

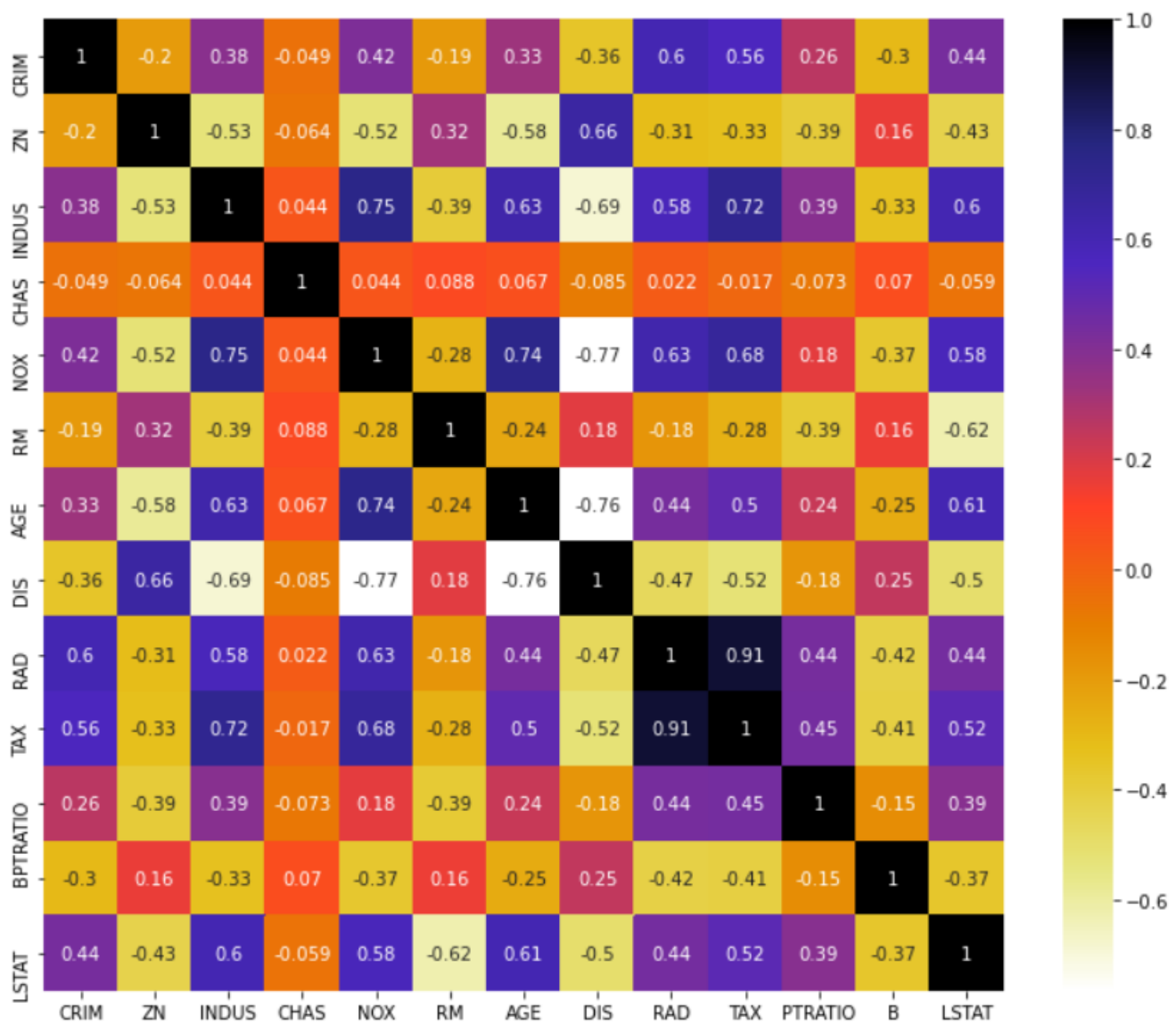
# Project Name: Entire Feature-Selection and Feature-Importance

## Table of Contents

Demo  
Overview  
Motivation  
Technical Aspect  
Installation  
Run/How to Use/Steps  
Directory Tree/Structure of Project  
To Do/Future Scope  
Technologies Used/System Requirements/Tech Stack  
Credits

## Demo

---



## Overview

---

This is diving into Feature Selection and Feature Importance Concept.

Feature Engineering is the act of extracting important features from raw data and transforming them into formats that are suitable for machine learning. Feature Engineering is a data preparation process.

Feature Selection: All features aren't equal. It is all about selecting a small subset of features from a large pool of features. We select those attributes which best explain the relationship of an independent variable with the target variable. There are certain features which are more important than other features to the accuracy of the model. It is different from dimensionality reduction because the dimensionality reduction method does so by combining existing attributes, whereas the feature selection method includes or excludes those features. The methods of Feature Selection are Chi-squared test, correlation coefficient scores, LASSO, Ridge regression etc.

Feature importance refers to a class of techniques for assigning scores to input features to a predictive model that indicates the relative importance of each feature when making a prediction.

This repository contains the code for Feature Selection and Feature Importance using python's various libraries.

It used pandas, matplotlib, seaborn and sklearn libraries.

These libraries help to perform individually one particular functionality.

Pandas objects rely heavily on Numpy objects.

Matplotlib is a plotting library.

Seaborn is data visualization library based on matplotlib.

Sklearn has 100 to 200 models.

The purpose of creating this repository is to gain insights into how to extract and transform data.

These python libraries raised knowledge in discovering these libraries with practical use of it.

It leads to growth in my ML repository.

This above screenshot will help you to understand flow of output.

## Motivation

---

The reason behind building this is, of course all businesses want that their final result should be accurate therefore need to learn Feature Engineering which has a major influence on the performance of machine learning models and even the quality of insights derived during EDA. The models take features as input. A feature is generally a numeric representation of an aspect of real-world phenomena or data. My job as a Data Scientist is to find a clear path to the end goal of insights. Feature engineering helps extract information from raw data, i.e., it has created a lot of features. This means we need to find the main features of the whole lot. This is also known as the Curse of Dimensionality. I sometimes try to relate things of both the world that is

IT and Life Philosophy, In a normal life also it is not possible for one person to study and grab all degrees in the world or all the knowledge in the world so people extract degrees as per their interests so similar analogy is here I have to extract relevant features from raw data as per requirement. Through feature engineering, you can isolate key information, highlight patterns, and bring in domain expertise. If I wish to convey my point to buyer than this can be helpful as I can highlight relevant pattern and isolate key information, which is manipulation of data and hence of decision which is a tricky strategy which is what is call for apart from just a normal coder. I believe that, it's not all always about coding but also about convincing or confusing or ponder them to imagination or influence which is also an must which can be achieved through Feature Engineering. Feature engineering, also known as feature creation, is the process of constructing new features from existing data to train a machine learning model. This step can be more important than the actual model used because a machine learning algorithm only learns from the data, we give it, and creating features that are relevant to a task is absolutely crucial. Statistically, 8 out of 10 people does not like to perform Feature Engineering as it is time-consuming but what fascinates me is, it is art of direction. You know that, master plan can make it to millions too. It enables the machine learning algorithm to train faster. It reduces the complexity of a model and makes it easier to interpret. It improves the accuracy of a model if the right subset is chosen. Feature importance refers to techniques that assign a score to input features based on how useful they are at predicting a target variable. The primary motivations are or should be to either mitigate a specific problem in the interplay between predictors and a model, or to reduce model complexity. RFE is a good example of a wrapper feature selection method. Wrapper methods evaluate multiple models using procedures that add and/or remove predictors to find the optimal combination that maximizes model performance. Feature importance is calculated as the decrease in node impurity weighted by the probability of reaching that node. The node probability can be calculated by the number of samples that reach the node, divided by the total number of samples. The higher the value the more important the feature. Feature Selection is the process where you automatically or manually select those features which contribute most to your prediction variable or output in which you are interested in. Having irrelevant features in your data can decrease the accuracy of the models. Hence, I continue to gain knowledge while practicing the same and spread literary wings in tech-heaven.

---

Pandas module mainly works with the tabular data. It contains Data Frame and Series. Pandas is 18 to 20 times slower than Numpy. Pandas is seriously a game changer when it comes to cleaning, transforming, manipulating and analyzing data.

Matplotlib is used for EDA. Visualization of graphs helps to understand data in better way than numbers in table format. Matplotlib is mainly deployed for basic plotting. It consists of bars, pies, lines, scatter plots and so on. Inline command display visualization inline within frontends like in Jupyter Notebook, directly below the code cell that produced it.

Seaborn provides a high-level interface for drawing attractive and informative statistical graphics. It provides a variety of visualization patterns and visualize random distributions.

Sklearn is known as scikit learn. It provides many ML libraries and algorithms for it. It provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python.

Need to train\_test\_split - Using the same dataset for both training and testing leaves room for miscalculations, thus increases the chances of inaccurate predictions.

The train\_test\_split function allows you to break a dataset with ease while pursuing an ideal model. Also, keep in mind that your model should not be overfitting or underfitting.

Feature selector that removes all low-variance features.

### All Features



### Feature Selection



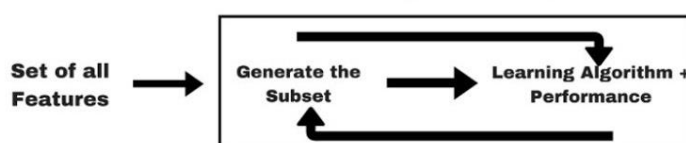
### Final Features

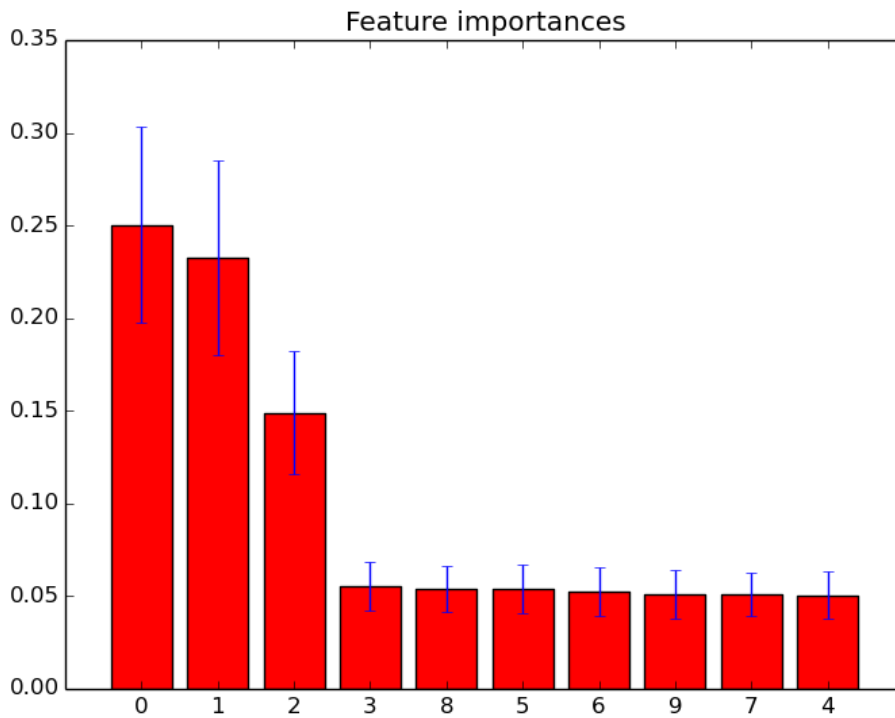


## Feature Selection



### Selecting the best subset





## Installation

---

Using intel core i5 9<sup>th</sup> generation with NVIDIA GFORCE GTX1650.

Windows 10 Environment Used.

Already Installed Anaconda Navigator for Python 3.x

The Code is written in Python 3.8.

If you don't have Python installed then please install Anaconda Navigator from its official site.

If you are using a lower version of Python you can upgrade using the pip package, ensuring you have the latest version of pip, *python -m pip install --upgrade pip and press Enter.*

## Run/How to Use/Steps

---

Keep your internet connection on while running or accessing files and throughout too.

Follow this when you want to perform from scratch.

Open Anaconda Prompt, Perform the following steps:

```
cd <PATH>
```

```
pip install pandas
```

```
pip install matplotlib
```

```
pip install seaborn
```

```
pip install sklearn
```

You can also create requirement.txt file as, `pip freeze > requirements.txt`  
run files.

Follow this when you want to just perform on local machine.

Download ZIP File.

Right-Click on ZIP file in download section and select Extract file option, which will unzip file.

Move unzip folder to desired folder/location be it D drive or desktop etc.

Open Anaconda Prompt, write `cd <PATH>` and press Enter.

eg: `cd C:\Users\Monica\Desktop\Projects\Python Projects`

`1\15)Feature_Engineering+Selection+Importance\`

`All_Feature_Selection_and_Feature_Importance`

In Anconda Prompt, `pip install -r requirements.txt` to install all packages.

Open in Jupyter Notebook, `<filename>.ipynb`

That is,

Open in Jupyter Notebook, `1)Entire_Feature_Selection_and_Feature_Importance.ipynb`

This takes `mobile_dataset.csv`, `santander.csv` file as input dataset.

Please be careful with spellings or numbers while typing filename and easier is just copy filename and then run it to avoid any silly errors.

Note: `cd <PATH>`

[Go to Folder where file is. Select the path from top and right-click and select copy option and paste it next to `cd` one space `<path>` and press enter, then you can access all files of that folder] [`cd` means change directory]

## Directory Tree/Structure of Project

---

Folder: `15)Feature_Engineering+Selection+Importance >`

`All_Feature_Selection_and_Feature_Importance`

`1)Entire_Feature_Selection_and_Feature_Importance.ipynb`

## To Do/Future Scope

---

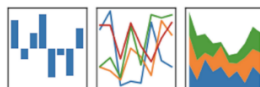
Can use other dataset.

## Technologies Used/System Requirements/Tech Stack

---



pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



matplotlib



## Credits

---

Krish Naik Channel