

# Non-parametric methods

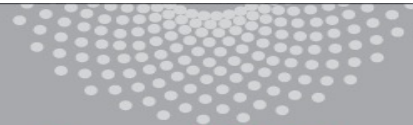
Aarti Singh

Machine Learning 10-315

Feb 23, 2022



**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

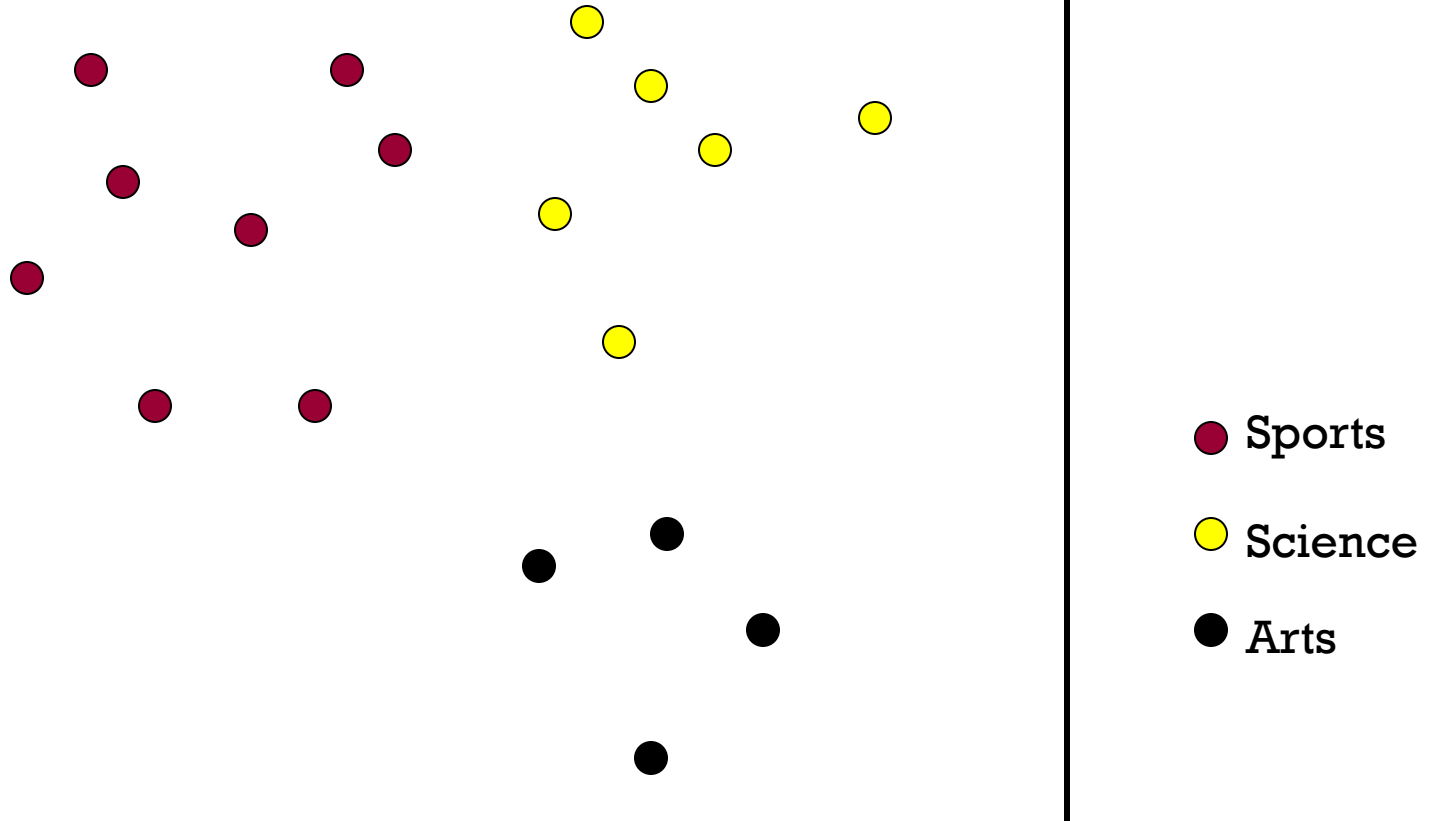
# Parametric methods

- 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 modeling assumptions, not satisfied in practice

# Non-Parametric methods

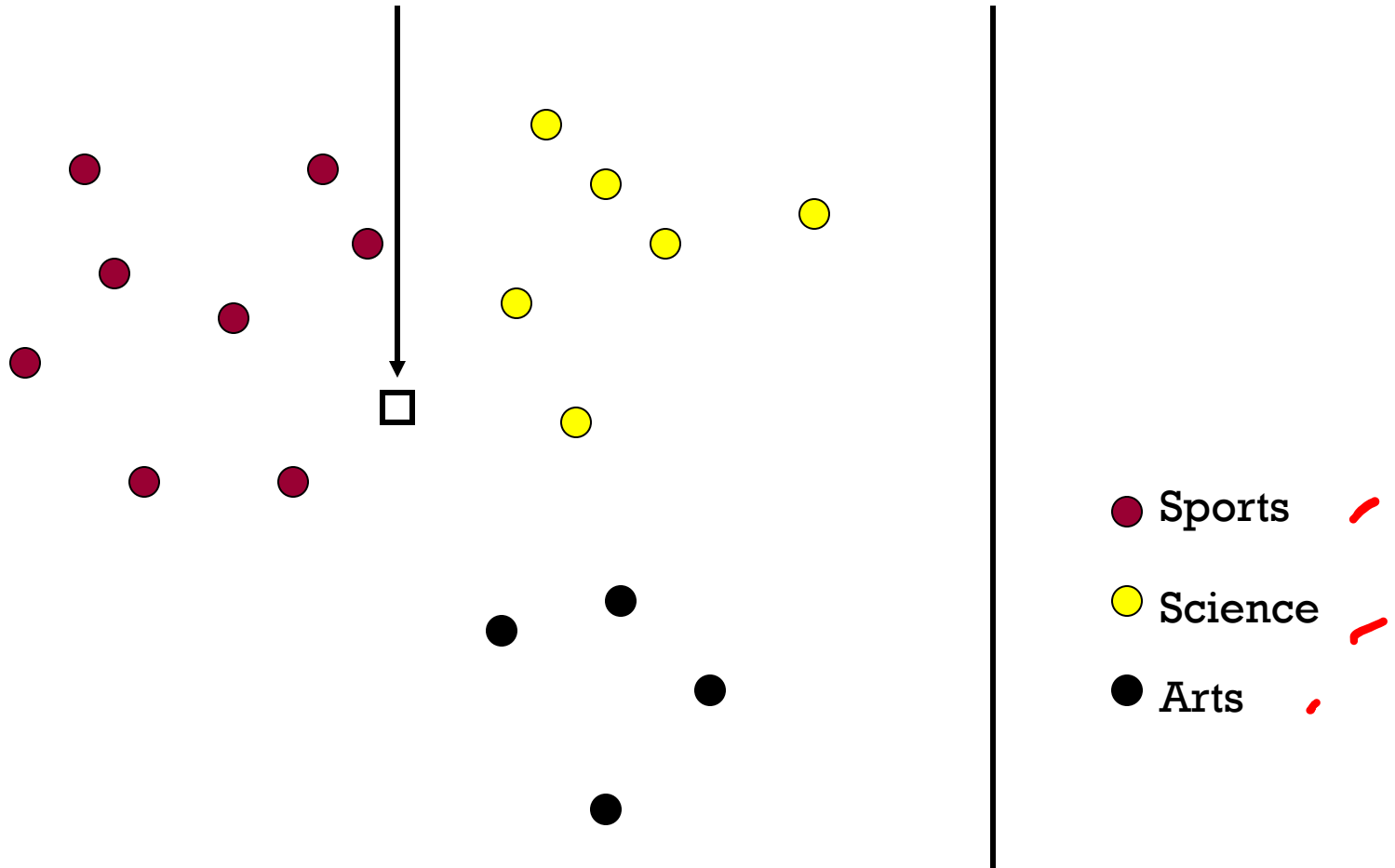
- Typically don't make any modeling assumptions
- As we have more data, we should be able to learn more complex models
- Let number of parameters scale with number of training data
- Some nonparametric methods
  - Classification:** Decision trees, k-NN (k-Nearest Neighbor) classifier
  - Density estimation:** k-NN, Histogram, Kernel density estimate
  - Regression:** Kernel regression

# k-NN classifier



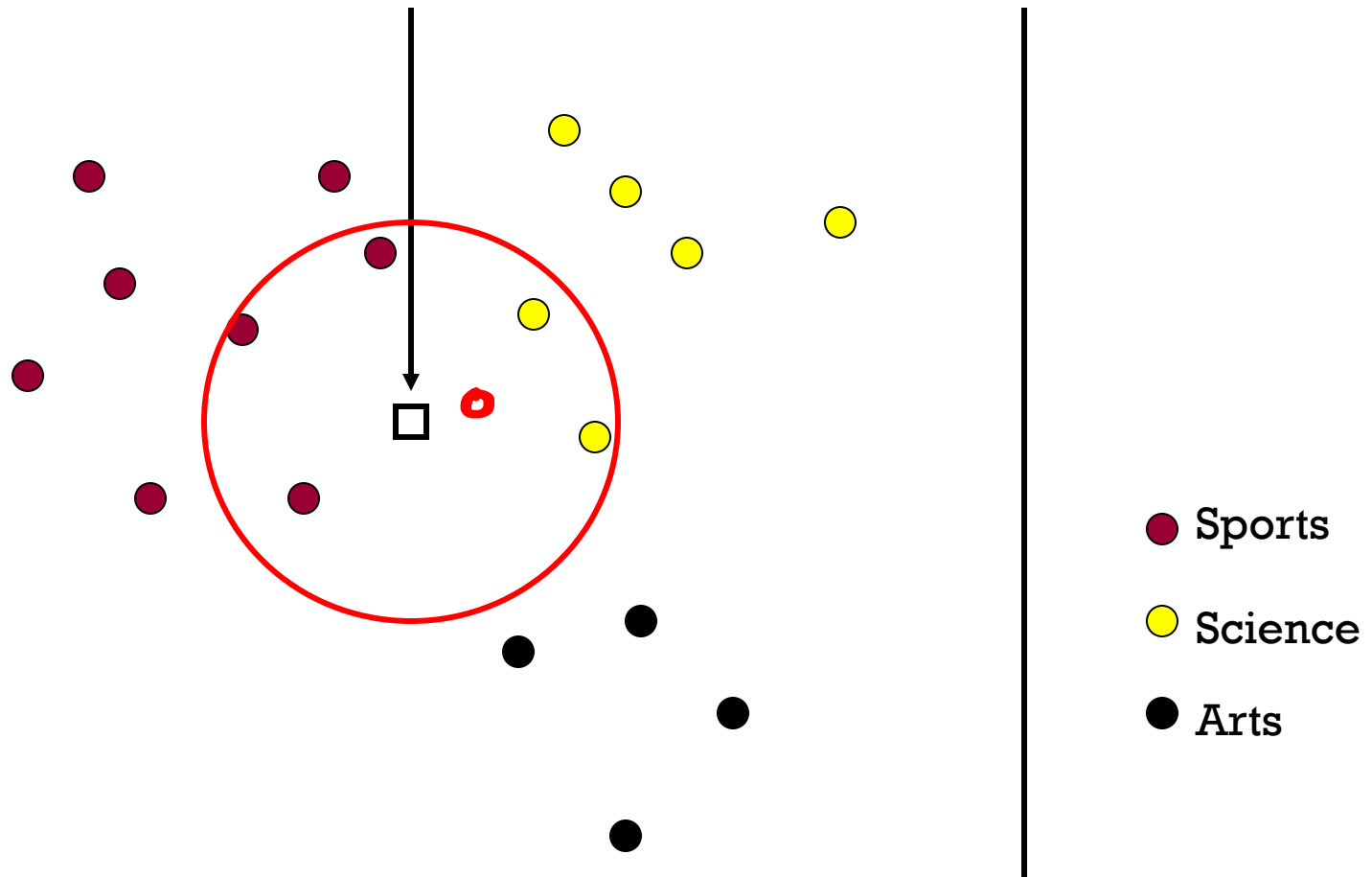
# k-NN classifier

Test document



# k-NN classifier (k=5)

Test document



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

# k-NN classifier

- Optimal Classifier:  $f^*(x) = \arg \max_y P(y|x)$   
 $= \arg \max_y \underline{P(x|y)} \underline{P(y)}$

- k-NN Classifier:  $\hat{f}_{kNN}(x) = \arg \max_y \hat{P}_{kNN}(x|y) \hat{P}(y)$   
 $= \arg \max_y \frac{k_y}{n_y} \frac{n_y}{n} = \frac{n_y}{n}$

$\rightarrow \hat{P}_{kNN}(x|y) = \frac{k_y}{n_y}$ 

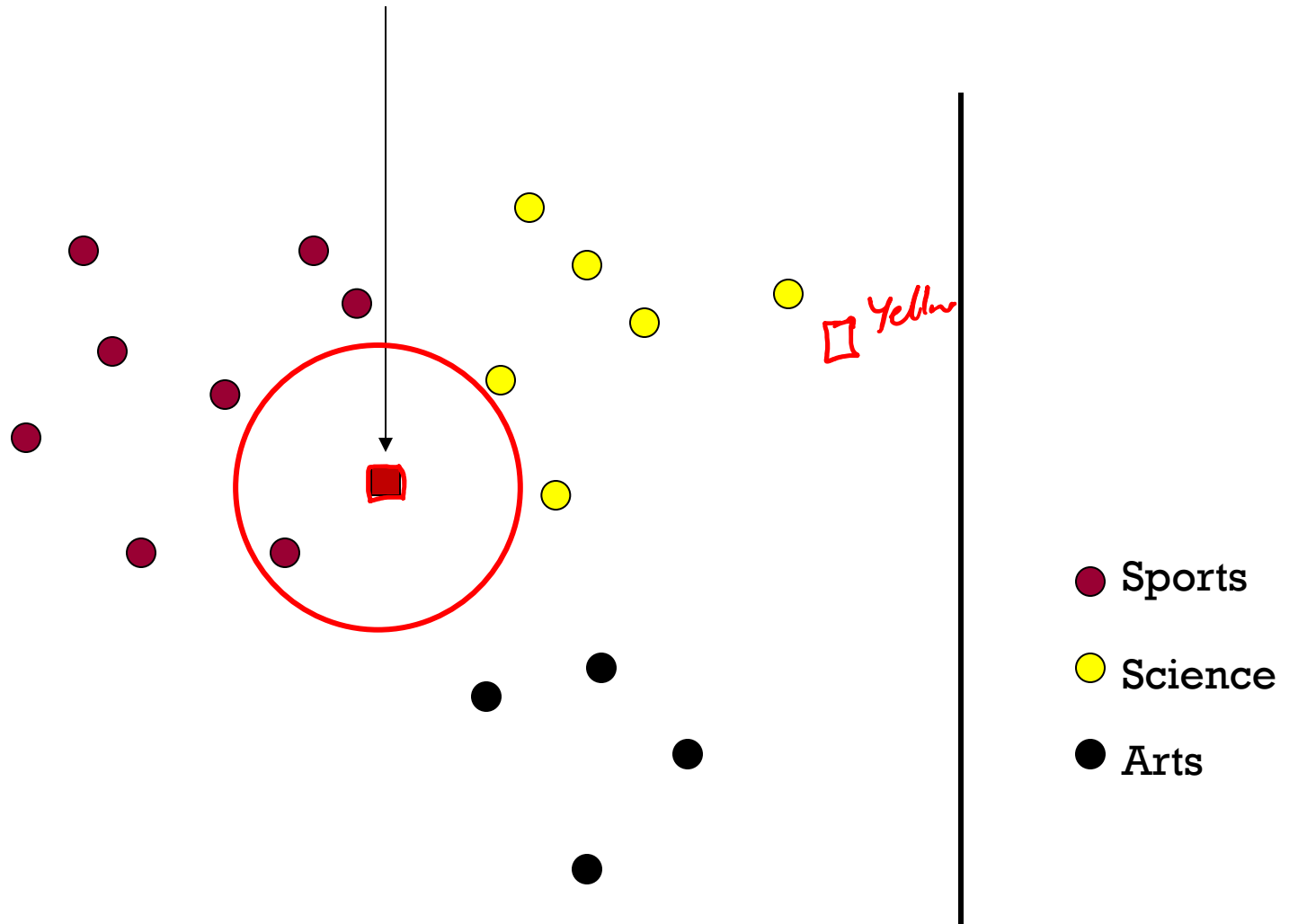
 $\rightarrow$  # training pts of class  $y$  amongst  $k$  NNs of  $x$ 

 $\sum_y k_y = k$

$\rightarrow$  # total training pts of class  $y$

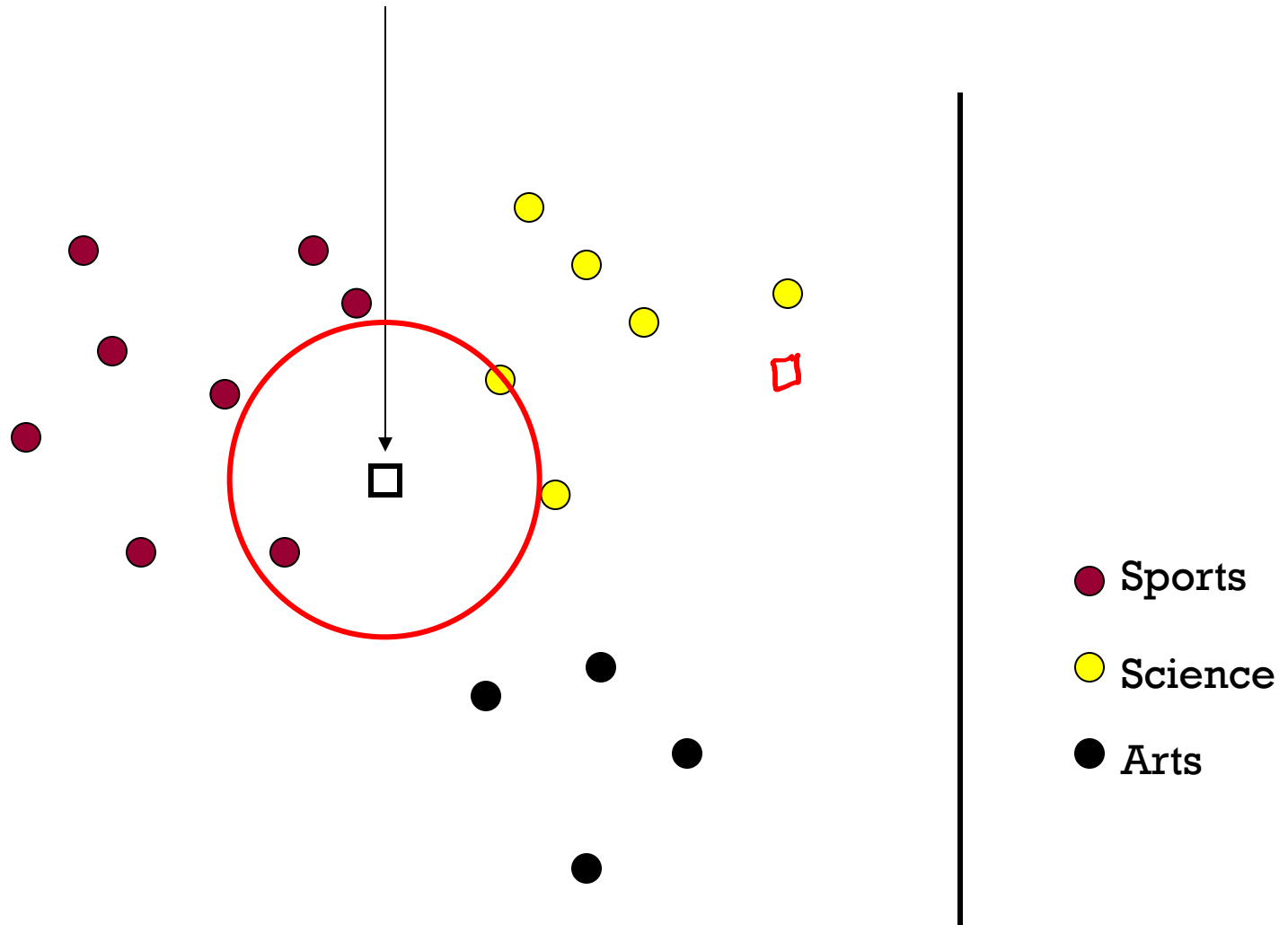
$$\hat{P}(y) = \frac{n_y}{n}$$

# 1-Nearest Neighbor (kNN) classifier

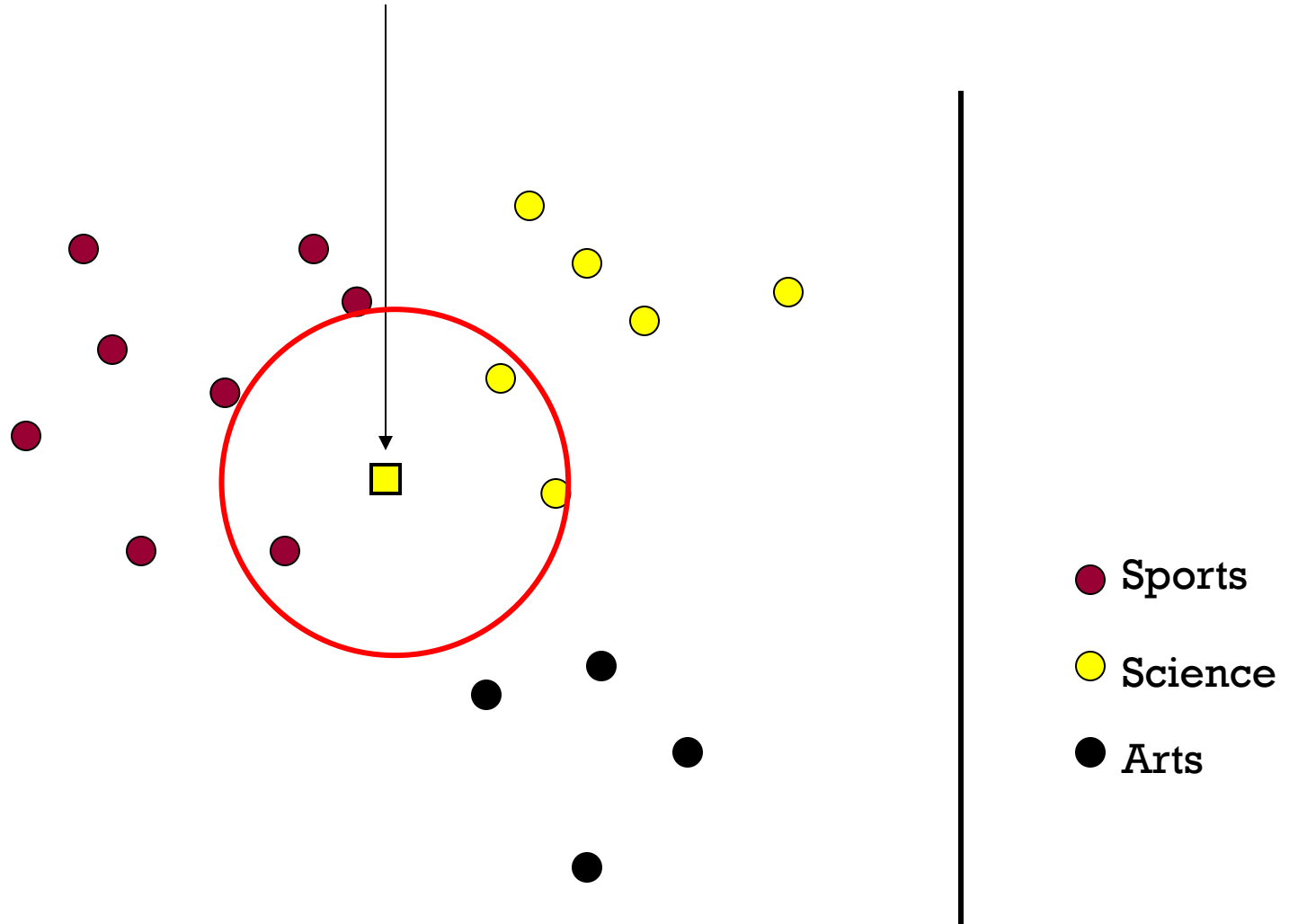




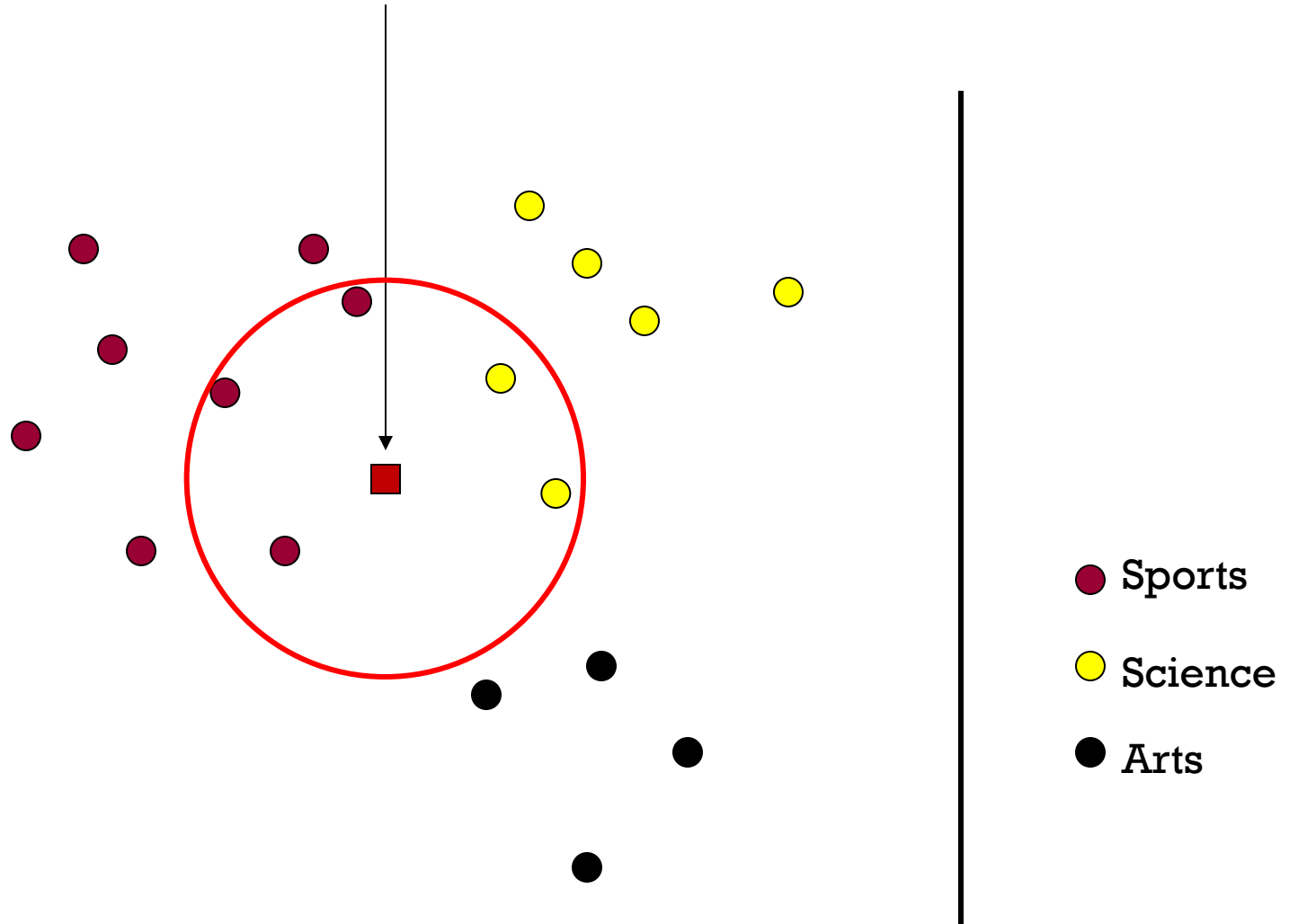
# 2-Nearest Neighbor (kNN) classifier



# 3-Nearest Neighbor (kNN) classifier

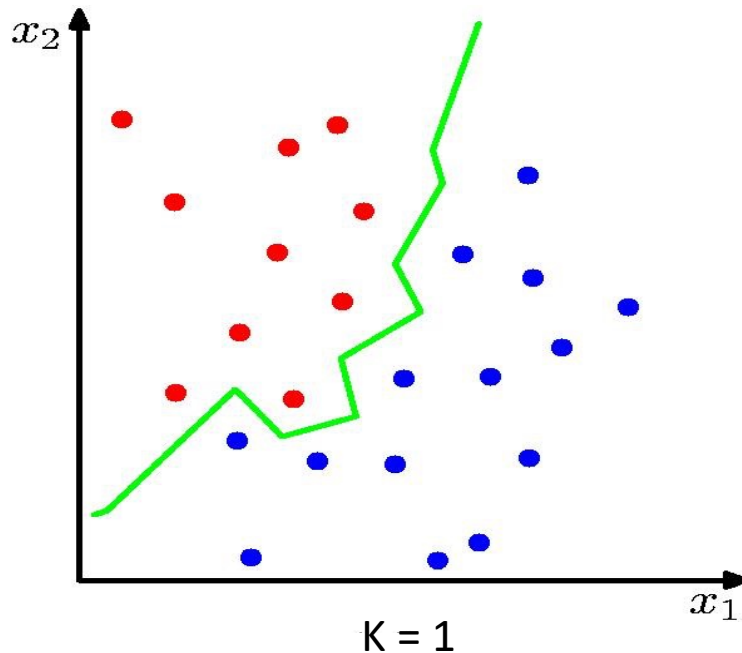


# 5-Nearest Neighbor (kNN) classifier

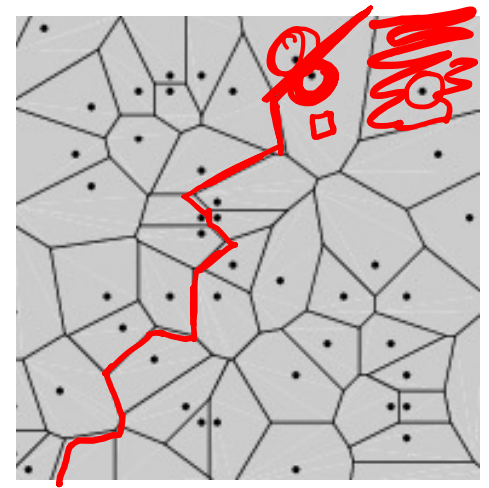


# What is the best k?

1-NN classifier decision boundary



Voronoi  
Diagram

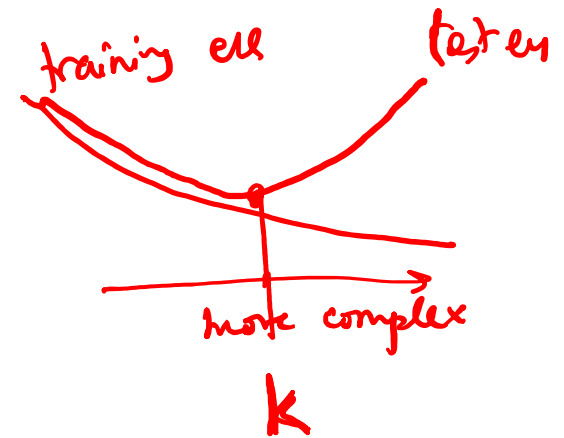


As k increases, boundary becomes smoother (less jagged).

# What is the best k?

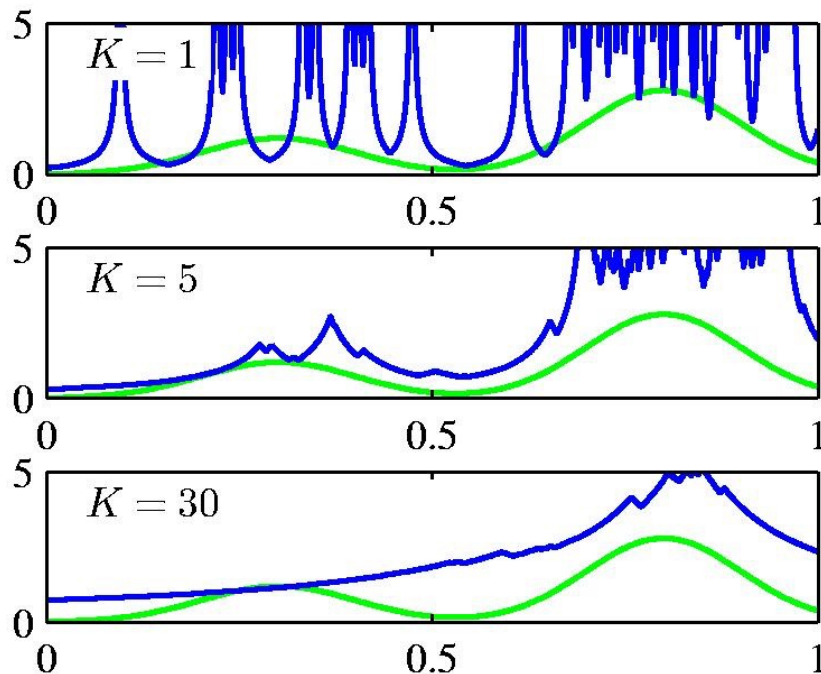
Approximation vs. Stability (aka Bias vs Variance) Tradeoff

- Larger  $K \Rightarrow$  predicted label is more stable
- Smaller  $K \Rightarrow$  predicted label can approximate best classifier well given enough data

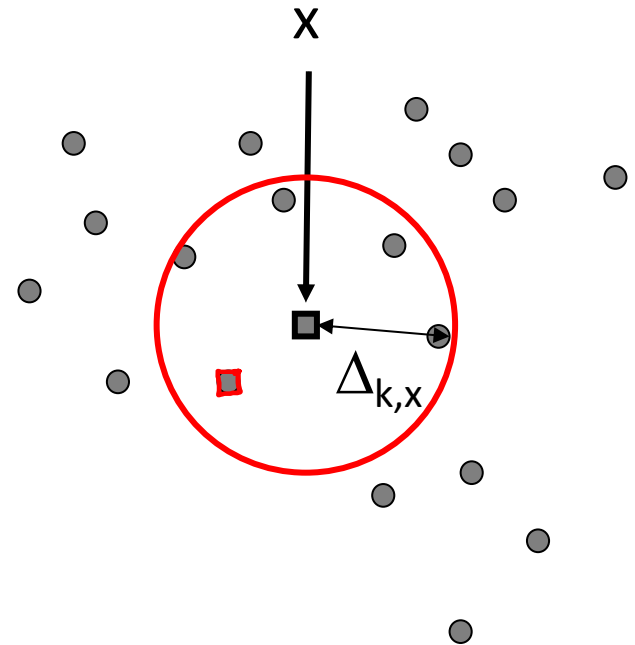


# k-NN density estimation

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



$k$  acts as a smoother.



Not very popular for density estimation – spiked estimates

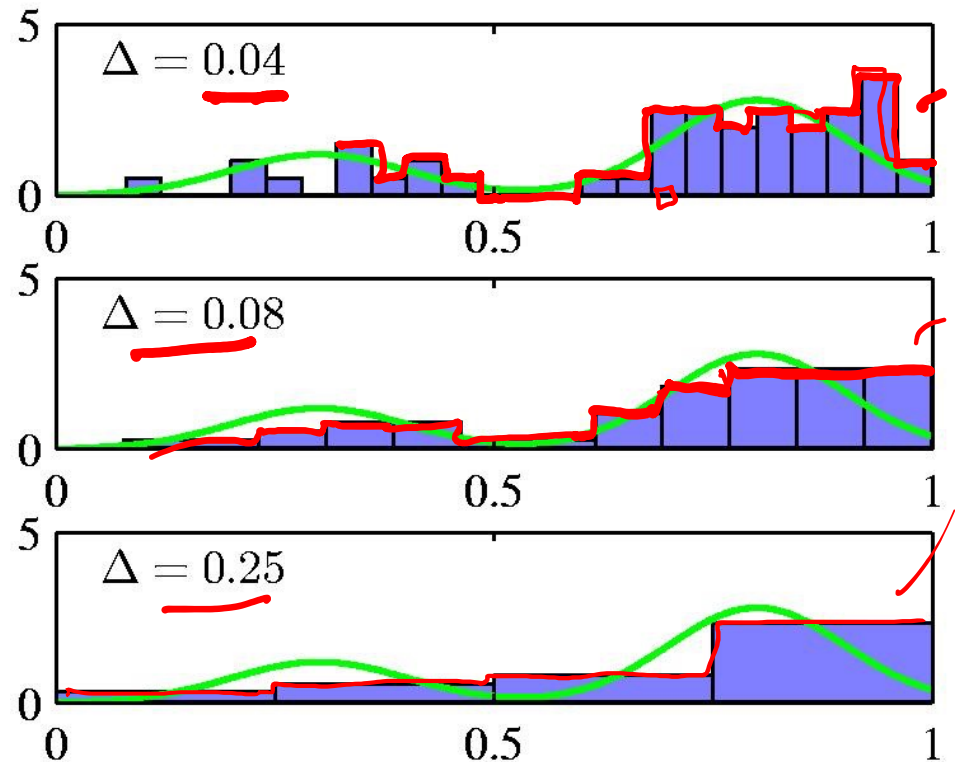
# Histogram density estimate

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

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

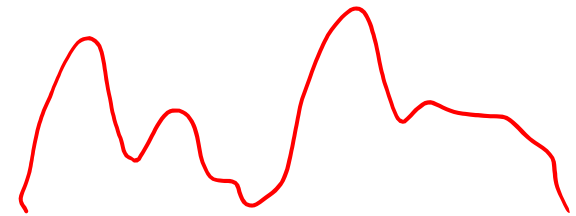
“Local relative frequency”

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

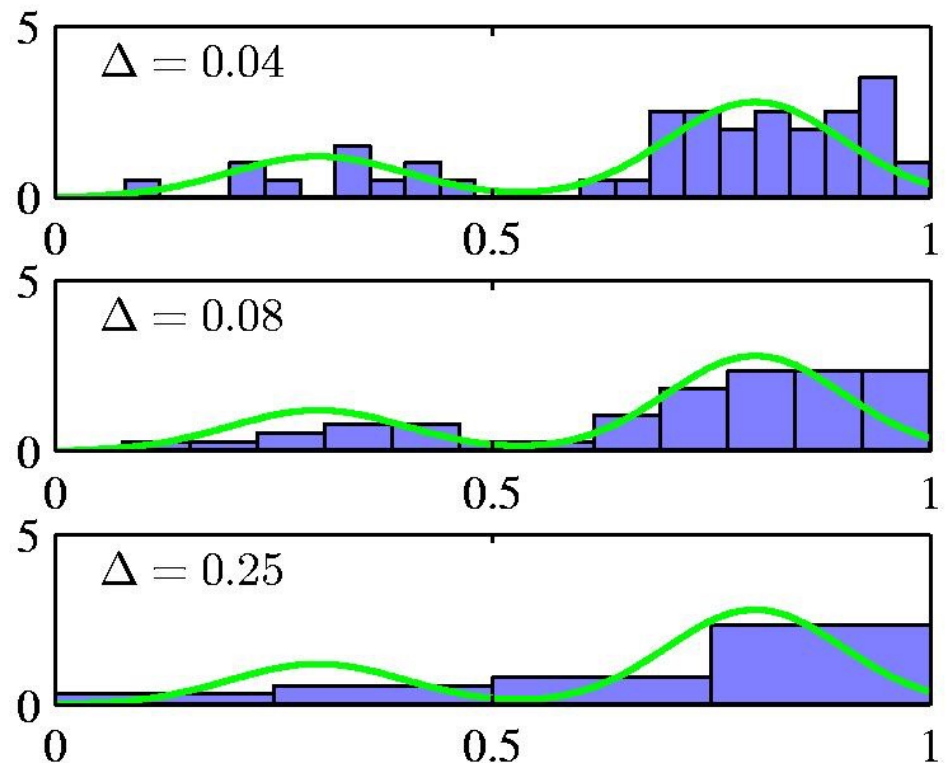


# Effect of histogram bin width

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



# bins =  $1/\Delta$



Small  $\Delta$ , large #bins  
Good fit but unstable  
(few points per bin)

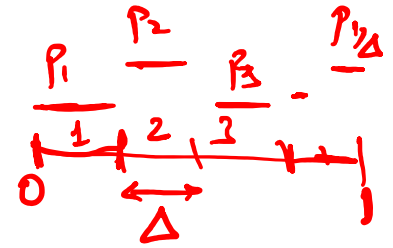
“**Small bias**, **Large variance**”

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  $\int p(x)dx = 1$

$$\sum_j p_j \Delta = 1$$

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

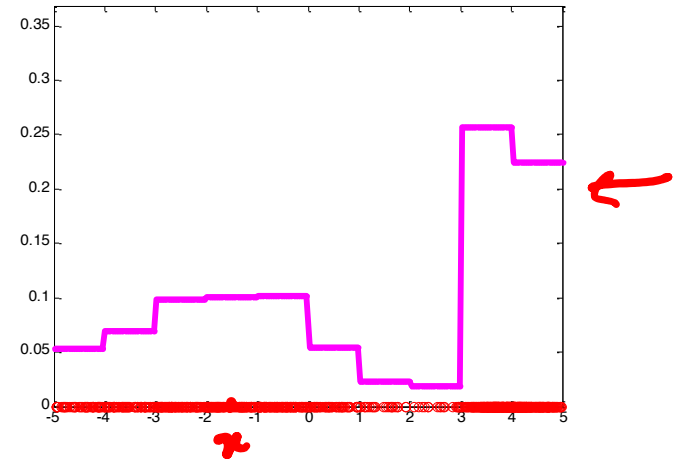
$$\hat{p}(x) = \arg \max_{\{p_j\}} P(\underbrace{X_1, \dots, X_n}_{\text{data}}; \underbrace{\{p_j\}_{j=1}^{1/\Delta}}_{\text{parameters}}) \quad \text{s.t.} \quad \sum_j p_j = 1/\Delta$$

- Show that histogram density estimate is MLE under this model = Categorical( $p_1\Delta$ ,  $p_2\Delta$ ,  $p_3\Delta$ ,...)

# Kernel density estimate

- Histogram – blocky estimate

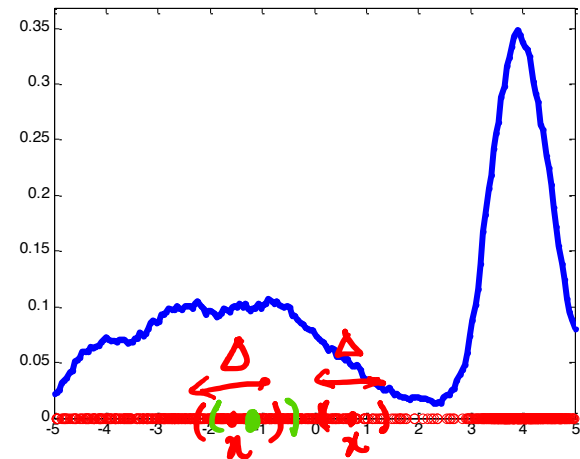
$$\hat{p}(x) = \frac{1}{\Delta} \frac{\sum_{j=1}^n \mathbf{1}_{X_j \in \text{Bin}_x}}{n} = \frac{n_x}{h\Delta}$$



- 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 ( $X_j, x, \Delta$ )



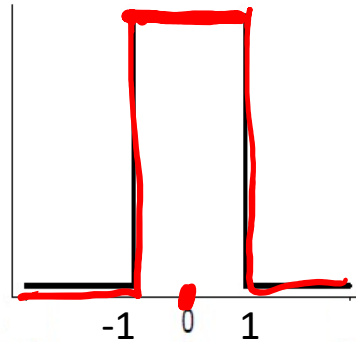
# Kernel density estimate

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

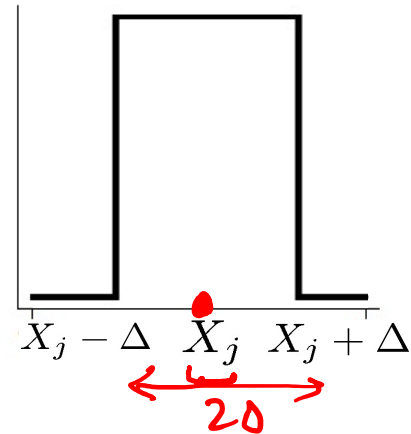
$\mathbb{1}_{\|X_j - x\| \leq \Delta}$

boxcar kernel :

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

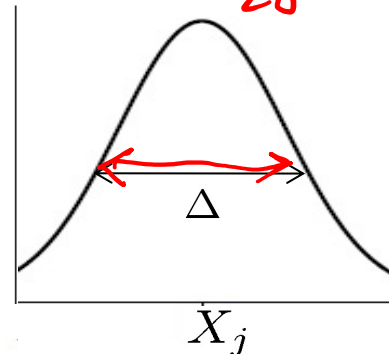
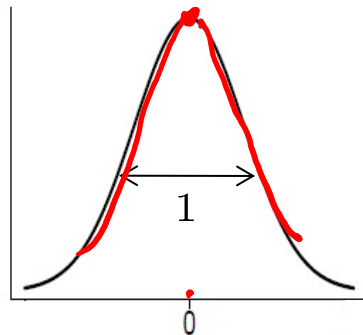


$$K\left(\frac{X_j - x}{\Delta}\right)$$



Gaussian kernel :

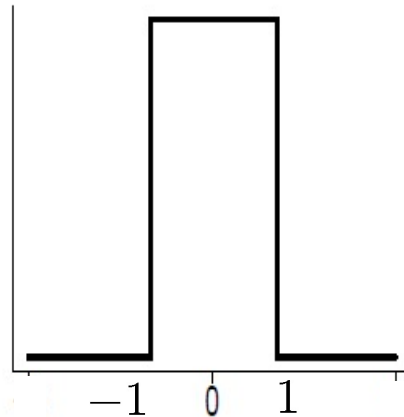
$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



# Kernels

boxcar kernel :

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



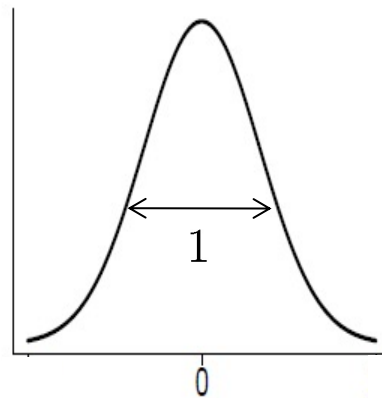
Any kernel function that satisfies

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

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

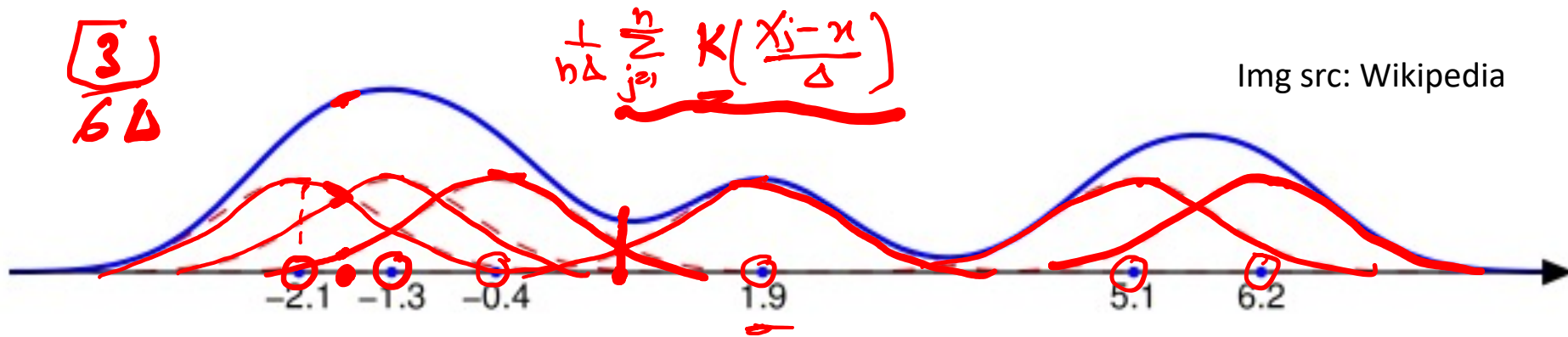
Gaussian kernel :

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



# Kernel density estimation

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



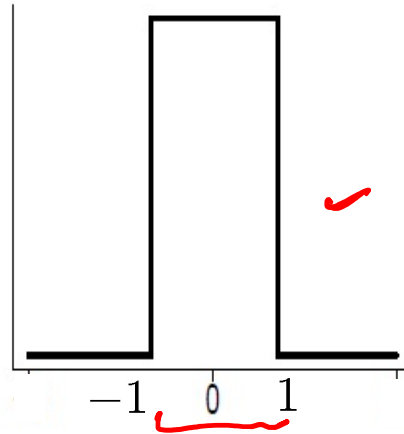
Gaussian bumps (red) around six data points and their sum (blue)

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

# Choice of Kernels

boxcar kernel :

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

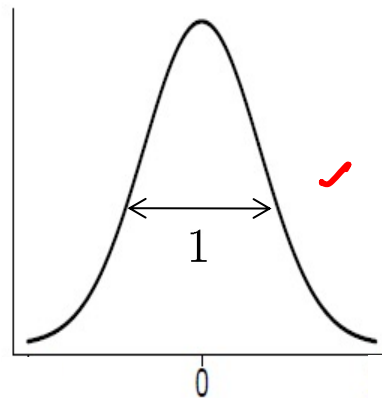


**Finite support**

- only need local points to compute estimate

Gaussian kernel :

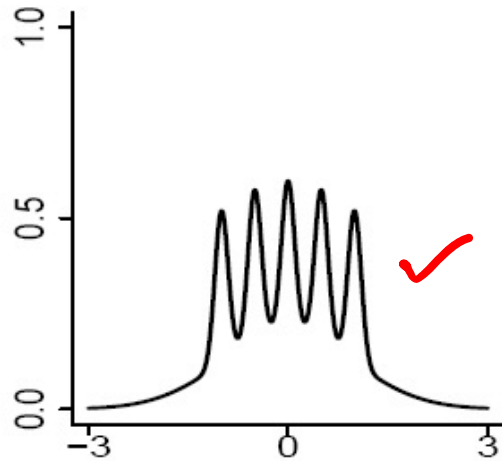
$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



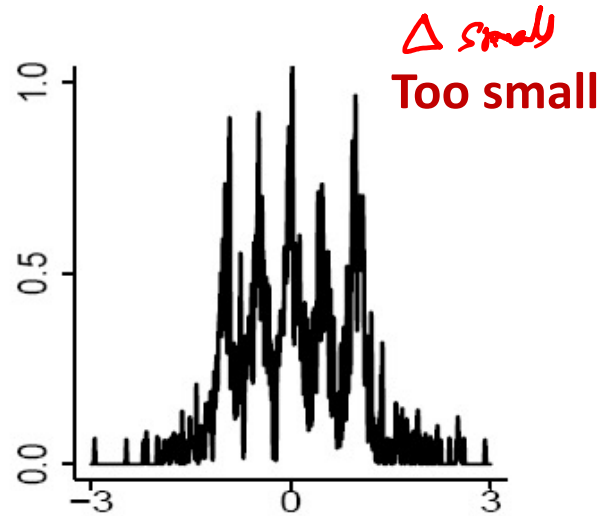
**Infinite support**

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

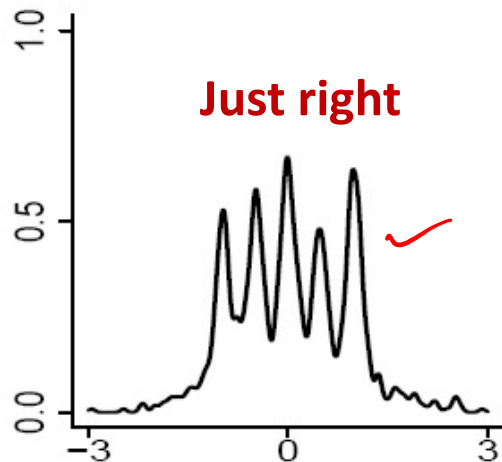
# Choice of kernel bandwidth



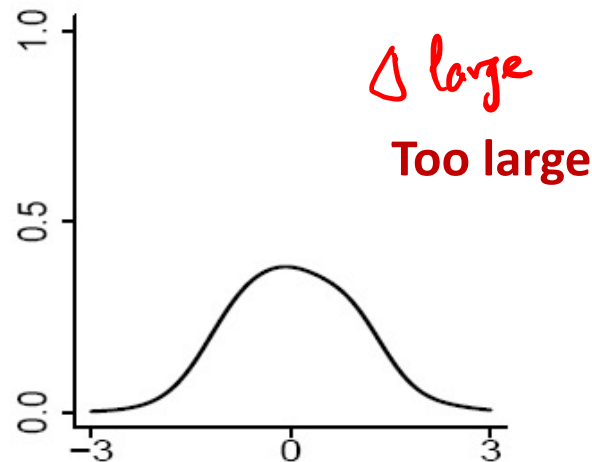
True Density



Undersmoothed



Just Right

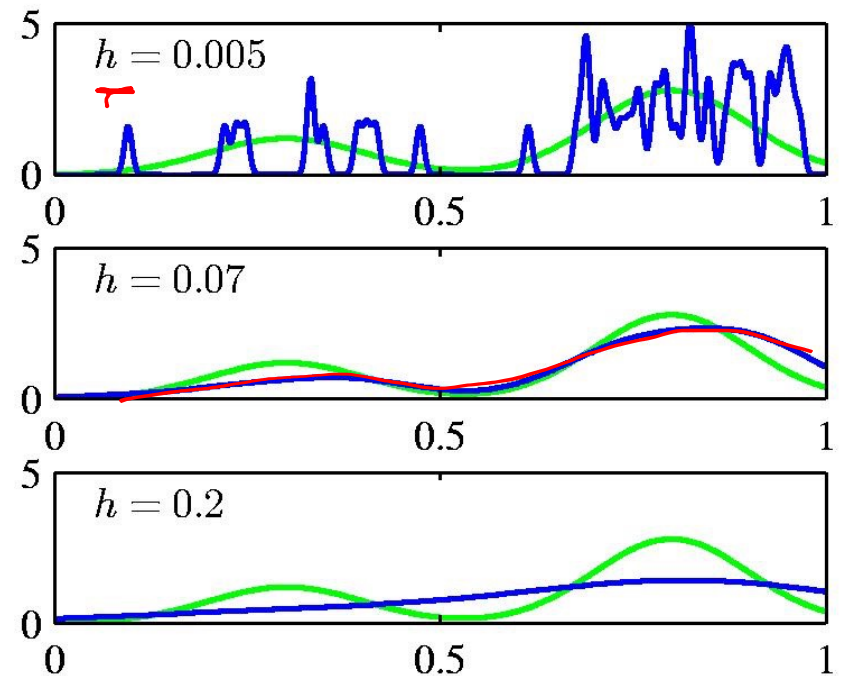
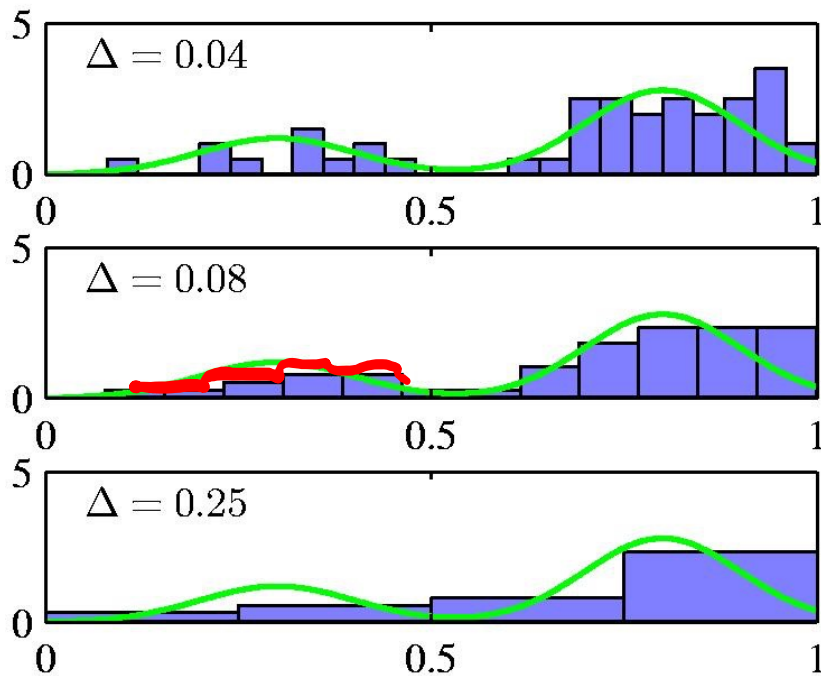


Oversmoothed

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

**Bart-Simpson  
Density**

# Histograms vs. Kernel density estimation



$\Delta = h$  acts as a smoother.





# Nonparametric density estimation

- Histogram

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

*non-overlapping  
bins*

- Kernel density est

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

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

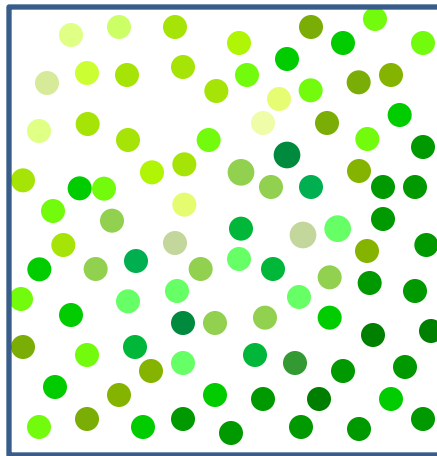
Fix  $n_x = k$ , estimate  $\Delta$  from data (volume of ball around  $x$  that contains  $k$  training pts)

- k-NN density est

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

# Local Kernel Regression

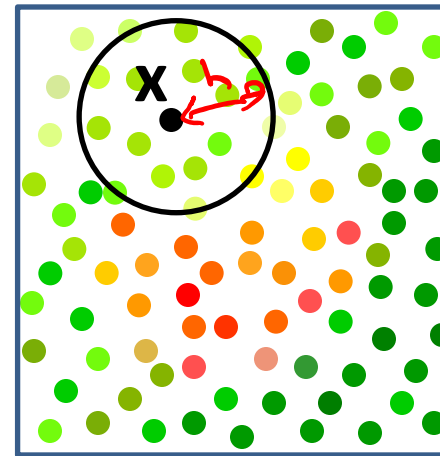
- What is the temperature in the room?



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

Average

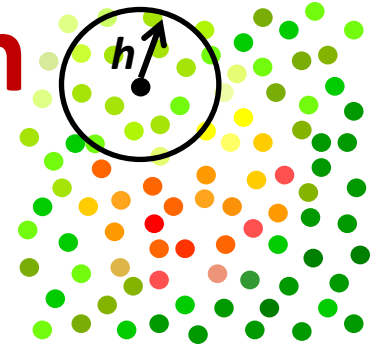
at location  $x$ ?



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

"Local" Average ✓

# Local Kernel Regression



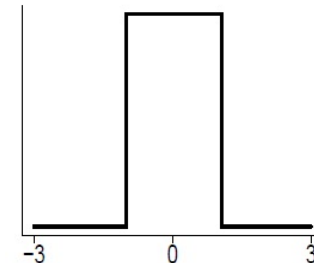
- Nonparametric estimator
- Nadaraya-Watson Kernel Estimator

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

- 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 $h$

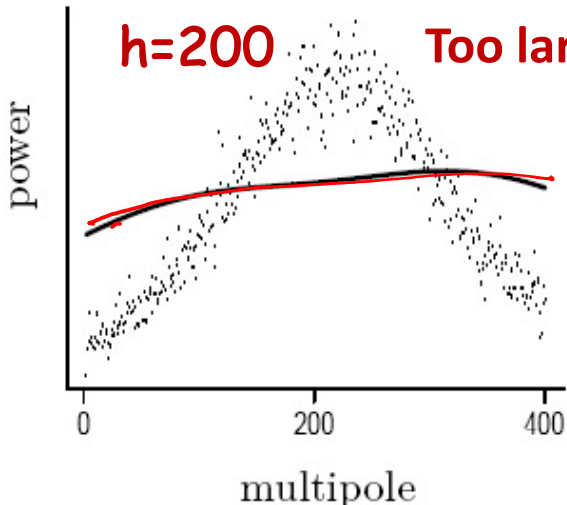
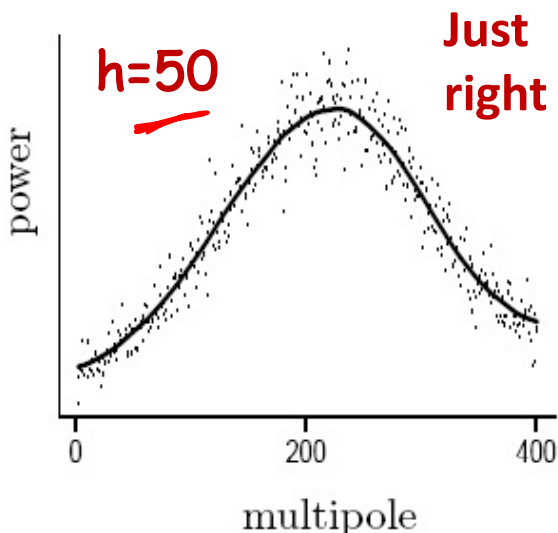
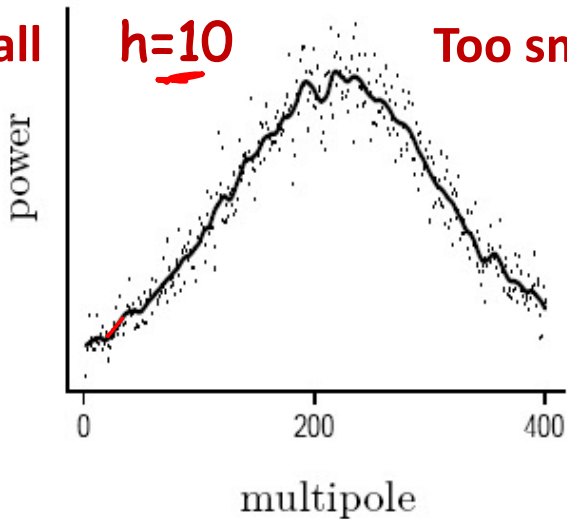
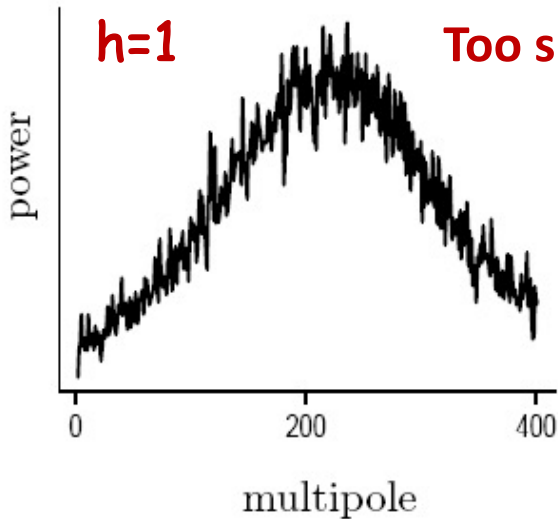
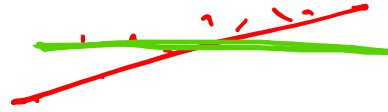


Image Source:  
Larry's book – All  
of Nonparametric  
Statistics

# Kernel Regression as Weighted Least

LS:  $\sum_{i=1}^n (f(x_i) - y_i)^2$   $f(x_i) = X_i \beta$  **Squares**



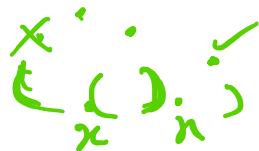
WLS:  $\min_f \sum_{i=1}^n w_i (f(X_i) - Y_i)^2$

Weighted Least Squares

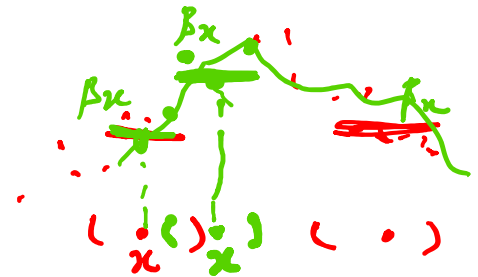
training pt.  $\uparrow$   
 $w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$   $\uparrow$   
 test point

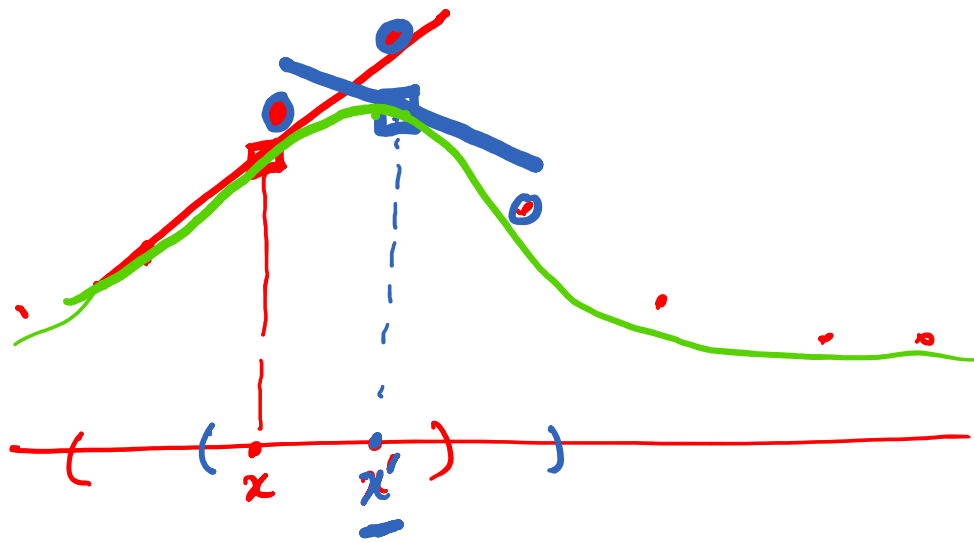
Kernel regression corresponds to locally constant estimator obtained from (locally) weighted least squares

i.e. set  $f(x_i) = \beta$  (a constant)  $\beta \equiv \beta(x)$



$f(x_i) = X_i \beta$





# Kernel Regression as Weighted Least Squares

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

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

$f(x_i)$   
 $\downarrow$   
 constant  
 $J(\beta)$

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

$$\beta \sum_{i=1}^n w_i = \sum_{i=1}^n w_i Y_i$$

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

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

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

$(X)$                        $(X)$

# Local Linear/Polynomial Regression

$$\min_f \sum_{i=1}^n \underbrace{w_i}_{\uparrow X} \underbrace{(f(X_i) - Y_i)^2}$$

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

Weighted Least Squares

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

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

(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,  $\text{dist}(x, X_i)$  ✓  $\|x - X_i\|$   
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) modeling assumptions

- Nonparametric models typically require storage and computation of the order of entire data set size.

Parametric models, once fitted, are much more efficient in terms of storage and computation.