

# How to learn a decision tree <sup>min</sup> $H(Y|X_i)$

- Top-down induction [ID3, C4.5, C5, ...]

$\max_i I(Y, X_i)$   
↑ feature

Main loop: C4.5

1.  $X \leftarrow$  the “best” decision feature for next *node*

2. Assign  $X$  as decision feature for *node*

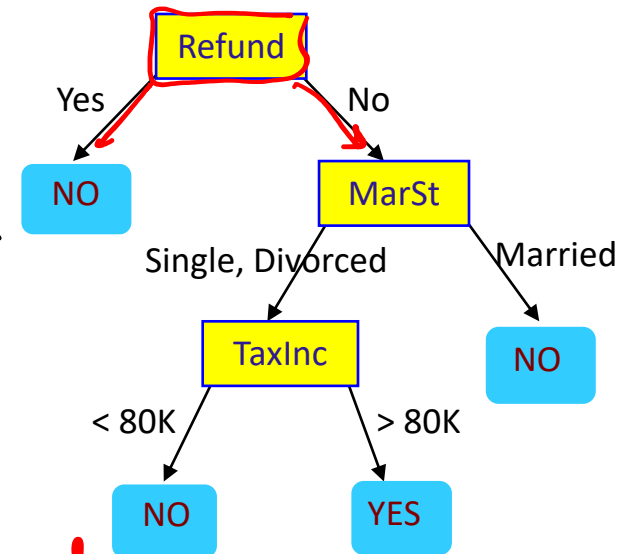
⇒ 3. For “best” split of  $X$ , create new descendants of *node*

4. Sort training examples to leaf nodes

⇒ 5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

⇒ 6. Prune back tree to reduce overfitting

7. Assign majority label to the leaf node



(don't discard previously used features)

# Handling continuous features (C4.5)

Convert continuous features into discrete by setting a threshold.

What threshold to pick?

Search for best one as per information gain. Infinitely many??

Don't need to search over more than  $\sim n$  (number of training data), e.g. say  $X_1$  takes values  $x_1^{(1)}, x_1^{(2)}, \dots, x_1^{(n)}$  in the training set. Then possible thresholds are

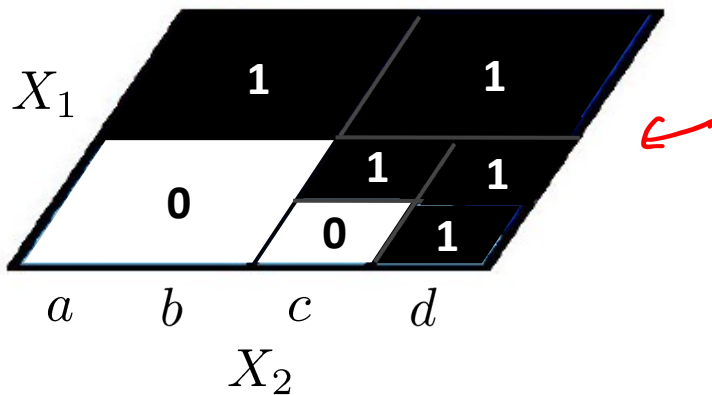
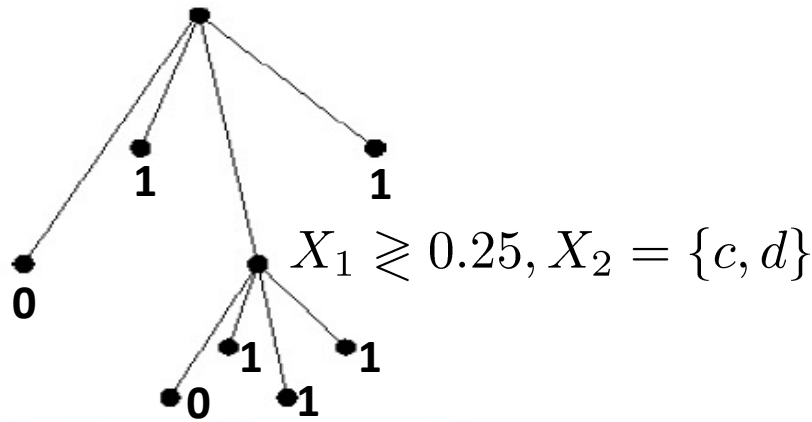
$$\underline{[x_1^{(1)} + x_1^{(2)}]/2}, \underline{[x_1^{(2)} + x_1^{(3)}]/2}, \dots, \underline{[x_1^{(n-1)} + x_1^{(n)}]/2}$$

# Poll

- What's the maximum depth of a decision tree learnt using ID3?
- What's the maximum depth of a decision tree learnt using C4.5?

# Decision Tree more generally...

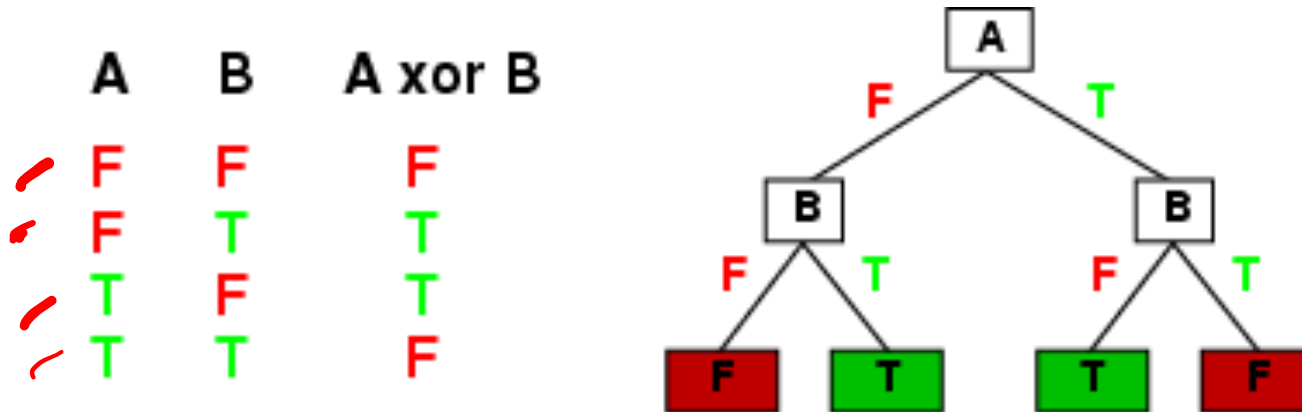
$$X_1 \geq 0.5, X_2 = \{a, b\} \text{ or } \{c, d\}$$



- Features can be discrete, continuous or categorical
- Each internal node: test some set of features  $\{X_i\}$
- Each branch from a node: selects a set of value for  $\{X_i\}$
- Each leaf node: prediction for  $Y$

# Expressiveness of Decision Trees

- Decision trees in general (without pruning) can express any function of the input features.
- E.g., for Boolean functions, truth table row  $\rightarrow$  path to leaf:



- There is a decision tree which perfectly classifies a training set with one path to leaf for each example - **overfitting**
- But it won't generalize well to new examples - prefer to find more **compact** decision trees

# Pruning the tree

- Many strategies for picking simpler trees:

- Pre-pruning

- Fixed depth (e.g. ID3)
- Fixed number of leaves

- Post-pruning

- Chi-square test

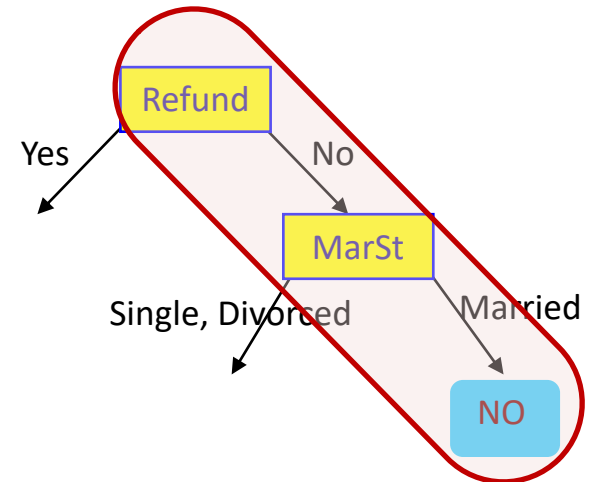
*– independence test*

- Convert decision tree to a set of rules

- Eliminate variable values in rules which are independent of label (using chi-square test for independence)

- Simplify rule set by eliminating unnecessary rules

- Information Criteria: MDL(Minimum Description Length)



# Information Criteria

- Penalize complex models by introducing cost

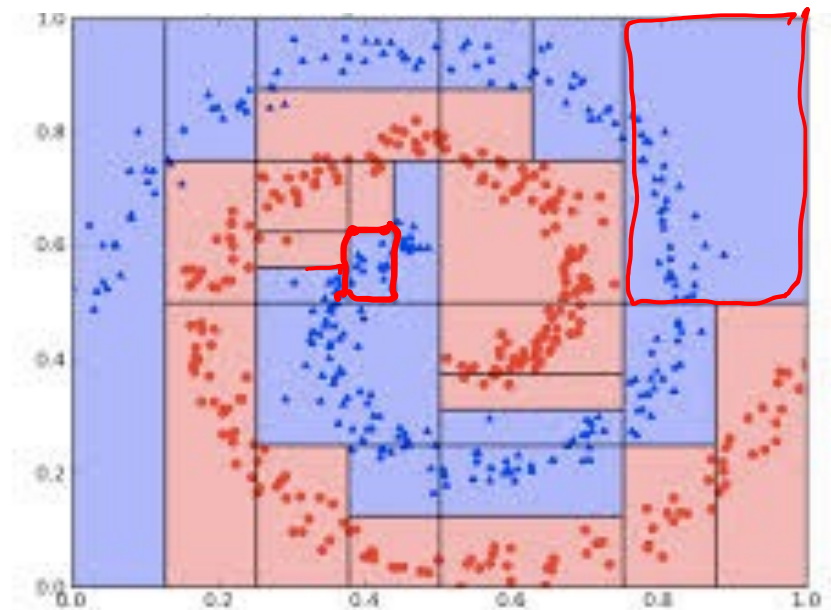
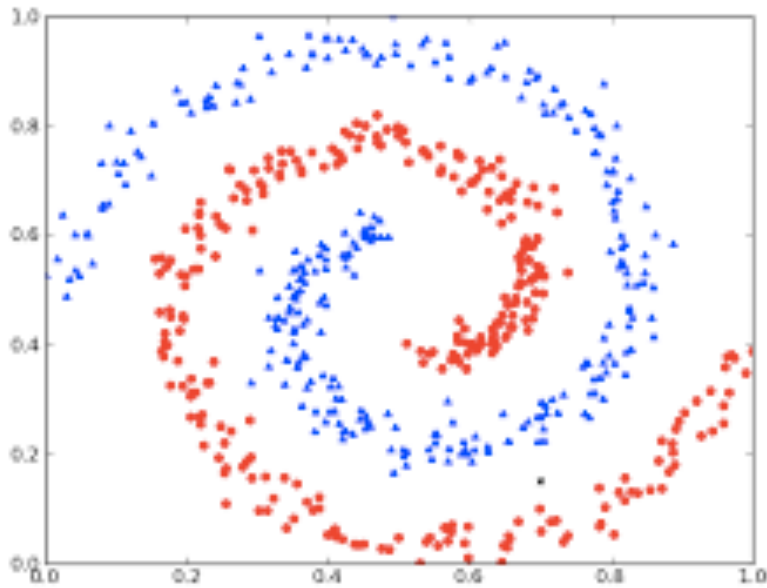
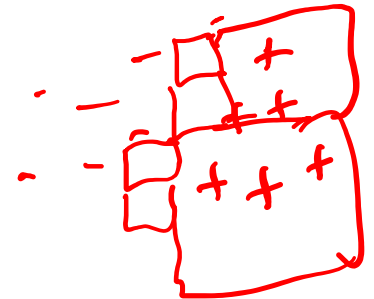
$$\hat{f} = \arg \min_T \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \text{loss}(\hat{f}_T(X_i), Y_i)}_{\text{log likelihood}} + \underbrace{\text{pen}(T)}_{\text{cost}} \right\}$$

$$\begin{aligned} \text{loss}(\hat{f}_T(X_i), Y_i) &= (\hat{f}_T(X_i) - Y_i)^2 && \text{regression} \quad \checkmark \\ &= \mathbf{1}_{\hat{f}_T(X_i) \neq Y_i} && \text{classification} \quad \checkmark \end{aligned}$$

$\text{pen}(T) \propto |T|$       penalize trees with more leaves

CART – optimization can be solved by dynamic programming

# Example of 2-feature decision tree classifier



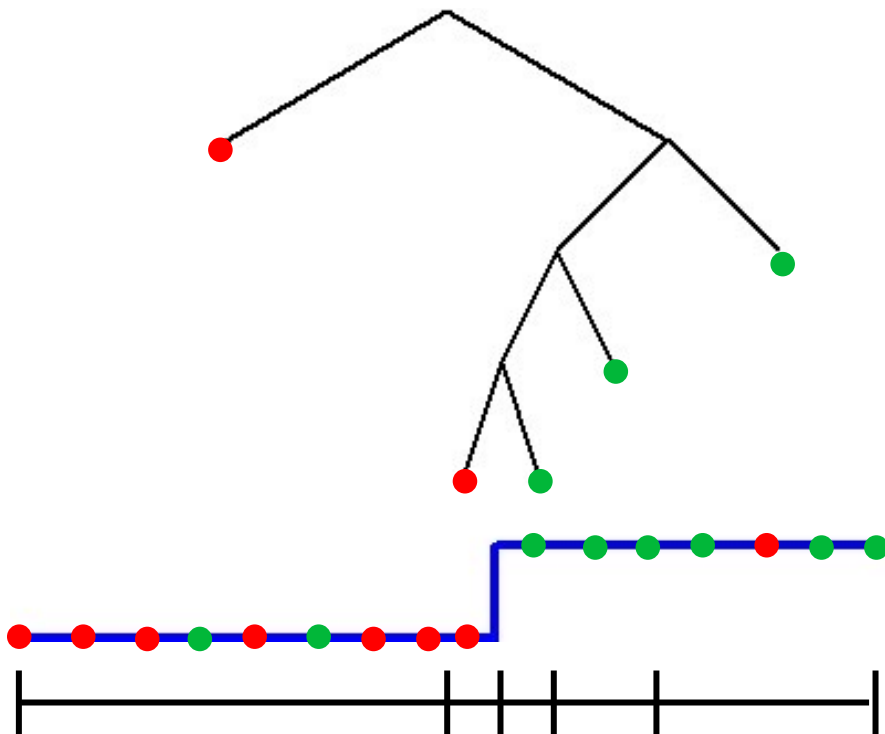
larger bins  $\Rightarrow$  low variance, high bias  
smaller bins  $\Rightarrow$  high variance, low bias



# How to assign label to each leaf

Classification – Majority vote

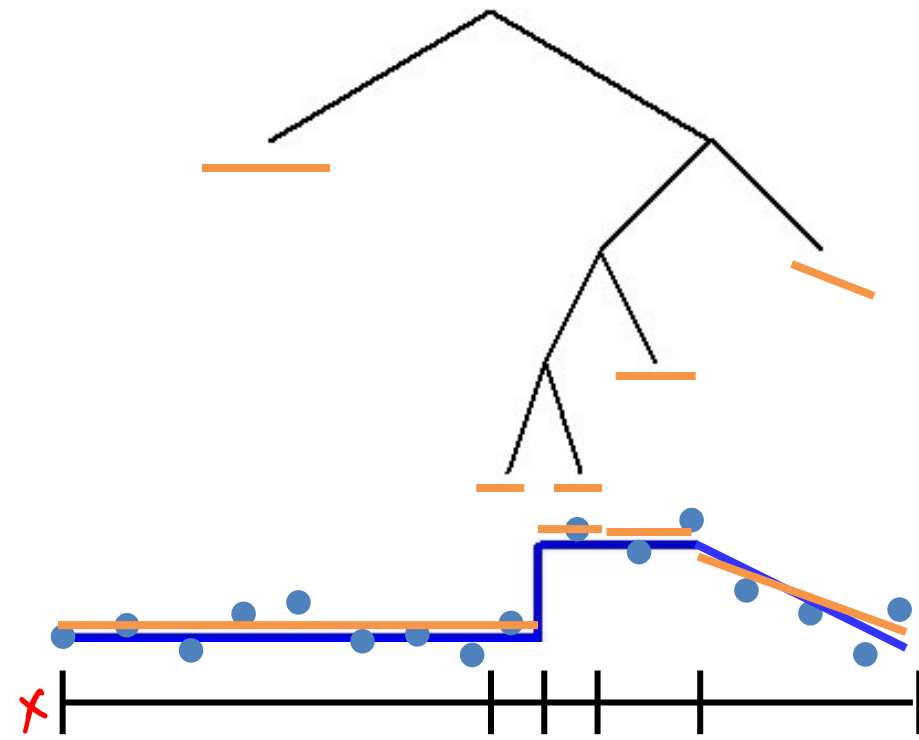
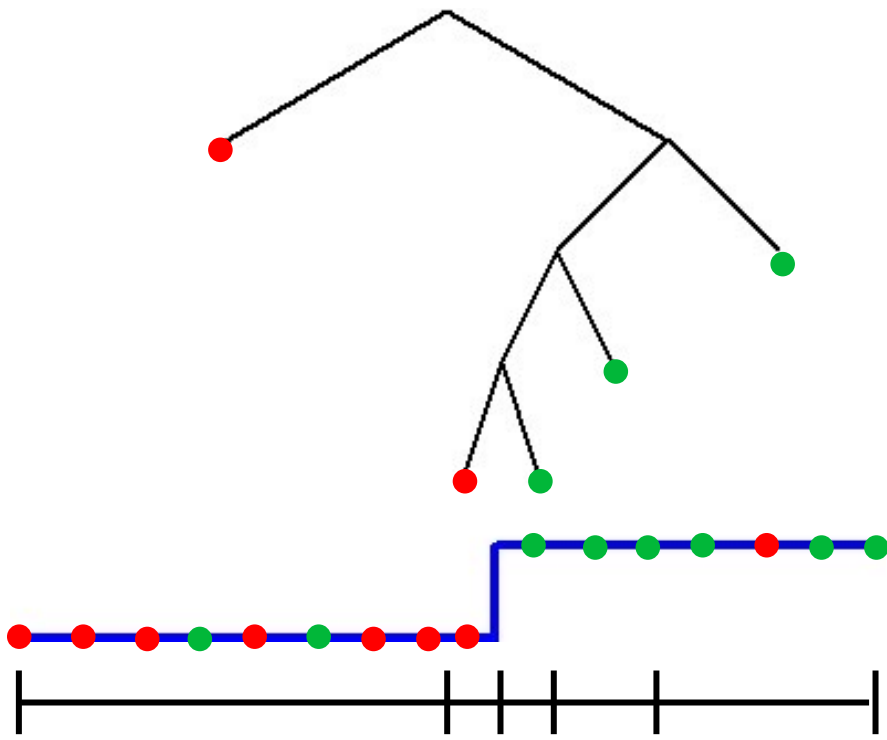
Regression – ?



# How to assign label to each leaf

Classification – Majority vote

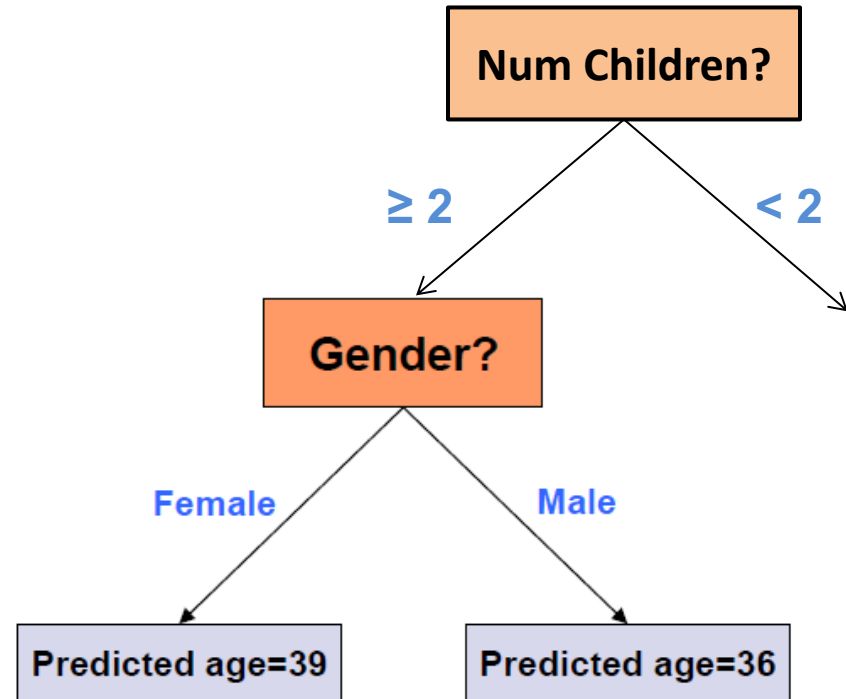
Regression – Constant/  
Linear/Poly fit



# Regression trees

$X^{(1)}$       ....       $X^{(p)}$      $Y$

Gender	Rich?	Num. Children	# travel per yr.	Age
F	No	2	5	38
M	No	0	2	25
M	Yes	1	0	72
:	:	:	:	:



Average (fit a constant ) using training data at the leaves

# What you should know

- Decision trees are one of the most popular data mining tools
  - Simplicity of design
  - Interpretability ↵
  - Ease of implementation
  - Good performance in practice (for small dimensions)
- Information gain to select attributes (ID3, C4.5,...) ✓
- Decision trees will overfit!!!
  - Must use tricks to find “simple trees”, e.g.,
    - Pre-Pruning: Fixed depth/Fixed number of leaves
    - Post-Pruning: Chi-square test of independence
    - Complexity Penalized/MDL model selection
- Can be used for classification, regression and density estimation too

*features  
≡ attributes*

# k-NN classifier ✓

## Nonparametric kernel regression

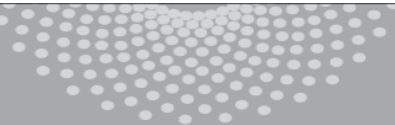
Aarti Singh

Feb 27, 2023

Machine Learning 10-701

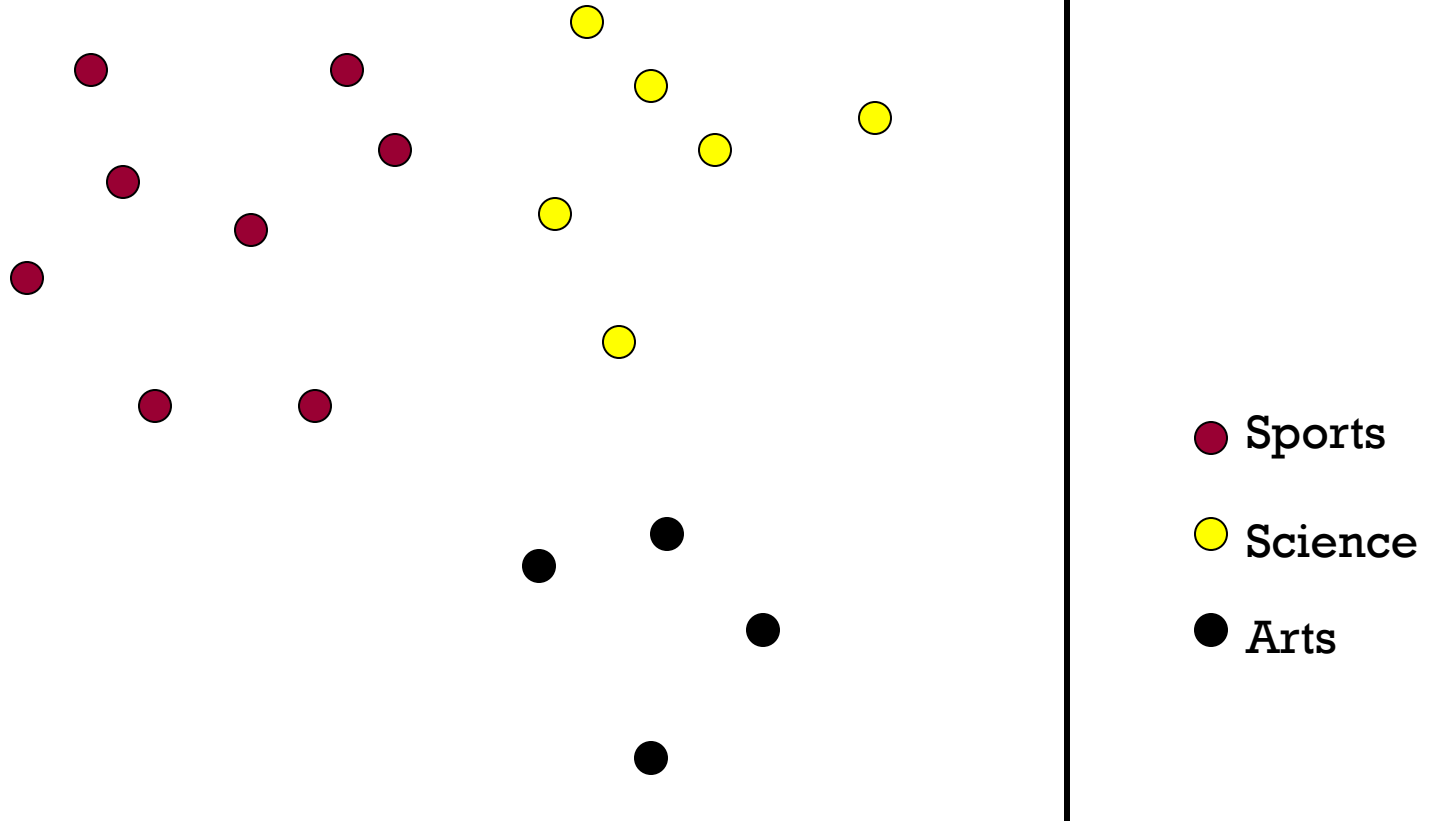


MACHINE LEARNING DEPARTMENT



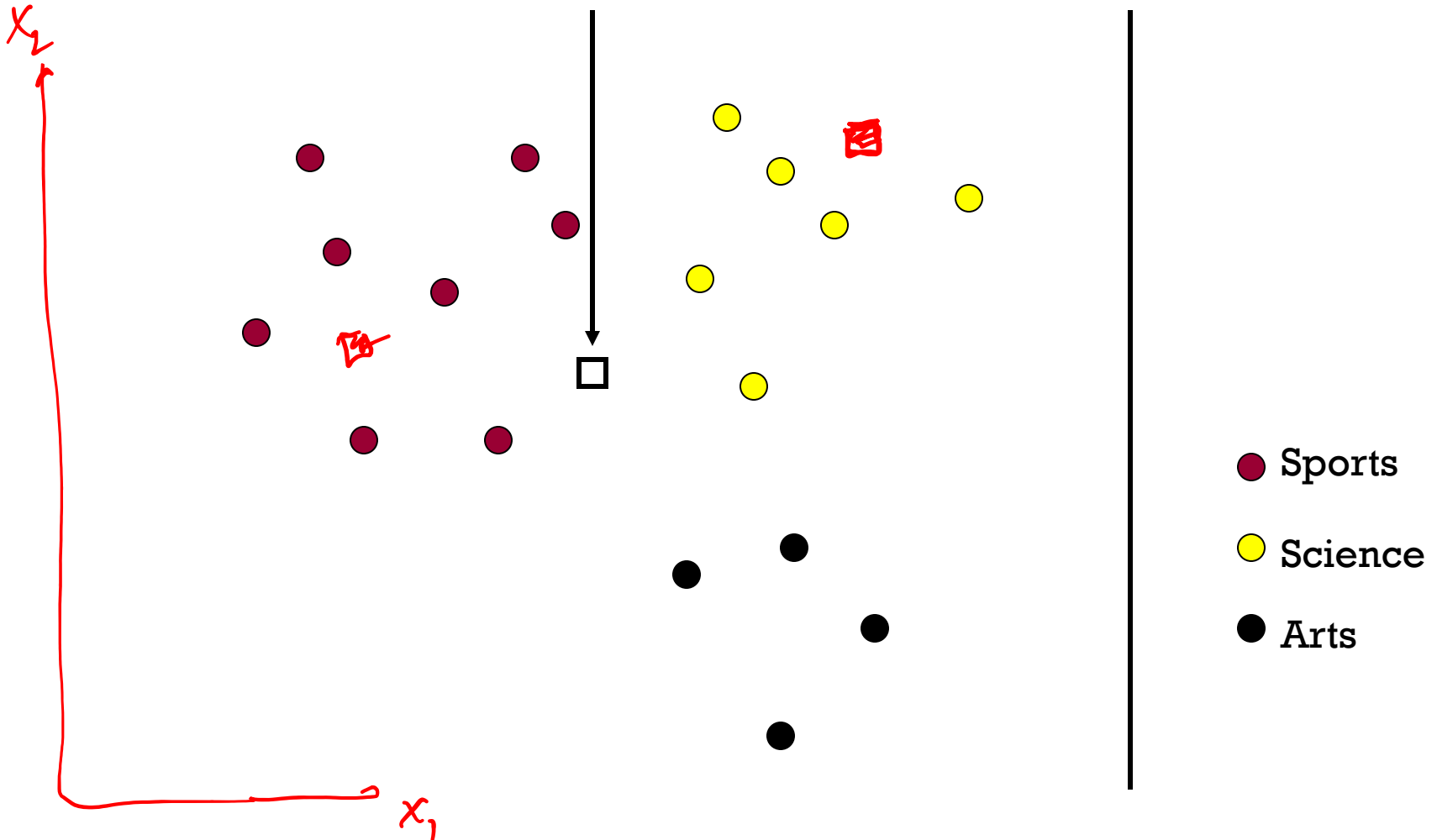
**Carnegie Mellon.**  
School of Computer Science

# k-NN classifier



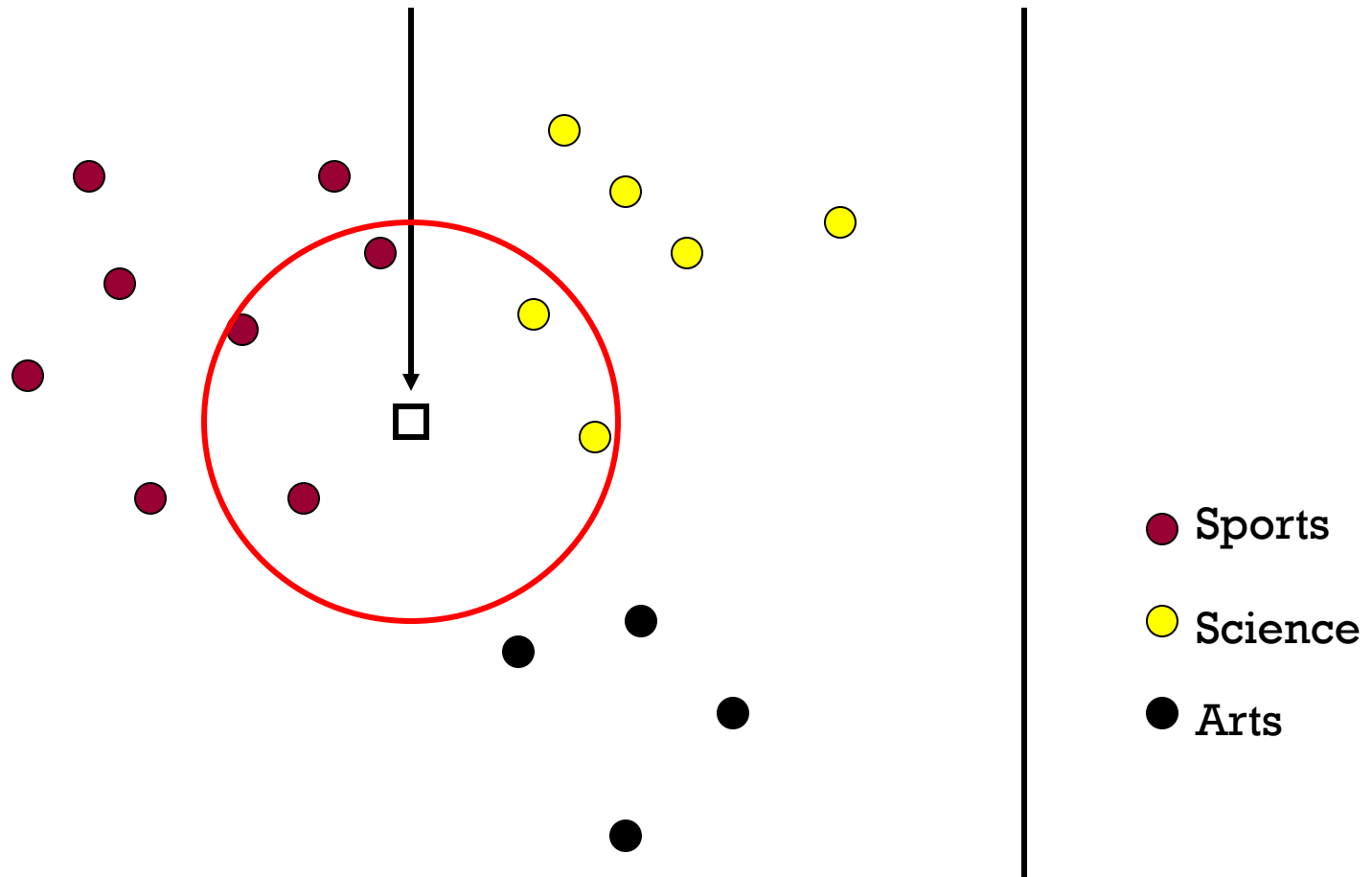
# k-NN classifier

Test document



# k-NN classifier (k=5)

Test document



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



# k-NN classifier

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

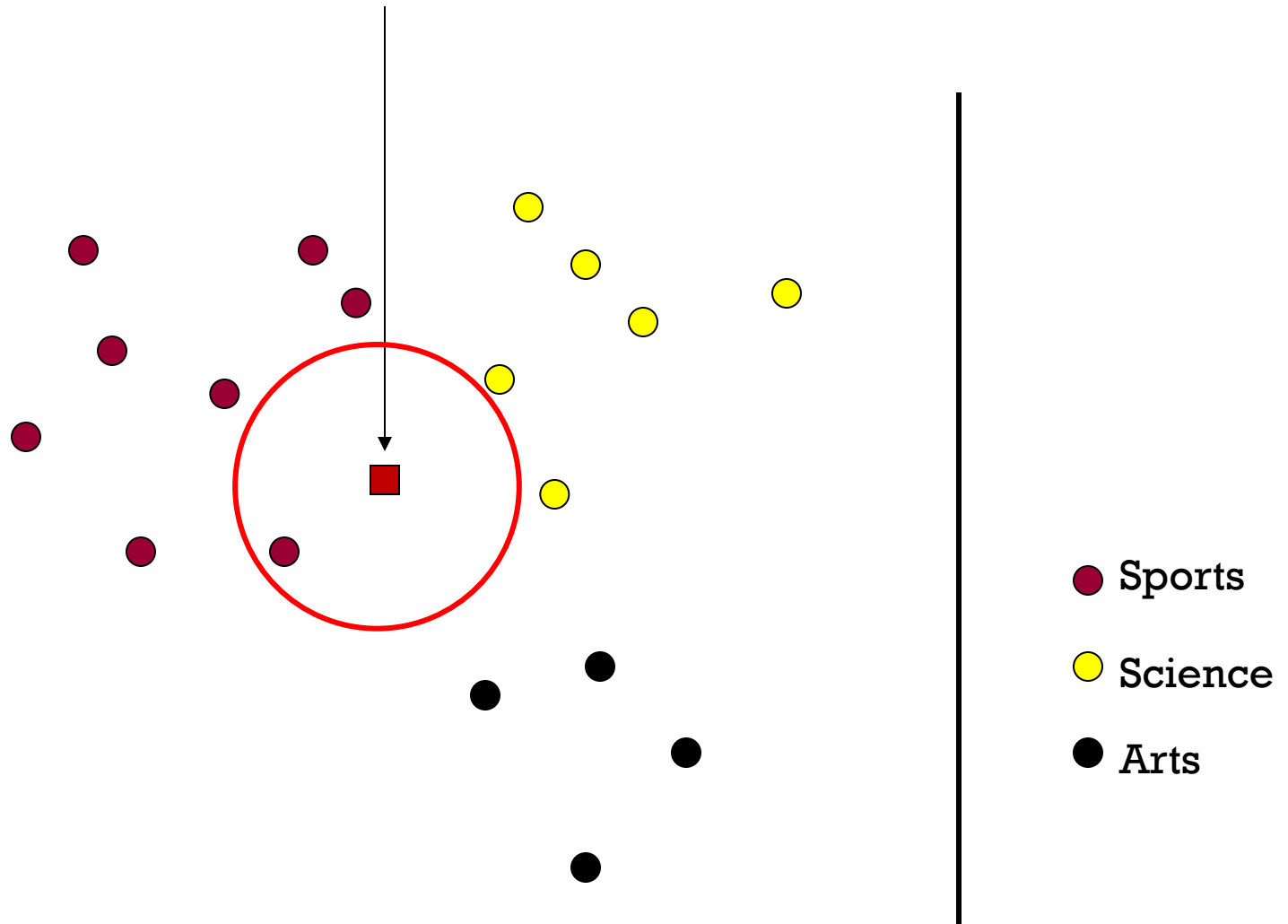
$$\hat{P}_{kNN}(x|y) = \frac{k_y}{n_y} \rightarrow \# \text{ training pts of class } y \text{ amongst } k \text{ NNs of } x \quad \checkmark \quad \sum_y k_y = k$$

$\frac{k_y}{n_y} \cdot \frac{n_y}{n}$

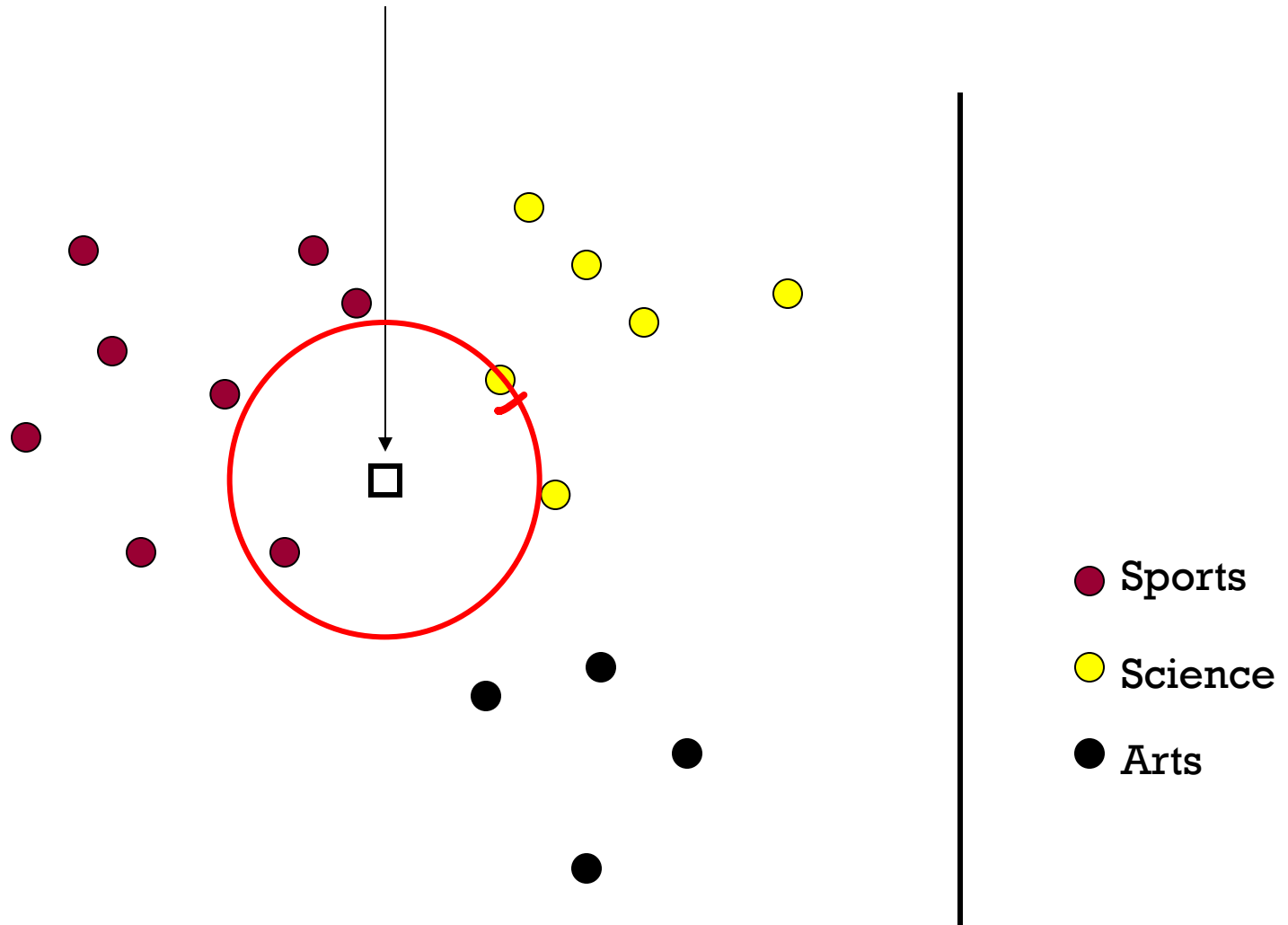
$\rightarrow \# \text{ total training pts of class } y \quad \checkmark$

$$\hat{P}(y) = \frac{n_y}{n} \leftarrow \begin{array}{l} \text{no. of training examples that have label } y \\ \text{total no. of training examples} \end{array}$$

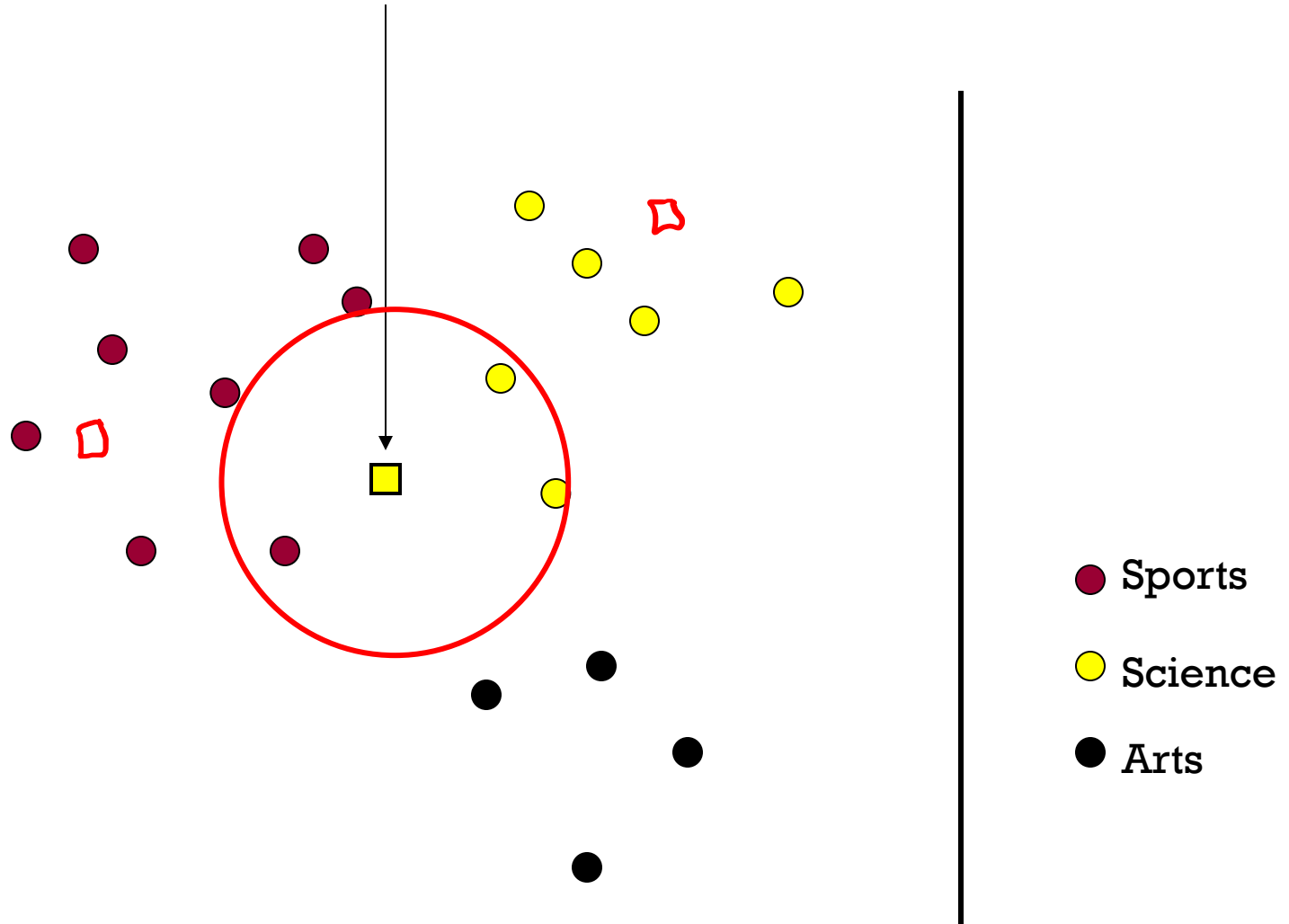
# 1-Nearest Neighbor (kNN) classifier



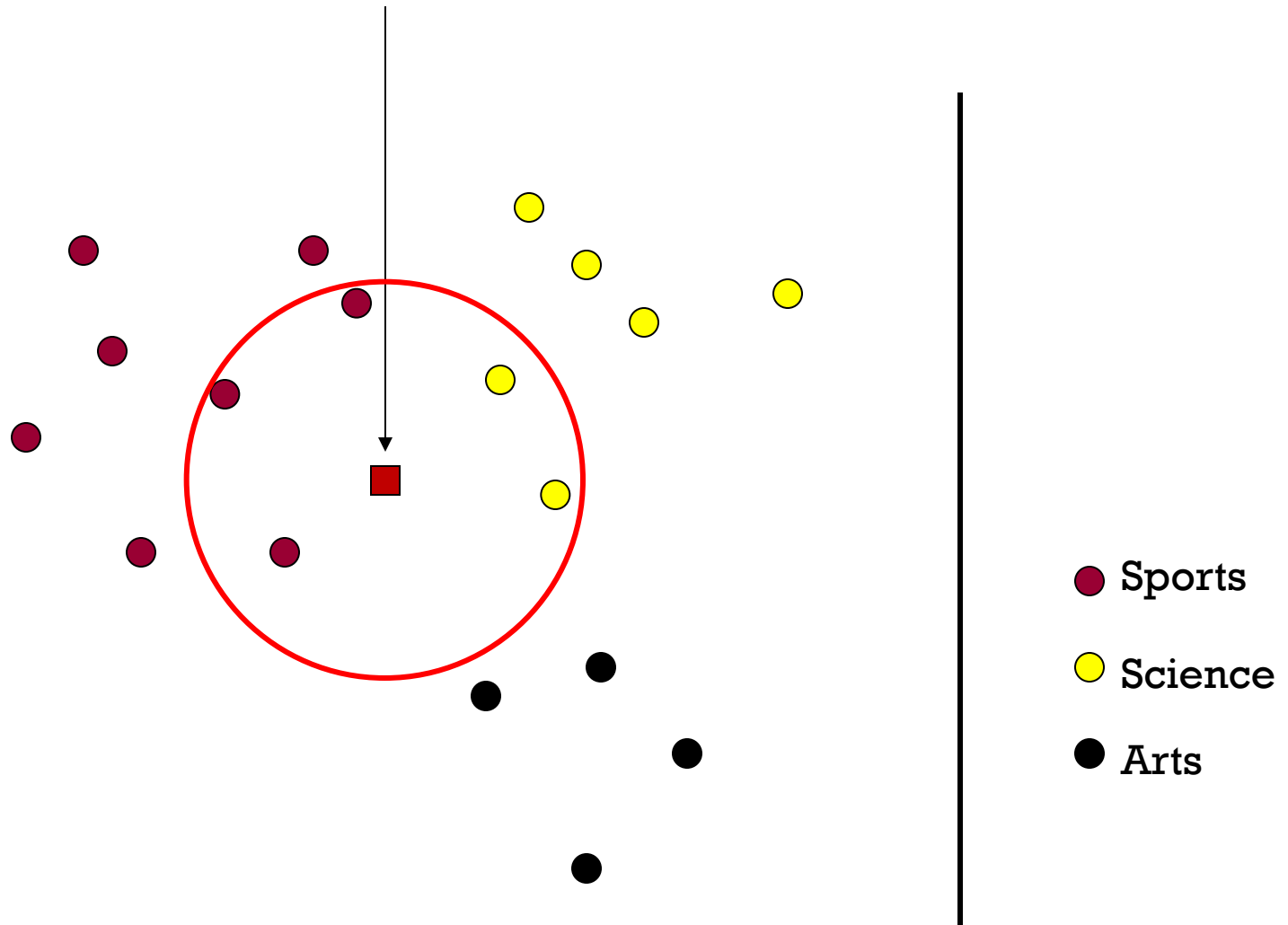
# 2-Nearest Neighbor (kNN) classifier



# 3-Nearest Neighbor (kNN) classifier

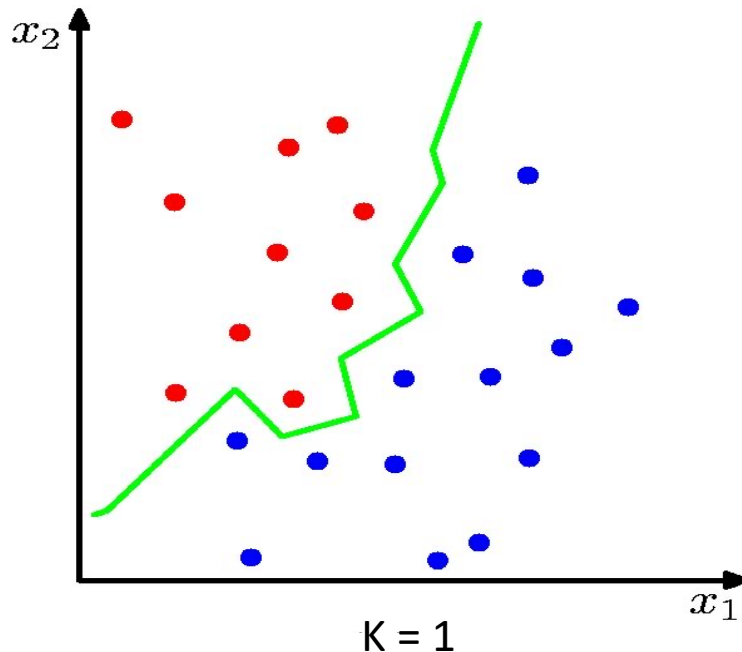


# 5-Nearest Neighbor (kNN) classifier

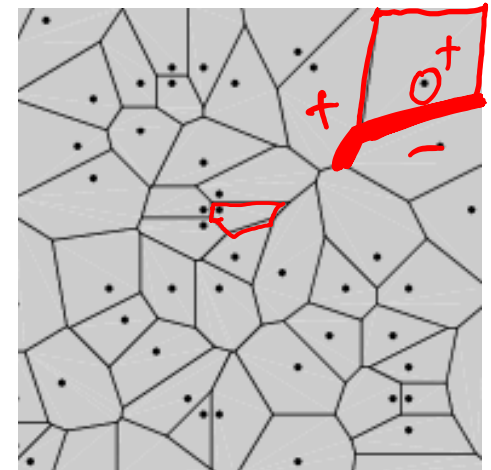


# What is the best k?

1-NN classifier decision boundary



Voronoi  
Diagram



*larger k  $\rightarrow$  smoother decision boundary*

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

# What is the best k?

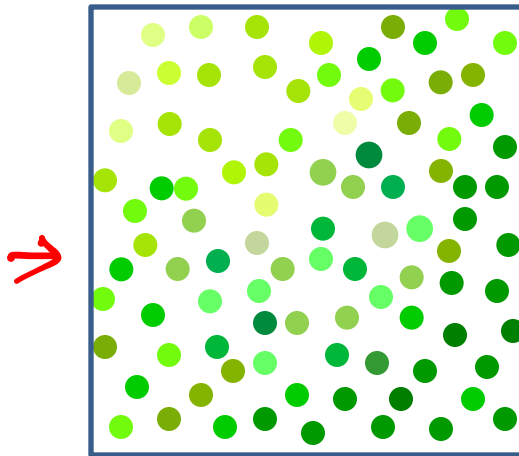
Approximation vs. Stability (aka Bias vs Variance) Tradeoff



- Larger  $K \Rightarrow$  predicted label is more stable (low variance) but potentially less accurate (high bias)
- Smaller  $K \Rightarrow$  predicted label can approximate best classifier well given enough data (low bias) but predict label is unstable (high variance)

# Local Kernel Regression

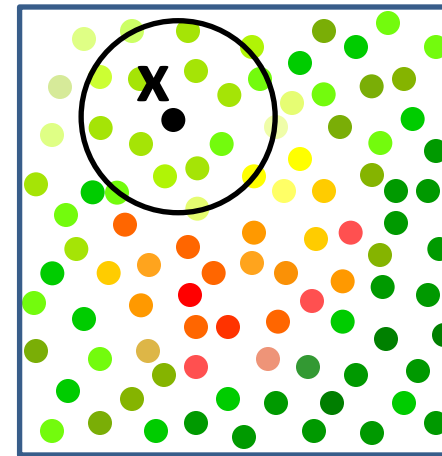
- What is the temperature in the room?



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

Average

at location  $x$ ?

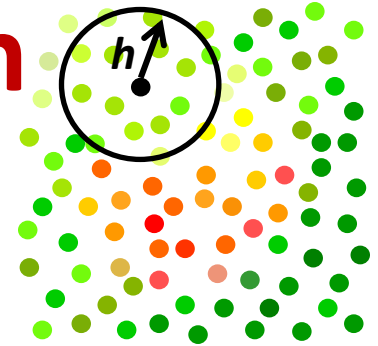


$$\hat{T}(x) = \frac{\sum_{i=1}^n Y_i \mathbf{1}_{\|X_i - x\| \leq h}}{\sum_{i=1}^n \mathbf{1}_{\|X_i - x\| \leq h}}$$

"Local" Average



# Local Kernel Regression



Global average  $w_i = \frac{1}{n} \Rightarrow \sum_i w_i = 1$

- Nonparametric estimator
- Nadaraya-Watson Kernel Estimator

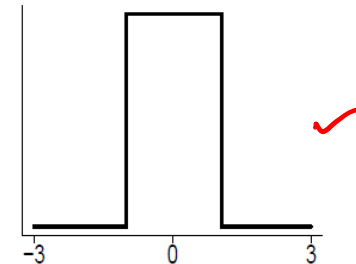
$$\hat{f}_n(X) = \sum_{i=1}^n w_i Y_i \quad \text{Where} \quad w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

$\Rightarrow \sum w_i = 1$

- 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 $h$

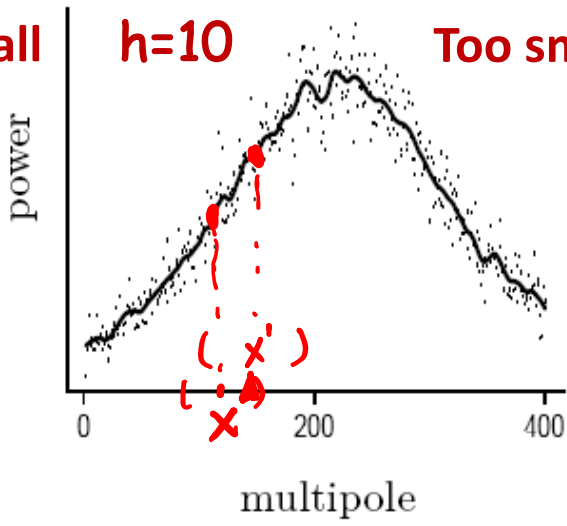
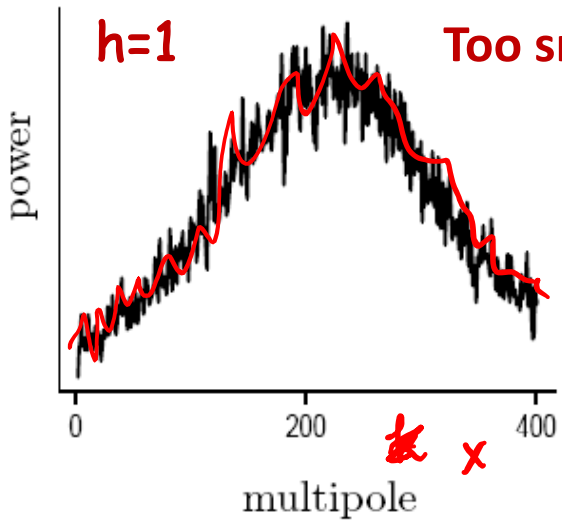
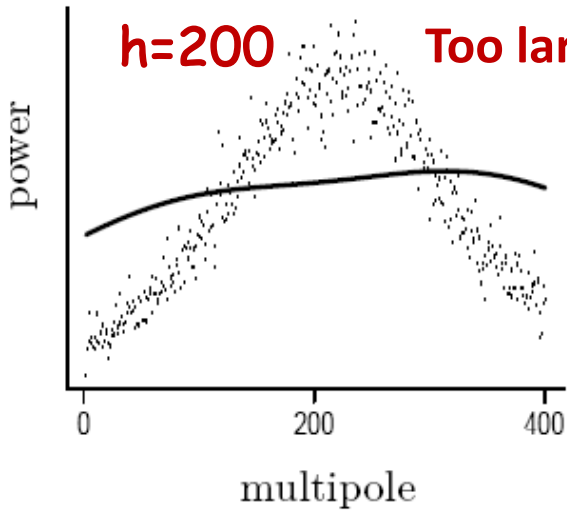
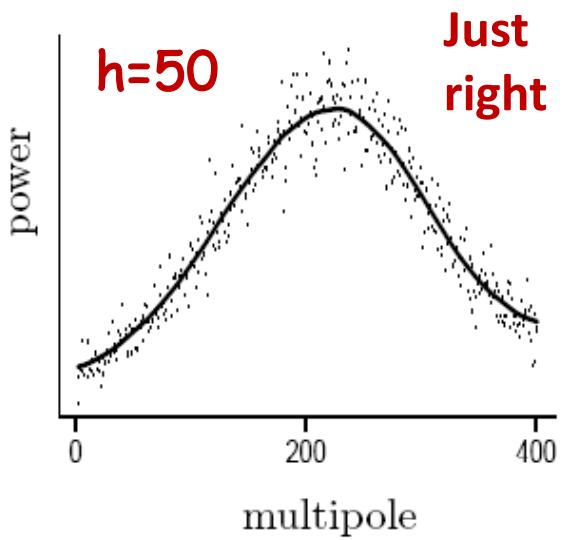


Image Source:  
Larry's book – All  
of Nonparametric  
Statistics



# Kernel Regression as Weighted Least Squares

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

=  $\uparrow \beta$

Weighted Least Squares

$$w_i(X) = \frac{e^{-\frac{(\|X-X_i\|^2)}{2h}} K\left(\frac{X-X_i}{h}\right) = K(X, X_i)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

if  $w_i = \frac{1}{n}$

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  $f(X_i) = \beta_x$  (a constant)

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

$\downarrow$   
 constant  
 $\underbrace{\hspace{10em}}$   
 $J(\beta)$

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

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

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

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

$$\Rightarrow \beta_x = \frac{\sum_i w_i Y_i}{\sum_i w_i}$$

# Local Linear/Polynomial Regression

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

Weighted Least Squares

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

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

(local polynomial of degree p around X)

$$f(x_i) = \beta$$

$$f(x_i) = \beta_0 + \beta_1(x_i - x)$$

# Summary

- Non-parametric approaches

**Four things make a nonparametric/memory/instance based/lazy 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) 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.