

```
In [2]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from dataprep.eda import create_report
from pandas_profiling import ProfileReport
from sklearn.compose import (
    ColumnTransformer,
    TransformedTargetRegressor,
    make_column_transformer,
)
from sklearn.dummy import DummyRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.impute import SimpleImputer
from sklearn.linear_model import Ridge, RidgeCV
from sklearn.metrics import make_scorer, mean_squared_error, r2_score
from sklearn.model_selection import (
    GridSearchCV,
    cross_val_score,
    cross_validate,
    train_test_split,
)
from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.preprocessing import OneHotEncoder, OrdinalEncoder, StandardScaler
from sklearn.tree import DecisionTreeRegressor
%matplotlib inline
```

```
Pandas backend loaded 1.4.3
Numpy backend loaded 1.23.2
Pyspark backend NOT loaded
Python backend loaded
C:\Users\Nandakumar\AppData\Local\Temp\ipykernel_16952\3904670446.py:5: Deprecatio
nWarning: `import pandas_profiling` is going to be deprecated by April 1st. Please
use `import ydata_profiling` instead.
    from pandas_profiling import ProfileReport
```

```
In [3]: import warnings

warnings.simplefilter(action="ignore", category=FutureWarning)
```

```
In [4]: data=pd.read_csv('query_result.csv')
```

## Exploratory Data Analysis

```
In [5]: profile = ProfileReport(data, title="Report")
profile

Summarize dataset:  0% | 0/5 [00:00<?, ?it/s]
Generate report structure:  0% | 0/1 [00:00<?, ?it/s]
Render HTML:  0% | 0/1 [00:00<?, ?it/s]
```

# Overview

## Dataset statistics

<b>Number of variables</b>	24
<b>Number of observations</b>	136119
<b>Missing cells</b>	2948
<b>Missing cells (%)</b>	0.1%
<b>Duplicate rows</b>	0
<b>Duplicate rows (%)</b>	0.0%
<b>Total size in memory</b>	24.9 MiB
<b>Average record size in memory</b>	192.0 B

## Variable types

<b>Categorical</b>	6
<b>Numeric</b>	18

## Alerts

<code>m1_num</code> has a high cardinality: 136119 distinct values	<span>High cardinality</span>
<code>date_start</code> has a high cardinality: 686 distinct values	<span>High cardinality</span>
<code>date_end</code> has a high cardinality: 366 distinct values	<span>High cardinality</span>

Out[5]:

## Key findings:

- columns named **Ip\_dol** has high correlation with our target variable
- **yr\_built** and **garage type** has some missing values which is less than 1% of the total dataset ( missing in random category)
- lat and lon does not have direct correlation between the target column (location might not be a huge factor as of now but lets dig deeper while making them a categorical

feature and checking later.)

- We can see that id\_community,id\_municipality,ml\_num are ids we have to cast it as objects so that they dont add in the Modelling part.
- ml\_num should not be considered in the future of modelling because it has all distinct values and hence making model more ambiguous.

The date column can be really helpful in upcoming predictions so converted into the right format

```
In [6]: data['date_start'] = pd.to_datetime(data['date_start'], format='%Y-%m-%d %H:%M:%S')
data['date_end'] = pd.to_datetime(data['date_end'], format='%Y-%m-%d %H:%M:%S')
data['diff_in_dates'] = pd.to_numeric(data['date_end'] - data['date_start'])
data['id_municipality'] = data['id_municipality'].astype(object)
data['id_community'] = data['id_community'].astype(object)
```

```
In [10]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 136119 entries, 0 to 136118
Data columns (total 25 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   ml_num            136119 non-null   object 
 1   property_type     136119 non-null   object 
 2   br                134949 non-null   float64
 3   br_plus           136119 non-null   int64  
 4   br_final          136119 non-null   float64
 5   bath_tot          136119 non-null   int64  
 6   taxes              136119 non-null   float64
 7   lp_dol             136119 non-null   int64  
 8   yr_built           136117 non-null   object 
 9   gar_type           134343 non-null   object 
 10  garage              136119 non-null   float64
 11  topHighschoolScore 136119 non-null   float64
 12  topBelowHighschoolScore 136119 non-null   float64
 13  geo_latitude       136119 non-null   float64
 14  geo_longitude      136119 non-null   float64
 15  lot_frontfeet      136119 non-null   float64
 16  lot_depthfeet      136119 non-null   float64
 17  lot_size             136119 non-null   float64
 18  sqft_numeric        136119 non-null   int64  
 19  id_community         136119 non-null   object 
 20  id_municipality     136119 non-null   object 
 21  date_start          136119 non-null   datetime64[ns]
 22  date_end            136119 non-null   datetime64[ns]
 23  price_sold           136119 non-null   int64  
 24  diff_in_dates        136119 non-null   int64  
dtypes: datetime64[ns](2), float64(11), int64(6), object(6)
memory usage: 26.0+ MB
```

```
In [7]: # splitting the data to maintain the golden rule
from sklearn.model_selection import cross_val_score, cross_validate, train_test_split
train_df,test_df = train_test_split(data,test_size=0.2 ,random_state=123)
train_df.shape
```

```
Out[7]: (108895, 25)
```

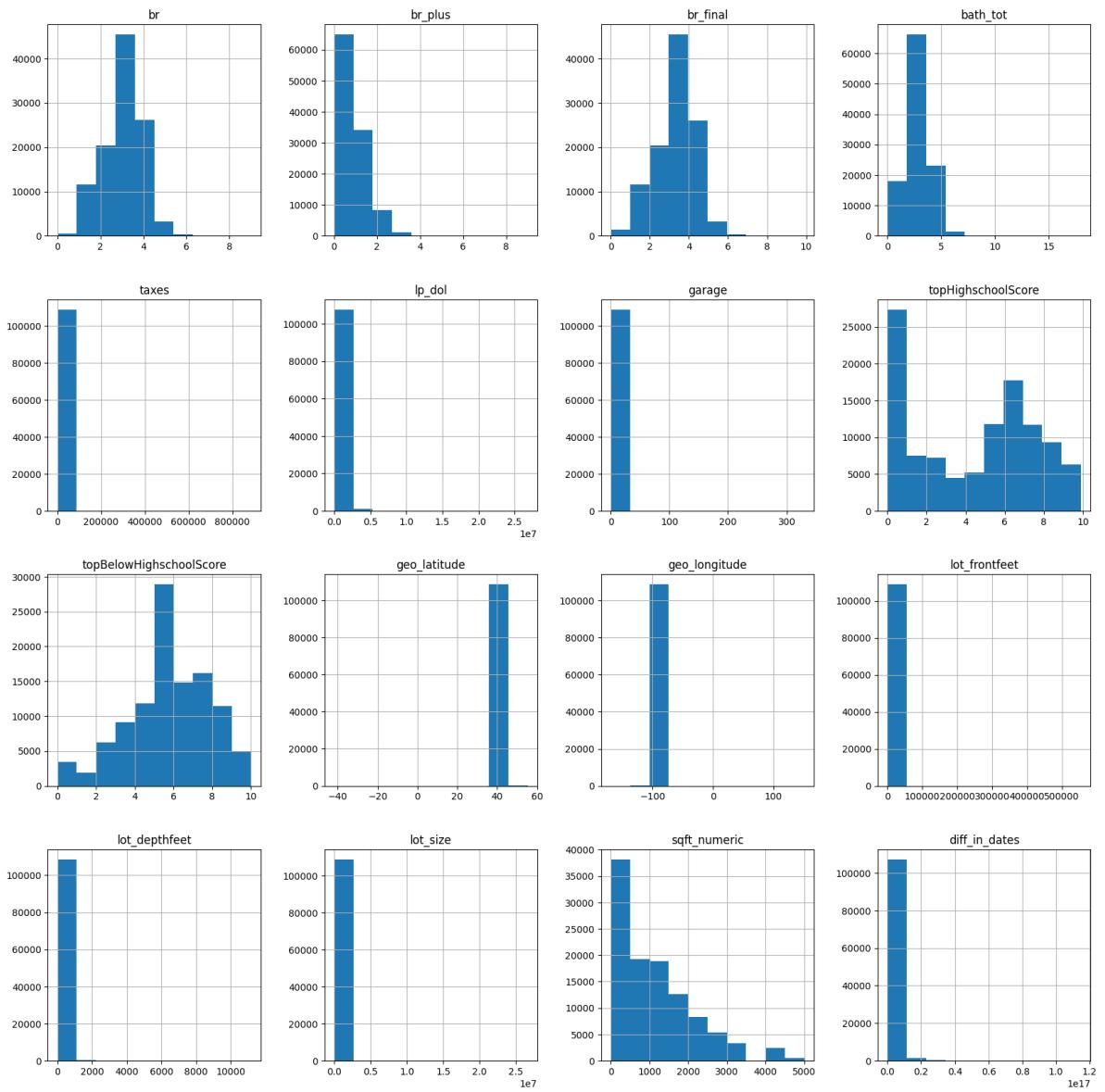
```
In [12]: #huge variability in the data these are all should be mostly closely located area in  
data['geo_latitude'].value_counts(normalize=True)
```

```
Out[12]: 43.777502    0.000984  
43.641579    0.000771  
43.641596    0.000676  
43.659329    0.000669  
43.636101    0.000654  
...  
44.403672    0.000007  
43.818745    0.000007  
43.701665    0.000007  
43.396994    0.000007  
43.894330    0.000007  
Name: geo_latitude, Length: 98896, dtype: float64
```

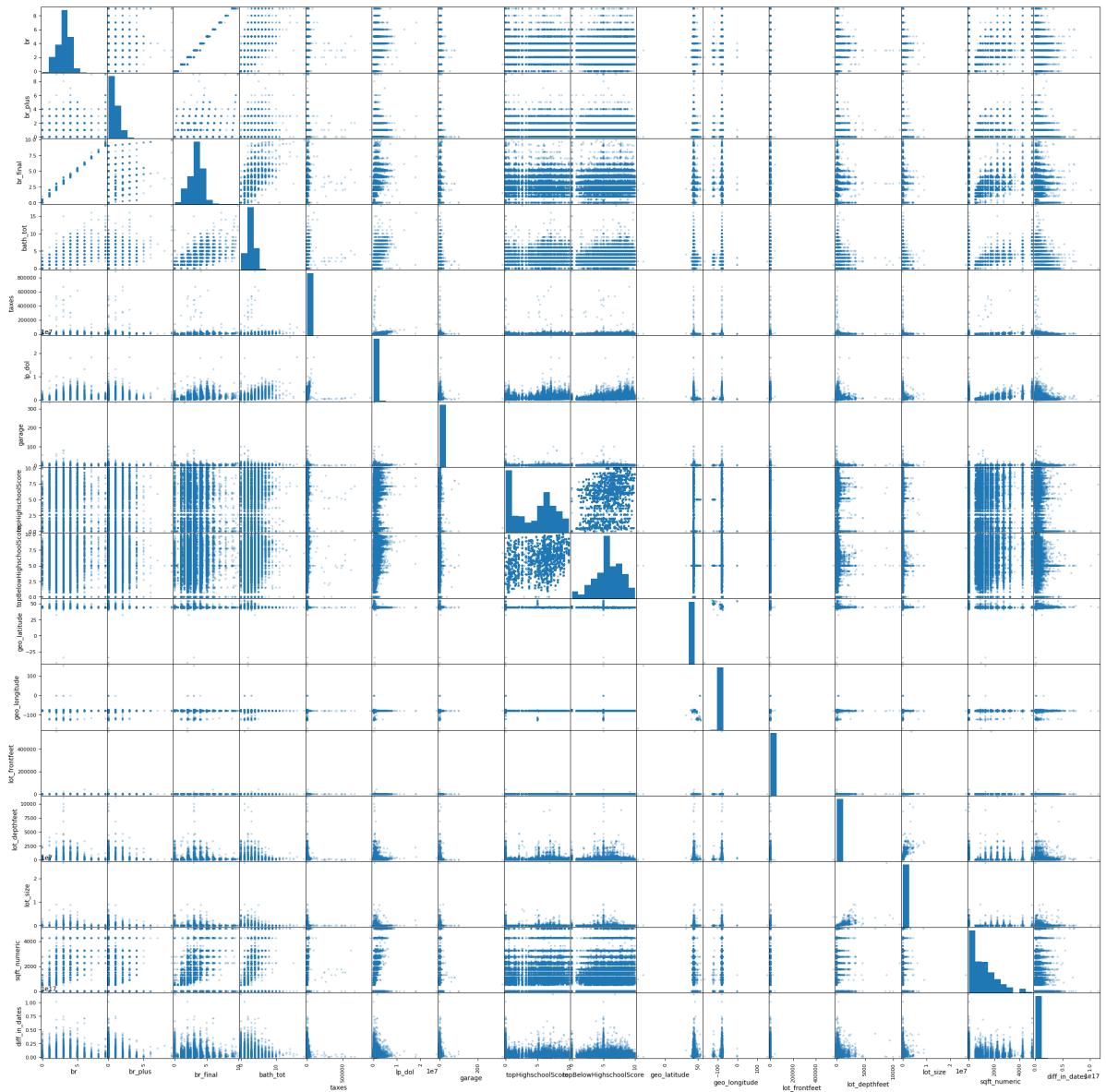
```
In [13]: train_df.describe().T
```

	count	mean	std	min	25%
<b>br</b>	107964.0	2.901078e+00	1.048421e+00	0.000000	2.000000e+00
<b>br_plus</b>	108895.0	5.094816e-01	7.061903e-01	0.000000	0.000000e+00
<b>br_final</b>	108895.0	2.964327e+00	1.042934e+00	0.000000	2.100000e+00
<b>bath_tot</b>	108895.0	2.630562e+00	1.207713e+00	0.000000	2.000000e+00
<b>taxes</b>	108895.0	3.781379e+03	6.478388e+03	0.000000	2.425000e+03
<b>lp_dol</b>	108895.0	6.735559e+05	5.135342e+05	1.000000	3.890000e+05
<b>garage</b>	108895.0	1.270563e+00	1.497978e+00	0.000000	1.000000e+00
<b>topHighschoolScore</b>	108895.0	4.352461e+00	3.189895e+00	0.000000	8.000000e-01
<b>topBelowHighschoolScore</b>	108895.0	5.651880e+00	2.137176e+00	0.000000	4.400000e+00
<b>geo_latitude</b>	108895.0	4.379830e+01	5.264347e-01	-41.872312	4.364237e+01
<b>geo_longitude</b>	108895.0	-7.950593e+01	1.650825e+00	-169.157540	-7.971425e+01
<b>lot_frontfeet</b>	108895.0	4.383721e+01	1.688723e+03	0.000000	0.000000e+00
<b>lot_depthfeet</b>	108895.0	9.012811e+01	1.526978e+02	0.000000	0.000000e+00
<b>lot_size</b>	108895.0	1.036533e+04	1.357734e+05	0.000000	0.000000e+00
<b>sqft_numeric</b>	108895.0	1.088336e+03	1.083436e+03	0.000000	0.000000e+00
<b>price_sold</b>	108895.0	6.889055e+05	5.156448e+05	1.000000	3.899000e+05
<b>diff_in_dates</b>	108895.0	1.808598e+15	2.964509e+15	0.000000	4.320000e+14

```
In [8]: import numpy as np  
import matplotlib.pyplot as plt  
%matplotlib inline  
numeric_data=train_df[['br','br_plus', 'br_final',  
                      'bath_tot', 'taxes',  
                      'lp_dol', 'garage', 'topHighschoolScore',  
                      'topBelowHighschoolScore', 'geo_latitude',  
                      'geo_longitude', 'lot_frontfeet', 'lot_depthfeet', 'lot_size'  
numeric_data.hist(figsize=(20,20));
```



```
In [13]: # we can see some correlated and some completely not corelated pattern here.
pd.plotting.scatter_matrix(numeric_data, alpha=.3, figsize=(30,30));
```



```
In [9]: #getting all the categorical data separately as they have to be treated differently
categorical_df=train_df.select_dtypes(include='object').columns
categorical_df
```

```
Out[9]: Index(['ml_num', 'property_type', 'yr_built', 'gar_type', 'id_community',
       'id_municipality'],
      dtype='object')
```

```
In [10]: from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.compose import make_column_transformer
from sklearn.preprocessing import OneHotEncoder, OrdinalEncoder, StandardScaler
```

```
In [11]: # there is one ordinal feature here year built as it has an inherent order to it
list(train_df['yr_built'].unique())
```

```
Out[11]: ['6-15',
          '0',
          '16-30',
          '0-5',
          'New',
          '31-50',
          '11-15',
          '100+',
          '51-99',
          '6-10',
          'nan']
```

```
In [12]: len(train_df['id_community'].unique())+len(train_df['property_type'].unique())+len(
```

```
Out[12]: 1658
```

```
In [13]: # as we can see below New ,0,0-5,6-15 ... up to 100+ lets create an order to it.
ordinal_feature_ordering = [
    ['New', '0', '0-5', '6-10', '6-15', '11-15', '16-30', '31-50', '51-99', '100+', 'nan']]
ordinal_features_oth = ['yr_built']
#re-initializing removing year built
categorical_features= ['property_type', 'gar_type', 'id_community',
                      'id_municipality']
# as this is id it can be dropped in the preparation step
drop_features = ['ml_num']
numeric_features = ['br', 'br_plus', 'br_final',
                     'bath_tot', 'taxes',
                     'lp_dol', 'garage', 'topHighschoolScore',
                     'topBelowHighschoolScore', 'geo_latitude',
                     'geo_longitude', 'lot_frontfeet', 'lot_depthfeet', 'lot_size'
# creating pipeline so that the same process is released in the test and dev as well
numeric_transformer = make_pipeline(SimpleImputer(strategy="median"), StandardScaler)
ordinal_transformer_oth = make_pipeline(
    SimpleImputer(strategy="most_frequent"),
    OrdinalEncoder(categories=ordinal_feature_ordering),
)
categorical_transformer = make_pipeline(
    SimpleImputer(strategy="constant", fill_value="missing"),
    OneHotEncoder(handle_unknown="ignore", sparse=False),
)
```

```
In [14]: X_train = train_df.drop(columns=['price_sold'])
X_test = test_df.drop(columns=['price_sold'])
y_train = train_df['price_sold']
y_test = test_df['price_sold']
```

```
In [15]: #preprocessor defined to handle all types of data
preprocessor = make_column_transformer(
    ("drop", drop_features),
    (numeric_transformer, numeric_features),
    (ordinal_transformer_oth, ordinal_features_oth),
    (categorical_transformer, categorical_features),
)
```

```
In [16]: preprocessor
```

```
Out[16]: ColumnTransformer(transformers=[('drop', 'drop', ['ml_num']),
                                         ('pipeline-1',
                                          Pipeline(steps=[('simpleimputer',
                                                          SimpleImputer(strategy='median'))]),
                                          ('standardscaler',
                                           StandardScaler()))]),
                                         ['br', 'br_plus', 'br_final', 'bath_tot',
                                          'taxes', 'lp_dol', 'garage',
                                          'topHighschoolScore',
                                          'topBelowHighschoolScore', 'geo_latitude',
                                          'geo_longitude', 'lot_frontfeet',
                                          'lot_depthfeet...'],
                                         OrdinalEncoder(categories=[[ 'New',
                                          '0',
                                          '0-5',
                                          '5-10',
                                          '10-15',
                                          '15-20',
                                          '20-25',
                                          '25-30',
                                          '30-35',
                                          '35-40',
                                          '40-45',
                                          '45-50',
                                          '50-55',
                                          '55-60',
                                          '60-65',
                                          '65-70',
                                          '70-75',
                                          '75-80',
                                          '80-85',
                                          '85-90',
                                          '90-95',
                                          '95-100',
                                          '+'],
                                          'na',
                                          'n']])))],
                                         ['yr_built']),
                                         ('pipeline-3',
                                          Pipeline(steps=[('simpleimputer',
                                                          SimpleImputer(fill_value='missing',
                                                          strategy='constant')),
                                                          ('onehotencoder',
                                                           OneHotEncoder(handle_unknown='ignore',
                                                           sparse=False))]),
                                          ['property_type', 'gar_type', 'id_community',
                                           'id_municipality']))]
```

```
In [17]: preprocessor.fit(X_train) # Calling fit to examine all the transformers.
preprocessor.named_transformers_
```

```
Out[17]: {'drop': 'drop',
  'pipeline-1': Pipeline(steps=[('simpleimputer', SimpleImputer(strategy='median')),
                                ('standardscaler', StandardScaler())]),
  'pipeline-2': Pipeline(steps=[('simpleimputer', SimpleImputer(strategy='most_frequent')),
                                ('ordinalencoder',
                                 OrdinalEncoder(categories=[[['New', '0', '0-5', '6-10', '6-15',
                                                               '11-15', '16-30', '31-50', '51-99',
                                                               '100+', 'nan']]]))],
  'pipeline-3': Pipeline(steps=[('simpleimputer',
                                 SimpleImputer(fill_value='missing', strategy='constant')),
                                ('onehotencoder',
                                 OneHotEncoder(handle_unknown='ignore', sparse=False))]),
  'remainder': 'drop'}
```

```
In [18]: ohe_columns = list(
    preprocessor.named_transformers_["pipeline-3"]
    .named_steps["onehotencoder"]
    .get_feature_names_out(categorical_features)
)
new_columns = (
    numeric_features + ordinal_features_oth + ohe_columns
)
```

```
In [19]: X_train_enc = pd.DataFrame(
    preprocessor.transform(X_train), index=X_train.index, columns=new_columns
)
X_train_enc.head()
```

Out[19]:

	<b>br</b>	<b>br_plus</b>	<b>br_final</b>	<b>bath_tot</b>	<b>taxes</b>	<b>lp_dol</b>	<b>garage</b>	<b>topHighschools</b>
<b>32534</b>	1.051833	-0.721454	0.993042	1.133916	0.089617	-0.182376	0.486950	-1.36
<b>51293</b>	-0.863941	-0.721454	-0.924633	-0.522115	-0.172718	-0.456751	-0.180620	-1.20
<b>59241</b>	-0.863941	0.694601	-0.445214	0.305901	-0.112186	0.051300	1.822089	0.98
<b>50702</b>	-1.821828	0.694601	-1.404052	-1.350130	-0.210964	-0.591504	-0.180620	0.70
<b>37089</b>	-0.863941	0.694601	-0.828749	-1.350130	-0.027998	-0.454998	-0.180620	-0.11

5 rows × 1675 columns

In [20]: X\_train.shape , X\_train\_enc.shape  
# there is jump of 1651 features

Out[20]: ((108895, 24), (108895, 1675))

## Model Building

1. Three machine learning models such as,
  - Regression based on k-nearest neighbors.
  - Ridge: Linear least squares with L2 regularization.
  - Random Forest Regressor: To improve score and avoid over-fitting.

```
In [30]: from sklearn.linear_model import Ridge
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
models = {
    "KNN_Reg": KNeighborsRegressor(),
    "Ridge": Ridge(),
    "RandomForest": RandomForestRegressor()
}
score_types_reg = {
    "neg_mean_squared_error": "neg_mean_squared_error",
    "neg_root_mean_squared_error": "neg_root_mean_squared_error",
    "neg_mape": "neg_mean_absolute_percentage_error",
    "r2": "r2",
}
```

```
In [32]: cross_val_results={}
for i in models:
    print(models[i])
    pipe=make_pipeline(preprocessor,models[i])
    cross_val_results[i] = pd.DataFrame(
        cross_validate(pipe, X_train, y_train, cv=2, return_train_score=True,
                      scoring=score_types_reg)).agg(['mean','std']).round(3).T
pd.concat(
    cross_val_results,
    axis='columns'
).xs(
    'mean',
    axis='columns',
    level=1
).style.format(
    precision=2
).background_gradient(
    axis=None
)
```

KNeighborsRegressor()  
Ridge()  
RandomForestRegressor()

```
Out[32]:
```

	KNN_Reg	Ridge	regr
<b>fit_time</b>	1.60	3.00	432.32
<b>score_time</b>	178.39	0.97	3.75
<b>test_neg_mean_squared_error</b>	-27204765429.02	-11801765631.02	-18768213966.71
<b>train_neg_mean_squared_error</b>	-17184316405.79	-10783325768.32	-937074347.40
<b>test_neg_root_mean_squared_error</b>	-163381.65	-103990.29	-136862.86
<b>train_neg_root_mean_squared_error</b>	-129525.85	-97981.02	-29882.99
<b>test_neg_mape</b>	-6.18	-6.81	-156.94
<b>train_neg_mape</b>	-6.90	-7.43	-2.89
<b>test_r2</b>	0.90	0.96	0.93
<b>train_r2</b>	0.94	0.96	1.00

- Random Forest regressor is Overfitting. Ridge regressor is working well for house prediction usecase based on Test\_r2 (r squared)

# Hyperparameter Tuning using Random Search Cross validation.

```
In [21]: from scipy.stats import loguniform
from sklearn.model_selection import RandomizedSearchCV
param_dist = {"ridge_alpha": 10.0 ** np.arange(-5, 5, 1)}
pipe = make_pipeline(preprocessor, Ridge())
search = RandomizedSearchCV(pipe, param_dist, return_train_score=True, cv=5, n_jobs=-1)
search.fit(X_train, y_train)
best_parameter = search.best_params_
best_parameter
```

```
Out[21]: {'ridge_alpha': 100.0}
```

```
In [24]: # Top 10 Important features interpretaion for numerical feature
coef=pd.DataFrame(best_model.named_steps.ridge.coef_,X_train_enc.columns,columns=[]
coef = coef.sort_values('importance', ascending=False).iloc[:10,:]
coef.style.format(
    precision=2
).background_gradient(
    axis=None
)
```

```
Out[24]:
```

	importance
lp_dol	468780.70
id_community_111	141065.02
id_municipality_10343	91293.35
id_community_705	89905.37
id_municipality_10234	79553.05
id_community_22	64720.36
id_community_42	61880.76
property_type_V.	60739.83
id_municipality_10280	52383.73
id_community_45	50727.03

```
In [28]: coef1=pd.DataFrame(best_model.named_steps.ridge.coef_,X_train_enc.columns)
```

```
In [36]: lr_coefs_property = coef1[coef1.index.str.startswith("property_type")]
lr_coefs_property
```

Out[36]:

0

<b>property_type_3.</b>	-6766.317240
<b>property_type_6.</b>	10162.402439
<b>property_type_9.</b>	-4431.448035
<b>property_type_@.</b>	-21069.249418
<b>property_type_A.</b>	525.650429
<b>property_type_C.</b>	-30942.148531
<b>property_type_D.</b>	24461.074285
<b>property_type_E.</b>	-20172.230901
<b>property_type_F.</b>	9814.998435
<b>property_type_G.</b>	12037.938570
<b>property_type_H.</b>	2186.819543
<b>property_type_J.</b>	1364.923725
<b>property_type_K.</b>	-6207.660579
<b>property_type_L.</b>	4656.073457
<b>property_type_M.</b>	-8851.487632
<b>property_type_N.</b>	-8009.714570
<b>property_type_O.</b>	-69.611331
<b>property_type_P.</b>	2289.727439
<b>property_type_R.</b>	11686.783217
<b>property_type_S.</b>	8207.671707
<b>property_type_T.</b>	-17472.892593
<b>property_type_V.</b>	60739.833706
<b>property_type_W.</b>	-27517.008843
<b>property_type_X.</b>	12978.909151
<b>property_type_Y.</b>	-9867.621679
<b>property_type_Z.</b>	264.585248

In [49]:

```
# Interpretation for categorical variable  
lr_coefs_proprty=lr_coefs_proprty.loc['property_type_A. ']
```

Out[49]:

0

<b>property_type_3.</b>	-7291.967670
<b>property_type_6.</b>	9636.752010
<b>property_type_9.</b>	-4957.098465
<b>property_type_@.</b>	-21594.899847
<b>property_type_A.</b>	0.000000
<b>property_type_C.</b>	-31467.798961
<b>property_type_D.</b>	23935.423856
<b>property_type_E.</b>	-20697.881330
<b>property_type_F.</b>	9289.348006
<b>property_type_G.</b>	11512.288141
<b>property_type_H.</b>	1661.169114
<b>property_type_J.</b>	839.273295
<b>property_type_K.</b>	-6733.311009
<b>property_type_L.</b>	4130.423028
<b>property_type_M.</b>	-9377.138061
<b>property_type_N.</b>	-8535.364999
<b>property_type_O.</b>	-595.261761
<b>property_type_P.</b>	1764.077010
<b>property_type_R.</b>	11161.132788
<b>property_type_S.</b>	7682.021278
<b>property_type_T.</b>	-17998.543023
<b>property_type_V.</b>	60214.183277
<b>property_type_W.</b>	-28042.659272
<b>property_type_X.</b>	12453.258722
<b>property_type_Y.</b>	-10393.272108
<b>property_type_Z.</b>	-261.065181

## Evaluating best model (Ridge with best parameter) using test dataset

In [23]:

```
best_model=search.best_estimator_
best_model.score(X_test,y_test)
```

Out[23]: 0.8129602321005545

- Model is achieving **81.3% R-squared value**. Which is a really good fit model for house price prediction data.

```
In [53]: import pickle
pickle.dump(best_model, open('Random_forest_regressor_model.pkl','wb'))
```

## Inference Pipeline: Predict house price prediction for new input

```
In [56]: # data=pd.read_csv('query_result.csv')
input_house_data = data.iloc[0,:]
new = pd.DataFrame(input_house_data).T
new['date_start'] = pd.to_datetime(new['date_start'], format='%Y-%m-%d %H:%M:%S')
new['date_end'] = pd.to_datetime(new['date_end'], format='%Y-%m-%d %H:%M:%S')
new['diff_in_dates'] =pd.to_numeric(new['date_end']- new ['date_start'])
new['id_municipality'] =new['id_municipality'].astype(object)
new['id_community'] = new['id_community'].astype(object)
new = new.drop(columns=['price_sold'])
```

```
In [55]: #predicted value for new input
model=pickle.load(open("Random_forest_regressor_model.pkl", "rb"))
print(model.predict(new))

[401622.83103571]
```

```
In [57]: #Actual value for the new input
input_house_data['price_sold']
```

```
Out[57]: 406400
```

**Model is able to predict the house price prediction using Ridge Regressor**

## Model deployment in Fast API

1. **app.py**: This is the main file for receiving required information for house price prediction through GUI or API calls and computing the predicted house price and returning it.
2. **index.html** Template folder: This folder contains an HTML template for user input, based on which the model will make house price predictions.
3. **requirements.txt**: This file provides packages to install for your web app to run.

### How to Run the model API

1. run the following code in command line. `python app.py`
2. Open the following link in the browser. `127.0.0.1:5000`
3. Provide the Input data for house price prediction and click on **predict house price** button.
4. The button will direct to the following link `127.0.0.1:5000/predict`
5. Predict page will provide the predicted house price at the bottom of the page.

**SNIPPETS OF THE MODEL AS API IS SHOWN BELOW**

House Price Prediction

ml\_num  
N3613770

property\_type  
A.

br  
4.0

br\_plus  
0

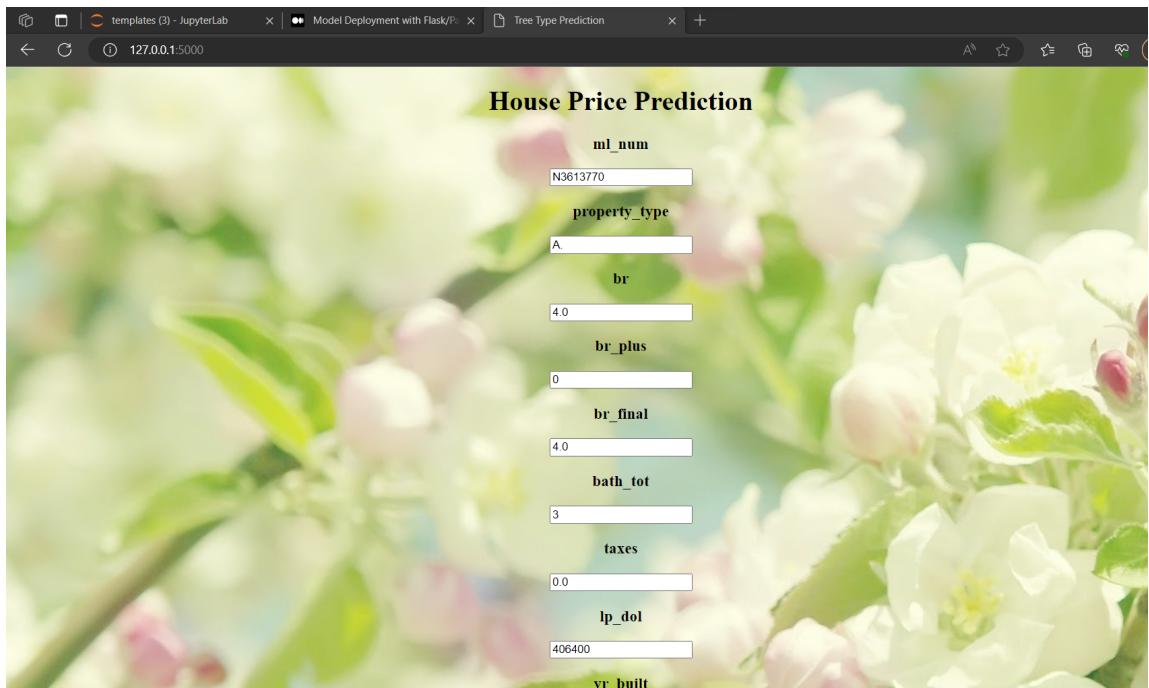
br\_final  
4.0

bath\_tot  
3

taxes  
0.0

lp\_dol  
406400

vr\_built



Model Deployment with Flask/Pa x Tree Type Prediction x +

lot\_frontfeet  
-79.857459

lot\_depthfeet  
30.02

sqft\_numeric  
117.29

id\_community  
sqft\_living  
1750

id\_municipality  
id\_neighborhood  
232

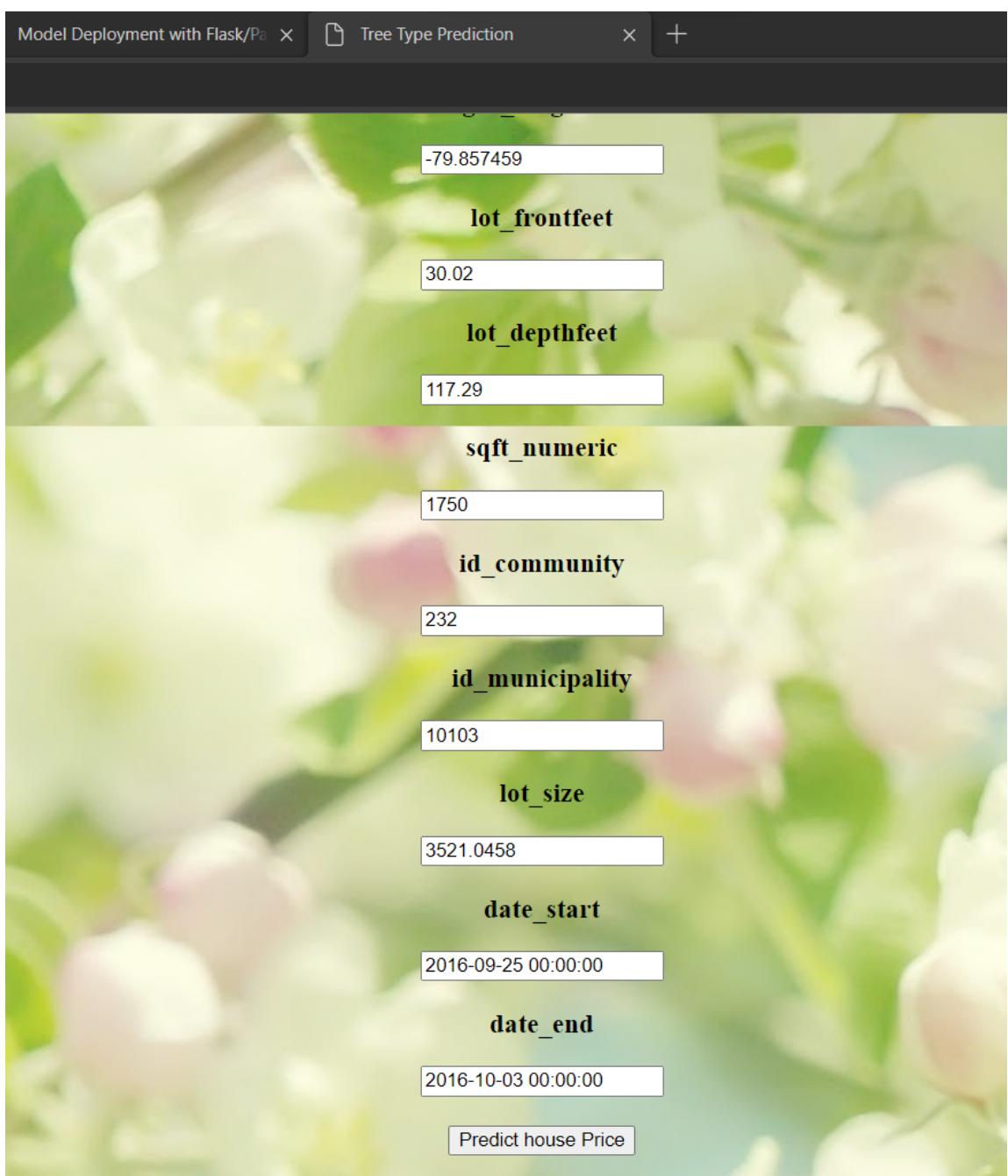
lot\_size  
date\_start  
10103

date\_end  
3521.0458

Predict house Price

2016-09-25 00:00:00

2016-10-03 00:00:00



**date\_end**

2016-10-03 00:00:00

[Predict house Price](#)

**Predicted House Price is 408714.3087299372**