

# k-近邻 (k-NN): 理论与实践

2025 年 9 月 10 日

## 1 引言

k-近邻 (k-Nearest Neighbors, k-NN) 是一种非参数、基于实例的“惰性学习”方法：预测时在训练集中查找与查询样本最近的  $k$  个邻居。其思想直观、实现简单，在低维、良好缩放的数据上常有竞争力，但对特征缩放敏感，并在高维下性能下降。

## 2 原理与公式

给定查询点  $\mathbf{x}$ ，在指定距离度量  $d(\cdot, \cdot)$ （如欧氏、曼哈顿）下找到其最近的  $k$  个邻居。分类任务采用多数表决（可选距离加权 `weights=distance` 使近邻权重更高）；回归任务取邻居目标值的平均（或距离加权平均）。

计算上，朴素查找每次预测代价为  $\mathcal{O}(nd)$ （ $n$  为样本数， $d$  为维度）。中等维度时可使用 KDTree/BallTree 提速。k-NN 受“维度灾难”影响，合适的特征缩放与度量选择至关重要。

## 3 应用与技巧

- **选择  $k$** : 通过交叉验证调参；二分类常取奇数  $k$  以减少平票。
- **缩放**: 对特征做标准化/归一化，或使用 Pipeline；距离对量纲非常敏感。
- **度量**: 尝试欧氏与曼哈顿；必要时考虑领域特定距离。
- **权重**: `uniform` 与 `distance` 可对类间重叠区产生不同效果。
- **复杂度**: 预测代价随数据量增长；大规模可考虑近似近邻检索。

## 4 Python 实战

在本章节目录运行下述命令，图片将保存到 figures/:

Listing 1: 生成 k-NN 配图

```
1 python gen_knn_figures.py
```

Listing 2: gen\_knn\_figures.py 源码

```
1 """
2 Figure generator for the k-NN chapter.
3
4 Generates illustrative figures and saves them into the chapter's '
   figures/'
5 folder next to this script, regardless of current working directory.
6
7 Requirements:
8 - Python 3.8+
9 - numpy, matplotlib, scikit-learn
10
11 Install (if needed):
12     pip install numpy matplotlib scikit-learn
13
14 This script avoids newer or experimental APIs for broader compatibility
15 .
16 """
17
18 from __future__ import annotations
19
20 import os
21 import numpy as np
22 import matplotlib.pyplot as plt
23 from matplotlib.colors import ListedColormap
24
25 try:
26     from sklearn.datasets import make_moons, make_regression,
27         make_classification
28     from sklearn.neighbors import KNeighborsClassifier,
29         KNeighborsRegressor
30     from sklearn.preprocessing import StandardScaler
31     from sklearn.pipeline import make_pipeline
32     from sklearn.model_selection import cross_val_score
33 except Exception:
34     raise SystemExit(
```

```
31         "Missing scikit-learn. Please install with: pip install scikit-
           learn"
32     )
33
34
35 def _ensure_figures_dir(path: str | None = None) -> str:
36     """Create figures directory under this chapter regardless of CWD.
       """
37     if path is None:
38         base = os.path.dirname(os.path.abspath(__file__))
39         path = os.path.join(base, "figures")
40     os.makedirs(path, exist_ok=True)
41     return path
42
43
44 def _plot_decision_boundary(ax, clf, X, y, title: str):
45     x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
46     y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
47     xx, yy = np.meshgrid(
48         np.linspace(x_min, x_max, 400), np.linspace(y_min, y_max, 400)
49     )
50     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
51     cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
52     cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
53     ax.contourf(xx, yy, Z, cmap=cmap_light, alpha=0.8, levels=np.unique(
54         Z).size)
55     ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s
56               =20)
57     ax.set_title(title)
58     ax.set_xlabel("Feature 1")
59     ax.set_ylabel("Feature 2")
60
61 def fig_knn_k_compare(out_dir: str) -> str:
62     np.random.seed(0)
63     X, y = make_moons(n_samples=500, noise=0.3, random_state=0)
64     models = [
65         (KNeighborsClassifier(n_neighbors=1), "k=1 (high variance)"),
66         (KNeighborsClassifier(n_neighbors=15), "k=15 (smoother)")
67     ]
68     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
69                           True, sharey=True)
70     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"), "
84          metric=euclidean"),
85         (KNeighborsClassifier(n_neighbors=11, metric="manhattan"), "
86          metric=manhattan"),
87     ]
88     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
89     True, sharey=True)
90     for ax, (m, title) in zip(axes, models):
91         m.fit(X, y)
92         _plot_decision_boundary(ax, m, X, y, f"k-NN: {title}")
93     fig.suptitle("Effect of distance metric")
94     out_path = os.path.join(out_dir, "knn_metric_compare.png")
95     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
96     fig.savefig(out_path)
97     plt.close(fig)
98     return out_path
99
100 def fig_knn_scaling_effect(out_dir: str) -> str:
101     np.random.seed(2)
102     X, y = make_classification(
103         n_samples=600,
104         n_features=2,
105         n_informative=2,
106         n_redundant=0,
107         n_clusters_per_class=1,
108         class_sep=1.0,
109         random_state=2,
110     )
```

```

109     # Impose different scales on features
110     X_scaled_variance = X.copy()
111     X_scaled_variance[:, 0] *= 8.0 # make feature 0 dominate distances
112
113     knn_raw = KNeighborsClassifier(n_neighbors=11)
114     knn_std = make_pipeline(StandardScaler(), KNeighborsClassifier(
115         n_neighbors=11))
116
117     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
118         True, sharey=True)
119     knn_raw.fit(X_scaled_variance, y)
120     _plot_decision_boundary(axes[0], knn_raw, X_scaled_variance, y, "
121         Without scaling")
122     knn_std.fit(X_scaled_variance, y)
123     _plot_decision_boundary(axes[1], knn_std, X_scaled_variance, y, "
124         With StandardScaler")
125     fig.suptitle("Feature scaling impact on k-NN")
126     out_path = os.path.join(out_dir, "knn_scaling_effect.png")
127     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
128     fig.savefig(out_path)
129     plt.close(fig)
130     return out_path
131
132 def fig_knn_weight_compare(out_dir: str) -> str:
133     np.random.seed(3)
134     X, y = make_moons(n_samples=500, noise=0.32, random_state=3)
135     models = [
136         (KNeighborsClassifier(n_neighbors=11, weights="uniform"), "
137             weights=uniform"),
138         (KNeighborsClassifier(n_neighbors=11, weights="distance"), "
139             weights=distance"),
140     ]
141     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
142         True, sharey=True)
143     for ax, (m, title) in zip(axes, models):
144         m.fit(X, y)
145         _plot_decision_boundary(ax, m, X, y, f"k-NN: {title}")
146     fig.suptitle("Uniform vs distance weighting")
147     out_path = os.path.join(out_dir, "knn_weight_compare.png")
148     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
149     fig.savefig(out_path)
150     plt.close(fig)

```

```

145     return out_path
146
147
148 def fig_knn_regression_curve(out_dir: str) -> str:
149     rng = np.random.RandomState(4)
150     # 1D regression:  $y = \sin(x) + \text{noise}$ 
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 结果

Effect of k on decision boundary

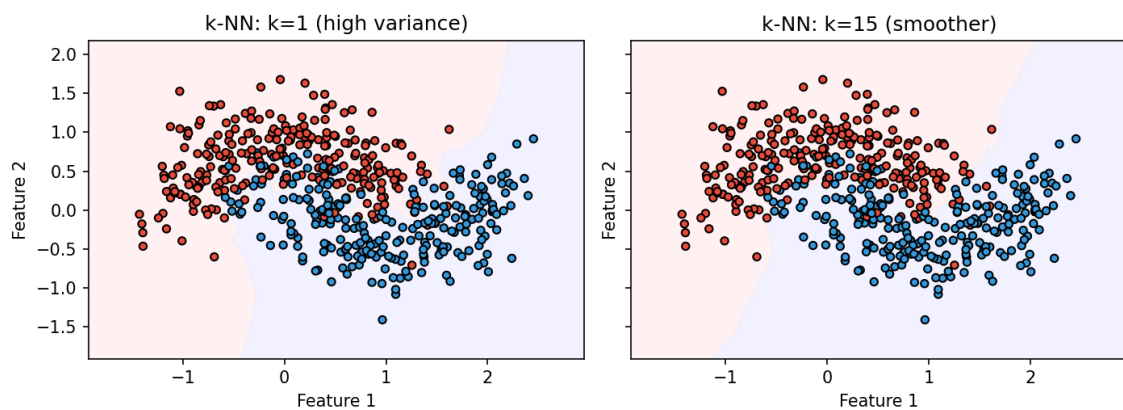


图 1: 不同 k (1 vs 15) 的决策边界对比。

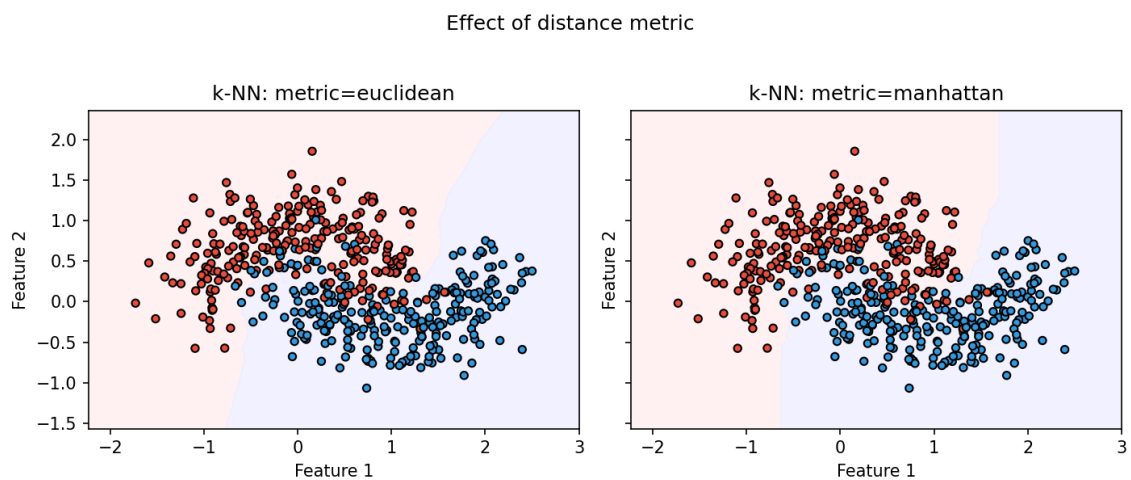


图 2: 不同距离度量: 欧氏 (Euclidean) vs 曼哈顿 (Manhattan)。

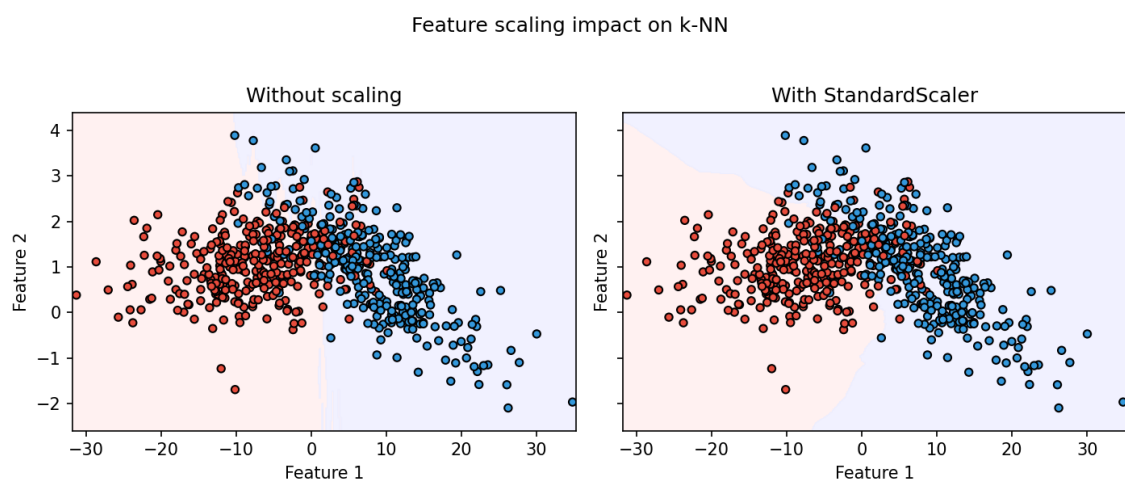


图 3: 特征缩放对决策边界的影响。

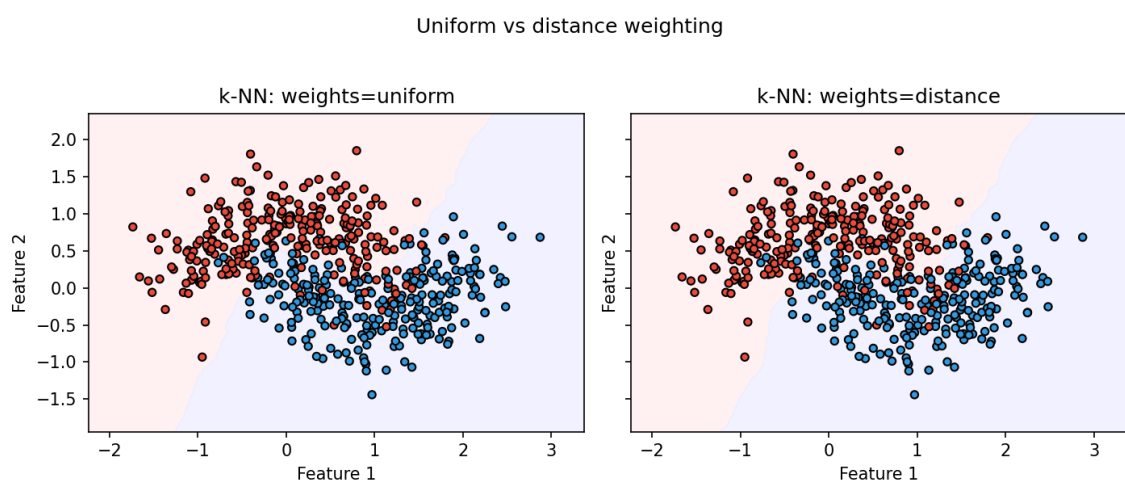


图 4: 均匀权重 vs 距离加权的对比。



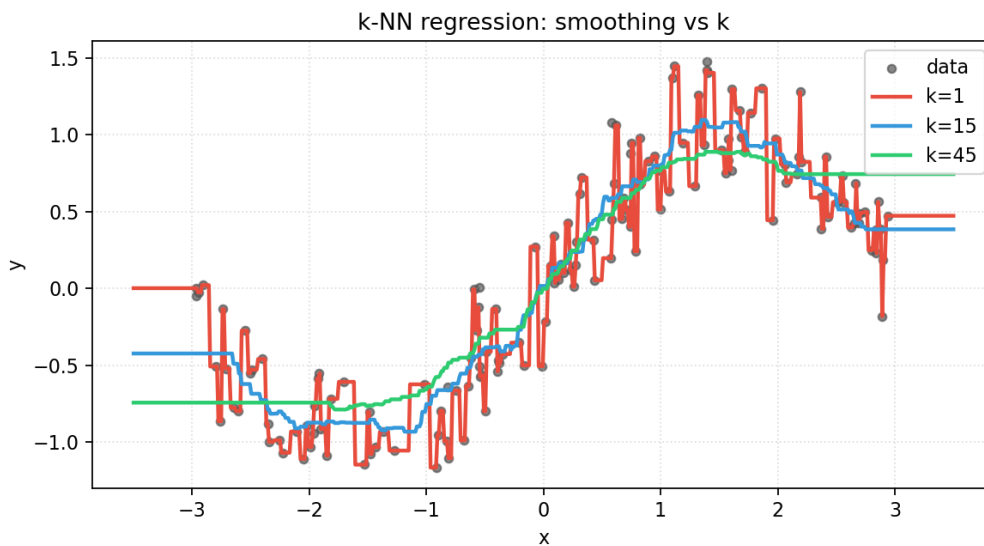


图 5: k-NN 回归：随  $k$  增大平滑程度的变化。

## 6 总结

k-NN 在特征良好缩放且维度适中的场景下是简洁有效的基线。通过验证选择合适的  $k$ 、距离度量与权重设置，并做好缩放预处理，可获得稳定表现。