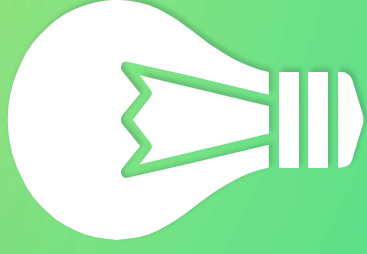


ADVANCED HOUSE PRICE PREDICTION USING MACHINE LEARNING





TEAM PRESENTATION

Sushith Vanga

CSU ID : 2830173

Srinidhi Muppidi

CSU ID : 2822455

Yagnesh

CSU ID : 2825787

Gamya

CSU ID :2829783

Vara Prasad

CSU ID : 2826343

Amulya

CSU ID : 2830012

OUTLINE

- PROJECT SUMMARY
- TECHNOLOGY AND TOOLS
- DATA COLLECTION / PREPROCESSING
- FEATURE ENGINEERING
- MODEL BUILDING
- WEB APPLICATION

PROJECT SUMMARY



PROJECT SUMMARY

- ▶ The objective of this project is predict the house prices using machine learning techniques.
- ▶ Features like number of bedrooms, locality, year of bought etc are dependent variables and SalesPrice is the target feature.



INTRODUCTION

TRAINING THE MODEL

The supervised machine learning model is trained on the preprocessed data.

Feature selection, Feature engineering, model development, evaluation are part of this.

DEVELOPMENT OF WEB APPLICATION

Creating a Graphical user interface for this trained model.

Streamlit API

TECHNOLOGY AND TOOLS

2



TECHNOLOGY

There are 3 types of machine learning techniques

SUPERVISED LEARNING

The machine learns by
using labeled data

Regression and
Classification

Ex : cat or dog, spam or
ham

UNSUPERVISED LEARNING

The machine is trained
on unlabeled data

Association and
Clustering

Ex : segmenting apples
and bananas

REINFORCEMENT LEARNING

An agent interacts with
its environment and
learns by doing

Reward - based

Ex : Self driving cars



TOOLS

- ▶ Python
- ▶ Pycharm
- ▶ Pandas
- ▶ Numpy

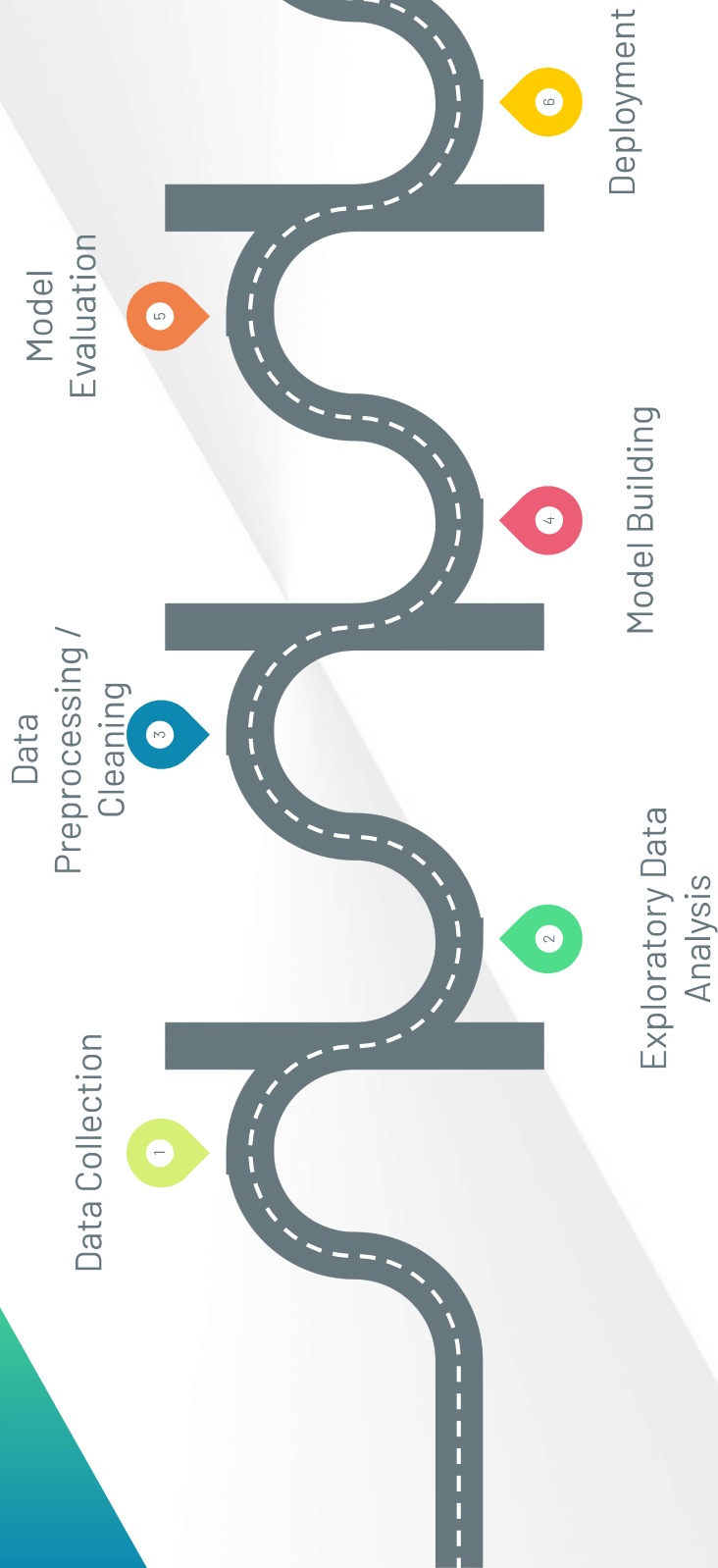
- ▶ Sklearn
- ▶ Kaggle
- ▶ Streamlit

DATA COLLECTION / PROCESSING

3



DATA SCIENCE PROJECT LIFECYCLE



DATA COLLECTION

Data downloaded from **Kaggle Competitions**

- Dataset prepared by *Dean De Cock*

Variables :

Total of 79 variables present in the dataset

- "SalePrice" represent final price at which the house was sold
- Remaining 78 variables represent different attributes like area, car parking, number of rooms etc

DATA PREPROCESSING

Numerical Variables 38

- *Discrete variables* - 17
- *Continuous variables* - 16
- *Temporal variables* - 4

Categorical Variables - 43

FEATURE ENGINEERING



DATA PREPROCESSING

MISSING VALUES

Detection and correction

Numerical missing values are treated with median.

Categorical missing values are treated with "Missing" string.

OUTLIERS

Boxplot of all the dependent features .

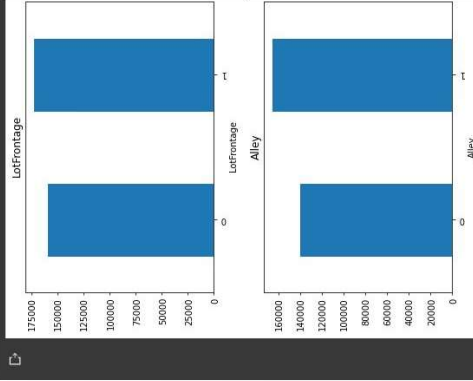
Dropping and replacing with mean, median and mode.

DATA TRANSFORMATION

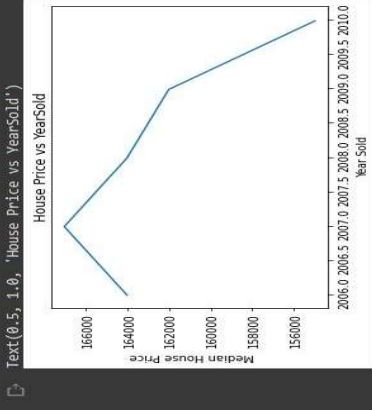
Applied normalization technique on train and test data.

DATA PREPROCESSING (EDA)

```
# for feature in features_with_na:  
    data = dataset.copy()  
    # let's make a variable that indicates 1 if the observation  
    data[feature] = np.where(data[feature].isnull(), 1, 0)  
    # let's calculate the mean SalePrice where the information  
    data.groupby(feature)['SalePrice'].median().plot.bar()  
    plt.title(feature)  
    plt.show()
```

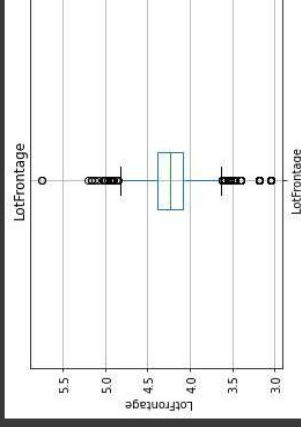


```
## Lets analyze the Temporal Datetime Variables  
## We will check whether there is a relation between year the house  
dataset.groupby('Year Sold')['SalePrice'].median().plot()  
plt.xlabel('Year Sold')  
plt.ylabel('Median House Price')  
plt.title('House Price vs YearSold')
```



Outliers

```
[ ] for feature in continuous_feature:  
    data=dataset.copy()  
    if 0 in data[feature].unique():  
        pass  
    else:  
        data[feature]=np.log(data[feature])  
        data.boxplot(column=feature)  
        plt.ylabel(feature)  
        plt.title(feature)  
        plt.show()
```



MODEL BUILDING

LS



MODEL BUILDING (DATA SPLITTING)





MODEL BUILDING (REGRESSION ANALYSIS)

LINEAR REGRESSION	RIDGE REGRESSION	KNN
MAE - 0.138	MAE - 0.030	MAE - 0.108
MSE - 0.031	MSE - 0.0018	MSE - 0.0195
RMSE - 0.178	RMSE - 0.043	RMSE - 0.1397
RMSE Logarithmic - 1.723	RMSE Logarithmic - 3.141	RMSE Logarithmic - 1.967

RANDOM FOREST

MAE - 0.0113
MSE - 0.0002
RMSE - 0.0171
RMSE Logarithmic - 4.0647

XGBOOST

MAE - 0.1080
MSE - 0.0195
RMSE - 0.1397
RMSE Logarithmic - 1.9675

6

WEB APPLICATION

WEB APPLICATION

- ▶ Built GUI using **streamlit**.
- ▶ Integrating all the functionalities into the GUI.

MODEL BUILDING

- ▶ Option to upload csv files.
- ▶ Display of basic stats by default.
- ▶ Flexibility on feature engineering, model building and metrics.

GUI for Machine Learning Model

Please upload the dataset

Choose a file

Drag and drop file here
Limit 200MB per file

Browse files

train.csv 449.3KB

X

Shape of dataset (1468, 81)

Column names

0	Id
1	MSSubClass
2	MSZoning
3	LotFrontage
4	LotArea
5	Street
6	Alley
7	LotShape
8	LandContour
9	Utilities

Dataset overview

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	La
0	1	60	RL	65.0000	8450	Pave	<NA>	Reg	LV
1	2	20	RL	80.0000	9600	Pave	<NA>	Reg	LV
2	3	60	RL	68.0000	11250	Pave	<NA>	IR1	LV
3	4	70	RL	60.0000	9550	Pave	<NA>	IR1	LV
4	5	60	RL	84.0000	14260	Pave	<NA>	IR1	LV

Nature of the data

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCo
count	1,460,0000	1,460,0000	1,201,0000	1,460,0000	1,460,0000	1,460,00
mean	730.5000	56.8973	70.0500	10,516.8281	6.0993	5.51
std	421.6100	42.3006	24.2848	9,981.2649	1.3830	1.11
min	1,0000	20,0000	21,0000	1,300,0000	1,0000	1,00
25%	365.7500	20,0000	59,0000	7,553,5000	5,0000	5,00
50%	730.5000	50,0000	69,0000	9,478,5000	6,0000	5,00
75%	1,095,2500	70,0000	80,0000	11,601,5000	7,0000	6,00
max	1,460,0000	100,0000	212,0000	212,345,0000	10,0000	0,00

Select the type of variable to analyse

Numerical variables

Number of numerical variables: 38

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt
0	1	60	65,0000	8450	7	5	2003
1	2	20	80,0000	9600	6	8	1976
2	3	60	68,0000	11250	7	5	2001
3	4	70	60,0000	9550	7	5	1915
4	5	60	84,0000	14260	8	5	2000

Outlier detection

Plot out the outliers in the data

Finding the relation between categorical and dependent features

Plots between categorical and dependent?

Feature Engineering

Display the missing values

Do you want to correct the missing values?

Display missing values in numerical data and correct it?

Data Transformation

Select the data transformation technique to be implemented

Standardization

Split the data

Yes

x_train shape (978, 83)

x_test shape (482, 83)

y_train shape (978, 1)

y_test shape (482, 1)

Model Building

Select the machine learning model to train on

Select

Select

Linear Regression

Lasso Regression

K Nearest Neighbours

Random Forest

THANKS!

