

# index

March 16, 2022

Link is [here](#) or you can use

<https://github.com/miladshiraniUCB/dsc-knn-with-scikit-learn-lab.git>

## 1 KNN with scikit-learn - Lab

### 1.1 Introduction

In this lab, you'll learn how to use scikit-learn's implementation of a KNN classifier on the classic Titanic dataset from Kaggle!

### 1.2 Objectives

In this lab you will:

- Conduct a parameter search to find the optimal value for K
- Use a KNN classifier to generate predictions on a real-world dataset
- Evaluate the performance of a KNN model

### 1.3 Getting Started

Start by importing the dataset, stored in the `titanic.csv` file, and previewing it.

[ ]:

```
[51]: # Your code here
      # Import pandas and set the standard alias

import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

# Import the data from 'titanic.csv' and store it in a pandas DataFrame
raw_df = pd.read_csv("titanic.csv")
```

```
# Print the head of the DataFrame to ensure everything loaded correctly
raw_df.head()
```

```
[51]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

  

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

  

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

Great! Next, you'll perform some preprocessing steps such as removing unnecessary columns and normalizing features.

## 1.4 Preprocessing the data

Preprocessing is an essential component in any data science pipeline. It's not always the most glamorous task as might be an engaging data visual or impressive neural network, but cleaning and normalizing raw datasets is very essential to produce useful and insightful datasets that form the backbone of all data powered projects. This can include changing column types, as in:

```
df['col_name'] = df['col_name'].astype('int')
```

Or extracting subsets of information, such as:

```
import re
df['street'] = df['address'].map(lambda x: re.findall('(.*)?\n', x)[0])
```

**Note:** While outside the scope of this particular lesson, **regular expressions** (mentioned above) are powerful tools for pattern matching! See the [regular expressions official documentation here](#).

Since you've done this before, you should be able to do this quite well yourself without much hand holding by now. In the cells below, complete the following steps:

1. Remove unnecessary columns ('PassengerId', 'Name', 'Ticket', and 'Cabin')
2. Convert 'Sex' to a binary encoding, where female is 0 and male is 1

3. Detect and deal with any missing values in the dataset:
  - For 'Age', replace missing values with the median age for the dataset
  - For 'Embarked', drop the rows that contain missing values
4. One-hot encode categorical columns such as 'Embarked'
5. Store the target column, 'Survived', in a separate variable and remove it from the DataFrame

```
[113]: # Drop the unnecessary columns
df = raw_df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], axis = 1)
df.head()
```

```
[113]:
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	0	3	male	22.0	1	0	7.2500	S
1	1	1	female	38.0	1	0	71.2833	C
2	1	3	female	26.0	0	0	7.9250	S
3	1	1	female	35.0	1	0	53.1000	S
4	0	3	male	35.0	0	0	8.0500	S

```
[53]: # Convert Sex to binary encoding
df['Sex'] = df['Sex'].map({"female":0, "male":1})
df.head()
```

```
[53]:
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	0	3	1	22.0	1	0	7.2500	S
1	1	1	0	38.0	1	0	71.2833	C
2	1	3	0	26.0	0	0	7.9250	S
3	1	1	0	35.0	1	0	53.1000	S
4	0	3	1	35.0	0	0	8.0500	S

```
[54]: # Find the number of missing values in each column
df.isna().sum()
```

```
[54]: Survived      0
Pclass          0
Sex             0
Age            177
SibSp           0
Parch           0
Fare           0
Embarked        2
dtype: int64
```

```
[57]: # Impute the missing values in 'Age'
df['Age'] = df["Age"].fillna(df["Age"].median())
df.isna().sum()
```

```
[57]: Survived    0
      Pclass     0
      Sex        0
      Age        0
      SibSp      0
      Parch      0
      Fare       0
      Embarked   2
      dtype: int64
```

```
[58]: # Drop the rows missing values in the 'Embarked' column
      df = df.dropna()
      df.isna().sum()
```

```
[58]: Survived    0
      Pclass     0
      Sex        0
      Age        0
      SibSp      0
      Parch      0
      Fare       0
      Embarked    0
      dtype: int64
```

```
[59]: # One-hot encode the categorical columns
      from sklearn.preprocessing import OneHotEncoder

      OHE = OneHotEncoder(categories='auto', sparse=False, handle_unknown='ignore')

      embarked = OHE.fit_transform(df[["Embarked"]])
      embarked_df = pd.DataFrame(embarked,
                                columns = OHE.categories_[0],
                                index = df.index)

      df.drop("Embarked", axis = 1, inplace = True)
      one_hot_df = pd.concat([df, embarked_df], axis = 1)
      one_hot_df.head()
```

```
[59]:   Survived  Pclass  Sex  Age  SibSp  Parch    Fare     C     Q     S
0         0       3    1  22.0     1     0   7.2500  0.0  0.0  1.0
1         1       1    0  38.0     1     0  71.2833  1.0  0.0  0.0
2         1       3    0  26.0     0     0   7.9250  0.0  0.0  1.0
3         1       1    0  35.0     1     0  53.1000  0.0  0.0  1.0
4         0       3    1  35.0     0     0   8.0500  0.0  0.0  1.0
```

```
[72]: # Assign the 'Survived' column to labels
labels = one_hot_df["Survived"]

# Drop the 'Survived' column from one_hot_df
one_hot_df.drop(columns = ["Survived"], axis = 1, inplace = True)
```

## 1.5 Create training and test sets

Now that you've preprocessed the data, it's time to split it into training and test sets.

In the cell below:

- Import `train_test_split` from the `sklearn.model_selection` module
- Use `train_test_split()` to split the data into training and test sets, with a `test_size` of 0.25. Set the `random_state` to 42

```
[73]: # Import train_test_split
from sklearn.model_selection import train_test_split

# Split the data
X_train, X_test, y_train, y_test = train_test_split(one_hot_df, labels,
                                                    test_size = 0.25,
                                                    random_state = 42)
```

## 1.6 Normalizing the data

The final step in your preprocessing efforts for this lab is to **normalize** the data. We normalize **after** splitting our data into training and test sets. This is to avoid information “leaking” from our test set into our training set (read more about data leakage [here](#) ). Remember that normalization (also sometimes called **Standardization** or **Scaling**) means making sure that all of your data is represented at the same scale. The most common way to do this is to convert all numerical values to z-scores.

Since KNN is a distance-based classifier, if data is in different scales, then larger scaled features have a larger impact on the distance between points.

To scale your data, use `StandardScaler` found in the `sklearn.preprocessing` module.

In the cell below:

- Import and instantiate `StandardScaler`
- Use the scaler's `.fit_transform()` method to create a scaled version of the training dataset
- Use the scaler's `.transform()` method to create a scaled version of the test dataset
- The result returned by `.fit_transform()` and `.transform()` methods will be numpy arrays, not a pandas DataFrame. Create a new pandas DataFrame out of this object called `scaled_df`. To set the column names back to their original state, set the `columns` parameter to `one_hot_df.columns`
- Print the head of `scaled_df` to ensure everything worked correctly

```
[75]: # Import StandardScaler
from sklearn.preprocessing import StandardScaler

# Instantiate StandardScaler
scaler = StandardScaler()

# Transform the training and test sets
scaled_data_train = scaler.fit_transform(X_train)
scaled_data_test = scaler.transform(X_test)

# # Convert into a DataFrame
scaled_df_train = pd.DataFrame(scaled_data_train, columns = X_train.columns)
scaled_df_train.head()
```

```
[75]:      Pclass      Sex      Age      SibSp      Parch      Fare      C \
0  0.815528 -1.390655 -0.575676 -0.474917 -0.480663 -0.500108 -0.483046
1 -0.386113 -1.390655  1.550175 -0.474917 -0.480663 -0.435393 -0.483046
2 -0.386113  0.719086 -0.120137 -0.474917 -0.480663 -0.644473 -0.483046
3 -1.587755  0.719086 -0.120137 -0.474917 -0.480663 -0.115799 -0.483046
4  0.815528 -1.390655 -1.107139  0.413551 -0.480663 -0.356656  2.070197

      Q      S
0 -0.311768  0.620174
1 -0.311768  0.620174
2 -0.311768  0.620174
3 -0.311768  0.620174
4 -0.311768 -1.612452
```

You may have noticed that the scaler also scaled our binary/one-hot encoded columns, too! Although it doesn't look as pretty, this has no negative effect on the model. Each 1 and 0 have been replaced with corresponding decimal values, but each binary column still only contains 2 values, meaning the overall information content of each column has not changed.

## 1.7 Fit a KNN model

Now that you've preprocessed the data it's time to train a KNN classifier and validate its accuracy.

In the cells below:

- Import `KNeighborsClassifier` from the `sklearn.neighbors` module
- Instantiate the classifier. For now, you can just use the default parameters
- Fit the classifier to the training data/labels
- Use the classifier to generate predictions on the test data. Store these predictions inside the variable `test_preds`

```
[78]: # Import KNeighborsClassifier
from sklearn.neighbors import KNeighborsClassifier
```

```

# Instantiate KNeighborsClassifier
clf = KNeighborsClassifier()

# Fit the classifier
clf.fit(scaled_data_train, y_train)

# Predict on the test set
test_preds = clf.predict(scaled_data_test)

```

## 1.8 Evaluate the model

Now, in the cells below, import all the necessary evaluation metrics from `sklearn.metrics` and complete the `print_metrics()` function so that it prints out **Precision, Recall, Accuracy, and F1-Score** when given a set of labels (the true values) and `preds` (the models predictions).

Finally, use `print_metrics()` to print the evaluation metrics for the test predictions stored in `test_preds`, and the corresponding labels in `y_test`.

```

[86]: # Your code here
      # Import the necessary functions

      from sklearn.metrics import precision_recall_fscore_support

      p, r, f, s = precision_recall_fscore_support(y_test, test_preds)
      print("precision score : ", p)
      print("recall score    : ", r)
      print("f1 score       : ", f)

```

```

precision score : [0.84057971 0.70588235]
recall score    : [0.82269504 0.73170732]
f1 score       : [0.83154122 0.71856287]

```

```

[83]: from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score

      print("precision score : ", precision_score(y_test, test_preds))
      print("recall score    : ", recall_score(y_test, test_preds))
      print("f1 score       : ", f1_score(y_test, test_preds))
      print("accuracy score  : ", accuracy_score(y_test, test_preds))

```

```

precision score : 0.7058823529411765
recall score    : 0.7317073170731707
f1 score       : 0.718562874251497
accuracy score  : 0.7892376681614349

```

```

[88]: from sklearn.metrics import classification_report
      print(classification_report(y_test, test_preds))

```

	precision	recall	f1-score	support
0	0.84	0.82	0.83	141
1	0.71	0.73	0.72	82
accuracy			0.79	223
macro avg	0.77	0.78	0.78	223
weighted avg	0.79	0.79	0.79	223

```
[96]: # Complete the function
def print_metrics(labels, preds):
    print("Precision Score: {}".format(precision_score(labels, preds)))
    print("Recall Score: {}".format(recall_score(labels, preds)))
    print("Accuracy Score: {}".format(accuracy_score(labels, preds)))
    print("F1 Score: {}".format(f1_score(labels, preds)))

print_metrics(y_test, test_preds)
```

```
Precision Score: 0.7058823529411765
Recall Score:    0.7317073170731707
Accuracy Score: 0.7892376681614349
F1 Score:       0.718562874251497
```

Interpret each of the metrics above, and explain what they tell you about your model's capabilities. If you had to pick one score to best describe the performance of the model, which would you choose? Explain your answer.

Write your answer below this line:

---

## 1.9 Improve model performance

While your overall model results should be better than random chance, they're probably mediocre at best given that you haven't tuned the model yet. For the remainder of this notebook, you'll focus on improving your model's performance. Remember that modeling is an *iterative process*, and developing a baseline out of the box model such as the one above is always a good start.

First, try to find the optimal number of neighbors to use for the classifier. To do this, complete the `find_best_k()` function below to iterate over multiple values of K and find the value of K that returns the best overall performance.

The function takes in six arguments: `* X_train * y_train * X_test * y_test * min_k` (default is 1) `* max_k` (default is 25)

**Pseudocode Hint:** 1. Create two variables, `best_k` and `best_score` 1. Iterate through every *odd number* between `min_k` and `max_k + 1`. 1. For each iteration: 1. Create a new KNN classifier, and set the `n_neighbors` parameter to the current value for k, as determined by the loop 1. Fit this classifier to the training data 1. Generate predictions for `X_test` using the fitted classifier 1. Calculate the *F1-score* for these



predictions 1. Compare this F1-score to `best_score`. If better, update `best_score` and `best_k` 1. Once all iterations are complete, print the best value for `k` and the F1-score it achieved

```
[97]: def find_best_k(X_train, y_train, X_test, y_test, min_k=1, max_k=25):
      # Your code here
      interval = range(1, 26, 2)
      best_score = 0
      best_k = 0
      for k in interval:
          KNN = KNeighborsClassifier(n_neighbors = k)
          KNN.fit(X_train, y_train)
          y_preds = KNN.predict(X_test)
          f1 = f1_score(y_test, y_preds)
          if f1 >= best_score:
              best_score = f1
              best_k = k
      print(f"best f1 score is: {best_score}")
      print(f"best k number is: {best_k}")
      return best_k
```

```
[98]: find_best_k(scaled_data_train, y_train, scaled_data_test, y_test)
      # Expected Output:

      # Best Value for k: 17
      # F1-Score: 0.7468354430379746
```

```
best f1 score is: 0.7468354430379746
best k number is: 17
```

```
[98]: 17
```

If all went well, you'll notice that model performance has improved by 3 percent by finding an optimal value for `k`. For further tuning, you can use scikit-learn's built-in `GridSearch()` to perform a similar exhaustive check of hyperparameter combinations and fine tune model performance. For a full list of model parameters, see the [sklearn documentation](#) !

## 1.10 (Optional) Level Up: Iterating on the data

As an optional (but recommended!) exercise, think about the decisions you made during the preprocessing steps that could have affected the overall model performance. For instance, you were asked to replace the missing age values with the column median. Could this have affected the overall performance? How might the model have fared if you had just dropped those rows, instead of using the column median? What if you reduced the data's dimensionality by ignoring some less important columns altogether?

In the cells below, revisit your preprocessing stage and see if you can improve the overall results of the classifier by doing things differently. Consider dropping certain columns, dealing with missing values differently, or using an alternative scaling function. Then see how these different preprocess-

ing techniques affect the performance of the model. Remember that the `find_best_k()` function handles all of the fitting; use this to iterate quickly as you try different strategies for dealing with data preprocessing!

### 1.10.1 We use mean values for age instead of median

```
[146]: ## Data Preperation, we use mean values for age instead of median

df_new = raw_df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], axis = 1)
df_new["Age"] = df_new["Age"].fillna(df_new["Age"].mean())
df_new.dropna( subset = ["Embarked"], inplace = True, axis = 0)
df_new['Sex'] = df_new['Sex'].map({"female":0, "male":1})
df_new.isna().sum()

### Using OneHotEncoder for Embarked
from sklearn.preprocessing import OneHotEncoder
OHE = OneHotEncoder(categories='auto', sparse=False, handle_unknown='ignore')

embarked_new = OHE.fit_transform(df_new[["Embarked"]])
embarked_df_new = pd.DataFrame(embarked_new,
                               columns = OHE.categories_[0],
                               index = df_new.index)

df_new.drop("Embarked", axis = 1, inplace = True)
one_hot_df_new = pd.concat([df_new, embarked_df_new], axis = 1)
one_hot_df_new.head()

#### defining X and y

y = one_hot_df_new["Survived"]
X = one_hot_df_new.drop(columns = ["Survived"], axis = 1)

### Creating train and test sets
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.25,
                                                    random_state = 42)

### Scaling X_train and X_tests
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

### Import KNeighborsClassifier
from sklearn.neighbors import KNeighborsClassifier
```

```

clf = KNeighborsClassifier()
clf.fit(X_train_scaled, y_train)
y_test_preds = clf.predict(X_test_scaled)

### Results
from sklearn.metrics import classification_report
print(classification_report(y_test, y_test_preds))

find_best_k(X_train_scaled, y_train, X_test_scaled, y_test)

```

	precision	recall	f1-score	support
0	0.83	0.82	0.83	141
1	0.70	0.72	0.71	82
accuracy			0.78	223
macro avg	0.77	0.77	0.77	223
weighted avg	0.79	0.78	0.79	223

best f1 score is: 0.7388535031847134  
best k number is: 17

[146]: 17

### 1.10.2 We drop missing values of age instead of median

```

[128]: df_new_II = raw_df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], axis = 1)
df_new_II.dropna(subset = ["Embarked", "Age"], inplace = True, axis = 0)
df_new_II['Sex'] = df_new_II['Sex'].map({"female":0, "male":1})
df_new_II.isna().sum()

### Using OneHotEncoder for Embarked
from sklearn.preprocessing import OneHotEncoder
OHE = OneHotEncoder(categories='auto', sparse=False, handle_unknown='ignore')

embarked_new_II = OHE.fit_transform(df_new_II[["Embarked"]])
embarked_df_new_II = pd.DataFrame(embarked_new_II,
                                  columns = OHE.categories_[0],
                                  index = df_new_II.index)

df_new_II.drop("Embarked", axis = 1, inplace = True)
one_hot_df_new_II = pd.concat([df_new_II, embarked_df_new_II], axis = 1)
one_hot_df_new_II.head()

#### defining X and y

```

```

y = one_hot_df_new_II["Survived"]
X = one_hot_df_new_II.drop(columns = ["Survived"], axis = 1)

### Creating train and test sets
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.25,
                                                    random_state = 42)

### Scaling X_train and X_tests
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

### Import KNeighborsClassifier
from sklearn.neighbors import KNeighborsClassifier

clf = KNeighborsClassifier()
clf.fit(X_train_scaled, y_train)
y_test_preds = clf.predict(X_test_scaled)

### Results
from sklearn.metrics import classification_report
print(classification_report(y_test, y_test_preds))

find_best_k(X_train_scaled, y_train, X_test_scaled, y_test)

```

	precision	recall	f1-score	support
0	0.78	0.84	0.81	99
1	0.78	0.71	0.74	79
accuracy			0.78	178
macro avg	0.78	0.77	0.78	178
weighted avg	0.78	0.78	0.78	178

best f1 score is: 0.7435897435897437

best k number is: 17

[128]: 17

### 1.10.3 We use media values for age but we use another metric for the distance

```
[145]: ## Data Preperation, we use media values for age but we use another metric

df_new = raw_df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], axis = 1)
df_new["Age"] = df_new["Age"].fillna(df_new["Age"].median())
df_new.dropna( subset = ["Embarked"], inplace = True, axis = 0)
df_new['Sex'] = df_new['Sex'].map({"female":0, "male":1})
df_new.isna().sum()

### Using OneHotEncoder for Embarked
from sklearn.preprocessing import OneHotEncoder
OHE = OneHotEncoder(categories='auto', sparse=False, handle_unknown='ignore')

embarked_new = OHE.fit_transform(df_new[["Embarked"]])
embarked_df_new = pd.DataFrame(embarked_new,
                               columns = OHE.categories_[0],
                               index = df_new.index)

df_new.drop("Embarked", axis = 1, inplace = True)
one_hot_df_new = pd.concat([df_new, embarked_df_new], axis = 1)
one_hot_df_new.head()

#### defining X and y

y = one_hot_df_new["Survived"]
X = one_hot_df_new.drop(columns = ["Survived"], axis = 1)

### Creating train and test sets
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.25,
                                                    random_state = 42)

### Scaling X_train and X_tests
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

### Import KNeighborsClassifier
from sklearn.neighbors import KNeighborsClassifier

clf = KNeighborsClassifier(p=1, metric='minkowski', n_neighbors=17)
clf.fit(X_train_scaled, y_train)
y_test_preds = clf.predict(X_test_scaled)
```

```

### Results
from sklearn.metrics import classification_report
print(classification_report(y_test, y_test_preds))

```

	precision	recall	f1-score	support
0	0.84	0.88	0.86	141
1	0.77	0.71	0.74	82
accuracy			0.82	223
macro avg	0.81	0.79	0.80	223
weighted avg	0.81	0.82	0.81	223

## 1.11 Summary

Well done! In this lab, you worked with the classic Titanic dataset and practiced fitting and tuning KNN classification models using scikit-learn! As always, this gave you another opportunity to continue practicing your data wrangling skills and model tuning skills using Pandas and scikit-learn!