

K-nearest neighbors

10/14

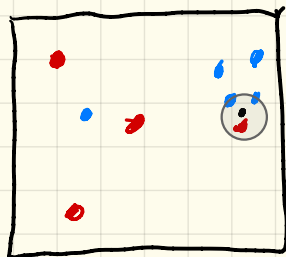
Let $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$ w/ $\vec{x}_i \in \mathbb{R}^d$ our training set
 $D_i \subseteq D$ pts w/ label c_i
labels $\{c_1, \dots, c_k\}$
 $n_i = |D_i|$

Given test point $\vec{x} \in \mathbb{R}^d$
and $k \in \mathbb{Z}^+$ # of neighbors to consider

Let $r \in \mathbb{R}$ be the dist of \vec{x} 's k -th nearest neighbor in D

$$B_d(\vec{x}, r) = \{p \in \mathbb{R}^d \mid \|p - \vec{x}\| \leq r\}$$

↑
Euclidean dist



K is the k -nearest neighbors of \vec{x}
 $K \subseteq B(\vec{x}, r)$

$$K_i = \{\vec{x}_j \in K \mid y_j = c_i\}$$

↙ all pts in radius r ball w/ label c_i

$$\text{and } n_i = |K_i|$$

Idea estimate $P(c_i | \vec{x})$
estimate the
conditional prob density at \vec{x}

- \vec{x}
- $c_1 = \text{red}$
- $c_2 = \text{blue pts}$
- $D_1 = \text{set of all red pts}$
- $D_2 = \text{set of all blue pts}$
- $k=3$
- $k_{\text{red}} = 1$
- $k_{\text{blue}} = 2$

Step 1: estimate density

Let $V = \text{vol}(B_d(\vec{x}, r))$ ↗ vol of the d-dim ball centered at \vec{x} w/ radius r

$$= \text{vol}(B_{d-2}(\vec{x}, r)) \times \frac{2\pi r^2}{d}$$

$$\text{w/ } \text{vol}(B_2(\vec{x}, r)) = \pi r^2$$

$$\text{vol}(B_3(\vec{x}, r)) = \frac{4}{3} \pi r^3$$

Estimate density at \vec{x}

$$f(\vec{x} | c_i) = \frac{\frac{k_i}{n_i}}{V} = \frac{k_i}{n_i V}$$

Step 2:

$$p(c_i | \vec{x}) = \frac{\hat{f}(\vec{x} | c_i) \hat{p}(c_i)}{\sum_{j=1}^K \hat{f}(\vec{x} | c_j) \hat{p}(c_j)}$$

← $\frac{n_i}{n}$

← $\frac{n_j}{n}$

$$\hat{f}(\vec{x} | c_i) \hat{p}(c_i) = \frac{k_i}{n_i V} \cdot \frac{n_i}{n} = \frac{k_i}{n V}$$

$$\hat{p}(c_i | \vec{x}) = \frac{\frac{k_i}{n V}}{\sum_{j=1}^K \frac{k_j}{n V}} = \frac{\frac{1}{n V} k_i}{\frac{1}{n V} \sum_{j=1}^K k_j} = \frac{k_i}{\sum_{j=1}^K k_j} = \frac{k_i}{K}$$

$$\hat{y} = \arg \max_{c_i} \{p(c_i | \vec{x})\} = \arg \max_{c_i} \left\{ \frac{k_i}{K} \right\} = \arg \max_{c_i} \{k_i\}$$

↗ pick label w/ majority vote from knn