# Feature Selection in Python

1. Poorya Mohammadi Nasab
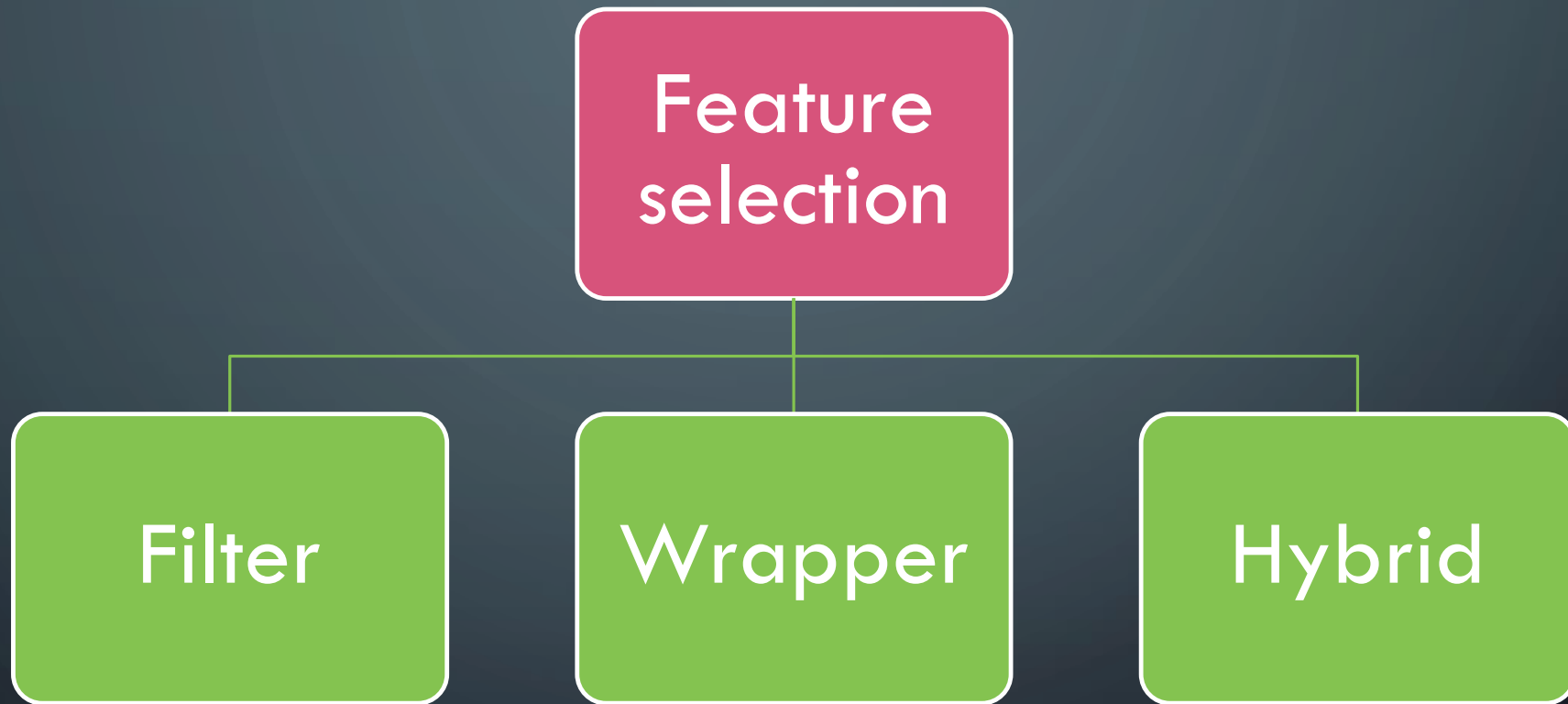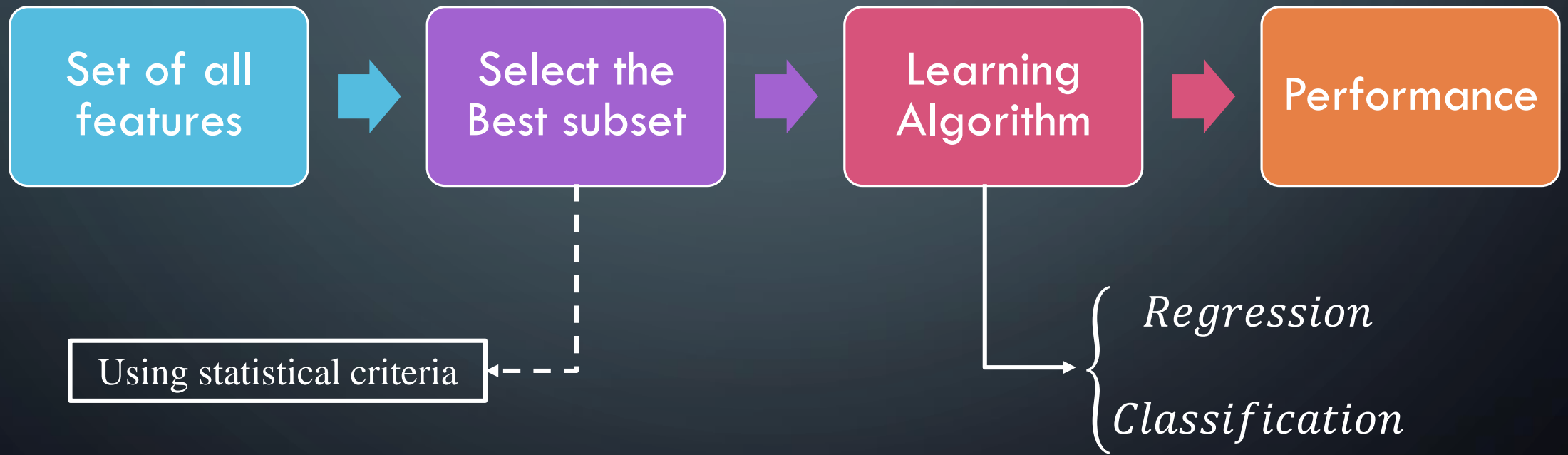2. Behzad Rabiei

Data Mining Course
Fall 2020

# What is Feature selection

Feature selection (FS) is an important step in machine learning and has a significant impact on data mining performance. In recent years, datasets became very massive in number of attributes and instances. By feature, we mean each column in the table and dimension is the number of dataset columns and instance is the number of records. Feature selection algorithms try to remove redundant, irrelevant and noisy data to reduce the size of the dataset and increase classification accuracy.

# Various types of FS
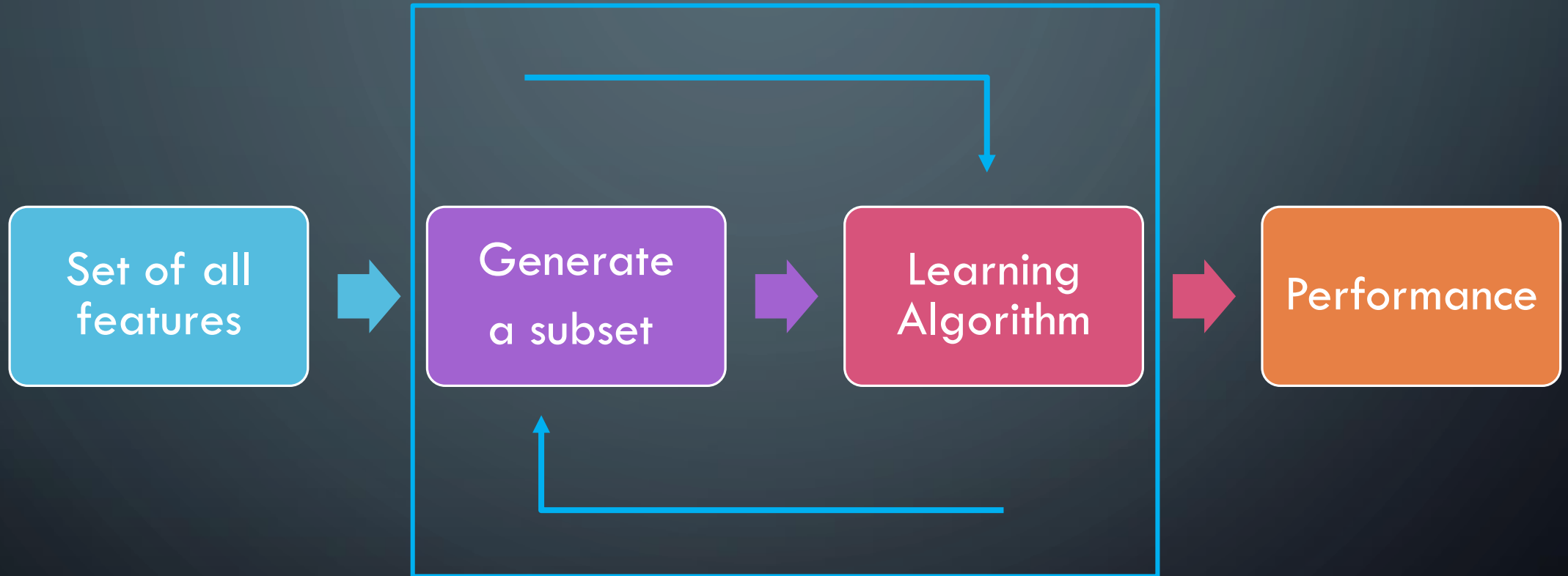
Feature selection

Filter

Wrapper

Hybrid

# Filter-Based Algorithms

Set of all features → Select the Best subset → Learning Algorithm → Performance

Using statistical criteria

$\begin{cases} Regression \\ Classification \end{cases}$

4

# Wrapper-Based Methods

# Crow Search Algorithm (CSA)

CSA algorithm is a meta-heuristic algorithm proposed by Askarzadeh in 2016. The main inspiration of this algorithm came from crow search mechanism for hiding their food. They are known to be thieves, as they are stealing the food of the other birds. The four main principles of CSA are defined as follows:

1) Crows live in the flock's form
2) Crows keep in their memory their hiding places of foods
3) Each member of crow follows each other while doing thievery
4) Crows are very cautions against thievery; they protect their caches by a probability

# Mobile Price Classification Dataset:

| | | | | |
|---|---|---|---|---|
| Battery Power (mAh) | *Bluetooth* | Processor Speed | *Dual Simcards* | Front Camera (mega pixels) |
| *4G* | Internal Memory (GB) | Mobile Depth (Cm) | Weight of mobile (gr) | Number of cores |
| Primary Camera (mega pixels) | Pixel Resolution Height | Pixel Resolution Width | RAM (MB) | Screen Height (Cm) |
| Screen Width (Cm) | Talk Time (h) | *3G* | *Touch Screen* | *WiFi* |

# Target Variable

Price range is the target variable with a value of:

- 0(low cost)
- 1(medium cost)
- 2(high cost)
- 3(very high cost)

# 1.Univariate Selection

$S$tatistical tests can be used to select those features that have the strongest relationship with the output variable. The scikit-learn library provides the SelectKBest class that can be used with a suite of different statistical tests to select a specific number of features. In our example we use the chi-squared (chi²) statistical test for non-negative features to select 10 of the best features from the Mobile Price Range Prediction Dataset.

# Using SelectKBest in Python

```python
>>> import pandas as pd
>>> import numpy as np
>>> from sklearn.feature_selection import SelectKBest
>>> from sklearn.feature_selection import chi2

>>> data = pd.read_csv('train.csv')
>>> X = data.iloc[:,0:20]   #independent columns
>>> y = data.iloc[:,-1]        #target column(price range)
    #apply SelectKBest class to extract top 10 best features
>>> bestfeatures = SelectKBest(score_func=chi2, k=10)
>>> fit = bestfeatures.fit(X,y)
>>> dfscores = pd.DataFrame(fit.scores_)
>>> dfcolumns = pd.DataFrame(X.columns)
    #concat two dataframes for better visualization
>>> featureScores = pd.concat([dfcolumns,dfscores],axis=1)
>>> featureScores.columns = ['Specs','Score']  #naming the dataframe columns
>>> print(featureScores.nlargest(10,'Score'))  #print 10 best features
```

# Output Data frame

Feature indexes

Calculated feature
score using Chi2

|     | Specs         | Score         |
|-----|---------------|---------------|
| 13  | ram           | 931267.519053 |
| 11  | px_height     | 17363.569536  |
| 0   | battery_power | 14129.866576  |
| 12  | px_width      | 9810.586750   |
| 8   | mobile_wt     | 95.972863     |
| 6   | int_memory    | 89.839124     |
| 15  | sc_w          | 16.480319     |
| 16  | talk_time     | 13.236400     |
| 4   | fc            | 10.135166     |
| 14  | sc_h          | 9.614878      |

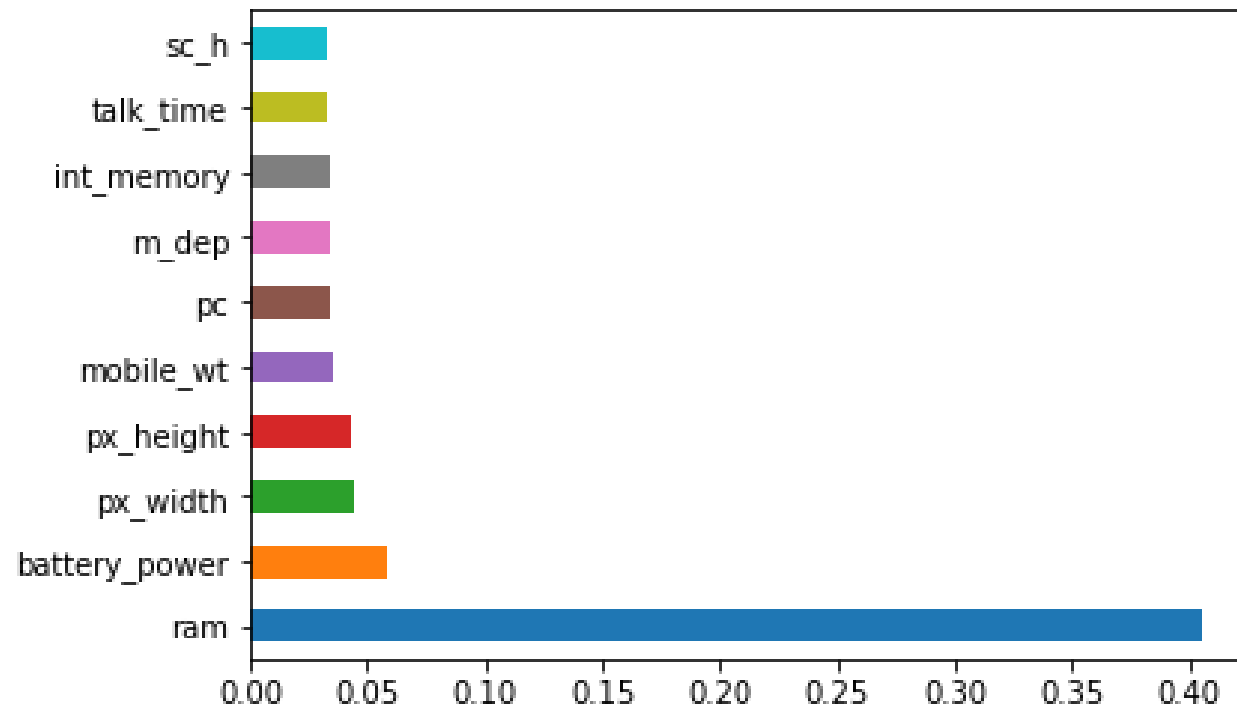Feature names

# 2.Feature Importance

Feature importance gives you a score for each feature of your data, the higher the score more important or relevant is the feature towards your output variable. Feature importance is an inbuilt class that comes with Tree Based Classifiers, we will be using Extra Tree Classifier for extracting the top 10 features for the dataset.

# Using Feature importance

```
>>> import pandas as pd
>>> import numpy as np
>>> import matplotlib.pyplot as plt
>>> from sklearn.ensemble import ExtraTreesClassifier

>>> data = pd.read_csv('train.csv')
>>> X = data.iloc[:,0:20]   #independent columns
>>> y = data.iloc[:,-1]       #target column(price range)
>>> model = ExtraTreesClassifier()
>>> model.fit(X,y)
>>> print(model.feature_importances_) #use inbuilt class feature_importances of tree based classifiers
        #plot graph of feature importances for better visualization
>>> feat_importances = pd.Series(model.feature_importances_, index=X.columns)
>>> feat_importances.nlargest(10).plot(kind='barh')
>>> plt.show()
```
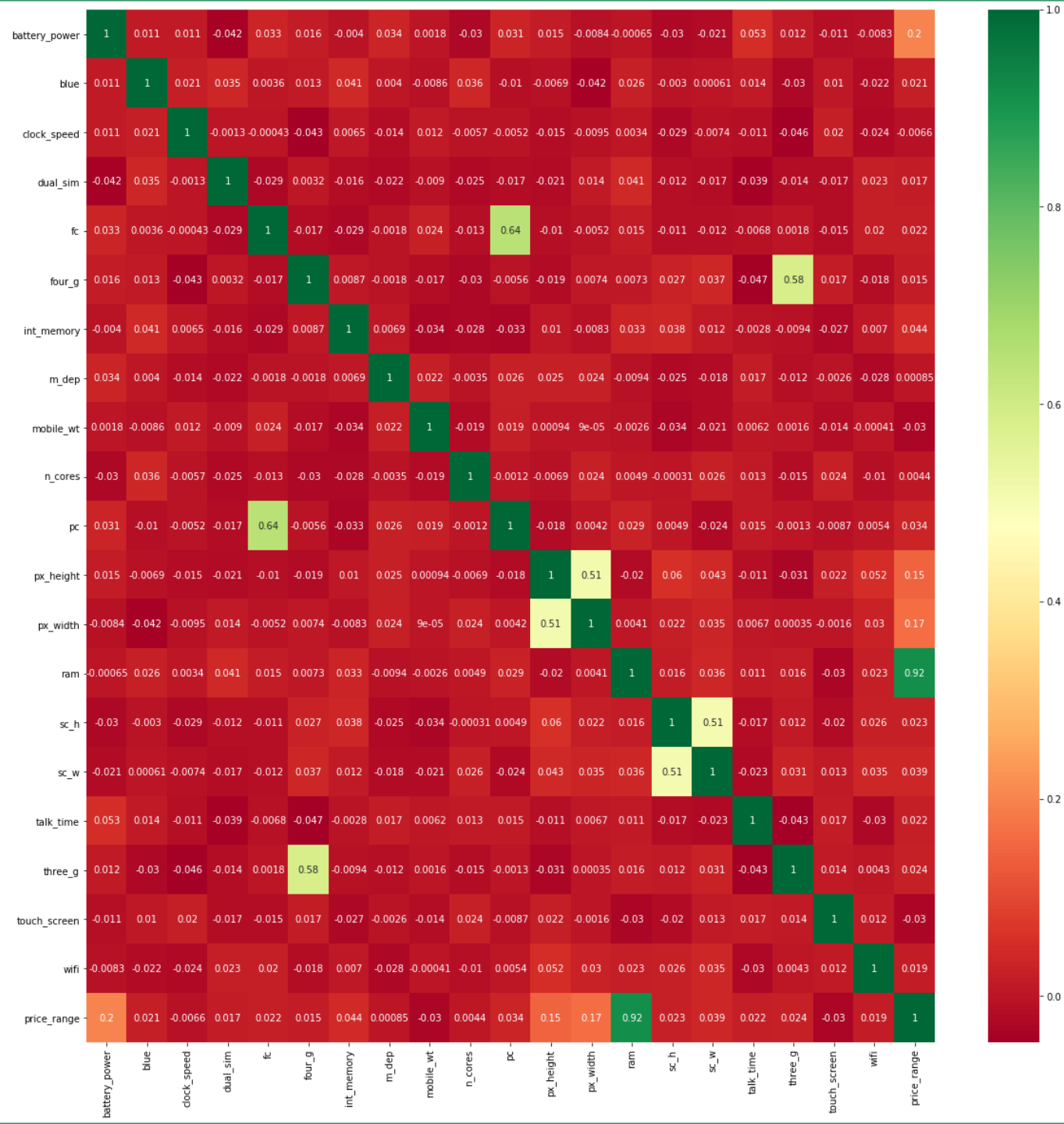
# Plt Output

# 3.Correlation Matrix with Heatmap

Correlation states how the features are related to each other or the target variable. Correlation can be positive (increase in one value of feature increases the value of the target variable) or negative (increase in one value of feature decreases the value of the target variable). Heatmap makes it easy to identify which features are most related to the target variable, we will plot heatmap of correlated features using the seaborn library.

# Plotting Heatmap

```
>>> import pandas as pd
>>> import numpy as np
>>> import seaborn as sns

>>> data = pd.read_csv('train.csv')
>>> X = data.iloc[:,0:20]   #independent columns
>>> y = data.iloc[:,-1]      #target column(price range)
    #get correlations of each features in dataset
>>> corrmat = data.corr()
>>> top_corr_features = corrmat.index
>>> plt.figure(figsize=(20,20))
    #plot heat map
>>> g=sns.heatmap(data[top_corr_features].corr(),annot=True,cmap="RdYlGn")
```