



# Nonparametric methods

---

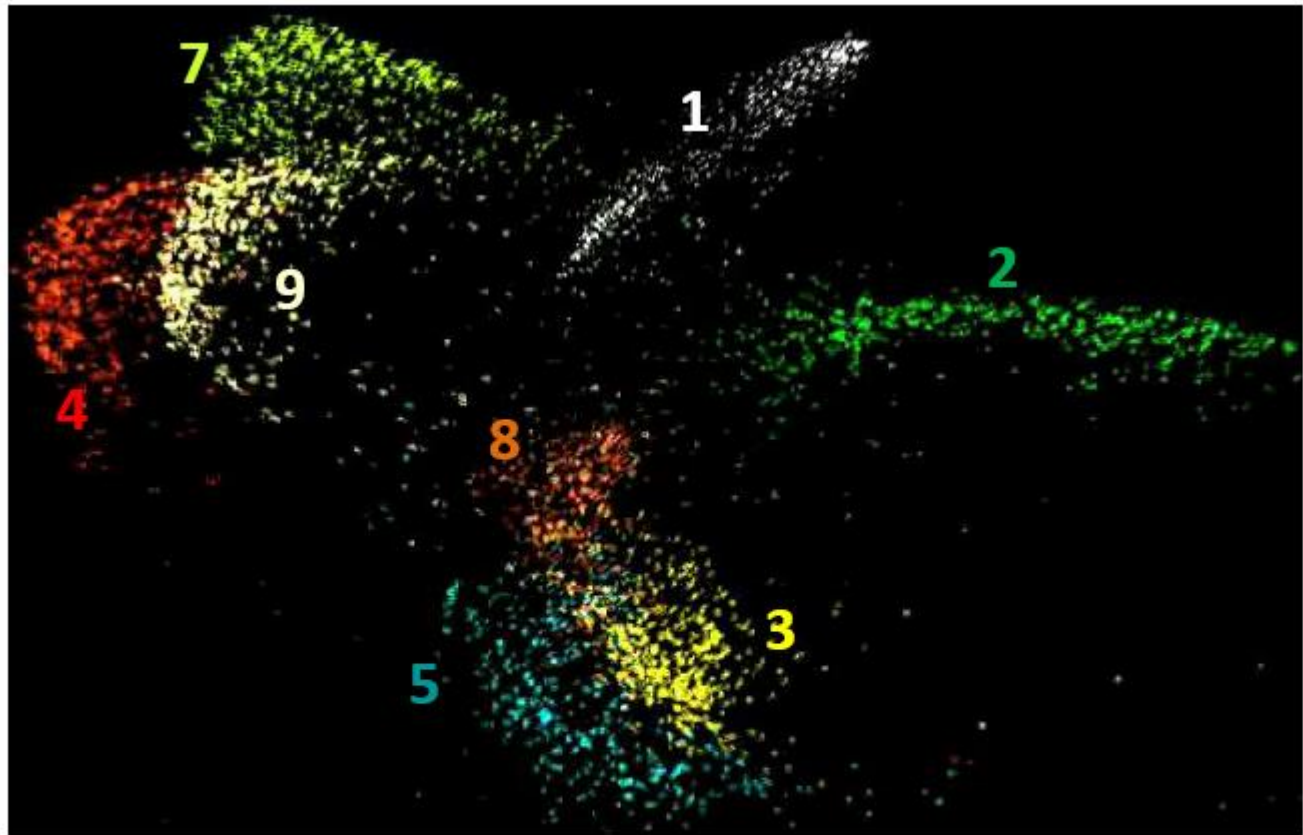


# Parametric methods

---

- Assume some functional form (Gaussian, Bernoulli, Multinomial, logistic, Linear) for
  - $P(X_i|Y)$  and  $P(Y)$  as in Naïve Bayes
  - $P(Y|X)$  as in Logistic regression
- 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

# Example



Hand-written digit images  
projected as points on a two-dimensional (nonlinear) feature spaces



# Non-Parametric methods

---

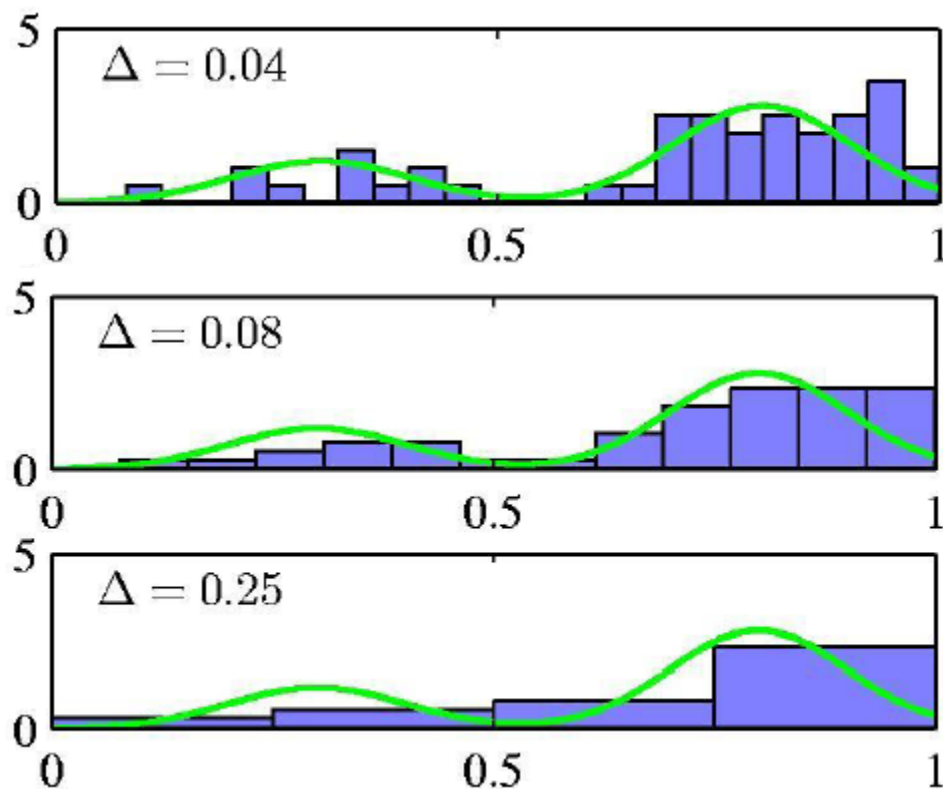
- Typically don't make any distributional assumptions
- As we have more data, we should be able to learn more complex models
- Let number of parameters scale with number of training data
- Today, we will see some nonparametric methods for
  - Density estimation
  - Classification

# Histogram density estimate

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

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

- Often, the same width is used for all bins,  $\Delta_i = \Delta$ .
- $\Delta$  acts as a smoothing parameter.

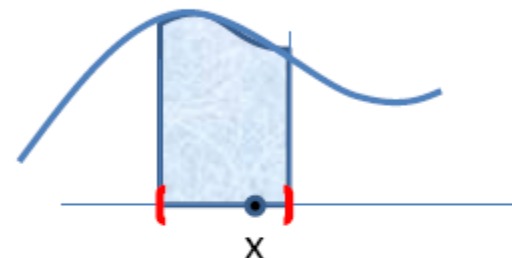


# Effect of histogram bin width

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

$$\# \text{ bins} = 1/\Delta$$

$$\hat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^n \mathbf{1}_{X_j \in \text{Bin}_x}}{n}$$



Bias of histogram density estimate:

$$\mathbb{E}[\hat{p}(x)] = \frac{1}{\Delta} P(X \in \text{Bin}_x) = \frac{1}{\Delta} \int_{z \in \text{Bin}_x} p(z) dz \approx \frac{p(x)\Delta}{\Delta} = p(x)$$



**Assuming density is roughly constant in each bin  
(holds true if  $\Delta$  is small)**



# Bias-Variance tradeoff

---

- Choice of #bins

$$\# \text{ bins} = 1/\Delta$$

$\mathbb{E}[\hat{p}(x)] \approx p(x)$  if  $\Delta$  is small      ( $p(x)$  approx constant per bin)

$\mathbb{E}[\hat{p}(x)] \approx \hat{p}(x)$  if  $\Delta$  is large      (more data per bin,  
stable estimate)

- Bias** – how close is the mean of estimate to the truth
- Variance** – how much does the estimate vary around mean

Small  $\Delta$ , large #bins  $\longleftrightarrow$  “**Small bias**, **Large variance**”

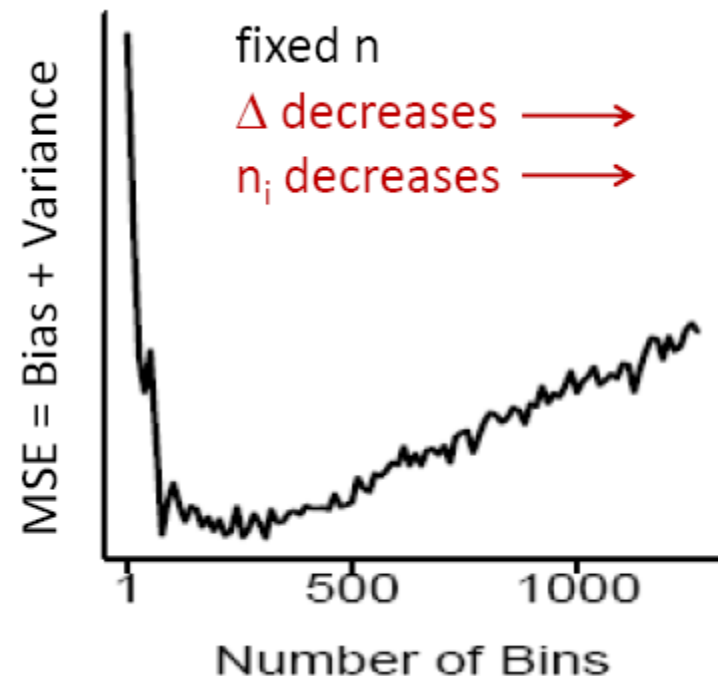
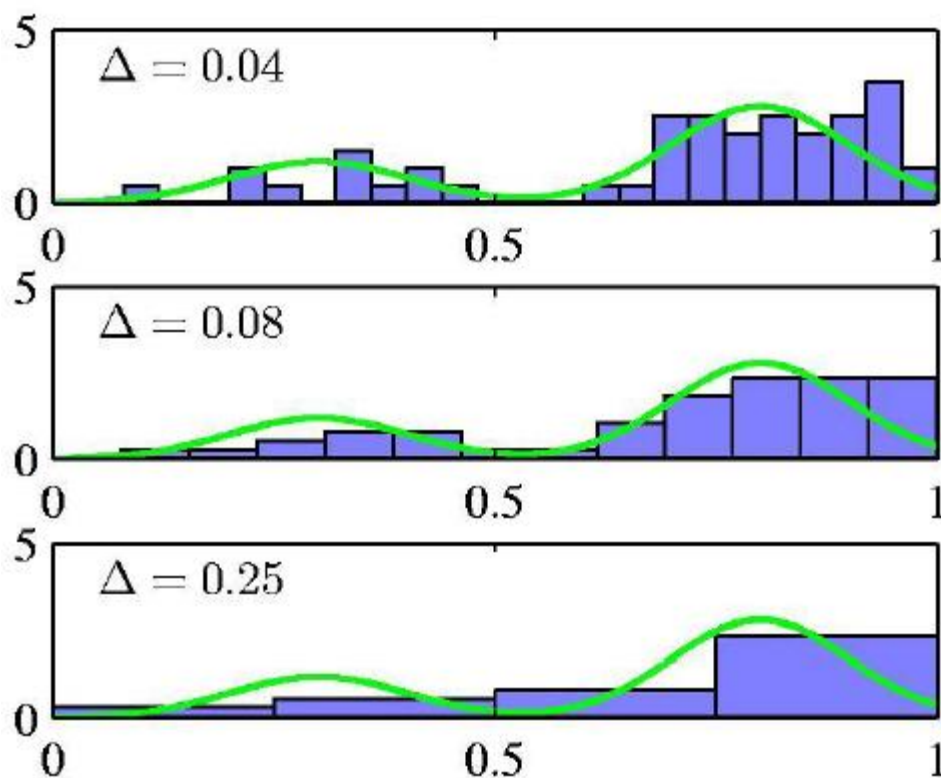
Large  $\Delta$ , small #bins  $\longleftrightarrow$  “**Large bias**, **Small variance**”

**Bias-Variance tradeoff**

# Choice of #bins

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

$$\# \text{ bins} = 1/\Delta$$







# Histogram: Conclusion

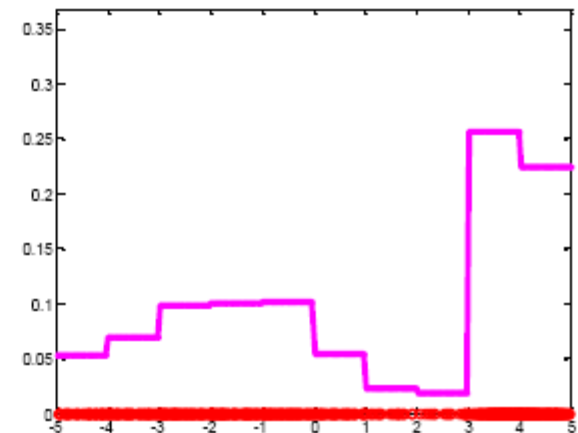
---

- **Pro**-once computed, the data set itself can be discarded, advantageous if the data set is large; easily applied if the data points are arriving sequentially; a quick visualization of data in one or two dims
- **Con**-discontinuities on bin edge;  $(\#bins)^D$
- Two hints: local; delta

# Kernel density estimate

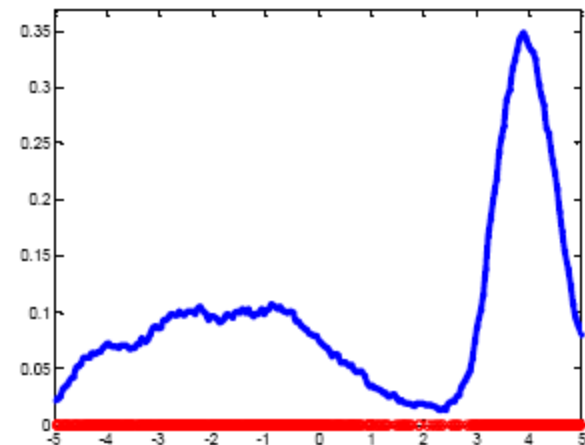
- Histogram – blocky estimate

$$\hat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^n \mathbf{1}_{X_j \in \text{Bin}_x}}{n}$$



- Kernel density estimate aka “Parzen/moving window method”

$$\hat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^n \mathbf{1}_{||X_j - x|| \leq \Delta}}{n}$$

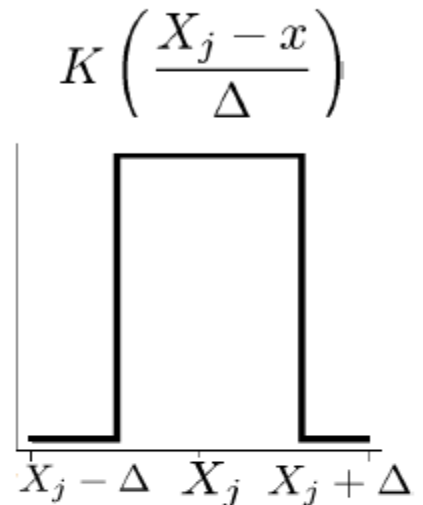
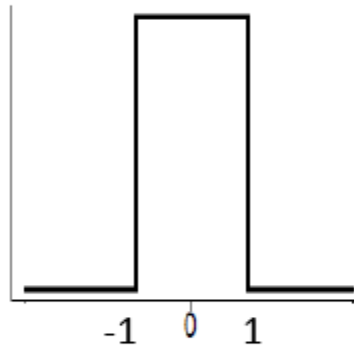


# Kernel density estimate

•  $\hat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^n K\left(\frac{X_j - x}{\Delta}\right)}{n}$  more generally

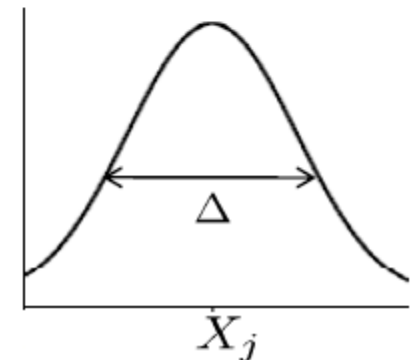
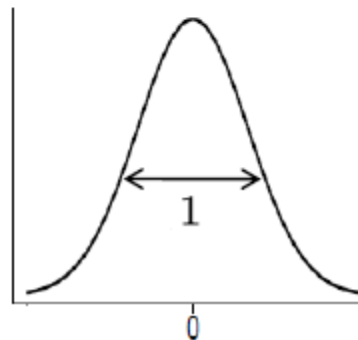
boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$



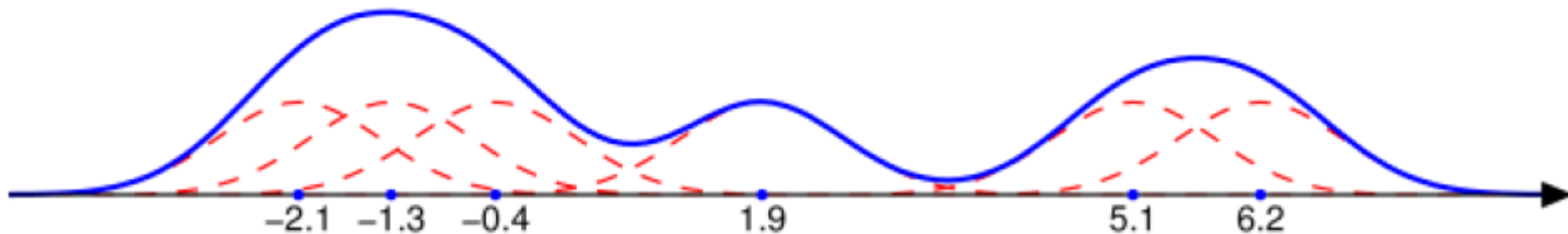
Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



# Kernel density estimate

- Place small "bumps" at each data point, determined by the kernel function.
- The estimator consists of a (normalized) "sum of bumps".



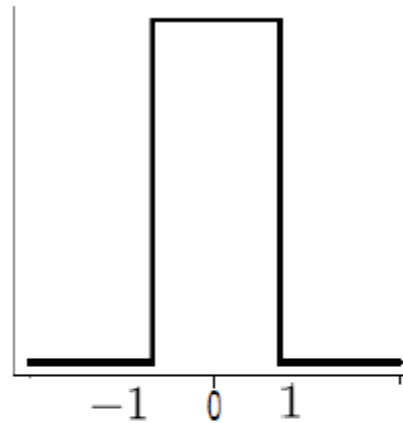
Gaussian bumps (red) around six data points and their sum (blue)

- Note that where the points are denser the density estimate will have higher values.

# Kernels

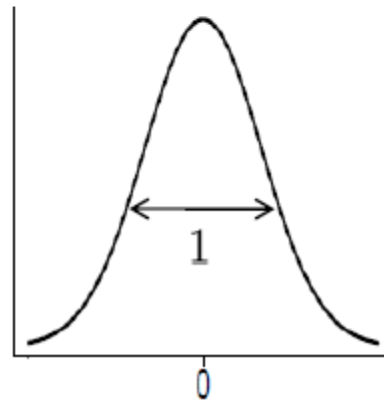
boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$



Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



Any kernel function that satisfies

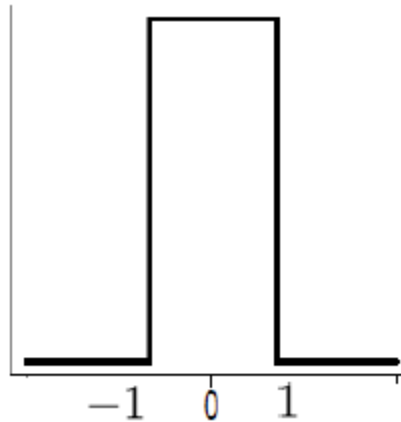
$$K(x) \geq 0,$$

$$\int K(x)dx = 1$$

# Kernels

boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$

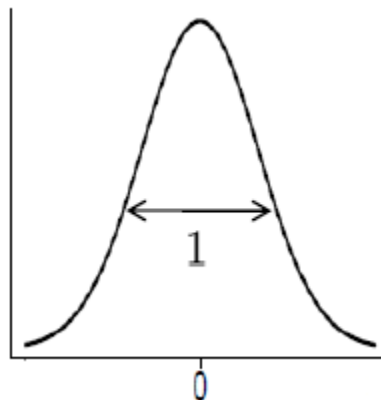


**Finite support**

- only need local points to compute estimate

Gaussian kernel :

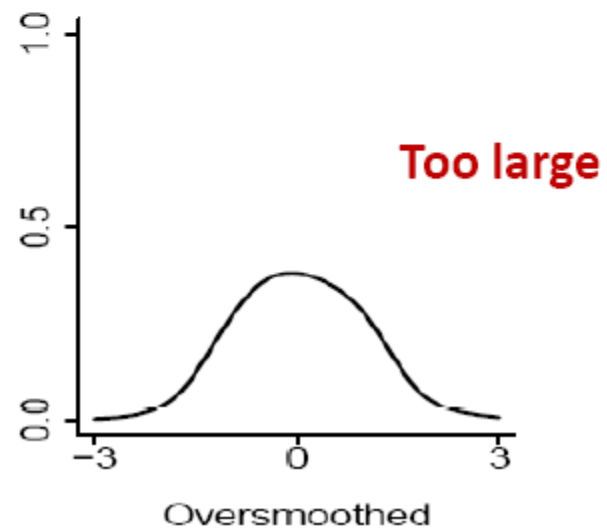
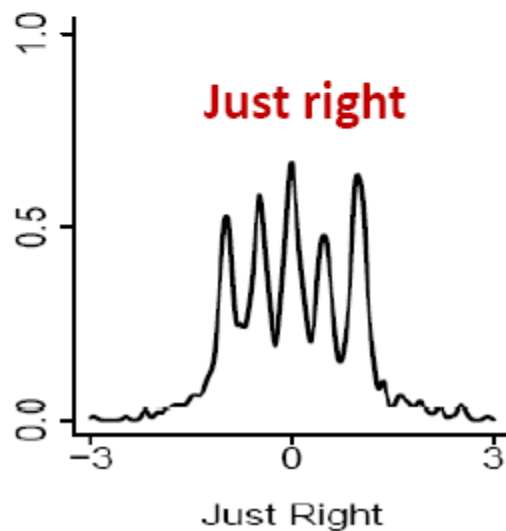
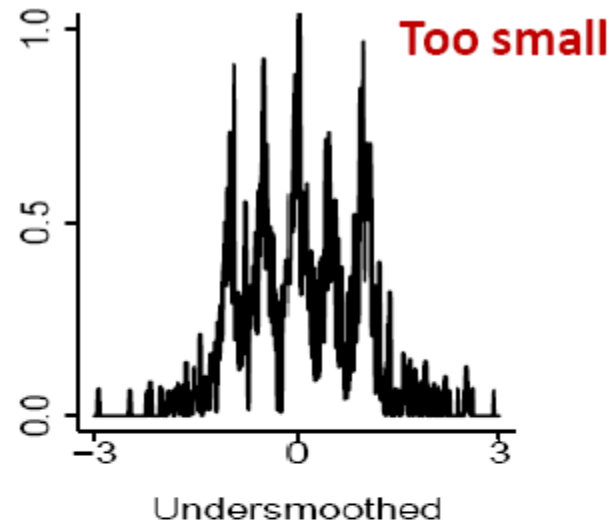
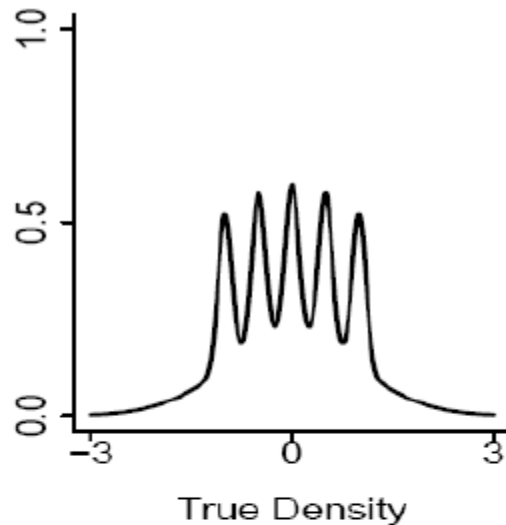
$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



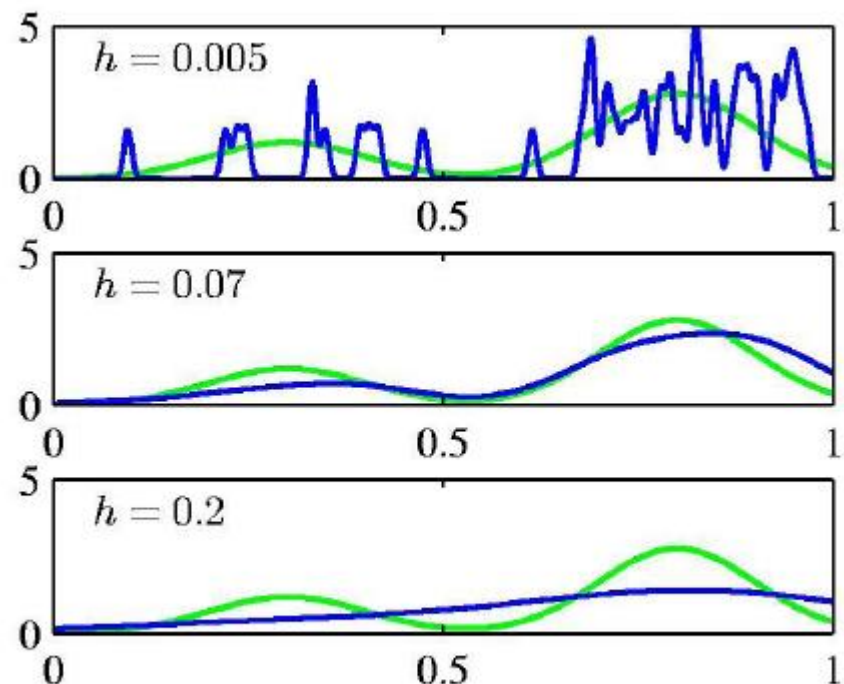
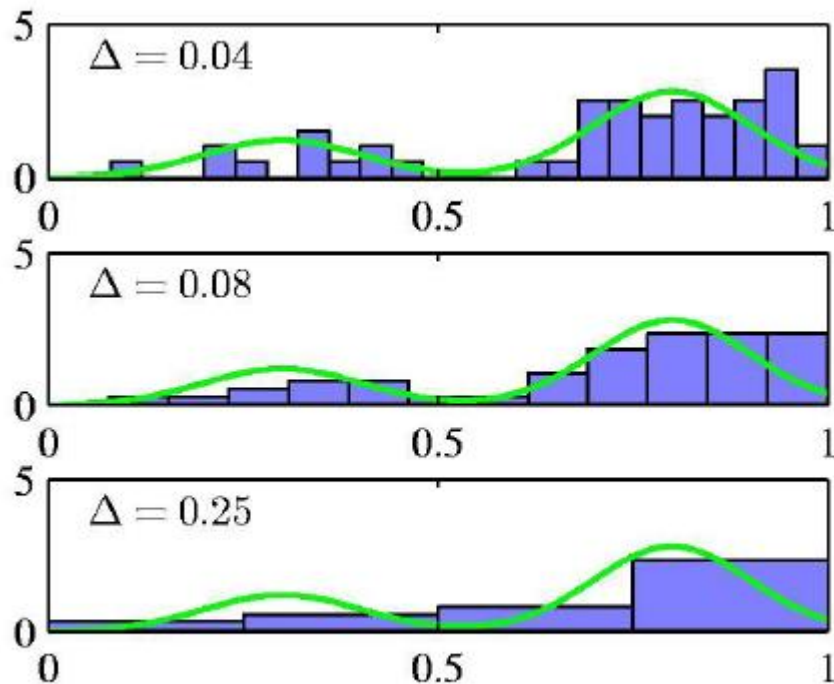
**Infinite support**

- need all points to compute estimate
- But quite popular since smoother

# Choice of kernel bandwidth



# Histograms vs. Kernel density estimation<sup>16</sup>



$\Delta = h$  acts as a smoother.



# K-NN(Nearest Neighbor) density estimation

- Histogram

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

- Kernel density est

$$\hat{p}(x) = \frac{n_x}{n\Delta}$$

Fix  $\Delta$ , estimate number of points within  $\Delta$  of  $x$  ( $n_i$  or  $n_x$ ) from data

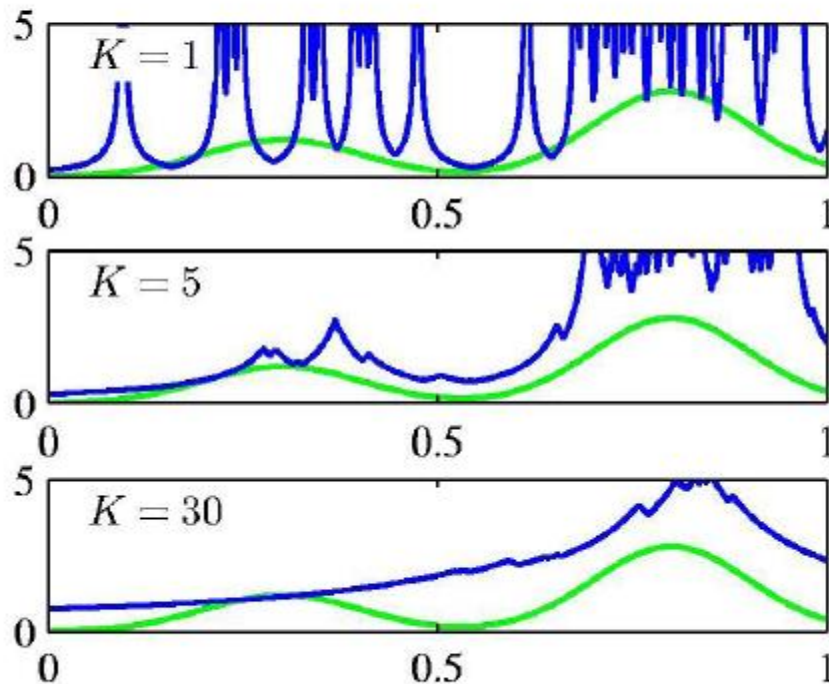
Fix  $n_x = k$ , estimate  $\Delta$  from data (volume of ball around  $x$  that contains  $k$  training pts)

- k-NN density est

$$\hat{p}(x) = \frac{k}{n\Delta_{k,x}}$$

# K-NN density estimation

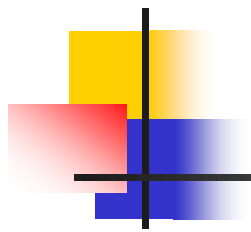
$$\hat{p}(x) = \frac{k}{n\Delta_{k,x}}$$



$k$  acts as a smoother.

Not very popular for density estimation - expensive to compute, bad estimates

But a related version for classification quite popular ...

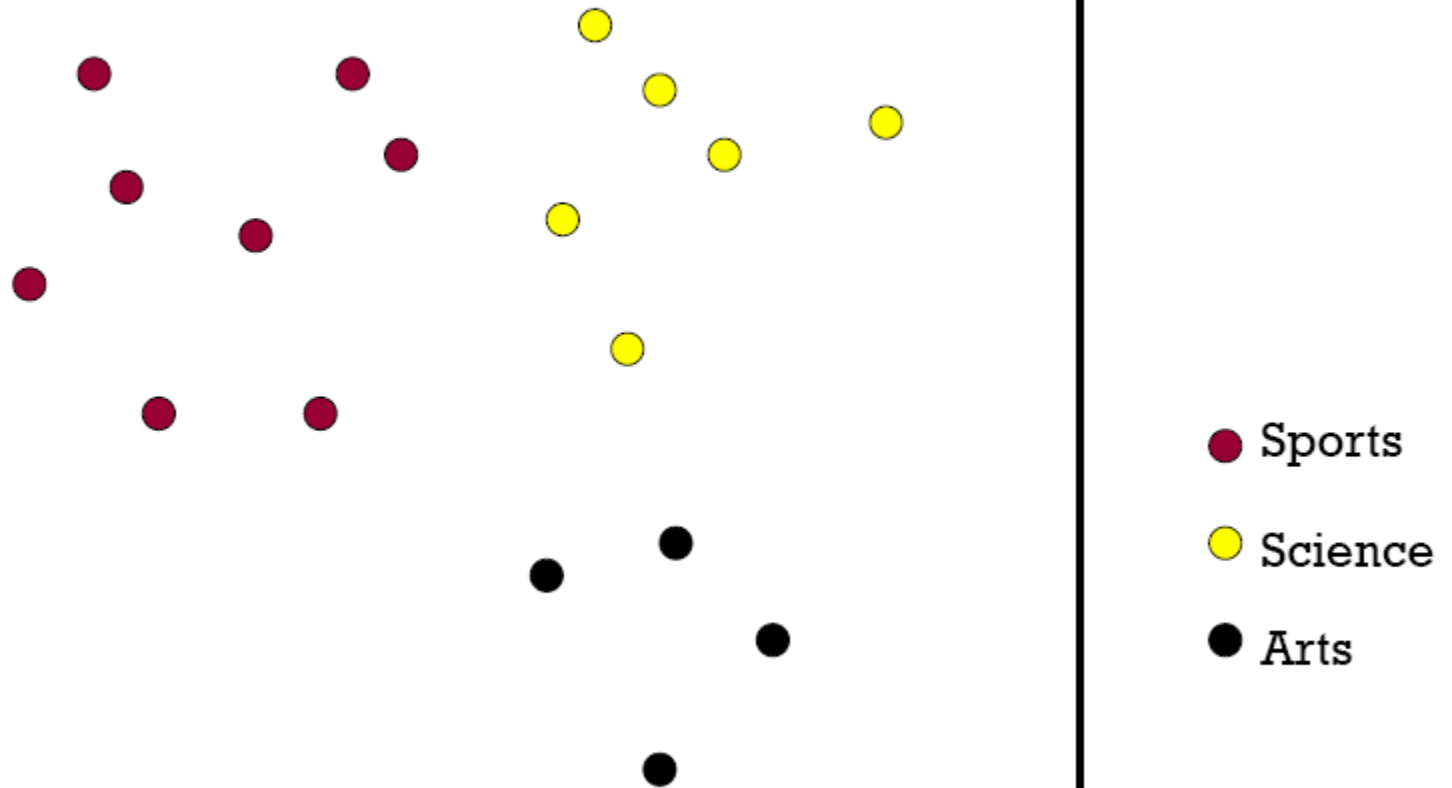


# From Density estimation to Classification



# K-NN classifier

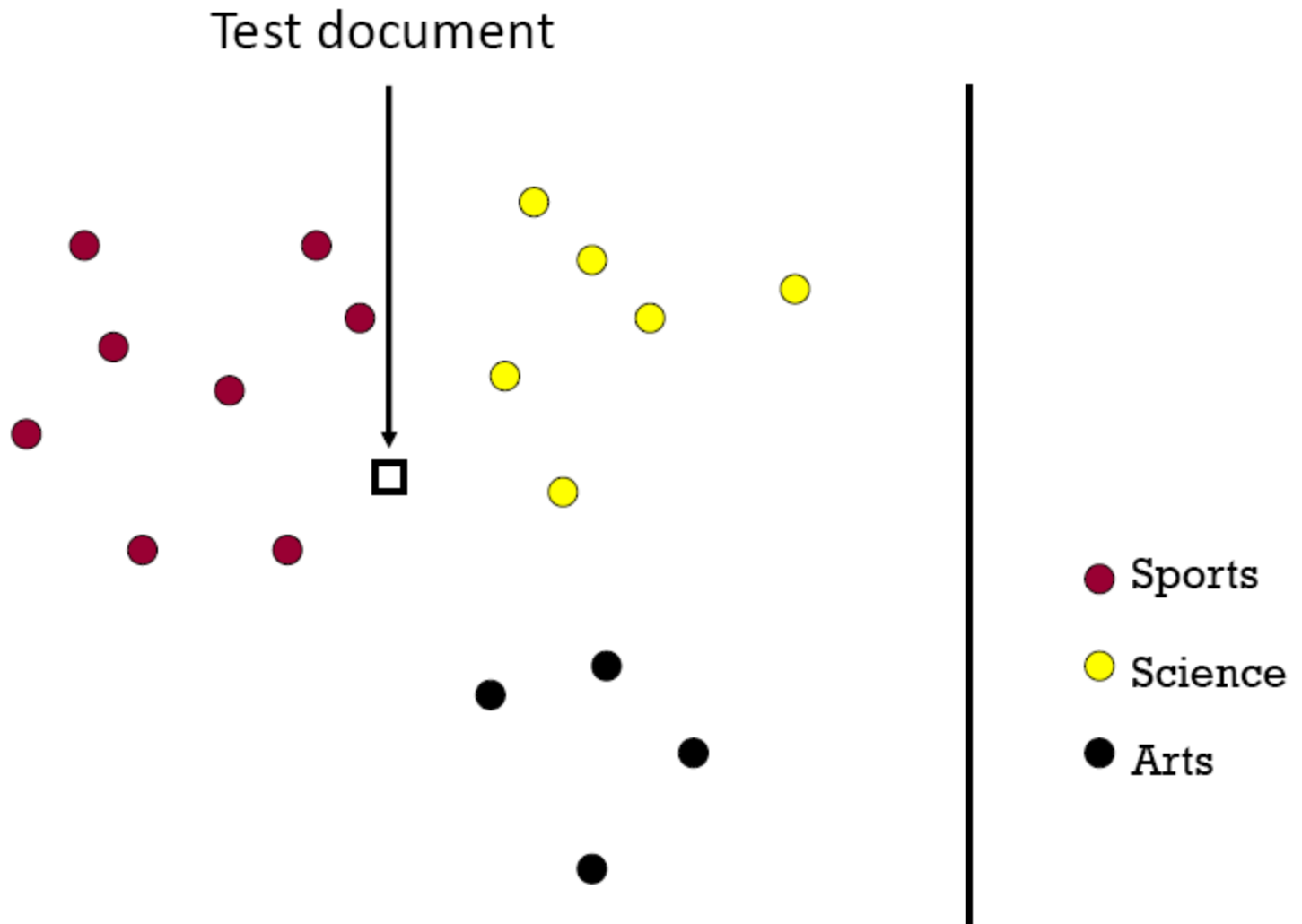
---



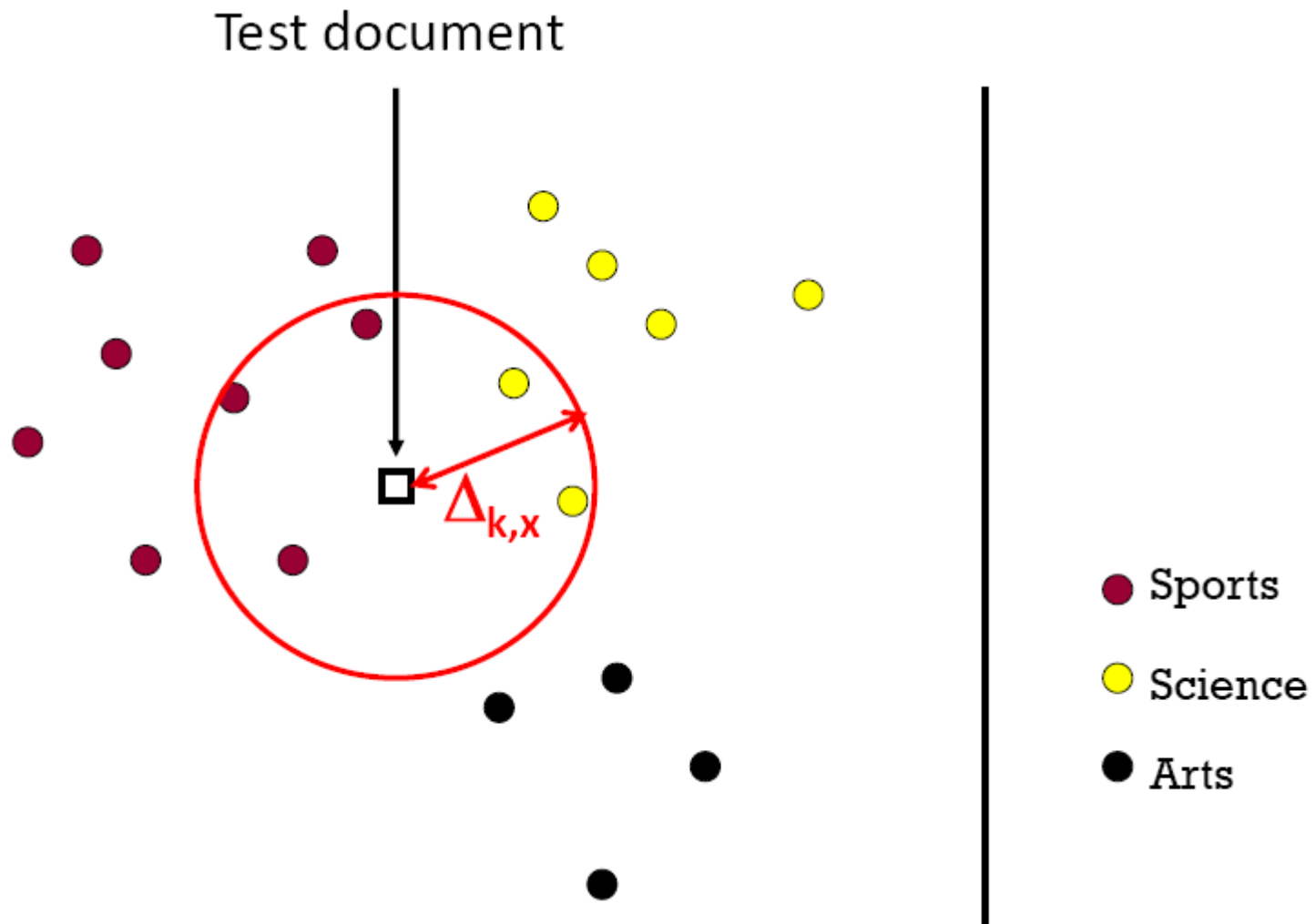


# K-NN classifier

---



# K-NN classifier (K=5)



What should we predict? ... Average? Majority? Why?

# K-NN classifier

- Optimal Classifier:  $f^*(x) = \arg \max_y P(y|x)$   
 $= \arg \max_y p(x|y)P(y)$
- k-NN Classifier:  $\hat{f}_{kNN}(x) = \arg \max_y \hat{p}_{kNN}(x|y)\hat{P}(y)$   
 $= \arg \max_y k_y$  **(Majority vote)**

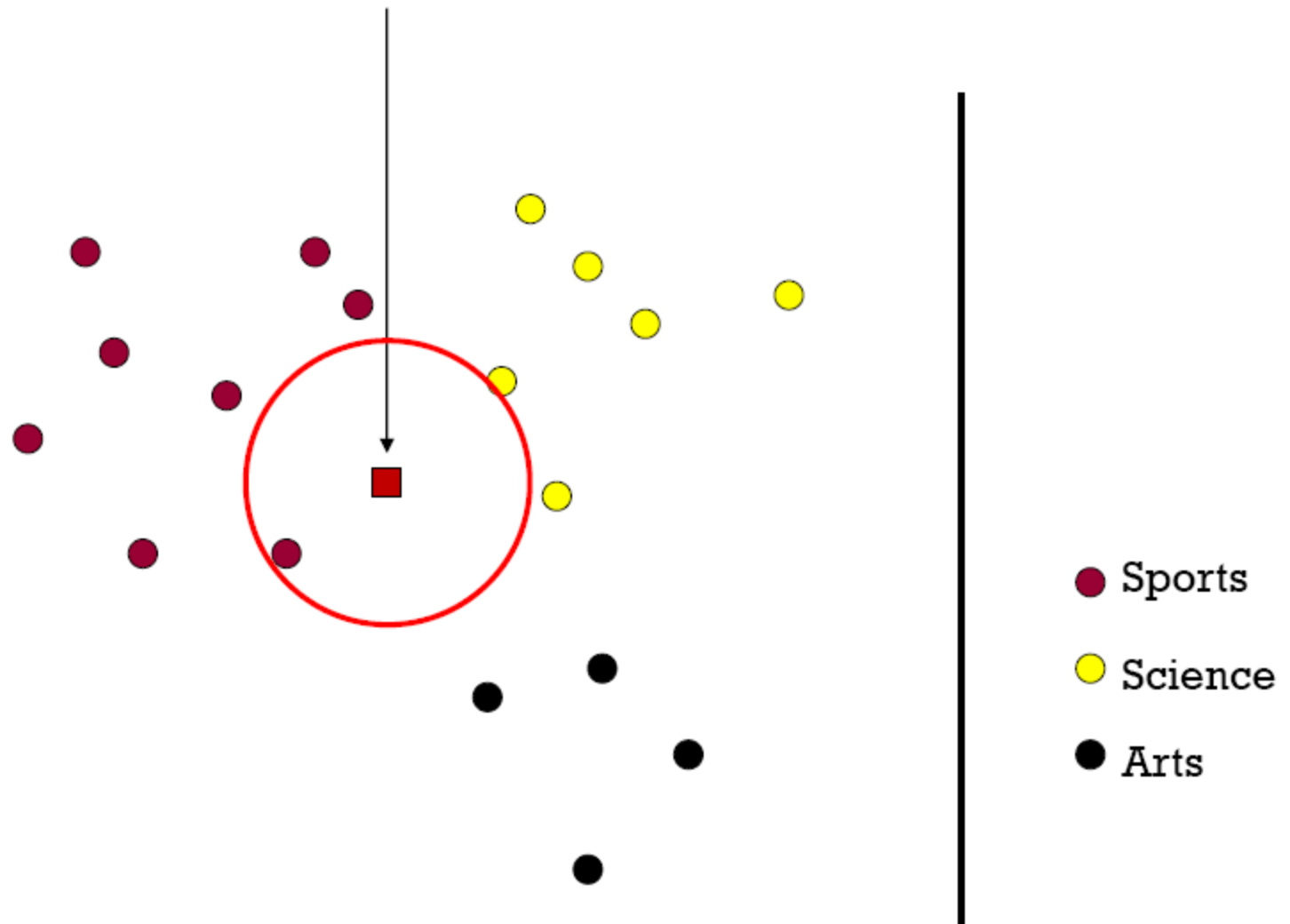
$$\hat{p}_{kNN}(x|y) = \frac{k_y}{n_y \Delta_{k,x}}$$

$\xrightarrow{\text{red arrow}} \# \text{ training pts of class } y \text{ that lie within } \Delta_k \text{ ball}$ 
 $\xrightarrow{\text{red arrow}} \# \text{ total training pts of class } y$

$\sum_y k_y = k$

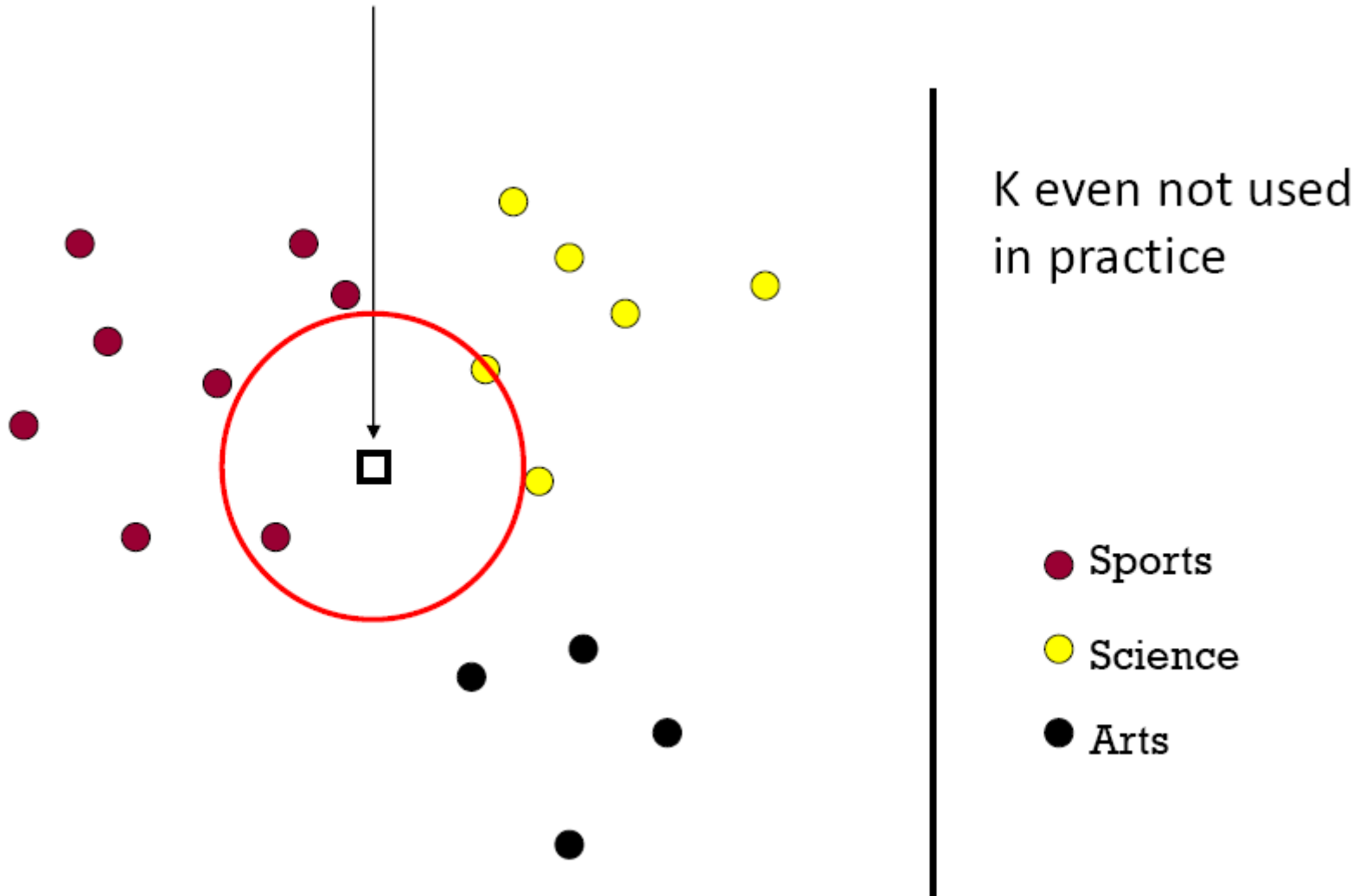
$$\hat{P}(y) = \frac{n_y}{n}$$

# 1-Nearest Neighbor (kNN) classifier

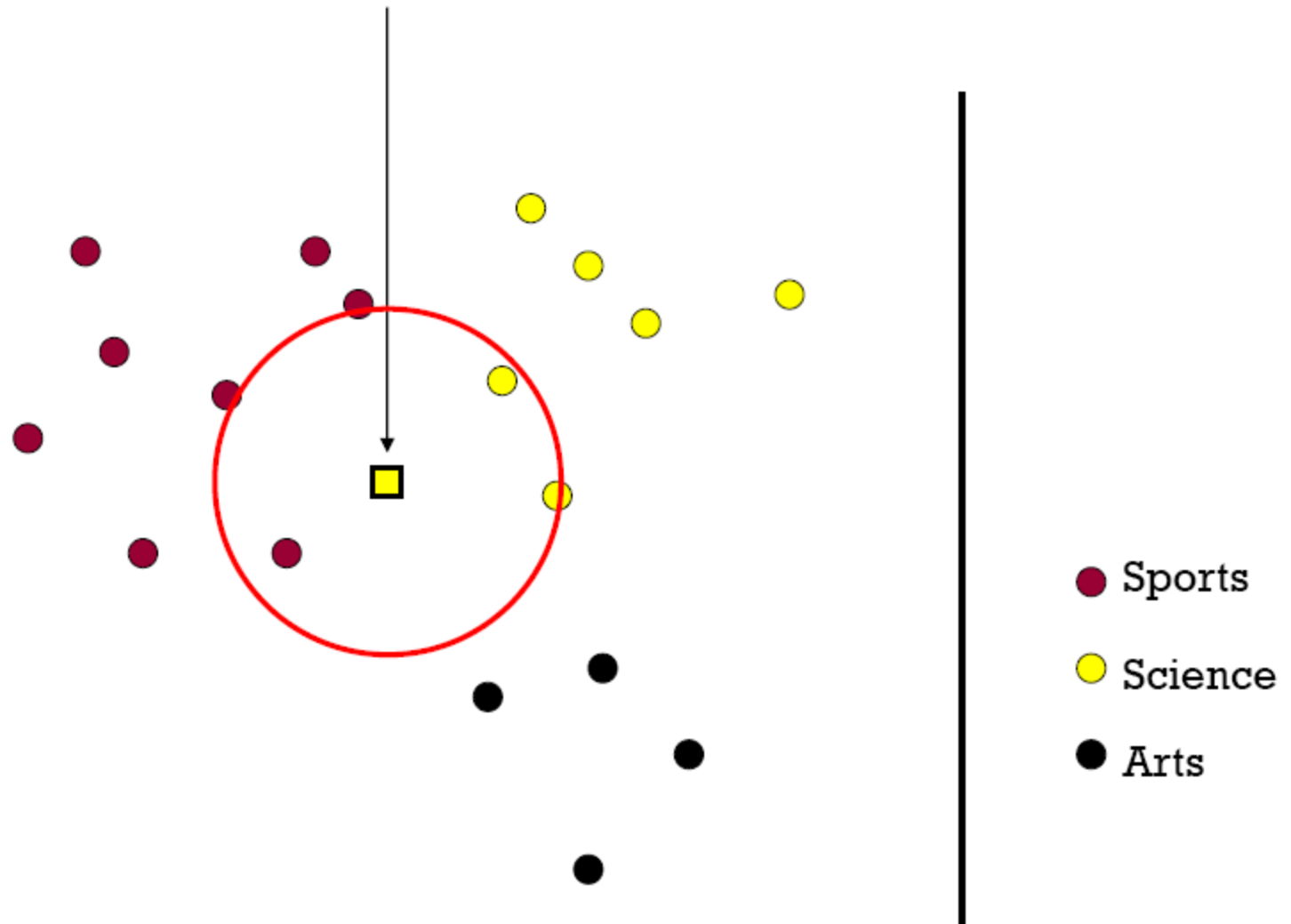




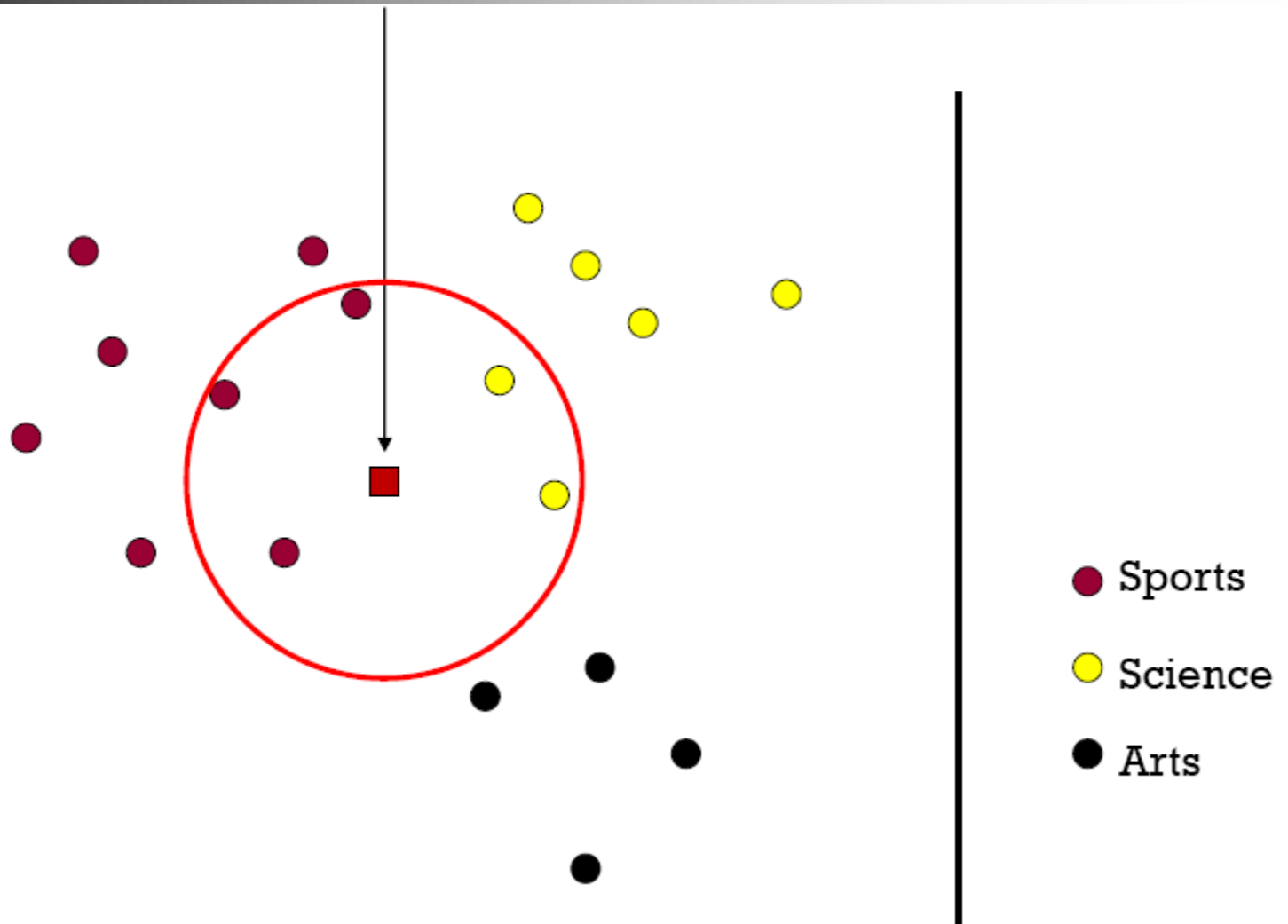
# 2-Nearest Neighbor (kNN) classifier



# 3-Nearest Neighbor (kNN) classifier



# 5-Nearest Neighbor (kNN) classifier





# What is the best K?

---

Bias-variance tradeoff

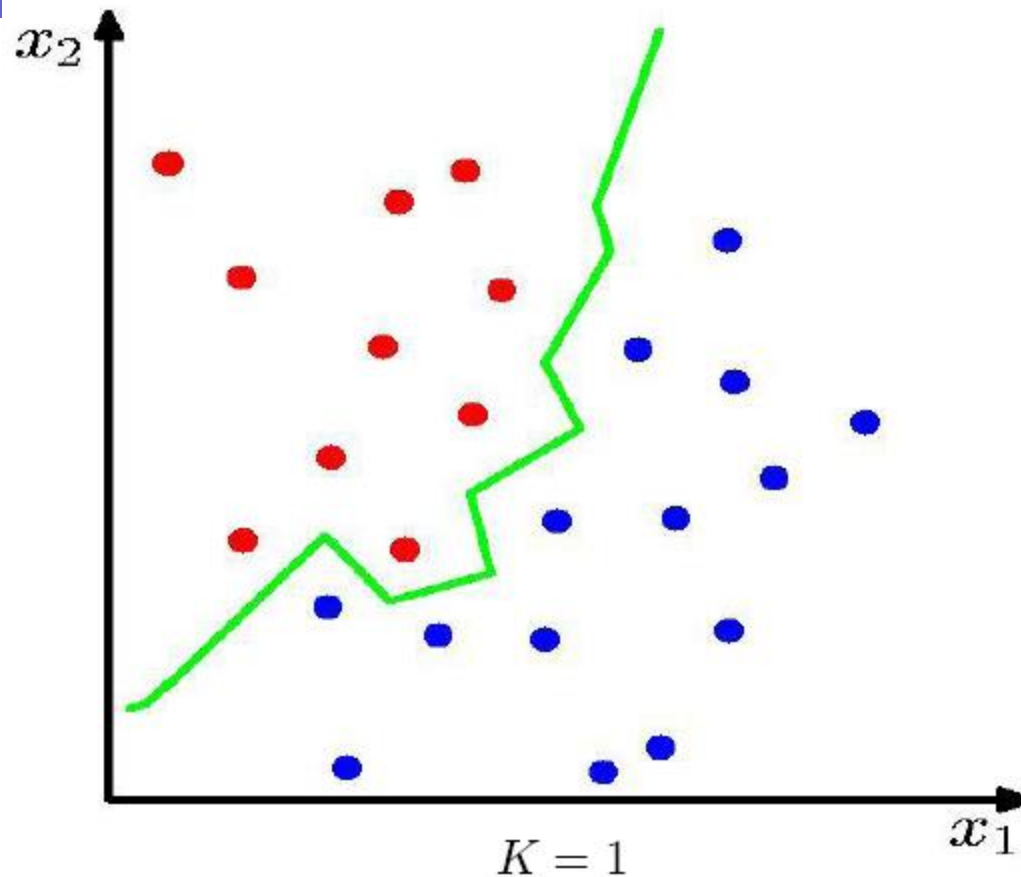
Larger  $K \Rightarrow$  predicted label is more stable

Smaller  $K \Rightarrow$  predicted label is more accurate

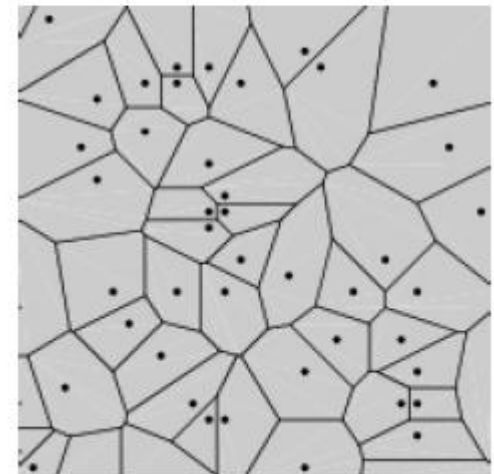
Similar to density estimation

Choice of  $K$  - in next class ...

# 1-NN classifier – decision boundary

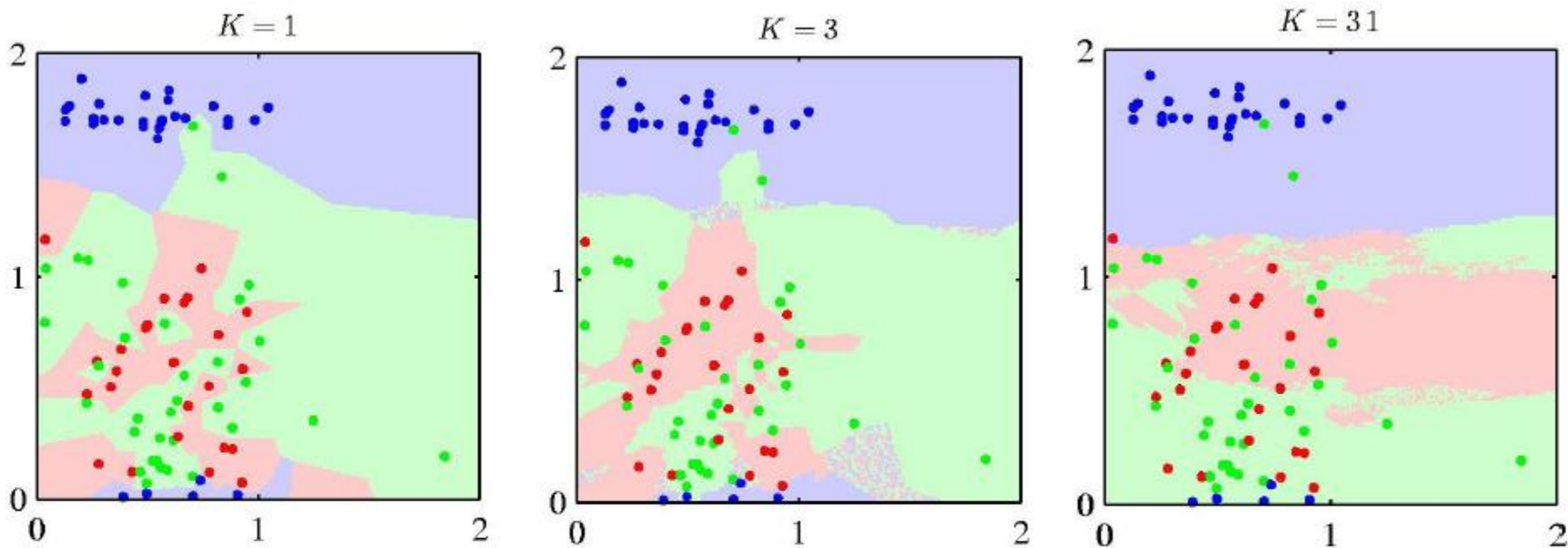
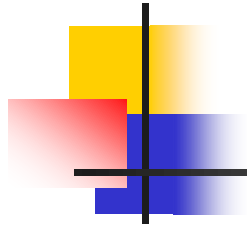


Voronoi  
Diagram



- Guarantee: For  $n \rightarrow \infty$ , the error rate of the 1-nearest-neighbour classifier is never more than twice the optimal error.

# k-NN classifier – decision boundary



- $K$  acts as a smoother (Bias-variance tradeoff)

# Case Study: kNN for Web Classification

- Dataset
  - 20 News Groups (20 classes)
  - Download :(<http://people.csail.mit.edu/jrennie/20Newsgroups/>)
  - 61,118 words, 18,774 documents
  - Class labels descriptions

comp.graphics comp.os.ms-windows.misc comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x	rec.autos rec.motorcycles rec.sport.baseball rec.sport.hockey	sci.crypt sci.electronics sci.med sci.space
misc.forsale	talk.politics.misc talk.politics.guns talk.politics.mideast	talk.religion.misc alt.atheism soc.religion.christian



# Experimental Setup

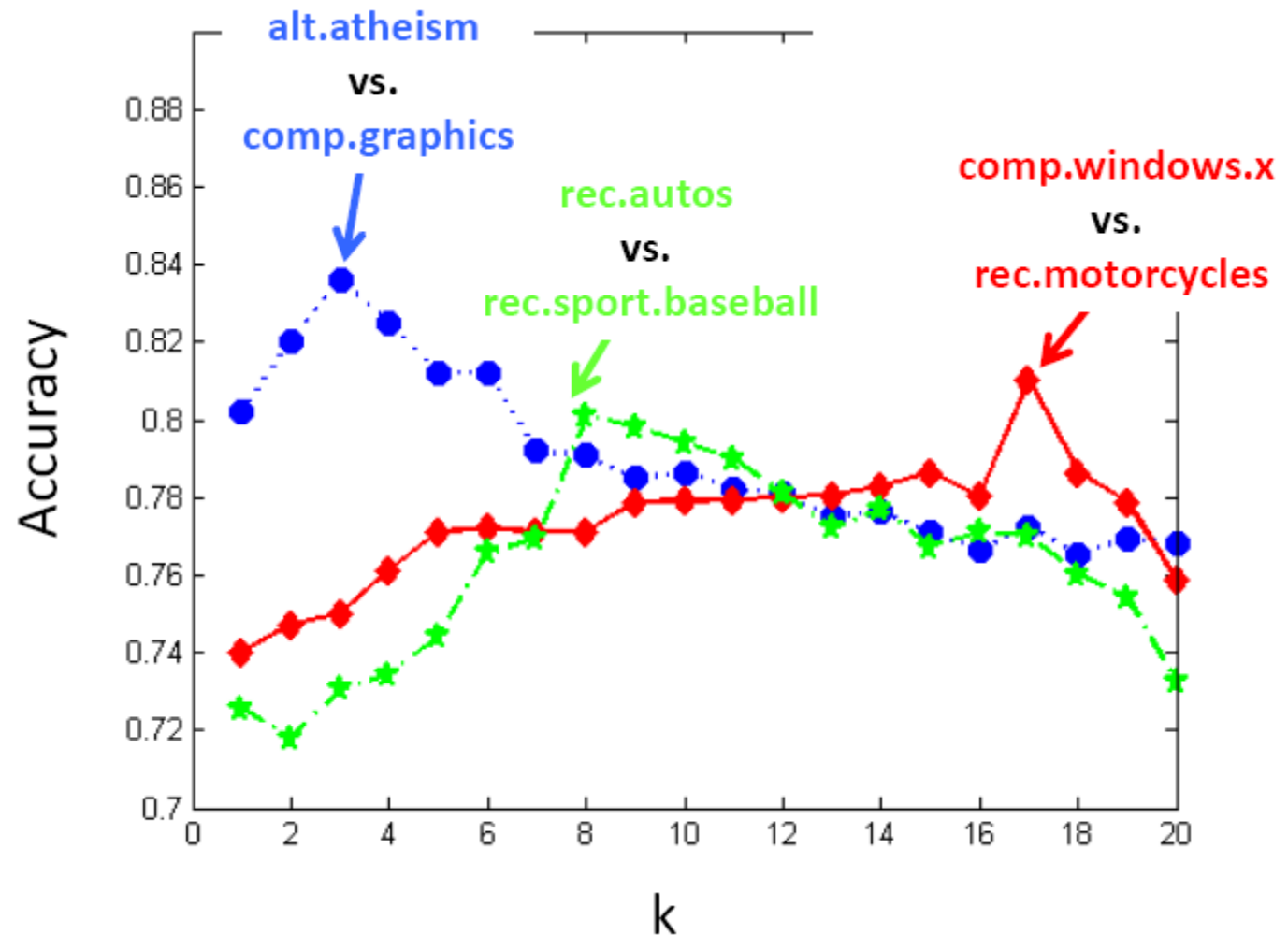
---

- Training/Test Sets:
  - 50%-50% randomly split.
  - 10 runs
  - report average results
- Evaluation Criteria:

$$Accuracy = \frac{\sum_{i \in \text{test set}} I(\text{predict}_i = \text{true label}_i)}{\# \text{ of test samples}}$$



# Results: Binary Classes





# Summary

---

- Instance based/non-parametric approaches

**Four things make a memory based learner:**

1. *A distance metric,  $\text{dist}(x, X_i)$*   
**Euclidean (and many more)**
2. *How many nearby neighbors/radius to look at?*  
 **$k, \Delta/h$**
3. *A weighting function (optional)*  
**W based on kernel K**
4. *How to fit with the local points?*  
**Average, Majority vote, Weighted average, Poly fit**



# Summary

---

- Parametric vs Nonparametric approaches
  - Nonparametric models place very mild assumptions on the data distribution and provide good models for complex data  
Parametric models rely on very strong (simplistic) distributional assumptions
  - Nonparametric models (not histograms) requires storing and computing with the entire data set.  
Parametric models, once fitted, are much more efficient in terms of storage and computation.



# What you should know...

---

- Histograms, Kernel density estimation
  - Effect of bin width/ kernel bandwidth
  - Bias-variance tradeoff
- K-NN classifier
  - Nonlinear decision boundaries