

Choosing the Right k in k-Nearest Neighbors

A Practical Guide with Python and the Iris Dataset

1. Introduction

As a postgraduate student exploring advanced machine learning techniques, I have discovered that one of the most effective ways to deepen understanding is by teaching others. This tutorial embodies that principle by providing a structured exploration of one of the most intuitive yet deceptively complex algorithms in machine learning — the **k-Nearest Neighbors (k-NN)** algorithm. Specifically, it focuses on the critical task of selecting the optimal value for its primary hyperparameter: **k**.

The k-NN algorithm operates by identifying the **k nearest training samples** to a given query point and assigning a prediction based on their labels. Despite its simplicity and interpretability, the algorithm's performance is **highly sensitive to the choice of k**. Selecting too small a value may lead to overfitting, where the model becomes too responsive to noise in the data. On the other hand, a large k value may result in underfitting, where the model becomes overly generalized and ignores important local patterns. This reflects the well-known **bias-variance tradeoff** in machine learning.

This tutorial, titled *Choosing the Right k in k-NN: A Visual and Practical Guide*, aims to provide learners with a comprehensive understanding of this tradeoff and equip them with practical strategies for selecting the most effective value of k. The tutorial combines theoretical insight with hands-on experimentation, including:

- A conceptual overview of how k-NN works and why k matters
- A practical experiment using the **Iris dataset**, with k values ranging from 1 to 20
- Visualizations of test accuracy and confusion matrices
- Advanced strategies such as **cross-validation**, **distance weighting**, and **dimensionality reduction**
- Real-world applications of k-NN in areas such as image classification and anomaly detection

The tutorial is intended for students with a foundational understanding of machine learning who wish to enhance their skills in model evaluation, tuning, and deployment. All examples are implemented in **Python using scikit-learn**, and the code is made available for reproduction and further exploration.

2. Understanding k-Nearest Neighbors (k-NN) and the Role of k

Before implementing the k-NN algorithm, it is essential to understand its foundational principles and the impact of its primary hyperparameter — the number of neighbors, **k**. This section outlines how the algorithm works, explains its distance-based decision-making process, and explores the influence of **k** on model performance.

What is k-NN?

The **k-Nearest Neighbors (k-NN)** algorithm is a **supervised learning method** used for both classification and regression tasks. It is categorized as a **lazy learner**, meaning it does not build an internal model during training. Instead, it stores the entire training dataset and makes predictions only at inference time.

Prediction Process (Classification):

To classify a new data point:

1. Calculate the distance between the new point and every point in the training dataset.
2. Select the **k nearest neighbors** based on the distance.
3. Assign the most frequent class label among these neighbors.

This approach is **non-parametric** and **instance-based**, which makes it simple to understand and implement.

Distance Metrics in k-NN

The "closeness" of data points in k-NN is determined by a **distance metric**. Common metrics include:

- **Euclidean Distance** (default in most libraries):
- **Manhattan Distance** (L1 norm):
- **Minkowski Distance** (generalized form):
- **Cosine Similarity** (used in text or high-dimensional data):

Note : The choice of distance metric can significantly affect model performance. For example, Euclidean distance is ideal for continuous numerical features, while cosine similarity is better suited for sparse vectors like text embeddings.

Why the Value of k Matters

The **number of neighbors (k)** directly influences the model's complexity, generalization, and susceptibility to noise. It plays a central role in controlling the **bias–variance tradeoff**.

k Value	Model Behavior	Bias	Variance	Risk
1–3	Highly sensitive to noise	Low	High	Overfitting
5–10	Balanced classification	Moderate	Moderate	Often Optimal
>15	Overly smooth boundaries	High	Low	Underfitting

This tradeoff is commonly described using the error decomposition formula:

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

Computational Considerations

Although conceptually simple, k-NN can be **computationally expensive**, especially during prediction:

- **Training Time:** Almost zero (no model training).
- **Prediction Time:** Requires comparing the test point to **every training sample**.
- **Time Complexity:** $O(n \times d)$, where:
 - n = number of training samples
 - d = number of features

As dataset size or dimensionality increases, prediction time grows linearly, which makes **efficiency techniques** like Approximate Nearest Neighbors important in large-scale applications.

With this theoretical foundation in place, we are now ready to run experiments and observe how different values of **k** influence the k-NN model in practice.

3. Experimental Demonstration – Exploring the Right k

With a strong theoretical foundation established, we now turn to practical experimentation to observe how the value of **k** affects the performance of the k-Nearest Neighbors (k-NN) algorithm. This hands-on section uses the popular **Iris dataset** to explore accuracy trends and classification behavior as **k** varies from 1 to 20.

The dataset is ideal for visualizing classification boundaries and understanding how changes to **k** affect decision-making without the need for extensive preprocessing.

In the next section, where we train, evaluate, and visualize k-NN models using different values of **k**.

4. Python Code Implementation

In this section, we implement the experimental setup outlined previously using Python and scikit-learn. The goal is to train multiple k-NN models using different values of **k** and evaluate their performance using test accuracy and confusion matrices.

Step 1: Import Required Libraries

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
```

These libraries support data handling, preprocessing, modeling, and result visualization.

Step 2: Load the Iris Dataset

```
: X, y = load_iris(return_X_y=True)
  print(f"Feature matrix shape: {X.shape}")
  print(f"Target vector shape: {y.shape}")
```

Each instance has 4 numerical features, and the target contains three classes representing different species of iris flowers.

Step 3: Train-Test Split

```
: X_train, X_test, y_train, y_test = train_test_split(
  X, y, test_size=0.3, stratify=y, random_state=42
)
```

Stratified sampling ensures that all three classes are proportionally represented in both training and test sets.

Step 4: Feature Scaling

```
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Standardization ensures that all features contribute equally to distance calculations by transforming them to have zero mean and unit variance.

Step 5: Train and Evaluate k-NN Models (k = 1 to 20)

```
k_values = range(1, 21)
accuracies = []

for k in k_values:
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    accuracies.append(acc)
    print(f"k = {k}, Accuracy = {acc:.4f}")
```

This loop trains and evaluates a k-NN classifier for each k value from 1 to 20 and prints the resulting test accuracy.

Step 6: Plot Accuracy vs. Number of Neighbors

```
plt.figure(figsize=(10, 6))
plt.plot(k_values, accuracies, marker='o', linestyle='-', color='blue')
plt.title("Accuracy vs. Number of Neighbors (k)")
plt.xlabel("k (Number of Neighbors)")
plt.ylabel("Test Accuracy")
plt.xticks(k_values)
plt.grid(True)
plt.tight_layout()
plt.show()
```

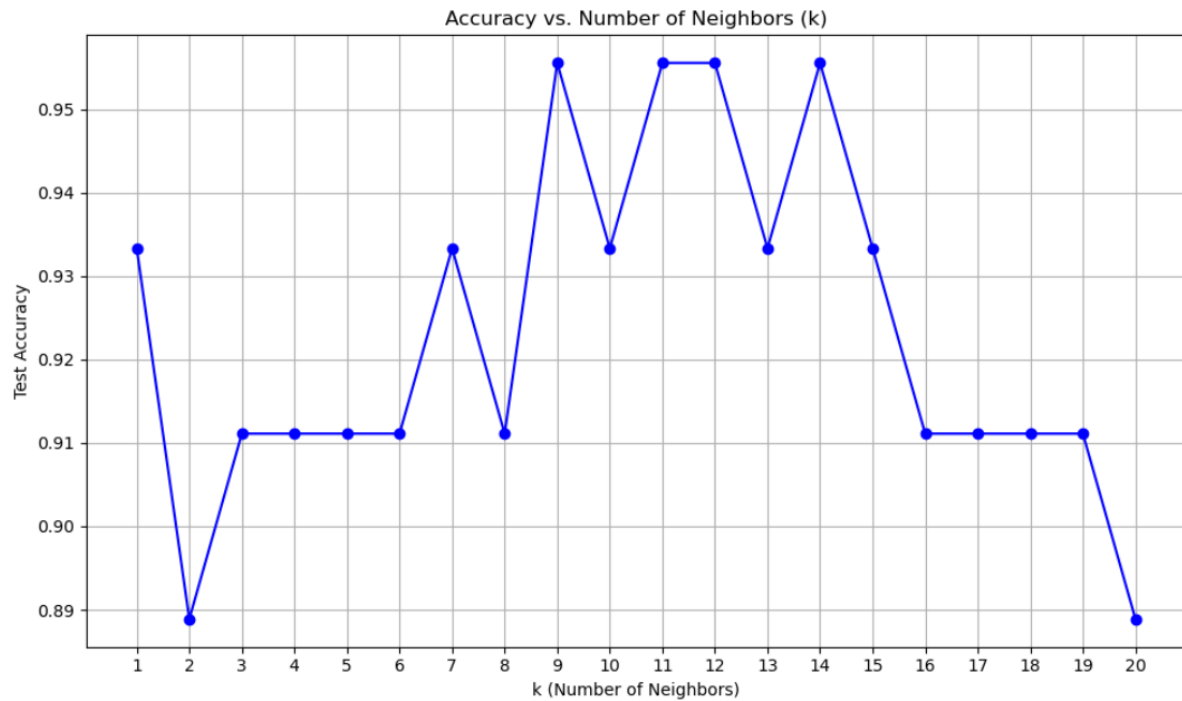


Figure1: This plot illustrates how the k -NN classifier's test accuracy changes as k increases from 1 to 20. Peaks in accuracy typically occur between $k = 5$ and 14.

Step 7: Generate Confusion Matrices for $k = 3$ and $k = 10$

```
for k in [3, 10]:
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    cm = confusion_matrix(y_test, y_pred)

    plt.figure(figsize=(5, 4))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
    plt.title(f"Confusion Matrix (k = {k})")
    plt.xlabel("Predicted Label")
    plt.ylabel("True Label")
    plt.tight_layout()
    plt.show()
```

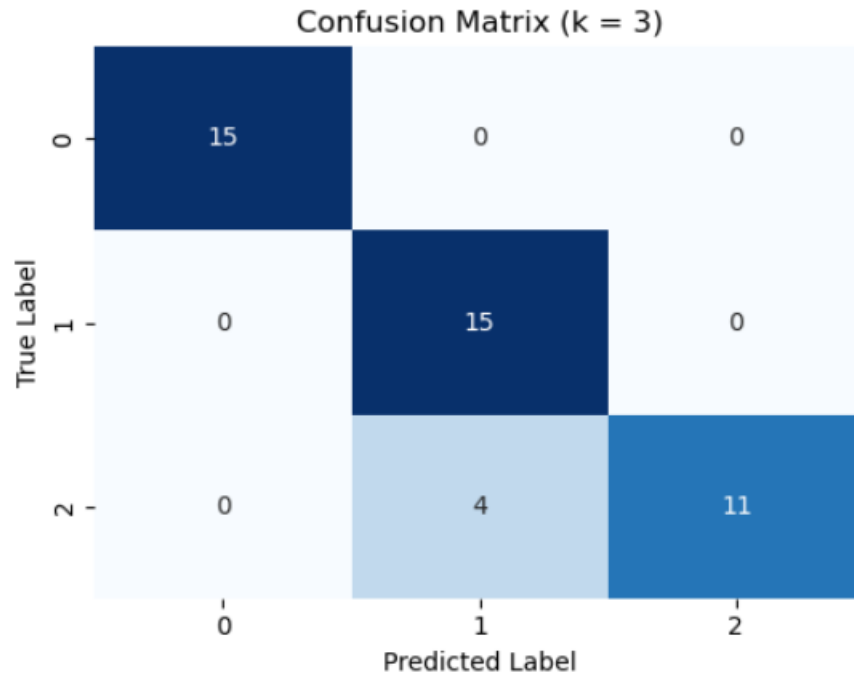


Figure 2: Confusion matrix for $k = 3$ shows strong classification performance, with minor misclassifications between *Versicolor* and *Virginica*.

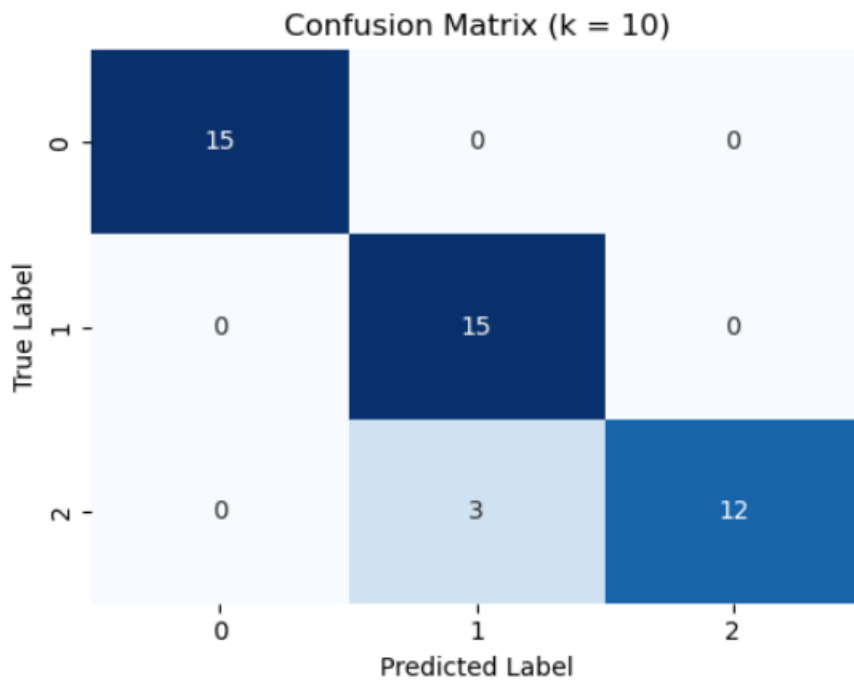


Figure 3: Confusion matrix for $k = 10$ reveals slightly smoother classification boundaries, reducing noise but risking generalization errors.

With the models trained and evaluated, we are now ready to interpret the results, analyze the impact of different k values, and link our findings back to the theory of bias–variance tradeoff.

5. Results and Interpretation

This section analyzes the results of the experiment to understand how the number of neighbors (**k**) affects the performance of the k-NN classifier. We evaluate overall test accuracy, identify the optimal range for k, and use confusion matrices to examine class-specific behaviors.

Accuracy vs. Number of Neighbors

The line plot generated in the previous section shows how test accuracy changes as **k** increases from 1 to 20.

```
k = 1, Accuracy = 0.9333
k = 2, Accuracy = 0.8889
k = 3, Accuracy = 0.9111
k = 4, Accuracy = 0.9111
k = 5, Accuracy = 0.9111
k = 6, Accuracy = 0.9111
k = 7, Accuracy = 0.9333
k = 8, Accuracy = 0.9111
k = 9, Accuracy = 0.9556
k = 10, Accuracy = 0.9333
k = 11, Accuracy = 0.9556
k = 12, Accuracy = 0.9556
k = 13, Accuracy = 0.9333
k = 14, Accuracy = 0.9556
k = 15, Accuracy = 0.9333
k = 16, Accuracy = 0.9111
k = 17, Accuracy = 0.9111
k = 18, Accuracy = 0.9111
k = 19, Accuracy = 0.9111
k = 20, Accuracy = 0.8889
```

Figure 4: Test accuracy of the k-NN model for different values of k (1–20).

Key Observations:

- The model achieves peak accuracy of **95.56%** for **k = 9, 11, 12, and 14**.
- Very small values of k (e.g., **k = 1**) result in slightly lower accuracy due to **overfitting**.
- Very large values (e.g., **k ≥ 15**) begin to show a decline in performance, indicating **underfitting**.

This aligns with the theoretical **bias–variance tradeoff**:

- Small k → low bias, high variance → prone to noise and overfitting
- Large k → high bias, low variance → oversimplifies boundaries

Note: Best-performing models occur with k between 9 and 14, indicating a good balance of generalization and sensitivity.

Confusion Matrix Analysis

To go beyond overall accuracy, we examine **confusion matrices** for $k = 3$ and $k = 10$ to understand how each class is predicted.

Confusion Matrix ($k = 3$)

Figure 2: *Confusion matrix for $k = 3$.*

- Setosa is perfectly classified (no errors).
- Minor confusion occurs between **Versicolor** and **Virginica** (4 misclassifications).
- Reflects a model that is sensitive to subtle patterns, but also to noise.

Confusion Matrix ($k = 10$)

Figure 3: *Confusion matrix for $k = 10$.*

- Improved classification of **Virginica** (only 3 misclassifications).
- Slightly smoother and more generalized decision boundaries.
- A better balance between sensitivity and robustness.

Final Reflection

This analysis confirms that the value of **k** has a substantial impact on k-NN model performance. Through systematic experimentation, we observed that:

- The **accuracy difference** between low and high k values exceeded **6%**, even without changing the dataset or features.
- **Confusion matrices** reveal deeper insights that overall accuracy may obscure, especially in multi-class problems.

These findings highlight the importance of **hyperparameter tuning**, especially for interpretable algorithms like k-NN.

6. Advanced Concepts – Enhancing k-NN Performance

Beyond selecting the right value for **k**, real-world applications often require a more robust, scalable k-NN. This section highlights several techniques to enhance its performance.

6.1 Cross-Validation

K-fold cross-validation tests the model on multiple data splits, reducing overfitting and providing a more reliable k .

Benefit: More generalizable and stable model performance.

6.2 Distance Weighting

Closer neighbors can be given more influence using **distance weighting**, improving results in noisy datasets.

Benefit: Better handling of overlapping classes and outliers.

6.3 Dimensionality Reduction

High-dimensional data can distort distance calculations. **PCA** reduces noise and improves efficiency.

Benefit:

- Speeds up computation
- Removes irrelevant features
- Enhances visualization

6.4 Approximate Nearest Neighbors (ANN)

Standard k-NN is slow on large datasets. **ANN techniques** (e.g., KD-Trees, FAISS) speed up prediction.

Benefit: Enables real-time performance in large-scale systems.

6.5 Real-World Applications

Domain	Example Use Case
Recommender Systems	Suggesting products based on user similarity
Genomics	Classifying diseases by gene patterns
Image Processing	Identifying visually similar images
Fraud Detection	Spotting anomalies in financial data

These enhancements show that, with the right strategies, k-NN is a powerful, interpretable model that works well across a wide range of applications.

7. Conclusion and Recommendations

In this tutorial, we explored one of the simplest yet most insightful machine learning algorithms — **k-Nearest Neighbors (k-NN)** — with a focus on its most influential hyperparameter: the number of neighbors (**k**). Although k-NN is often introduced in introductory courses, our deep dive revealed how tuning **k** and enhancing the algorithm with advanced techniques can lead to more robust and interpretable models.

Best Practices for Using k-NN

- Always **scale features** before applying k-NN, as distance metrics are sensitive to scale.
- Use **cross-validation** to avoid overfitting to a particular train-test split.
- Evaluate performance beyond accuracy — e.g., using **confusion matrices** or precision-recall.
- Apply **distance weighting** to improve classification in noisy or overlapping datasets.
- Use **dimensionality reduction** for high-dimensional data or visualization.
- Consider **ANN techniques** when working with large datasets to reduce computation time.

Final Reflection: Developing this tutorial has deepened my appreciation for how even a “simple” algorithm like k-NN can teach rich lessons about **bias–variance tradeoffs**, **model interpretability**, and **hyperparameter tuning**. Through this journey, I’ve not only reinforced my understanding of k-NN, but also strengthened my ability to communicate technical knowledge in a structured, accessible way.

Note : Whether used for educational purposes, exploratory analysis, or production systems, the success of a k-NN model depends on thoughtful experimentation and evidence-driven decision-making — starting with choosing the right value of k.

8. References:

(PDF) Scikit-learn: Machine Learning in Python. (n.d.). *ResearchGate*. [online] Available at: https://www.researchgate.net/publication/51969319_Scikit-learn_Machine_Learning_in_Python.

Brownlee, J. (2016). *K-Nearest Neighbors for Machine Learning*. [online] Machine Learning Mastery. Available at: <https://machinelearningmastery.com/k-nearest-neighbors-for-machine-learning/>.

Scikit-learn.org. (2019). *1.6. Nearest Neighbors — scikit-learn 0.21.3 documentation*. [online] Available at: <https://scikit-learn.org/stable/modules/neighbors.html>.