

Kernels

Background: local learning
aka instance-based "
memory-based "
non-parametric "

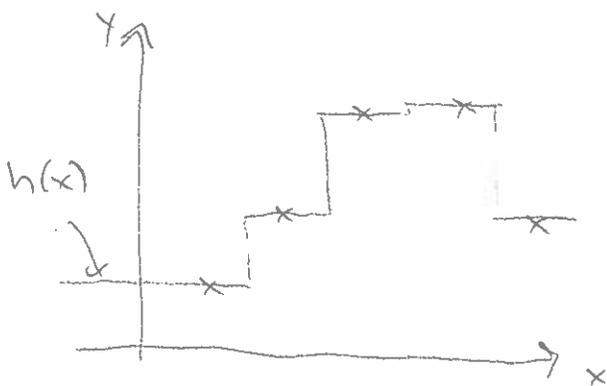
R&N 18.8

Nearest-neighbor method

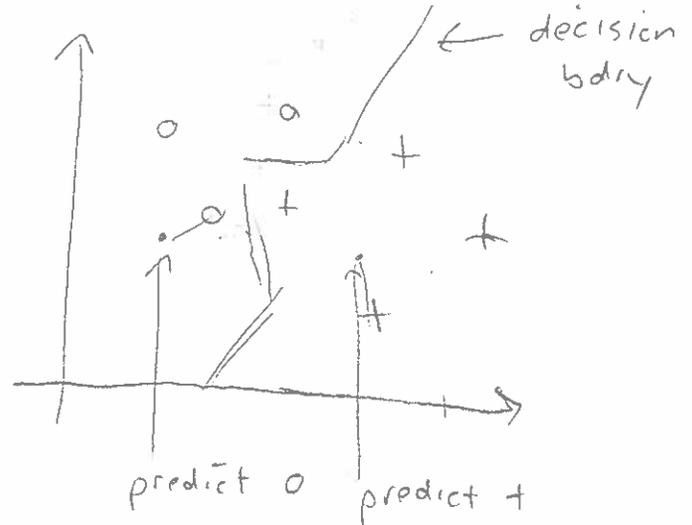
- Training examples $(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)$
- No training phase
- On input \vec{x} , predict output of closest training example to \vec{x} :

$$\begin{cases} j = \operatorname{argmin}_i \|x - x_i\|^2 \\ h(x) = y_j \end{cases}$$

E.g. 1D regression

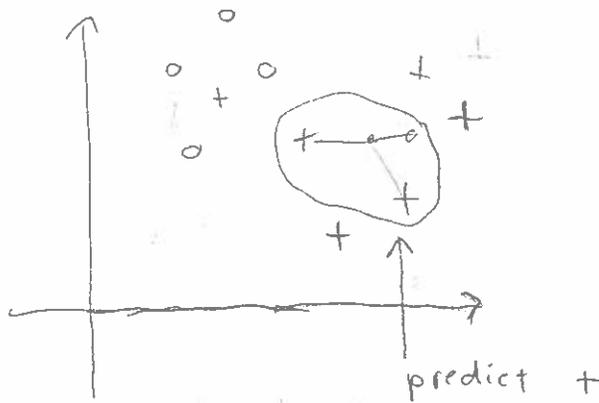


2D classification



Improvements

#1) Use K nearest neighbors instead of one
(avoid overfitting; see R&N Fig. 18.26)



e.g. $K=3$

Take majority vote
of 3 closest neighbors

Formally: on input \vec{x} ,

- Let j_1, j_2, \dots, j_k be set of k closest examples,
(smallest distance $\|\vec{x} - \vec{x}_{j_i}\|$)

$$h(\vec{x}) = \frac{1}{k} \sum_{p=1}^k Y_{j_p} \quad \leftarrow \text{average over neighbors}$$

- Classification

$h(\vec{x}) =$ fraction of K neighbors with $y=1$

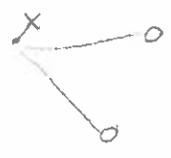
$$\text{predict} \begin{cases} 1 & h(\vec{x}) \geq \frac{1}{2} \\ 0 & h(\vec{x}) < \frac{1}{2} \end{cases}$$

(Select K odd to avoid ties)

#2) Kernel regression (classification)

: KNN downsides:

- all neighbors weighted equally
- how to choose k ?



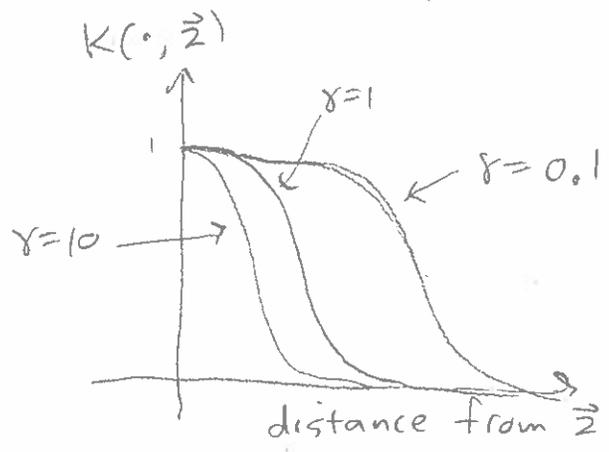
Idea: use all neighbors, weight by distance, using kernel function

$$K(\vec{x}, \vec{z}) \equiv \text{similarity between } \vec{x} \text{ + } \vec{z}$$

(later: view as a dot product)

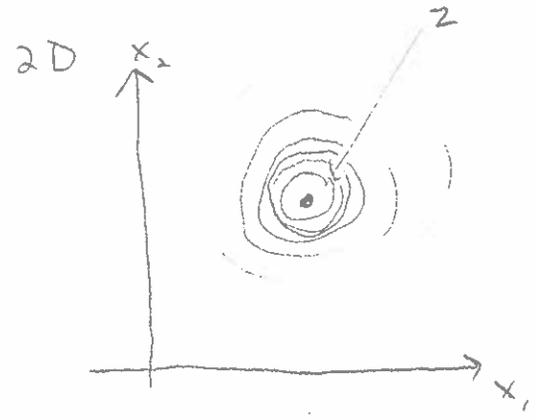
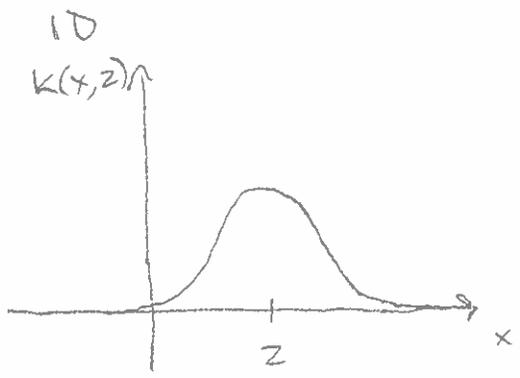
E.g. Gaussian Kernel

$$K(\vec{x}, \vec{z}) = \exp(-\gamma \|\vec{x} - \vec{z}\|^2)$$



- closer points are most similar
- exponential dropoff with distance

Similarity to a fixed point \vec{z}

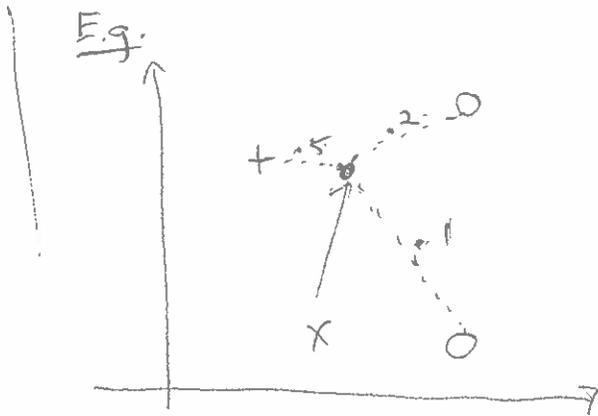


Kernel regression algorithm

$$h(\vec{x}) = \frac{\sum_{i=1}^N K(\vec{x}, \vec{x}_i) \cdot y_i}{\sum_{i=1}^N K(\vec{x}, \vec{x}_i)}$$

← weight of \vec{x}_i

← total weight



$$h(x) = \frac{(0.5)(1) + (0.2)(0) + (0.1)(0)}{0.5 + 0.2 + 0.1}$$

$$= \frac{0.5}{0.8}$$

$$= \frac{5}{8} > \frac{1}{2}$$

⇒ predict '+'