

1. Clustering

Sex

Values 1: Male 2: Female

Age

How have you been feeling over the last two weeks

q1 - I have felt cheerful and in good spirits

q2 - I have felt calm and relaxed

q3 - I have felt active and vigorous

q4 - I woke up feeling fresh and rested

q5 - My daily life has been filled with things that interest me

Values:

1. All of the time.
2. Most of the time
3. More than half of the time
4. Less than half of the time
5. Some of the time
6. At no time

Please tell me how often you feel this way

w1 - At my work I feel full of energy

w2 - I am enthusiastic about my job

w3 - Time flies when I am working

w4 - In my opinion, I am good at my job

Values for variables Q00a to Q00f

1. Always
2. Most of the time
3. Sometimes
4. Rarely
5. Never

PCA, Elbow method Clustering, K-means clustering

```
In [204]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

plt.style.use('seaborn')
plt.rc('font', size=14)
plt.rc('figure', titlesize=18)
plt.rc('axes', labelsiz=15)
plt.rc('axes', titlesize=18)

init_df = pd.read_csv('EWCS_2016.csv')
init_df = init_df.replace(-999, np.nan)
init_df = init_df.dropna()
init_df.columns = ["sex", "age", "q1", "q2", "q3", "q4", "q5", "w1", "w2", "w3", "w4"]
init_df.head()
```

```
Out[204]:
```

	sex	age	q1	q2	q3	q4	q5	w1	w2	w3	w4
0	1.0	63.0	3.0	3.0	3.0	3.0	3.0	2.0	2.0	2.0	2.0
1	2.0	58.0	2.0	3.0	2.0	3.0	2.0	2.0	3.0	2.0	2.0
2	2.0	32.0	2.0	3.0	2.0	3.0	2.0	2.0	2.0	2.0	2.0
3	1.0	35.0	3.0	2.0	2.0	3.0	2.0	2.0	2.0	2.0	2.0
4	2.0	27.0	2.0	2.0	3.0	3.0	2.0	2.0	4.0	2.0	2.0

PCA

```
In [554]: from sklearn.decomposition import PCA
from sklearn import preprocessing

df = init_df.copy(deep=True)

def do_pca(init_df, delta=None, devide=True):
    return_pca = PCA()
    pca_data = pca.transform(df)

    per_var = np.round(pca.explained_variance_ratio_ * 100, 1)

    labels = ['PC' + str(i) for i in range(1, len(per_var) + 1)]

    # Plot
    # Scree plot
    if return_scree:
        plt.bar(range(len(per_var)+1), height=per_var, tick_label=labels)
        plt.ylabel('Percentage of Explained Variance')
        plt.xlabel('Principal Component')
        plt.title('Scree Plot')
        plt.show()

    # PC1 vs PC2
    pca_df = pd.DataFrame(pca_data, index=init_df.index, columns=labels)

    plt.scatter(pca_df.PC1, pca_df.PC2)
    plt.title('PCA Graph')
    plt.xlabel('PC1 - %s' % per_var[0])
    plt.ylabel('PC2 - %s' % per_var[1])

    if devide:
        seen = set()
        clusters = []
        for c in range(1, len(per_var)+1):
            d1 = (pca_df.PC1.max() - pca_df.PC1.min()) / 10
            d2 = (pca_df.PC2.max() - pca_df.PC2.min()) / 10
            delta = (d1 + d2) / 2

            for sample in pca_df.index:
                x, y = pca_df.PC1.loc[sample], pca_df.PC2.loc[sample]

                found = False
                for s in seen:
                    x1, y1 = pca_df.PC1.loc[s], pca_df.PC2.loc[s]
                    if sum(np.power(np.array(x1, y1) - np.array(x, y), 2)) < delta**2:
                        found = True
                        break

                if found:
                    clusters[s].append(sample)
                else:
                    clusters[sample] = [sample]
                    seen.add(sample)

            clusts = []
            for c in range(1, len(per_var)+1):
                i = len(clusters[c])
                if i > len(df)/20:
                    clusts.append(c)
                    x, y = pca_df.PC1.loc[c], pca_df.PC2.loc[c]
                    plt.annotate('C' + str(i), [x, y])

            if return_info:
                return('Centers of clusters are represented by the following samples: %s' % list(clusts.keys()))

            return clusts

    return None
```

Raw Data

```
In [555]: df = init_df.copy(deep=True)
df.head()
```

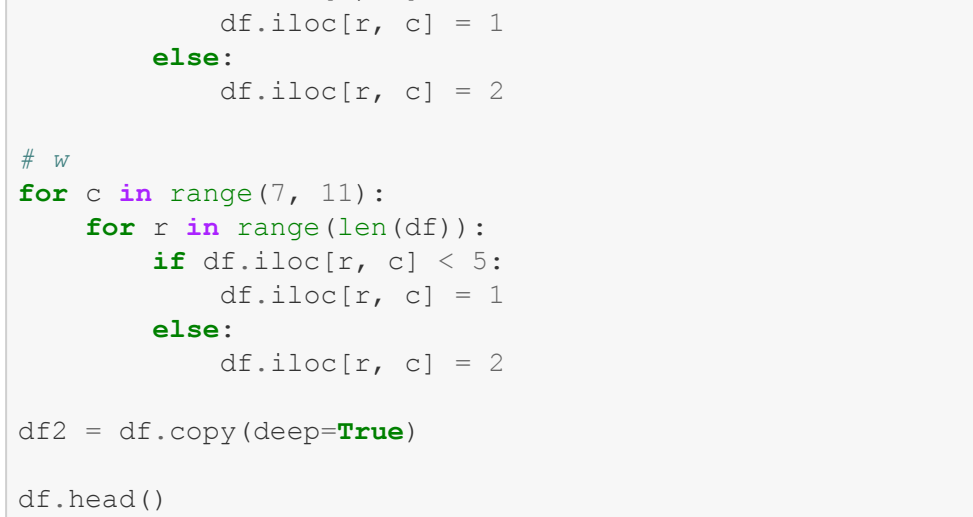
```
Out[555]:
```

	sex	age	q1	q2	q3	q4	q5	w1	w2	w3	w4
0	1.0	63.0	3.0	3.0	3.0	3.0	3.0	2.0	2.0	2.0	2.0
1	2.0	58.0	2.0	3.0	2.0	3.0	2.0	2.0	3.0	2.0	2.0
2	2.0	32.0	2.0	3.0	2.0	3.0	2.0	2.0	2.0	2.0	2.0
3	1.0	35.0	3.0	2.0	2.0	3.0	2.0	2.0	2.0	2.0	2.0
4	2.0	27.0	2.0	2.0	3.0	3.0	2.0	2.0	4.0	2.0	2.0

```
In [556]: clusts = do_pca(df, delta=8)
```



Centers of clusters are represented by the following samples:



Data-frame v2

1. For q:

1,2,3 = 1

4,5,6 = 2

2. For w:

1,2,3,4 = 1

5,6 = 2

```
In [210]: df = init_df.copy(deep=True)

# q
for c in range(2, 7):
    for r in range(len(df)):
        if df.iloc[r, c] < 4:
            df.iloc[r, c] = 1
        else:
            df.iloc[r, c] = 2

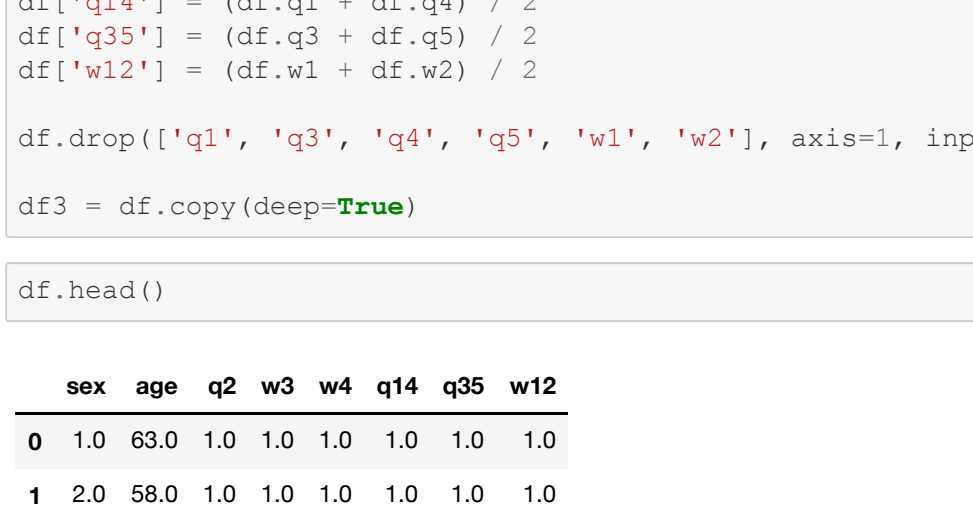
# w
for c in range(7, 11):
    for r in range(len(df)):
        if df.iloc[r, c] < 5:
            df.iloc[r, c] = 1
        else:
            df.iloc[r, c] = 2

df2 = df.copy(deep=True)
df2.head()
```

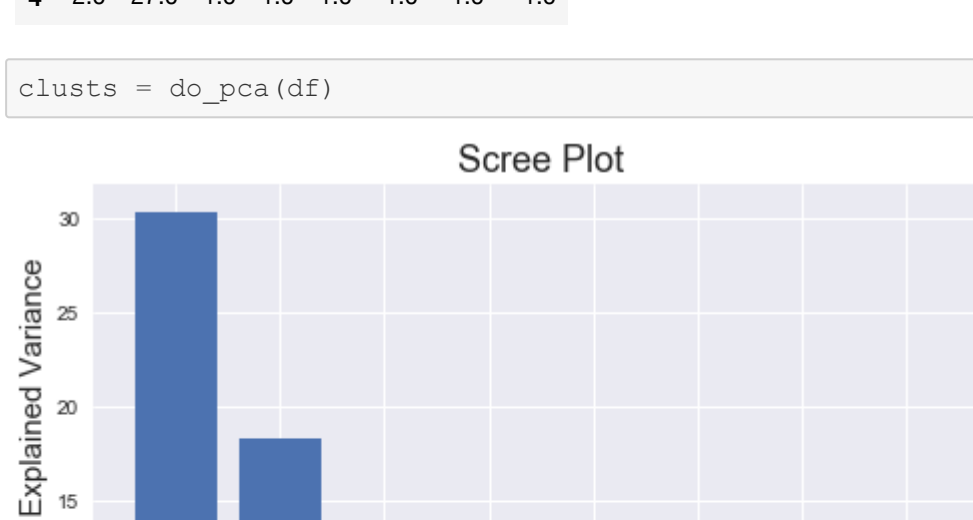
```
Out[210]:
```

	sex	age	q1	q2	q3	q4	q5	w1	w2	w3	w4
0	1.0	63.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
1	2.0	58.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	2.0	32.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
3	1.0	35.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
4	2.0	27.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

```
In [208]: clusts = do_pca(df2)
```



Centers of clusters are represented by the following samples:



Data-frame v3

Modified DF v2

1. Join q1 and q4

2. Join q3 and q5

3. Join w1 and w2

```
In [211]: df['q1q4'] = (df.q1 + df.q4) / 2
df['q3q5'] = (df.q3 + df.q5) / 2
df['w1w2'] = (df.w1 + df.w2) / 2

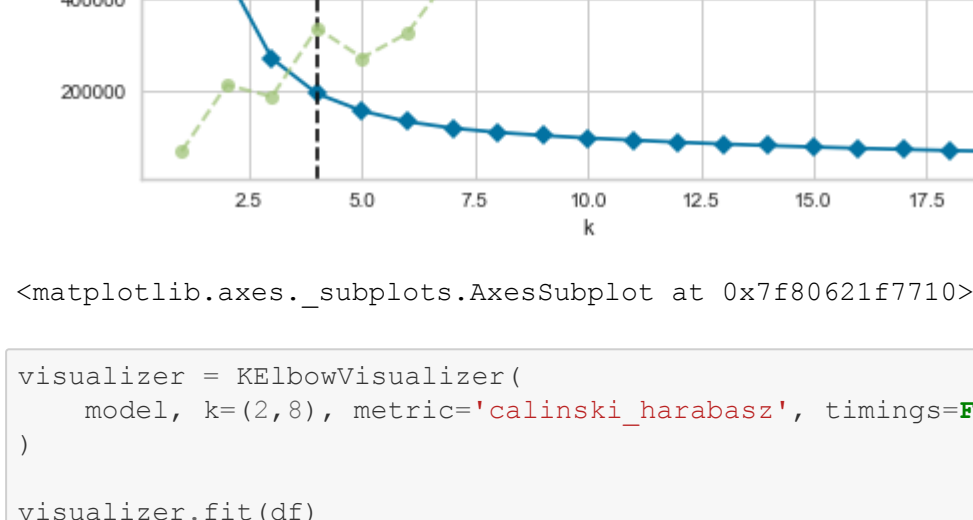
df.drop(['q1', 'q3', 'q4', 'q5', 'w1', 'w2'], axis=1, inplace=True)
df3 = df.copy(deep=True)
df3.head()
```

```
In [212]: df.head()
```

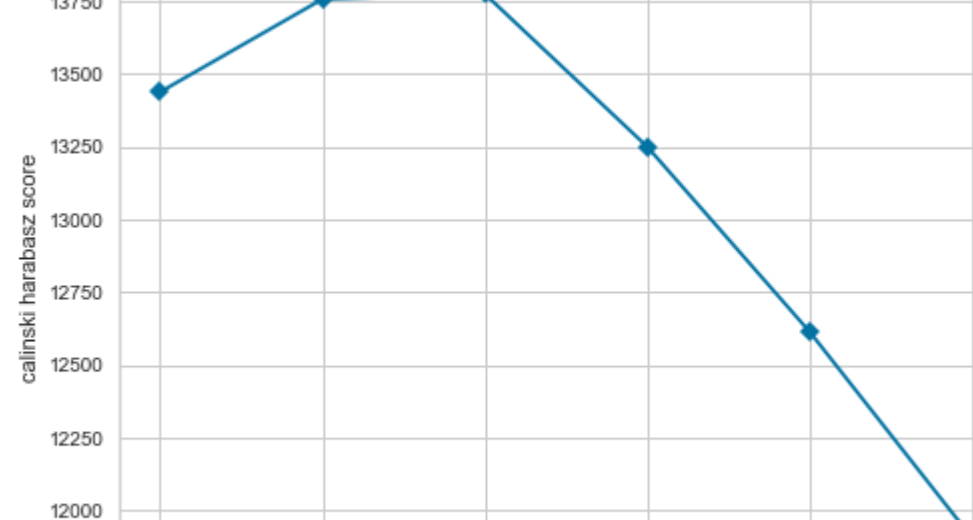
```
Out[212]:
```

	sex	age	q2	w3	w4	q1q4	q3q5	w1w2
0	1.0	63.0	1.0	1.0	1.0	1.0	1.0	1.0
1	2.0	58.0	1.0	1.0	1.0	1.0	1.0	1.0
2	2.0	32.0	1.0	1.0	1.0	1.0	1.0	1.0
3	1.0	35.0	1.0	1.0	1.0	1.0	1.0	1.0
4	2.0	27.0	1.0	1.0	1.0	1.0	1.0	1.0

```
In [213]: clusts = do_pca(df3)
```



Centers of clusters are represented by the following samples:

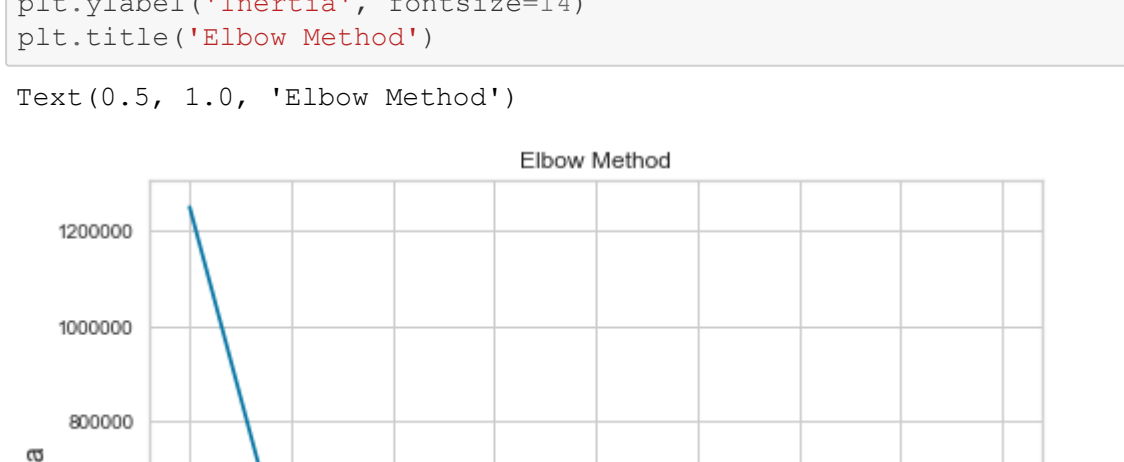


K-means Clustering

```
In [563]: from sklearn.cluster import KMeans
from yellowbrick.cluster import KElbowVisualizer
```

```
model = KMeans()
visualizer = KElbowVisualizer(model, k=(1,20))
```

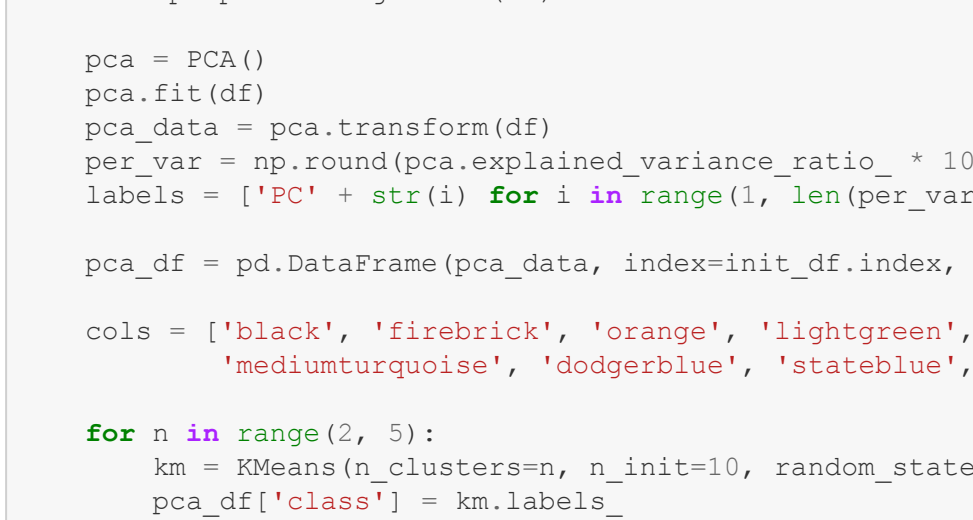
```
visualizer.fit(df)
visualizer.show()
```



```
Out[563]: <matplotlib.axes._subplots.AxesSubplot at 0x7f80621f7710>
```

```
In [574]: visualizer = KElbowVisualizer(
    model, k=(2,8), metric='calinski_harabasz', timings=False, locate_elbow=False)

visualizer.fit(df)
visualizer.show()
```



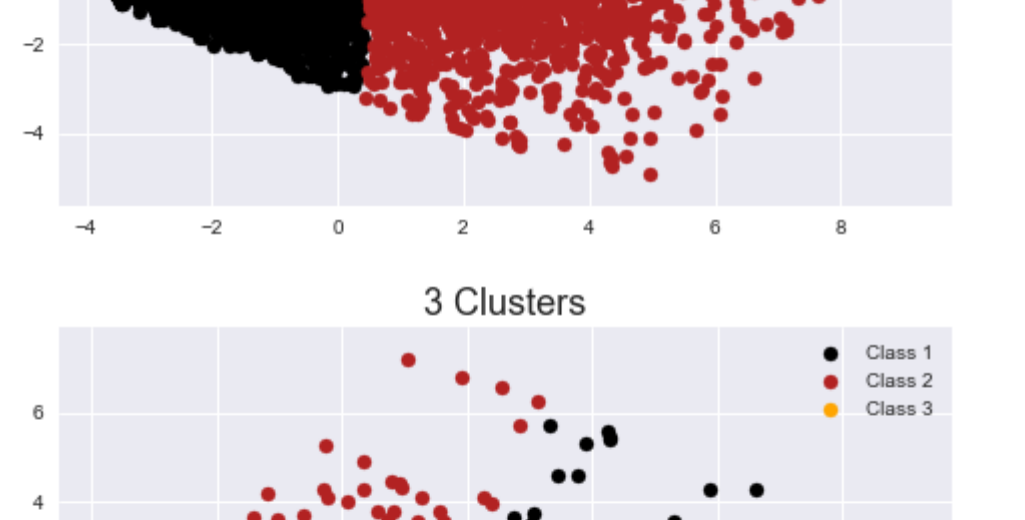
```
Out[574]: <matplotlib.axes._subplots.AxesSubplot at 0x7f80a1a94e50>
```

```
In [570]: iner = []
```

```
for k in range(1, 10, 1):
    km = KMeans(n_clusters=k).fit(df)
    iner.append(km.inertia_)

plt.plot(range(1, 10), iner)
plt.xlabel('k', fontsize=18)
plt.ylabel('Inertia', fontsize=14)
plt.title('Elbow Method')
```

```
Out[570]: Text(0.5, 1.0, 'Elbow Method')
```



```
In [557]: def show_clusts(init_df):
    df = init_df.copy(deep=True)
    df = preprocessing.scale(df)

    pca = PCA()
    pca.fit(df)
    pca_data = pca.transform(df)
    per_var = np.round(pca.explained_variance_ratio_ * 100, 1)
    labels = ['PC' + str(i) for i in range(1, len(per_var) + 1)]
    pca_df = pd.DataFrame(pca_data, index=init_df.index, columns=labels)

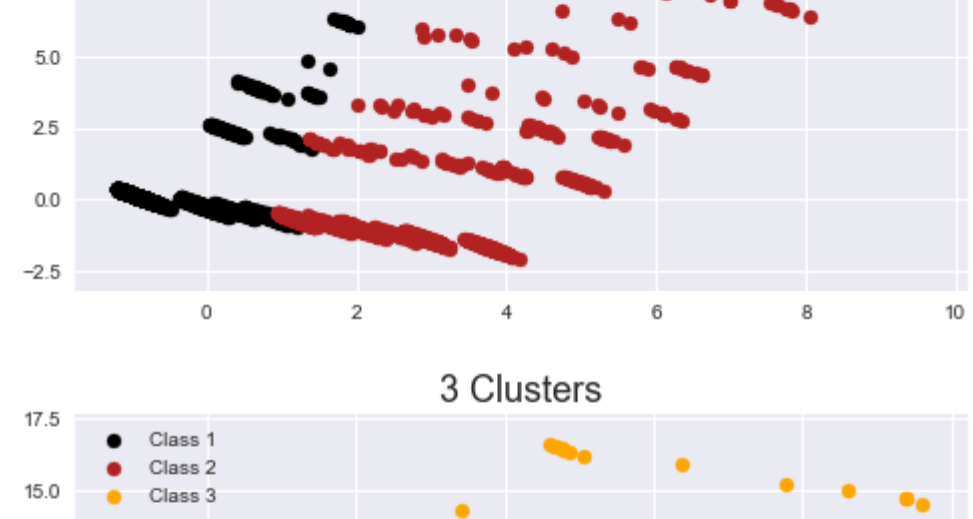
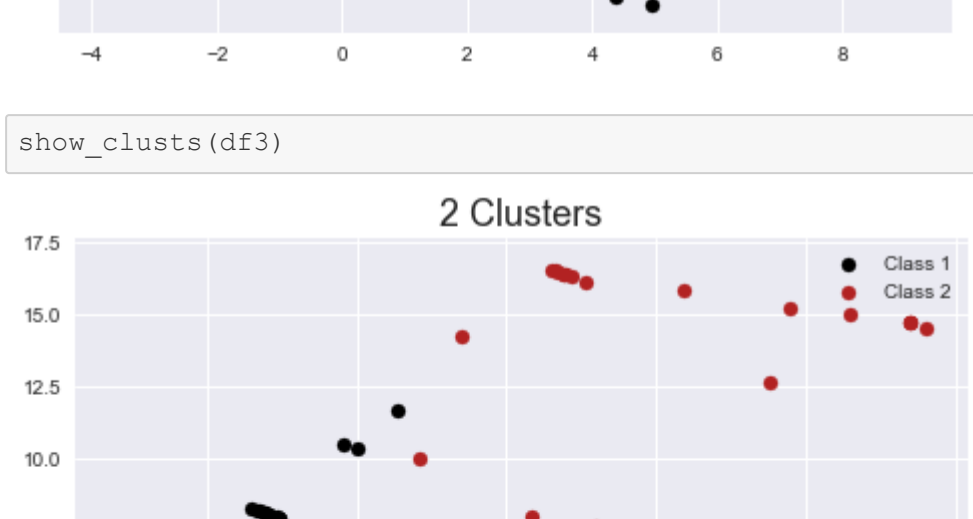
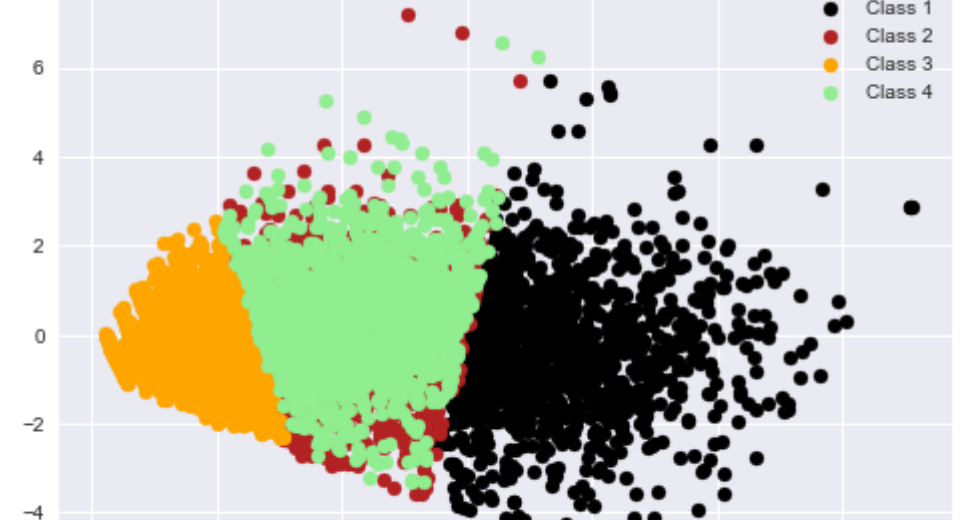
    cols = ['black', 'firebrick', 'orange', 'lightgreen',
            'mediumslateblue', 'dodgerblue', 'slateblue', 'violet']

    for n in range(2, 5):
        km = KMeans(n_clusters=n, n_init=10, random_state=0).fit(df)
        pca_df['class'] = km.labels_

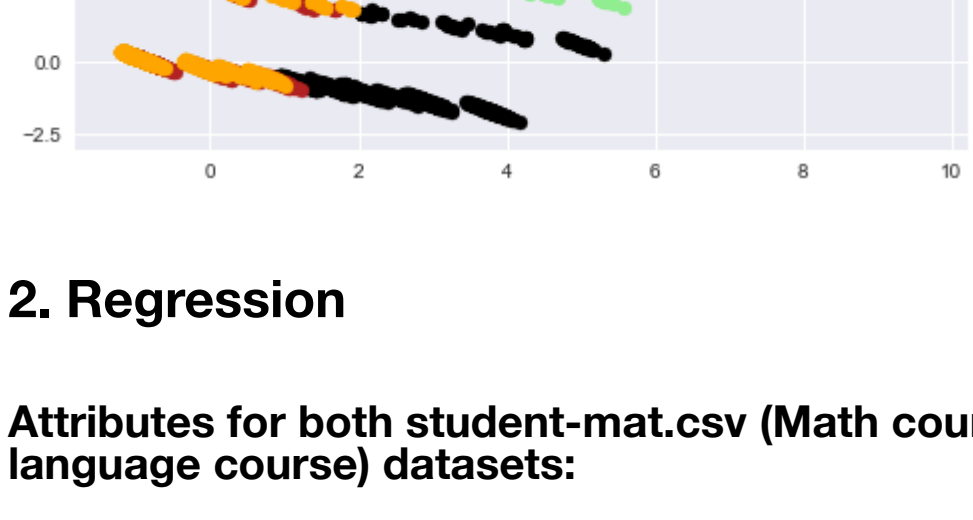
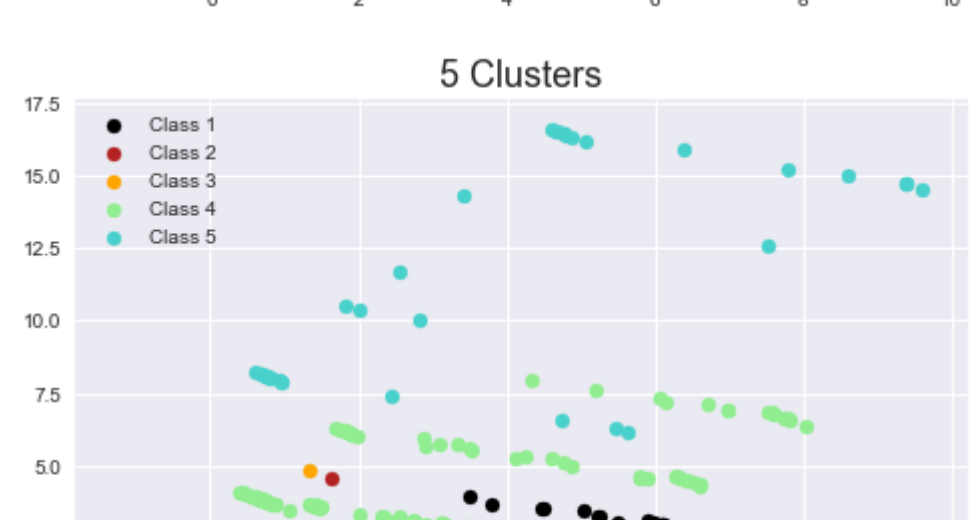
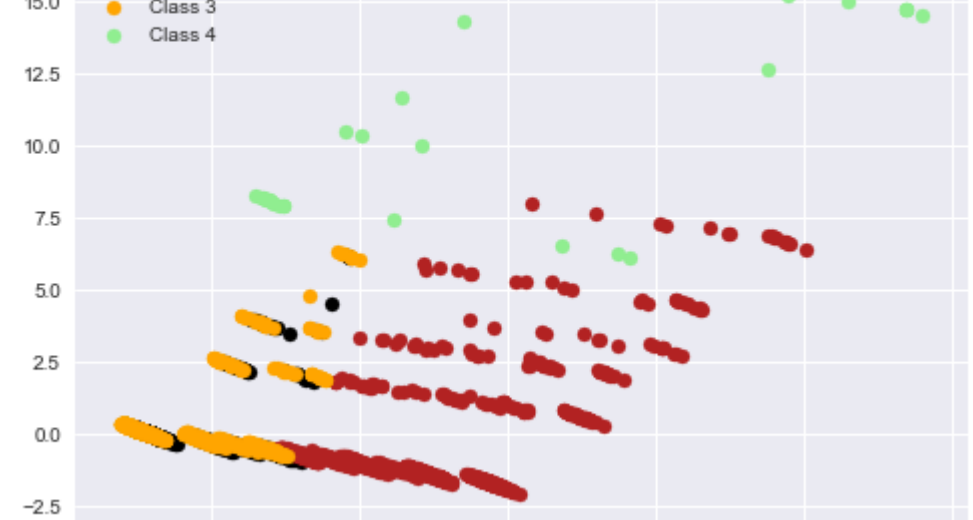
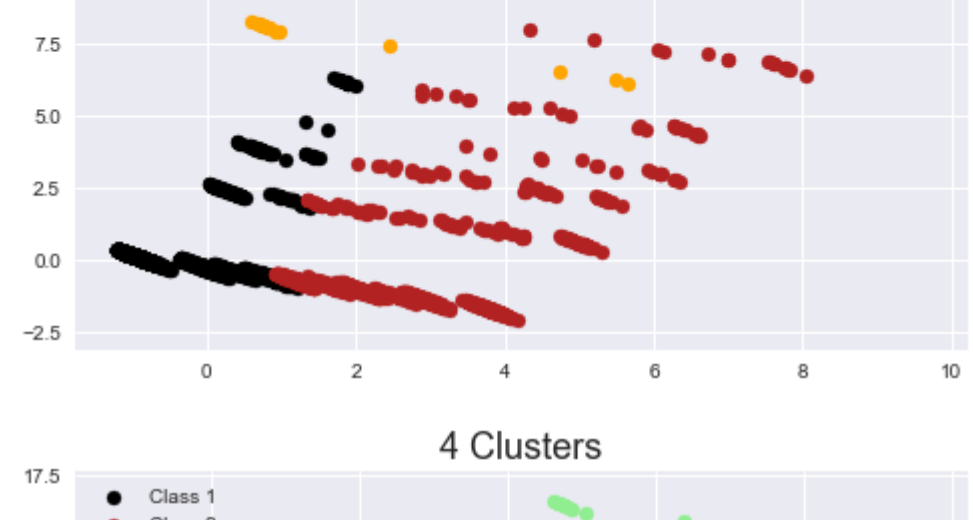
        for cl in np.unique(km.labels_):
            plt.scatter(pca_df[pca_df['class'] == cl].PC1,
                        pca_df[pca_df['class'] == cl].PC2,
                        c=cols[cl], label=f'Class {cl+1}')

        plt.title(f'({n}) Clusters')
        plt.legend()
        plt.show()
```

```
In [558]: show_clusts(init_df)
```



```
In [216]: show_clusts(df3)
```



2. Regression

Attributes for both student-mat.csv (Math course) and student-por.csv (Portuguese language course) datasets:

1 school - student's school (binary: 'GP' - Gabriel Pereira or 'MS' - Mousinho da Silveira)

2 sex - student's sex (binary: 'F' - female or 'M' - male)

3 age - student's age (numeric: from 15 to 22)

4 address - student's home address type (binary: 'U' - urban or 'R' - rural)

5 famsize - family size (binary: 'LE3' - less or equal to 3 or 'GT3' - greater than 3)

6 Pstatus - parent's cohabitation status (binary: 'T' - living together or 'A' - apart)

7 Medu - mother's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)

8 Fedu - father's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)

9 Mjob - mother's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at_home' or 'other')

10 Fjob - father's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at_home' or 'other')

11 reason - reason to choose this school (nominal: close to 'home', school 'reputation', 'course' preference or 'other')

12 guardian - student's guardian (nominal: 'mother', 'father' or 'other')

13 traveltime - home to school travel time (numeric: 1 - <15 min., 2 - 15 to 30 min., 3 - 30 min. to 1 hour, or 4 - >1 hour)

14 studytime - weekly study time (numeric: 1 - <2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4 - >10 hours)

15 failures - number of past class failures (numeric: n if 1<n<3, else 4)

16 schoolsup - extra educational support (binary: yes or no)

17 famsup - family educational support (binary: yes or no)

18 paid - extra paid classes within the course subject (Math or Portuguese) (binary: yes or no)

19 activities - extra-curricular activities (binary: yes or no)

20 nursery - attended nursery school (binary: yes or no)

21 higher - wants to take higher education (binary: yes or no)

22 internet - Internet access at home (binary: yes or no)

23 romantic - with a romantic relationship (binary: yes or no)

24 famrel - quality of family relationships (numeric: from 1 - very bad to 5 - excellent)

25 freetime - free time after school (numeric: from 1 - very low to 5 - very high)

26 goout - going out with friends (numeric: from 1 - very low to 5 - very high)

27 Dalc - workday alcohol consumption (numeric: from 1 - very low to 5 - very high)

28 Walc - weekend alcohol consumption (numeric: from 1 - very low to 5 - very high)

29 health - current health status (numeric: from 1 - very bad to 5 - very good)

30 absences - number of school absences (numeric: from 0 to 93)

these grades are related with the course subject, Math or Portuguese:

31 G1 - first period grade (numeric: from 0 to 20)

31 G2 - second period grade (numeric: from 0 to 20)

32 G3 - final grade (numeric: from 0 to 20, output target)

Preprocessing

```
In [651]: def prep(df, dummy=True):
    df.dropna(inplace=True)

    # Binary to 0/1
    df['school'] = np.where(df['school'] == 'GP', 0, 1)
    df['sex'] = np.where(df['sex'] == 'F', 0, 1)
    df['address'] = np.where(df['address'] == 'U', 1, 0)
    df['famsize'] = np.where(df['famsize'] == 'LE3', 0, 1)
    df['Pstatus'] = np.where(df['Pstatus'] == 'T', 1, 0)

    # Binary with yes/no to 0/1
    features = ['schoolsup', 'famsup', 'paid', 'activities',
               'nursery', 'higher', 'internet', 'romantic']
    for f in features:
        df[f] = np.where(df[f] == 'yes', 1, 0)

    if dummy:
        # Dummy vars for nominal features
        features = ['Mjob', 'Fjob', 'reason', 'guardian']
        df = pd.get_dummies(df[features])
        df.drop(features, axis=1, inplace=True)
        df = df.join(mf)

    return df
```

```
In [652]: init_df_mat = pd.read_csv('student-mat.csv', sep=';')
init_df_por = pd.read_csv('student-por.csv', sep=';')
```

```
In [683]: # Read and preprocess data
df_mat = prep(df_mat, dummy=True)
df_por = prep(df_por, dummy=True)

df_mat.head()
```

school	sex	age	address	famsize	Pstatus	Medu	Fedu	traveltime	studytime	...	Fjob	other	Fjob	services	Fjob	teacher	reason
0	1	0	0	18	1	1	1	1	1	2	...	0	0	0	0	0	1
1	0	0	17	1	1	1	1	1	1	1	2	...	1	0	0	0	0
2	0	0	15	1	0	1	1	1	1	1	1	2	...	1	0	0	0
3	0	0	15	1	1	1	4	2	1	3	...	0	0	1	0	0	0
4	0	0	16	1	1	1	3	3	1	2	...	1	0	0	0	0	0

5 rows × 48 columns


```
In [684]: df_pqr.head()

Out[684]:
  school sex age address famsize Pstatus Medu Fedu travelttime studytime failures schoolsup famsup paid
0      1  0  0  18      1      1      1      1      1      2      0      0      0      0      0
1      1  0  0  17      1      1      1      1      1      2      0      1      0      0      0
2      3  0  0  15      1      0      1      1      1      1      3      0      0      0      0
3      3  0  0  15      1      1      1      4      2      1      2      0      1      0      0
4      4  0  0  16      1      1      1      3      3      1      2      0      1      0      0

5 rows x 18 columns

Models

OLS

In [685]: import statsmodels.api as sm
import statsmodels.formula.api as smf
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error as mse

In [686]: df = df_pqr.copy(deep=True)
# df['const'] = np.ones(len(df))
df.const() = style.background_gradient(cmap='coolwarm').format({'0':1,1:F})

Out[686]:
  school sex age address famsize Pstatus Medu Fedu travelttime studytime failures schoolsup famsup paid
0      1  1.0 -0.1 0.1 -0.4 -0.0 -0.1 0.1 -0.3 -0.2 0.3 -0.1 0.1 -0.1 -0.1 -0.1
1      1  0.1 -1.0 -0.0 0.0 -0.0 -0.1 0.1 0.1 0.0 -0.2 0.1 -0.1 -0.1 -0.1 -0.0
2      3  0.1 -0.0 1.0 -0.0 -0.0 -0.0 -0.1 0.1 0.0 -0.0 0.3 -0.2 -0.1 -0.1 -0.0
3      3  0.4 0.0 -0.0 1.0 -0.0 -0.1 0.2 0.1 -0.3 0.1 -0.1 -0.1 0.0 -0.0 -0.0
4      3  0.0 0.1 0.0 -0.0 1.0 0.0 -0.1 0.2 0.0 -0.0 0.0 0.1 0.1 0.0 0.1
5      4  0.0 0.1 0.0 -0.0 1.0 0.0 0.0 -0.1 0.2 0.0 -0.0 0.0 -0.1 0.1 0.0 0.1
6      1  0.3 0.1 -0.1 0.1 0.2 0.1 -0.1 0.6 0.6 -0.3 0.1 -0.2 -0.0 0.1 0.1
7      1  0.2 0.1 -0.1 0.1 0.1 0.0 -0.0 0.8 1.0 -0.2 0.1 -0.2 0.0 0.1 0.1
8      3  0.3 0.0 0.0 0.0 -0.3 -0.0 -0.0 -0.3 -0.2 -0.1 -0.1 0.1 -0.0 -0.0 -0.1
9      3  0.1 0.1 0.3 0.3 -0.1 0.1 -0.0 -0.2 0.1 -0.1 -0.1 -0.1 0.1 0.1 -0.0
10     1  0.1 0.1 -0.2 0.0 0.1 -0.0 -0.0 0.0 -0.0 0.1 -0.0 -0.1 -0.0 1.0 0.1 0.0
11     1  0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.1 1.0 0.1
12     1  0.1 0.1 -0.0 -0.0 0.1 0.0 0.1 -0.0 -0.0 -0.0 0.1 0.0 0.1 0.0 0.1 1.0
13     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
14     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
15     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
16     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
17     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
18     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
19     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
20     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
21     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
22     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
23     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
24     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
25     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
26     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
27     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
28     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
29     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
30     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
31     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
32     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
33     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
34     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
35     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
36     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
37     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
38     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
39     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
40     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
41     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
42     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
43     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
44     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
45     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
46     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
47     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
48     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
49     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
50     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
51     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
52     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
53     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
54     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
55     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
56     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
57     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
58     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
59     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
60     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
61     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
62     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
63     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
64     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
65     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
66     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
67     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
68     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
69     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
70     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
71     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
72     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
73     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
74     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
75     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
76     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
77     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
78     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
79     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
80     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
81     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
82     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
83     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
84     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
85     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
86     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
87     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
88     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
89     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
90     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
91     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
92     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
93     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
94     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
95     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
96     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
97     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
98     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
99     1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0
100    1  0.1 0.1 -0.1 -0.1 0.0 0.0 0.0 0.1 0.1 -0.0 0.1 -0.0 0.1 -0.0 0.0 0.0

In [687]: # all columns except from G3
pred_cols_all = [i for i in df.columns if i != 'G3']
# all columns except from G1, G2, G3
pred_cols = [i for i in df.columns if i not in ['G1', 'G2', 'G3']]

# y/X
y = df['G3']
X = df[pred_cols]
X_all = df[pred_cols_all]

# train/test split
x_train, x_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, random_state=1)

x_train_all, x_test_all, y_train_all, y_test_all = train_test_split(
    X_all, y, test_size=0.25, random_state=1)

Model without G1 and G2

In [688]: mod = sm.OLS(y_train, x_train).fit()
mod.summary()

Out[688]:
OLS Regression Results

Dep. Variable:      G3      R-squared:      0.402
Model:              OLS      Adj. R-squared:  0.350
Method:             Least Squares      F-statistic:  7.686
Date:               Thu, 01 Apr 2021      Prob (F-statistic):  1.95e-30
Time:              19:46:39      Log-Likelihood: -1121.8
No. Observations:   486      AIC:      2324.
Df Residuals:      446      BIC:      2491.
Df Model:          39
Covariance Type:    nonrobust

coef    std err    t    P>[ ]    [0.025    0.975]
school   -1.7784    0.304   -5.867    0.000   -1.774   -0.579
sex       -0.5888    0.284  -2.081    0.038   -1.147   -0.033
age       0.1281    0.117    1.093    0.275   -0.102    0.359
address   0.3605    0.302    1.195    0.233   -0.232    0.953
famsize   -0.4531    0.275  -1.650    0.100   -0.993    0.087
Pstatus   0.1327    0.379    0.350    0.726   -0.612    0.878
Medu      0.2499    0.155    1.612    0.108   -0.055    0.555
Fedu      0.1808    0.182    0.993    0.321   -0.177    0.539
studytime 0.3780    0.154    2.462    0.014    0.076    0.680
failures  -1.3380    0.216  -6.193    0.000   -1.763    -0.913
schoolsup -1.3227    0.395  -3.352    0.001   -2.098    -0.547
famsup     0.0912    0.255    0.357    0.721   -0.410    0.593
paid       -0.6197    0.521  -1.188    0.235   -1.645    0.405
activities 0.4781    0.248    1.928    0.055   -0.009    0.965
nursery    0.0308    0.302    0.102    0.919   -0.563    0.625
higher     1.8337    0.414    4.427    0.000    1.020    2.648
internet   0.5399    0.306    1.766    0.078   -0.061    1.141
romantic   -0.4005    0.256  -1.564    0.118   -0.904    0.103
famelrel   0.2394    0.131    1.830    0.068   -0.018    0.496
freetime   -0.1631    0.124  -1.312    0.180   -0.407    0.081
goutout    -0.0926    0.120  -0.771    0.441   -0.339    0.145
Dalc       -0.0588    0.165  -0.356    0.722   -0.383    0.266
Walc       -0.1649    0.131  -1.256    0.210   -0.423    0.093
health     -0.1549    0.086  -1.809    0.071   -0.323    0.013
absences   -0.0301    0.027  -1.101    0.272   -0.084    0.024
Mjob_at_home 1.9848    0.599    3.299    0.001    0.817    3.172
Mjob_health 1.5642    0.507    3.087    0.002    0.569    2.560
Mjob_other 1.6755    0.546    3.066    0.002    0.602    2.749
Mjob_teacher 2.1950    0.611    3.143    0.002    0.719    3.189
Fjob_at_home 2.0185    0.642    3.145    0.002    0.757    3.310
Fjob_health 1.3289    0.686    1.937    0.053   -0.020    2.678
Fjob_other 1.7755    0.512    3.471    0.001    0.770    2.781
Fjob_services 1.3929    0.519    2.684    0.008    0.373    2.413
Fjob_teacher 2.0698    0.683    3.032    0.003    0.728    3.411
reason_course 2.0593    0.598    3.446    0.001    0.885    3.234
reason_home 2.2429    0.620    3.618    0.000    1.025    3.461
reason_other 1.5938    0.649    2.455    0.014    0.318    2.870
reason_reputation 2.6887    0.645    4.169    0.000    1.422    3.958
guardian_father 2.7958    0.769    3.634    0.000    1.281    4.300
guardian_mother 2.7978    0.752    3.723    0.000    1.321    4.275
guardian_other 2.9974    0.951    3.153    0.002    1.129    4.866

Omnibus: 77.462   Durbin-Watson:      2.107
Prob(Omnibus): 0.000   Jarque-Bera (JB): 230.745
Skew: -0.748      Prob(JB): 7.84e-51
Kurtosis: 6.026      Cond. No. 1.22e+16

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The smallest eigenvalue is 1.23e-27. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

In [689]: y_pred = mod.predict(x_test)
print(f'MSE: {mse(y_pred, y_test)}')

MSE: 9.52399245086358

Model with G1 and G2

In [690]: mod_all = sm.OLS(y_train_all, x_train_all).fit()
mod_all.summary()

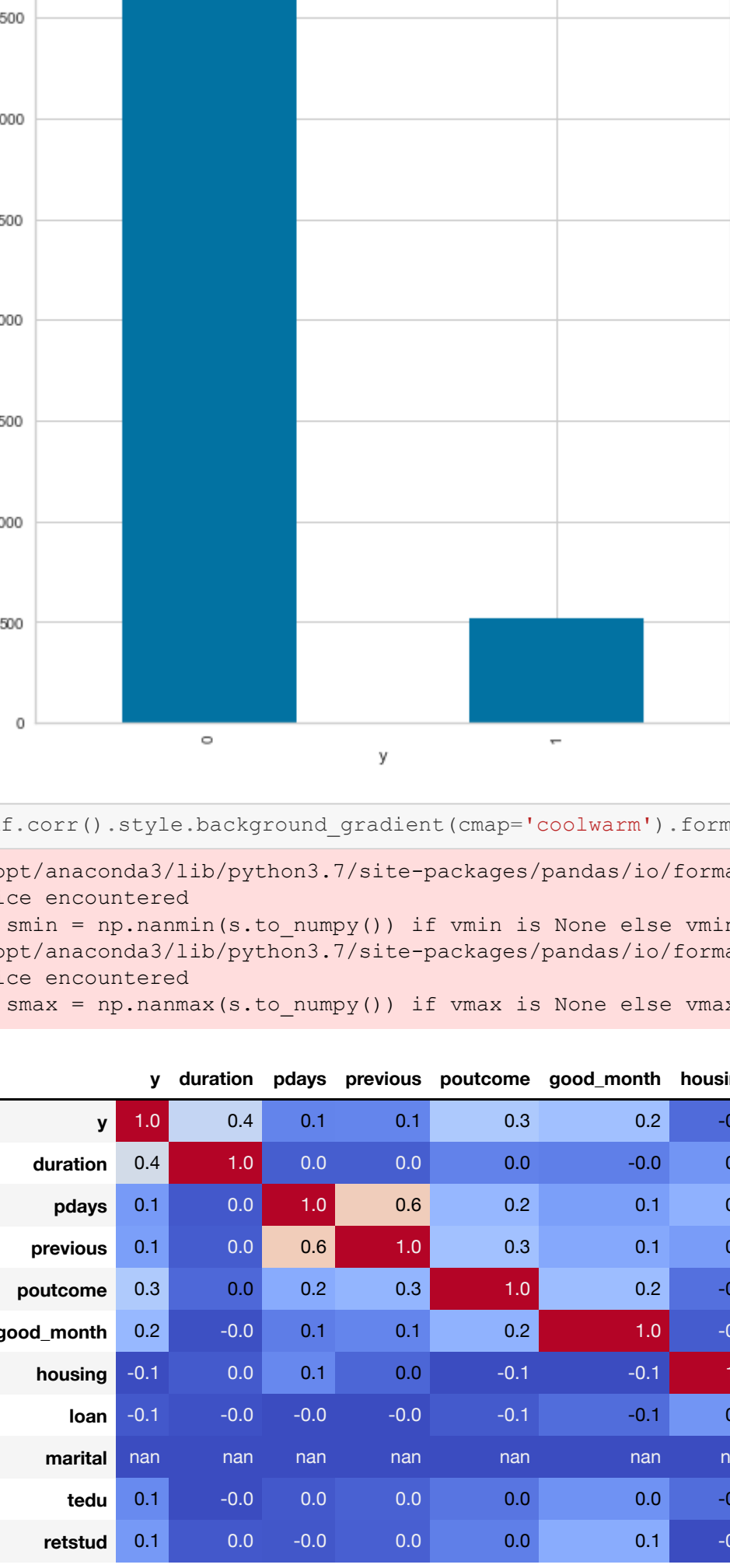
Out[690]:
OLS Regression Results

Dep. Variable:      G3      R-squared:      0.859
Model:              OLS      Adj. R-squared:  0.845
Method:             Least Squares      F-statistic:  65.72
Date:               Thu, 01 Apr 2021      Prob (F-statistic):  3.67e-162
Time:              19:46:55      Log-Likelihood: -771.46
No. Observations:   486      AIC:      1627.
Df Residuals:      444      BIC:      1803.
Df Model:          41
Covariance Type:    nonrobust

coef    std err    t    P>[ ]    [0.025    0.975]
school   -0.2543    0.152  -1.673    0.095   -0.553    0.044
sex      -0.1100    0.139  -0.782    0.429   -0.383    0.163
age       0.0069    0.058    0.120    0.905   -0.106    0.120
address   0.1373    0.147    0.933    0.352   -0.152    0.427
famsize   -0.0580    0.135  -0.431    0.667   -0.322    0.207
Pstatus   -0.0778    0.185  -0.421    0.674   -0.441    0.285
Medu      0.0848    0.078    1.118    0.264   -0.064    0.233
Fedu      0.1021    0.089    1.150    0.251   -0.072    0.277
studytime 0.0654    0.076    0.864    0.388   -0.083    0.214
failures  -0.2051    0.110  -1.870    0.062   -0.421    0.010
schoolsup -0.2280    0.195  -1.167    0.244   -0.612    0.166
famsup     0.1147    0.124    0.922    0.357   -0.130    0.359
paid      -0.3151    0.256  -1.231    0.219   -0.818    0.188
activities 0.0677    0.121    0.557    0.577   -0.171    0.306
nursery    0.0274    0.147  -0.186    0.853   -0.317    0.262
higher     0.3829    0.206    1.856    0.064   -0.023    0.788
internet   0.1543    0.149    1.033    0.302   -0.139    0.448
romantic   -0.0491    0.125  -0.392    0.695   -0.295    0.197
romantic   -0.0543    0.065  -0.840    0.402   -0.181    0.073
famelrel   -0.0547    0.061  -0.900    0.369   -0.174    0.065
freetime   -0.0413    0.059  -0.703    0.482   -0.157    0.074
goutout    -0.0480    0.081    0.596    0.551   -0.110    0.206
Dalc       -0.0567    0.064  -0.917    0.360   -0.185    0.067
Walc       -0.0524    0.042  -1.251    0.212   -0.135    0.030
health     0.0124    0.013    0.923    0.356   -0.014    0.039
absences   0.1255    0.047    2.654    0.008    0.033    0.218
G1        0.8591    0.043   19.852    0.000    0.774    0.944
G2        0.8591    0.043   19.852    0.000    0.774    0.944
Mjob_at_home 0.0644    0.276    0.234    0.815   -0.477    0.606
Mjob_health 0.3355    0.297    1.129    0.260   -0.249    0.920
Mjob_other 0.0045    0.255    0.018    0.986   -0.496    0.505
Mjob_services 0.0679    0.273    0.249    0.804   -0.469    0.604
Mjob_teacher 0.2084    0.304    0.686    0.493   -0.398    0.805
Fjob_at_home 0.5388    0.317    1.700    0.090   -0.084    1.162
Fjob_health -0.0467    0.339  -0.138    0.890   -0.712    0.619
Fjob_other 0.2034    0.256    0.794    0.427   -0.300    0.707
Fjob_services 0.1763    0.259    0.680    0.497   -0.333    0.685
Fjob_teacher -0.1910    0.343  -0.557    0.578   -0.885    0.483
reason_course 0.3906    0.298    1.311    0.195   -0.195    0.976
reason_home 0.1955    0.312    0.626    0.531   -0.418    0.809
reason_other -0.1690    0.325  -0.581    0.561   -0.828    0.450
reason_reputation 0.2935    0.325    0.872    0.384   -0.355    0.923
guardian_father 0.2204    0.388    0.568    0.570   -0.542    0.962
guardian_mother 0.1802    0.380    0.475    0.635   -0.566    0.926
guardian_other 0.
```



```
In [1016]: df.groupby('y')['l'].count().plot.bar()
Out[1016]: <matplotlib.axes._subplots.AxesSubplot at 0x7f7fa70f5fd0>
```



```
In [995]: df.corr().style.background_gradient(cmap='coolwarm').format({'0':,1F'})

/opt/anaconda3/lib/python3.7/site-packages/pandas/io/formats/style.py:1264: RuntimeWarning: All-NaN s
lice encountered
  vmin = np.nanmin(s.to_numpy()) if vmin is None else vmin
/opt/anaconda3/lib/python3.7/site-packages/pandas/io/formats/style.py:1265: RuntimeWarning: All-NaN s
lice encountered
  smax = np.nanmax(s.to_numpy()) if vmax is None else vmax
```

Out[995]:

	y	duration	pdays	previous	poutcome	good_month	housing	loan	marital	tedu	retstud
duration	1.0	0.4	0.1	0.1	0.3	0.2	-0.1	-0.1	nan	0.1	0.1
pdays	0.1	0.0	1.0	0.6	0.2	-0.0	0.0	-0.0	nan	-0.0	0.0
previous	0.1	0.0	0.6	1.0	0.3	0.1	0.1	-0.0	nan	0.0	-0.0
poutcome	0.3	0.0	0.2	0.3	1.0	0.2	-0.1	-0.1	nan	0.0	0.0
good_month	0.2	-0.0	0.1	0.1	0.2	1.0	-0.1	-0.1	nan	0.0	0.1
housing	-0.1	0.0	0.1	0.0	-0.1	-0.1	1.0	0.0	nan	-0.1	-0.2
loan	-0.1	-0.0	-0.0	-0.0	-0.1	-0.1	0.0	1.0	nan	-0.0	-0.0
marital	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
tedu	0.1	-0.0	0.0	0.0	0.0	0.0	-0.1	-0.0	nan	1.0	-0.1
retstud	0.1	0.0	-0.0	0.0	0.0	0.1	-0.2	-0.0	nan	-0.1	1.0

```
In [1047]: pred_cols = [i for i in df.columns if i != 'y']
y = df['y']
X = df[pred_cols]

# train/test
x_train, x_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, random_state=1)

train_df = x_train.join(y_train)
test_df = x_test.join(y_test)

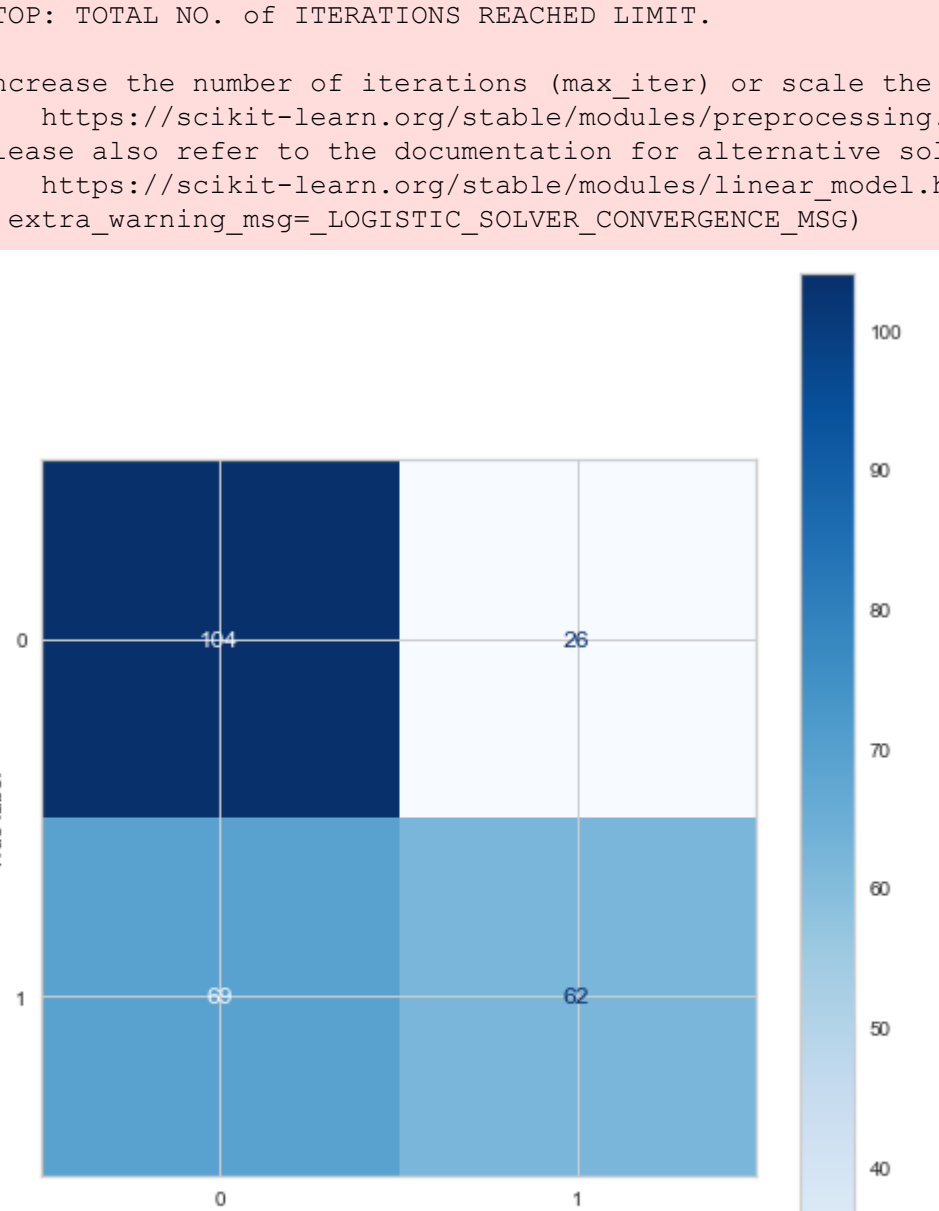
mod = LogisticRegression().fit(x_train, y_train)
plot_confusion_matrix(mod, x_test, y_test, cmap=plt.cm.Blues);

y_pred = mod.predict(x_test)
print('ROC: ', roc_auc_score(y_pred, y_test))
print('F1: ', f1_score(y_test, y_pred, zero_division=1))

ROC: 0.7658725368900692
F1: 0.3804878048780488

/opt/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:765: ConvergenceWarning:
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
```



```
In [1051]: # undersampling
df = init_df[['y', 'pdays', 'previous']].copy(deep=True)

df['poutcome'] = np.where(init_df['poutcome'] == 'success', 1, 0)
df['good_month'] = np.where(init_df['month'].isin(['dec', 'mar', 'oct', 'sep']), 1, 0)
df['housing'] = np.where(init_df['housing'] == 'yes', 1, 0)
df['loan'] = np.where(init_df['loan'] == 'yes', 1, 0)
df['marital'] = np.where(init_df['marital'] == 'married', 1, 0)
df['tedu'] = np.where(init_df['education'] == 'tertiary', 1, 0)
df['retstud'] = np.where(init_df['job'].isin(['retired', 'student']), 1, 0)

pred_cols = [i for i in df.columns if i != 'y']

c11 = df[df['y'] == 0].sample(n=521)
c12 = df[df['y'] == 1]

df = c11.append(c12)
y = df['y']
X = df[pred_cols]

# train/test
x_train, x_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, random_state=1)

train_df = x_train.join(y_train)
test_df = x_test.join(y_test)

mod = LogisticRegression().fit(x_train, y_train)
plot_confusion_matrix(mod, x_test, y_test, cmap=plt.cm.Blues);

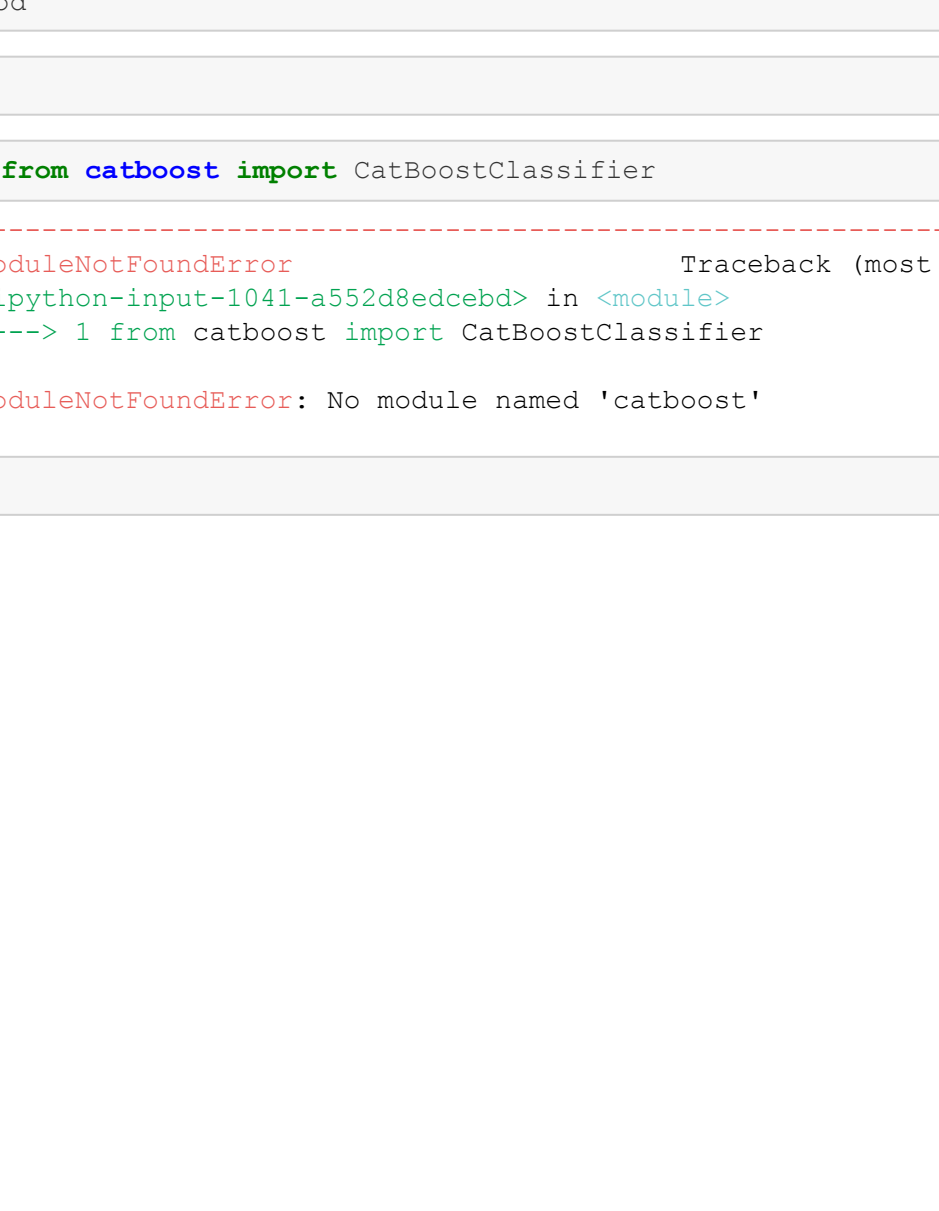
y_pred = mod.predict(X)

print('ROC: ', roc_auc_score(y_pred, y))
print('F1: ', f1_score(y, y_pred, zero_division=1))
print('Acc: ', accuracy_score(y_pred, y))

ROC: 0.681032851745193
F1: 0.6098654708520179
Acc: 0.6660268714011516

/opt/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:765: ConvergenceWarning:
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
```



Trees

```
In [1043]: from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import cross_val_score

mod = DecisionTreeClassifier()
cross_val_score(mod, x_train, y_train, cv=10)

Out[1043]: array([0.75949367, 0.75641026, 0.73076923, 0.73076923, 0.74358974,
0.71794872, 0.82051282, 0.65384615, 0.69230769, 0.74358974])
```

```
In [ ]: mod

In [ ]:
```

```
In [1041]: from catboost import CatBoostClassifier

ModuleNotFoundError Traceback (most recent call last)
<ipython-input-1041-c552d8ede6bd> in <module>
----> 1 from catboost import CatBoostClassifier

ModuleNotFoundError: No module named 'catboost'

In [ ]:
```