# Institute of Information Technology (IIT)
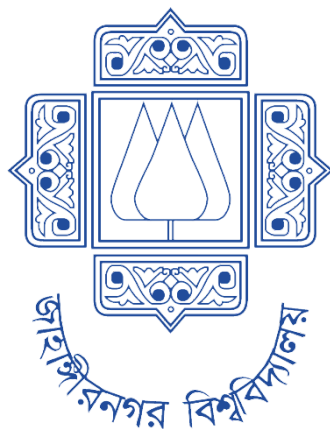
## Jahangirnagar University

**Lab Report: 06**

Submitted by:

**Name:** MD Zuleyenine Ibne Noman

**Roll No:** 2002

**Lab Date:** 22 August,2023

**Submission Date:** 25 August,2023

# Dataset related information:

- filename : survey lung cancer.csv
- Link of dataset : https://www.kaggle.com/code/kelvinfoo123/predicting-lung-cancer-with-knn/input (https://www.kaggle.com/code/kelvinfoo123/predicting-lung-cancer-with-knn/input)
- following are the Dataset characterizations:
    - Total no. of attributes:16
    - No .of instances:284 Attribute information: 1. Gender: M(male), F(female) 2. Age: Age of the patient 3. Smoking: YES=2 , NO=1. 4. Yellow fingers: YES=2 , NO=1. 5. Anxiety: YES=2 , NO=1. 6. Peer_pressure: YES=2 , NO=1. 7. Chronic Disease: YES=2 , NO=1. 8. Fatigue: YES=2 , NO=1. 9. Allergy: YES=2 , NO=1. 10. Wheezing: YES=2 , NO=1. 11. Alcohol: YES=2 , NO=1. 12. Coughing: YES=2 , NO=1. 13. Shortness of Breath: YES=2 , NO=1. 14. Swallowing Difficulty: YES=2 , NO=1. 15. Chest pain: YES=2 , NO=1. 16. Lung Cancer: YES , NO.

# All needed installations:

In [4]:
```
!pip install numpy
```

Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: numpy in /home/eyenine/.local/lib/python3.10/site-packages (1.23.5)

In [5]:
```
!pip install pandas
```

Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: pandas in /home/eyenine/.local/lib/python3.10/site-packages (2.0.3)
Requirement already satisfied: numpy>=1.21.0 in /home/eyenine/.local/lib/python3.10/site-packages (from pandas) (1.23.5)
Requirement already satisfied: tzdata>=2022.1 in /home/eyenine/.local/lib/python3.10/site-packages (from pandas) (2023.3)
Requirement already satisfied: python-dateutil>=2.8.2 in /home/eyenine/.local/lib/python3.10/site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/lib/python3/dist-packages (from pandas) (2022.1)
Requirement already satisfied: six>=1.5 in /usr/lib/python3/dist-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)

In [6]:
```
!pip install seaborn
```

```
Defaulting to user installation because normal site-packages is not writ
eable
Requirement already satisfied: seaborn in /home/eyenine/.local/lib/py
thon3.10/site-packages (0.12.2)
Requirement already satisfied: numpy!=1.24.0,>=1.17 in /home/eyenine/.l
ocal/lib/python3.10/site-packages (from seaborn) (1.23.5)
Requirement already satisfied: pandas>=0.25 in /home/eyenine/.local/li
b/python3.10/site-packages (from seaborn) (2.0.3)
Requirement already satisfied: matplotlib!=3.6.1,>=3.1 in /home/eyenine
/.local/lib/python3.10/site-packages (from seaborn) (3.6.2)
Requirement already satisfied: kiwisolver>=1.0.1 in /home/eyenine/.loc
al/lib/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seabor
n) (1.4.4)
Requirement already satisfied: pillow>=6.2.0 in /usr/lib/python3/dist-
packages (from matplotlib!=3.6.1,>=3.1->seaborn) (9.0.1)
Requirement already satisfied: packaging>=20.0 in /home/eyenine/.local
/lib/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn)
(21.3)
Requirement already satisfied: cycler>=0.10 in /home/eyenine/.local/li
b/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (0.1
1.0)
Requirement already satisfied: pyparsing>=2.2.1 in /usr/lib/python3/di
st-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (2.4.7)
Requirement already satisfied: contourpy>=1.0.1 in /home/eyenine/.loca
l/lib/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn)
(1.0.6)
Requirement already satisfied: fonttools>=4.22.0 in /home/eyenine/.loca
l/lib/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn)
(4.38.0)
Requirement already satisfied: python-dateutil>=2.7 in /home/eyenine/.
local/lib/python3.10/site-packages (from matplotlib!=3.6.1,>=3.1->seab
orn) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/lib/python3/dist-p
ackages (from pandas>=0.25->seaborn) (2022.1)
Requirement already satisfied: tzdata>=2022.1 in /home/eyenine/.local/
lib/python3.10/site-packages (from pandas>=0.25->seaborn) (2023.3)
Requirement already satisfied: six>=1.5 in /usr/lib/python3/dist-pack
ages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.1->seaborn) (1.1
6.0)
```

## Import Libraries

In [7]:
```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

## Read the file /dataset

In [8]:
```
cancer=pd.read_csv("survey lung cancer.csv")
```

In [9]: `cancer`

Out[9]:

| | GENDER | AGE | SMOKING | YELLOW_FINGERS | ANXIETY | PEER_PRESSURE | CHRONIC DISEASE |
|---|---|---|---|---|---|---|---|
| 0 | M | 69 | 1 | 2 | 2 | 1 | 1 |
| 1 | M | 74 | 2 | 1 | 1 | 1 | 2 |
| 2 | F | 59 | 1 | 1 | 1 | 2 | 1 |
| 3 | M | 63 | 2 | 2 | 2 | 1 | 1 |
| 4 | F | 63 | 1 | 2 | 1 | 1 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 304 | F | 56 | 1 | 1 | 1 | 2 | 2 |
| 305 | M | 70 | 2 | 1 | 1 | 1 | 1 |
| 306 | M | 58 | 2 | 1 | 1 | 1 | 1 |
| 307 | M | 67 | 2 | 1 | 2 | 1 | 1 |
| 308 | M | 62 | 1 | 1 | 1 | 2 | 1 |

309 rows × 16 columns

# Exploring the dataset

In [10]: `cancer.head()`

Out[10]:

| | GENDER | AGE | SMOKING | YELLOW_FINGERS | ANXIETY | PEER_PRESSURE | CHRONIC DISEASE | F/ |
|---|---|---|---|---|---|---|---|---|
| 0 | M | 69 | 1 | 2 | 2 | 1 | 1 | |
| 1 | M | 74 | 2 | 1 | 1 | 1 | 2 | |
| 2 | F | 59 | 1 | 1 | 1 | 2 | 1 | |
| 3 | M | 63 | 2 | 2 | 2 | 1 | 1 | |
| 4 | F | 63 | 1 | 2 | 1 | 1 | 1 | |

In [11]: `cancer.shape`

Out[11]: (309, 16)

In [13]: `cancer.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 309 entries, 0 to 308
Data columns (total 16 columns):
 #  Column               Non-Null Count Dtype
--- ------               -------------- -----
 0  GENDER               309 non-null   object
 1  AGE                  309 non-null   int64
 2  SMOKING              309 non-null   int64
 3  YELLOW_FINGERS       309 non-null   int64
 4  ANXIETY              309 non-null   int64
 5  PEER_PRESSURE        309 non-null   int64
 6  CHRONIC DISEASE      309 non-null   int64
 7  FATIGUE              309 non-null   int64
 8  ALLERGY              309 non-null   int64
 9  WHEEZING             309 non-null   int64
 10 ALCOHOL CONSUMING    309 non-null   int64
 11 COUGHING             309 non-null   int64
 12 SHORTNESS OF BREATH  309 non-null   int64
 13 SWALLOWING DIFFICULTY 309 non-null  int64
 14 CHEST PAIN           309 non-null   int64
 15 LUNG_CANCER          309 non-null   object
dtypes: int64(14), object(2)
memory usage: 38.8+ KB
```

In [14]: `cancer.describe()`

Out[14]:

|  | AGE | SMOKING | YELLOW_FINGERS | ANXIETY | PEER_PRESSURE | CHRONIC DISEASE |
|---|---|---|---|---|---|---|
| count | 309.000000 | 309.000000 | 309.000000 | 309.000000 | 309.000000 | 309.000000 |
| mean | 62.673139 | 1.563107 | 1.569579 | 1.498382 | 1.501618 | 1.504854 |
| std | 8.210301 | 0.496806 | 0.495938 | 0.500808 | 0.500808 | 0.500787 |
| min | 21.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 25% | 57.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 50% | 62.000000 | 2.000000 | 2.000000 | 1.000000 | 2.000000 | 2.000000 |
| 75% | 69.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 |
| max | 87.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 |

## Dataset preprocessing

```
!pip install sklearn
!pip install scikit-learn
```

In [22]: 
```python
from sklearn.preprocessing import LabelEncoder
lb=LabelEncoder()
cancer["LUNG_CANCER"]=lb.fit_transform(cancer["LUNG_CANCER"])
```
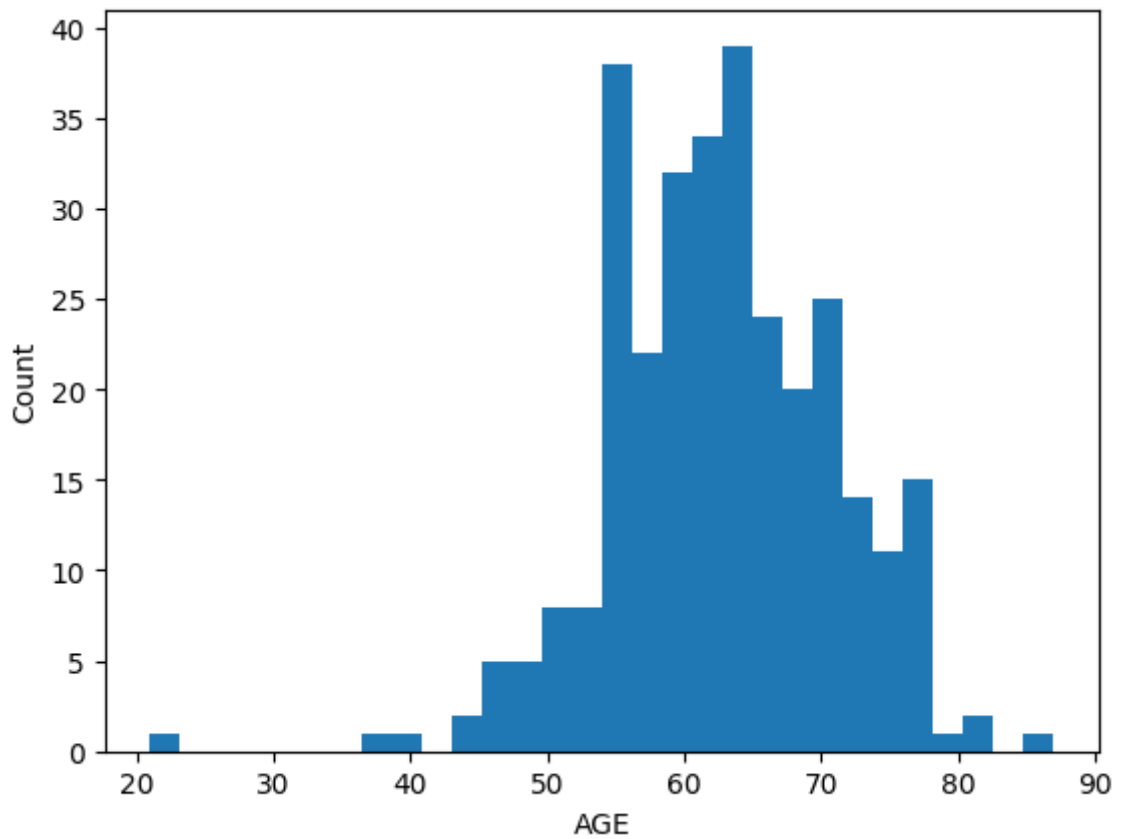
In [23]: 
```python
cancer["GENDER"]=lb.fit_transform(cancer["GENDER"])
```

In [24]: `cancer["LUNG_CANCER"].value_counts()`

Out[24]: 
```
LUNG_CANCER
1   270
0    39
Name: count, dtype: int64
```

In [25]: 
```python
plt.hist(cancer["AGE"],bins=30)
plt.xlabel("AGE")
plt.ylabel("Count")
plt.show()
```



In [26]: `cancer=cancer[cancer.AGE>30]`

In [27]:
```python
cancer_without_age = cancer.drop(["AGE"], axis = 1)
for i in cancer_without_age.columns:
    print(cancer_without_age[i].value_counts())
```

```
GENDER
1   162
0   146
Name: count, dtype: int64
SMOKING
2   173
1   135
Name: count, dtype: int64
YELLOW_FINGERS
2   176
1   132
Name: count, dtype: int64
ANXIETY
2   154
1   154
Name: count, dtype: int64
PEER_PRESSURE
2   155
1   153
Name: count, dtype: int64
CHRONIC DISEASE
2   155
1   153
Name: count, dtype: int64
FATIGUE
2   207
1   101
Name: count, dtype: int64
ALLERGY
2   171
1   137
Name: count, dtype: int64
WHEEZING
2   172
1   136
Name: count, dtype: int64
ALCOHOL CONSUMING
2   172
1   136
Name: count, dtype: int64
COUGHING
2   179
1   129
Name: count, dtype: int64
SHORTNESS OF BREATH
2   197
1   111
Name: count, dtype: int64
SWALLOWING DIFFICULTY
1   163
2   145
Name: count, dtype: int64
CHEST PAIN
2   172
```

```
1   136
Name: count, dtype: int64
LUNG_CANCER
1   270
0   38
Name: count, dtype: int64
```
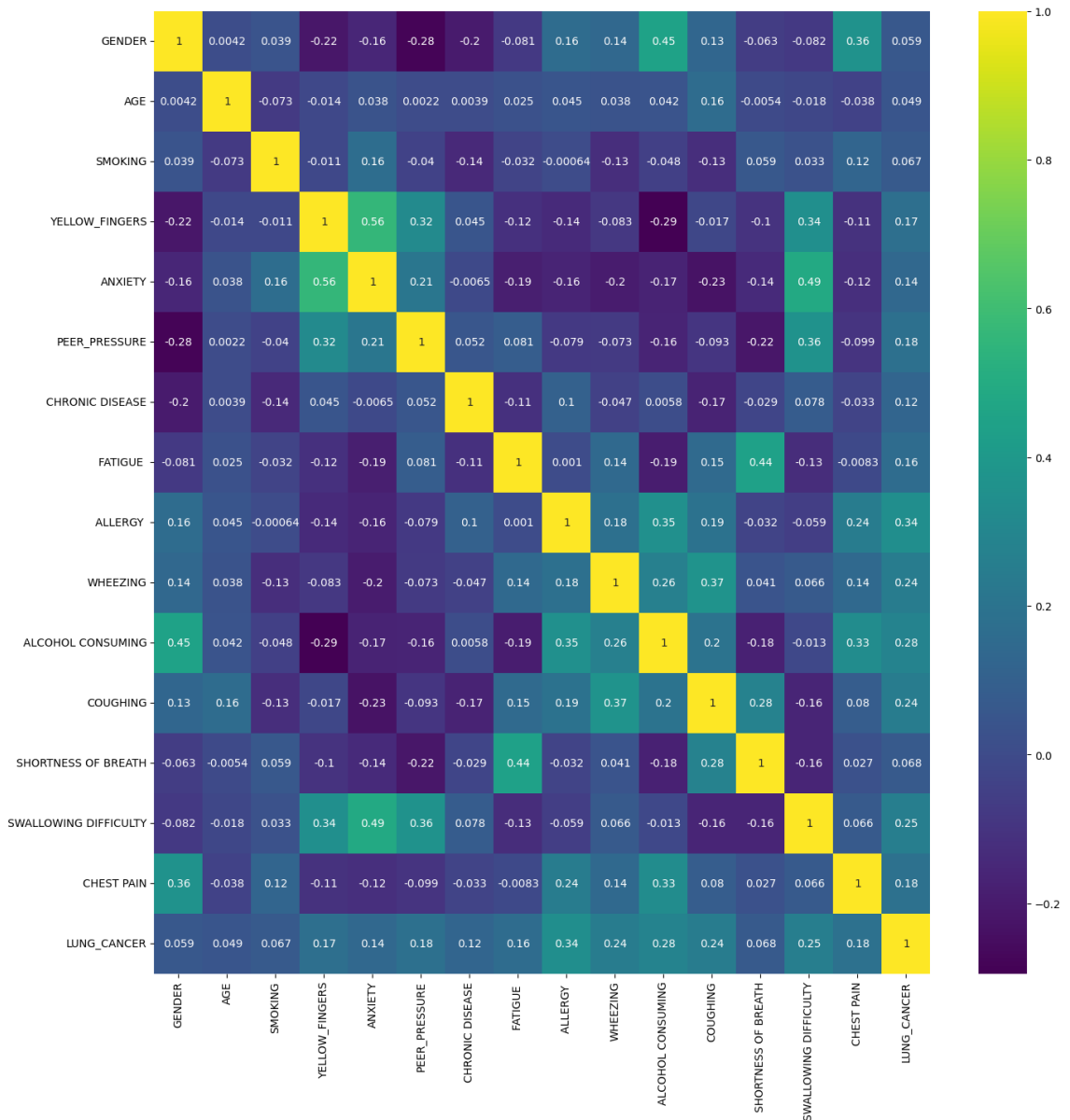
In [28]:
```python
plt.figure(figsize=(16,16))
sns.heatmap(cancer.corr(),annot=True,cmap="viridis")
```

Out[28]: <AxesSubplot: >



# KNN model Training and Evaluation

```python
from sklearn.model_selection import KFold
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.neighbors import KNeighborsClassifier
```

In [33]: 
```python
x_train, x_test, y_train, y_test = train_test_split(cancer.drop('LUNG_
```

In [34]: 
```python
knn = KNeighborsClassifier(n_neighbors = 11)
knn.fit(x_train,y_train)
```

Out[34]: 
```
▼          KNeighborsClassifier
KNeighborsClassifier(n_neighbors=11)
```

In [35]: 
```python
predictions=knn.predict(x_test)
```

In [36]: 
```python
print(confusion_matrix(y_test,predictions))
```
```
[[ 0 10]
 [ 1 82]]
```

In [37]: 
```python
print(accuracy_score(y_test,predictions))
```
```
0.8817204301075269
```

In [38]: 
```python
print(classification_report(y_test, predictions))
```
```
              precision   recall  f1-score   support

           0       0.00     0.00      0.00        10
           1       0.89     0.99      0.94        83

    accuracy                          0.88        93
   macro avg       0.45     0.49      0.47        93
weighted avg       0.80     0.88      0.84        93
```

In [40]: 
```python
print(precision_score(y_test, predictions))
```
```
0.8913043478260869
```

In [41]: 
```python
print(recall_score(y_test, predictions))
```
```
0.9879518072289156
```

In [42]: 
```python
print(f1_score(y_test, predictions))
```
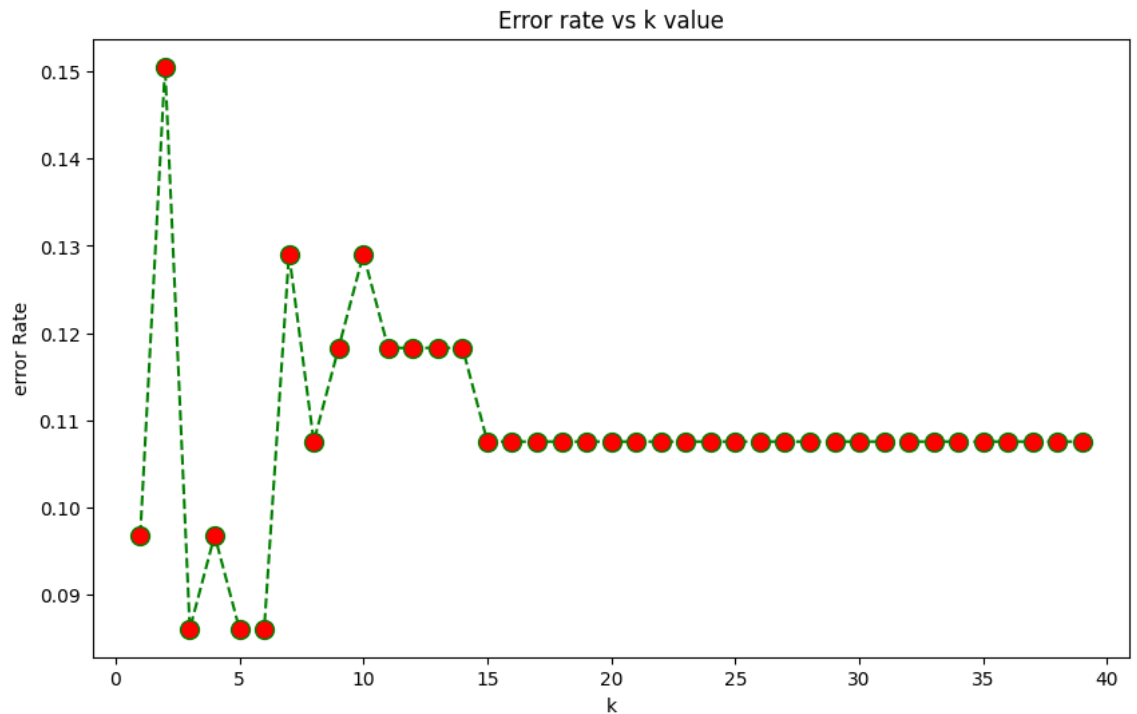```
0.9371428571428572
```

In [43]: 
```python
print(confusion_matrix(y_test, predictions))
```
```
[[ 0 10]
 [ 1 82]]
```

In [45]: 
```python
errorRate=[]
for i in range(1,40):
  knn=KNeighborsClassifier(n_neighbors=i)
  knn.fit(x_train,y_train)
  predictions_i=knn.predict(x_test)
  errorRate.append(np.mean(predictions_i!=y_test))
```

In [50]:
```python
plt.figure(figsize=(10,6))
plt.plot(range(1,40),errorRate,color='green',linestyle='dashed',marker='o',r
plt.title('Error rate vs k value')
plt.xlabel('k')
plt.ylabel('error Rate')
```
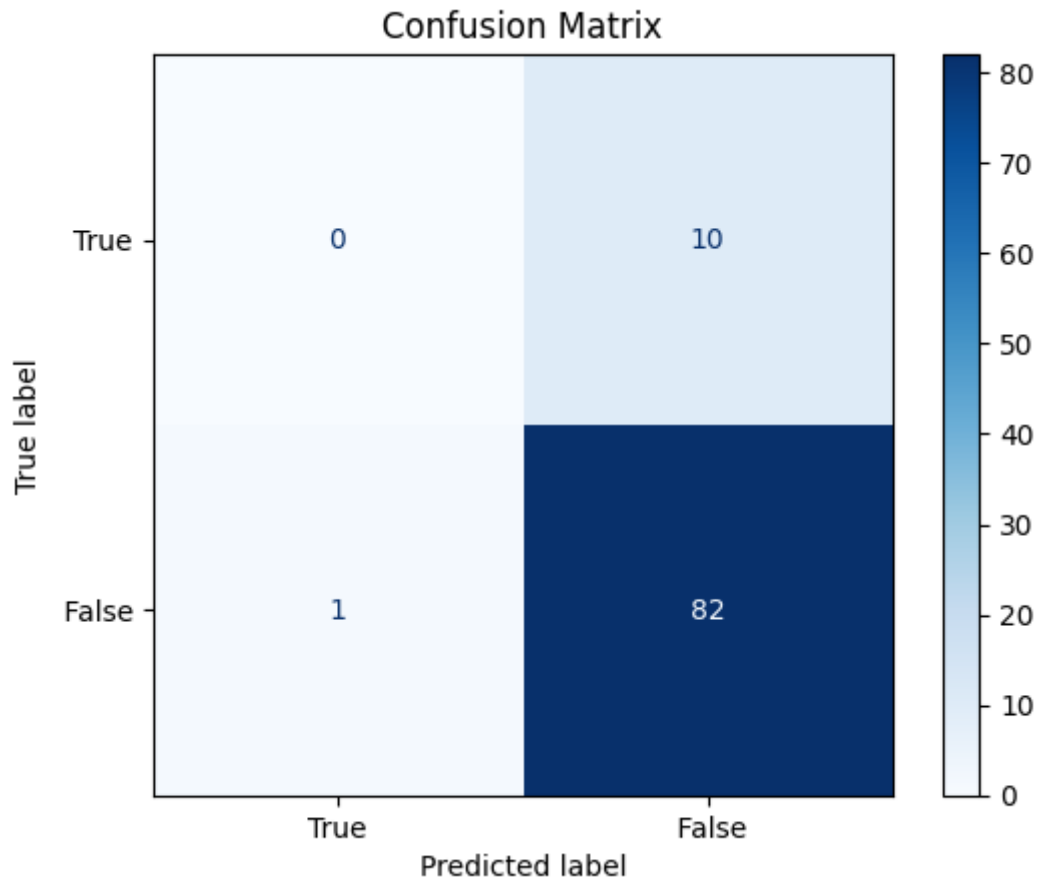
Out[50]: Text(0, 0.5, 'error Rate')



In [59]:
```python
knn = KNeighborsClassifier(n_neighbors = 13)
knn.fit(x_train,y_train)
predictions=knn.predict(x_test)
```

In [60]:
```python
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay


cm = confusion_matrix(y_test, predictions)


disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["T
disp.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix")
plt.show()
```



Confusion Matrix

In [55]: `print(classification_report(y_test,predictions))`

```
              precision    recall  f1-score   support

           0       0.00      0.00      0.00        10
           1       0.89      1.00      0.94        83

    accuracy                           0.89        93
   macro avg       0.45      0.50      0.47        93
weighted avg       0.80      0.89      0.84        93
```

```
/home/eyenine/.local/lib/python3.10/site-packages/sklearn/metrics/
_classification.py:1469: UndefinedMetricWarning: Precision and F-score
are ill-defined and being set to 0.0 in labels with no predicted samples. U
se `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/home/eyenine/.local/lib/python3.10/site-packages/sklearn/metrics/
_classification.py:1469: UndefinedMetricWarning: Precision and F-score
are ill-defined and being set to 0.0 in labels with no predicted samples. U
se `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/home/eyenine/.local/lib/python3.10/site-packages/sklearn/metrics/
_classification.py:1469: UndefinedMetricWarning: Precision and F-score
are ill-defined and being set to 0.0 in labels with no predicted samples. U
se `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

In [ ]: