Non-parametric methods

Aarti Singh

Machine Learning 10-315 Oct 11, 2021





Parametric methods

- Assume some model (Gaussian, Bernoulli, Multinomial, logistic, network of logistic units, Linear, Quadratic) with fixed number of parameters
 - Gaussian Bayes, Naïve Bayes, Logistic Regression, Neural Networks
- Estimate parameters $(\mu, \sigma^2, \theta, w, \beta)$ using MLE/MAP and plug in
- Pro need few data points to learn parameters
- Con Strong distributional assumptions, not satisfied in practice

Non-Parametric methods

modeling

- Typically don't make any distributional assumptions
- As we have more data, we should be able to learn more complex models



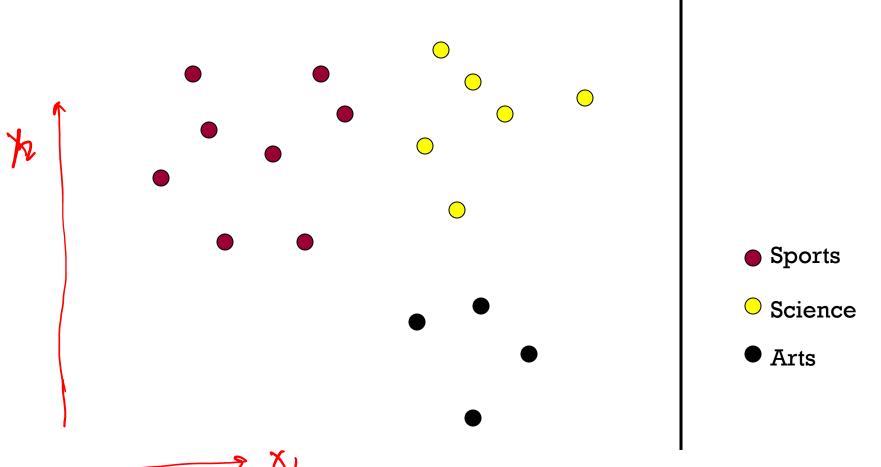
- Let number of parameters scale with number of training data
- Some nonparametric methods

Classification: Decision trees, k-NN (k-Nearest Neighbor) classifier

Density estimation: k-NN, Histogram, Kernel density estimate

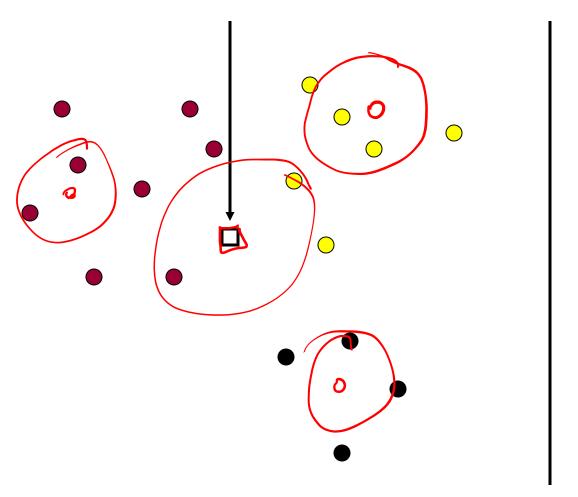
Regression: Kernel regression

k-NN classifier



k-NN classifier

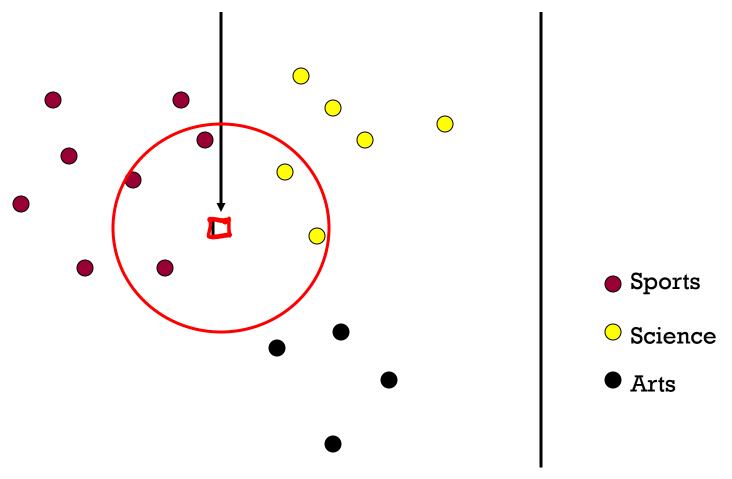
Test document



- Sports
- Science
- Arts

k-NN classifier (k=5)

Test document



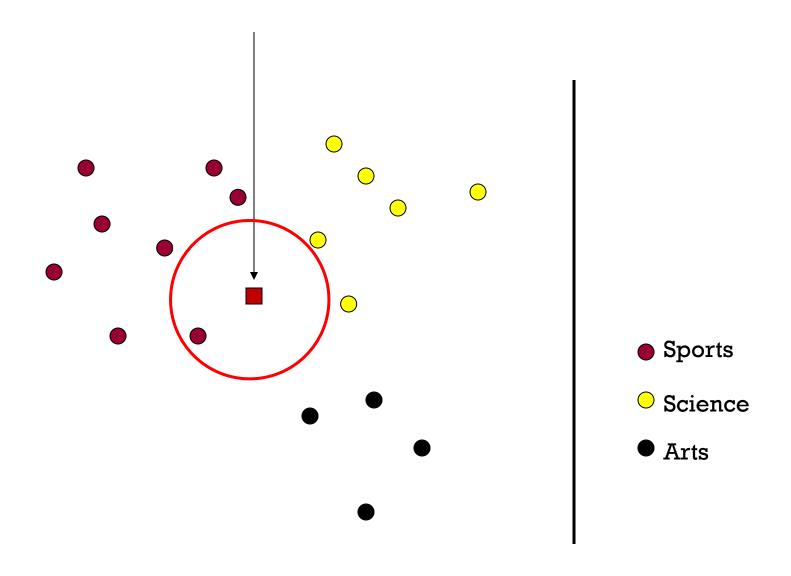
k-NN classifier

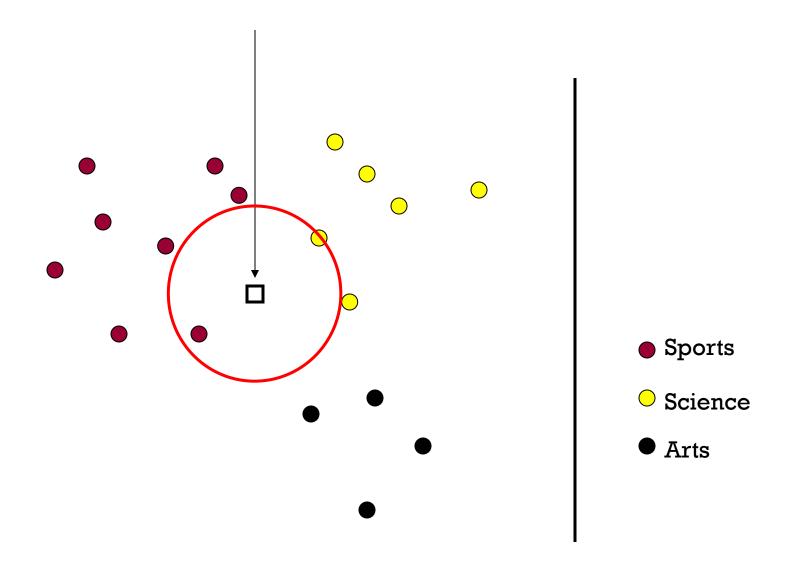
- Optimal Classifier: $f^*(x) = \arg\max_y P(y|x)$ = $\arg\max_y P(x|y)P(y)$
- k-NN Classifier: $\widehat{f}_{kNN}(x) = \arg\max_{y} \widehat{P}_{kNN}(x|y)\widehat{P}(y)$ $= \arg\max_{y} k_{y}$

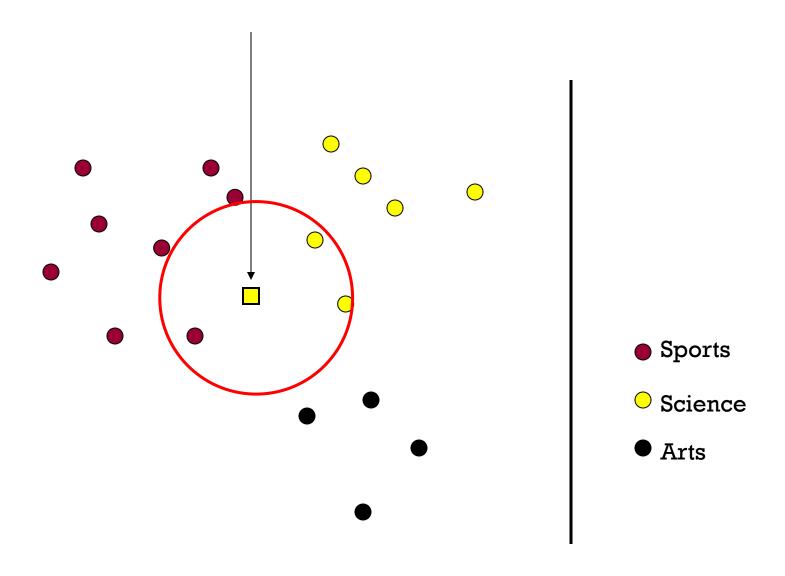
$$\widehat{P}_{kNN}(x|y) = \frac{k_y}{n_y} \longrightarrow \text{\# training pts of class y} \\ \text{amongst k NNs of x} \\ \text{\searrow} k_y = k$$

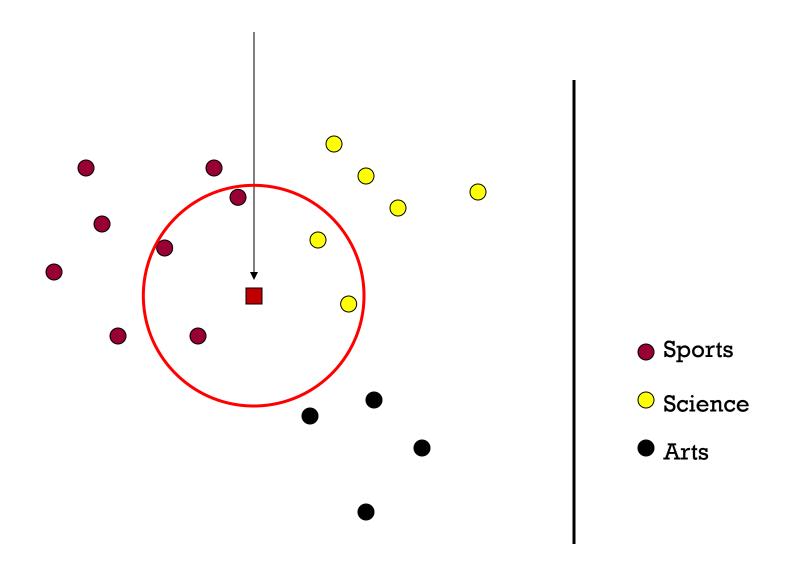
$$\text{\# total training pts of class y}$$

$$\widehat{P}(y) = \frac{n_y}{n}$$



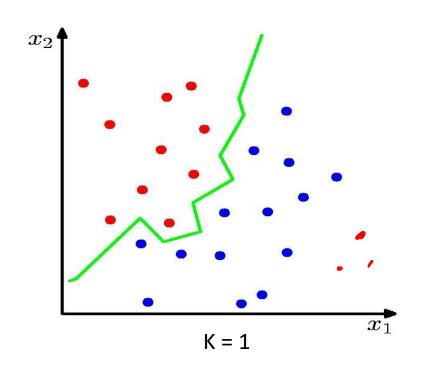






What is the best k?

1-NN classifier decision boundary



Voronoi Diagram

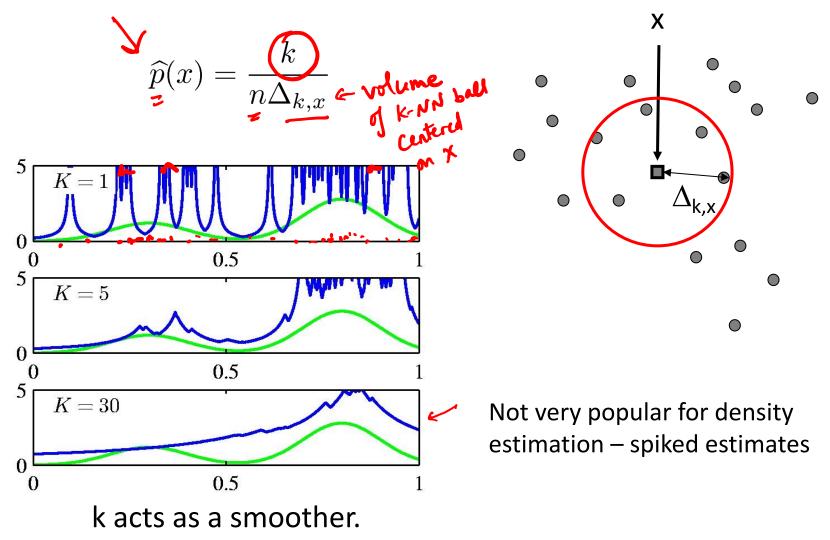
As k increases, boundary becomes smoother (less jagged).

What is the best k?

Approximation vs. Stability (aka Bias vs Variance) Tradeoff

- Larger K => predicted label is more stable
- Smaller K => predicted label can approximate best classifier well given enough data

k-NN density estimation



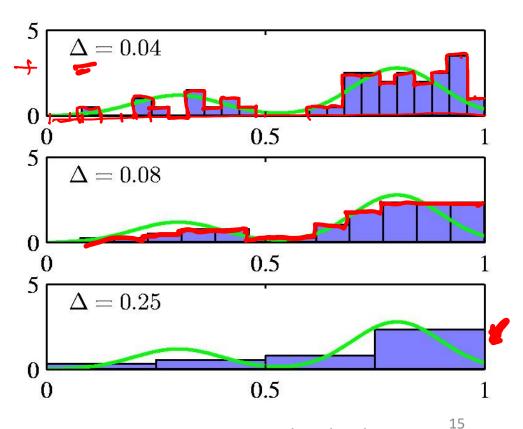
Histogram density estimate

Partition the feature space into distinct bins with widths Δ_i and count the number of observations, n_i , in each bin.

$$\widehat{p}(x) = \frac{n_i}{n\Delta_i} \mathbf{1}_{x \in \text{Bin}_i}$$

"Local relative frequency"

- Often, the same width is used for all bins, $\Delta_i = \Delta$.
- \(\frac{1}{2}\) acts as a smoothing parameter.

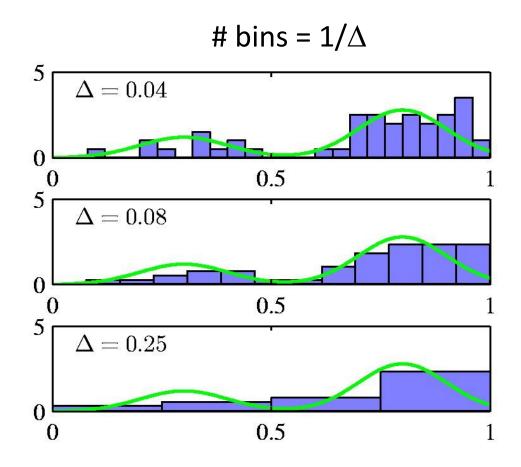


Effect of histogram bin width

$$\widehat{p}(x) = \frac{n_i}{n\Delta} \mathbf{1}_{x \in \text{Bin}_i}$$

Small △, large #bins
Good fit but unstable
(few points per bin)
"Small bias, Large variance"

Large △, small #bins
Poor fit but stable
(many points per bin)
"Large bias, Small variance"



Histogram as MLE



Underlying model – density is constant on each bin P(XEBin;) = Pja Parameters p_j : density in bin j

Note
$$\sum_{j} p_{j} = 1/\Delta$$
 since $\int p(x)dx = 1$

Maximize likelihood of data under probability model with $parameters \ p_j$

parameters
$$\mathbf{p}_j$$

$$\hat{p}(x) = \arg\max_{\{p_j\}} P(X_1, \dots, X_n; \{p_j\}_{j=1}^{1/\Delta}) \quad \text{s.t.} \quad \sum_j p_j = 1/\Delta$$

Show that histogram density estimate is MLE under this model

Histogram as MLE Pinus 12

$$\hat{p}(x) = \arg\max_{\{p_j\}} P(X_1, \dots, X_n; \{p_j\}_{j=1}^{1/\Delta}) \quad \text{s.t.} \quad \sum_j p_j = 1/\Delta$$

$$A \quad \prod_{j=1}^{1/\Delta} n^{n_j} \text{ where } n_j = number \text{ of } data \text{ in him } i$$

A.
$$\prod_{j=1}^{n_j} p_j^{n_j} \text{ where } n_j - number \text{ of data in bin } j$$

$$B. \prod_{i=1}^{n} p_i$$

B.
$$\prod_{j=1}^n p_j^{1/\Delta}$$

arg max
$$log(\Pi p_i^n) = \frac{1}{2} n_i log p_i$$
 c.t. $Zp_i = 1/2$
 Sp_i^n

$$\frac{\partial}{\partial P_i} = \frac{n_i}{P_i} - \lambda = 0$$

Irg max log (
$$\prod_{j=1}^{n} P_{j}^{i}$$
) = $\lim_{j \to \infty} \frac{1}{2} \prod_{j=1}^{n} \frac{1}{2} \prod_{j=1$