

SML_Assignment_3

Name: Ayush Patel

Collabration: Abhinav Nippani, Adwait Patil

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import itertools
4 import os
5 import sklearn
6 import matplotlib.pyplot as plt
7 from sklearn import tree
8 from sklearn.model_selection import train_test_split
9 from sklearn.linear_model import LinearRegression
10 from sklearn.model_selection import cross_val_score
11 import statsmodels.api as sm
12 from sklearn.metrics import r2_score
13 from sklearn.metrics import mean_squared_error
14 from sklearn.preprocessing import StandardScaler
15 from sklearn import linear_model
16 from sklearn.decomposition import PCA
17 from sklearn.cross_decomposition import PLSRegression
18 from sklearn.preprocessing import PolynomialFeatures
19 from sklearn.tree import DecisionTreeRegressor
20 from sklearn.ensemble import BaggingRegressor
21 from sklearn.ensemble import RandomForestRegressor
22 from sklearn.metrics import mean_squared_error
23 from sklearn.model_selection import StratifiedKFold
24 from sklearn.utils import resample
25 from sklearn.tree import plot_tree
26 import seaborn as sns
27 import warnings
28 warnings.filterwarnings("ignore")
29 from sklearn.model_selection import GridSearchCV
30 from sklearn.utils import resample
31 from sklearn.tree import DecisionTreeClassifier
32 from scipy import stats
33 from sklearn.metrics import f1_score
34 from sklearn.metrics import roc_auc_score
35 from sklearn.ensemble import RandomForestClassifier
36 import random
37 import scipy
38 from sklearn.preprocessing import OneHotEncoder
```

```
In [2]: 1 path = os.getcwd()
```

Problem 1

```
In [3]: 1 np.random.seed(123)
2 x = np.random.normal(0, 1, (200))
3 y = x + 2 * x**2 - 2 * x**3 + np.random.normal(0, 1, (200))
```

Bias column to ensure the columns are named properly

```
In [4]: 1 for polynomial in range(1,11):
2     polynomial_features = PolynomialFeatures(degree = polynomial, include_bias = True)
3     x_polynomial = polynomial_features.fit_transform(x.reshape(-1,1))
```

Removing the bias column

In [5]:

```
1 x_polynomial = pd.DataFrame(x_polynomial).iloc[:,1:]
2 display(x_polynomial)
```

	1	2	3	4	5	6	7	8	9	10
0	-1.085631	1.178594	-1.279518	1.389083	-1.508031	1.637165	-1.777356e+00	1.929553e+00	-2.094781e+00	2.274159e+00
1	0.997345	0.994698	0.992057	0.989424	0.986798	0.984178	9.815655e-01	9.789598e-01	9.763611e-01	9.737693e-01
2	0.282978	0.080077	0.022660	0.006412	0.001815	0.000513	1.453028e-04	4.111758e-05	1.163539e-05	3.292565e-06
3	-1.506295	2.268924	-3.417668	5.148015	-7.754428	11.680454	-1.759421e+01	2.650206e+01	-3.991991e+01	6.013115e+01
4	-0.578600	0.334778	-0.193703	0.112076	-0.064847	0.037521	-2.170953e-02	1.256114e-02	-7.267877e-03	4.205195e-03
...
195	-3.231055	10.439716	-33.731298	108.987680	-352.145189	1137.800476	-3.676296e+03	1.187831e+04	-3.837949e+04	1.240062e+05
196	-0.269293	0.072519	-0.019529	0.005259	-0.001416	0.000381	-1.027025e-04	2.765711e-05	-7.447880e-06	2.005666e-06
197	-0.110851	0.012288	-0.001362	0.000151	-0.000017	0.000002	-2.056694e-07	2.279860e-08	-2.527241e-09	2.801465e-10
198	-0.341262	0.116460	-0.039743	0.013563	-0.004628	0.001580	-5.390301e-04	1.839503e-04	-6.277520e-05	2.142277e-05
199	-0.217946	0.047501	-0.010353	0.002256	-0.000492	0.000107	-2.335856e-05	5.090910e-06	-1.109545e-06	2.418211e-07

200 rows × 10 columns

a)

In [6]:

```
1 list_of_columns = list(x_polynomial.columns)
2
3 subset_selection = {"Features" : [], "AIC" : [], "BIC" : [], "Adjusted_R_Square" : []}
4
5 for A in range(len(list_of_columns) + 1):
6     for subset in itertools.combinations(list_of_columns, A):
7
8         if(list(subset)!=[]):
9
10            x_polynomial_subset = x_polynomial[list(subset)]
11            subset_selection["Features"].append(list(subset))
12
13            regr = sm.OLS(y, sm.add_constant(x_polynomial_subset)).fit()
14
15            aic = regr.aic
16            subset_selection["AIC"].append(aic)
17            bic = regr.bic
18            subset_selection["BIC"].append(bic)
19            adj_r_square = regr.rsquared_adj
20            subset_selection["Adjusted_R_Square"].append(adj_r_square)
```

In [7]:

```
1 subset_selection = pd.DataFrame(subset_selection)
2 display(subset_selection)
```

	Features	AIC	BIC	Adjusted_R_Square
0	[1]	1356.768181	1363.364816	0.421163
1	[2]	1387.524884	1394.121519	0.324938
2	[3]	1065.501335	1072.097970	0.865079
3	[4]	1330.869591	1337.466225	0.491468
4	[5]	1063.848982	1070.445617	0.866189
...
1018	[1, 2, 3, 5, 6, 7, 8, 9, 10]	552.240533	585.223707	0.990030
1019	[1, 2, 4, 5, 6, 7, 8, 9, 10]	558.460353	591.443527	0.989715
1020	[1, 3, 4, 5, 6, 7, 8, 9, 10]	558.104425	591.087599	0.989733
1021	[2, 3, 4, 5, 6, 7, 8, 9, 10]	557.049808	590.032982	0.989787
1022	[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]	552.185498	588.466989	0.990079

1023 rows × 4 columns

Finding features with minimum AIC.

In [8]:

```
1 subset_selection[subset_selection["AIC"] == subset_selection["AIC"].min()]
```

Out[8]:

	Features	AIC	BIC	Adjusted_R_Square
55	[1, 2, 3]	542.275105	555.468374	0.990236

Finding features with minimum BIC.

In [9]: 1 subset_selection[subset_selection["BIC"] == subset_selection["BIC"].min()]

Out[9]:

	Features	AIC	BIC	Adjusted_R_Square
55	[1, 2, 3]	542.275105	555.468374	0.990236

Finding features with maximum Adjusted R square.

In [10]: 1 subset_selection[subset_selection["Adjusted_R_Square"] == subset_selection["Adjusted_R_Square"].max()]

Out[10]:

	Features	AIC	BIC	Adjusted_R_Square
55	[1, 2, 3]	542.275105	555.468374	0.990236

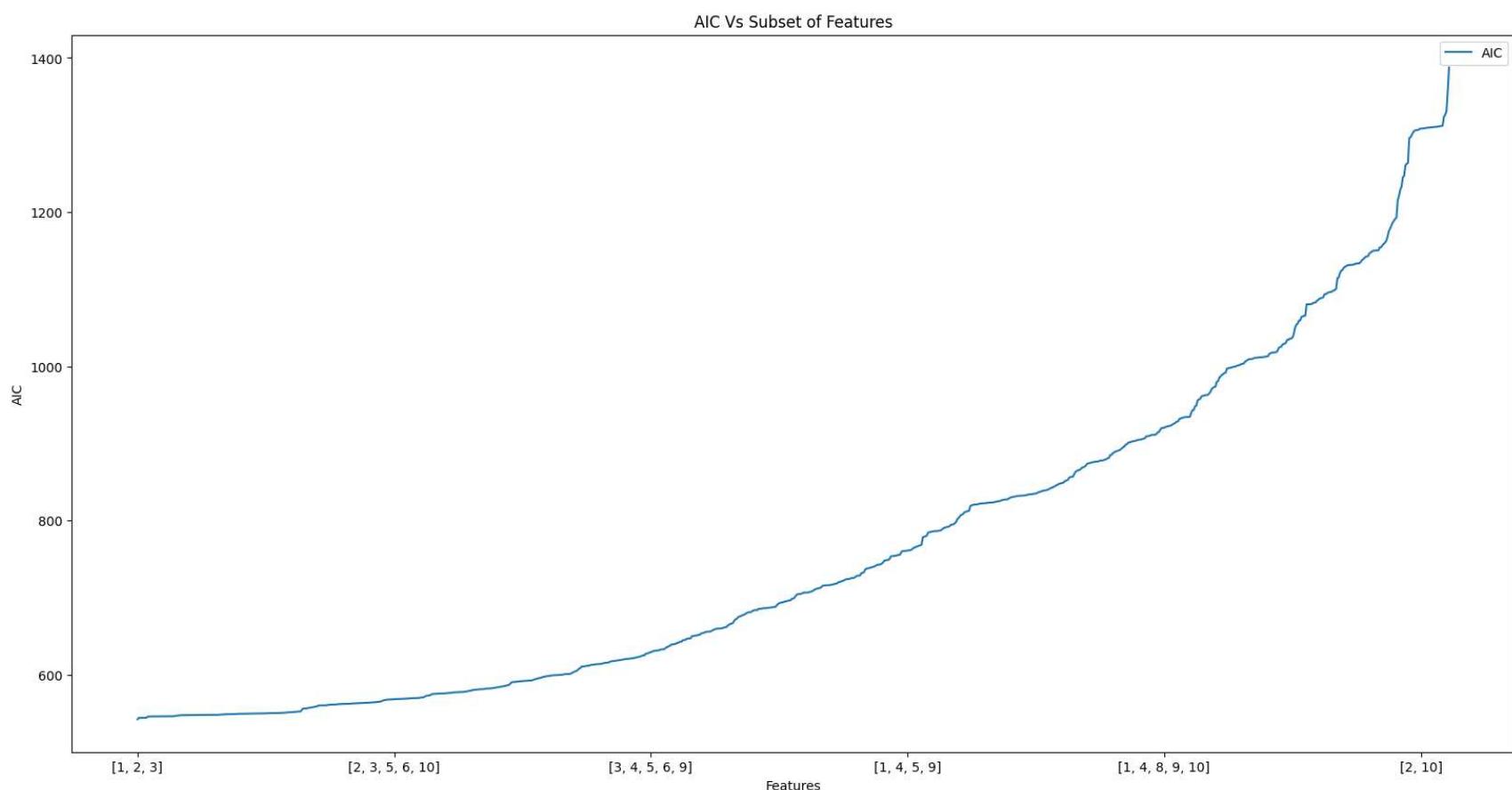
Plotting a graph of AIC vs Subset of Features

In [11]:

```

1 subset_selection.sort_values(["AIC"]).plot(x="Features",y="AIC",figsize=(20,10))
2 plt.xlabel("Features")
3 plt.ylabel("AIC")
4 plt.title("AIC Vs Subset of Features")
5 plt.show()

```



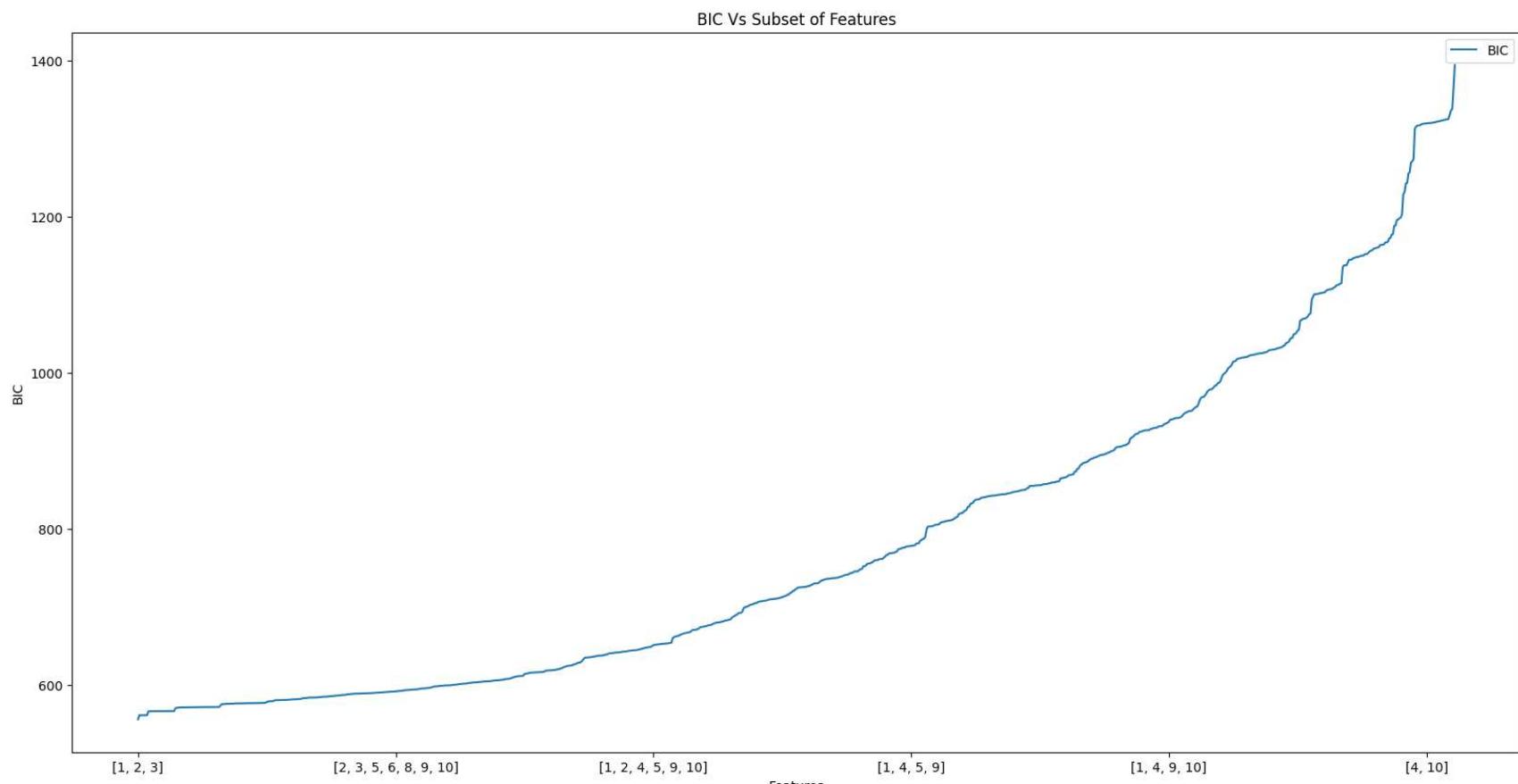
Plotting a graph of BIC vs Subset of Features

In [12]:

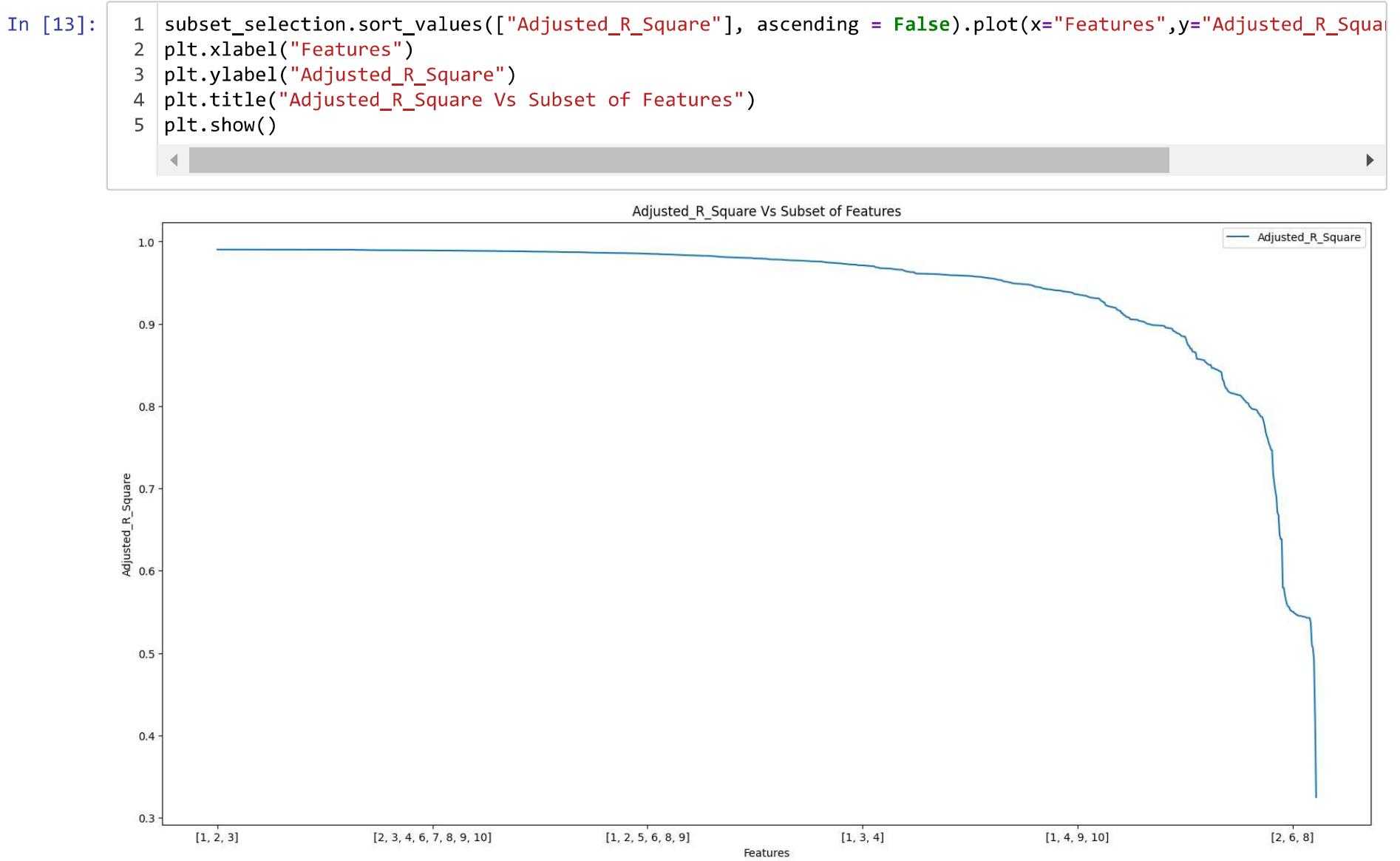
```

1 subset_selection.sort_values(["BIC"]).plot(x="Features",y="BIC",figsize=(20,10))
2 plt.xlabel("Features")
3 plt.ylabel("BIC")
4 plt.title("BIC Vs Subset of Features")
5 plt.show()

```



Plotting a graph of Adjusted_R_Square vs Subset of Features

**Observation:**

All the above graphs plotted indicate that the optimal features are [1,2,3]

Filtering the optimal features, then training the model using OLS

In [14]:

```

1 best_features = [1,2,3]
2 best_x_polynomial = x_polynomial[best_features]
3 regr = sm.OLS(y, sm.add_constant(best_x_polynomial)).fit()

```

Coefficients of the best model

In [15]:

```
1 pd.DataFrame(regr.params).T
```

Out[15]:

	const	1	2	3
0	-0.181156	0.914572	2.06306	-1.970858

b)

In [16]:

```

1 list_of_columns = list(x_polynomial.columns)
2 final_column = []
3 best_adj_r_square = -np.inf
4 for i in range(len(list_of_columns)):
5     for column in list_of_columns:
6
7         if(column not in final_column):
8             x_polynomial_subset = x_polynomial[final_column + [column]]
9             regr = sm.OLS(y, sm.add_constant(x_polynomial_subset)).fit()
10            adj_r_square = regr.rsquared_adj
11
12            if(adj_r_square > best_adj_r_square):
13                #print(adj_r_square, column)
14                best_adj_r_square = adj_r_square
15                best_column = column
16
17            if(best_column not in final_column):
18                final_column.append(best_column)

```

In [17]:

```
1 final_column
```

Out[17]:

[5, 2, 3, 1]

Observation:

We can see that after every iteration, one feature that contributes the most is being added to the selection which can be seen with the increase in adjusted R² values at every step. The selection stops after 4 iterations since there is no improvement in the model performance after 4 features i.e; [5,2,3,1]

Although, from subset selection, we know that it is not the most optimal model, there is not a significant difference in the performance between the models trained on the obtained subsets (best subset selection and forward subset selection). However, the computational time significantly reduces in the case of forward subset selection

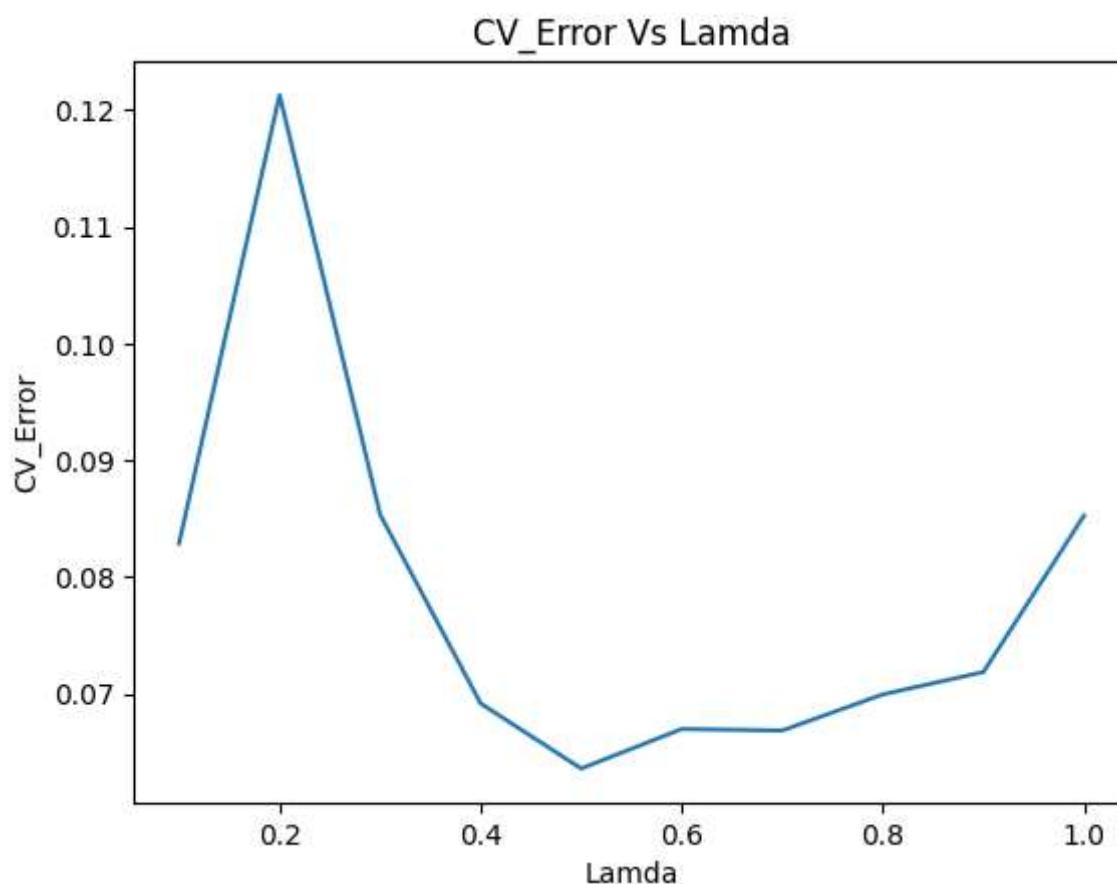
c)

```
In [18]: 1 cv_error_data = {"Lamda":[], "CV_Error":[]}
2
3
4 for alpha in np.arange(0.1,1.1,0.1):
5     #print(alpha)
6     clf = linear_model.Lasso(alpha=alpha,max_iter = 100000000)
7
8     best_cv_error = 1 - cross_val_score(clf, x_polynomial, y, scoring = "r2").mean()
9
10    cv_error_data["Lamda"].append(alpha)
11    cv_error_data["CV_Error"].append(best_cv_error)
12
13
14 cv_error_data = pd.DataFrame(cv_error_data)
```

```
In [19]: 1 cv_error_data=cv_error_data["CV_Error"] == cv_error_data["CV_Error"].min()
```

```
Out[19]: Lamda  CV_Error
4      0.5  0.063629
```

```
In [20]: 1 plt.plot(cv_error_data["Lamda"],cv_error_data["CV_Error"])
2 plt.xlabel("Lamda")
3 plt.ylabel("CV_Error")
4 plt.title("CV_Error Vs Lamda")
5 plt.show()
```



Filtering the best features and printing their co-efficients

```
In [21]: 1 regr = linear_model.Lasso(alpha=0.5,max_iter = 1000000).fit(x_polynomial,y)
2 coefficient = pd.DataFrame(regr.coef_).T
3 coefficient.columns = x_polynomial.columns
4 coefficient
```

```
Out[21]: 1  2   3   4   5   6   7   8   9   10
0  -0.0  0.0 -0.0  0.670943 -0.635998  0.0  0.048937 -0.021227  0.00071  0.001987
```

```
In [22]: 1 regr.intercept_
```

```
Out[22]: 0.5804372836802458
```

Problem 2

```
In [23]: 1 train = pd.read_csv('D:\SML\Train_Data.csv')
2 test = pd.read_csv('D:\SML\Test_Data.csv')
```

In [24]: 1 train

Out[24]:

	age	sex	bmi	smoker	region	children	charges
0	21.000000	male	25.745000	no	northeast	2	3279.868550
1	36.976978	female	25.744165	yes	southeast	3	21454.494239
2	18.000000	male	30.030000	no	southeast	1	1720.353700
3	37.000000	male	30.676891	no	northeast	3	6801.437542
4	58.000000	male	32.010000	no	southeast	1	11946.625900
...
3625	48.820767	female	41.426984	no	northwest	4	10987.324964
3626	38.661977	female	26.202557	no	southeast	2	11735.844352
3627	56.000000	male	40.300000	no	southwest	0	10602.385000
3628	48.061207	female	34.930624	no	southeast	1	8976.140452
3629	37.598865	female	25.219233	no	northeast	3	7027.698968

3630 rows × 7 columns

In [25]: 1 X = train.drop(columns=["charges"])
2 y = train[["charges"]]

Performing One Hot Encoding

In [26]:

```

1 def encoding(df,col_name):
2     dummy = pd.get_dummies(df[col_name]).iloc[:,1:]
3     dummy.columns = [col_name+"_"+str(x) for x in dummy.columns]
4
5     df = df.drop(col_name, axis = 1)
6
7     df = df.join(dummy)
8
9
10    return df

```

In [27]: 1 X = encoding(X,"sex")
2 X = encoding(X,"smoker")
3 X = encoding(X,"region")

a)

In [28]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)

b)

In [29]: 1 regr = sm.OLS(y_train, X_train).fit()
2 y_pred = regr.predict(X_test)

In [30]: 1 print("The mean squared error(test error) is: {}".format(mean_squared_error(y_pred, y_test, squared=False)))

The mean squared error(test error) is: 5860.094769311387

c)

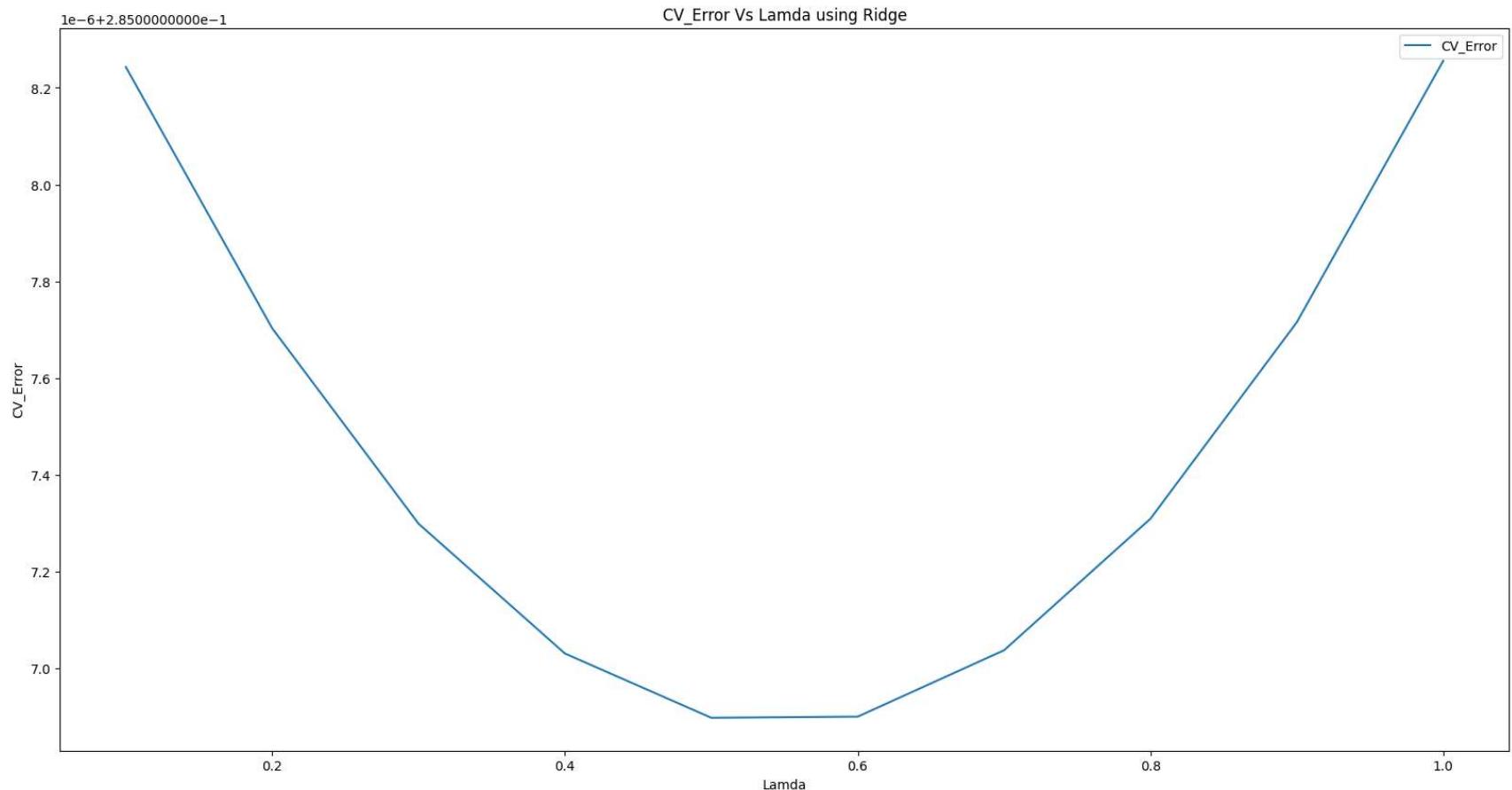
```
In [31]: 1 cross_val = {"Lamda": [], "CV_Error": []}
2
3 for lamda in np.arange(0.1, 1.1, 0.1):
4
5     regr = linear_model.Ridge(alpha = lamda)
6
7     cv_error = 1 - cross_val_score(regr, X_train, y_train, scoring="r2").mean()
8
9     cross_val["Lamda"].append(lamda)
10    cross_val["CV_Error"].append(cv_error)
11
12 cross_val = pd.DataFrame(cross_val)
13 cross_val
```

Out[31]:

	Lamda	CV_Error
0	0.1	0.285008
1	0.2	0.285008
2	0.3	0.285007
3	0.4	0.285007
4	0.5	0.285007
5	0.6	0.285007
6	0.7	0.285007
7	0.8	0.285007
8	0.9	0.285008
9	1.0	0.285008

Plotting the graph to find the minimal cv error

```
In [32]: 1 cross_val.plot(x="Lamda", y="CV_Error", figsize=(20,10))
2 plt.xlabel("Lamda")
3 plt.ylabel("CV_Error")
4 plt.title("CV_Error Vs Lamda using Ridge")
5 plt.show()
```



```
In [33]: 1 min_cv_error = cross_val[cross_val["CV_Error"] == cross_val["CV_Error"].min()]
2 min_cv_error
```

Out[33]:

	Lamda	CV_Error
4	0.5	0.285007

The optimal value of lamda obtained is 0.5

```
In [34]: 1 regr = linear_model.Ridge(alpha=min_cv_error["Lamda"].values[0]).fit(X_train,y_train)
2
3 y_pred = regr.predict(X_test)
```

```
In [35]: 1 print("The mean squared error(test error) for Ridge is: {}".format(mean_squared_error(y_pred, y_test, square
```

The mean squared error(test error) for Ridge is: 5671.667330077637

d)

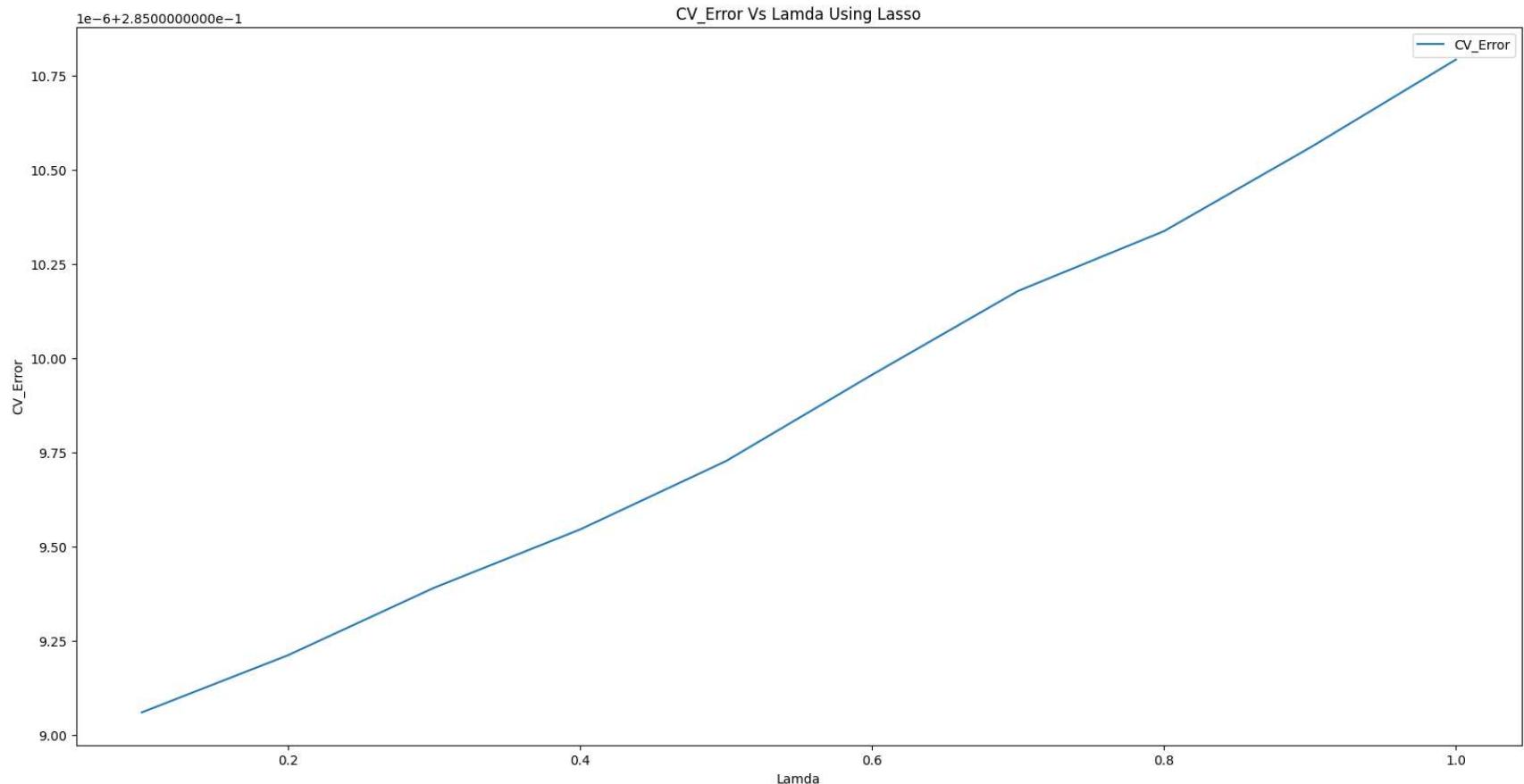
```
In [36]: 1 cross_val = {"Lamda":[], "CV_Error":[]}
2
3 for lamda in np.arange(0.1, 1.1, 0.1):
4     regr = linear_model.Lasso(alpha = lamda)
5
6     cv_error = 1 - cross_val_score(regr, X_train, y_train, scoring="r2").mean()
7
8     cross_val["Lamda"].append(lamda)
9     cross_val["CV_Error"].append(cv_error)
10
11 cross_val = pd.DataFrame(cross_val)
12 cross_val
```

Out[36]: Lamda CV_Error

	Lamda	CV_Error
0	0.1	0.285009
1	0.2	0.285009
2	0.3	0.285009
3	0.4	0.285010
4	0.5	0.285010
5	0.6	0.285010
6	0.7	0.285010
7	0.8	0.285010
8	0.9	0.285011
9	1.0	0.285011

Plotting the graph to find the minimal cv error

```
In [37]: 1 cross_val.plot(x="Lamda", y="CV_Error", figsize=(20,10))
2 plt.xlabel("Lamda")
3 plt.ylabel("CV_Error")
4 plt.title("CV_Error Vs Lamda Using Lasso")
5 plt.show()
```



```
In [38]: 1 min_cv_error = cross_val[cross_val["CV_Error"] == cross_val["CV_Error"].min()]
2 min_cv_error
```

Out[38]: Lamda CV_Error

	Lamda	CV_Error
0	0.1	0.285009

The optimal value of lamda obtained is 0.1

```
In [39]: 1 regr = linear_model.Lasso(alpha=min_cv_error["Lamda"].values[0]).fit(X_train,y_train)
2
3 y_pred = regr.predict(X_test)
```

```
In [40]: 1 print("The mean squared error(test error) for Lasso is: {}".format(mean_squared_error(y_pred, y_test, square
```

The mean squared error(test error) for Lasso is: 5670.879834633087

```
In [41]: 1 coefficient = pd.DataFrame(regr.coef_).
2 coefficient.columns = X_train.columns
3 coefficient
```

```
Out[41]:      age        bmi   children  sex_male  smoker_yes  region_northwest  region_southeast  region_southwest
0  235.507643  292.260956  507.388153  960.624686  22623.46324       -2276.372029       -1908.611075       -2319.871851
```

There are 8 non-zero coefficient estimates

```
In [42]: 1 regr.intercept_
```

```
Out[42]: array([-9038.37916747])
```

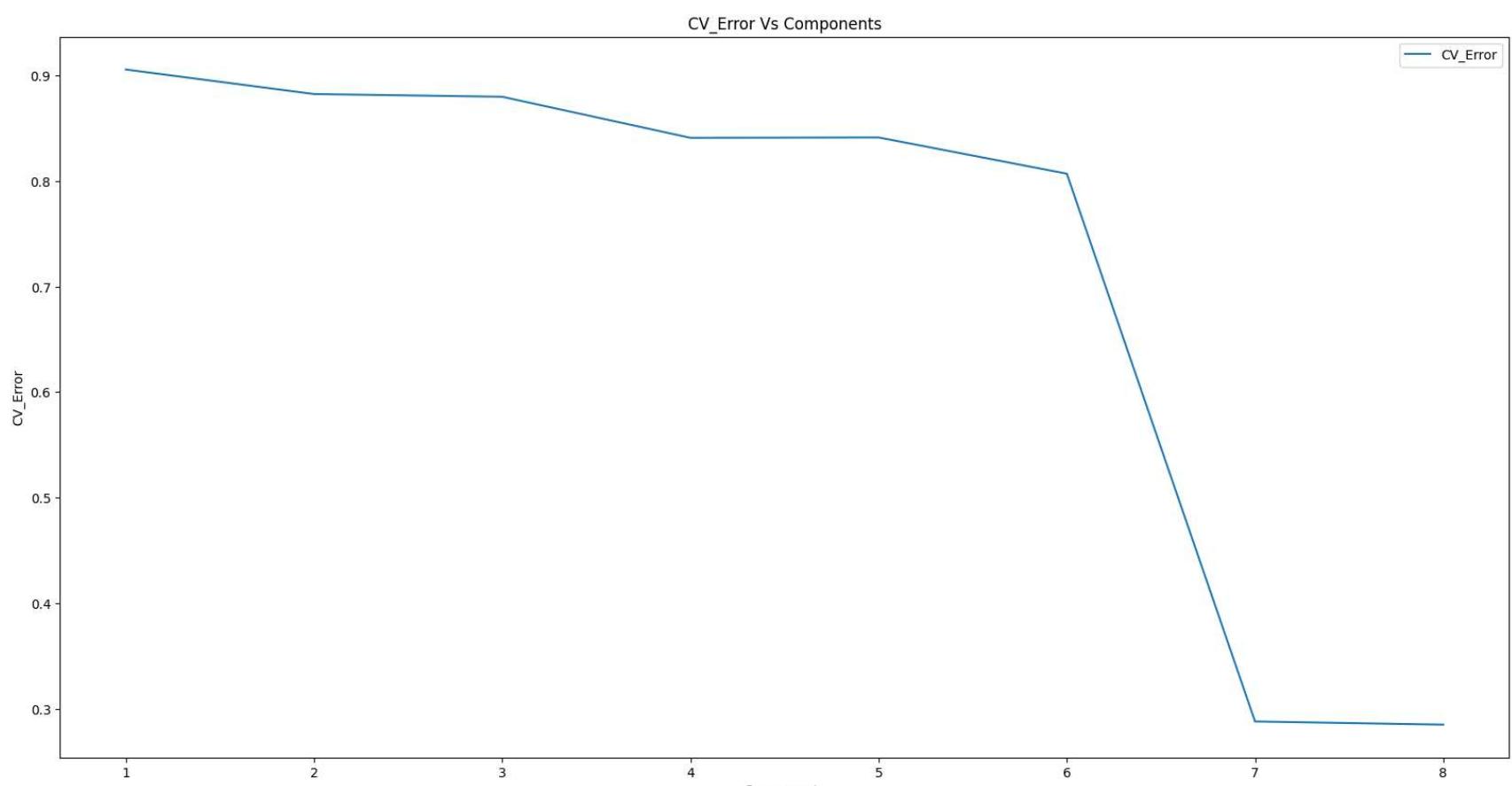
e)

```
In [43]: 1 cross_val = {"Components":[], "CV_Error":[]}
2
3 for c in range(1,X_train.shape[1]+1):
4     pca = PCA(n_components = c)
5     X_train_reduced = pca.fit_transform(X_train)
6     regr = linear_model.LinearRegression()
7     cv_error = 1 - cross_val_score(regr,X_train_reduced,y_train, scoring = "r2").mean()
8
9     cross_val["Components"].append(c)
10    cross_val["CV_Error"].append(cv_error)
11
12 cross_val = pd.DataFrame(cross_val)
13 cross_val
```

	Components	CV_Error
0	1	0.905684
1	2	0.882465
2	3	0.879893
3	4	0.840989
4	5	0.841323
5	6	0.806988
6	7	0.288013
7	8	0.285009

Plotting the graph to find the minimal cv error

```
In [44]: 1 cross_val.plot(x="Components", y="CV_Error", figsize=(20,10))
2 plt.xlabel("Components")
3 plt.ylabel("CV_Error")
4 plt.title("CV_Error Vs Components")
5 plt.show()
```



By plotting a graph between the Number of Principal Components and the RMSE Score we can see that from 6 onwards there is a significant drop in RMSE and from 7 and 8 the value does not change much so we can safely say that it does not vary much when it is at 7 or 8.

```
In [45]: 1 min_cv_error = cross_val[cross_val["CV_Error"] == cross_val["CV_Error"].min()]
2 min_cv_error
```

```
Out[45]: Components CV_Error
7 8 0.285009
```

The optimal component obtained is 8.

```
In [46]: 1 pca = PCA(n_components = min_cv_error["Components"].values[0])
2 X_train_reduced = pca.fit_transform(X_train)
3 X_test_reduced = pca.transform(X_test)
4 regr = linear_model.LinearRegression().fit(X_train_reduced, y_train)
5 y_pred = regr.predict(X_test_reduced)
```

```
In [47]: 1 print("The mean squared error(test error) for PCA is: {}".format(mean_squared_error(y_pred, y_test, squared=True)))
```

The mean squared error(test error) for PCA is: 5670.859881750769

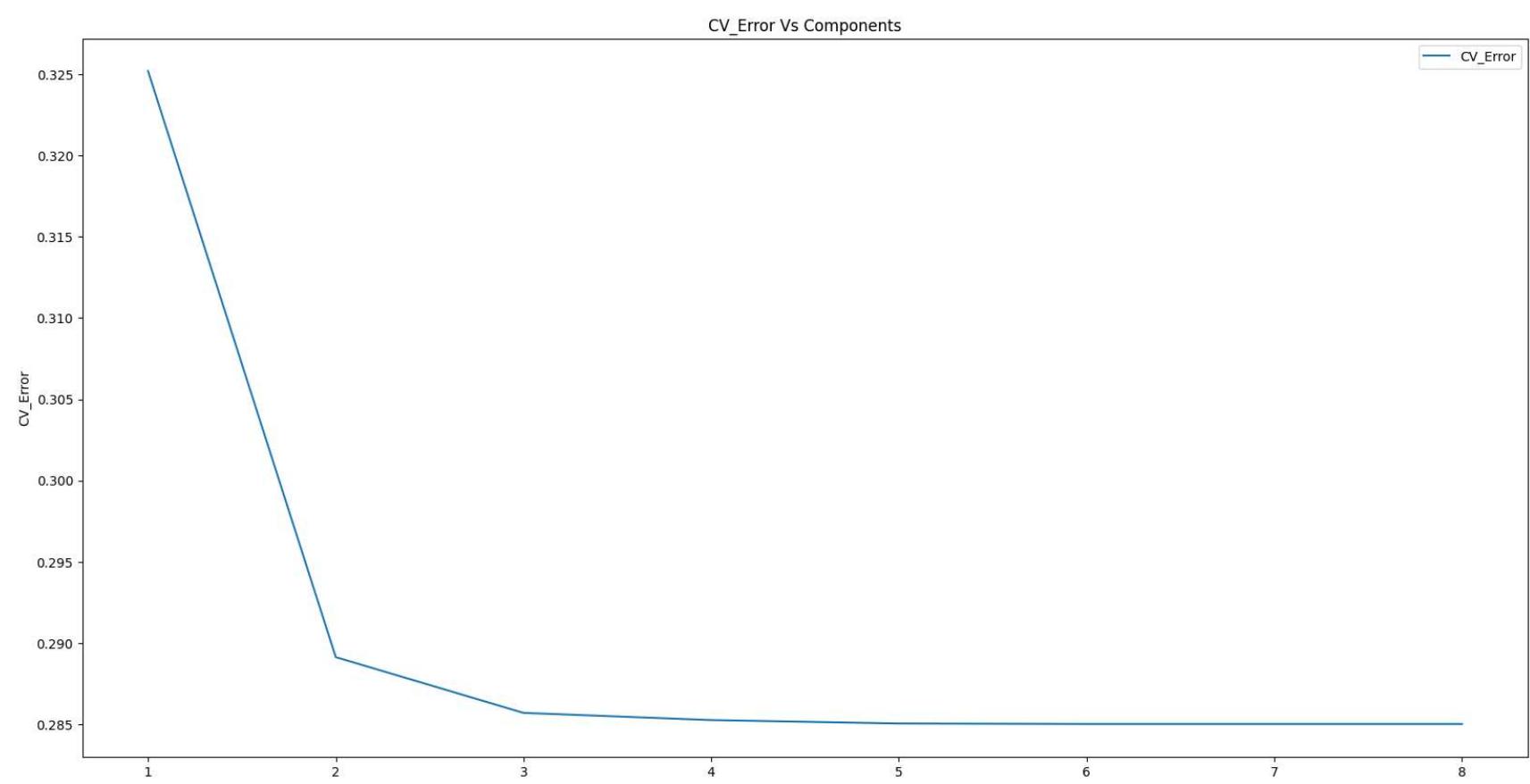
f)

```
In [48]: 1 cross_val = {"Components": [], "CV_Error": []}
2
3 # Iterating through no of components
4 for c in range(1, X_train.shape[1]+1):
5     pls = PLSRegression(n_components = c)
6     cv_error = 1 - cross_val_score(pls, X_train, y_train, scoring = "r2").mean()
7
8     cross_val["Components"].append(c)
9     cross_val["CV_Error"].append(cv_error)
10
11 cross_val = pd.DataFrame(cross_val)
12 cross_val
```

```
Out[48]: Components CV_Error
0 1 0.325200
1 2 0.289125
2 3 0.285693
3 4 0.285250
4 5 0.285040
5 6 0.285010
6 7 0.285009
7 8 0.285009
```

Plotting the graph to find the minimal cv error

```
In [49]: 1 cross_val.plot(x="Components", y="CV_Error", figsize=(20,10))
2 plt.xlabel("Components")
3 plt.ylabel("CV_Error")
4 plt.title("CV_Error Vs Components")
5 plt.show()
```



```
In [50]: 1 min_cv_error = cross_val[cross_val["CV_Error"] == cross_val["CV_Error"].min()]
2 min_cv_error
```

Out[50]: Components CV_Error

	Components	CV_Error
7	8	0.285009

The optimal component obtained is 8.

```
In [51]: 1 pls = PLSRegression(n_components = min_cv_error["Components"].values[0])
2 regr.fit(X_train,y_train)
3
4 y_pred = regr.predict(X_test)
```

```
In [52]: 1 print("The mean squared error(test error) for PLSRegression is: {}".format(mean_squared_error(y_pred, y_te
```

The mean squared error(test error) for PLSRegression is: 5670.859881750768

Problem 3

```
In [53]: 1 df = pd.read_csv("D:\SML\College.csv")
```

```
In [54]: 1 df['Accept.Rate'] = df['Accept'] / df['Apps']
```

```
In [55]: 1 df.head()
```

Out[55]:

	Unnamed: 0	Private	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	I
0	Abilene Christian University	Yes	1660	1232	721	23	52	2885	537	7440	3300	450	2200	
1	Adelphi University	Yes	2186	1924	512	16	29	2683	1227	12280	6450	750	1500	
2	Adrian College	Yes	1428	1097	336	22	50	1036	99	11250	3750	400	1165	
3	Agnes Scott College	Yes	417	349	137	60	89	510	63	12960	5450	450	875	
4	Alaska Pacific University	Yes	193	146	55	16	44	249	869	7560	4120	800	1500	

```
In [56]: 1 df = pd.get_dummies(df, columns = ['Private'])
```

a)

```
In [57]: 1 X = df.drop(['Accept.Rate', 'Apps', 'Accept', 'Unnamed: 0', 'Private_No'], axis = 1)
2 y = df['Accept.Rate']
```

```
In [58]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
```

```
In [59]: 1 X_train
```

Out[59]:

	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ratio	perc.alu
739	478	12	25	2138	227	4470	2890	600	1210	33	33	16.3	
133	249	23	57	1698	894	9990	5666	800	1500	66	71	14.3	
234	198	7	20	545	42	11750	2700	400	850	77	83	14.0	
55	156	25	55	421	27	6500	2700	500	1000	76	76	14.3	
639	588	56	86	1846	154	9843	3180	600	1500	74	78	14.6	
...	
71	313	71	95	1088	16	18165	6750	500	1200	100	100	12.3	
106	288	55	82	943	7	11850	4270	600	900	95	99	11.4	
270	471	55	86	1818	23	14360	4090	400	650	77	92	12.9	
435	351	22	44	1419	228	6400	3150	500	1900	58	64	16.2	
102	902	6	24	6394	3881	5962	4444	500	985	69	73	16.7	

621 rows × 16 columns

b)

```
In [60]: 1 model = DecisionTreeRegressor()
2 model.fit(X_train, y_train)
```

Out[60]: DecisionTreeRegressor()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

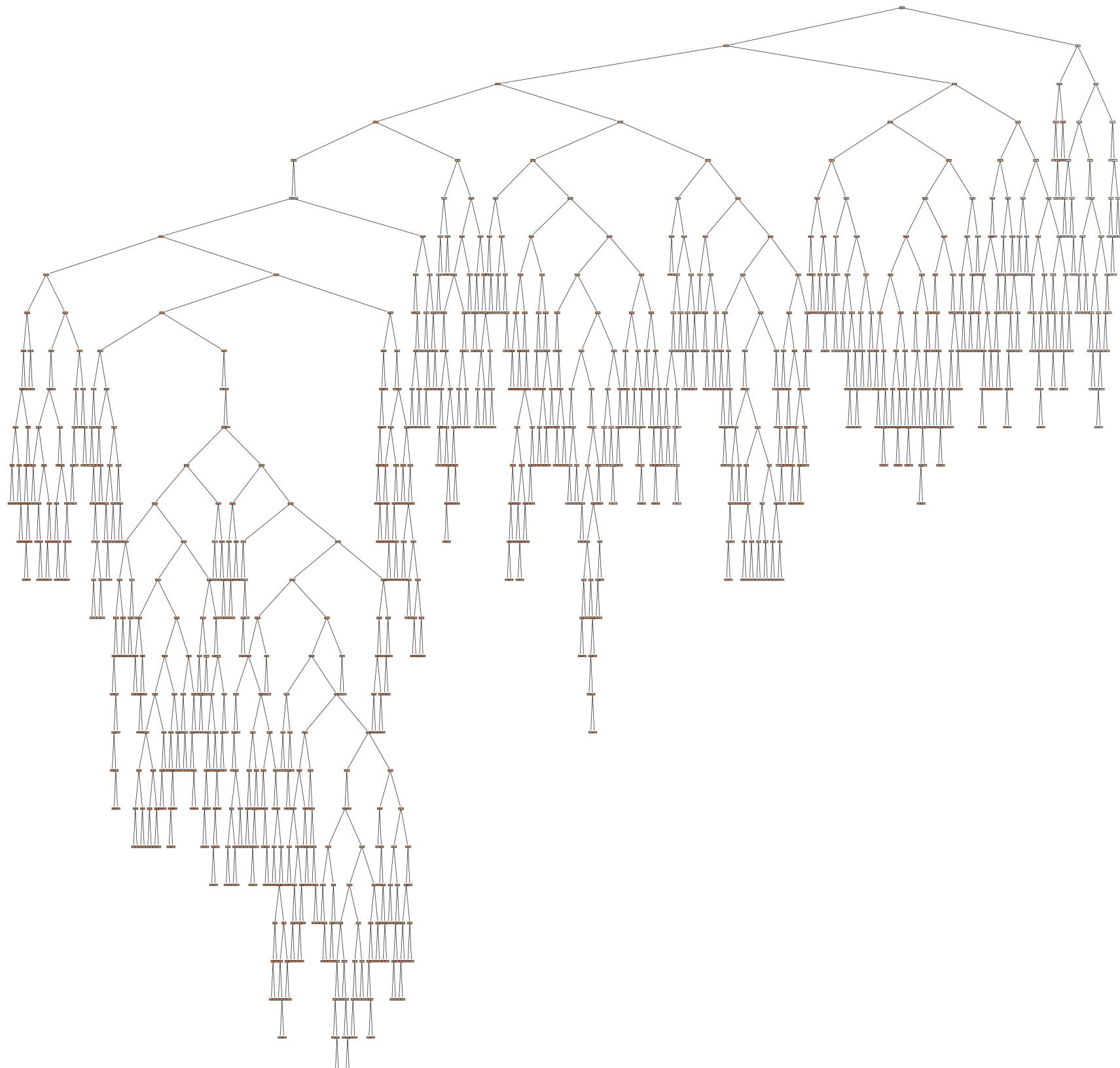
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [61]: 1 y_predictions = model.predict(X_test)
```

```
In [62]: 1 print("The mean squared error for Decision Tree Regressor is: {}".format(mean_squared_error(y_predictions, y_test)))
```

The mean squared error for Decision Tree Regressor is: 0.027824697997187

```
In [63]: 1 fig = plt.figure(figsize = (50, 50))
2 plot_tree(model, feature_names = X_train.columns, filled = True)
3 plt.show()
```



```
In [64]: 1 print("The mean squared error of the test set is: {}".format(mean_squared_error(y_test, y_predictions, squared=False)))
```

The mean squared error of the test set is: 0.027824697997187

Observation:

The tree plotted above has a large depth which makes it difficult for visualization.

Since, it is a big tree, there are a large number of parameters which might mean that the model is overfitting on the data. By reducing the depth, the performance of the model might improve further.

When we just pick a decision tree without specifying the depth, it would go till the maximum depth until values don't split further. As a result, we see that the decision tree is fully grown and it takes into account many features before coming to a conclusion for the acceptance rate.

The mean squared error of the test set is 0.0278246 respectively

c)

```
In [65]: 1 parameters = {'max_depth':range(1,30), 'min_samples_split':range(3,30), 'min_samples_leaf':range(3,30)}
2 clf = GridSearchCV(model, parameters, scoring='neg_mean_squared_error')
3 clf.fit(X_train, y_train)
```

```
Out[65]: GridSearchCV(estimator=DecisionTreeRegressor(),
param_grid={'max_depth': range(1, 30),
'min_samples_leaf': range(3, 30),
'min_samples_split': range(3, 30)},
scoring='neg_mean_squared_error')
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [66]: 1 clf.best_params_
```

```
Out[66]: {'max_depth': 3, 'min_samples_leaf': 26, 'min_samples_split': 22}
```

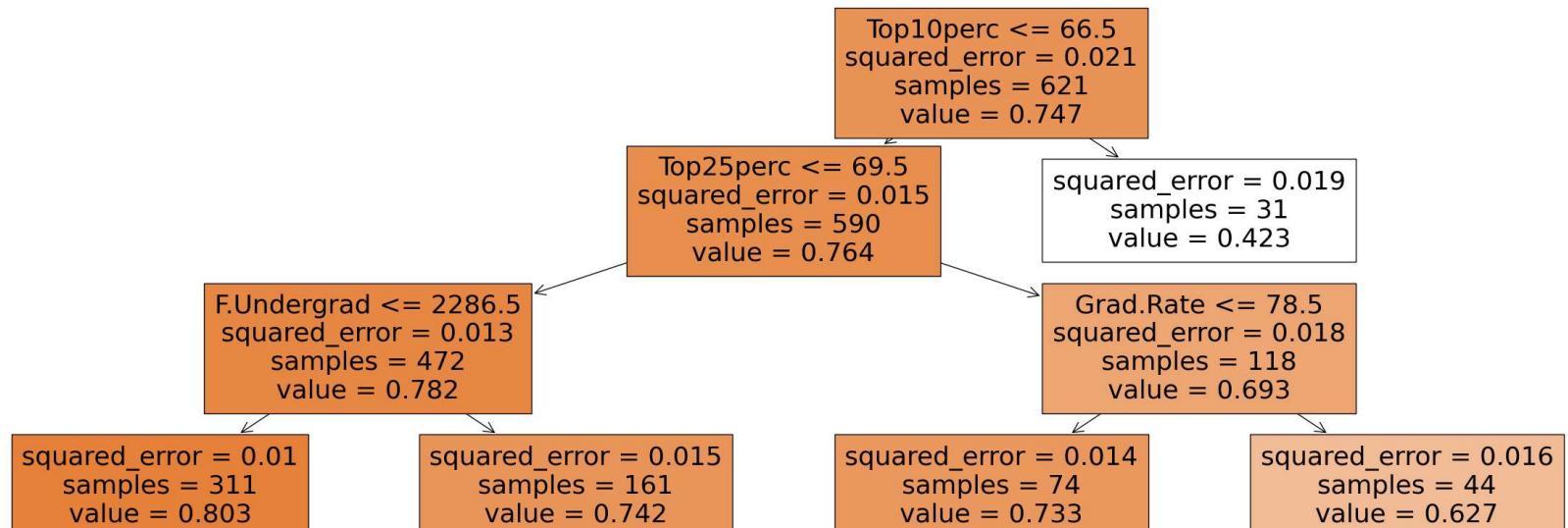
```
In [67]: 1 model = DecisionTreeRegressor(min_samples_split = clf.best_params_['min_samples_leaf'], min_samples_leaf =
```

```
In [68]: 1 model.fit(X_train, y_train)
2 y_predictions = model.predict(X_test)
```

```
In [69]: 1 print("The mean squared error of the best tree is: {}".format(mean_squared_error(y_predictions, y_test)), sq
```

The mean squared error of the best tree is: 0.02005733000440076

```
In [70]: 1 fig = plt.figure(figsize = (30, 10))
2 plot_tree(model, feature_names = X_train.columns, filled = True)
3 plt.show()
```



Observation:

Based on the results from getting the best tree, we see that the mean square error reduced from 0.0278 to 0.020 with the best values of max_depth, min_samples_split and min_samples_leaf respectively.

The plot that we see above is quite readable and it points to the best features that were important to split the data into categories for better interpretability. The plot that we've seen in b) model takes a lot of parameters and it reaches the full depth of the tree.

d)

```
In [71]: 1 training_data = pd.concat([X_train, y_train], axis = 1, ignore_index = True)
2 testing_data = pd.concat([X_test, y_test], axis = 1, ignore_index = True)
```

```
In [72]: 1 training_data_len = len(training_data)
```

```
In [73]: 1 testing_data_len = len(testing_data)
```

```
In [74]: 1 bagging = BaggingRegressor(base_estimator = DecisionTreeRegressor(max_depth = clf.best_params_['max_depth'],
2 min_samples_leaf = clf.best_params_['min_samples_leaf'],
3 min_samples_split = clf.best_params_['min_samples_split']),
4 n_estimators = 20)
5
```

```
In [75]: 1 bagging.fit(X_train,y_train.values.ravel())
2 y_predictions = bagging.predict(X_test)
```

```
In [76]: 1 print("The Mean Squared Error (MSE) of Bagged Decision Tree is: {}".format(mean_squared_error(y_predictions,
```

The Mean Squared Error (MSE) of Bagged Decision Tree is: 0.01952223508092106

```
In [77]: 1 imp_df = pd.DataFrame(model.feature_importances_)
2 imp_df.columns = ["Feature Importances"]
3 imp_df["Features"] = list(X_train.columns)
4 imp_df = imp_df[["Features", "Feature Importances"]]
5 imp_df = imp_df.sort_values(["Feature Importances"], ascending=False)
6 imp_df
```

Out[77]:

	Features	Feature Importances
1	Top10perc	0.702555
2	Top25perc	0.152808
3	F.Undergrad	0.081307
14	Grad.Rate	0.063330
0	Enroll	0.000000
4	P.Undergrad	0.000000
5	Outstate	0.000000
6	Room.Board	0.000000
7	Books	0.000000
8	Personal	0.000000
9	PhD	0.000000
10	Terminal	0.000000
11	S.F.Ratio	0.000000
12	perc.alumni	0.000000
13	Expend	0.000000
15	Private_Yes	0.000000

Observation:

We see that the features Top25perc, Top10perc and F.Undergrad are important for our decision trees as can be seen above. The mean squared error obtained as a result of bagging is about 0.019 respectively.

e)

```
In [78]: 1 metrics_df = {}
2 metrics_df["Estimators"] = []
3 metrics_df["Max Features"] = []
4 metrics_df["Error"] = []
5 for estimators in [1,20,50,100,150,200,300]:
6     for max_feat in range(1,X_train.shape[1] + 1):
7         rfr = RandomForestRegressor(n_estimators = estimators, max_features = max_feat)
8         rfr.fit(X_train,y_train.values.ravel())
9         y_prediction = rfr.predict(X_test)
10        error = mean_squared_error(y_test, y_prediction, squared = True)
11        metrics_df["Estimators"].append(estimators)
12        metrics_df["Max Features"].append(max_feat)
13        metrics_df["Error"].append(error)
14 metrics_df = pd.DataFrame(metrics_df)
15 metrics_df
```

Out[78]:

	Estimators	Max Features	Error
0	1	1	0.031652
1	1	2	0.032917
2	1	3	0.030255
3	1	4	0.028350
4	1	5	0.030153
...
107	300	12	0.016030
108	300	13	0.016236
109	300	14	0.016052
110	300	15	0.015960
111	300	16	0.016168

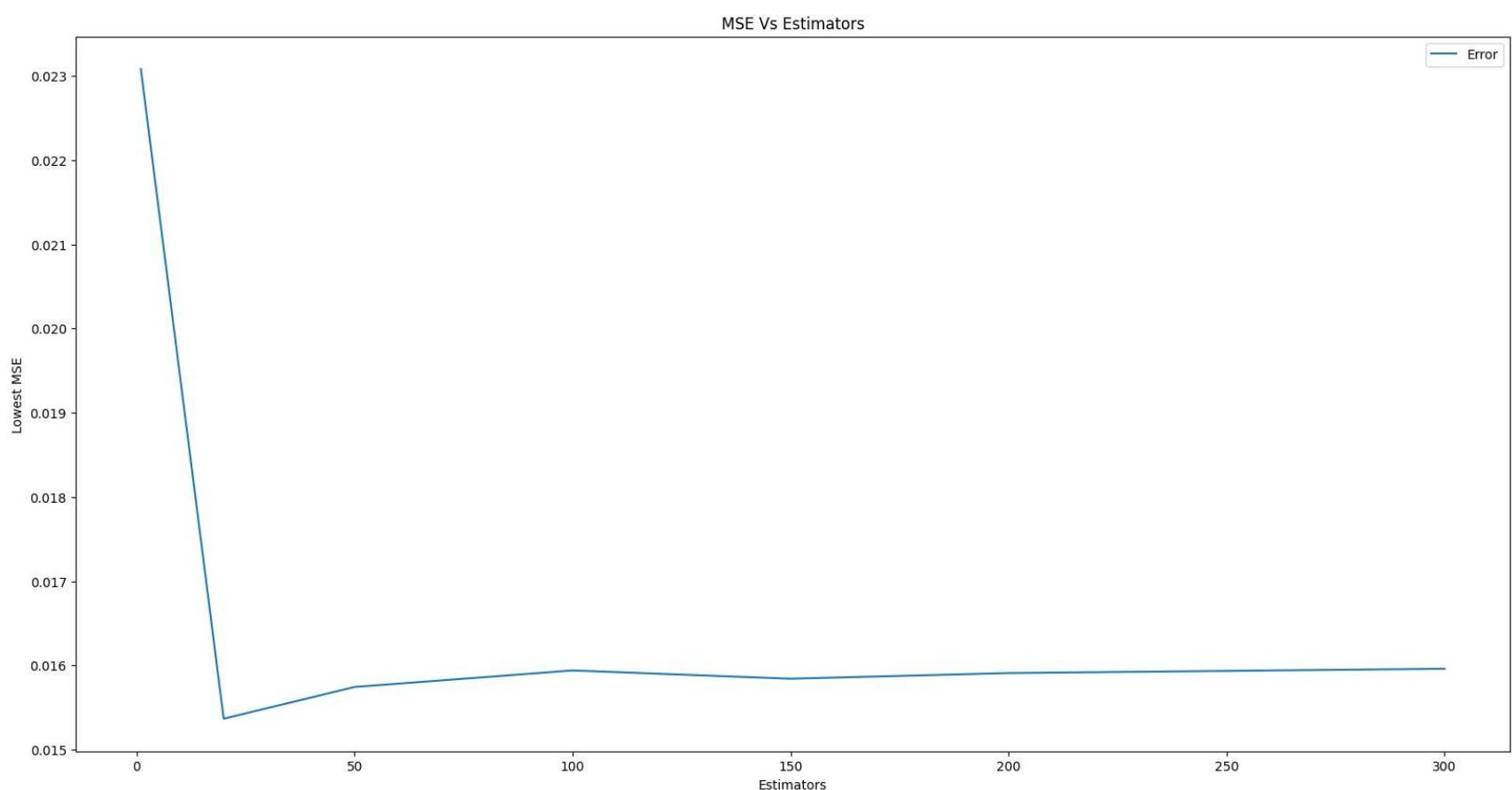
112 rows × 3 columns

```
In [79]: 1 best_params_df = metrics_df[metrics_df["Error"] == metrics_df["Error"].min()]
2 best_params_df
```

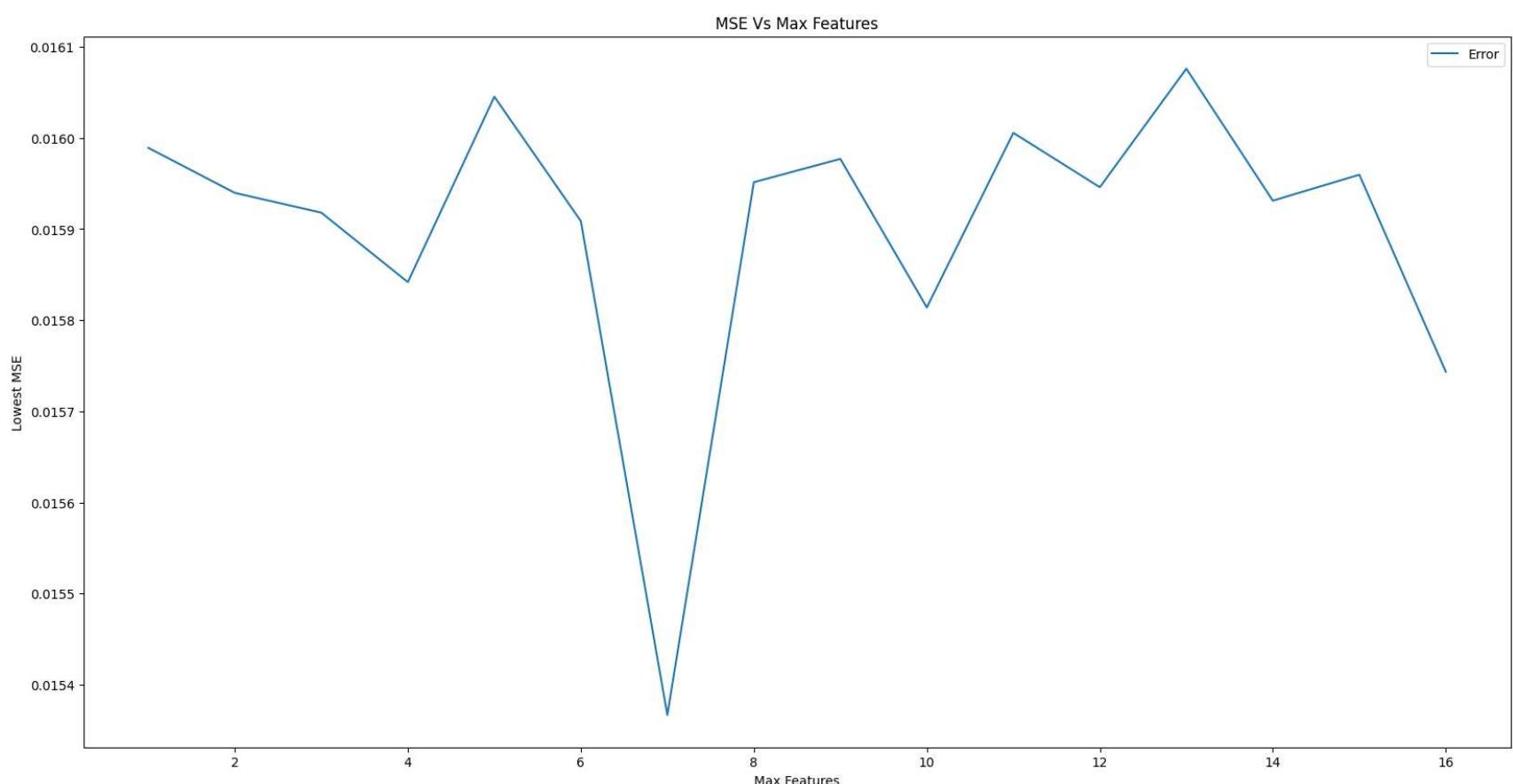
Out[79]:

Estimators	Max Features	Error
22	20	7 0.015367

```
In [80]: 1 estimators_df = metrics_df.groupby(["Estimators"],as_index=False)[["Error"]].min()
2 estimators_df.plot(x="Estimators",y="Error",kind="line",figsize=(20,10))
3 plt.xlabel("Estimators")
4 plt.ylabel("Lowest MSE")
5 plt.title("MSE Vs Estimators")
6 plt.show()
```



```
In [81]: 1 max_features_df = metrics_df.groupby(["Max Features"],as_index=False)[["Error"]].min()
2 max_features_df.plot(x="Max Features",y="Error",kind="line",figsize=(20,10))
3 plt.xlabel("Max Features")
4 plt.ylabel("Lowest MSE")
5 plt.title("MSE Vs Max Features")
6 plt.show()
7
```



```
In [82]: 1 rfr = RandomForestRegressor(n_estimators=best_params_df["Estimators"].values[0],
2                               max_features = best_params_df["Max Features"].values[0])
3 rfr.fit(X_train, y_train.values.ravel())
```

Out[82]: RandomForestRegressor(max_features=7, n_estimators=20)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [83]: 1 imp_df = pd.DataFrame(rfr.feature_importances_)
2 imp_df.columns = ["Feature Importances"]
3 imp_df["Features"] = list(X_train.columns)
4 imp_df = imp_df[["Features", "Feature Importances"]]
5 imp_df = imp_df.sort_values(["Feature Importances"], ascending=False)
6 imp_df
```

Out[83]:

	Features	Feature Importances
2	Top25perc	0.159545
14	Grad.Rate	0.125136
1	Top10perc	0.124239
5	Outstate	0.081853
13	Expend	0.077530
3	F.Undergrad	0.065681
11	S.F.Ratio	0.054187
6	Room.Board	0.050510
7	Books	0.048101
4	P.Undergrad	0.041976
8	Personal	0.035931
12	perc.alumni	0.034739
0	Enroll	0.034697
10	Terminal	0.030795
9	PhD	0.029050
15	Private_Yes	0.006028

In []:

1

```
In [84]: 1 print("The Mean Squared Error of the best Random Forest Regressor is: {}".format(mean_squared_error(y_predi
```

◀

▶

The Mean Squared Error of the best Random Forest Regressor is: 0.01952223508092106

Observation:

The most important features according to our random forest regressor are Top25perc, Top10perc and Grad.Rate respectively.

Problem 4

```
In [85]: 1 penguin_size = pd.read_csv('D:\SML\penguins_size.csv')
2 penguin_lter = pd.read_csv('D:\SML\penguins_lter.csv')
```

```
In [86]: 1 df = pd.concat([penguin_lter, penguin_size], axis=1)
2 df.head()
```

Out[86]:

	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)	...	Delta 15 N (o/oo)	Delta 13 C (o/oo)	Comment
0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1	...	NaN	NaN	Not enough blood for isotope
1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5	...	8.94956	-24.69454	Na
2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3	...	8.36821	-25.33302	Na
3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN	...	NaN	NaN	Adult not sampled
4	PAL0708	5	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7	...	8.76651	-25.32426	Na

5 rows × 24 columns

In [87]: 1 df.columns

Out[87]: Index(['studyName', 'Sample Number', 'Species', 'Region', 'Island', 'Stage', 'Individual ID', 'Clutch Completion', 'Date Egg', 'Culmen Length (mm)', 'Culmen Depth (mm)', 'Flipper Length (mm)', 'Body Mass (g)', 'Sex', 'Delta 15 N (o/oo)', 'Delta 13 C (o/oo)', 'Comments', 'species', 'island', 'culmen_length_mm', 'culmen_depth_mm', 'flipper_length_mm', 'body_mass_g', 'sex'], dtype='object')

In [88]: 1 # Removing unnecessary columns
2 df = df.drop(columns=["Sample Number", "Individual ID", "Date Egg", "Stage", "Region", "Comments"])
3
4 # Removing duplicate columns
5 df = df.drop(columns=["Island", "Culmen Length (mm)", "Culmen Depth (mm)", "Flipper Length (mm)", "Body Mass (g)"])
6
7 df

Out[88]:

	studyName	Clutch Completion	Delta 15 N (o/oo)	Delta 13 C (o/oo)	species	island	culmen_length_mm	culmen_depth_mm	flipper_length_mm	body_ma
0	PAL0708	Yes	NaN	NaN	Adelie	Torgersen	39.1	18.7	181.0	37
1	PAL0708	Yes	8.94956	-24.69454	Adelie	Torgersen	39.5	17.4	186.0	38
2	PAL0708	Yes	8.36821	-25.33302	Adelie	Torgersen	40.3	18.0	195.0	39
3	PAL0708	Yes	NaN	NaN	Adelie	Torgersen	NaN	NaN	NaN	NaN
4	PAL0708	Yes	8.76651	-25.32426	Adelie	Torgersen	36.7	19.3	193.0	34
...
339	PAL0910	No	NaN	NaN	Gentoo	Biscoe	NaN	NaN	NaN	NaN
340	PAL0910	Yes	8.41151	-26.13832	Gentoo	Biscoe	46.8	14.3	215.0	48
341	PAL0910	Yes	8.30166	-26.04117	Gentoo	Biscoe	50.4	15.7	222.0	51
342	PAL0910	Yes	8.24246	-26.11969	Gentoo	Biscoe	45.2	14.8	212.0	52
343	PAL0910	Yes	8.36390	-26.15531	Gentoo	Biscoe	49.9	16.1	213.0	54

344 rows × 11 columns

In [89]: 1 df = df[~df["Delta 15 N (o/oo)"].isnull()].reset_index(drop = True)
2 df = df[~df["Delta 13 C (o/oo)"].isnull()].reset_index(drop = True)
3 df = df[~df["sex"].isnull()].reset_index(drop = True)

In [90]: 1 df["sex"].unique()

Out[90]: array(['FEMALE', 'MALE', '.'], dtype=object)

In [91]: 1 df = df[~(df["sex"] == ".")].reset_index(drop = True)

In [92]: 1 df = encoding(df, "island")
2 df = encoding(df, "sex")
3 df = encoding(df, "Clutch Completion")
4 df = encoding(df, "studyName")

In [93]: 1 df["species"].unique()

Out[93]: array(['Adelie', 'Chinstrap', 'Gentoo'], dtype=object)

In [94]: 1 df["species"] = df["species"].replace("Adelie", 1)
2 df["species"] = df["species"].replace("Chinstrap", 2)
3 df["species"] = df["species"].replace("Gentoo", 3)

In [95]: 1 X = df.drop(columns=["species"])
2 Y = df[["species"]]

In [96]: 1 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 0.2, random_state = 42)

a)

In [97]: 1 X_train_bootstrap, y_train_bootstrap = resample(X_train, y_train, n_samples = 100)

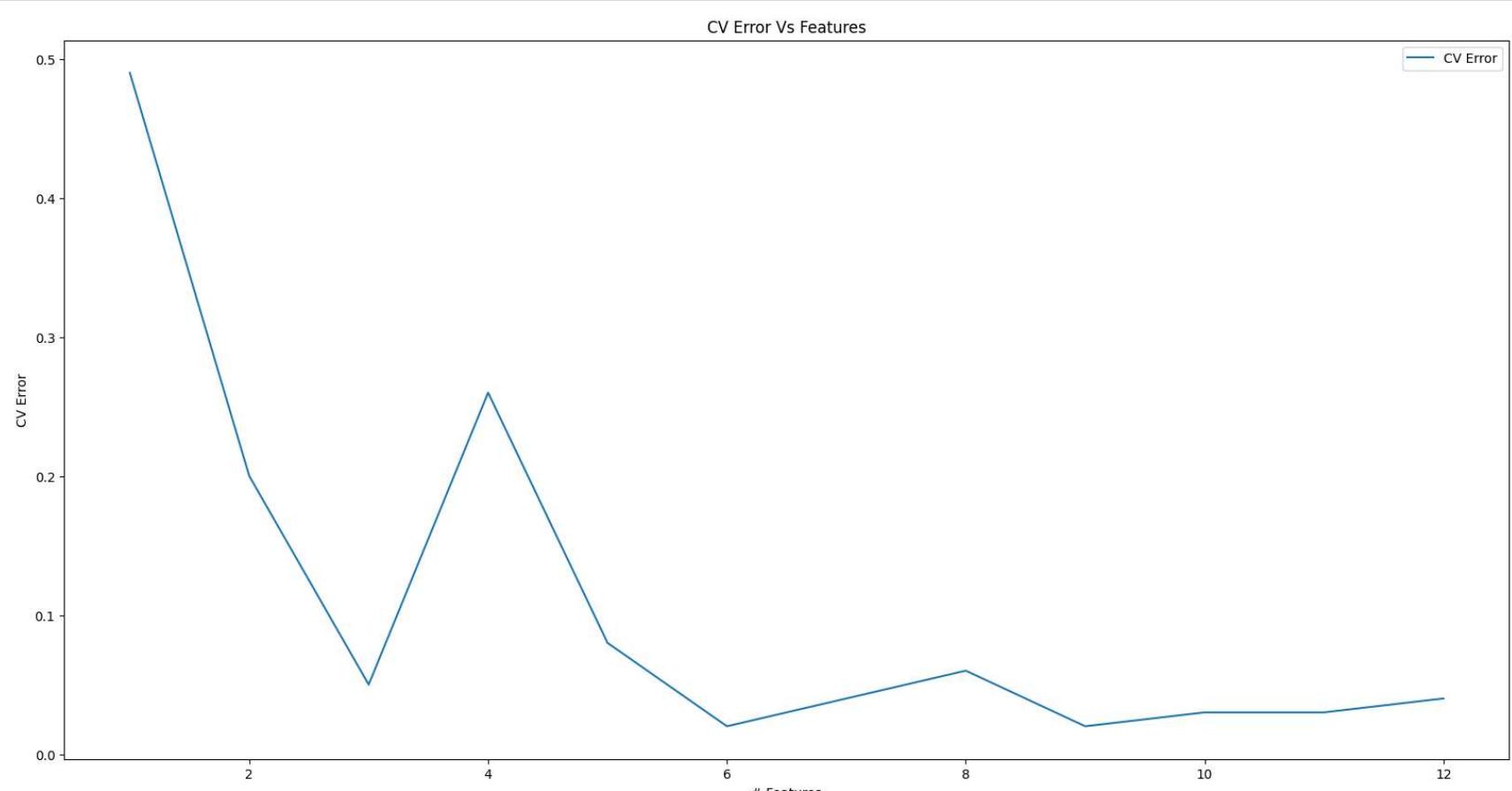
```
In [98]: 1 column_list = list(X_train_bootstrap.columns)
2
3 cv_df = {}
4 cv_df["Features"] = []
5 cv_df["# Features"] = []
6 cv_df["CV Error"] = []
7
8 # Iterating through all the combinations
9 for p in range(len(column_list) + 1):
10     for subset in random.sample(list(itertools.combinations(column_list, p)), 1):
11         # Removing the empty subset
12         if(list(subset)!=[]):
13
14             # Filtering the subsetted dataframe
15             X_subset = X_train_bootstrap[list(subset)]
16
17             # Training a Decision Tree Classifier with max depth 6
18             clf = DecisionTreeClassifier(max_depth = 6).fit(X_subset,y_train_bootstrap)
19
20             # Finding the CV Error
21             cv_error = 1 - cross_val_score(clf, X_subset, y_train_bootstrap, scoring = "accuracy").mean()
22
23             cv_df["Features"].append(list(subset))
24             cv_df["# Features"].append(len(list(subset)))
25             cv_df["CV Error"].append(cv_error)
26
27 cv_df = pd.DataFrame(cv_df)
28 cv_df
```

Out[98]:

	Features	# Features	CV Error
0	[studyName_PAL0910]	1	0.49
1	[culmen_length_mm, island_Torgersen]	2	0.20
2	[Delta 15 N (o/oo), Delta 13 C (o/oo), body_ma...	3	0.05
3	[culmen_depth_mm, body_mass_g, sex_MALE, study...	4	0.26
4	[Delta 13 C (o/oo), culmen_depth_mm, island_To...	5	0.08
5	[Delta 15 N (o/oo), Delta 13 C (o/oo), flipper...	6	0.02
6	[culmen_length_mm, flipper_length_mm, body_mas...	7	0.04
7	[Delta 15 N (o/oo), Delta 13 C (o/oo), culmen_...	8	0.06
8	[Delta 15 N (o/oo), Delta 13 C (o/oo), culmen_...	9	0.02
9	[culmen_length_mm, culmen_depth_mm, flipper_le...	10	0.03
10	[Delta 15 N (o/oo), Delta 13 C (o/oo), culmen_...	11	0.03
11	[Delta 15 N (o/oo), Delta 13 C (o/oo), culmen_...	12	0.04

In [99]:

```
1 cv_df.plot(x="# Features",y="CV Error",figsize=(20,10))
2 plt.xlabel("# Features")
3 plt.ylabel("CV Error")
4 plt.title("CV Error Vs Features")
5 plt.show()
```



```
In [100]: 1 min_cv_df = cv_df[cv_df["CV_Error"] == cv_df["CV_Error"].min()].reset_index(drop=True)
2 display(min_cv_df)
3 print("\nTherefore, the most optimal p obtained is:",min_cv_df["# Features"].values[0])
4
```

	Features	# Features	CV Error
0	[Delta 15 N (o/oo), Delta 13 C (o/oo), flipper...	6	0.02
1	[Delta 15 N (o/oo), Delta 13 C (o/oo), culmen...	9	0.02

Therefore, the most optimal p obtained is: 6

b)

```
In [101]: 1 y_pred_final_train = {}
2 y_pred_final_test = {}
3
4 # Calculating predictions for each tree
5 for n_trees in [1,50,100,150,200,300,400]:
6     final_y_pred_train = []
7     final_y_pred_test = []
8
9     for i in range(n_trees):
10
11         # Iterating through all the combinations of the best feature subset (p) obtained above
12         for subset in random.sample(list(itertools.combinations(column_list, min_cv_df["# Features"].values
13
14             # Filtering the subsetted dataframes
15             X_train_bootstrap_subset = X_train_bootstrap[list(subset)]
16             X_train_subset = X_train[list(subset)]
17             X_test_subset = X_test[list(subset)]
18
19             clf = DecisionTreeClassifier(max_depth = 6)
20             clf.fit(X_train_bootstrap_subset, y_train_bootstrap)
21
22             y_pred_train = clf.predict(X_train_subset)
23             y_pred_test = clf.predict(X_test_subset)
24
25             final_y_pred_train.append(y_pred_train)
26             final_y_pred_test.append(y_pred_test)
27
28         # Finding the mode of the predictions obtained
29         y_pred_final_train[n_trees] = stats.mode(np.array(final_y_pred_train), keepdims = True)[0][0]
30         y_pred_final_test[n_trees] = stats.mode(np.array(final_y_pred_test), keepdims = True)[0][0]
31
32 y_pred_final_train = pd.DataFrame(y_pred_final_train)
33 y_pred_final_test = pd.DataFrame(y_pred_final_test)
34
35 print("Reporting the final predictions of the Training Set for each tree in [1,50,100,150,200,300,400]:")
36 y_pred_final_train
```

Reporting the final predictions of the Training Set for each tree in [1,50,100,150,200,300,400]:

Out[101]:

	1	50	100	150	200	300	400
0	1	1	1	1	1	1	1
1	2	2	2	2	2	2	2
2	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1
4	2	2	2	2	2	2	2
...
254	2	2	2	2	2	2	2
255	1	1	1	1	1	1	1
256	1	1	1	1	1	1	1
257	3	3	3	3	3	3	3
258	1	1	1	1	1	1	1

259 rows × 7 columns

c)

```
In [102]: 1 metrics = {}
2 metrics["No of Trees"] = []
3 metrics["Train Zero One Test Error"] = []
4 metrics["Test Zero One Test Error"] = []
5 metrics["F1 Score"] = []
6 metrics["AUC"] = []
7
8 # For every T
9 for col in y_pred_final_train:
10     metrics["No of Trees"].append(col)
11
12 # Calculating the Train Error
13 zero_one_error_train = 1 - (np.array(y_train["species"]) == np.array(y_pred_final_train[col])).sum()/len(y_train)
14 metrics["Train Zero One Test Error"].append(zero_one_error_train)
15
16 # Calculating the Test Error
17 zero_one_error_test = 1 - (np.array(y_test["species"]) == np.array(y_pred_final_test[col])).sum()/len(y_test)
18 metrics["Test Zero One Test Error"].append(zero_one_error_test)
19
20 # Calculating the F1 Score (One Vs ALL - 3 values for each class)
21 f1score = f1_score(y_test["species"],y_pred_final_test[col],average = None)
22 f1score = [round(x,3) for x in f1score]
23 metrics["F1 Score"].append(f1score)
24
25 # Calculating the AUC (One Vs ALL - 3 values for each class)
26 y_pred_class = pd.DataFrame(OneHotEncoder().fit_transform(y_pred_final_test[[col]]).toarray()).astype(int)
27 y_test_class = pd.DataFrame(OneHotEncoder().fit_transform(y_test[["species"]]).toarray())
28
29 auc_list = []
30 for col1 in y_test_class.columns:
31     auc = round(roc_auc_score(list(y_test_class[col1]),list(y_pred_class[col1])),3)
32     auc_list.append(auc)
33
34 metrics["AUC"].append(auc_list)
35
36 print("Reporting the Final Metrics for every tree(T): ")
37 pd.DataFrame(metrics)
```

Reporting the Final Metrics for every tree(T):

	No of Trees	Train Zero One Test Error	Test Zero One Test Error	F1 Score	AUC
0	1	0.046332	0.046154	[0.968, 0.923, 0.952]	[0.969, 0.981, 0.965]
1	50	0.038610	0.061538	[0.933, 0.923, 0.955]	[0.938, 0.981, 0.977]
2	100	0.038610	0.046154	[0.951, 0.96, 0.955]	[0.953, 0.991, 0.977]
3	150	0.023166	0.015385	[0.984, 1.0, 0.977]	[0.984, 1.0, 0.989]
4	200	0.027027	0.015385	[0.984, 1.0, 0.977]	[0.984, 1.0, 0.989]
5	300	0.027027	0.015385	[0.984, 1.0, 0.977]	[0.984, 1.0, 0.989]
6	400	0.027027	0.015385	[0.984, 1.0, 0.977]	[0.984, 1.0, 0.989]

d)

```
In [103]: 1 rf_metrics = {}
2 rf_metrics["No of Trees"] = []
3 rf_metrics["Train Zero One Error"] = []
4 rf_metrics["Test Zero One Error"] = []
5 rf_metrics["F1 Score Train"] = []
6 rf_metrics["F1 Score Test"] = []
7 rf_metrics["AUC Train"] = []
8 rf_metrics["AUC Test"] = []
9
10 feature_imp = {}
11
12 for n_trees in [10,50,100]:
13
14     rf_metrics["No of Trees"].append(n_trees)
15
16     # Training a Random Forest Classifier
17     clf = RandomForestClassifier(n_estimators=n_trees).fit(X_train_bootstrap, y_train_bootstrap.values.ravel())
18
19     # Predicting the Classes
20     y_pred_train = clf.predict(X_train)
21     y_pred_test = clf.predict(X_test)
22
23     # Computing the Zero One Test Error
24     rf_metrics["Train Zero One Error"].append(1 - (np.array(y_train["species"]) == y_pred_train).sum()/len(y_train))
25     rf_metrics["Test Zero One Error"].append(1 - (np.array(y_test["species"]) == y_pred_test).sum()/len(y_test))
26
27     # Computing the F1 Score
28     f1score = f1_score(y_train["species"],y_pred_train,average=None)
29     f1score = [round(x,3) for x in f1score]
30     rf_metrics["F1 Score Train"].append(f1score)
31
32     f1score = f1_score(y_test["species"],y_pred_test,average=None)
33     f1score = [round(x,3) for x in f1score]
34     rf_metrics["F1 Score Test"].append(f1score)
35
36     # Computing the AUC - Train
37     y_pred_class = pd.DataFrame(OneHotEncoder().fit_transform(y_pred_train.reshape(-1,1)).toarray()).astype(bool)
38     y_train_class = pd.DataFrame(OneHotEncoder().fit_transform(y_train[["species"]]).toarray())
39
40     auc_list = []
41     for col1 in y_train_class.columns:
42         auc = round(roc_auc_score(list(y_train_class[col1]),list(y_pred_class[col1])),3)
43         auc_list.append(auc)
44
45     rf_metrics["AUC Train"].append(auc_list)
46
47     # Computing the AUC - Test
48     y_pred_class = pd.DataFrame(OneHotEncoder().fit_transform(y_pred_test.reshape(-1,1)).toarray()).astype(bool)
49     y_test_class = pd.DataFrame(OneHotEncoder().fit_transform(y_test[["species"]]).toarray())
50
51     auc_list = []
52     for col1 in y_test_class.columns:
53         auc = round(roc_auc_score(list(y_test_class[col1]),list(y_pred_class[col1])),3)
54         auc_list.append(auc)
55
56     rf_metrics["AUC Test"].append(auc_list)
57
58     # Calculating the Top Ten Features wrt their importances
59     imp_df = pd.DataFrame(clf.feature_importances_)
60     imp_df.columns = ["Feature Importances"]
61     imp_df["Features"] = list(X_train.columns)
62     imp_df = imp_df[["Features","Feature Importances"]]
63
64     top10_imp_df = imp_df.sort_values(["Feature Importances"],ascending=False).reset_index(drop=True)[:10]
65
66     feature_imp[n_trees] = top10_imp_df.copy()
67
68 print("Reporting the Final Metrics for every tree(T): ")
69 pd.DataFrame(rf_metrics)
```

Reporting the Final Metrics for every tree(T):

	No of Trees	Train Zero One Error	Test Zero One Error	F1 Score Train	F1 Score Test	AUC Train	AUC Test
0	10	0.023166	0.000000	[0.972, 0.945, 1.0]	[1.0, 1.0, 1.0]	[0.976, 0.965, 1.0]	[1.0, 1.0, 1.0]
1	50	0.015444	0.000000	[0.981, 0.964, 1.0]	[1.0, 1.0, 1.0]	[0.984, 0.977, 1.0]	[1.0, 1.0, 1.0]
2	100	0.019305	0.015385	[0.977, 0.955, 1.0]	[0.984, 1.0, 0.977]	[0.979, 0.974, 1.0]	[0.984, 1.0, 0.989]

```
In [104]: 1 for n_trees in feature_imp.keys():
2
3     print(f"\nTop 10 Features wrt Feature Importances of a Random Forest Classifier of {n_trees} trees:")
4     display(feature_imp[n_trees])
```

Top 10 Features wrt Feature Importances of a Random Forest Classifier of 10 trees:

	Features	Feature Importances
0	flipper_length_mm	0.218543
1	Delta 13 C (o/oo)	0.209097
2	culmen_depth_mm	0.188981
3	culmen_length_mm	0.144862
4	body_mass_g	0.097088
5	Delta 15 N (o/oo)	0.078094
6	island_Dream	0.055412
7	island_Torgersen	0.005335
8	sex_MALE	0.001599
9	studyName_PAL0910	0.000988

Top 10 Features wrt Feature Importances of a Random Forest Classifier of 50 trees:

	Features	Feature Importances
0	culmen_length_mm	0.306287
1	flipper_length_mm	0.150637
2	Delta 13 C (o/oo)	0.146349
3	body_mass_g	0.112374
4	Delta 15 N (o/oo)	0.096862
5	culmen_depth_mm	0.096744
6	island_Dream	0.071028
7	studyName_PAL0809	0.006501
8	studyName_PAL0910	0.004701
9	sex_MALE	0.003870

Top 10 Features wrt Feature Importances of a Random Forest Classifier of 100 trees:

	Features	Feature Importances
0	culmen_length_mm	0.239834
1	flipper_length_mm	0.232335
2	Delta 13 C (o/oo)	0.175041
3	body_mass_g	0.109158
4	culmen_depth_mm	0.106433
5	Delta 15 N (o/oo)	0.067624
6	island_Dream	0.048551
7	studyName_PAL0809	0.006786
8	studyName_PAL0910	0.006202
9	sex_MALE	0.004824

Problem 5

a)

Solution:

Let h_1 and h_2 represent the nodes in the hidden layer,

From the given neural network architecture:

$$h_1 = c(w_1 + w_3 \cdot x_1 + w_5 \cdot x_2)$$

$$h_2 = c(w_2 + w_4 \cdot x_1 + w_6 \cdot x_2)$$

Let z represent the value of the output node before applying the activation function,

$$z = w_7 + h_1 \cdot w_8 + h_2 \cdot w_9$$

$$z = w_7 + c \cdot w_8 \cdot (w_1 + w_3 \cdot x_1 + w_5 \cdot x_2) + c \cdot w_9 \cdot (w_2 + w_4 \cdot x_1 + w_6 \cdot x_2)$$

$$z = (w_7 + c \cdot w_8 \cdot w_1 + c \cdot w_9 \cdot w_2) + (c \cdot w_8 \cdot w_3 + c \cdot w_9 \cdot w_4) \cdot x_1 + (c \cdot w_8 \cdot w_5 + c \cdot w_9 \cdot w_6) \cdot x_2$$

Now applying the activation function (sigmoid), we get:

$$P(y = 1) = \frac{1}{(1+e^{-z})}$$

$$\Rightarrow P(y = 1) = \frac{1}{(1+e^{-(w_7+c.w_8.w_1+c.w_9.w_2)+(c.w_8.w_3+c.w_9.w_4).x_1+(c.w_8.w_5+c.w_9.w_6).x_2)})}$$

The final classification boundary is when:

$$P(y=1) = 0.5$$

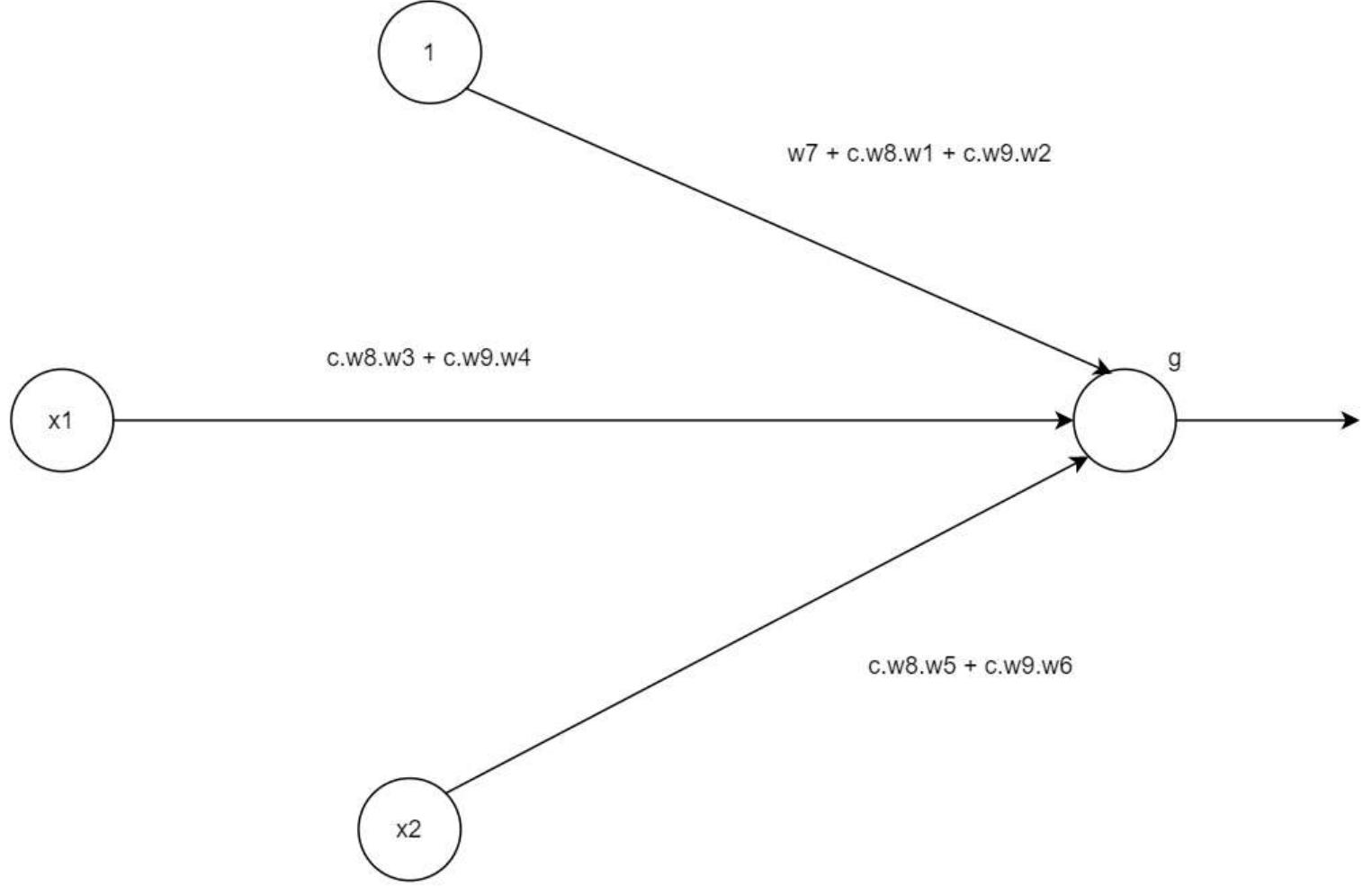
$$\Rightarrow \frac{1}{(1+e^{-z})} = 0.5 \Rightarrow z = 0$$

Therefore, the classification boundary is:

$$(w_7 + c \cdot w_8 \cdot w_1 + c \cdot w_9 \cdot w_2) + (c \cdot w_8 \cdot w_3 + c \cdot w_9 \cdot w_4) \cdot x_1 + (c \cdot w_8 \cdot w_5 + c \cdot w_9 \cdot w_6) \cdot x_2 = 0$$

We obtain a linear classification boundary in terms of x

b)



Let the new weights be w_1, w_2, w_3

$$w_1 = w_7 + c \cdot w_8 \cdot w_1 + c \cdot w_9 \cdot w_2$$

$$w_2 = c \cdot w_8 \cdot w_3 + c \cdot w_9 \cdot w_4$$

$$w_3 = c \cdot w_8 \cdot w_5 + c \cdot w_9 \cdot w_6$$

$$P(y = 1) = g(w_1 + w_2 \cdot x_1 + w_3 \cdot x_2)$$

Yes. If linear activation functions are used for all the hidden units, output from hidden units will be written as linear combination of input features. Since these intermediate output serves as input for the final output layer, we can always find an equivalent neural net which does not have any hidden layer as seen in the example above.