# Non-parametric methods

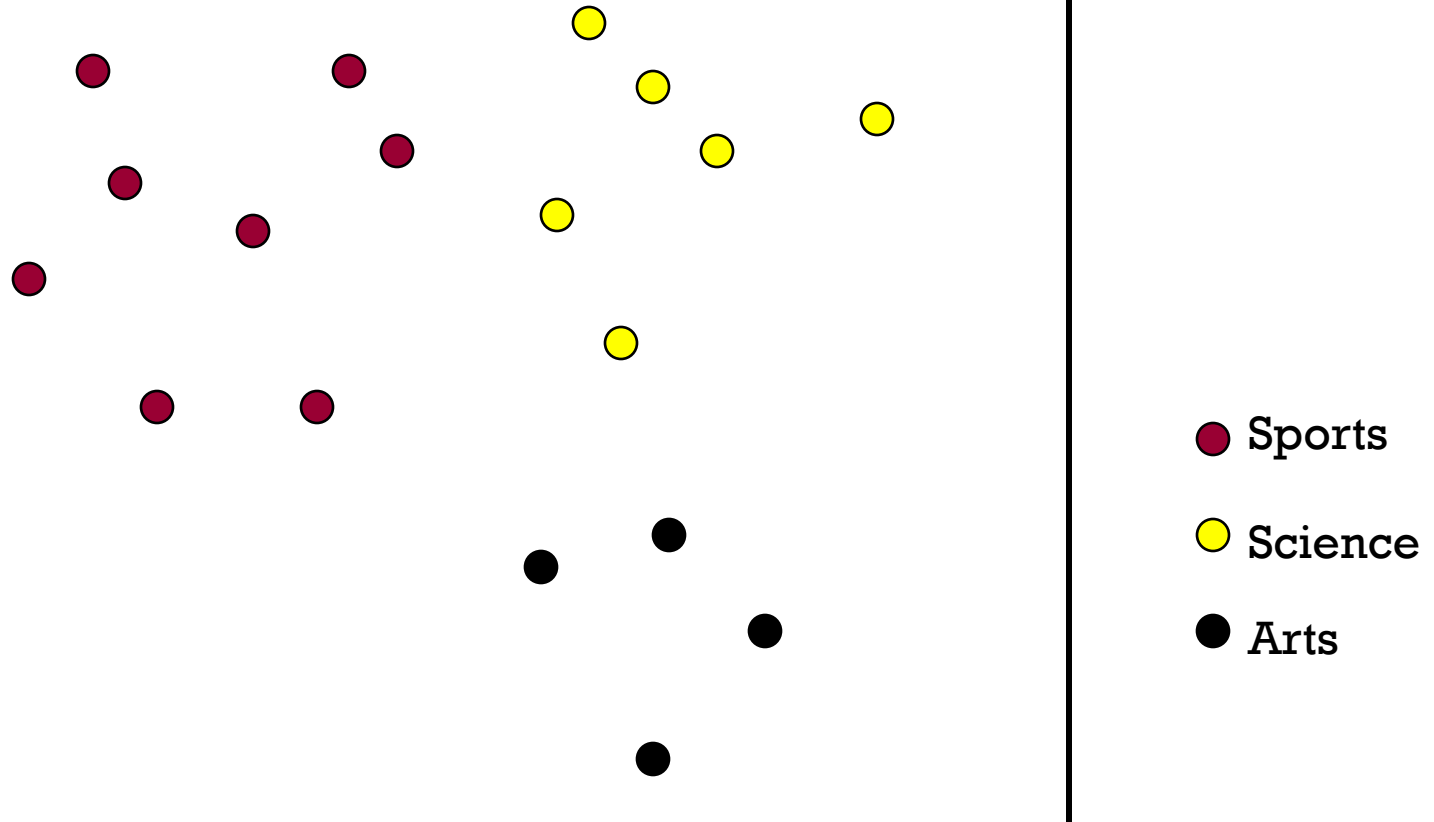Aarti Singh

Machine Learning 10-315
Oct 12, 2020

# 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
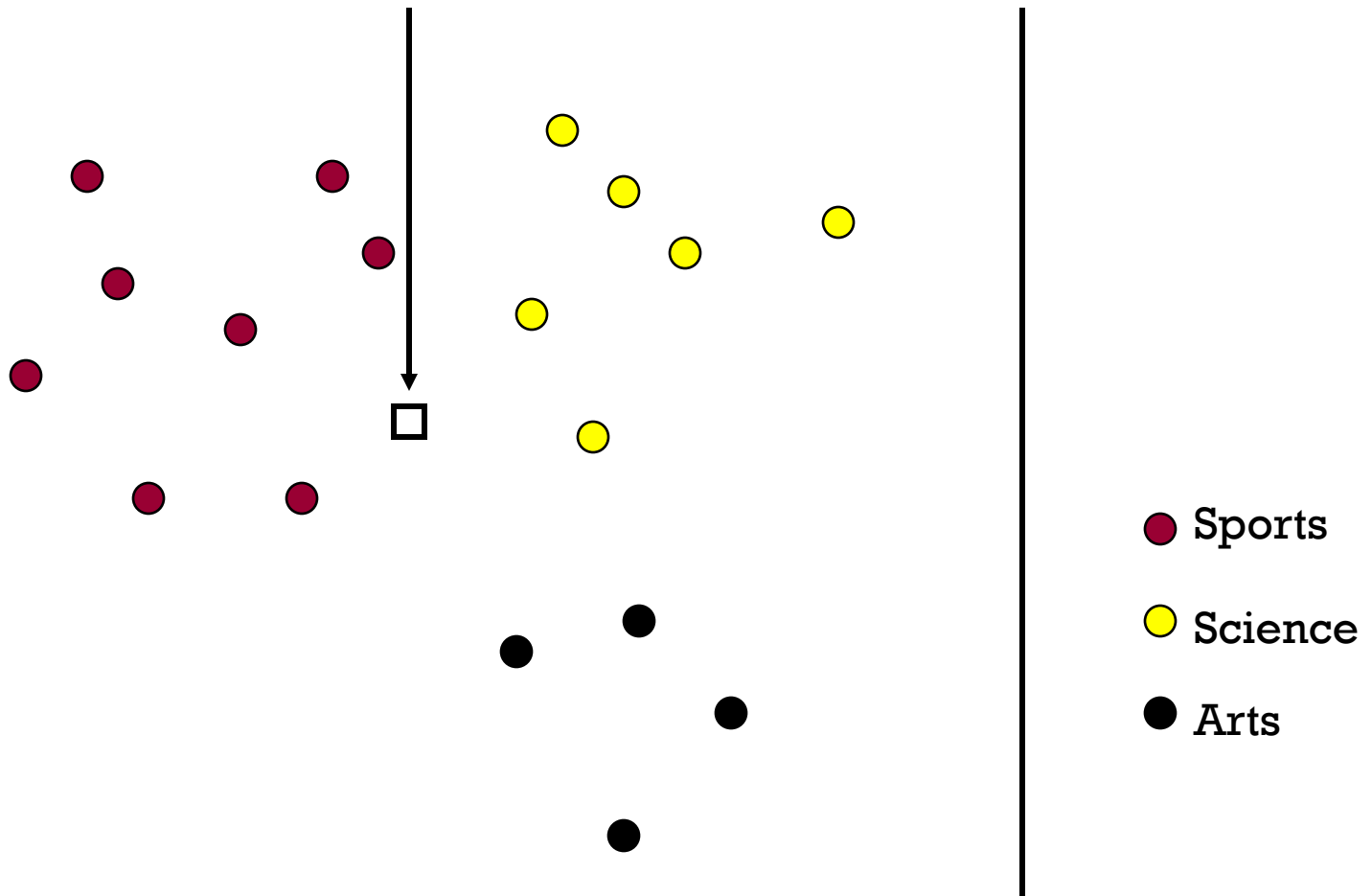
# Non-Parametric methods

- 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 today (more later)

  **Classification:** k-NN (k-Nearest Neighbor) classifier

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

  **Regression:** Kernel regression
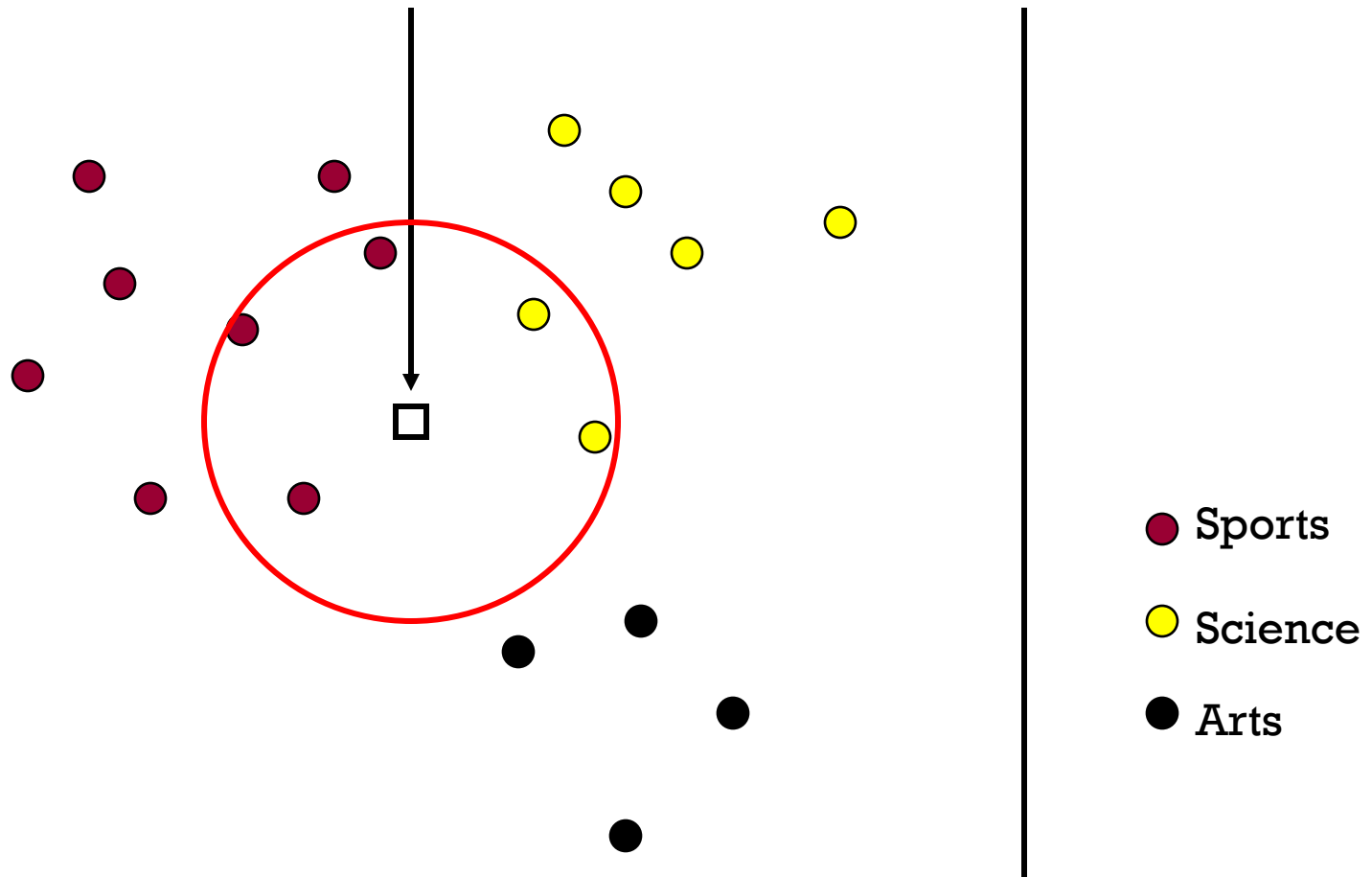
# k-NN classifier



Sports

Science

Arts

# k-NN classifier

Test document

# k-NN classifier (k=5)



Test document

Sports
Science
Arts

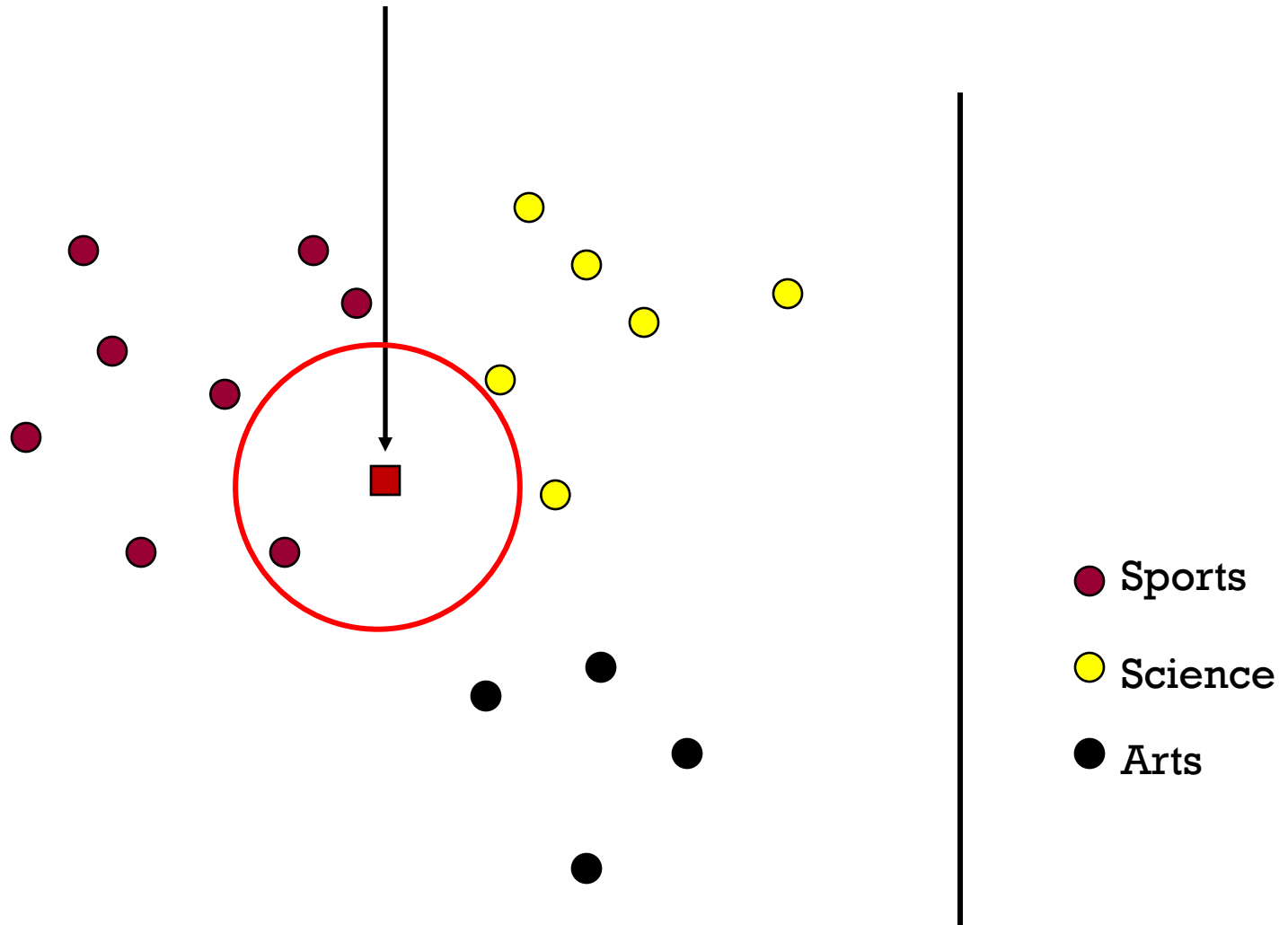What should we predict? … Average? Majority? Why?

# k-NN classifier

- Optimal Classifier:

$$f^*(x) = \arg\max_y P(y|x)$$

$$= \arg\max_y P(x|y)P(y)$$

- k-NN Classifier:

$$\widehat{f}_{kNN}(x) = \arg\max_y \widehat{P}_{kNN}(x|y)\widehat{P}(y)$$

$$= \arg\max_y k_y$$

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

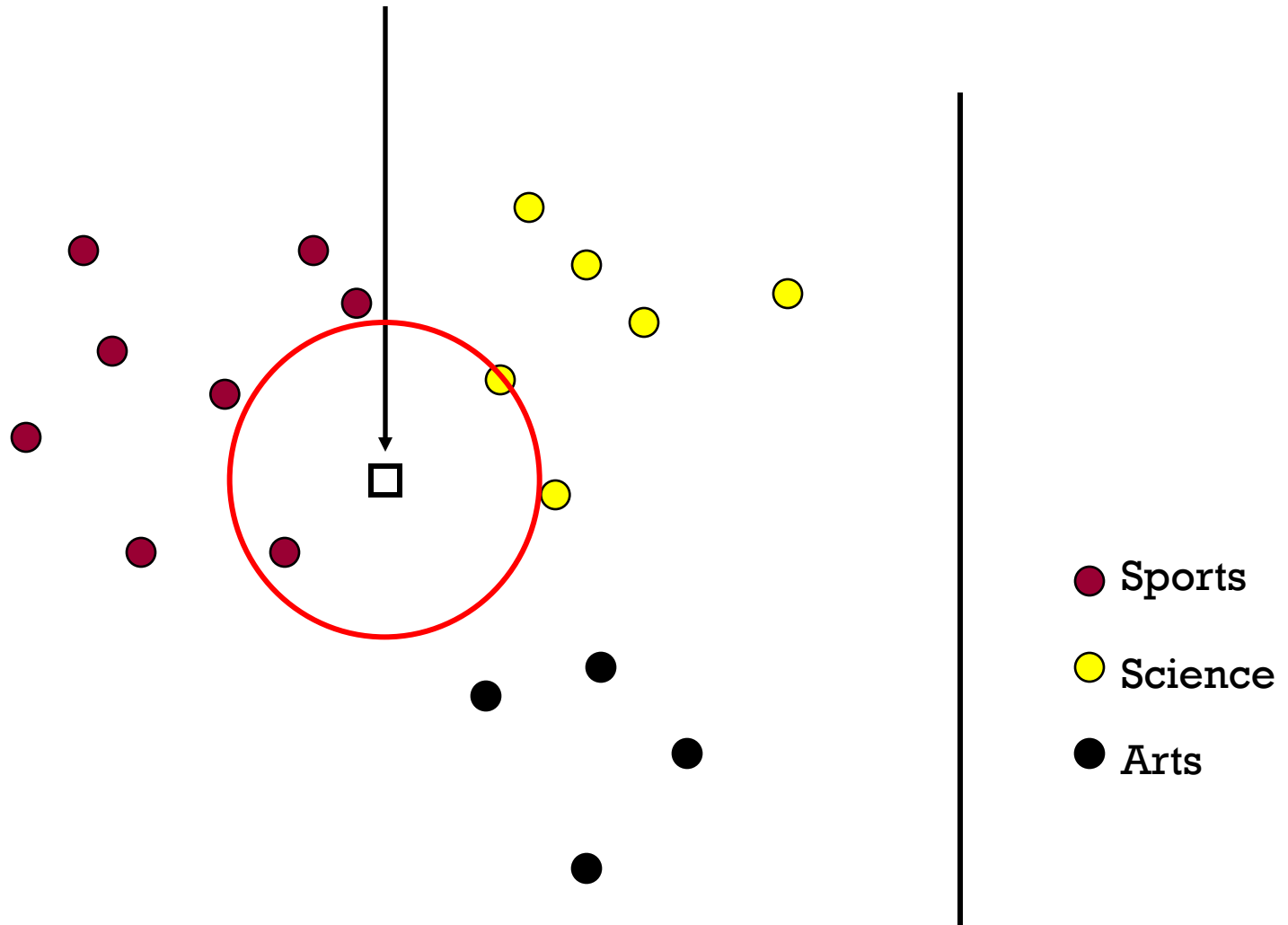$$\longrightarrow \textbf{\# total training pts of class y}$$
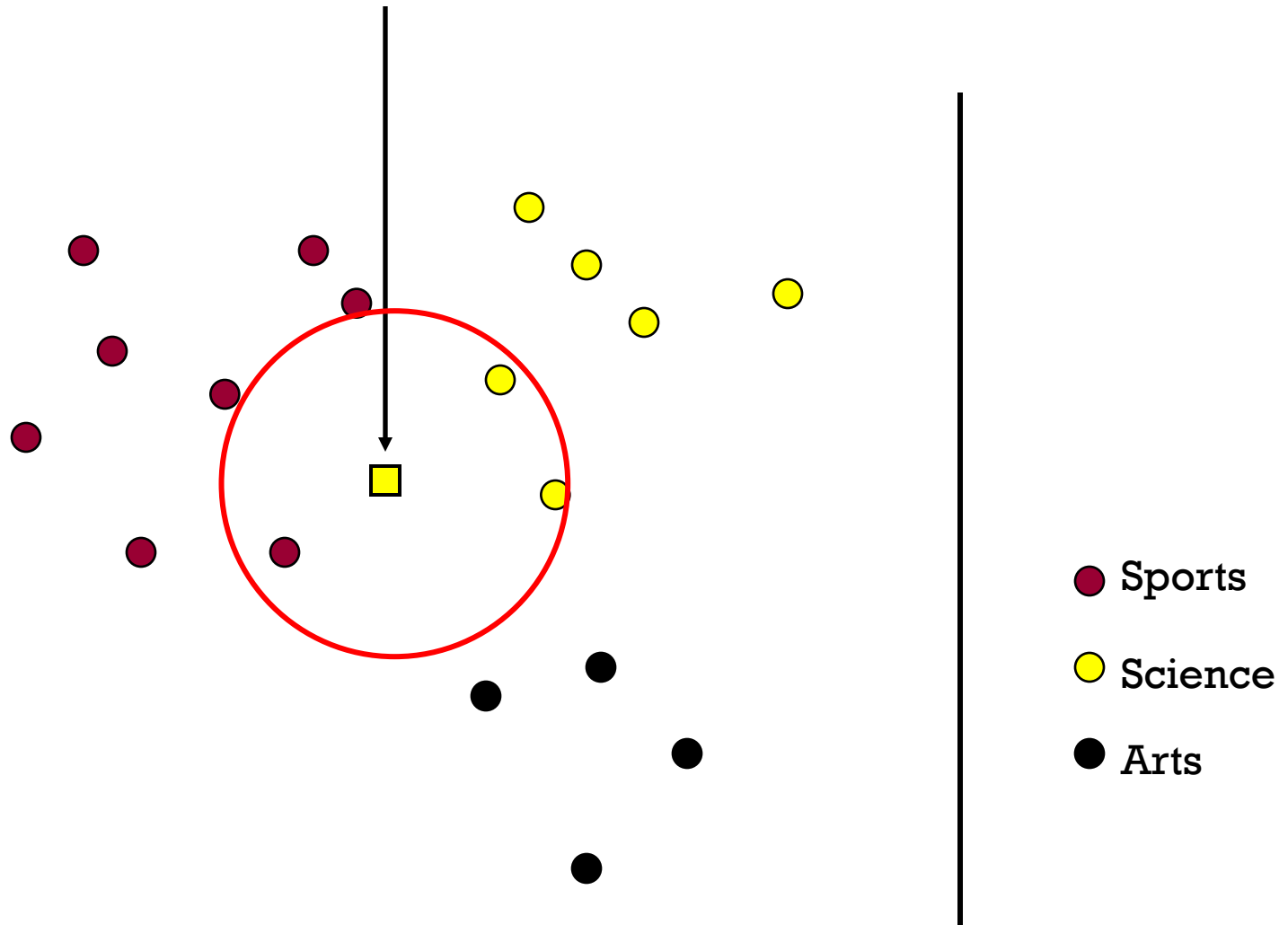
$$\widehat{P}(y) = \frac{n_y}{n}$$
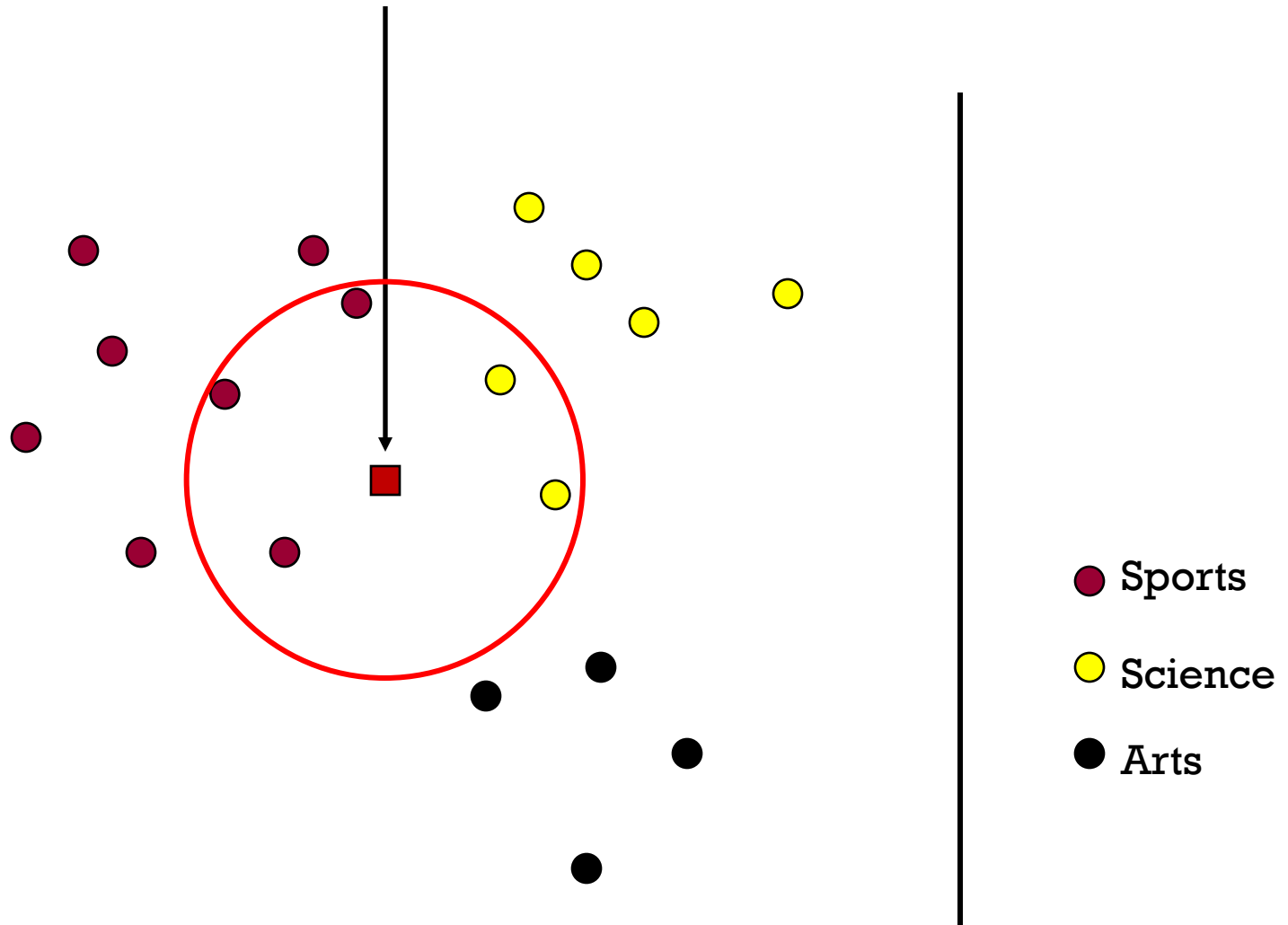
# 1-Nearest Neighbor (kNN) classifier



Sports

Science

Arts

# 2-Nearest Neighbor (kNN) classifier

Sports

Science

Arts

# 3-Nearest Neighbor (kNN) classifier



Sports

Science

Arts

# 5-Nearest Neighbor (kNN) classifier

Sports
Science
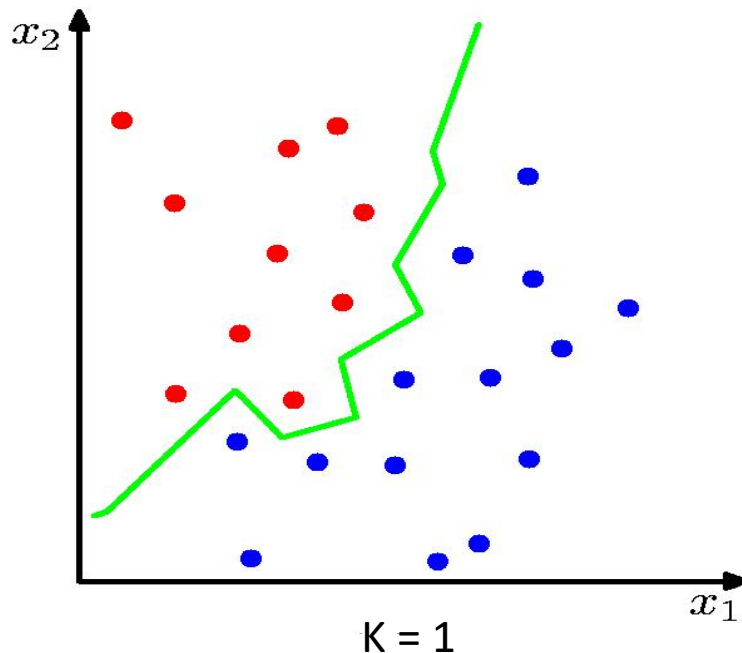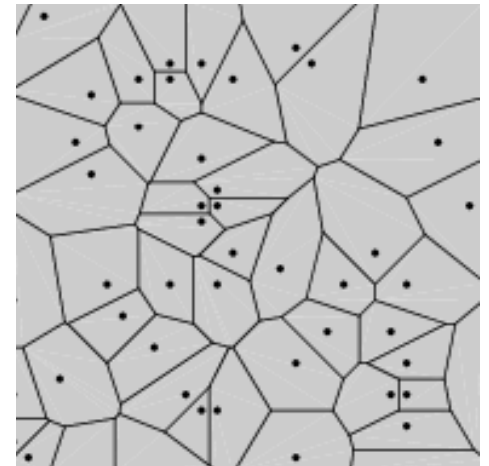Arts

# What is the best k?

1-NN classifier decision boundary

Voronoi Diagram

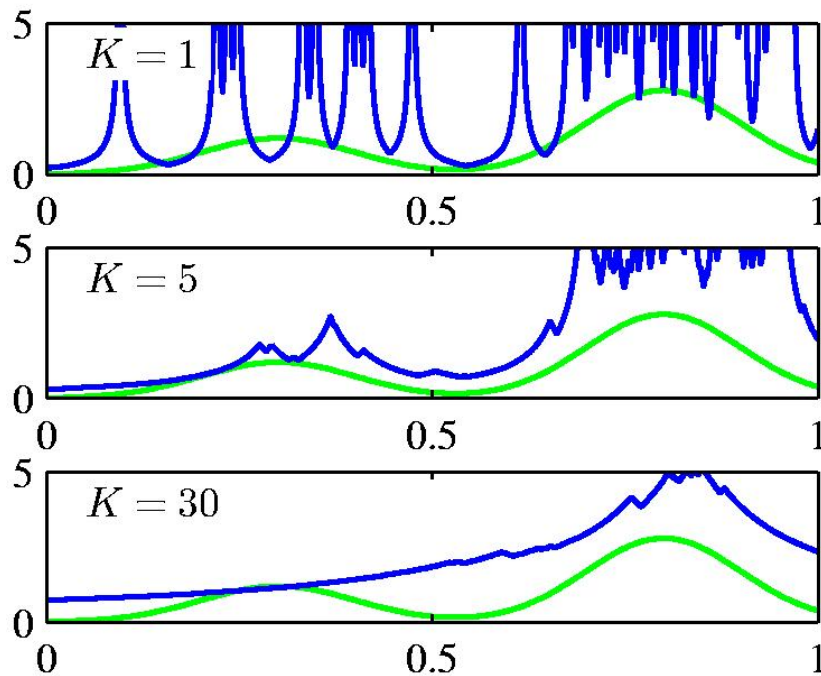

K = 1

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

# What is the best k?

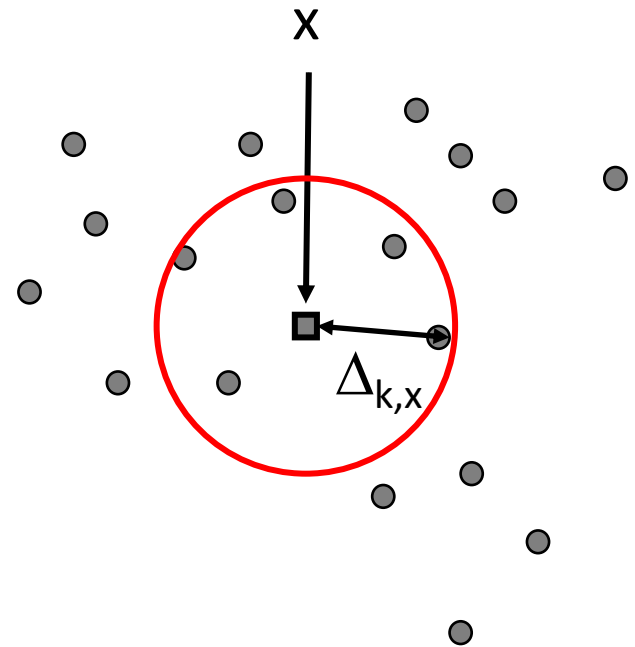Approximation vs. Stability (aka Bias vs Variance) Tradeoff

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

# k-NN density estimation

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



x

$\Delta_{k,x}$



$K = 1$

$K = 5$

$K = 30$

k acts as a smoother.

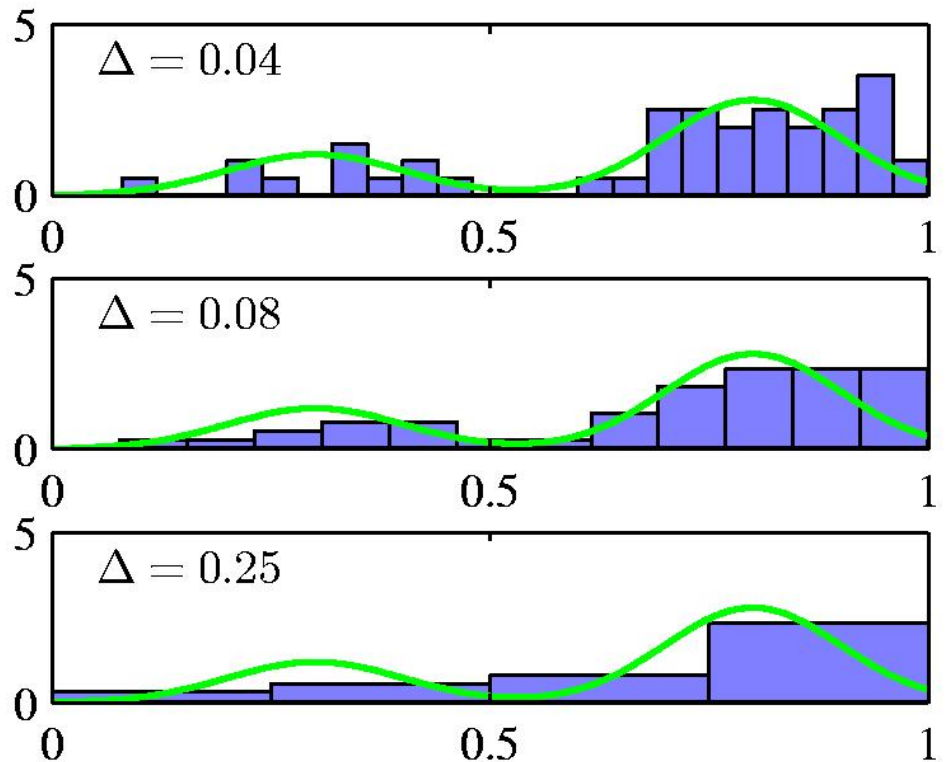Not very popular for density estimation – spiked estimates

# 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.

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

"Local relative frequency"

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



Image src: Bishop book

15

# Effect of histogram bin width

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

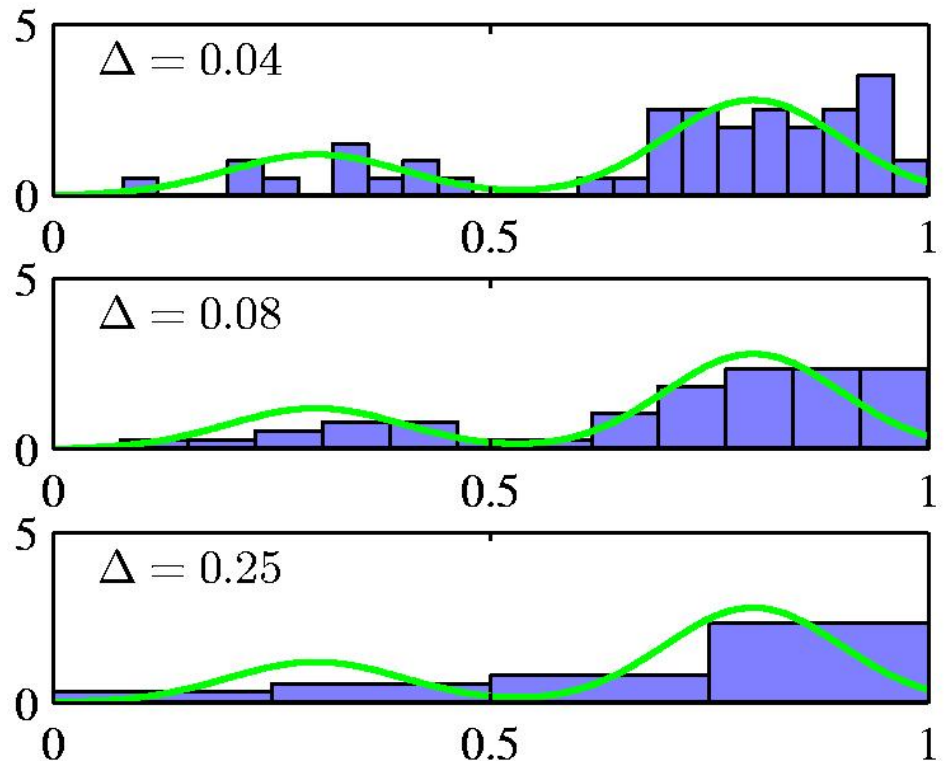# bins = $1/\Delta$

Small $\Delta$, large #bins
Good fit but unstable
(few points per bin)
"**Small bias**, **Large variance**"

Large $\Delta$, 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

  Parameters $p_j$ : density in bin j

  Note $\sum_j p_j = 1/\Delta$ since $\int p(x)dx = 1$

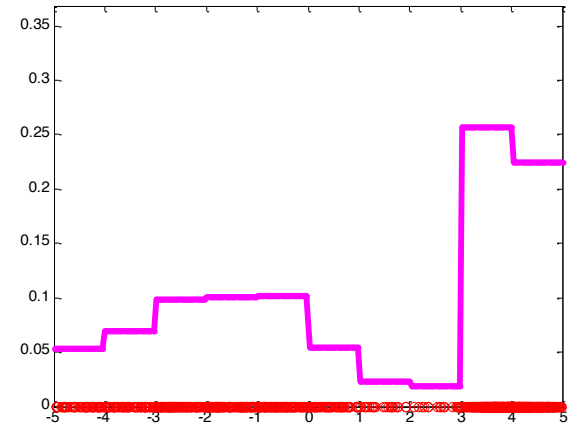- Maximize likelihood of data under probability model with parameters $p_j$

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

- Show that histogram density estimate is MLE under this model
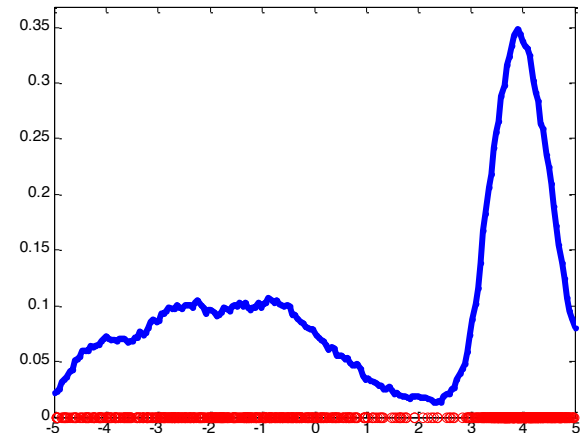
# Kernel density estimate

- Histogram – blocky estimate

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



- Kernel density estimate aka "Parzen/moving window method"

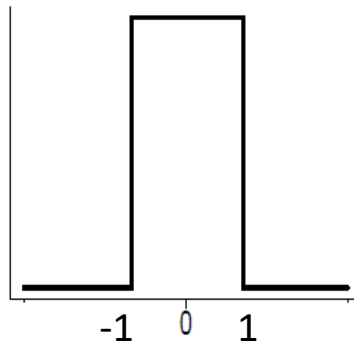$$\widehat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^{n} \mathbf{1}_{||X_j - x|| \leq \Delta}}{n}$$
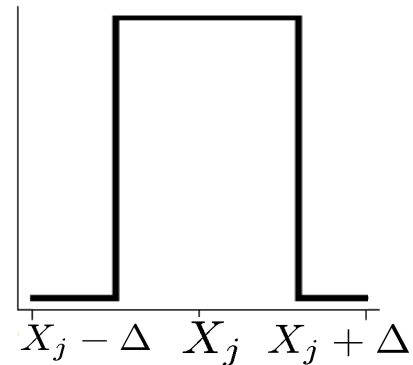
# Kernel density estimate

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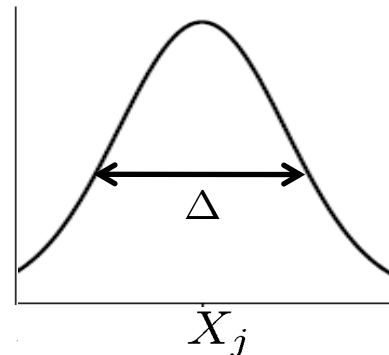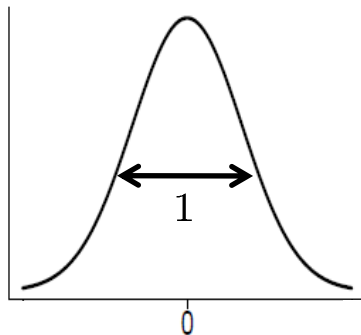
$$K\left(\frac{X_j - x}{\Delta}\right)$$

boxcar kernel :

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

Gaussian kernel :

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

# Kernels

boxcar kernel :

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

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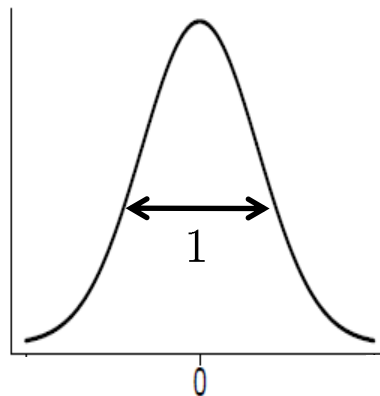
$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$

Any kernel
function that
satisfies

$$K(x) \geq 0,$$
$$\int K(x)dx = 1$$

# Kernel density estimation

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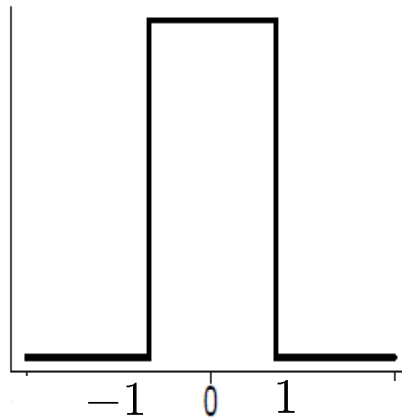
Img src: Wikipedia



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.

# Choice of Kernels

boxcar kernel :

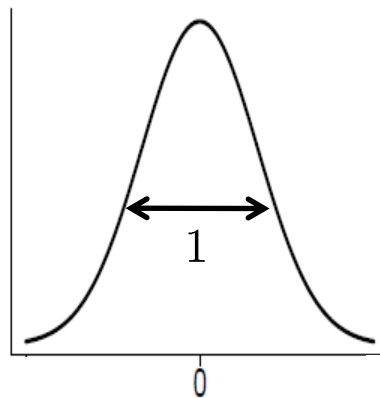$$K(x) = \frac{1}{2}I(x),$$



Finite support
– only need local
  points to compute
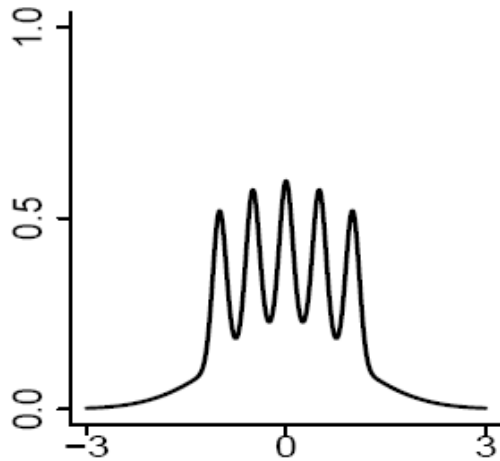  estimate

Gaussian kernel :

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



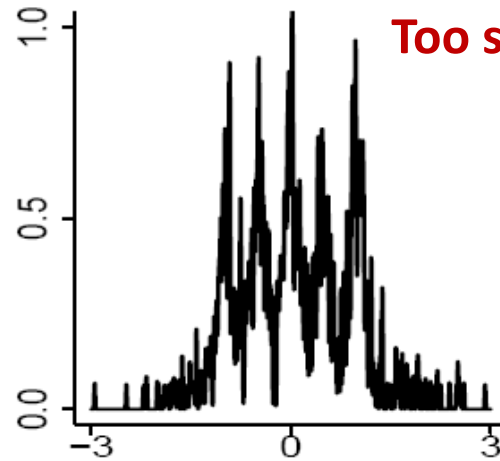Infinite support
- need all points to
  compute estimate
-But quite popular
  since smoother
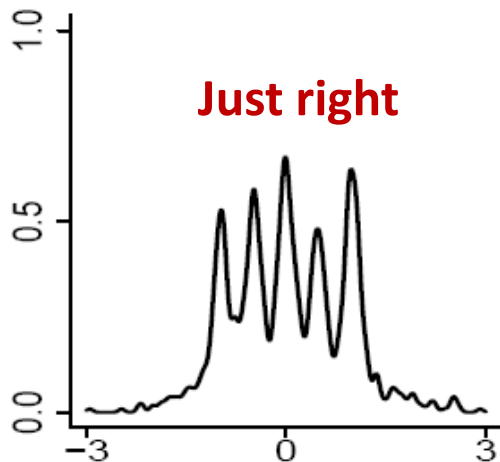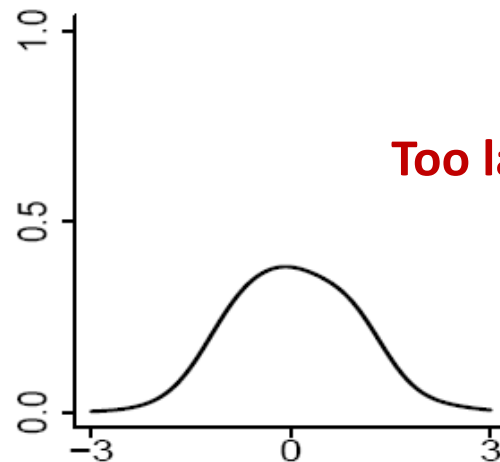
# Choice of kernel bandwidth



**Too small**

Image Source:
Larry's book – All
of Nonparametric
Statistics

**Bart-Simpson
Density**

**Just right**

**Too large**

# Histograms vs. Kernel density estimation



$\Delta$ = h acts as a smoother.

# Nonparametric density estimation

- Histogram

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

- Kernel density est

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

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

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

# Local Kernel Regression

- What is the temperature

  in the room?                                    at location x?



$$\widehat{T} = \frac{1}{n} \sum_{i=1}^{n} Y_i$$

**Average**

$$\widehat{T}(x) = \frac{\sum_{i=1}^{n} Y_i \mathbf{1}_{||X_i - x|| \le h}}{\sum_{i=1}^{n} \mathbf{1}_{||X_i - x|| \le h}}$$

**"Local" Average**

# Local Kernel Regression

- Nonparametric estimator
- Nadaraya-Watson Kernel Estimator

$$\widehat{f}_n(X) = \sum_{i=1}^{n} w_i Y_i \quad \text{Where} \quad w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^{n} K\left(\frac{X-X_i}{h}\right)}$$

- Weight each training point based on distance to test point
- Boxcar kernel yields local average

boxcar kernel :

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

# Choice of kernel bandwidth h



Image Source: Larry's book – All of Nonparametric Statistics

# Kernel Regression as Weighted Least Squares

$$\min_f \sum_{i=1}^{n} w_i (f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^{n} K\left(\frac{X - X_i}{h}\right)}$$

<span style="color:red">Weighted Least Squares</span>

Kernel regression corresponds to locally constant estimator obtained from (locally) weighted least squares

i.e. set   $f(X_i) = \beta$   (a constant)

# Kernel Regression as Weighted Least Squares

set $f(X_i) = \beta$ (a constant)

$$\min_{\beta} \sum_{i=1}^{n} w_i(\beta - Y_i)^2$$

<span style="color:red">constant</span>

$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^{n} K\left(\frac{X - X_i}{h}\right)}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^{n} w_i(\beta - Y_i) = 0$$

Notice that $\displaystyle\sum_{i=1}^{n} w_i = 1$

$$\Rightarrow \widehat{f}_n(X) = \widehat{\beta} = \sum_{i=1}^{n} w_i Y_i$$

# Local Linear/Polynomial Regression

$$\min_f \sum_{i=1}^{n} w_i(f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^{n} K\left(\frac{X-X_i}{h}\right)}$$

**Weighted Least Squares**

Local Polynomial regression corresponds to locally polynomial estimator obtained from (locally) weighted least squares

i.e. set $f(X_i) = \beta_0 + \beta_1(X_i - X) + \frac{\beta_2}{2!}(X_i - X)^2 + \cdots + \frac{\beta_p}{p!}(X_i - X)^p$

**(local polynomial of degree p around X)**

# Summary

- Non-parametric approaches

**Four things make a nonparametric/memory/instance based/lazy learner:**

1. *A distance metric, 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.