# STUDENTS PERFORMANCE PREDICTION

# Introduction

Education is an important element of the society, every government and country in the world work so hard to improve this sector. With the corona-virus outbreak that has disrupted life around the globe in 2020, the educational systems have been affected in many ways; studies show that student's performance has decreased since then, which highlights the need to deal with this problem more seriously and try to find effective solutions, as well as the influencing factors.

# Motivation

The educational systems need, at this specific time, innovative ways to improve quality of education to achieve the best results and decrease the failure rate.

As students of the IT department who studied in the last month a little about machine learning, we know that in order for an institute to provide quality education to learners, deep analysis of previous records of the learners can play a vital role, and wanted to work on this challenging task.

# Objectives

In this notebook, we will:

- **Predict whether or not a student will pass the final exam based on certain information given**
- **Compare the three learning algorithms**
- **Find out what most affects student achievement**
- **Find the best algorithm with high accuracy**

We will be using three learning algorithms:

- **Logistic regression**
- **Supported vector machine**
- **KNN**

# Problem Statement:

As already mentioned, with the help of the old students records, we can came up with a model that can let us help students improve their performance in exams by predicting the student success. So, it is obvious it's a problem of classification , and we will classify a student based on his given informations, and we will also use different classifiers such as KNN or SVM classifier and compare between them. Many factors affect a student performance in exams like family problems or alcohol consumption, and by using our skills in machine learning we want to:

   1) predict whether a student will pass his final exam or not.
   2) came up with the best classifier that is more accurate and avoid
   overfitting and underfitting by using simple techniques.
   3) know what the most factors affect a student performance.

So, teachers and parents will be able to intervene before students reach the exam stage and solve the problems.

# Dataset:

Dataset name: Student.csv

Source: Kaggle

# EDA:

```
df.shape

(395, 31)
```

```
df.dropna().shape # their is no null value "fortunately:)"

(395, 31)
```

```
df.columns

Index(['school', 'sex', 'age', 'address', 'famsize', 'Pstatus', 'Medu', 'Fedu',
       'Mjob', 'Fjob', 'reason', 'guardian', 'traveltime', 'studytime',
       'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery',
       'higher', 'internet', 'romantic', 'famrel', 'freetime', 'goout', 'Dalc',
       'Walc', 'health', 'absences', 'passed'],
      dtype='object')
```
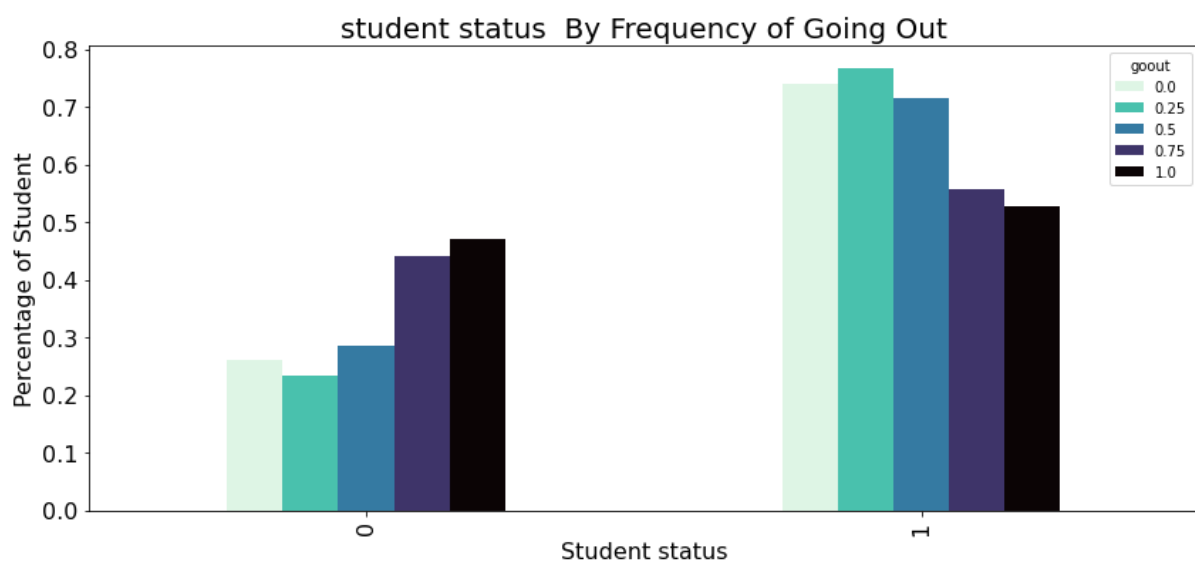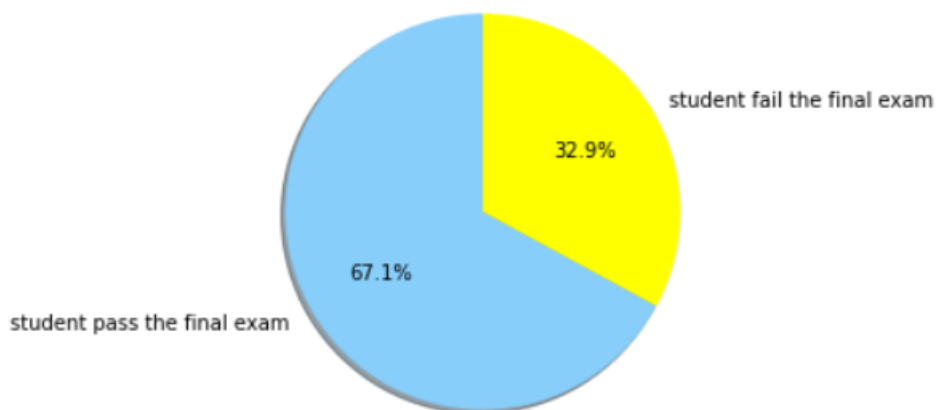
```
features=['school', 'sex', 'age', 'address', 'famsize', 'Pstatus', 'Medu', 'Fedu',
       'Mjob', 'Fjob', 'reason', 'guardian', 'traveltime', 'studytime',
       'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery',
       'higher', 'internet', 'romantic', 'famrel', 'freetime', 'goout', 'Dalc',
       'Walc', 'health', 'absences']
```

```python
#plot of student status
dfv['passed'].value_counts()
```
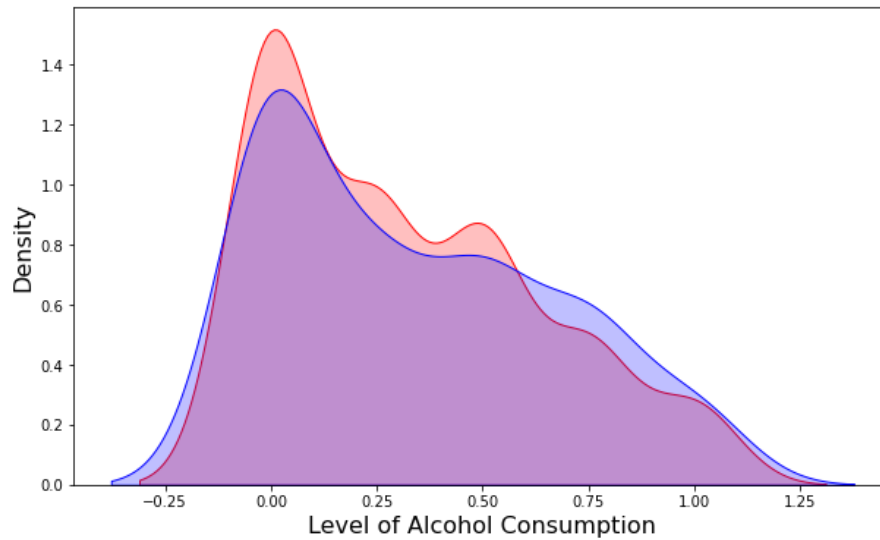
```
yes      265
no       130
Name: passed, dtype: int64
```

```python
df["goout"].unique()
```

```
array([0.75, 0.5 , 0.25, 0.  , 1.  ])
```

Good Performance vs. Poor Performance Student Weekend Alcohol Consumption

# Pre-processing (dimensionality reduction):

- Scaling:

```
# Let's scal our features
feature_scaling(df)

# Now we are ready for models training
df
```

| | school | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob | Fjob | ... | internet | romantic | famrel | freetime | goout | Dalc | Walc | health | absences | passed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 1.0 | 0.059264 | 0.0 | 1.0 | 1.0 | 1.00 | 1.00 | 0.75 | 0.00 | ... | 0.0 | 0.0 | 0.75 | 0.50 | 0.75 | 0.00 | 0.00 | 0.50 | 0.003882 | 0.0 |
| 1 | 0.0 | 1.0 | 0.013809 | 0.0 | 1.0 | 0.0 | 0.25 | 0.25 | 0.75 | 1.00 | ... | 1.0 | 0.0 | 1.00 | 0.50 | 0.50 | 0.00 | 0.00 | 0.50 | -0.022785 | 0.0 |
| 2 | 0.0 | 1.0 | -0.077100 | 0.0 | 0.0 | 0.0 | 0.25 | 0.25 | 0.75 | 1.00 | ... | 1.0 | 0.0 | 0.75 | 0.50 | 0.25 | 0.25 | 0.50 | 0.50 | 0.057215 | 1.0 |
| 3 | 0.0 | 1.0 | -0.077100 | 0.0 | 1.0 | 0.0 | 1.00 | 0.50 | 0.25 | 0.50 | ... | 1.0 | 1.0 | 0.50 | 0.25 | 0.25 | 0.00 | 0.00 | 1.00 | -0.049451 | 1.0 |
| 4 | 0.0 | 1.0 | -0.031646 | 0.0 | 1.0 | 0.0 | 0.75 | 0.75 | 1.00 | 1.00 | ... | 0.0 | 0.0 | 0.75 | 0.50 | 0.25 | 0.00 | 0.25 | 1.00 | -0.022785 | 1.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 390 | 1.0 | 0.0 | 0.150173 | 0.0 | 0.0 | 1.0 | 0.50 | 0.50 | 0.50 | 0.50 | ... | 0.0 | 0.0 | 1.00 | 1.00 | 0.75 | 0.75 | 1.00 | 0.75 | 0.070549 | 0.0 |
| 391 | 1.0 | 0.0 | 0.013809 | 0.0 | 0.0 | 0.0 | 0.75 | 0.25 | 0.50 | 0.50 | ... | 1.0 | 0.0 | 0.25 | 0.75 | 1.00 | 0.50 | 0.75 | 0.25 | -0.036118 | 1.0 |
| 392 | 1.0 | 0.0 | 0.195627 | 1.0 | 1.0 | 0.0 | 0.25 | 0.25 | 1.00 | 1.00 | ... | 0.0 | 0.0 | 1.00 | 1.00 | 0.50 | 0.50 | 0.50 | 0.50 | -0.036118 | 0.0 |
| 393 | 1.0 | 0.0 | 0.059264 | 1.0 | 0.0 | 0.0 | 0.75 | 0.50 | 0.50 | 1.00 | ... | 1.0 | 0.0 | 0.75 | 0.75 | 0.00 | 0.50 | 0.75 | 1.00 | -0.076118 | 1.0 |
| 394 | 1.0 | 0.0 | 0.104718 | 0.0 | 0.0 | 0.0 | 0.25 | 0.25 | 1.00 | 0.75 | ... | 1.0 | 0.0 | 0.50 | 0.25 | 0.50 | 0.50 | 0.50 | 1.00 | -0.009451 | 0.0 |

395 rows × 31 columns

- Encoding:

```python
def numerical_data():
    df['school'] = df['school'].map({'GP': 0, 'MS': 1})
    df['sex'] = df['sex'].map({'M': 0, 'F': 1})
    df['address'] = df['address'].map({'U': 0, 'R': 1})
    df['famsize'] = df['famsize'].map({'LE3': 0, 'GT3': 1})
    df['Pstatus'] = df['Pstatus'].map({'T': 0, 'A': 1})
    df['Mjob'] = df['Mjob'].map({'teacher': 0, 'health': 1, 'services': 2, 'at_home': 3, 'other': 4})
    df['Fjob'] = df['Fjob'].map({'teacher': 0, 'health': 1, 'services': 2, 'at_home': 3, 'other': 4})
    df['reason'] = df['reason'].map({'home': 0, 'reputation': 1, 'course': 2, 'other': 3})
    df['guardian'] = df['guardian'].map({'mother': 0, 'father': 1, 'other': 2})
    df['schoolsup'] = df['schoolsup'].map({'no': 0, 'yes': 1})
    df['famsup'] = df['famsup'].map({'no': 0, 'yes': 1})
    df['paid'] = df['paid'].map({'no': 0, 'yes': 1})
    df['activities'] = df['activities'].map({'no': 0, 'yes': 1})
    df['nursery'] = df['nursery'].map({'no': 0, 'yes': 1})
    df['higher'] = df['higher'].map({'no': 0, 'yes': 1})
    df['internet'] = df['internet'].map({'no': 0, 'yes': 1})
    df['romantic'] = df['romantic'].map({'no': 0, 'yes' : 1})
    df['passed'] = df['passed'].map({'no': 0, 'yes': 1})
    # reorder dataframe columns :
    col = df['passed']
    del df['passed']
    df['passed'] = col
```

**SOURCE CODE & OUTPUT:**

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from time import time
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split,GridSearchCV
from sklearn.metrics import confusion_matrix, roc_curve, accuracy_score, f1_score, roc_auc_score, classification_report
from astropy.table import Table
from sklearn.metrics import roc_auc_score

df = pd.read_csv('student.csv')
dfv = pd.read_csv('student.csv')
```

```python
df
```

| | school | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob | Fjob | ... | internet | romantic | famrel | freetime | goout | Dalc | Walc | health | absences | passed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | GP | F | 18 | U | GT3 | A | 4 | 4 | at_home | teacher | ... | no | no | 4 | 3 | 4 | 1 | 1 | 3 | 6 | no |
| 1 | GP | F | 17 | U | GT3 | T | 1 | 1 | at_home | other | ... | yes | no | 5 | 3 | 3 | 1 | 1 | 3 | 4 | no |
| 2 | GP | F | 15 | U | LE3 | T | 1 | 1 | at_home | other | ... | yes | no | 4 | 3 | 2 | 2 | 3 | 3 | 10 | yes |
| 3 | GP | F | 15 | U | GT3 | T | 4 | 2 | health | services | ... | yes | yes | 3 | 2 | 2 | 1 | 1 | 5 | 2 | yes |
| 4 | GP | F | 16 | U | GT3 | T | 3 | 3 | other | other | ... | no | no | 4 | 3 | 2 | 1 | 2 | 5 | 4 | yes |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 390 | MS | M | 20 | U | LE3 | A | 2 | 2 | services | services | ... | no | no | 5 | 5 | 4 | 4 | 5 | 4 | 11 | no |
| 391 | MS | M | 17 | U | LE3 | T | 3 | 1 | services | services | ... | yes | no | 2 | 4 | 5 | 3 | 4 | 2 | 3 | yes |

```python
def numerical_data():
    df['school'] = df['school'].map({'GP': 0, 'MS': 1})
    df['sex'] = df['sex'].map({'M': 0, 'F': 1})
    df['address'] = df['address'].map({'U': 0, 'R': 1})
    df['famsize'] = df['famsize'].map({'LE3': 0, 'GT3': 1})
    df['Pstatus'] = df['Pstatus'].map({'T': 0, 'A': 1})
    df['Mjob'] = df['Mjob'].map({'teacher': 0, 'health': 1, 'services': 2, 'at_home': 3, 'other': 4})
    df['Fjob'] = df['Fjob'].map({'teacher': 0, 'health': 1, 'services': 2, 'at_home': 3, 'other': 4})
    df['reason'] = df['reason'].map({'home': 0, 'reputation': 1, 'course': 2, 'other': 3})
    df['guardian'] = df['guardian'].map({'mother': 0, 'father': 1, 'other': 2})
    df['schoolsup'] = df['schoolsup'].map({'no': 0, 'yes': 1})
    df['famsup'] = df['famsup'].map({'no': 0, 'yes': 1})
    df['paid'] = df['paid'].map({'no': 0, 'yes': 1})
    df['activities'] = df['activities'].map({'no': 0, 'yes': 1})
    df['nursery'] = df['nursery'].map({'no': 0, 'yes': 1})
    df['higher'] = df['higher'].map({'no': 0, 'yes': 1})
    df['internet'] = df['internet'].map({'no': 0, 'yes': 1})
    df['romantic'] = df['romantic'].map({'no': 0, 'yes' : 1})
    df['passed'] = df['passed'].map({'no': 0, 'yes': 1})
    # reorder dataframe columns :
    col = df['passed']
    del df['passed']
    df['passed'] = col

    # feature scaling will allow the algorithm to converge faster, large data will have same scal
def feature_scaling(df):
    for i in df:
        col = df[i]
        # let's choose columns that have large values
        if(np.max(col)>6):
            Max = max(col)
            Min = min(col)
            mean = np.mean(col)
            col = (col-mean)/(Max)
            df[i] = col
```

```python
# All values in numerical after calling numerical_data() function
numerical_data()
df
```
Python

|  | school | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob | Fjob | ... | internet | romantic | famrel | freetime | goout | Dalc | Walc | health | absences | passed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 18 | 0 | 1 | 1 | 4 | 4 | 3 | 0 | ... | 0 | 0 | 4 | 3 | 4 | 1 | 1 | 3 | 6 | 0 |
| 1 | 0 | 1 | 17 | 0 | 1 | 0 | 1 | 1 | 3 | 4 | ... | 1 | 0 | 5 | 3 | 3 | 1 | 1 | 3 | 4 | 0 |
| 2 | 0 | 1 | 15 | 0 | 0 | 0 | 1 | 1 | 3 | 4 | ... | 1 | 0 | 4 | 3 | 2 | 2 | 3 | 3 | 10 | 1 |
| 3 | 0 | 1 | 15 | 0 | 1 | 0 | 4 | 2 | 1 | 2 | ... | 1 | 1 | 3 | 2 | 2 | 1 | 1 | 5 | 2 | 1 |
| 4 | 0 | 1 | 16 | 0 | 1 | 0 | 3 | 3 | 4 | 4 | ... | 0 | 0 | 4 | 3 | 2 | 1 | 2 | 5 | 4 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 390 | 1 | 0 | 20 | 0 | 0 | 1 | 2 | 2 | 2 | 2 | ... | 0 | 0 | 5 | 5 | 4 | 4 | 5 | 4 | 11 | 0 |
| 391 | 1 | 0 | 17 | 0 | 0 | 0 | 3 | 1 | 2 | 2 | ... | 1 | 0 | 2 | 4 | 5 | 3 | 4 | 2 | 3 | 1 |
| 392 | 1 | 0 | 21 | 1 | 1 | 0 | 1 | 1 | 4 | 4 | ... | 0 | 0 | 5 | 5 | 3 | 3 | 3 | 3 | 3 | 0 |
| 393 | 1 | 0 | 18 | 1 | 0 | 0 | 3 | 2 | 2 | 4 | ... | 1 | 0 | 4 | 4 | 1 | 3 | 4 | 5 | 0 | 1 |
| 394 | 1 | 0 | 19 | 0 | 0 | 0 | 1 | 1 | 4 | 3 | ... | 1 | 0 | 3 | 2 | 3 | 3 | 3 | 5 | 5 | 0 |

395 rows × 31 columns

```python
# Let's scal our features
feature_scaling(df)

# Now we are ready for models training
df
```
Python

|  | school | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob | Fjob | ... | internet | romantic | famrel | freetime | goout | Dalc | Walc | health | absences | passed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 1.0 | 0.059264 | 0.0 | 1.0 | 1.0 | 1.00 | 1.00 | 0.75 | 0.00 | ... | 0.0 | 0.0 | 0.75 | 0.50 | 0.75 | 0.00 | 0.00 | 0.50 | 0.003882 | 0.0 |
| 1 | 0.0 | 1.0 | 0.013809 | 0.0 | 1.0 | 0.0 | 0.25 | 0.25 | 0.75 | 1.00 | ... | 1.0 | 0.0 | 1.00 | 0.50 | 0.50 | 0.00 | 0.00 | 0.50 | -0.022785 | 0.0 |
| 2 | 0.0 | 1.0 | -0.077100 | 0.0 | 0.0 | 0.0 | 0.25 | 0.25 | 0.75 | 1.00 | ... | 1.0 | 0.0 | 0.75 | 0.50 | 0.25 | 0.25 | 0.50 | 0.50 | 0.057215 | 1.0 |
| 3 | 0.0 | 1.0 | -0.077100 | 0.0 | 1.0 | 0.0 | 1.00 | 0.50 | 0.25 | 0.50 | ... | 1.0 | 1.0 | 0.50 | 0.25 | 0.25 | 0.00 | 0.00 | 1.00 | -0.049451 | 1.0 |
| 4 | 0.0 | 1.0 | -0.031646 | 0.0 | 1.0 | 0.0 | 0.75 | 0.75 | 1.00 | 1.00 | ... | 0.0 | 0.0 | 0.75 | 0.50 | 0.25 | 0.00 | 0.25 | 1.00 | -0.022785 | 1.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 390 | 1.0 | 0.0 | 0.150173 | 0.0 | 0.0 | 1.0 | 0.50 | 0.50 | 0.50 | 0.50 | ... | 0.0 | 0.0 | 1.00 | 1.00 | 0.75 | 0.75 | 1.00 | 0.75 | 0.070549 | 0.0 |
| 391 | 1.0 | 0.0 | 0.013809 | 0.0 | 0.0 | 0.0 | 0.75 | 0.25 | 0.50 | 0.50 | ... | 1.0 | 0.0 | 0.25 | 0.75 | 1.00 | 0.50 | 0.75 | 0.25 | -0.036118 | 1.0 |
| 392 | 1.0 | 0.0 | 0.195627 | 1.0 | 1.0 | 0.0 | 0.25 | 0.25 | 1.00 | 1.00 | ... | 0.0 | 0.0 | 1.00 | 1.00 | 0.50 | 0.50 | 0.50 | 0.50 | -0.036118 | 0.0 |
| 393 | 1.0 | 0.0 | 0.059264 | 1.0 | 0.0 | 0.0 | 0.75 | 0.50 | 0.50 | 1.00 | ... | 1.0 | 0.0 | 0.75 | 0.75 | 0.00 | 0.50 | 0.75 | 1.00 | -0.076118 | 1.0 |
| 394 | 1.0 | 0.0 | 0.104718 | 0.0 | 0.0 | 0.0 | 0.25 | 0.25 | 1.00 | 0.75 | ... | 1.0 | 0.0 | 0.50 | 0.25 | 0.50 | 0.50 | 0.50 | 1.00 | -0.009451 | 0.0 |

395 rows × 31 columns

```python
df.shape
```
Python

(395, 31)

```
    df.dropna().shape # their is no null value "fortunately:)"
```

```
(395, 31)
```

```
    df.columns
```

```
Index(['school', 'sex', 'age', 'address', 'famsize', 'Pstatus', 'Medu', 'Fedu',
       'Mjob', 'Fjob', 'reason', 'guardian', 'traveltime', 'studytime',
       'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery',
       'higher', 'internet', 'romantic', 'famrel', 'freetime', 'goout', 'Dalc',
       'Walc', 'health', 'absences', 'passed'],
      dtype='object')
```

```
    features=['school', 'sex', 'age', 'address', 'famsize', 'Pstatus', 'Medu', 'Fedu',
             'Mjob', 'Fjob', 'reason', 'guardian', 'traveltime', 'studytime',
             'failures', 'schoolsup', 'famsup', 'paid', 'activities', 'nursery',
             'higher', 'internet', 'romantic', 'famrel', 'freetime', 'goout', 'Dalc',
             'Walc', 'health', 'absences']
```

```
    #plot of student status
    dfv['passed'].value_counts()
```

```
yes    265
no     130
Name: passed, dtype: int64
```

```python
labels = 'student pass the final exam ', 'student fail the final exam'
sizes = [265, 130]
colors=['lightskyblue','yellow']
fig1, ax1 = plt.subplots()
ax1.pie(sizes,  labels=labels, autopct='%1.1f%%',colors=colors,
        shadow=True, startangle=90)
ax1.axis('equal')  # Equal aspect ratio ensures that pie is drawn as a circle.
plt.show()
```
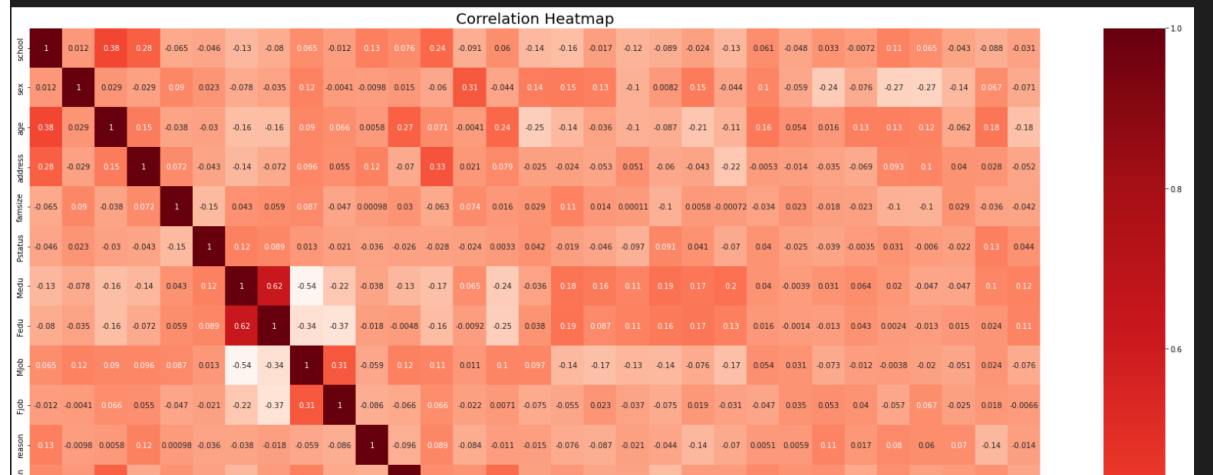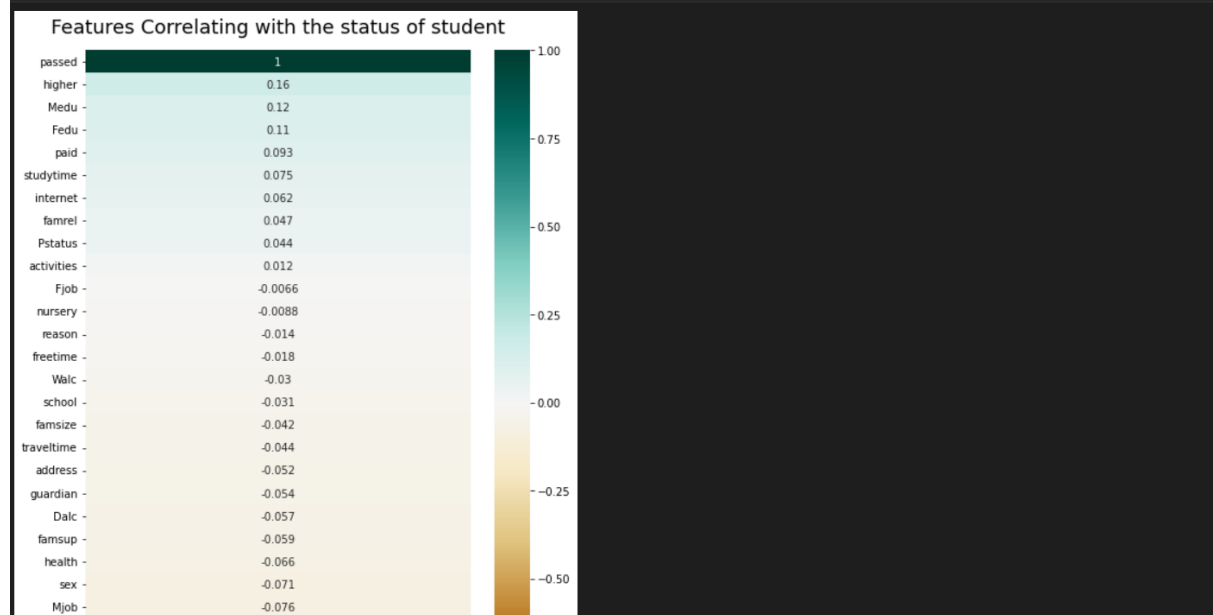
```python
# see correlation between variables through a correlation heatmap
corr = df.corr()
plt.figure(figsize=(30,30))
sns.heatmap(corr, annot=True, cmap="Reds")
plt.title('Correlation Heatmap', fontsize=20)
```
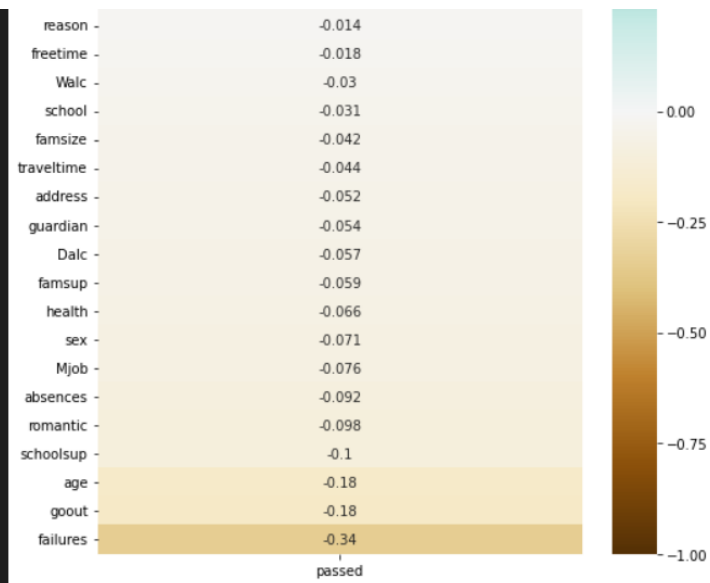
```
Text(0.5, 1.0, 'Correlation Heatmap')
```



```python
plt.figure(figsize=(8, 12))
heatmap = sns.heatmap(df.corr()[['passed']].sort_values(by='passed', ascending=False), vmin=-1, vmax=1, annot=True, cmap='BrBG')
heatmap.set_title('Features Correlating with the status of student', fontdict={'fontsize':18}, pad=16);
```
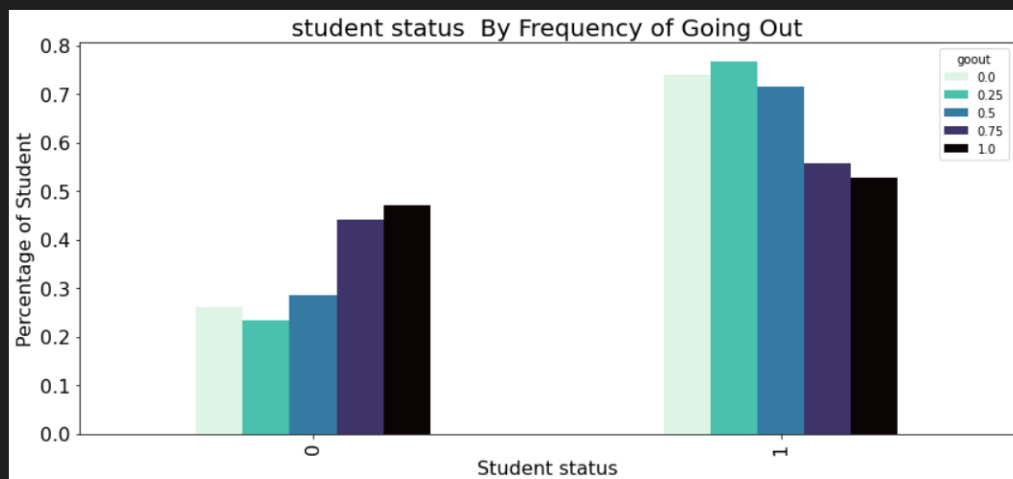
```
reason     -0.014
freetime   -0.018
Walc       -0.03
school     -0.031
famsize    -0.042
traveltime -0.044
address    -0.052
guardian   -0.054
Dalc       -0.057
famsup     -0.059
health     -0.066
sex        -0.071
Mjob       -0.076
absences   -0.092
romantic   -0.098
schoolsup  -0.1
age        -0.18
goout      -0.18
failures   -0.34
           passed
```

```python
df["goout"].unique()
```

```
array([0.75, 0.5 , 0.25, 0.  , 1.  ])
```

```python
# going out
perc = (lambda col: col/col.sum())
index = [0,1]
out_tab = pd.crosstab(index=df.passed, columns=df.goout)
out_perc = out_tab.apply(perc).reindex(index)
out_perc.plot.bar(colormap="mako_r", fontsize=16, figsize=(14,6))
plt.title('student status  By Frequency of Going Out', fontsize=20)
plt.ylabel('Percentage of Student', fontsize=16)
plt.xlabel('Student status', fontsize=16)
```
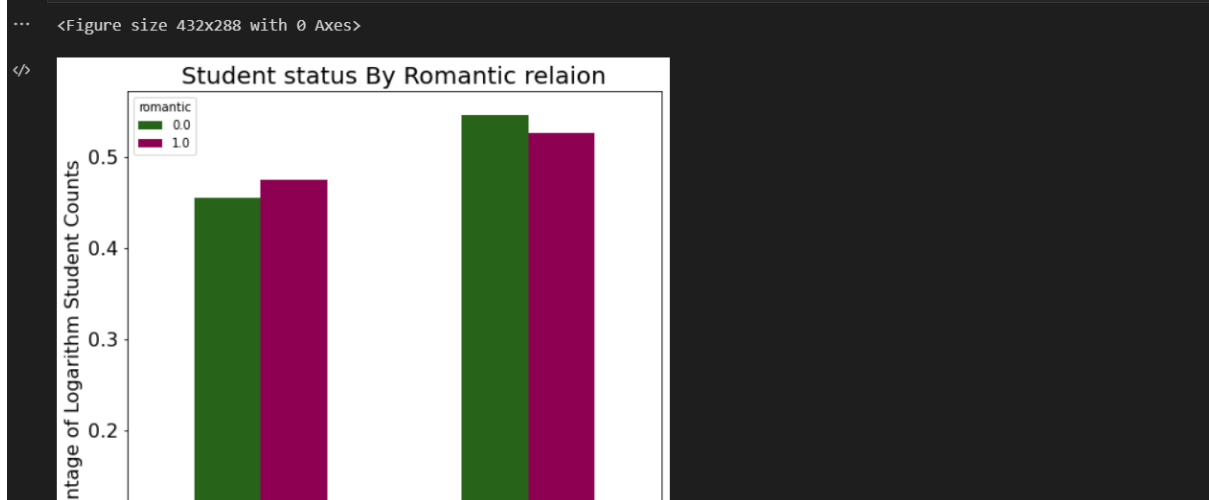
```
Text(0.5, 0, 'Student status')
```

```
# romantic status
romance_tab1 = pd.crosstab(index=df.passed, columns=df.romantic)
romance_tab = np.log(romance_tab1)
romance_perc = romance_tab.apply(perc).reindex(index)
plt.figure()
romance_perc.plot.bar(colormap="PiYG_r", fontsize=16, figsize=(8,8))
plt.title('Student status By Romantic relaion', fontsize=20)
plt.ylabel('Percentage of Logarithm Student Counts ', fontsize=16)
plt.xlabel('Student status', fontsize=16)
plt.show()
# 0 in romantic mean no romantic relation
```
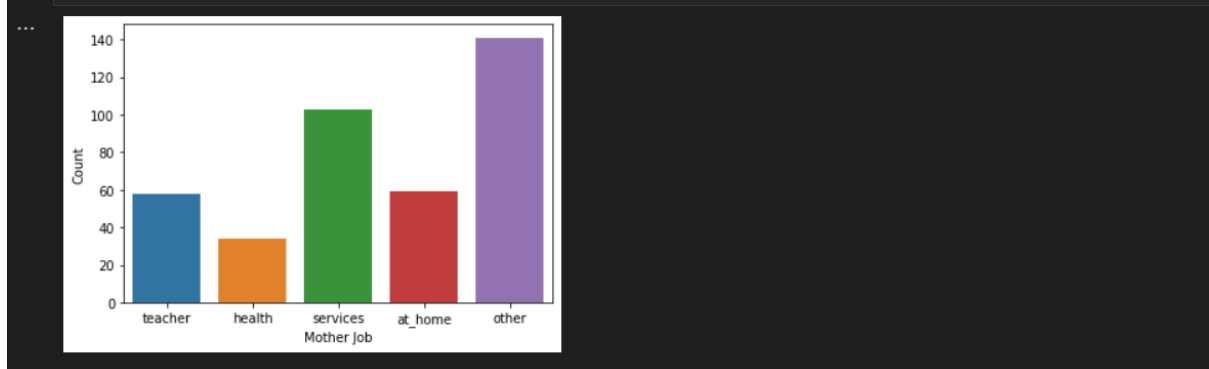
<Figure size 432x288 with 0 Axes>



```
# 1) mother job
# Mjob distribution
f, fx = plt.subplots()
figure = sns.countplot(x = 'Mjob', data=dfv, order=['teacher','health','services','at_home','other'])
fx = fx.set(ylabel="Count", xlabel="Mother Job")
figure.grid(False)
```
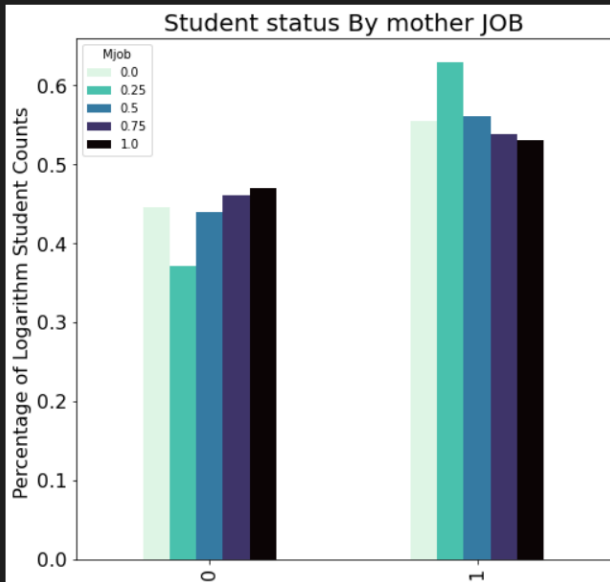
```
mjob_perc.plot.bar(colormap="mako_r", fontsize=16, figsize=(8,8))
plt.title('Student status By mother JOB', fontsize=20)
plt.ylabel('Percentage of Logarithm Student Counts ', fontsize=16)
plt.xlabel('Student status', fontsize=16)
plt.show()
#'teacher': 0, 'health': 1, 'services': 2, 'at_home': 3, 'other': 4
```

<Figure size 432x288 with 0 Axes>
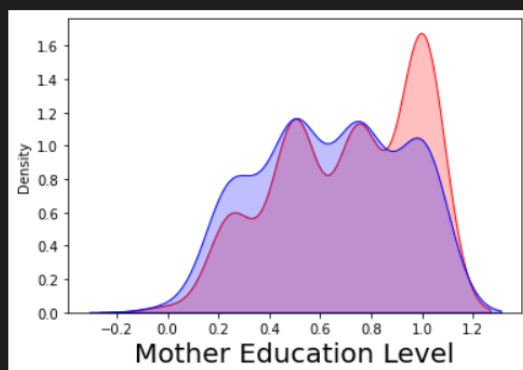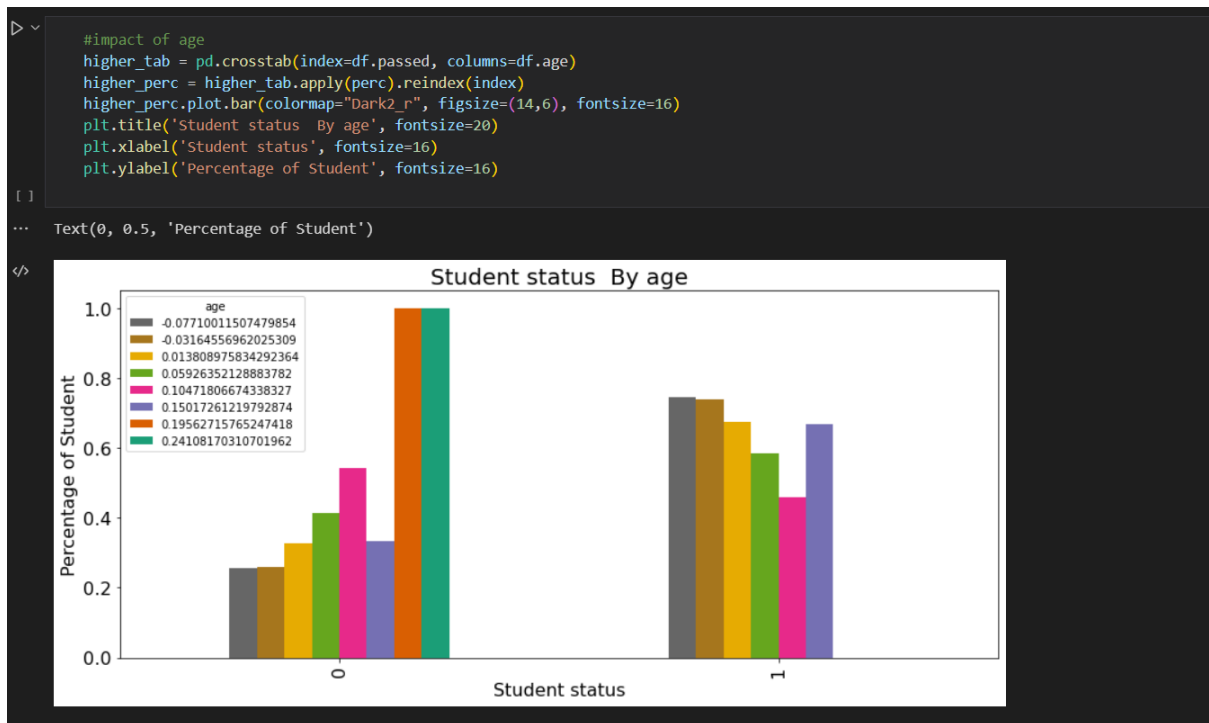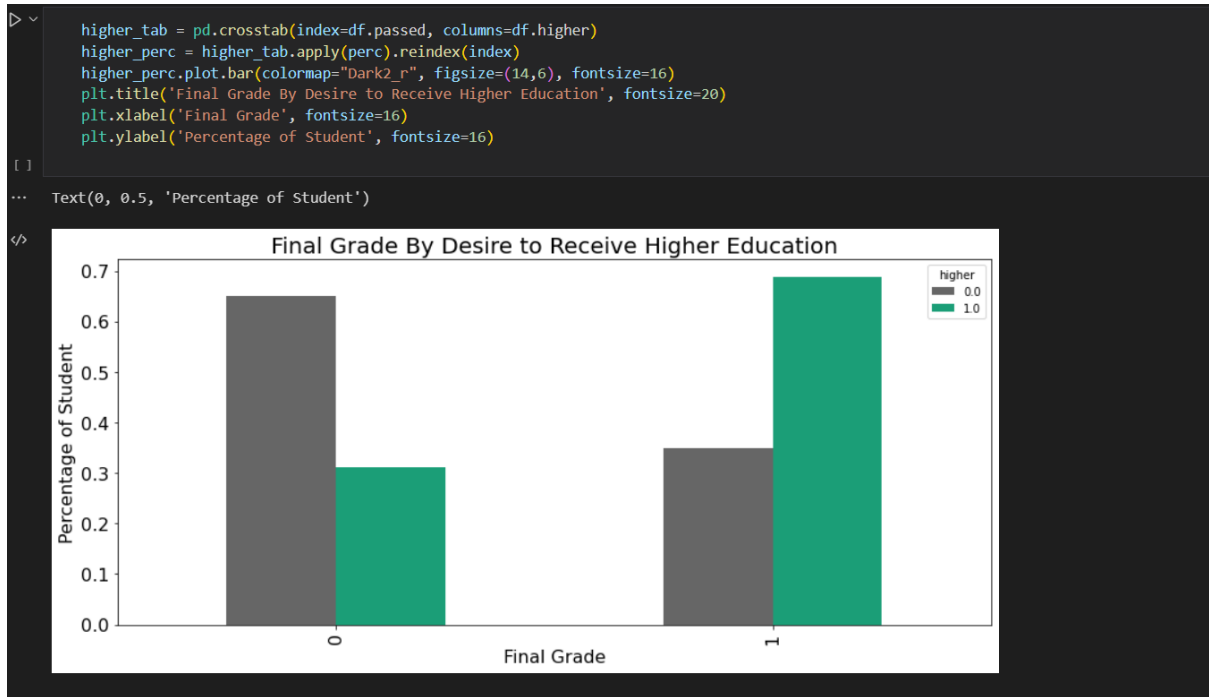


```
#Mother education:
good = df.loc[df.passed==1]
poor=df.loc[df.passed==0]
good['good_student_mother_education'] = good.Medu
poor['poor_student_mother_education'] = poor.Medu
plt.figure(figsize=(6,4))
p=sns.kdeplot(good['good_student_mother_education'], shade=True, color="r")#good_student in red
p=sns.kdeplot(poor['poor_student_mother_education'], shade=True, color="b")#poor_student in blue
plt.xlabel('Mother Education Level', fontsize=20)
```

```
higher_tab = pd.crosstab(index=df.passed, columns=df.higher)
higher_perc = higher_tab.apply(perc).reindex(index)
higher_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('Final Grade By Desire to Receive Higher Education', fontsize=20)
plt.xlabel('Final Grade', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

Text(0, 0.5, 'Percentage of Student')



```
#impact of age
higher_tab = pd.crosstab(index=df.passed, columns=df.age)
higher_perc = higher_tab.apply(perc).reindex(index)
higher_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('Student status  By age', fontsize=20)
plt.xlabel('Student status', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

Text(0, 0.5, 'Percentage of Student')

```
fail_tab = pd.crosstab(index=df.passed, columns=df.failures)
fail_perc = fail_tab.apply(perc).reindex(index)
fail_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('student status By failures', fontsize=20)
plt.xlabel('Final Grade', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

Text(0, 0.5, 'Percentage of Student')



```
#first let's see the distribution of students who live in urban or rural area
f, fx = plt.subplots()
figure = sns.countplot(x = 'address', data=dfv, order=['U','R'])
fx = fx.set(ylabel="Count", xlabel="address")
figure.grid(False)
plt.title('Address Distribution')
```

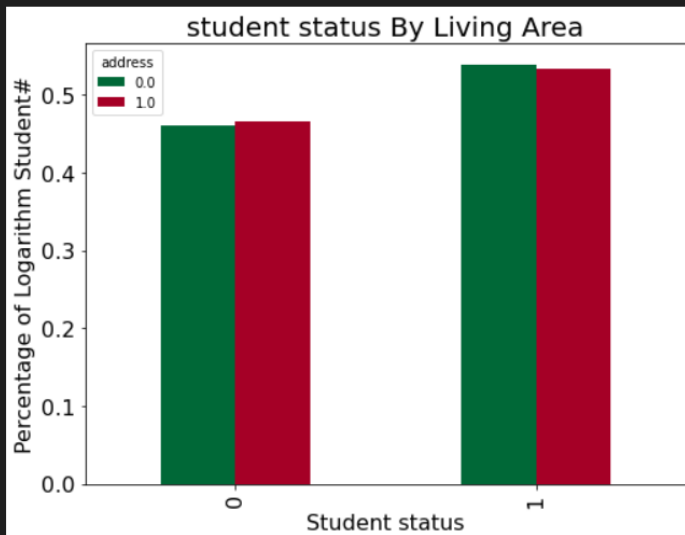Text(0.5, 1.0, 'Address Distribution')

```
ad_tab1 = pd.crosstab(index=df.passed, columns=df.address)
ad_tab = np.log(ad_tab1)
ad_perc = ad_tab.apply(perc).reindex(index)
ad_perc.plot.bar(colormap="RdYlGn_r", fontsize=16, figsize=(8,6))
plt.title('student status By Living Area', fontsize=20)
plt.ylabel('Percentage of Logarithm Student#', fontsize=16)
plt.xlabel('Student status', fontsize=16)
```
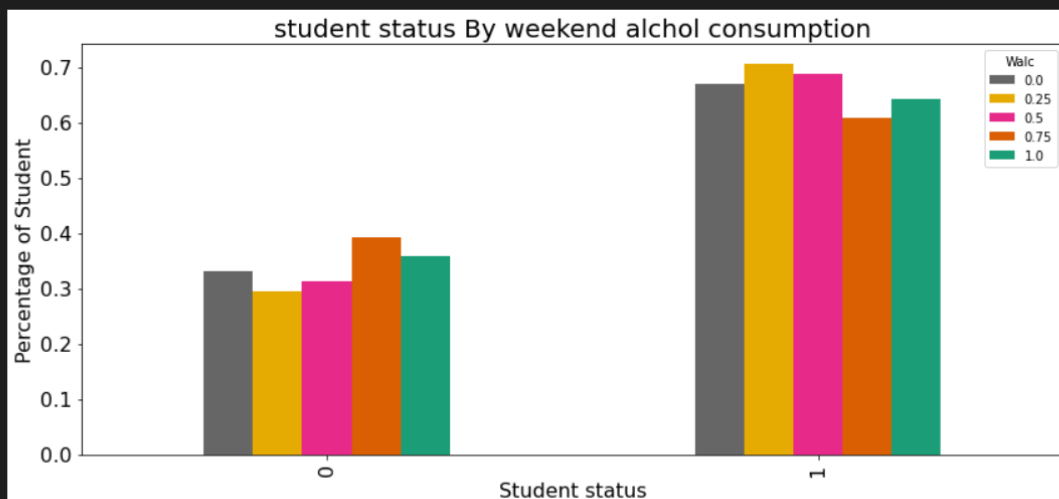
Text(0.5, 0, 'Student status')



```
#impact of weekend alcohol consumption in student performance
alc_tab = pd.crosstab(index=df.passed, columns=df.Walc)
alc_perc = alc_tab.apply(perc).reindex(index)
alc_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('student status By weekend alchol consumption', fontsize=20)
plt.xlabel('Student status', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```
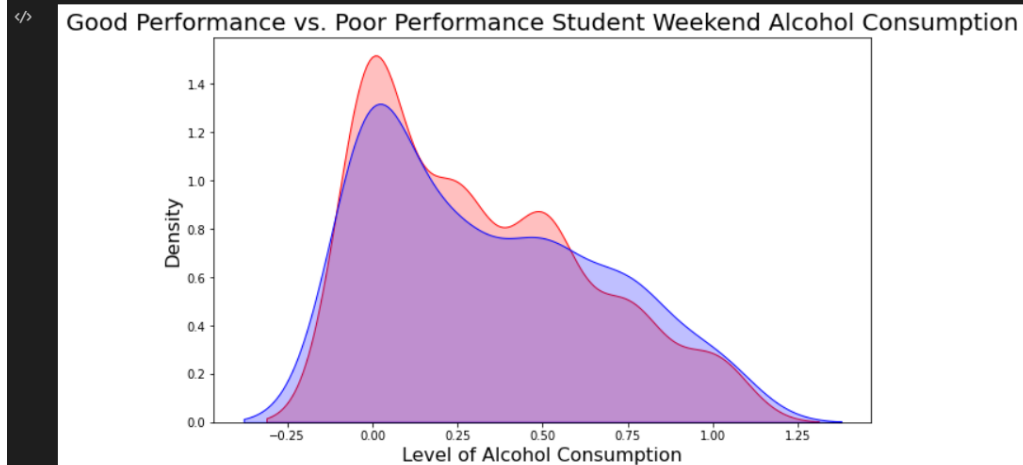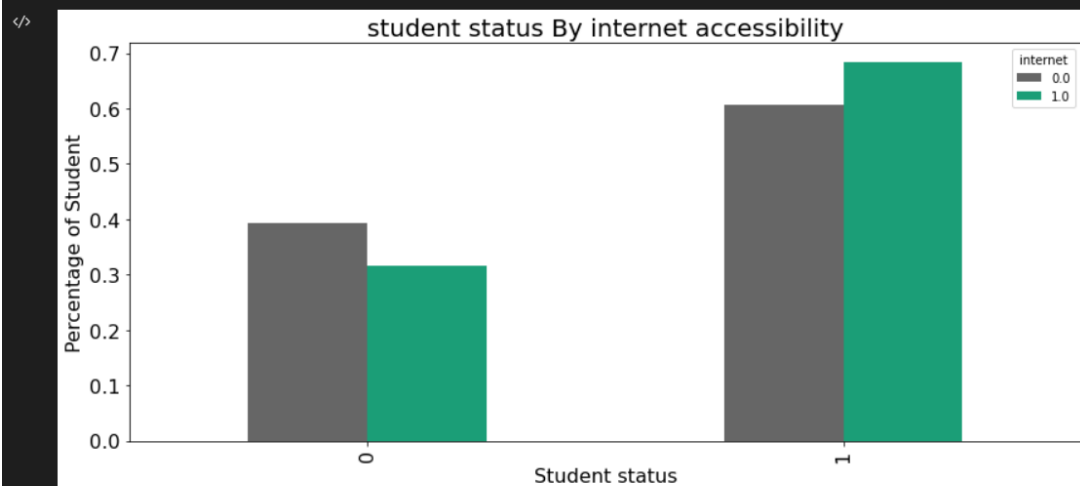
Text(0, 0.5, 'Percentage of Student')

```python
# weekend alcohol consumption
# create good student dataframe
good = df.loc[df.passed == 1]
good['good_alcohol_usage']=good.Walc
# create poor student dataframe
poor = df.loc[df.passed == 0]
poor['poor_alcohol_usage']=poor.Walc
plt.figure(figsize=(10,6))
p1=sns.kdeplot(good['good_alcohol_usage'], shade=True, color="r")
p1=sns.kdeplot(poor['poor_alcohol_usage'], shade=True, color="b")
plt.title('Good Performance vs. Poor Performance Student Weekend Alcohol Consumption', fontsize=20)
plt.ylabel('Density', fontsize=16)
plt.xlabel('Level of Alcohol Consumption', fontsize=16)
```



```python
alc_tab = pd.crosstab(index=df.passed, columns=df.internet)
alc_perc = alc_tab.apply(perc).reindex(index)
alc_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('student status By internet accessibility', fontsize=20)
plt.xlabel('Student status', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

```
Text(0, 0.5, 'Percentage of Student')
```
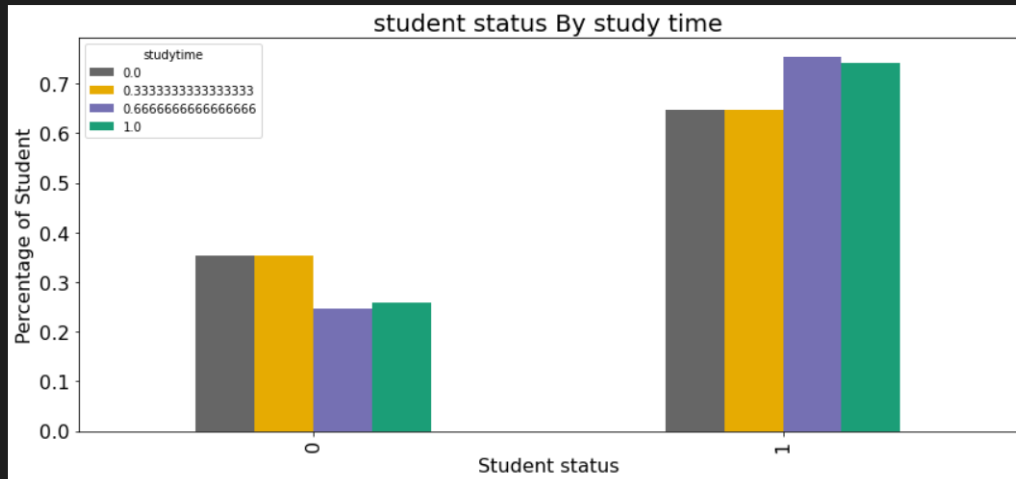
```python
stu_tab = pd.crosstab(index=df.passed, columns=df.studytime)
stu_perc = stu_tab.apply(perc).reindex(index)
stu_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('student status By study time', fontsize=20)
plt.xlabel('Student status', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

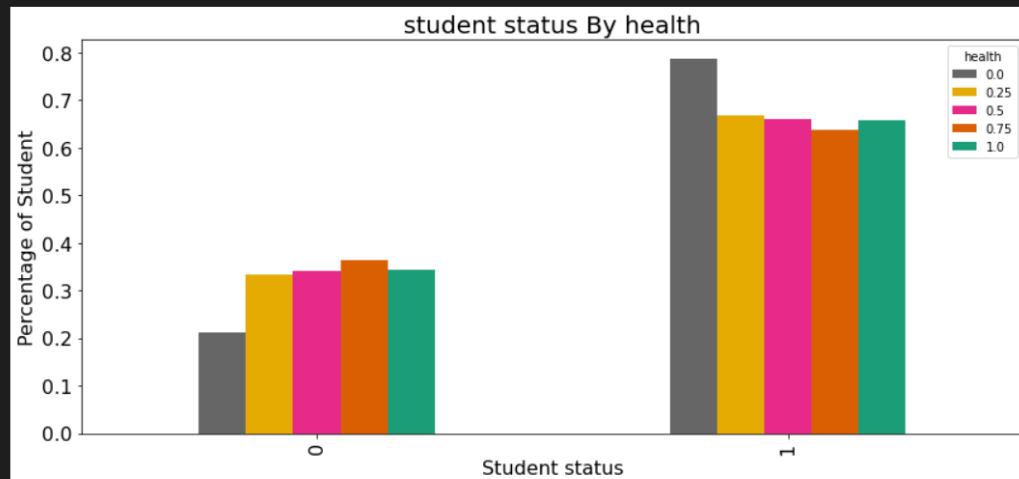Text(0, 0.5, 'Percentage of Student')



```python
he_tab = pd.crosstab(index=df.passed, columns=df.health)
he_perc = he_tab.apply(perc).reindex(index)
he_perc.plot.bar(colormap="Dark2_r", figsize=(14,6), fontsize=16)
plt.title('student status By health', fontsize=20)
plt.xlabel('Student status', fontsize=16)
plt.ylabel('Percentage of Student', fontsize=16)
```

Text(0, 0.5, 'Percentage of Student')

```python
# split data train 70 % and test 30 %

data = df.to_numpy()
n = data.shape[1]
x = data[:,0:n-1]
y = data[:,n-1]
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=0)

# Once our data is split, we can forget about x_test and y_test until we define our model.
#x_train and y_train are the samples we will use to train the model
```

```python
# let's create a model and train it

logisticRegr = LogisticRegression(C=1)
```

```python
#and now let's do the training

logisticRegr.fit(x_train,y_train)
```

```
LogisticRegression(C=1)
```

```python
y_pred=logisticRegr.predict(x_test)
y_pred
```

```
array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
       1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 0., 0., 1., 0.,
       1., 0., 1., 1., 1., 1., 1., 0., 0., 1., 0., 1., 1., 1., 0., 1.,
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1.,
       1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1.,
       1., 0., 1., 1., 1., 1., 0., 0., 1., 1., 1., 1., 0., 1., 1., 1., 1.])
```

```python
#let's have a look at the accuracy of the model

Sctest=logisticRegr.score(x_test,y_test)
Sctrain=logisticRegr.score(x_train,y_train)

print('#Accuracy test is: ',Sctest)
print('#Accuracy train is: ',Sctrain)


f1 = f1_score(y_test, y_pred, average='macro')

print('\n#f1 score is: ',f1)
```

```
#Accuracy test is:  0.6386554621848739
#Accuracy train is:  0.7463768115942029

#f1 score is:  0.5533734834598935
```

```python
#let's have a look at the accuracy of the model

Sctest=logisticRegr.score(x_test,y_test)
Sctrain=logisticRegr.score(x_train,y_train)

print('Accuracy test is: ',Sctest)
print('Accuracy train is: ',Sctrain)
```

```
Accuracy test is:  0.6386554621848739
Accuracy train is:  0.7463768115942029
```

```python
#now, we can get the confusion matrix with confusion_matrix():

confusion_matrix(y_test, y_pred)
```

```
array([[12, 38],
       [ 5, 64]])
```

```python
#let's visualize the confusion matrix:
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm,annot=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fa89987ac90>
```

```python
#import classification_report
print(classification_report(y_test, y_pred))
```

```
              precision    recall  f1-score   support

         0.0       0.71      0.24      0.36        50
         1.0       0.63      0.93      0.75        69

    accuracy                           0.64       119
   macro avg       0.67      0.58      0.55       119
weighted avg       0.66      0.64      0.58       119
```

```python
#ploting the roc_curve

fpositif, tpositif, thresholds = roc_curve(y_test, y_pred)
plt.plot([0,1],[0,1],'--')
plt.plot(fpositif,tpositif, label='LogisticRegr')
plt.xlabel('false positif')
plt.ylabel('true positif')
plt.title('LogisticRegr ROC curve')
p=plt.show()
```

```python
max_iteration = 0
maxF1 = 0
maxAccuracy = 0
optimal_state = 0
import random
for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+ ', Current f1 : '+str(maxF1), end="\r")
    split_state = np.random.randint(1,100000000)-1
    x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=split_state)
    logisticRegr = LogisticRegression(C=1)
    logisticRegr.fit(x_train,y_train)
    y_pred=logisticRegr.predict(x_test)
    f1 = f1_score(y_test, y_pred, average='macro')
    accuracy = accuracy_score(y_test, y_pred)*100

    if (accuracy>maxAccuracy and f1>maxF1):
        maxF1 = f1
        maxAccuracy = accuracy
        optimal_state = split_state


optimal_state = 85491961
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=optimal_state)
logisticRegr = LogisticRegression(C=1)
logisticRegr.fit(x_train,y_train)
y_pred=logisticRegr.predict(x_test)
f1 = f1_score(y_test, y_pred, average='macro')
accuracy = accuracy_score(y_test, y_pred)*100
print('\n\n\n*Accuracy is: '+str(accuracy)+'\n*f1 score is: ',f1)

yt_lg,yp_lg = y_test,y_pred
#ploting the roc_curve

print ( '\n\n *the ROC curve: ')
```

```python
    fpositif, tpositif, thresholds = roc_curve(y_test, y_pred)
    plt.plot([0,1],[0,1],'--')
    plt.plot(fpositif,tpositif, label='LogisticRegr')
    plt.xlabel('false positif')
    plt.ylabel('true positif')
    plt.title('LogisticRegr ROC curve')
    p=plt.show()


    #visualizig the confusion matrix:

    print (' *the confusion matrix ')

    cm = confusion_matrix(y_test, y_pred)
    sns.heatmap(cm,annot=True)
```

```
*Accuracy is: 80.67226890756302
*f1 score is:  0.7408389357068459
```
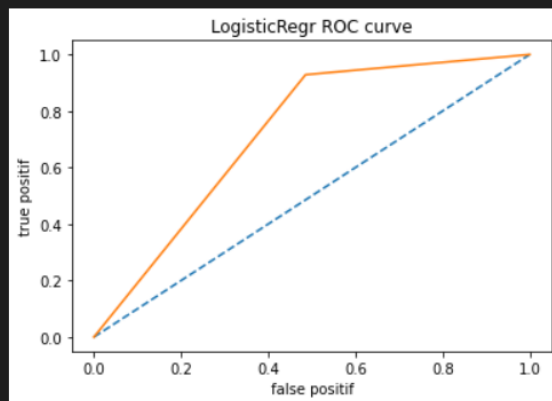
*the ROC curve:



*the confusion matrix

<matplotlib.axes._subplots.AxesSubplot at 0x7fa8997bf050>



```python
#define data
y=df.passed
target=["passed"]
x = df.drop(target,axis = 1 )
```

```python
max_iteration = 0
maxF1 = 0
maxAccuracy = 0
optimal_state = 0
for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+ ', Current f1 : '+str(maxF1), end="\
    split_state = np.random.randint(1,100000000)-1
    x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=split_state)
    KNN = KNeighborsClassifier()
    KNN.fit(x_train,y_train)
    y_pred=KNN.predict(x_test)
    f1 = f1_score(y_test, y_pred, average='macro')
    accuracy = accuracy_score(y_test, y_pred)*100

    if (accuracy>maxAccuracy and f1>maxF1):
        maxF1 = f1
        maxAccuracy = accuracy
        optimal_state = split_state

optimal_state = 71027464
```

```python
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=optimal_state)
KNN= KNeighborsClassifier()
KNN.fit(x_train,y_train)
y_pred=KNN.predict(x_test)
f1 = f1_score(y_test, y_pred, average='macro')
accuracy = accuracy_score(y_test, y_pred)*100
print('\n\n\n*Accuracy is: '+str(accuracy)+'\n*f1 score is: ',f1)

print ('random_state is ',optimal_state)


#ploting the roc_curve

print ( '\n\n *the ROC curve: ')

fpositif, tpositif, thresholds = roc_curve(y_test, y_pred)
plt.plot([0,1],[0,1],'--')
plt.plot(fpositif,tpositif, label='knn')
plt.xlabel('false positif')
plt.ylabel('true positif')
plt.title('KNN ROC curve')
p=plt.show()

yt_knn,yp_knn= y_test,y_pred
#visualizig the confusion matrix:

print (' *the confusion matrix ')

cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm,annot=True)
```

```
*Accuracy is: 78.15126050420169
*f1 score is:  0.7102996254681648
random_state is  71027464



 *the ROC curve:
```

*the confusion matrix

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fa899538f10>
```



```python
#Setup arrays to store training and test accuracies
neighbors= np.arange(1,20)
train_accuracy =np.empty(19)
test_accuracy = np.empty(19)

for i,k in enumerate(neighbors):
    #Setup a knn classifier with k neighbors
    knn = KNeighborsClassifier(n_neighbors=k)

    #Fit the model
    knn.fit(x_train, y_train)

    #Compute accuracy on the training set
    train_accuracy[i] = knn.score(x_train, y_train)

    #Compute accuracy on the test set
    test_accuracy[i] = knn.score(x_test, y_test)

#  Plotting the curv
plt.title('k-NN Varying number of neighbors')
plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
plt.plot(neighbors, train_accuracy, label='Training accuracy')
plt.legend()
plt.xlabel('Number of neighbors')
plt.ylabel('Accuracy')
plt.show()
```

k-NN Varying number of neighbors

```
#In case of classifier like knn the parameter to be tuned is n_neighbors
param_grid = {'n_neighbors':np.arange(1,20)}
knn = KNeighborsClassifier()
knn_cv= GridSearchCV(knn,param_grid,cv=5)
knn_cv.fit(x_train,y_train)
#best score\n",
knn_cv.best_score_
```

```
0.6449350649350649
```

```
knn_cv.best_params_
```

```
{'n_neighbors': 19}
```

```
param_grid = {'n_neighbors':np.arange(1,20)}
knn = KNeighborsClassifier()
knn_cv= GridSearchCV(knn,param_grid,cv=5)
knn_cv.fit(x,y)
#best score\n",
knn_cv.best_score_
```

```
0.6734177215189873
```

```
knn_cv.best_params_
```

```
{'n_neighbors': 7}
```

```python
params = {"n_neighbors":[7,19] , "metric":["euclidean", "manhattan", "chebyshev"]}
acc = {}

for m in params["metric"]:
    acc[m] = []
    for k in params["n_neighbors"]:
        print("Model_{} metric: {}, n_neighbors: {}".format(i, m, k))
        i += 1
        t = time()
        knn = KNeighborsClassifier(n_neighbors=k, metric=m)
        knn.fit(x_train,y_train)
        pred = knn.predict(x_test)
        print("Time: ", time() - t)
        acc[m].append(accuracy_score(y_test, y_pred))
        print("Acc: ", acc[m][-1])
```

```
Model_18 metric: euclidean, n_neighbors: 7
Time:   0.012510061264038086
Acc:   0.7815126050420168
Model_19 metric: euclidean, n_neighbors: 19
Time:   0.011599302291870117
Acc:   0.7815126050420168
Model_20 metric: manhattan, n_neighbors: 7
Time:   0.012638568878173828
Acc:   0.7815126050420168
Model_21 metric: manhattan, n_neighbors: 19
Time:   0.012667417526245117
Acc:   0.7815126050420168
Model_22 metric: chebyshev, n_neighbors: 7
Time:   0.011996269226074219
Acc:   0.7815126050420168
```

```python
max_iteration = 0
maxF1 = 0
maxAccuracy = 0
optimal_state = 0
f1 = 0
accuracy = 0
True60 = False
for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+ ', Current f1 : '+str(maxF1), end="\r")
    split_state = np.random.randint(1,100000000)-1
    x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=split_state)
    KNN = KNeighborsClassifier(n_neighbors=7,metric='chebyshev')
    KNN.fit(x_train,y_train)
    y_pred=KNN.predict(x_test)
    f1 = f1_score(y_test, y_pred, average='macro')
    accuracy = accuracy_score(y_test, y_pred)*100

    if accuracy>maxAccuracy and f1>=0.5:
        maxF1 = f1
        maxAccuracy = accuracy
        optimal_state = split_state
        if maxAccuracy>79:
            break

optimal_state = 29300362
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=optimal_state)
KNN_f= KNeighborsClassifier(n_neighbors=7,metric='chebyshev')
KNN_f.fit(x_train,y_train)
y_pred=KNN_f.predict(x_test)
f1 = f1_score(y_test, y_pred, average='macro')
accuracy = accuracy_score(y_test, y_pred)*100
print('\n\n\n*Accuracy is: '+str(accuracy)+'\n*f1 score is: ',f1)
```

```
    yt_knn,yp_knn= y_test,y_pred
```

```
*Accuracy is: 69.74789915966386
*f1 score is:  0.47959183673469385
random_state is  29300362
```

```
ac = accuracy_score(yt_knn,yp_knn)
print('Accuracy is: ',ac)
cm= confusion_matrix(yt_knn,yp_knn)
sns.heatmap(cm,annot=True)
yt_knn,yp_knn = y_test,y_pred
```

Accuracy is:  0.6974789915966386



```
print(classification_report(y_test,y_pred))
```

```
              precision    recall  f1-score   support

         0.0       0.19      0.12      0.14        26
         1.0       0.78      0.86      0.82        93

    accuracy                           0.70       119
   macro avg       0.48      0.49      0.48       119
weighted avg       0.65      0.70      0.67       119
```

```
#ploting the roc_curve

print ( ' the ROC curve: ')

fpositif, tpositif, thresholds = roc_curve(y_test, y_pred)
plt.plot([0,1],[0,1],'--')
plt.plot(fpositif,tpositif, label='final knn model')
plt.xlabel('false positif')
plt.ylabel('true positif')
plt.title('knn_f ROC curve')
p=plt.show()
```

the ROC curve:



knn_f ROC curve

```python
def showResults(accuracy, trainingTime, y_pred,model):

    print('---------------------------------------------Results :',model,'-----------------------------------------
    confusionMatrix = confusion_matrix(y_test, y_pred)
    print('\n The ROC curve is :\n')
    fig, _ = plt.subplots()
    fpr,tpr,thresholds=roc_curve(y_test,y_pred)
    plt.plot([0, 1],[0, 1],'--')
    plt.plot(fpr,tpr,label=model)
    plt.xlabel('false positive')
    plt.ylabel('false negative')
    plt.legend()
    fig.suptitle('ROC curve: '+str(model))
    plt.show()

    print('---------------------------------------------')
    print('The model  accuracy:', round(accuracy),'%')
    print('---------------------------------------------')
    print('The training time is: ',trainingTime)
    print('---------------------------------------------')
    print('The f1 score is :',round(100*f1_score(y_test, y_pred, average='macro'))/100)
    print('---------------------------------------------')
    print('The roc_auc_score is :',round(100*roc_auc_score(y_test, y_pred))/100)
    print('---------------------------------------------')
    print('The confusion matrix is :\n')
    ax = plt.axes()
    sns.heatmap(confusionMatrix,annot=True)
```

```python
# Optimal C
def optimal_C_value():
    Ci = np.array(( 0.0001,0.001,0.01,0.05,0.1,4,10,40,100))
    minError = float('Inf')
    optimal_C = float('Inf')

    for c in Ci:
        clf = SVC(C=c,kernel='linear')
        clf.fit(X_train, y_train)
        predictions = clf.predict(X_val)
        error = np.mean(np.double(predictions != y_val))
        if error < minError:
            minError = error
            optimal_C = c
    return optimal_C

# Optimal C and the degree of the polynomial
def optimal_C_d_values():
    Ci = np.array(( 0.0001,0.001,0.01,0.05,0.1,4,10,40,100))
    Di = np.array(( 2, 5, 10, 15, 20, 25, 30))
    minError = float('Inf')
    optimal_C = float('Inf')
    optimal_d = float('Inf')

    for d in Di:
        for c in Ci:
            clf = SVC(C=c,kernel='poly', degree=d)
            clf.fit(X_train, y_train)
            predictions = clf.predict(X_val)
            error = np.mean(np.double(predictions != y_val))
            if error < minError:
                minError = error
                optimal_C = c
                optimal_d = d
    return optimal_C,optimal_d
```

```python
# Optimal C and gamma
def optimal_C_gamma_values():
    Ci = np.array(( 0.0001,0.001,0.01,0.05,0.1,4,10,40,100))
    Gi = np.array(( 0.000001,0.00001,0.01,1,2,3,5,20,70,100,500,1000))
    minError = float('Inf')
    optimal_C = float('Inf')
    optimal_g = float('Inf')

    for g in Gi:
        for c in Ci:
            clf = SVC(C=c,kernel='rbf', gamma=g)
            clf.fit(X_train, y_train)
            predictions = clf.predict(X_val)
            error = np.mean(np.double(predictions != y_val))
            if error < minError:
                minError = error
                optimal_C = c
                optimal_g = g
    return optimal_C,optimal_g
# ----------------------------------------------------------------------------------
# Compare the three kernels


def compare_kernels():
    X_train1,X_val1,X_test1,y_train1,y_val1,y_test1 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state1)
    X_train2,X_val2,X_test2,y_train2,y_val2,y_test2 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state2)
    X_train3,X_val3,X_test3,y_train3,y_val3,y_test3 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state3)
    print('--------------------------------------------- Comparison ---------------------------------------------')
    print('\n')
    f11 = "{:.2f}".format(f1_score(y_test1, y_linear, average='macro'))
    f22 = "{:.2f}".format(f1_score(y_test2, y_poly, average='macro'))
    f33 = "{:.2f}".format(f1_score(y_test3, y_gauss, average='macro'))
    roc1 = "{:.2f}".format(roc_auc_score(y_test1, y_linear))
    roc2 = "{:.2f}".format(roc_auc_score(y_test2, y_poly))
    roc3 = "{:.2f}".format(roc_auc_score(y_test3, y_gauss))
```

```python
        a1,a2 = confusion_matrix(y_test1, y_linear)[0],confusion_matrix(y_test1, y_linear)[1]
        b1,b2 = confusion_matrix(y_test2, y_poly)[0],confusion_matrix(y_test2, y_poly)[1]
        c1,c2 = confusion_matrix(y_test3, y_gauss)[0],confusion_matrix(y_test3, y_gauss)[1]
        data_rows = [('training time',time1, time2, time3),
                      ('','','',''),
                       ('accuracy %',linear_accuracy, poly_accuracy, gauss_accuracy),
                      ('','','',''),
                      ('confusion matrix',a1, b1, c1),
                      ('',a2,b2,c2),
                      ('','','',''),
                      ('f1 score',f11,f22,f33),
                      ('','','',''),
                      ('roc_auc_score',roc1,roc2,roc3)]
        t = Table(rows=data_rows, names=('metric','Linear kernel', 'polynomial kernel', 'gaussian kernel'))
        print(t)
        print('\n\n')
        print('The Roc curves :\n')
        y_pred1 = y_linear
        y_pred2 = y_poly
        y_pred3 = y_gauss
        fig, _ = plt.subplots()
        fig.suptitle('Comparison of three ROC curves')
        fpr,tpr,thresholds=roc_curve(y_test1,y_pred1)
        plt.plot([0, 1],[0, 1],'--')
        plt.plot(fpr,tpr,label='Linear kernel :'+str(roc1))
        plt.xlabel('false positive')
        plt.ylabel('false negative')
        fpr,tpr,thresholds=roc_curve(y_test2,y_pred2)
        plt.plot(fpr,tpr,label='Polynomial kernel :'+str(roc2))
        fpr,tpr,thresholds=roc_curve(y_test3,y_pred3)
        plt.plot(fpr,tpr,label='Gaussian kernel :'+str(roc3))
        plt.legend()
        plt.show()
```

```python
def best_kernel(kernel):
    X_train1,X_val1,X_test1,y_train1,y_val1,y_test1 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state1)
    X_train2,X_val2,X_test2,y_train2,y_val2,y_test2 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state2)
    X_train3,X_val3,X_test3,y_train3,y_val3,y_test3 = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state3)

    time = 0
    f1 = 0
    accuracy = 0
    rc = 0
    y = 0
    if kernel == 'linear kernel':
        time = time1
        f1 = "{:.2f}".format(f1_score(y_test1, y_linear, average='macro'))
        accuracy = round(100*linear_accuracy)/100
        rc = round(100*roc_auc_score(y_test1, y_linear))/100
        y_test = y_test1
        y = y_linear
    elif kernel == 'polynomial kernel':
        time = time2
        f1 = "{:.2f}".format(f1_score(y_test2, y_poly, average='macro'))
        accuracy = round(100*poly_accuracy)/100
        rc = round(100*roc_auc_score(y_test2, y_poly))/100
        y_test = y_test2
        y = y_poly
    else :
        time = time3
        f1 = "{:.2f}".format(f1_score(y_test3, y_gauss, average='macro'))
        accuracy = round(100*gauss_accuracy)/100
        rc = round(100*roc_auc_score(y_test3, y_gauss))/100
        y_test = y_test3
        y = y_gauss

    # used for comparing three classfiers(knn, logistic regression and svm)
    yt_svm,yp_svm = y_test, y
```

```python
    print('The choosen kernel :',kernel)
    print('the training :',time)
    print('the accuracy :',round(accuracy),'%')
    print('the f1 score :',f1)
    print('The roc_auc_score is :',rc)
    print('--------------------------------------\nThe ROC curve :')
    fig, _ = plt.subplots()
    fpr,tpr,thresholds=roc_curve(y_test,y)
    plt.plot([0, 1],[0, 1],'--')
    plt.plot(fpr,tpr,label=kernel+': '+str(rc))
    plt.xlabel('false positive')
    plt.ylabel('false negative')
    plt.legend()
    plt.show()
    confusionMatrix = confusion_matrix(y_test, y)
    print('--------------------------------------\nThe confusion matrix is  :')
    ax = plt.axes()
    sns.heatmap(confusionMatrix,annot=True)
    ax.set_title('Confusion matrix of SVM '+str(kernel))
    return yt_svm,yp_svm


# ----------------------------------------------------------------------------------------------------
# svm factor : factor affecting students performance, later on on this Ipython notebook  we will explain how we will do this


# 1) factor as svm coefficients
def factors(array, K, max_or_min, df):

    n = array.shape[1]
    array = array.reshape(n,1)
    my_list = array.tolist()
    if max_or_min == 'max':
        temp = sorted(my_list)[-K:]
        res = []
        for ele in temp:
```

```python
    n = array.shape[1]
    array = array.reshape(n,1)
    my_list = array.tolist()
    if max_or_min == 'max':
        temp = sorted(my_list)[-K:]
        res = []
        for ele in temp:
            res.append(my_list.index(ele))
        return(get_factors(res, df))


    elif max_or_min == 'min':
        temp = sorted(my_list, reverse=True)[-K:]
        temp = temp = np.array(temp).reshape(K,1)
        res = []
        for ele in temp:
            if ele<0:
                res.append(my_list.index(ele))
        return(get_factors(res, df))


    else:
        return


# 2) converts those factors to dataset columns name
def get_factors(index, df):
    f = []
    for i in index:
        f.append(df.columns[i])
    return f
```

```python
# 3) Convert column names to understandable string

columns_name = {'famsize': 'family size', 'Pstatus': "parent's cohabitation status ", 'Medu': "mother's education",
                'Fedu': "father's education", 'Mjob': "mother's job", 'Fjob': "father's job",
                'reason': 'reason to choose this school ','schoolsup': 'extra educational support', 'famsup': 'family educational support',
                'paid': 'extra paid classes within the course subject', 'higher': 'wants to take higher education',
                'romantic': 'with a romantic relationship ', 'famrel': 'quality of family relationships', 'goout': 'going out with friends',
                'Dalc': 'workday alcohol consumption', 'Walc': 'weekend alcohol consumption'}


def column_to_string(fcts,max_or_min):

    if max_or_min == 'max':
        print('-----------------------------------------------------------------------------')
        print('Factors helping students succeed :')
    else:
        print('-----------------------------------------------------------------------------')
        print('-----------------------------------------------------------------------------')
        print('Factors leading students to failure')

    for fct in fcts:
        if fct in columns_name:
            print(columns_name[fct])
        else:
            print(fct)
```

```python
    def split(df,rest_size,test_size,randomState):
        data = df.to_numpy()
        n = data.shape[1]
        x = data[:,0:n-1]
        y = data[:,n-1]
        if(randomState):
            X_train,X_rest,y_train,y_rest = train_test_split(x,y,test_size=rest_size,random_state=randomState)
            X_val,X_test,y_val,y_test = train_test_split(X_rest,y_rest,test_size=test_size,random_state=randomState)
        else:
            X_train,X_rest,y_train,y_rest = train_test_split(x,y,test_size=rest_size,random_state=0)
            X_val,X_test,y_val,y_test = train_test_split(X_rest,y_rest,test_size=test_size,random_state=0)

        return X_train,X_val,X_test,y_train,y_val,y_test
    # We will use the three different svm classifier kernels
    # Linear kernel, polynomial kernel and gaussian kernel and we will choose the most accurate
```

```python
########################################################## Linear kernel ##########################################################
optimal_split_state1 = 0
maxAccuracy = 0
maxF1 = 0

# We already tune parameters, we do not need to loop over all the hyperparamters again,
# if you want to do so just set max_iteration to 2000 for example
# and remove the line 'optimal_split_state = 388628375' at the bottom of this cell.

max_iteration = 0
if max_iteration != 0:
    print ('-----------------------------------------Hyperparameters tunning starts-----------------------------------------\n\n')

for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+' Current f1 '+str(maxF1), end="\r")
    # Let's get the optimal C value for the linear kernal
    split_state = np.random.randint(1,1000000000)-1
    X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=split_state)
    optimal_C = optimal_C_value()


    # Now let's use the optimal C value
    linear_clf = SVC(C=optimal_C,kernel='linear')

    # Let's train the model with the optimal C value and calculate the training time
    tic = time()
    linear_clf.fit(X_train, y_train)
    toc = time()
    time1 = str(round(1000*(toc-tic))) + "ms"
    y_linear = linear_clf.predict(X_test)
    linear_f1 = f1_score(y_test, y_linear, average='macro')
    linear_accuracy = accuracy_score(y_test, y_linear)*100
```

```python
        if linear_accuracy>maxAccuracy and linear_f1>maxF1:
            maxAccuracy = linear_accuracy
            maxF1 = linear_f1
            optimal_split_state1 = split_state
        if maxAccuracy>86 and maxF1>80:
            break;
# We've already tuned our hyperparameters, we will not repeat that again as it takes soo long.
# The optimal split state for linear kernel is 388628375
# Let's try that split state
optimal_split_state1 = 388628375
X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state1)
optimal_C = optimal_C_value()


# Now let's use the optimal C value
linear_clf = SVC(C=optimal_C,kernel='linear')

# Let's train the model with the optimal C value and calculate the training time
tic = time()
linear_clf.fit(X_train, y_train)
toc = time()
time1 = str(round(1000*(toc-tic))) + "ms"
y_linear = linear_clf.predict(X_test)
linear_accuracy = accuracy_score(y_test, y_linear)*100
if max_iteration != 0:
    print('\n\n\n                              ---------------------------process ended'\
          '-----------------------------                      \n\n\n')

# Let's show the resuls
showResults(linear_accuracy, time1, y_linear,'SVM linear kernel')
```

The ROC curve is :



ROC curve: SVM linear kernel

----------------------------------------------
The model   accuracy: 84 %
----------------------------------------------
The training time is:  10ms
----------------------------------------------
The f1 score is : 0.82
----------------------------------------------
The roc_auc_score is : 0.8
----------------------------------------------

The confusion matrix is :



```python
######################################################## Polynomial kernel ########################################################
optimal_split_state2 = 0
maxAccuracy = 0
maxF1 = 0


# We already tune parameters, we do not need to loop over all the hyperparamters again,
# if you want to do so just set max_iteration to 500 for example
# and remove the line 'optimal_split_state2 = 7070621' at the bottom of this cell.

max_iteration = 0
if max_iteration != 0:
    print ('----------------------------------------Hyperparameters tunning starts----------------------------------------\n\n')
for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+', Current f1 '+str(maxF1), end="\r")

    split_state = np.random.randint(1,100000000)-1
    X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=split_state)

    # Let's get the optimal C and the degree value for the polynomial kernal
    optimal_C, optimal_d = optimal_C_d_values()

    # Now let's use the optimal c value and the optimal degree value
    poly_clf = SVC(C=optimal_C,kernel='poly', degree=optimal_d)

    # Let's train the model with the optimal C value
    poly_clf.fit(X_train, y_train)
    y_poly = poly_clf.predict(X_test)
    poly_f1 = f1_score(y_test, y_poly, average='macro')
    poly_accuracy = accuracy_score(y_test, y_poly)*100
```

```python
    if poly_accuracy>maxAccuracy and poly_f1>maxF1:
        maxAccuracy = poly_accuracy
        maxF1 = poly_f1
        optimal_split_state2 = split_state
# We've already tuned our hyperparameters, we will not repeat that again as it takes soo long.
# The optimal split state for polynomial kernel is 7070621
# Let's try that split state
optimal_split_state2 = 7070621

X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state2)

optimal_C, optimal_d = optimal_C_d_values()


# Now let's use the optimal C value
poly_clf = SVC(C=optimal_C,kernel='poly', degree=optimal_d)

# Let's train the model and calculate the training time
tic = time()
poly_clf.fit(X_train, y_train)
toc = time()
time2 = str(round(1000*(toc-tic))) + "ms"
y_poly = poly_clf.predict(X_test)
poly_accuracy = accuracy_score(y_test, y_poly)*100
if max_iteration != 0:
    print('\n\n\n                             ---------------------------process ended'\
        '----------------------------------                     \n\n\n')

# Let's show the resuls
showResults(poly_accuracy, time2, y_poly,'SVM polynomial kernel')
```
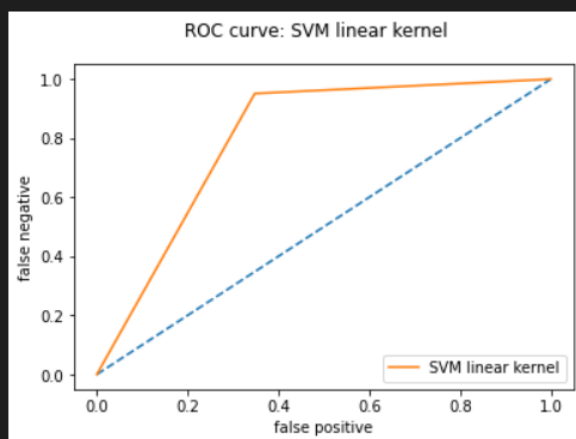
```
...    -----------------------------------------Results : SVM polynomial kernel -----------------------------------------

       The ROC curve is :
```



```
       -----------------------------------------
       The model   accuracy: 78 %
       -----------------------------------------
       The training time is:  8ms
       -----------------------------------------
       The f1 score is : 0.74
       -----------------------------------------
       The roc_auc_score is : 0.73
       -----------------------------------------
```

```
       The confusion matrix is :
```



```python
########################################################## Gaussian kernel ##########################################################
optimal_split_state3 = 0
maxAccuracy = 0
maxF1 = 0


# We already tune parameters, we do not need to loop over all the hyperparamters again,
# if you want to do so just set max_iteration to 500 for example
# and remove the line 'optimal_split_state3 = 93895097' at the bottom of this cell.

max_iteration = 0
if max_iteration != 0:
    print ('-----------------------------------------------Hyperparameters tunning starts'\
           '-----------------------------------------------\n\n')
for k in range(max_iteration):
    print ('Iteration :'+str(k)+', Current accuracy: '+str(maxAccuracy)+', Current f1 '+str(maxF1), end="\r")

    split_state = np.random.randint(1,100000000)-1
    X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=split_state)

    # Let's get the optimal C and the degree value for the polynomial kernal
    optimal_C, optimal_gamma = optimal_C_gamma_values()

    # Now let's use the optimal c value and the optimal degree value
    gauss_clf = SVC(C=optimal_C,kernel='rbf',gamma=optimal_gamma)

    # Let's train the model with the optimal C value
    gauss_clf.fit(X_train, y_train)
    y_gauss = gauss_clf.predict(X_test)
    gauss_f1 = f1_score(y_test, y_gauss, average='macro')
    gauss_accuracy = accuracy_score(y_test, y_gauss)*100
```

```python
        if gauss_accuracy>maxAccuracy and gauss_f1>maxF1:
            maxAccuracy = gauss_accuracy
            maxF1 = gauss_f1
            optimal_split_state3 = split_state

    # We've already tuned our hyperparameters, we will not repeat that again as it takes soo long.
    # The optimal split state for polynomial kernel is 93895097
    # Let's try that split state
    optimal_split_state3 = 93895097
    X_train,X_val,X_test,y_train,y_val,y_test = split(df,rest_size=0.4,test_size=0.4,randomState=optimal_split_state3)

    optimal_C, optimal_gamma = optimal_C_gamma_values()


    # Now let's use the optimal C value
    gauss_clf = SVC(C=optimal_C,kernel='rbf',gamma=optimal_gamma)

    # Let's train the model and calculate the training time
    tic = time()
    gauss_clf.fit(X_train, y_train)
    toc = time()
    time3 = str(round(1000*(toc-tic))) + "ms"
    y_gauss = gauss_clf.predict(X_test)
    gauss_accuracy = (accuracy_score(y_test, y_gauss)*100)

    if max_iteration != 0:
        print('\n\n\n                             ---------------------------process ended'\
              '-----------------------------                  \n\n\n')

    # Let's show the resuls
    showResults(gauss_accuracy, time3, y_gauss,'SVM gaussian kernel')
```
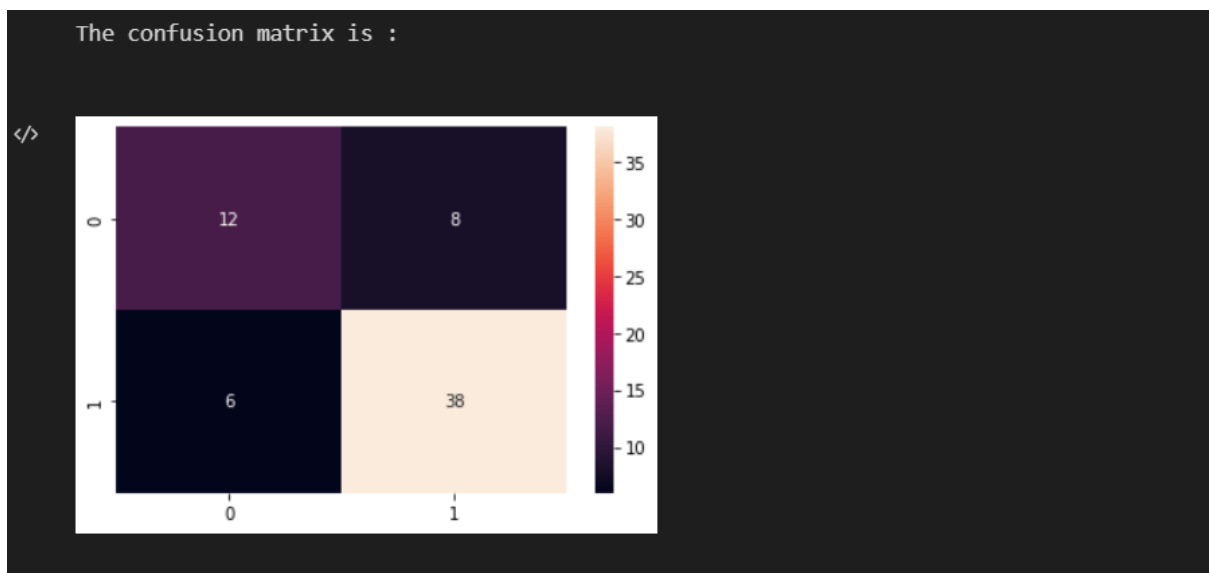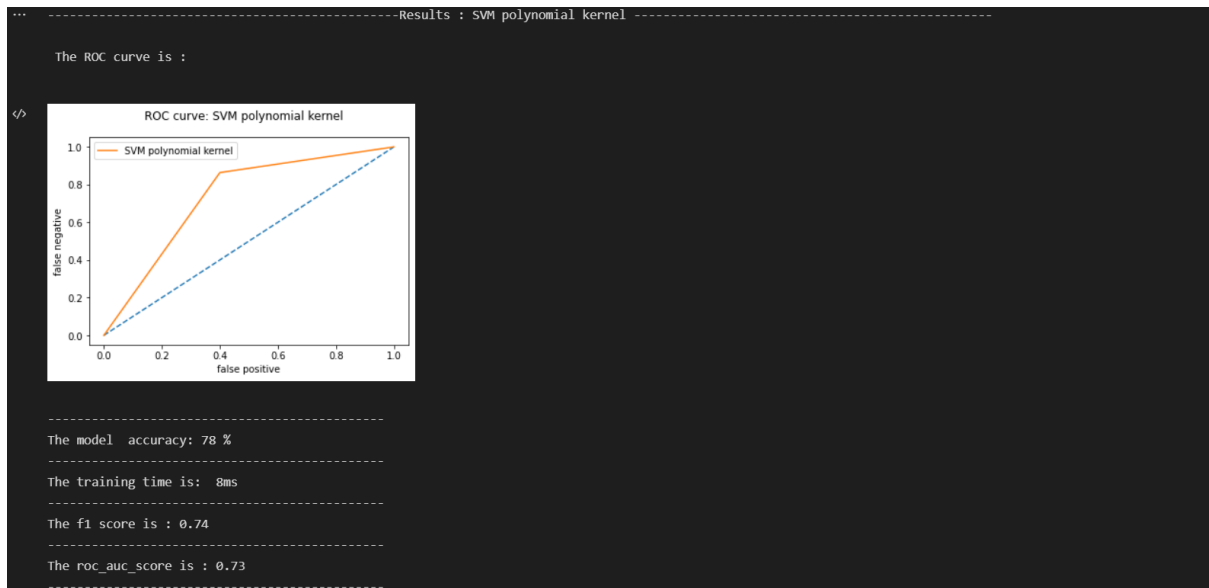
```
--------------------------------------------------Results : SVM gaussian kernel --------------------------------------------------

    The ROC curve is :
```



```
------------------------------------------------
The model  accuracy: 83 %
------------------------------------------------
The training time is:  5ms
------------------------------------------------
The f1 score is : 0.77
------------------------------------------------
```

The confusion matrix is :

```
compare_kernels()
```

```
[ ]

...    ---------------------------------------------- Comparison ----------------------------------------


         metric      Linear kernel polynomial kernel gaussian kernel
    ---------------- ------------- ----------------- ---------------
      training time       10ms                8ms              5ms

        accuracy %       84.375            78.125          82.8125

   confusion matrix      [15  8]            [12  8]          [10 11]
                         [ 2 39]            [ 6 38]          [ 0 43]

          f1 score         0.82               0.74             0.77

     roc_auc_score         0.80               0.73             0.74
```

The Roc curves :



Comparison of three ROC curves

```
yt_svm,yp_svm = best_kernel('linear kernel')
```

```
The choosen kernel : linear kernel
the training : 10ms
the accuracy : 84 %
the f1 score : 0.82
The roc_auc_score is : 0.8
----------------------------------------
The ROC curve :
```



```
----------------------------------------
The confusion matrix is  :
```

```python
    # Get svm parameters
    coefs = linear_clf.coef_

    # factors helping students to succeed
    column_to_string(factors(coefs, 5, 'max', df),'max')

    # factors leading students to failure
    column_to_string(factors(coefs, 5, 'min', df), 'min')
```

```
...    --------------------------------------------------------------------------------
       Factors helping students succeed :
       father's education
       guardian
       wants to take higher education
       studytime
       father's job
       --------------------------------------------------------------------------------
       --------------------------------------------------------------------------------
       Factors leading students to failure
       age
       health
       going out with friends
       absences
       failures
```

```python
# Function to compare the three classifiers (Logistic regression, KNN and SVM) performances :

def compare_lg_knn_svm(yt_knn,yp_knn,yt_lg,yp_lg,yt_svm,yp_svm):
    #F1 score
    f1_lg = round(f1_score(yt_lg, yp_lg, average='macro')*100)
    f1_knn = round(f1_score(yt_knn, yp_knn, average='macro')*100)
    f1_svm = round(f1_score(yt_svm, yp_svm, average='macro')*100)

    #Accuracy score
    acc_lg = round(accuracy_score(yt_lg, yp_lg)*100)
    acc_knn = round(accuracy_score(yt_knn, yp_knn)*100)
    acc_svm = round(accuracy_score(yt_svm, yp_svm)*100)

    #Confusion matrix
    conf_lg = confusion_matrix(yt_lg, yp_lg)
    conf_knn = confusion_matrix(yt_knn, yp_knn)
    conf_svm = confusion_matrix(yt_svm, yp_svm)

    #ROC score
    roc_c_lg = round(roc_auc_score(yt_lg, yp_lg)*100)
    roc_c_knn = round(roc_auc_score(yt_knn, yp_knn)*100)
    roc_c_svm = round(roc_auc_score(yt_svm, yp_svm)*100)

    #ROC curve thresholds
    roc_knn = roc_curve(yt_knn,yp_knn)
    roc_lg = roc_curve(yt_lg,yp_lg)
    roc_svm = roc_curve(yt_svm,yp_svm)
```

```python
    # Table of metrics
    print('--------------------------------Table of metrics---------------------------------------\n\n')
    data_rows = [('f1 score',f1_lg,f1_knn,f1_svm),
                 ('','','',''),
                  ('accuracy %',acc_lg,acc_knn,acc_svm),
                 ('','','',''),
                 ('confusion matrix',conf_lg[0], conf_knn[0], conf_svm[0]),
                 ('',conf_lg[1], conf_knn[1], conf_svm[1]),
                 ('','','',''),
                 ('ROC score',roc_c_lg,roc_c_knn,roc_c_svm)]
    t = Table(rows=data_rows, names=('metric','Logistic regression', 'KNN', 'SVM'))
    print(t)

    #Plot ROC curve
    print('\n\n---------------------------ROC curves--------------------------------------\n\n')
    fig, _ = plt.subplots()
    fig.suptitle('Comparison of three ROC curves')
    fpr,tpr,thresholds=roc_lg
    plt.plot([0, 1],[0, 1],'--')
    plt.plot(fpr,tpr,label='Logistic regression :'+str(roc_c_lg))
    plt.xlabel('false positive')
    plt.ylabel('false negative')
    fpr,tpr,thresholds=roc_knn
    plt.plot(fpr,tpr,label='KNN :'+str(roc_c_knn))
    fpr,tpr,thresholds=roc_svm
    plt.plot(fpr,tpr,label='SVM :'+str(roc_c_svm))
    plt.legend()
    plt.show()
```

```python
    # Maximum metrics
    print('----------------------------Max of metrics--------------------------------------\n\n')
    data_rows = [('max f1 score',algo_with_max_metric(f1_lg,f1_knn,f1_svm)),
                 ('','','',''),
                  ('max accuracy %',algo_with_max_metric(acc_lg,acc_knn,acc_svm)),
                 ('','','',''),
                 ('max ROC score',algo_with_max_metric(roc_c_lg,roc_c_knn,roc_c_svm))]
    t = Table(rows=data_rows, names=('metric','Learning algorithm winnig'))
    print(t)

# Function returning name of winnig algorithm based on a single metric
def algo_with_max_metric(a,b,c):
    max_metric = max(a,b,c)
    if max_metric == a:
        return 'Logistic regression'
    elif max_metric == b:
        return 'KNN'
    else:
        return 'SVM'
```
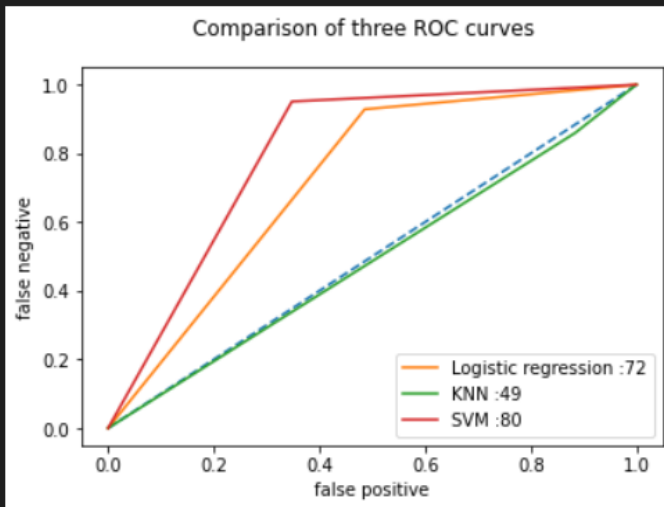
```python
    compare_lg_knn_svm(yt_knn,yp_knn,yt_lg,yp_lg,yt_svm,yp_svm)
```

```
...    ----------------------------Table of metrics--------------------------------------


        metric      Logistic regression   KNN     SVM
    ---------------  -------------------  ------- -------
          f1 score                   74      48      82

       accuracy %                    81      70      84

  confusion matrix               [18 17] [ 3 23] [15  8]
                                 [ 6 78] [13 80] [ 2 39]

         ROC score                   72      49      80
```

```
-------------------------------ROC curves-------------------------------------
```

Comparison of three ROC curves

Logistic regression :72
KNN :49
SVM :80

```
---------------------------Max of metrics-------------------------------------
```

| metric | Learning algorithm winnig |
|--------|--------------------------|
| max f1 score | SVM |
| max accuracy % | SVM |
| max ROC score | SVM |

**INFERENCE**:

- Here we have used three models such as logistic regression, KNN and SVM to check the performance of those models and identify which model performs well with better accuracy.
- The models are hyper tuned to increase their accuracy.

- From this it is identified that SVM provides greater accuracy compared to all other models
- In SVM we identified the accuracy of three different kernels such as linear, polynomial and gaussian kernel among the three the linear kernel of SVM provides greater accuracy of 84%