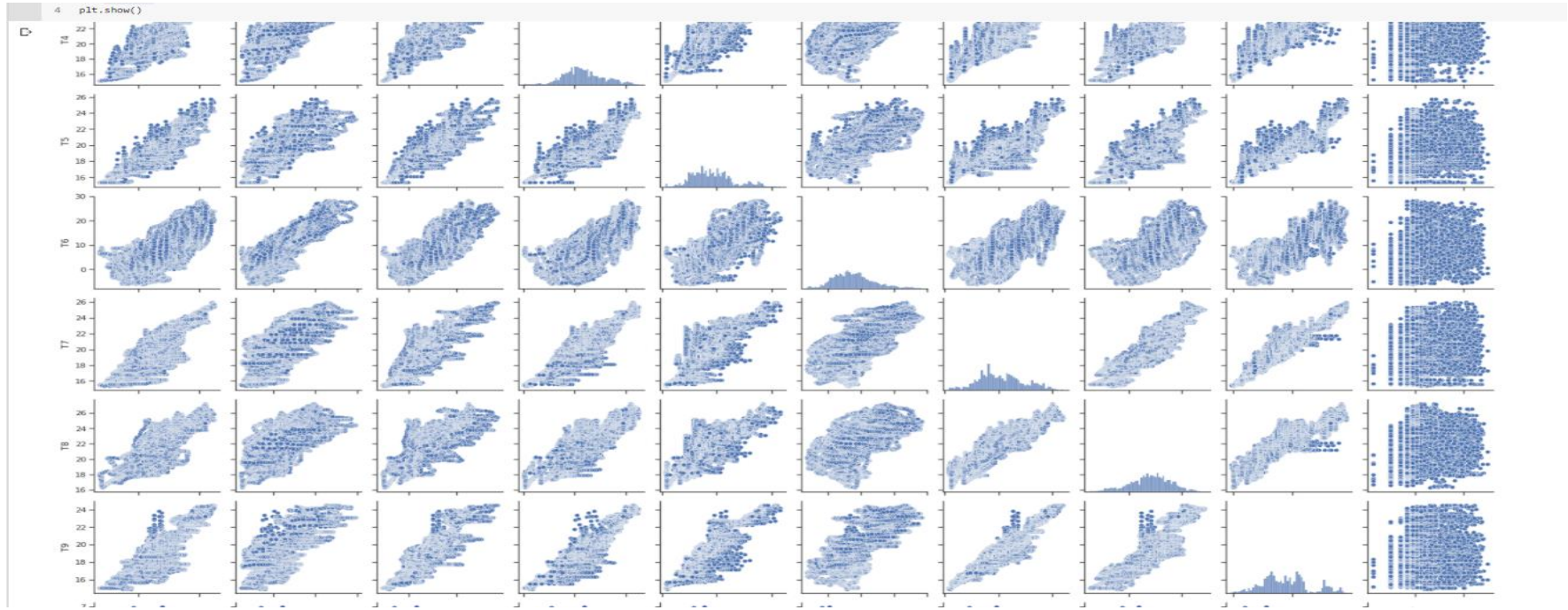
5. EXTRACTING NEW FEATURES FROM DATE COLUMN

```
1 energy_df['week_day'] = ((pd.to_datetime(energy_df['date']).dt.dayofweek) // 5 == 1).astype(float)
2 energy_df['date'] = pd.to_datetime(energy_df['date'])
3
4 energy_df['hours'] = energy_df['date'].dt.hour
5 energy_df['month'] = energy_df['date'].dt.month
6 energy_df['day'] = energy_df['date'].dt.day
7 energy_df['week_of_month'] = (energy_df['date'].dt.day // 7) + 1
```

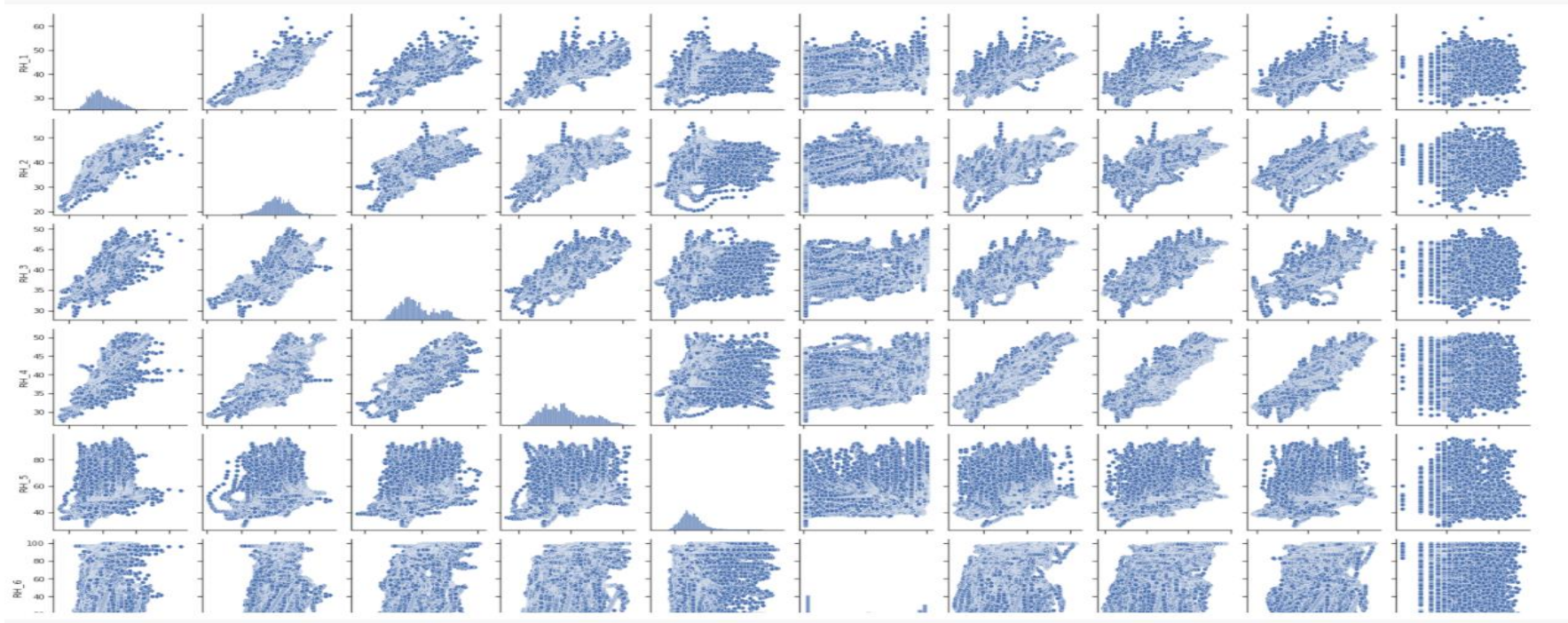

EDA (Seasonal.....I guess not)



EDA (temperatures)



EDA (Humidities)



EDA

↳ `sm.pairplot(met_weather, diag_kind='kde', corner_kde_kwargs={'bw': 0.2})`



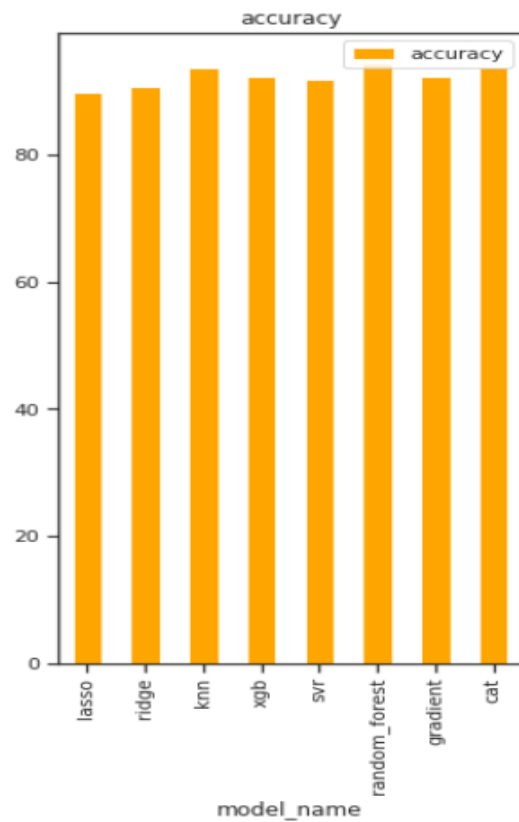
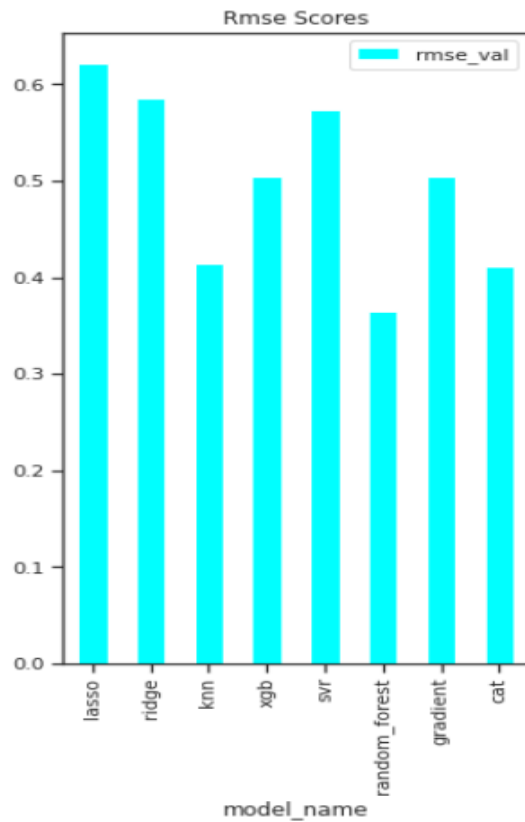
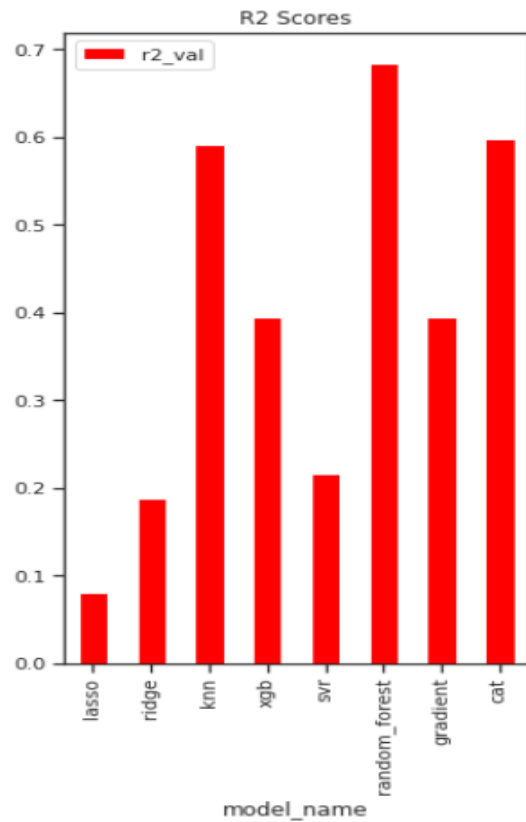
Models

- Lasso
- Ridge
- Knn
- SVM
- Random Forest
- Gradient Boosting
- Catboost

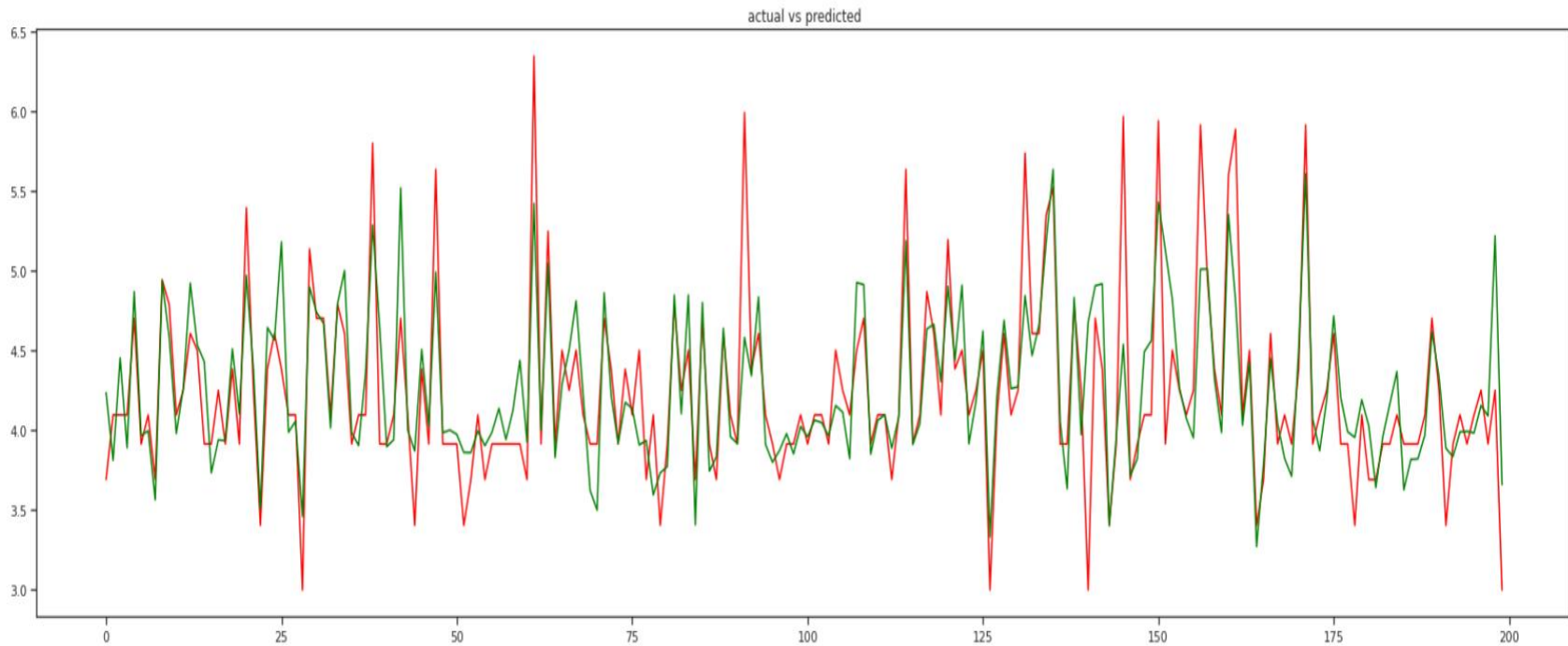
Validation and selection



	<code>model_name</code>	<code>accuracy</code>	<code>r2_val</code>	<code>rmse_val</code>
0	lasso	89.741500	0.080544	0.622252
1	ridge	90.694278	0.187442	0.584963
2	knn	93.697068	0.591920	0.414547
3	xgb	92.242783	0.394802	0.504836
4	svr	91.763680	0.216791	0.574301
5	random_forest	94.471283	0.684807	0.364325
6	gradient	92.235615	0.394756	0.504854
7	cat	93.607803	0.598229	0.411330



Actual vs Predicted



Hyperparameters

```
[ ] 1  from sklearn.model_selection import RandomizedSearchCV
    2  # Number of trees in random forest
    3  n_estimators = [int(x) for x in np.linspace(start = 50, stop = 200, num = 10)]
    4  # Number of features to consider at every split
    5  max_features = ['auto', 'sqrt']
    6  # Maximum number of levels in tree
    7  max_depth = [int(x) for x in np.linspace(10, 110, num = 11)]
    8  max_depth.append(None)
    9  # Minimum number of samples required to split a node
   10  min_samples_split = [2, 5, 10]
   11  # Minimum number of samples required at each leaf node
   12  min_samples_leaf = [1, 2, 4]
   13  # Method of selecting samples for training each tree
   14  bootstrap = [True, False]
   15  # Create the random grid
   16  random_grid = {'n_estimators': n_estimators,
   17                 |         |         |         |         |         |         |
   18                 |         |         |         |         |         |         |
   19                 |         |         |         |         |         |         |
   20                 |         |         |         |         |         |         |
   21                 |         |         |         |         |         |         |
                 'max_features': max_features,
                 'max_depth': max_depth,
                 'min_samples_split': min_samples_split,
                 'min_samples_leaf': min_samples_leaf,
                 'bootstrap': bootstrap}
```

Slight Improvement

	NAME OF MODEL	R2 SCORES	ACCURACIES	RMSE
0	random_forest	0.684807	94.471283	0.364325
1	random_forest_after_tuning	0.709293	94.724074	0.349888

Other things I tried

- Removing Outliers – Didn't work well
- Including variables with very less correlation with target variable gave less accuracy and r^2 scores
- Ridge and Lasso didn't cross 91% accuracy

Conclusion

- Final accuracy was 94.7 with random forest
- Catboost also worked well
- Randomized SearchCV had slight improvement of 0.4%
- Independent Variable selection worked well
- It can improve with further tuning

Challenges

- Feature selection – so many features and their correlation was confusing
- Computation time in randomized search CV

Q & A