# औद्योगिक प्रशिक्षण के लिए राष्ट्रीय संस्थान

## National Institute for Industrial Training

One Premier Organization with Non Profit Status | Registered Under Govt. of WB
Empanelled Under Planning Commission Govt. of India
Inspired By: National Task Force on IT & SD Government of India

National Institute for Industrial Training- One Premier Organization with Non Profit Status Registered Under Govt.of West Bengal,Empanelled Under Planning Commission Govt.of India , Empanelled Under Central Social Welfare Board Govt. of India , Registered with National Career Services , Registered with National Employment Services.

**SUBJECT: PYTHON WITH DATA SCIENCE**

**Submitted by: PREETI SHARMA**

**Submitted to: SOUMATANU MAJUMDAR**

**Name:** Preeti Sharma

**College:** Kurukshetra University(U.I.E.T.)

**Course:** Biotech Engineering

**Year:** 2nd

**Year of passing:** 2023
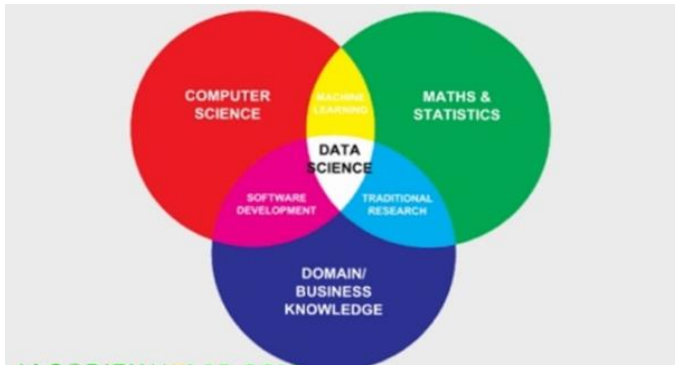
**E-mail ID:** preetishr27@gmail.com

# Contents

# Acknowledgment

 I am privileged and grateful to acknowledge my knowledge to all those who have guided me to put these ideas, well above the level of simplicity and into something concrete. I express my warm gratitude to National Institute for Industrial Training for their constant guidance and supervision as well as for providing necessary information regarding the project. I would like to express my sole thanks of gratitude to (Soumotanu Majumdar) for his support, co-operation and encouragement which helped me in the completion of this project.

My journey of completing the project wouldn't would remain incomplete if I don't extend my gratitude to my parents for their love , affection , support and constant guidance.
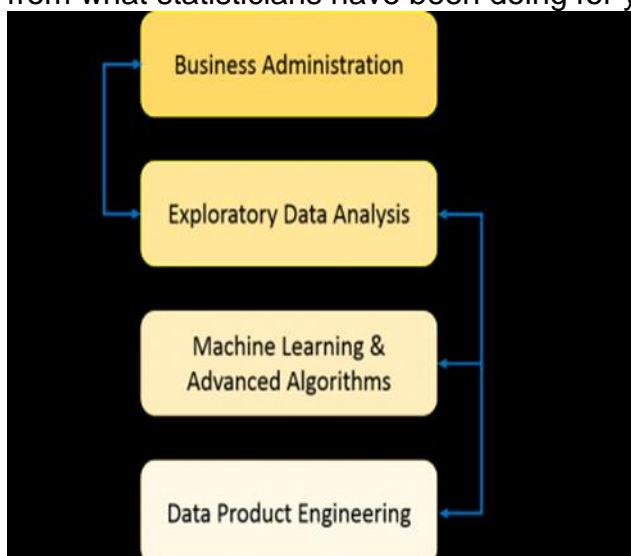
# Introduction

 "A combination of information technology , modelling , and business management"- Dr. Thomas Miller of Northwestern University.

In today's world Data science and Machine Learning are one of the most prominent topics in the arena of technology. Data science and machine learning are two distinguishable topics although they look analogous to each other. Combining , computing , comparing and concluding the insights of datas are known as data science while on the flip side machine learning is the processes by which a result can be obtained from the given data.

## What is Data Science ?

Have anyone wondered why our life in this era of iron age has become so comfortable and scientific? It's only because of computers , smartphones , tablets , laptops and many more electronic devices which have completely digitalized our life and consequently resulted in huge amount data. This data needs to be processed , organized and coordinated. The phenomenon of mastering these process or the study of these data sets are basically categorized under data science.

Data Science is a blend of various tools, algorithms, and machine learning principles with the goal to discover hidden patterns from the raw data. But how is this different from what statisticians have been doing for years?

 The answer lies in the difference

between explaining and predicting. The principal purpose of Data Science is to find patterns within data. It uses various statistical techniques to analyze and draw insights from the data. From data extraction, wrangling and pre-processing, a Data Scientist must scrutinize the data thoroughly. Then, he has the responsibility of making predictions from the data. The goal of a Data Scientist is to derive conclusions from the data. Through these conclusions, he is able to assist companies in making smarter business decisions. We will divide this blog into various sections to understand the role of a Data Scientist in more detail. Industries need data to help them make careful decisions. Data Science churns raw data into meaningful insights. Therefore, industries need data science. A Data Scientist is a wizard who knows how to create magic using data. A skilled Data Scientist will know how to dig out meaningful information with whatever data he comes across. He helps the company in the right direction. The company requires strong data-driven decisions at which he's an expert. The Data Scientist is an expert in various underlying fields of Statistics and Computer Science. He uses his analytical aptitude to solve business problems. Data Scientist is well versed with problem-solving and is assigned to find patterns in data. His goal is to recognize redundant samples and draw insights from it. Data science requires a variety of tools to extract information from the data. A Data Scientist is responsible for collecting, storing and maintaining the structured and unstructured form of data.

While the role of Data Science focuses on the analysis and management of data, it is dependent on the area that the company is specialized in. This requires the Data Scientist to have domain knowledge of that particular industry.

# Advantages

☐**A great library ecosystem:** A great choice of libraries is one of the main reasons Python is the most popular programming language used for AI. A library is a module or a group of modules published by different sources like Py Pi which include a pre-written piece of code that allows users to reach some functionality or perform different actions. Python libraries provide base level items so developers don't have to code them from the very beginning every time.

☐**A low entry barrier:** Working in the ML and AI industry means dealing with a bunch of data that you need to process in the most convenient and effective way. The low entry barrier allows more data scientists to quickly pick up Python and start using it for AI development without wasting too much effort on learning the language.

☐**Flexibility:** Python for machine learning is a great choice, as this language is very flexible. It offers an option to choose either to use OOPs or scripting. There's also no need to recompile the source code, developers can implement any changes and quickly see the results. Programmers can combine Python and other languages to reach their goals. The flexibility factor decreases the possibility of errors, as programmers have a chance to take the situation under control and work in a comfortable environment.

☐**Platform independence:** Python is not only comfortable to use and easy to learn but also very versatile. What we mean is that Python for machine learning development can run on any platform including Windows, MacOS, Linux, UNIX, and twenty-one others. To transfer the process from one platform to another, developers need to implement several small-scale changes and modify some lines of code to create an executable form of code for the chosen platform. Developers can use packages like Py Installer to prepare their code for running on different platforms. Again, this saves time and money for tests on various platforms and makes the overall process more simple and convenient.

☐**Readability**: Python is very easy to read so every Python developer can understand the code of their peers and change, copy or share it. There's no confusion, errors or conflicting paradigms, and this leads to a more efficient exchange of algorithms, ideas, and tools between AI and ML professionals.

**Good visualization options:** We've already mentioned that Python offers a variety of libraries, and some of them are great visualization tools. However, for AI developers, it's important to highlight that inartificial intelligence, deep learning, and machine learning, it's vital to be able to represent data in a human-readable format. Libraries like Matplotlib allow data scientists to build charts, histograms, and plots for better data comprehension, effective presentation, and visualization. Different application programming interfaces also simplify the visualization process and make it easier to create clear reports.

**Community support:** It's always very helpful when there's strong community support built around the programming language. Python is an open-source language which means that there's a bunch of resources open for programmers starting from beginners and ending with pros. A lot of Python documentation is available online as well as in Python communities and forums, where programmers and machine learning developers discuss errors, solve problems, and help each other out. Python programming language is absolutely free as is the variety of useful libraries and tools.

**Growing Popularity:** As a result of the advantages discussed above, Python is becoming more and more popular among data scientists. This means it's easier to search for developers and replace team players if required. Also, the cost of their work may be not as high as when using a less popular programming language.

# Future Scope



Let's dig deeper and see how Data Science is being used in various domains.

- How about if you could understand the precise requirements of your customers from the existing data like the customer's past browsing history, purchase history, age and income. No doubt you had all this data earlier too, but now with the vast amount and variety of data, you can train models more effectively and recommend the product to your customers with more precision. Wouldn't it be amazing as it will bring more business to your organization?

Let's take a different scenario to understand the role of Data Science in decision making. How about if your car had the intelligence to drive you home? The self-driving cars collect live data from sensors, including radars, cameras, and lasers to create a map of its surroundings. Based

- on this data, it takes decisions like when to speed up, when to speed down, when to overtake, where to take a turn – making use of advanced machine learning algorithms.
- Let's see how Data Science can be used in predictive analytics. Let's take weather forecasting as an example. Data from ships, aircraft, radars, satellites can be collected and analyzed to build models. These models will not only forecast the weather but also help in predicting the occurrence of any natural calamities. It will help you to take appropriate measures beforehand and save many precious lives.

# System Requirements

☐**Operating system:** Windows 7 or newer, 64-bitmacOS 10.13+, or Linux, including Ubuntu, RedHat, CentOS 6+, and others.

☐**System architecture:** Windows- 64-bit x86, 32-bitx86; MacOS- 64-bit x86; Linux- 64-bit x86, 64-bitPower8/Power9.

☐**Disk Space:** Minimum 5 GB disk space to download and install anaconda distribution.

☐**RAM**: 2 GB RAM recommended.

☐**Graphics:** For neural networks in Machine Learning, better graphics cards will yield faster results with some high-end graphics cards created especially for Machine Learning purposes. However, Google Colab or some similar website can be used to perform same tasks using cloud computing, without needing a graphics card.

# Objectives

The goal of this project is to write a python program to implement a few Machine Learning algorithms, namely, Linear Regression, Logistic Regression, Decision Tree and Random Forest, with a graphical user interface, using Tkinter, such that users without much knowledge of programming or Machine Learning can use these methods with ease. The objectives are as follows:

☐The Tkinter application should allow the user to select .csv(comma separated values) files to use as the data set.

☐The user should be able to view the file as a table in order to properly select the variables.

☐The users should be able to select the independent variables as well as the dependent variable.

☐The program should allow the user to enter the test size and random state for the train test split operation.

☐The program should allow the user to obtain the coefficients and scatter plot in case of Linear Regression.

☐The user should be able to view the confusion matrix and classification report for Logistic, Regression, Decision Tree and Random Forest.

☐The user should be able to view the mean absolute error, mean squared error and root mean squared error in all the cases.

☐The program is to be written in such a way that other Machine Learning algorithms can be added without much difficulty.

# SOURCE CODE

```python
1. import numpy as np
2. import pandas as pd
3. import matplotlib.pyplot as plt
4. import seaborn as sns
5. import warnings
6. warnings.filterwarnings('ignore') from sklearn.ensemble
   import RandomForestClassifier
7. from sklearn.svm import SVC
8. import tkinter as tk
9. from matplotlib.backends.backend_tkagg import
   FigureCanvasTkAgg
10.  from sklearn.linear_model import SGDClassifier
11.  from sklearn.metrics import confusion_matrix,
   classification_report
12.  from sklearn.preprocessing import StandardScaler,
   LabelEncoder
13.  from sklearn.model_selection import train_test_split,
   GridSearchCV, cross_val_score
14.  from imblearn.over_sampling import SMOTE
15.  !pip install imblearn
16.  df = pd.read_csv("winequality-red.csv")
17.  df.columns
18.  plt.figure(figsize=(10, 6))
         sns.countplot(df["quality"], palette="muted")
         df["quality"].value_counts()
19.  print("Rows, columns: " + str(df.shape))
20.  df.head()
21.  df.isnull().sum()
22.  df.info()
23.  df.describe()
24.  fig, ax = plt.subplots(ncols=6, nrows=2,figsize=(20,20))
25.  index = 0
26.  ax = ax.flatten()
27.  for col, value in df.items():
28.  if col != 'type':
```

```python
29.   sns.boxplot(y=col, data=df, ax=ax[index]) index += 1
30.   plt.tight_layout(pad=0.5, w_pad=0.7, h_pad=5.0)


31.   fig, ax = plt.subplots(ncols=6, nrows=2, figsize=(20,10))
32.   index = 0
33.   ax = ax.flatten()
34.   for col, value in df.items():
35.   if col != 'type':
36.   sns.distplot(value, ax=ax[index]) index += 1
37.   plt.tight_layout(pad=0.5, w_pad=0.7, h_pad=5.0)
38.   #Here we see that fixed acidity does not give any
      specification to classify the quality
39.    fig = plt.figure(figsize = (10,6))
40.   sns.barplot(x = 'quality', y = 'fixed acidity', data = df)
41.   #Here we see that its quite a downing trend in the volatile
      acidity as we go higher the
42.   fig = plt.figure(figsize = (10,6))
43.   sns.barplot(x = 'quality', y = 'volatile acidity', data = df)
44.   #Composition of citric acid go higher as we go higher in
      the quality of the wine
45.   fig = plt.figure(figsize = (10,6))
46.   sns.barplot(x = 'quality', y = 'citric acid', data = df)
47.   fig = plt.figure(figsize = (10,6))
48.   sns.barplot(x = 'quality', y = 'residual sugar', data = df)
49.   #Composition of chloride also go down as we go higher in
      the quality of the wine
50.   fig = plt.figure(figsize = (10,6))
51.   sns.barplot(x = 'quality', y = 'chlorides', data = df)
52.   fig = plt.figure(figsize = (10,6))
53.   sns.barplot(x = 'quality', y = 'free sulfur dioxide', data = df)
54.    fig = plt.figure(figsize = (10,6))


55.   sns.barplot(x = 'quality', y = 'total sulfur dioxide', data = df)
56.   #Sulphates level goes higher with the quality of wine
57.   fig = plt.figure(figsize = (10,6))
```

```python
58.  sns.barplot(x = 'quality', y = 'sulphates', data = df)
59.  #Alcohol level also goes higher as te quality of wine
     increases
60.  fig = plt.figure(figsize = (10,6))
61.  sns.barplot(x = 'quality', y = 'alcohol', data = df)
62.  #Making binary classificaion for the response variable.
63.  #Dividing wine as good and bad by giving the limit for the
     quality
64.  bins = (2, 6.5, 8)
65.  group_names = ['bad', 'good']
66.  df['quality'] = pd.cut(df['quality'], bins = bins, labels =
     group_names)
67.  #Now lets assign a labels to our quality variable
68.  label_quality = LabelEncoder()
69.  #Bad becomes 0 and good becomes 1
70.  df['quality'] = label_quality.fit_transform(df['quality'])
71.  df['quality'].value_counts()
72.  sns.countplot(df['quality'])
73.  sns.scatterplot(x='fixed acidity', y='density', data=df)
74.  plt.figure(figsize=(12,8),dpi=200)
75.  sns.scatterplot(x='fixed acidity',y='density',
     data=df,hue='quality', palette='viridis')
76.  sns.scatterplot(x='volatile acidity', y='density', data=df)
77.  plt.figure(figsize=(12,8),dpi=200)
78.  sns.scatterplot(x='volatile acidity',y='density',
     data=df,hue='quality', palette='viridis')
79.  sns.scatterplot(x='citric acid', y='density', data=df)
80.  plt.figure(figsize=(12,8),dpi=200)
81.  sns.scatterplot(x='citric acid',y='density',
     data=df,hue='quality', palette='viridis')
82.  sns.scatterplot(x='alcohol', y='density', data=df)
83.  plt.figure(figsize=(12,8),dpi=200)
84.  sns.scatterplot(x='alcohol',y='density',
     data=df,hue='quality', palette='viridis')
85.  sns.scatterplot(x='residual sugar', y='density', data=df)
86.  plt.figure(figsize=(12,8),dpi=200)
```

```python
87.  sns.scatterplot(x='residual sugar',y='density',
     data=df,hue='quality', palette='viridis')
88.  corr = df.corr()
89.  #Let's look at the correlation among the variables using
     Correlation chart
90.  colormap = plt.cm.viridis
91.  plt.figure(figsize=(12,12))
92.  plt.title('Correlation of Features', y=1.05, size=15)
93.  sns.heatmap(df.astype(float).corr(),linewidths=0.1,vmax=1
     .0, square=True,
94.  linecolor='white', annot=True)
95.  Data Transformation We want to transfer the score(num)
     to low-medium-high quality
96.  level(categorical) by:
97.  3,4 -> low
98.  5,6 -> medium
99.  7,8,9 -> high
100. quality = df["quality"].values
101. category = []
102. for num in quality:
103. if num < 5:
104. category.append("Low")
105. elif num > 6:
106. category.append("High")
107. else:
108. category.append("Medium")
109.  [(i, category.count(i)) for i in set(category)]
110. #Now seperate the dataset as response variable and
     feature variabes
111. X = df.drop('quality', axis = 1)
112. y = df['quality']
113. #Train and Test splitting of data
114. X_train, X_test, y_train, y_test = train_test_split(X, y,
     test_size = 0.2, random_state
115. #Applying Standard scaling to get optimized result
116. sc = StandardScaler()
117. X_train = sc.fit_transform(X_train)
```

```
118. X_test = sc.fit_transform(X_test)
119. rfc = RandomForestClassifier(n_estimators=200)
120. rfc.fit(X_train, y_train)
121. pred_rfc = rfc.predict(X_test)
122. #Let's see how our model performed
123. print(classification_report(y_test, pred_rfc))
124. #Confusion matrix for the random forest classification
125. print(confusion_matrix(y_test, pred_rfc))
126. sgd = SGDClassifier(penalty=None)
127. sgd.fit(X_train, y_train)
128. pred_sgd = sgd.predict(X_test)
129. svc = SVC()
130. svc.fit(X_train, y_train)
131. pred_svc = svc.predict(X_test)
132. print(classification_report(y_test, pred_svc))
133. # @hidden_cell
134. # relabel back : 0 means good, 1 for low, 2 for medium for
     better visualization
135. y_test_re = list(y_test)
136. for i in range(len(y_test_re)):
137. if y_test_re[i] == 0:
138. y_test_re[i] = "good"
139. if y_test_re[i] == 1:
140. y_test_re[i] = "low"
141. if y_test_re[i] == 2:
142. y_test_re[i] = "medium"
143. pred_sgd_re = list(pred_sgd)
144. for i in range(len(pred_sgd_re)):
145. if pred_sgd_re[i] == 0:
146. pred_sgd_re[i] = "good"
147. if pred_sgd_re[i] == 1:
148. pred_sgd_re[i] = "low"
149. if pred_sgd_re[i] == 2:
150. pred_sgd_re[i] = "medium"
151. y_actu = pd.Series(y_test_re, name='Actual')
152. y_pred = pd.Series(pred_sgd_re, name='Predicted')
153. svm_confusion = pd.crosstab(y_actu, y_pred)
```

```
154. svm_confusion
155. Grid Search CV
156. #Finding best parameters for our SVC model
157. param = {
158. 'C': [0.1,0.8,0.9,1,1.1,1.2,1.3,1.4],
159. 'kernel':['linear', 'rbf'],
160. 'gamma' :[0.1,0.8,0.9,1,1.1,1.2,1.3,1.4]
161. }
162. grid_svc = GridSearchCV(svc, param_grid=param,
     scoring='accuracy', cv=10)
163. grid_svc.fit(X_train, y_train)
164. #Best parameters for our svc model
165. grid_svc.best_params_
166. #Let's run our SVC again with the best parameters.
167. svc2 = SVC(C = 1.2, gamma = 0.9, kernel= 'rbf')
168. svc2.fit(X_train, y_train)
169. pred_svc2 = svc2.predict(X_test)
170. print(classification_report(y_test, pred_svc2))
171. grid_svc.best_score_
172. df.columns[:-1]
173. #Now lets try to do some evaluation for random forest
     model using cross validation.
174. rfc_eval = cross_val_score(estimator = rfc, X = X_train, y
     = y_train, cv = 10)
175. rfc_eval.mean()
176. from sklearn.ensemble import AdaBoostClassifier
177. model3 = AdaBoostClassifier(random_state=1)
178. model3.fit(X_train, y_train)
179. y_pred3 = model3.predict(X_test)
180. print(classification_report(y_test, y_pred3))
181. from sklearn.ensemble import GradientBoostingClassifier
182. model4 = GradientBoostingClassifier(random_state=1)
183. model4.fit(X_train, y_train)
184. y_pred4 = model4.predict(X_test)
185. print(classification_report(y_test, y_pred4))
186. # Filtering df for only good quality
187. df_temp = df[df['quality']==1]
```

```
188. df_temp.describe()
189. # Filtering df for only bad quality
190. df_temp2 = df[df['quality']==0]
191. df_temp2.describe()
192. def showQuality():
193. new =
     np.array([[float(e1.get()),float(e2.get()),float(e3.get()),float(e4.
     get()),float(e5.get())])
194. Ans = RF_clf.predict(new)
195. fin=str(Ans)[1:-1]
196. #IT WILL remove[ ]
197. quality.insert(0, fin)
198. #Train and evaluate the Random Forest Classifier with
     Cross Validation
199. # Instantiate the Random Forest Classifier
200. RF_clf = RandomForestClassifier(random_state=0)
201. # Compute k-fold cross validation on training dataset and
     see mean accuracy score
202. cv_scores = cross_val_score(RF_clf,X_train, y_train,
     cv=10, scoring='accuracy')
203. #Perform predictions
204. RF_clf.fit(X_train, y_train)
205. pred_RF = RF_clf.predict(X_test)
206. master = tk.Tk()
207. tk.Label(master, text="Fixed Acidity", anchor="nw",
     width=15).grid(row=0)
208. tk.Label(master, text="Volatile Acidity", anchor="nw",
     width=15).grid(row=1)
209. tk.Label(master, text="Citric Acid", anchor="nw",
     width=15).grid(row=2)
210. tk.Label(master, text="Residual Sugar", anchor="nw",
     width=15).grid(row=3)
211. tk.Label(master, text="Chlorides", anchor="nw",
     width=15).grid(row=4)
212. tk.Label(master, text="Sulfur Dioxide", anchor="nw",
     width=15).grid(row=5)
```

```python
213. tk.Label(master, text="Total Sulfur Dioxide", anchor="nw",
     width=15).grid(row=6)
214. tk.Label(master, text="Density", anchor="nw",
     width=15).grid(row=7)
215. tk.Label(master, text="pH", anchor="nw",
     width=15).grid(row=8)
216. tk.Label(master, text="Sulphates", anchor="nw",
     width=15).grid(row=9)
217. tk.Label(master, text="Alcohol", anchor="nw",
     width=15).grid(row=10)
218. tk.Label(master, text = "Quality", anchor="nw",
     width=15).grid(row=13)
219. e1 = tk.Entry(master)
220. e2 = tk.Entry(master)
221. e3 = tk.Entry(master)
222. e4 = tk.Entry(master)
223. e5 = tk.Entry(master)
224. e6 = tk.Entry(master)
225. e7 = tk.Entry(master)
226. e8 = tk.Entry(master)
227. e9 = tk.Entry(master)
228. e10 = tk.Entry(master)
229. e11 = tk.Entry(master)
230. quality = tk.Entry(master)
231. e1.grid(row=0, column=1)
232. e2.grid(row=1, column=1)
233. e3.grid(row=2, column=1)
234. e4.grid(row=3, column=1)
235. e5.grid(row=4, column=1)
236. e6.grid(row=5, column=1)
237. e7.grid(row=6, column=1)
238. e8.grid(row=7, column=1)
239. e9.grid(row=8, column=1)
240. e10.grid(row=9, column=1)
241. e11.grid(row=10, column=1)
242. quality.grid(row=13, column=1)
```

243. tk.Button(master, text='Quit',
    command=master.destroy,width=15).grid(row=11, column=0,
244. tk.Button(master, text='Find Quality',
    command=showQuality,width=17).grid(row=11, colum
245. tk.Button(master, text='Project By',width=15).grid(row=14,
    column=0, pady=4)
246. tk.Button(master, text='Preeti
    Sharma',width=17).grid(row=14, column=1, pady=4)
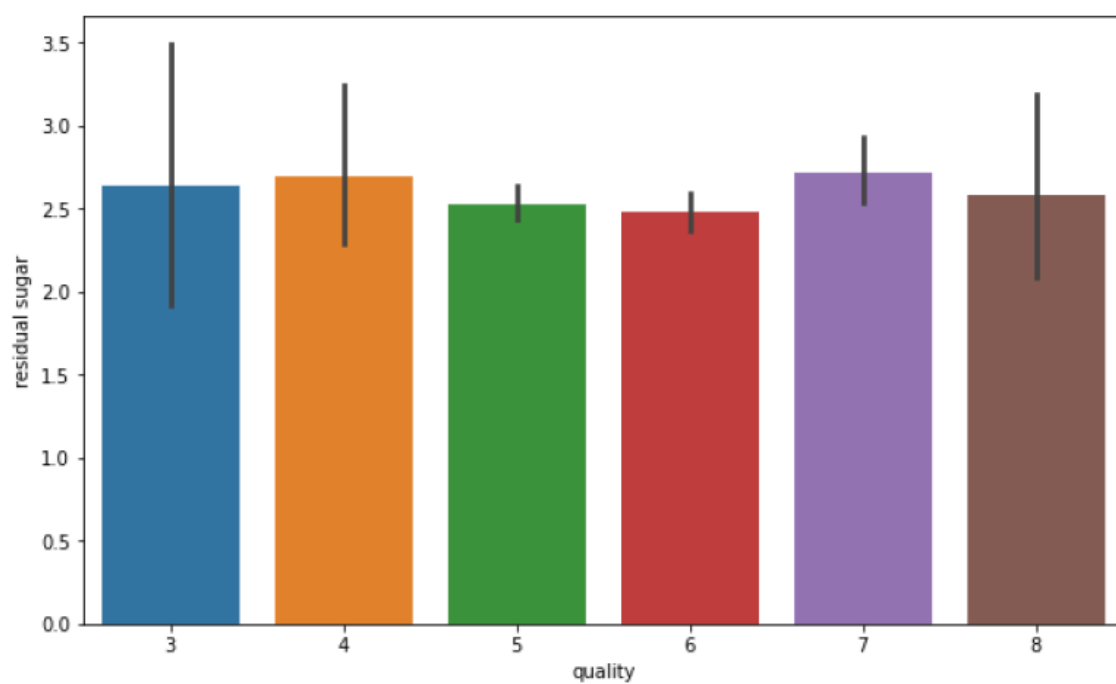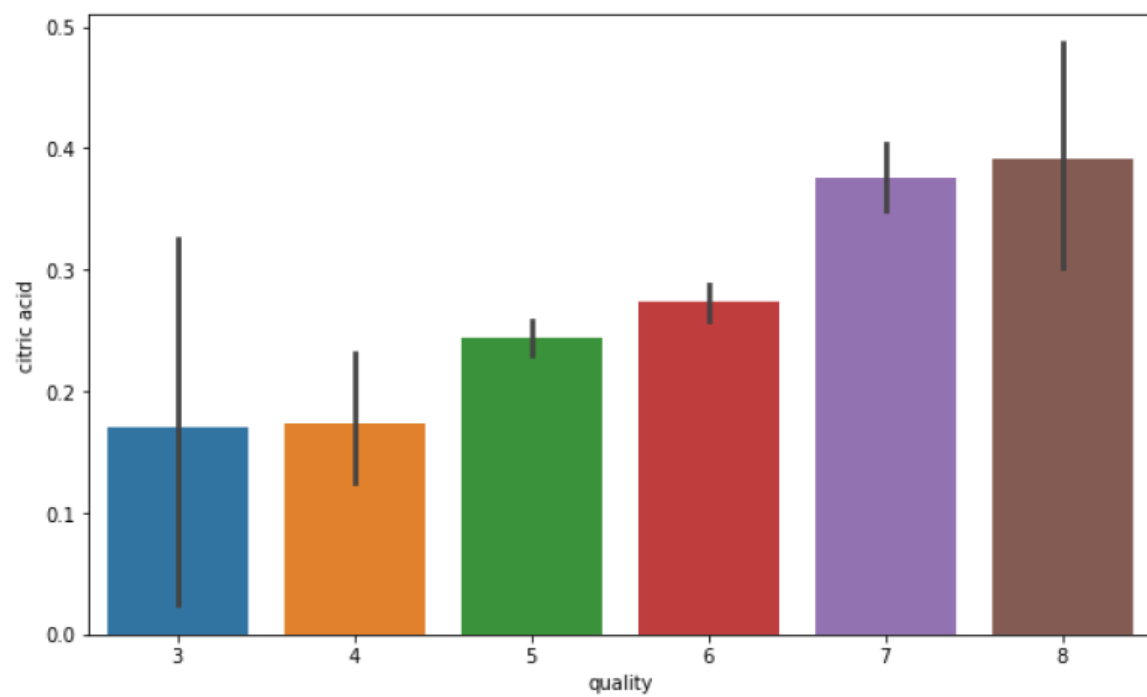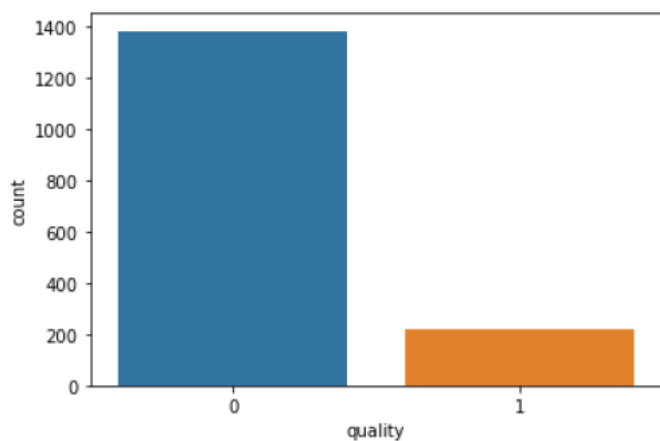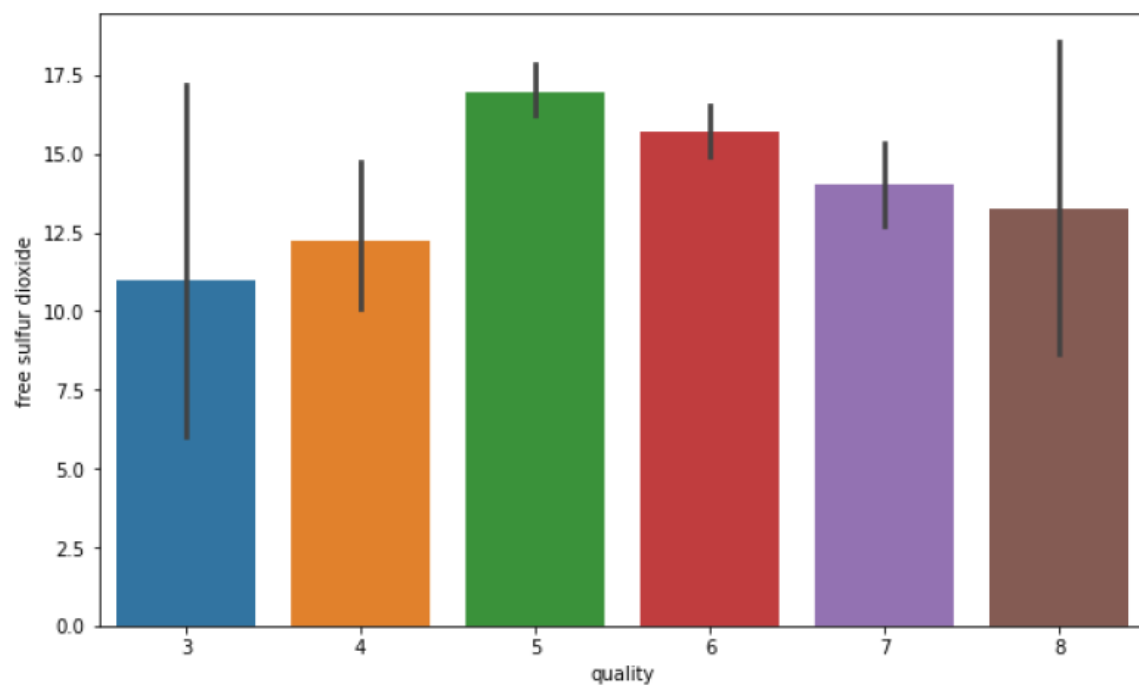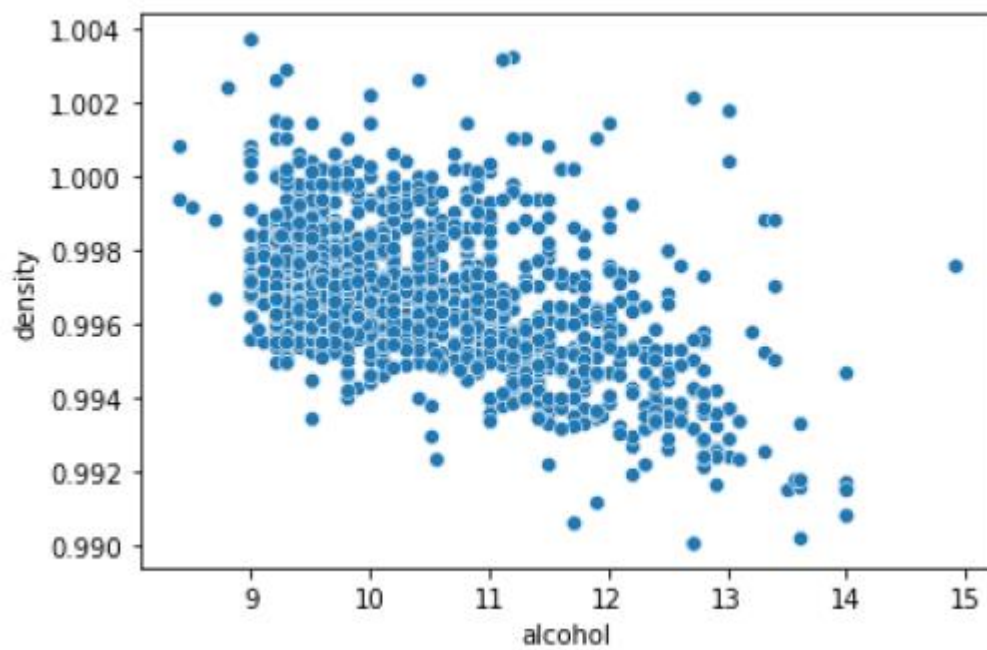247. master.mainloop()
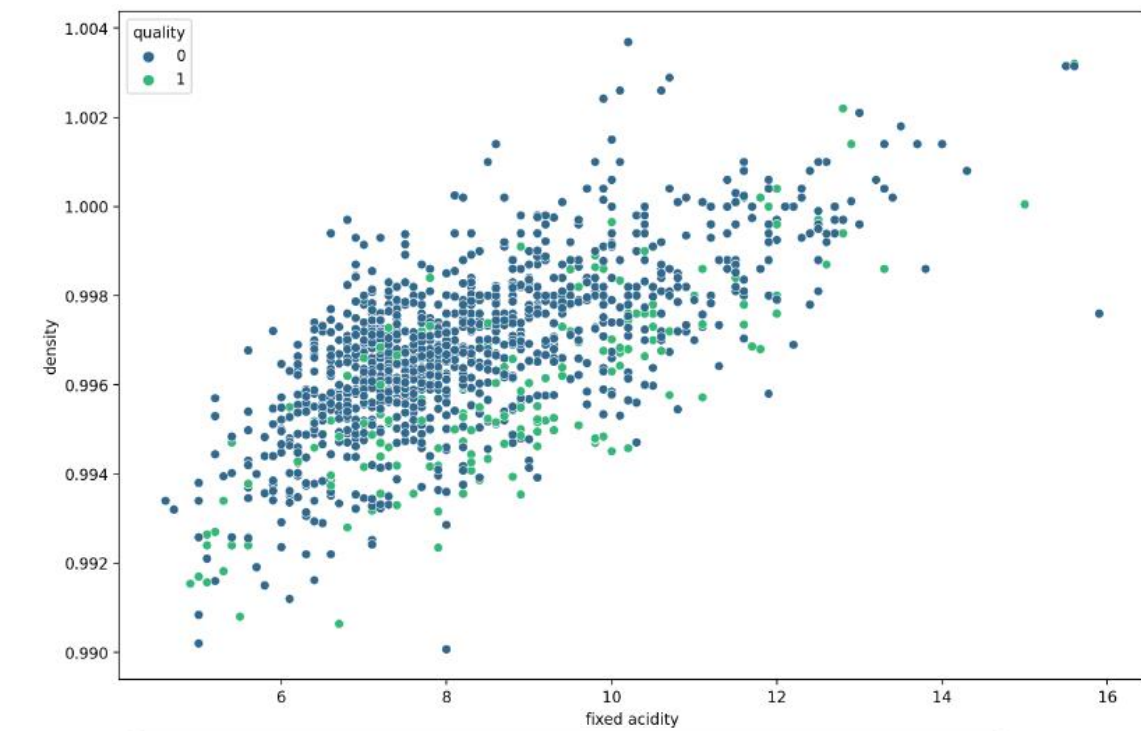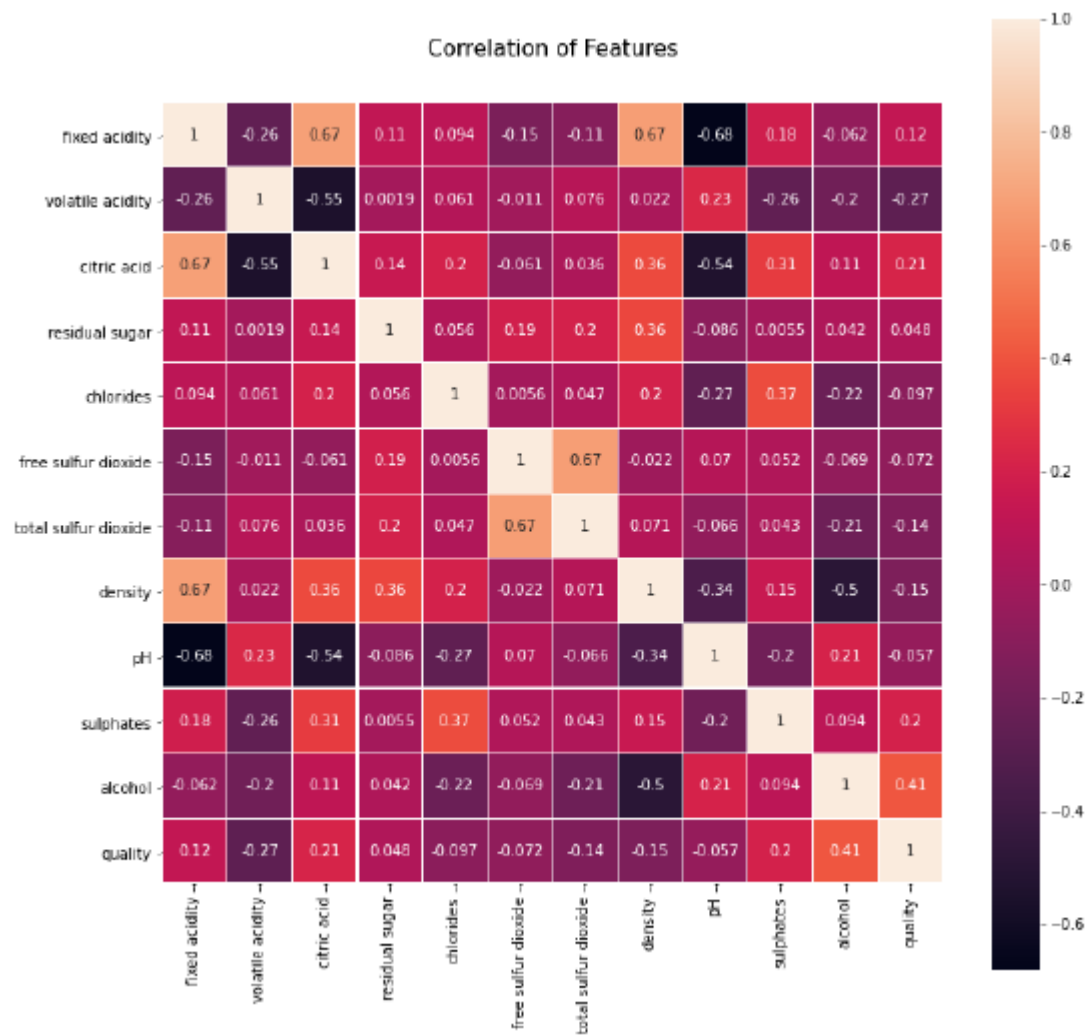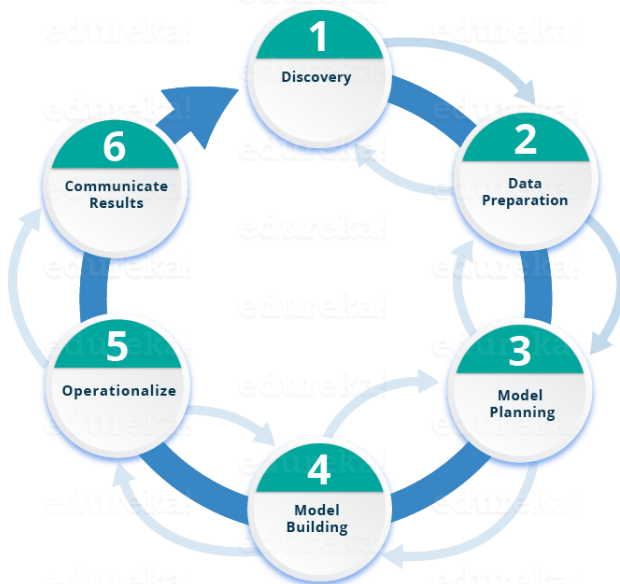
Correlation of Features

https://drive.google.com/file/d/1cDUxQXLCGLoTDxGWumK67doRO--fEfbN/view?usp=sharing

# Conclusion

Main phases of data science life cycle:-



Data science is emerging as a field that is revolutionizing science and industries alike. Work across nearly all domains is becoming more data driven, affecting both the jobs that are available and the skills that are required. As more data and ways of analyzing them become available, more aspects of the economy, society, and daily life will become dependent on data. It is imperative that educators, administrators, and students begin today to consider how to best prepare for and keep pace with this data-driven era of tomorrow. Undergraduate teaching, in particular, offers a critical link in offering more data science exposure to students and expanding the supply of data science talent. Customer data is key to making their lives better. Healthcare industries use the data available to them to assist their customers in their everyday life. Data Scientists in these type of industries have the purpose of analyzing the personal data, health history and create products that tackle the problems faced by customers.

From the above instances of data-centric companies, it is clear that each company uses data differently. The use of data varies as per company requirements. Therefore, the purpose of Data Scientists depends on the interests of the company.  the purpose of Data Science, we conclude that Data Scientists are the backbone of data-intensive companies. The purpose of Data Scientists is to extract, preprocess and analyze data. Through this, companies can make better decisions. Various companies have their own requirements and use data accordingly. In the end, the goal of Data Scientist to make businesses grow better. With the decisions and insights provided, the companies can adopt appropriate strategies and customize themselves for enhanced customer experience

# Bibliography

• Python 3.8.1 documentation, https://docs.python.org/3/

• Graphical User Interfaces with Tk, https://docs.python.org/3/library/tk.html

• Stack Overflow, https://stackoverflow.com

• Why Python is Good for Artificial Intelligence and Machine Learning, retrieved from https://djangostars.com/blog/

• Scikit-learn User Guide, https://scikit-learn.org/stable/user_guide.html