

K-Nearest Neighbors (KNN) Classifier

1. Pendahuluan

K-Nearest Neighbors (KNN) adalah algoritma non-parametrik untuk klasifikasi yang berdasarkan pada prinsip bahwa suatu data x akan diprediksi sebagai kelas mayoritas dari k tetangga terdekatnya dalam ruang fitur. Algoritma ini tidak melakukan proses pelatihan eksplisit, melainkan menyimpan data pelatihan dan melakukan prediksi dengan membandingkan data baru terhadap seluruh data pelatihan.

2. Formulasi Matematis

Diberikan:

- Dataset pelatihan: $D = \{(x_i, y_i)\}_{i=1}^n$
- Titik input baru: $x \in \mathbb{R}^d$
- Fungsi jarak: $d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Langkah prediksi:

1. Hitung jarak $d(x, x_i)$ untuk seluruh $i \in \{1, \dots, n\}$
2. Ambil $N_k(x) \subset D$, yaitu k data pelatihan dengan jarak terkecil terhadap x
3. Prediksi kelas berdasarkan mayoritas:

$$\hat{y} = \arg \max_{c \in \mathcal{C}} \sum_{(x_i, y_i) \in N_k(x)} 1[y_i = c]$$

3. Parameter dan Pengaturan

- k : Jumlah tetangga yang dipertimbangkan dalam prediksi.
- **distance_metric**: Jenis fungsi jarak yang digunakan:
 - **Euclidean**: $d(x, x') = \sqrt{\sum_j (x_j - x'_j)^2}$
 - **Manhattan**: $d(x, x') = \sum_j |x_j - x'_j|$
 - **Minkowski (p=2)**: $d(x, x') = (\sum_j |x_j - x'_j|^p)^{1/p}$

4. Algoritma

4.1. Training (fit)

Fungsi **fit(X, y)** menyimpan data pelatihan tanpa proses parameter learning. Data disimpan dalam:

$$X_{\text{train}} \in \mathbb{R}^{n \times d}, \quad y_{\text{train}} \in \mathbb{R}^n$$

4.2. Prediksi

Untuk setiap sampel uji x :

1. Hitung jarak terhadap semua $x_i \in X_{\text{train}}$
2. Pilih k sampel dengan jarak terkecil
3. Ambil label terbanyak sebagai hasil prediksi \hat{y}

$$\hat{y}(x) = \text{mode}(y_i \mid (x_i, y_i) \in N_k(x))$$

4.3. Tiebreaking

Jika terjadi jumlah kelas yang sama, digunakan aturan `np.argmax` pada `np.unique(..., return_counts=True)` untuk mengambil label dengan urutan indeks terkecil dalam `unique_labels`.

5. Kompleksitas

- **Training time:** $O(1)$
- **Prediction time:** $O(n \cdot d)$ per sampel uji
- **Memory usage:** $O(n \cdot d)$