```
In [33]:    1  import pandas as pd
            2  import numpy as np
```

```
In [34]:    1  data = pd.read_csv("diabetes.csv")
            2  data
```

Out[34]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFur |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | |
| ... | ... | ... | ... | ... | ... | ... | |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | |

768 rows × 9 columns

◀ ▶

```
In [35]:    1  df = pd.DataFrame(data)
            2  df.head()
```

Out[35]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunct |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0. |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0. |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0. |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0. |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2. |

◀ ▶

```
In [36]:    1  df.isnull().sum()
```

Out[36]:
```
Pregnancies                 0
Glucose                     0
BloodPressure               0
SkinThickness               0
Insulin                     0
BMI                         0
DiabetesPedigreeFunction    0
Age                         0
Outcome                     0
dtype: int64
```
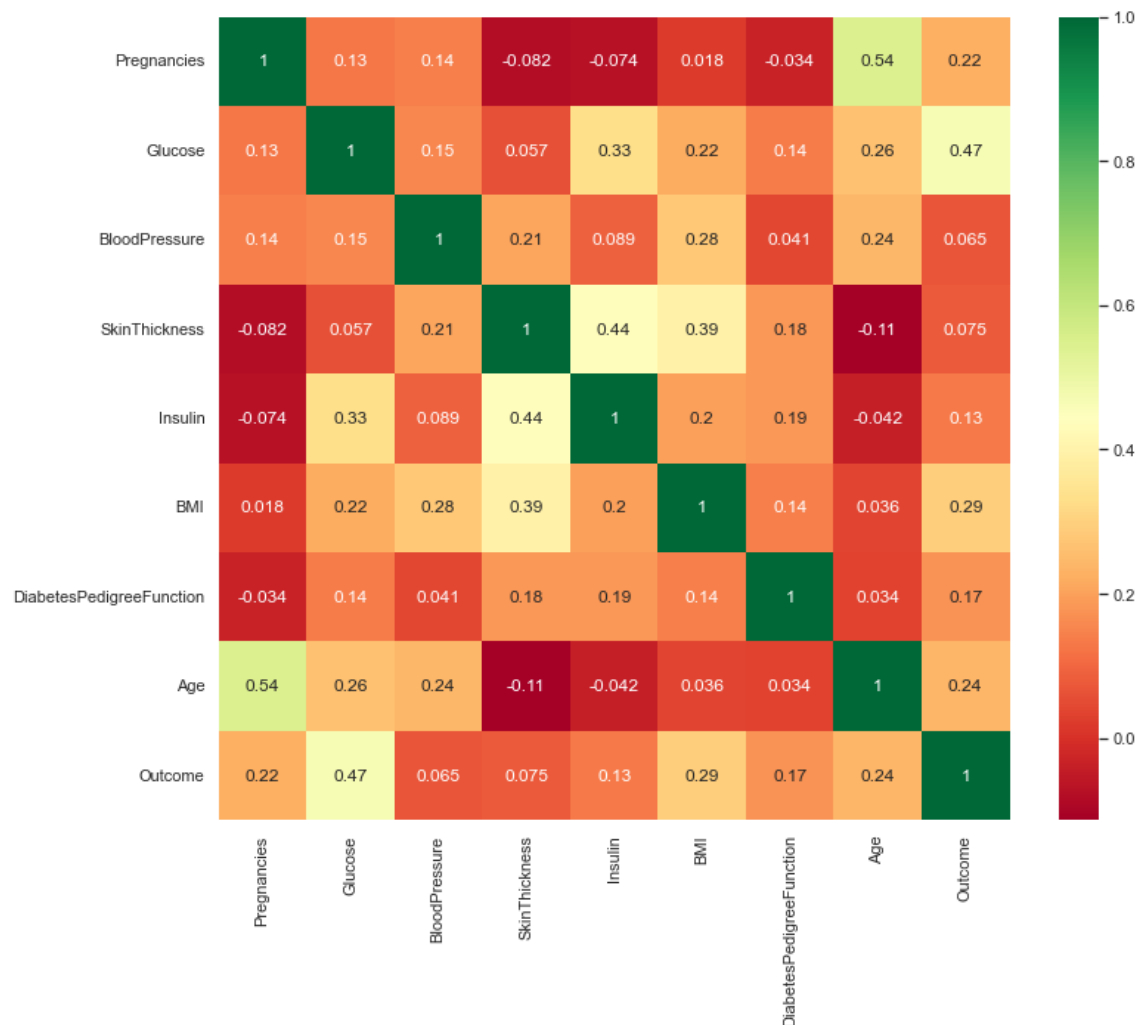
```
In [37]:  1  plt.figure(figsize=(12,10))  # on this line I just set the size of figu
          2  p=sns.heatmap(df.corr(), annot=True,cmap ='RdYlGn')  # seaborn has very
```



# Manipulating and Cleaning our dataset

```
In [38]:  1  cols_clean = ['Glucose','BloodPressure','SkinThickness','Insulin','BMI'
          2  for i in cols_clean:
          3      df[i] = df[i].replace(0,np.NaN)
          4      cols_mean = int(df[i].mean(skipna=True))
          5      df[i] = df[i].replace(np.NaN, cols_mean)
          6  data1 = df
          7  data1.head().style.highlight_max(color="lightblue").highlight_min(color
```

Out[38]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesP |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148.000000 | 72.000000 | 35.000000 | 155.000000 | 33.600000 | |
| 1 | 1 | 85.000000 | 66.000000 | 29.000000 | 155.000000 | 26.600000 | |
| 2 | 8 | 183.000000 | 64.000000 | 29.000000 | 155.000000 | 23.300000 | |
| 3 | 1 | 89.000000 | 66.000000 | 23.000000 | 94.000000 | 28.100000 | |
| 4 | 0 | 137.000000 | 40.000000 | 35.000000 | 168.000000 | 43.100000 | |

```
In [39]:  1  print(data1.describe())
```
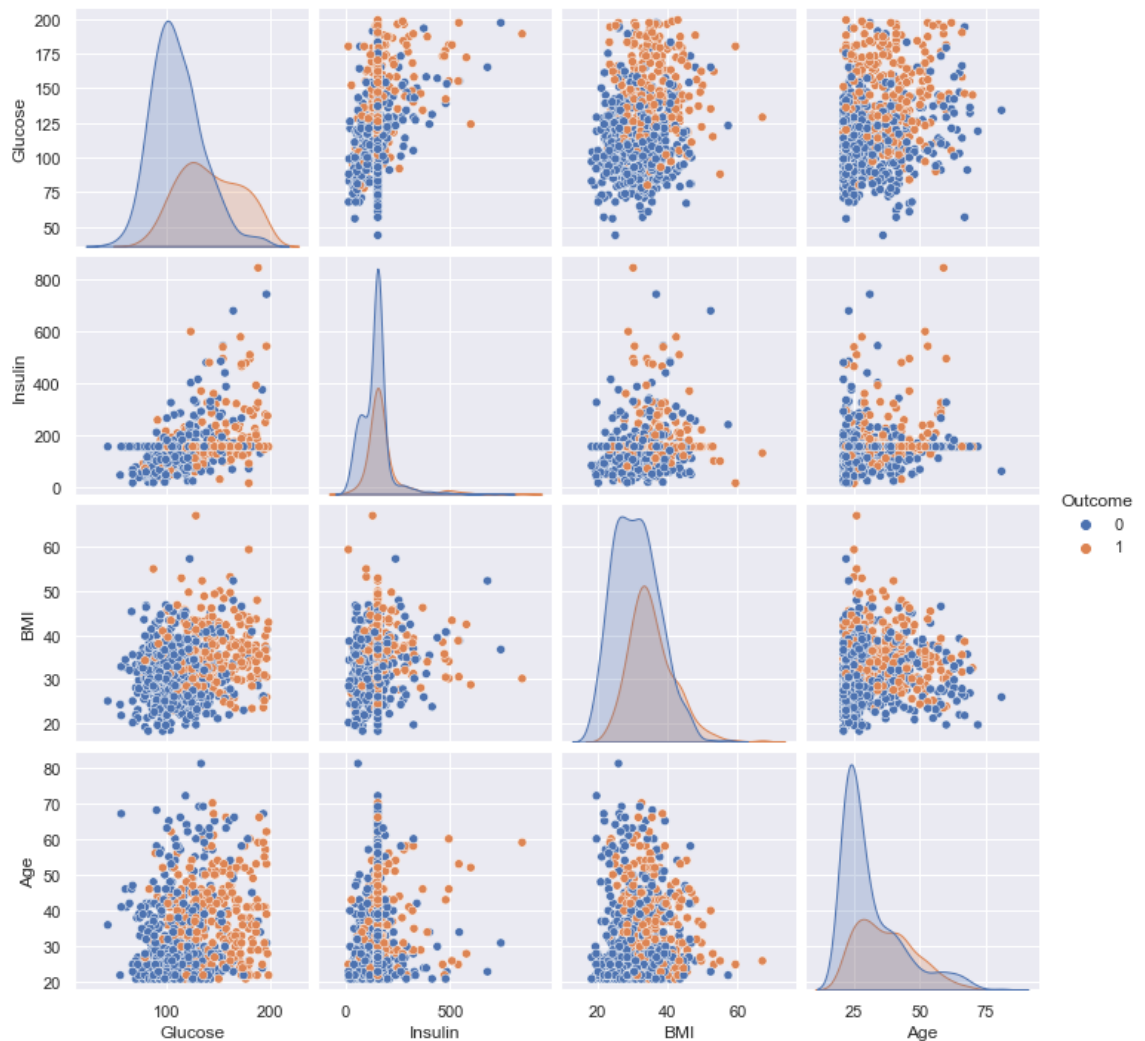
```
       Pregnancies      Glucose  BloodPressure  SkinThickness     Insulin  \
count   768.000000   768.000000     768.000000     768.000000   768.00000
mean      3.845052   121.682292      72.386719      29.108073   155.28125
std       3.369578    30.435999      12.096642       8.791221    85.02155
min       0.000000    44.000000      24.000000       7.000000    14.00000
25%       1.000000    99.750000      64.000000      25.000000   121.50000
50%       3.000000   117.000000      72.000000      29.000000   155.00000
75%       6.000000   140.250000      80.000000      32.000000   155.00000
max      17.000000   199.000000     122.000000      99.000000   846.00000

              BMI  DiabetesPedigreeFunction         Age     Outcome
count  768.000000                768.000000  768.000000  768.000000
mean    32.450911                  0.471876   33.240885    0.348958
std      6.875366                  0.331329   11.760232    0.476951
min     18.200000                  0.078000   21.000000    0.000000
25%     27.500000                  0.243750   24.000000    0.000000
50%     32.000000                  0.372500   29.000000    0.000000
75%     36.600000                  0.626250   41.000000    1.000000
max     67.100000                  2.420000   81.000000    1.000000
```

```
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3  %matplotlib inline
4
5  # graph = ['Glucose','Insulin','BMI','Age','Outcome']
6  sns.set()
7  # print(sns.pairplot(data1[graph],hue='Outcome', diag_kind='kde'))
8  print(sns.pairplot(data1[graph],hue='Outcome', diag_kind='kde'))
```

<seaborn.axisgrid.PairGrid object at 0x000002A3E0B306A0>

```
1  # for the purpose of simplicity and analysing the most relevent  data ,
2  # Glucose , Insulin and BMI
3  # defining variables and features for the dataset for splitting
4  # q_cols = ['Glucose','Insulin','BMI','Outcome']
5  q_cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', '
6
7  df = data1[q_cols]
8  print(df.head(2))
```

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI |
|---|---|---|---|---|---|---|
| 0 | 6 | 148.0 | 72.0 | 35.0 | 155.0 | 33.6 |
| 1 | 1 | 85.0 | 66.0 | 29.0 | 155.0 | 26.6 |

|   | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|
| 0 | 0.627 | 50 | 1 |
| 1 | 0.351 | 31 | 0 |

```python
In [47]:   1  # # let's split the data into training and testing datasets
           2  # split = 0.75 # 75% train and 25% test dataset
           3  # total_len = len(df)
           4  # split_df = int(total_len*split)
           5  # train, test = df.iloc[:split_df,0:4],df.iloc[split_df:,0:4]
           6  # train_x = train[['Glucose','Insulin','BMI']]
           7  # train_y = train['Outcome']
           8  # test_x = test[['Glucose','Insulin','BMI']]
           9  # test_y = test['Outcome']
          10
          11  # Split the data into training and testing datasets
          12  split = 0.75   # 75% train and 25% test dataset
          13  total_len = len(df)
          14  split_df = int(total_len * split)
          15  train, test = df.iloc[:split_df], df.iloc[split_df:]
          16
          17  # Select the columns specified in q_cols for training and testing
          18  train_x = train[q_cols[:-1]]   # Exclude the 'Outcome' column from featu
          19  train_y = train['Outcome']     # Target variable
          20  test_x = test[q_cols[:-1]]     # Exclude the 'Outcome' column from featu
          21  test_y = test['Outcome']       # Target variable
          22
```

```python
In [48]:   1  a = len(train_x)
           2  b = len(test_x)
           3  print(' Training data =',a,'\n','Testing data =',b,'\n','Total data ler
```

```
Training data = 576
Testing data = 192
Total data length =  768
```

```python
In [49]:   1  from sklearn.neighbors import KNeighborsClassifier
           2  from sklearn import metrics
           3
           4  def knn(x_train, y_train, x_test, y_test,n):
           5      n_range = range(1, n)
           6      results = []
           7      for n in n_range:
           8          knn = KNeighborsClassifier(n_neighbors=n)
           9          knn.fit(x_train, y_train)
          10          #Predict the response for test dataset
          11          predict_y = knn.predict(x_test)
          12          accuracy = metrics.accuracy_score(y_test, predict_y)
          13          #matrix = confusion_matrix(y_test,predict_y)
          14          #seaborn_matrix = sns.heatmap(matrix, annot = True, cmap="Blues
          15          results.append(accuracy)
          16      return results
```
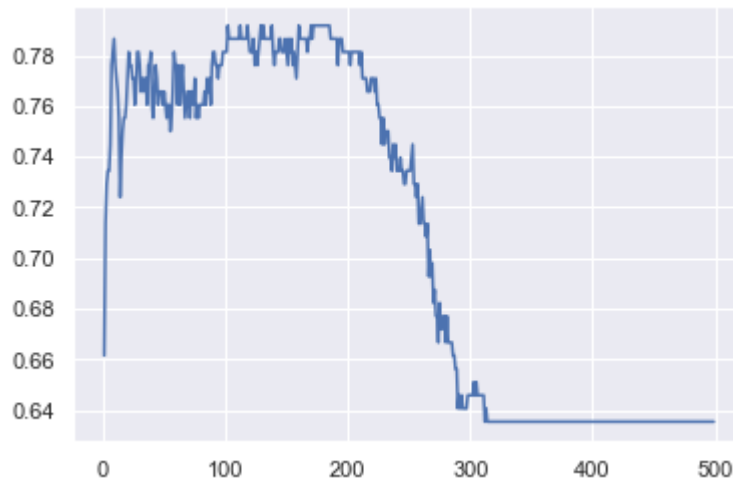
In [50]:
```
1  n= 500
2  output = knn(train_x,train_y,test_x,test_y,n)
3  n_range = range(1, n)
4  plt.plot(n_range, output)
```

Out[50]: [<matplotlib.lines.Line2D at 0x2a3ec0c6100>]



In [51]:
```
1  # best k that could optimize this model is between 100 to 200 offering
2  # ideal k value for this dataset should be 150 give or take
```

```
In [52]:    1  from sklearn.metrics import confusion_matrix
            2  from sklearn.metrics import accuracy_score, precision_score, recall_sco
            3  y_pred = knn(train_x,train_y,test_x,test_y,n)
            4  cnf_matrix = confusion_matrix(test_y, y_pred)
```

---------------------------------------------------------------------------
-
ValueError                                Traceback (most recent call las
t)
~\AppData\Local\Temp/ipykernel_18924/597529570.py in <module>
      2 from sklearn.metrics import accuracy_score, precision_score, recal
l_score, f1_score, fbeta_score
      3 y_pred = knn(train_x,train_y,test_x,test_y,n)
----> 4 cnf_matrix = confusion_matrix(test_y, y_pred)

c:\users\kedar\appdata\local\programs\python\python39\lib\site-packages\sk
learn\metrics\_classification.py in confusion_matrix(y_true, y_pred, label
s, sample_weight, normalize)
    305         (0, 2, 1, 1)
    306     """
--> 307     y_type, y_true, y_pred = _check_targets(y_true, y_pred)
    308     if y_type not in ("binary", "multiclass"):
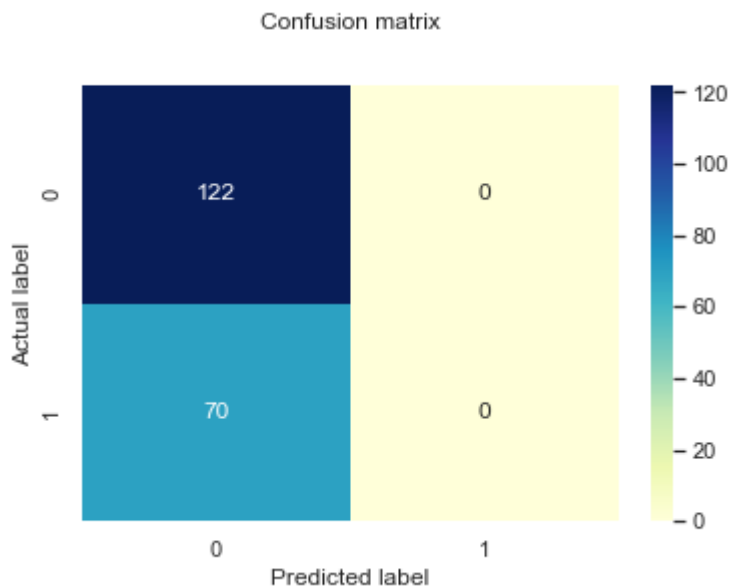    309         raise ValueError("%s is not supported" % y_type)

c:\users\kedar\appdata\local\programs\python\python39\lib\site-packages\sk
learn\metrics\_classification.py in _check_targets(y_true, y_pred)
     82     y_pred : array or indicator matrix
     83     """
---> 84     check_consistent_length(y_true, y_pred)
     85     type_true = type_of_target(y_true, input_name="y_true")
     86     type_pred = type_of_target(y_pred, input_name="y_pred")

c:\users\kedar\appdata\local\programs\python\python39\lib\site-packages\sk
learn\utils\validation.py in check_consistent_length(*arrays)
    385     uniques = np.unique(lengths)
    386     if len(uniques) > 1:
--> 387         raise ValueError(
    388             "Found input variables with inconsistent numbers of sa
mples: %r"
    389             % [int(l) for l in lengths]

ValueError: Found input variables with inconsistent numbers of samples: [1
92, 499]

```
In [53]:    1  p = sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu" ,fm
            2  plt.title('Confusion matrix', y=1.1)
            3  plt.ylabel('Actual label')
            4  plt.xlabel('Predicted label')
```

Out[53]:  Text(0.5, 12.5, 'Predicted label')



```
In [54]:    1  # Define your KNN function to return predictions
            2  def knn2(x_train, y_train, x_test, y_test, n):
            3      knn = KNeighborsClassifier(n_neighbors=n)
            4      knn.fit(x_train, y_train)
            5      # Predict the response for the test dataset
            6      predict_y = knn.predict(x_test)
            7      return predict_y
```

```
In [55]:    1  n = 500
            2  y_pred = knn2(train_x, train_y, test_x, test_y, n)
            3  cnf_matrix = confusion_matrix(test_y, y_pred)
```

```
In [56]:    1  # Now you can calculate other metrics like accuracy, precision, recall,
            2  accuracy = accuracy_score(test_y, y_pred)
            3  precision = precision_score(test_y, y_pred)
            4  recall = recall_score(test_y, y_pred)
            5  f1 = f1_score(test_y, y_pred)
            6  fbeta = fbeta_score(test_y, y_pred, beta=0.5)
```

```
c:\users\kedar\appdata\local\programs\python\python39\lib\site-packages\sk
learn\metrics\_classification.py:1327: UndefinedMetricWarning: Precision i
s ill-defined and being set to 0.0 due to no predicted samples. Use `zero_
division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

```
In [57]:   1  # Print the confusion matrix and other metrics
           2  print("Confusion Matrix:\n", cnf_matrix)
           3  print("Accuracy:", accuracy)
           4  print("Precision:", precision)
           5  print("Recall:", recall)
           6  print("F1 Score:", f1)
           7  print("F-beta Score:", fbeta)
```

```
Confusion Matrix:
 [[122    0]
 [ 70    0]]
Accuracy: 0.6354166666666666
Precision: 0.0
Recall: 0.0
F1 Score: 0.0
F-beta Score: 0.0
```
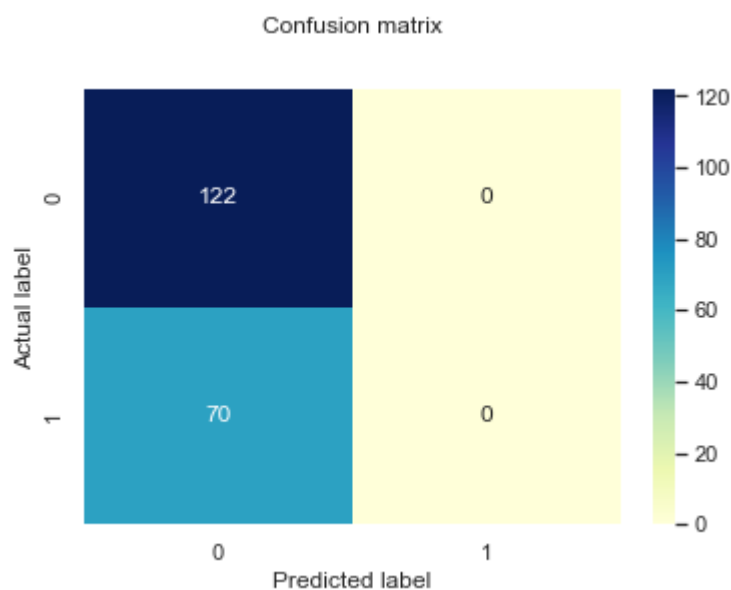
```
In [58]:   1  p = sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu" ,fm
           2  plt.title('Confusion matrix', y=1.1)
           3  plt.ylabel('Actual label')
           4  plt.xlabel('Predicted label')
```

Out[58]:  Text(0.5, 12.5, 'Predicted label')



```
In [60]:   1  from sklearn.model_selection import train_test_split
           2  from sklearn.preprocessing import StandardScaler
           3  from sklearn.neighbors import KNeighborsClassifier
           4  from sklearn.metrics import accuracy_score, confusion_matrix, classific
           5
```

```
In [61]:   1
           2  # Load your dataset
           3  # Replace 'your_dataset.csv' with the actual file path to your dataset
           4  df = pd.read_csv('diabetes.csv')
           5
           6  # Define your feature columns and target column
           7  q_cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', '
           8  target_col = 'Outcome'
           9
          10  # Split the data into features (X) and target (y)
          11  X = df[q_cols]
          12  y = df[target_col]
          13
          14  # Split the data into training and testing datasets
          15  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
          16
          17  # Perform feature scaling (standardization) on the features
          18  scaler = StandardScaler()
          19  X_train_scaled = scaler.fit_transform(X_train)
          20  X_test_scaled = scaler.transform(X_test)
          21
          22  # Create and train a K-nearest neighbors (KNN) classifier
          23  k = 5  # You can adjust the value of k
          24  knn_classifier = KNeighborsClassifier(n_neighbors=k)
          25  knn_classifier.fit(X_train_scaled, y_train)
          26
          27  # Make predictions on the test data
          28  y_pred = knn_classifier.predict(X_test_scaled)
          29
          30  # Evaluate the model
          31  accuracy = accuracy_score(y_test, y_pred)
          32  conf_matrix = confusion_matrix(y_test, y_pred)
          33  classification_rep = classification_report(y_test, y_pred)
          34
          35  # Print the results
          36  print(f"Accuracy: {accuracy}")
          37  print("Confusion Matrix:")
          38  print(conf_matrix)
          39  print("Classification Report:")
          40  print(classification_rep)
```

```
Accuracy: 0.6822916666666666
Confusion Matrix:
[[94 29]
 [32 37]]
Classification Report:
              precision    recall  f1-score   support

           0       0.75      0.76      0.76       123
           1       0.56      0.54      0.55        69

    accuracy                           0.68       192
   macro avg       0.65      0.65      0.65       192
weighted avg       0.68      0.68      0.68       192
```
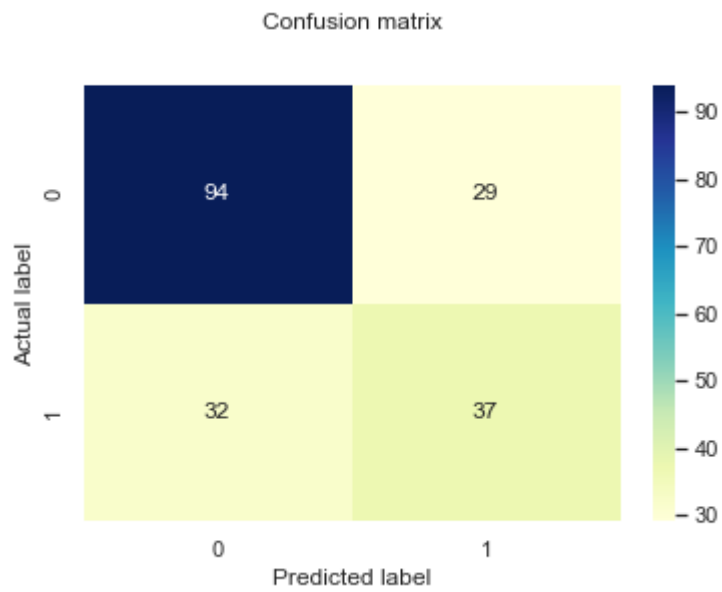
In [63]:
```
1 p = sns.heatmap(pd.DataFrame(conf_matrix), annot=True, cmap="YlGnBu" ,f
2 plt.title('Confusion matrix', y=1.1)
3 plt.ylabel('Actual label')
4 plt.xlabel('Predicted label')
```

Out[63]: Text(0.5, 12.5, 'Predicted label')



Confusion matrix

In [ ]:
```
1
```