OPTION #1: KNN Classifier with Iris Data

Scott Miner

Colorado State University – Global Campus

Abstract

This paper presents a k-Nearest Neighbors classifier (k-NN-C) implemented in Python, which achieves a mean accuracy of 96.67% on the Iris Dataset using 5-fold cross-validation, with the number of nearest neighbors set to 10. The Iris Dataset, comprising three classes and four attributes, serves as a foundation for understanding classification techniques. The k-NN-C model is a lazy learner that leverages Euclidean distance to identify similarities between observation pairs and makes predictions based on the k nearest neighbors. By using cross-validation, the classifier's performance is assessed on mutually exclusive data splits, ensuring a robust evaluation. The implemented Python program reads the Iris Dataset, preprocesses the data, and applies the k-NN-C algorithm to make predictions. Additionally, the program accepts user input for previously unseen iris plant features and generates class predictions based on the model. This work demonstrates the effectiveness of the k-NN-C model on a widely recognized dataset and lays the groundwork for future research in feature normalization and model optimization.

```
KNN.py
normalize_dataset
 1  import pandas as pd
 2  from sklearn.neighbors import KNeighborsClassifier
 3  from sklearn.model_selection import train_test_split
 4  # k-nearest neighbors on the Iris Flowers Dataset
 5  from random import seed
 6  from random import randrange
 7  from csv import reader
 8  from math import sqrt
 9  import re
10
11  # Load a CSV file
12  def load_csv(filename):
13      dataset = list()
14      with open(filename, 'r') as file:
15          csv_reader = reader(file)
16          for row in csv_reader:
17              if not row:
18                  continue
19              dataset.append(row)
20      return dataset
21
22  # Convert string column to float
23  def str_column_to_float(dataset, column):
24      for row in dataset:
25          row[column] = float(row[column].strip())
26
27  # Convert string column to integer
28  def str_column_to_int(dataset, column):
29      class_values = [row[column] for row in dataset]
30      unique = set(class_values)
31      lookup = dict()
32      # create dictionary
33      for i, value in enumerate(unique):
34          lookup[value] = i
35      # convert dataset column
36      for row in dataset:
37          row[column] = lookup[row[column]]
38      return lookup
39
40  # Find the min and max values for each column
41  def dataset_minmax(dataset):
42      minmax = list()
43      for i in range(len(dataset[0])):
44          col_values = [row[i] for row in dataset]
45          value_min = min(col_values)
46          value_max = max(col_values)
47          minmax.append([value_min, value_max])
48      return minmax
49
50  # Rescale dataset columns to the range 0-1
51  def normalize_dataset(dataset, minmax):
52      for row in dataset:
53          for i in range(len(row)):
54              row[i] = (row[i] - minmax[i][0]) / (minmax[i][1] - minmax[i][0])
55
56  # Split a dataset into k folds
57  def cross_validation_split(dataset, n_folds):
58      dataset_split = list()
59      dataset_copy = list(dataset)
60      fold_size = int(len(dataset) / n_folds)
61      for _ in range(n_folds):
62          fold = list()
63          while len(fold) < fold_size:
64              index = randrange(len(dataset_copy))
65              fold.append(dataset_copy.pop(index))
66          dataset_split.append(fold)
67      return dataset_split
```

*Figure 1.* Python code to implement a *k*-NN-C from scratch (part 1)

```
KNN.py
evaluate_algorithm
 67        return dataset_split
 68
 69   # Calculate accuracy percentage
 70   def accuracy_metric(actual, predicted):
 71        correct = 0
 72        for i in range(len(actual)):
 73            if actual[i] == predicted[i]:
 74                correct += 1
 75        return correct / float(len(actual)) * 100.0
 76
 77   # Evaluate an algorithm using a cross validation split
 78   def evaluate_algorithm(dataset, algorithm, n_folds, *args):
 79        folds = cross_validation_split(dataset, n_folds)
 80        scores = list()
 81        for fold in folds:
 82            train_set = list(folds)
 83            # create hold out set
 84            train_set.remove(fold)
 85            #combine train sets
 86            train_set = sum(train_set, [])
 87            # create test set on new hold
 88            test_set = list()
 89            for row in fold:
 90                row_copy = list(row)
 91                test_set.append(row_copy)
 92                # remove prediction from hold out set
 93                row_copy[-1] = None
 94            predicted = algorithm(train_set, test_set, *args)
 95            actual = [row[-1] for row in fold]
 96            accuracy = accuracy_metric(actual, predicted)
 97            scores.append(accuracy)
 98        return scores
 99
100   # Calculate the Euclidean distance between two vectors
101   def euclidean_distance(row1, row2):
102        distance = 0.0
103        for i in range(len(row1) - 1):
104            distance += (row1[i] - row2[i]) ** 2
105        return sqrt(distance)
106
107   # Locate the most similar neighbors
108   def get_neighbors(train, test_row, num_neighbors):
109        distances = list()
110        for train_row in train:
111            dist = euclidean_distance(test_row, train_row)
112            distances.append((train_row, dist))
113        distances.sort(key=lambda tup: tup[1])
114        neighbors = list()
115        for i in range(num_neighbors):
116            neighbors.append(distances[i][0])
117        return neighbors
118
119   # Make a prediction with neighbors
120   def predict_classification(train, test_row, num_neighbors):
121        neighbors = get_neighbors(train, test_row, num_neighbors)
122        output_values = [row[-1] for row in neighbors]
123        prediction = max(set(output_values), key=output_values.count)
124        return prediction
125
126   # kNN Algorithm
127   def k_nearest_neighbors(train, test, num_neighbors):
128        predictions = list()
129        for row in test:
130            output = predict_classification(train, row, num_neighbors)
131            predictions.append(output)
132        return(predictions)
```

*Figure 2.* Python code to implement a *k*-NN-C from scratch (part 2)

```
133
134    # Test the kNN on the Iris Flowers dataset
135    seed(1)
136    filename = 'data/iris.txt'
137    dataset = load_csv(filename)
138    for i in range(len(dataset[0]) - 1):
139        str_column_to_float(dataset[1:], i)
140    # convert class column to integers
141    # versicolor : 0
142    # virginica: 1
143    # setosa: 2
144    lookup = str_column_to_int(dataset[1:], len(dataset[0]) - 1)
145
146    # evaluate algorithm
147    n_folds = 5
148    num_neighbors = 10
149    scores = evaluate_algorithm(dataset[1:], k_nearest_neighbors, n_folds, num_neighbors)
150    print(f'********************************************************************************')
151    print(f'*')
152    print(f'*   K-Nearest Neighbor (KNN) algorithm with {num_neighbors} neighbors trained on a ')
153    print(f'*   dataset containing  {len(dataset)-1} rows and {len(dataset[1])-1} features, using {n_folds}-fold cross validation.')
154    print(f'*')
155    print(f'*   Users can adjust the n_folds and num_neighbors variables in the script.')
156    print(f'*')
157    print(f'********************************************************************************')
158    print()
159    print(f'Accuracy per fold: {scores}')
160    print(f'Mean Accuracy: {sum(scores) / float(len(scores)):.3f}')
161
162    while True:
163        try:
164            sl, sw, pl, pw = [float(x) for x in re.split(r'\, | |\,', input('\nPlease input four floating point numbers, representing, respectively,\n\
165    sepal length, sepal width, petal length, and petal width \n(e.g., 5.1, 3.5, 1.4, 0.2).  The model will guess the flower category\n\
166    (i.e., setosa, versicolor, or virginica) based on your input (CTRL-C to Exit): '))]
167            test_row = [sl, sw, pl, pw]
168            test_row
169            prediction = predict_classification(train=dataset[1:], test_row=test_row, num_neighbors=num_neighbors)
170            for key, value in lookup.items():
171
172                if prediction == value:
173                    prediction = key
174            print()
175            print(f'Prediction: {prediction}')
176        except ValueError:
177            print('\nNote: wrong input format.')
178
```

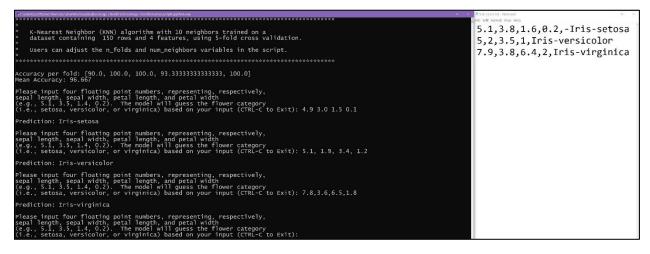*Figure 3.* Python code to implement a *k*-NN-C from scratch (part 3)



*Figure 4.* Program output and predictions based on user input similar to the training examples found on the right-hand side of the image

## Table of Contents

**OPTION #1: KNN Classifier with Iris Data**

This paper presents a k-Nearest Neighbors classifier (k-NN-C) built from scratch in Python, achieving a notable mean accuracy of 96.67% on the renowned Iris Dataset using 5-fold cross-validation (CV). The model's number of nearest neighbors is set to 10. Fenner (2019) discusses the Iris Dataset, often referred to as Fisher's Iris Dataset, named after the eminent mid-20th-century statistician, Sir Ronald Fisher, who pioneered its use in one of the earliest academic papers on classification. The dataset comprises three distinct classes, with 50 instances per class, representing the types of iris plants: (a) Setosa, (b) Versicolor, and (c) Virginica. Four descriptive attributes measure these plants in centimeters: (a) sepal length, (b) sepal width, (c) petal length, and (d) petal width.

**User Input and Model Predictions**

Upon training the k-NN-C model on the Iris Dataset, the program invites users to input four floating-point numbers, separated by commas or spaces, representing the features of previously unencountered iris plants. Subsequently, the program generates the model's class prediction based on this user input, demonstrating its practical utility in real-world scenarios.

**Introduction to the k-Nearest Neighbor Classifier (k-NN-C)**

Fenner (2019) highlights k-NN-C as a relatively simple yet effective machine-learning model designed to make predictions using labeled datasets. At its core, k-NN-C assesses similarities between pairs of observations, selects a predefined number of the most similar instances, and combines these findings to generate a single output prediction. While the Euclidean distance is frequently employed to measure similarities among features, alternative metrics like Minkowski and Hamming distances can also be utilized. In k-NN-C, the variable 'k'

denotes the number of nearest neighbors that the model relies on for its predictions. Typically, practitioners experiment with values such as 1, 3, 10, and 20 to find the optimal setting.

*Understanding the k-Nearest Neighbor Classifier (k-NN-C) Algorithm and Its Unique Characteristics*

Fenner (2019) emphasizes the unique characteristics of k-NN-C models in comparison to other machine learning approaches, particularly their reliance on the entirety of the training data when making predictions for new test cases. Consequently, removing any training observations could result in inaccurate predictions, as these records might have been crucial in determining the nearest neighbors for specific test cases. This property classifies k-NN-C as a lazy learner, which, unlike eager learners such as logistic regression, does not have a dedicated training phase. Instead, it retains all training data for use during the prediction phase, making this process more computationally intensive than that of eager learners (*KNN Algorithm - Finding Nearest Neighbors*, n.d.; *Why Is Nearest Neighbor a Lazy Algorithm?*, 2021).

**Program and Model Architecture**

Figures 1 – 3 illustrate the Python code used to create the k-NN-C model, with much of the code adapted from Brownlee (2019). The program starts by reading the Iris Dataset from the data/iris.txt file using the *load_csv()* function (lines 12 – 20, Figure 1). To ensure the dataset's features are in the correct format, *str_column_to_float()* (lines 28 – 38, Figure 1) converts the features from strings to floating-point numbers, while *str_column_to_int()* (lines 23 – 35, Figure 1) converts the string representation of each iris plant category to integers for easier processing by the k-NN-C model.

**Cross-Validation Process**

The *cross_validation_split()* method (lines 56 – 67, Figure 1) divides the dataset into five cross-validation (CV) splits. Barrow and Crone (2013) describe CV as a technique to estimate the expected accuracy of a predictive algorithm by averaging predictive errors across mutually exclusive subsamples of the data (p. 1). The function divides the dataset into k mutually exclusive groups of roughly equal size, as determined by the user-defined k value. One fold is reserved for the validation dataset, while the remaining folds form the training data. The model is trained on the training data and evaluated using the validation dataset (Brownlee, 2019).

Barrow and Crone (2013) explain that the CV process is repeated k times, with each fold serving as the validation dataset exactly once. The program records the accuracy scores for each of the k evaluation sessions and calculates the final measure of the model's predictive accuracy by averaging these scores across the k folds. The accuracy_metric() function (lines 70 – 75, Figure 2) calculates the model's accuracy for each fold using the formula: $\frac{\text{\# correct predictions}}{\text{\# total predictions}}$. The advantages of *k*-fold CV are that it uses all observations for both training and validation datasets, uses all training observations with equal weight, and uses each observation for validation exactly once.

**Generating Predictions and Calculating Euclidean Distance**

Since k-NN-C is a lazy learning algorithm, it doesn't have a specialized training phase but instead uses all its training data to generate predictions during the evaluation phase. The *evaluate_algorithm()* method (lines 78 – 98, Figure 2) calls *k_nearest_neighbors()*, which iterates over the test dataset and calls *predict_classification()* for each row of test data. *Predict_classification()* then calls *get_neighbors()*, which in turn calls the *euclidean_distance()* function (lines 101 – 105, Figure 2) to calculate the Euclidean distance between each training example and a given test example using the formula:

$$d = \sqrt{(sl_{train} - sl_{test})^2 + (sw_{train} - sw_{test})^2 + (pl_{train} - pl_{test})^2 + (pw_{train} - pw_{test})^2}.$$

The variables sl, sw, pl, and pw represent sepal length, sepal width, petal length, and petal width

for each training and test example (Brownlee, 2019).

**Model Evaluation, User Input, and Future Research**

The k-NN-C model achieves a mean accuracy of 96.67% on the Iris Dataset using 5-fold

CV, with the number of nearest neighbors set to 10 (Figure 4). The program then allows users to

input four floating-point numbers representing the sepal length, sepal width, petal length, and

petal width of previously unseen iris plants. For example, by inputting features similar to, but not

identical with, the training examples shown in Figure 4, the model accurately predicts each iris

plant category. The driver code for user input is shown in lines 134 – 177 of Figure 3. Users

must input four valid numerical features, separated by spaces or commas, for the model to

generate a valid prediction. If the input is incorrect, the program notifies the user and prompts for

new input. Additional functions, like the *normalize_dataset()* function, can be used in future

research to assess how normalizing the input features of the Iris Dataset affects the classifier's

mean accuracy (Brownlee, 2019).

<div align="center">

**Conclusion**

</div>

In conclusion, this paper presented an overview of a *k*-Nearest Neighbor classifier (k-

NN-C) built from scratch in Python, achieving 96.67% accuracy on the Iris Dataset using 5-fold

cross-validation (CV), with the model's number of neighbors set to 10. The paper provided an

overview of the Iris Dataset, k-NN-C models, and CV. The classifier's inner workings were also

explained, including the calculation of Euclidean distance for each training example per test

example, sorting of results, and summation of frequencies of the *k* nearest neighbor classes to

produce predictions. The program also accepts user input for iris plant features previously unseen

by the model and generates predictions based on the training data.

References

Barrow, D. K., & Crone, S. F. (2013). Crogging (cross-validation aggregation) for forecasting—

    A novel algorithm of neural network ensembles on time series subsamples. *The 2013*

    *International Joint Conference on Neural Networks (IJCNN)*, 1–8.

Brownlee, J. (2019, October 23). Develop k-Nearest Neighbors in Python from Scratch. *Machine*

    *Learning Mastery*. https://machinelearningmastery.com/tutorial-to-implement-k-nearest-

    neighbors-in-python-from-scratch/

Fenner, M. (2019). *Machine Learning in Python for Everyone*. Addison-Wesley.

*KNN Algorithm—Finding Nearest Neighbors*. (n.d.). Retrieved August 28, 2021, from

    https://www.tutorialspoint.com/machine_learning_with_python/machine_learning_with_

    python_knn_algorithm_finding_nearest_neighbors.htm

*Why is Nearest Neighbor a Lazy Algorithm?* (2021, August 25). Dr. Sebastian Raschka.

    https://sebastianraschka.com/faq/docs/lazy-knn.html