# k-Nearest Neighbors: Theory and Practice

September 10, 2025

## 1 Introduction

k-Nearest Neighbors (k-NN) is a non-parametric, instance-based learner: predictions are made by looking up the closest training samples. It is simple and often competitive on low-dimensional, well-scaled data, but can be sensitive to feature scaling and suffers in high dimensions.

## 2 Theory and Formulas

Given a query $\mathbf{x}$, find its $k$ nearest neighbors under a distance metric $d(\cdot, \cdot)$ (e.g., Euclidean, Manhattan). For classification, predict by majority vote among the labels of these neighbors; with distance weighting (`weights=distance`), closer neighbors have higher influence. For regression, predict the average (or distance-weighted average) of neighbor targets.

Computationally, naive search costs $\mathcal{O}(nd)$ per query with $n$ samples and $d$ features. Tree-based indices (KDTree/BallTree) can accelerate queries in moderate dimensions. k-NN is affected by the curse of dimensionality; proper feature scaling and metric choice are critical.

## 3 Applications and Tips

- **Choose k:** tune via cross-validation; odd $k$ helps avoid ties in binary classification.

- **Scaling:** standardize features or use pipelines; scale-sensitive.

- **Metric:** try Euclidean vs Manhattan; consider domain-specific distances.

- **Weights:** `uniform` vs `distance`; weighting can help with class overlap.

- **Complexity:** prediction cost grows with data size; consider approximate neighbors for large datasets.

## 4 Python Practice

Run the script in this chapter directory to generate figures into `figures/`.

Listing 1: Generate k-NN figures

```
python gen_knn_figures.py
```

Listing 2: gen_knn_figures.py

```python
"""
Figure generator for the k-NN chapter.

Generates illustrative figures and saves them into the chapter's 'figures/'
folder next to this script, regardless of current working directory.

Requirements:
- Python 3.8+
- numpy, matplotlib, scikit-learn

Install (if needed):
    pip install numpy matplotlib scikit-learn

This script avoids newer or experimental APIs for broader compatibility.
"""
from __future__ import annotations

import os
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap

try:
    from sklearn.datasets import make_moons, make_regression,
        make_classification
    from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
    from sklearn.preprocessing import StandardScaler
    from sklearn.pipeline import make_pipeline
    from sklearn.model_selection import cross_val_score
except Exception:
    raise SystemExit(
        "Missing scikit-learn. Please install with: pip install scikit-learn"
    )


def _ensure_figures_dir(path: str | None = None) -> str:
    """Create figures directory under this chapter regardless of CWD."""
    if path is None:
        base = os.path.dirname(os.path.abspath(__file__))
        path = os.path.join(base, "figures")
    os.makedirs(path, exist_ok=True)
    return path


def _plot_decision_boundary(ax, clf, X, y, title: str):
    x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
    y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
    xx, yy = np.meshgrid(
        np.linspace(x_min, x_max, 400), np.linspace(y_min, y_max, 400)
    )
    Z = clf.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
    cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
```

```python
52          cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
53          ax.contourf(xx, yy, Z, cmap=cmap_light, alpha=0.8, levels=np.unique(Z).
                size)
54          ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s=20)
55          ax.set_title(title)
56          ax.set_xlabel("Feature 1")
57          ax.set_ylabel("Feature 2")
58
59
60      def fig_knn_k_compare(out_dir: str) -> str:
61          np.random.seed(0)
62          X, y = make_moons(n_samples=500, noise=0.3, random_state=0)
63          models = [
64              (KNeighborsClassifier(n_neighbors=1), "k=1 (high variance)"),
65              (KNeighborsClassifier(n_neighbors=15), "k=15 (smoother)")
66          ]
67          fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
                sharey=True)
68          for ax, (m, title) in zip(axes, models):
69              m.fit(X, y)
70              _plot_decision_boundary(ax, m, X, y, f"k-NN: {title}")
71          fig.suptitle("Effect of k on decision boundary")
72          out_path = os.path.join(out_dir, "knn_k_compare.png")
73          fig.tight_layout(rect=[0, 0.03, 1, 0.95])
74          fig.savefig(out_path)
75          plt.close(fig)
76          return out_path
77
78
79      def fig_knn_metric_compare(out_dir: str) -> str:
80          np.random.seed(1)
81          X, y = make_moons(n_samples=500, noise=0.28, random_state=1)
82          models = [
83              (KNeighborsClassifier(n_neighbors=11, metric="euclidean"), "metric=
                    euclidean"),
84              (KNeighborsClassifier(n_neighbors=11, metric="manhattan"), "metric=
                    manhattan"),
85          ]
86          fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
                sharey=True)
87          for ax, (m, title) in zip(axes, models):
88              m.fit(X, y)
89              _plot_decision_boundary(ax, m, X, y, f"k-NN: {title}")
90          fig.suptitle("Effect of distance metric")
91          out_path = os.path.join(out_dir, "knn_metric_compare.png")
92          fig.tight_layout(rect=[0, 0.03, 1, 0.95])
93          fig.savefig(out_path)
94          plt.close(fig)
95          return out_path
96
97
98      def fig_knn_scaling_effect(out_dir: str) -> str:
99          np.random.seed(2)
100         X, y = make_classification(
```

```python
        n_samples=600,
        n_features=2,
        n_informative=2,
        n_redundant=0,
        n_clusters_per_class=1,
        class_sep=1.0,
        random_state=2,
    )
    # Impose different scales on features
    X_scaled_variance = X.copy()
    X_scaled_variance[:, 0] *= 8.0  # make feature 0 dominate distances

    knn_raw = KNeighborsClassifier(n_neighbors=11)
    knn_std = make_pipeline(StandardScaler(), KNeighborsClassifier(n_neighbors
        =11))

    fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
        sharey=True)
    knn_raw.fit(X_scaled_variance, y)
    _plot_decision_boundary(axes[0], knn_raw, X_scaled_variance, y, "Without
        scaling")
    knn_std.fit(X_scaled_variance, y)
    _plot_decision_boundary(axes[1], knn_std, X_scaled_variance, y, "With
        StandardScaler")
    fig.suptitle("Feature scaling impact on k-NN")
    out_path = os.path.join(out_dir, "knn_scaling_effect.png")
    fig.tight_layout(rect=[0, 0.03, 1, 0.95])
    fig.savefig(out_path)
    plt.close(fig)
    return out_path


def fig_knn_weight_compare(out_dir: str) -> str:
    np.random.seed(3)
    X, y = make_moons(n_samples=500, noise=0.32, random_state=3)
    models = [
        (KNeighborsClassifier(n_neighbors=11, weights="uniform"), "weights=
            uniform"),
        (KNeighborsClassifier(n_neighbors=11, weights="distance"), "weights=
            distance"),
    ]
    fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
        sharey=True)
    for ax, (m, title) in zip(axes, models):
        m.fit(X, y)
        _plot_decision_boundary(ax, m, X, y, f"k-NN: {title}")
    fig.suptitle("Uniform vs distance weighting")
    out_path = os.path.join(out_dir, "knn_weight_compare.png")
    fig.tight_layout(rect=[0, 0.03, 1, 0.95])
    fig.savefig(out_path)
    plt.close(fig)
    return out_path
```

4

```python
148  def fig_knn_regression_curve(out_dir: str) -> str:
149      rng = np.random.RandomState(4)
150      # 1D regression: y = sin(x) + noise
151      X = np.sort(rng.uniform(-3.0, 3.0, size=150)).reshape(-1, 1)
152      y = np.sin(X).ravel() + rng.normal(scale=0.25, size=X.shape[0])
153
154      grid = np.linspace(-3.5, 3.5, 600).reshape(-1, 1)
155      models = [
156          (KNeighborsRegressor(n_neighbors=1), "k=1"),
157          (KNeighborsRegressor(n_neighbors=15), "k=15"),
158          (KNeighborsRegressor(n_neighbors=45), "k=45"),
159      ]
160      fig, ax = plt.subplots(figsize=(7.5, 4.2), dpi=160)
161      ax.scatter(X[:, 0], y, s=18, c="#555", alpha=0.7, label="data")
162      colors = ["#E74C3C", "#3498DB", "#2ECC71"]
163      for (m, title), col in zip(models, colors):
164          m.fit(X, y)
165          y_pred = m.predict(grid)
166          ax.plot(grid[:, 0], y_pred, color=col, lw=2, label=title)
167      ax.set_title("k-NN regression: smoothing vs k")
168      ax.set_xlabel("x")
169      ax.set_ylabel("y")
170      ax.legend()
171      ax.grid(True, linestyle=":", alpha=0.4)
172      out_path = os.path.join(out_dir, "knn_regression_curve.png")
173      fig.tight_layout()
174      fig.savefig(out_path)
175      plt.close(fig)
176      return out_path
177
178
179  def main():
180      out_dir = _ensure_figures_dir(None)
181      generators = [
182          fig_knn_k_compare,
183          fig_knn_metric_compare,
184          fig_knn_scaling_effect,
185          fig_knn_weight_compare,
186          fig_knn_regression_curve,
187      ]
188      print("Generating figures into:", os.path.abspath(out_dir))
189      for gen in generators:
190          try:
191              p = gen(out_dir)
192              print("Saved:", p)
193          except Exception as e:
194              print("Failed generating", gen.__name__, ":", e)
195
196
197  if __name__ == "__main__":
198      main()
```
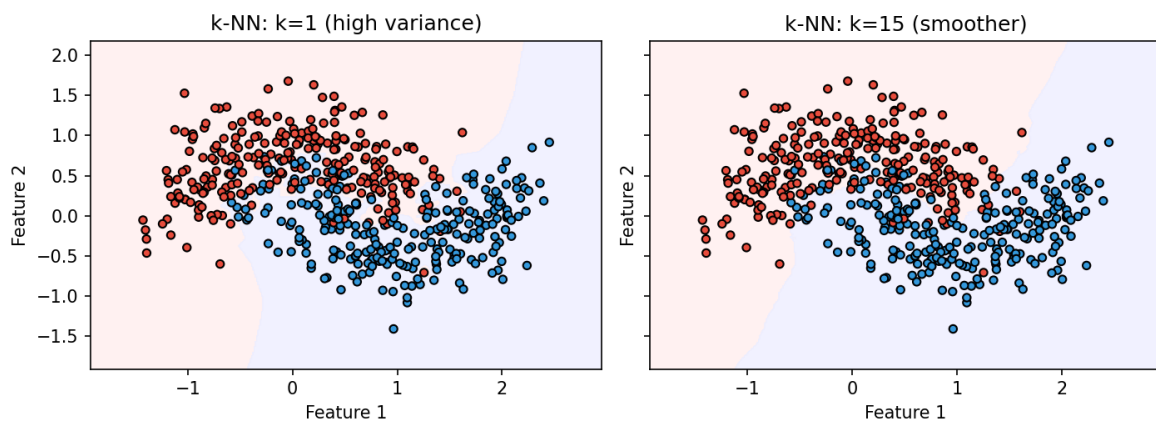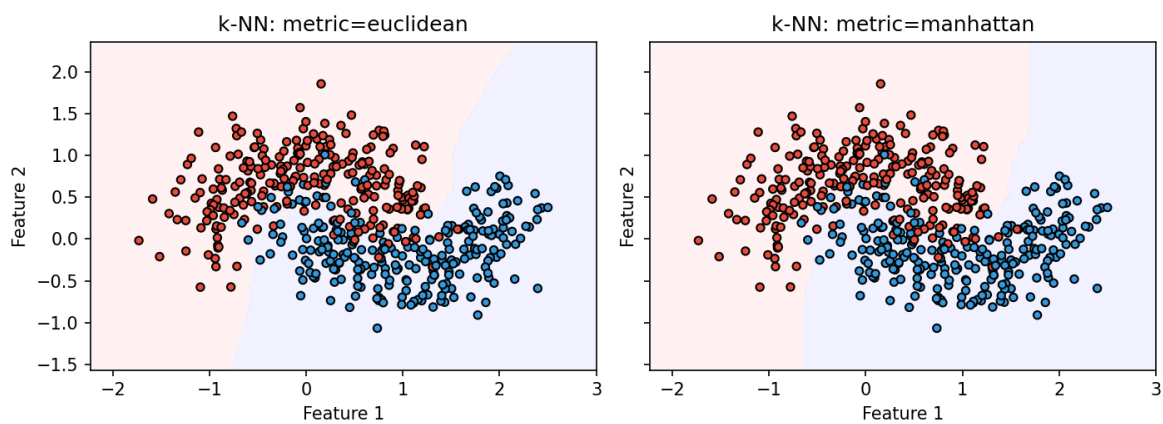
# 5 Result



Figure 1: Decision boundaries for different k (1 vs 15).



Figure 2: Effect of metric: Euclidean vs Manhattan.
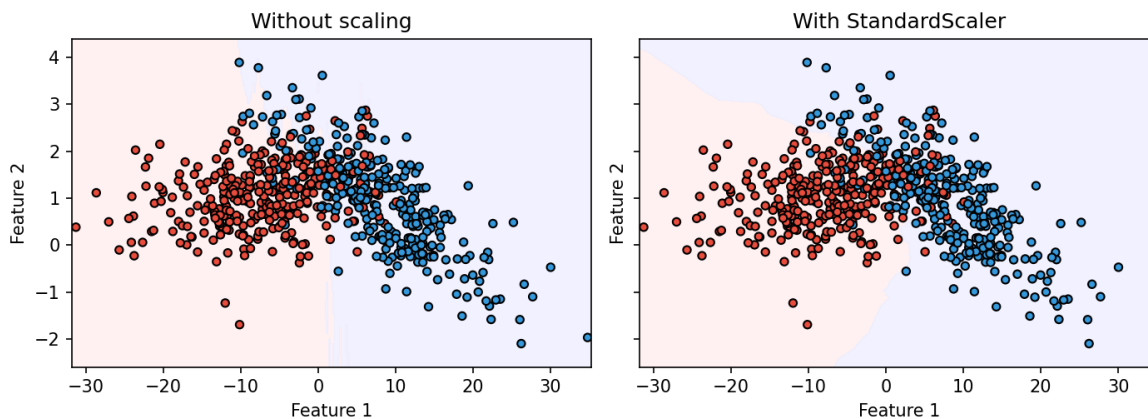
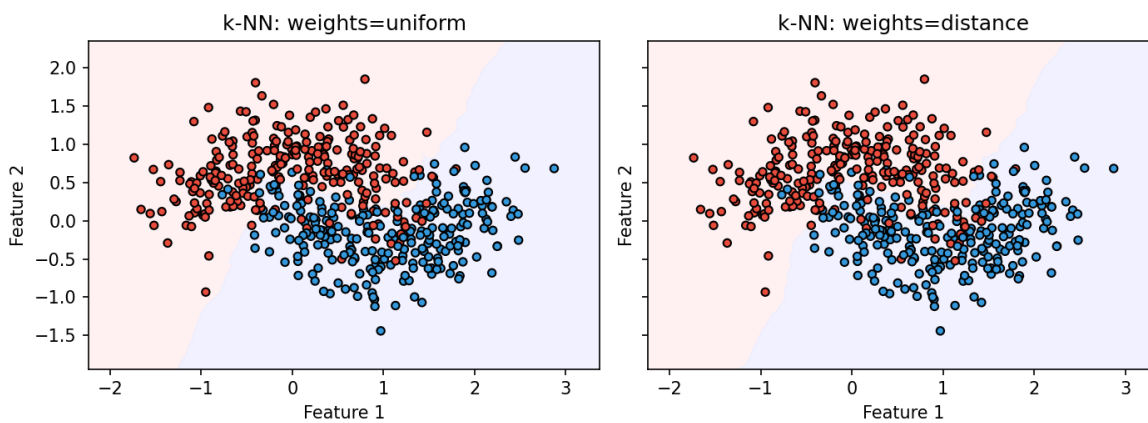Figure 3: Impact of feature scaling on decision boundary.



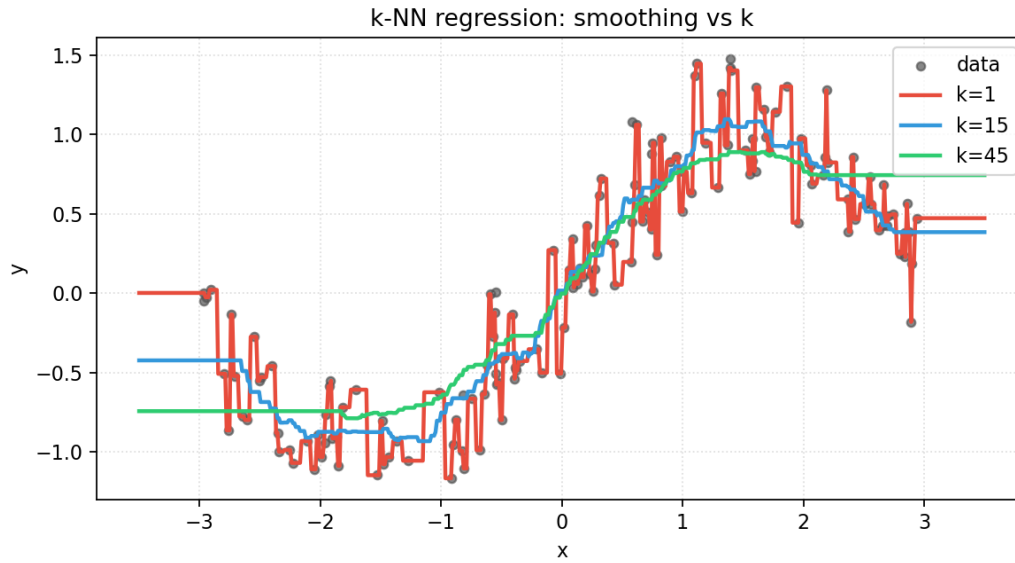Figure 4: Uniform vs distance weighting in k-NN.

Figure 5: k-NN regression: smoothing effect as k increases.

# 6 Summary

k-NN is a simple yet powerful baseline for both classification and regression when features are well-scaled and dimensionality is moderate. Select $k$, metric, and weighting via validation, and use scaling to ensure distances are meaningful.