

# CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



**HVERFORD**  
COLLEGE

Sit somewhere new!

# Admin

- Sorelle office hours **TODAY: 4-5pm in H110**
- **Lab 2** due Thursday Feb 8 (week from today)

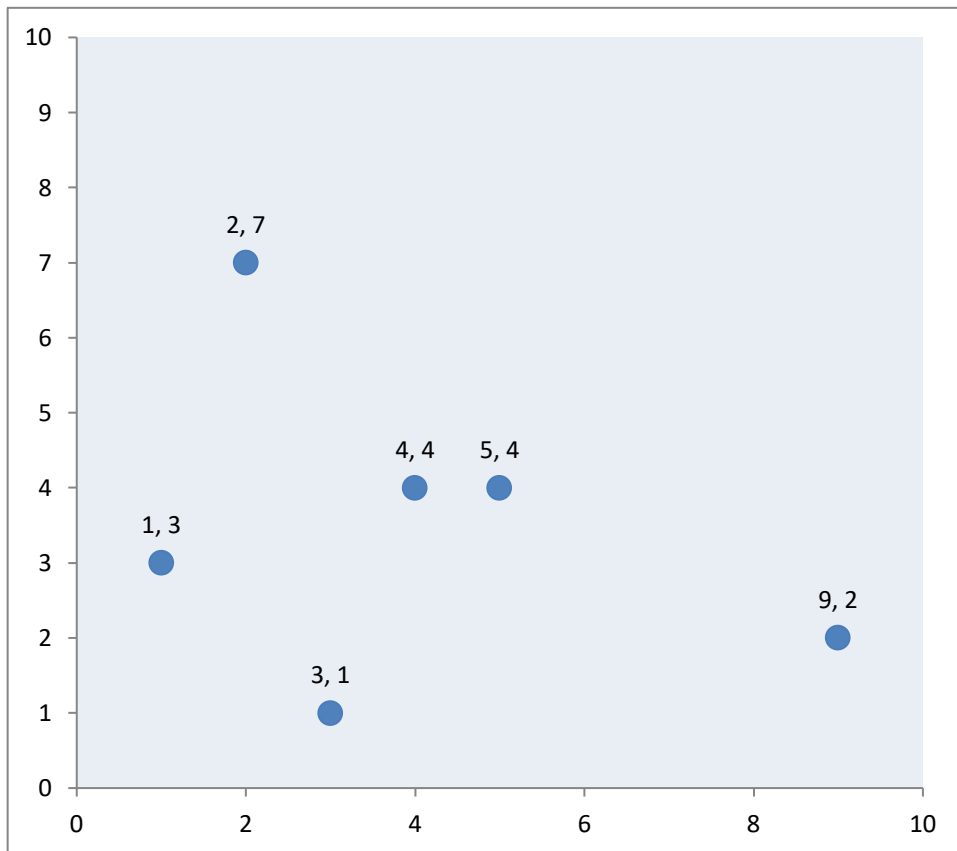
# Outline for Feb 1

- Finish KD Trees
  - Nearest neighbor algorithm
  - Extending to  $k > 1$
- Evaluation metrics beyond CS260
  - AUC
  - Precision/recall curves

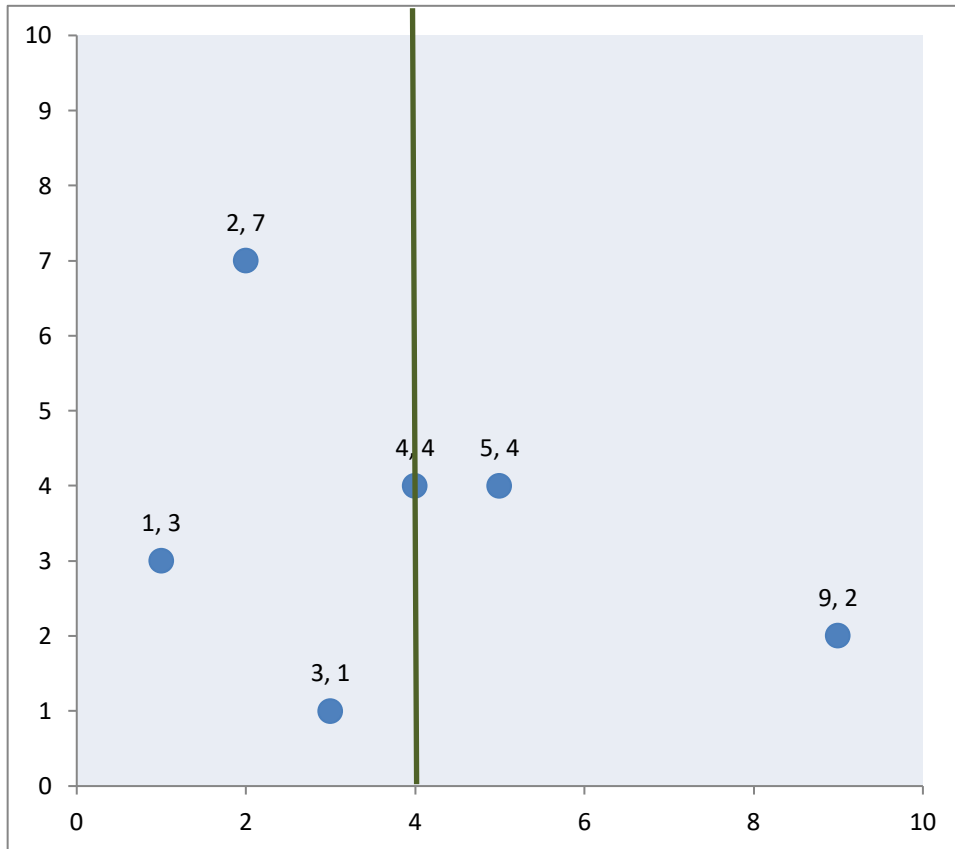
# Outline for Feb 1

- Finish KD Trees
  - Nearest neighbor algorithm
  - Extending to  $k > 1$
- Evaluation metrics beyond CS260
  - AUC
  - Precision/recall curves

# Making a kd-tree

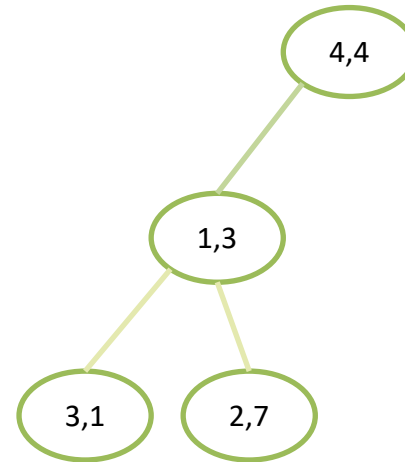
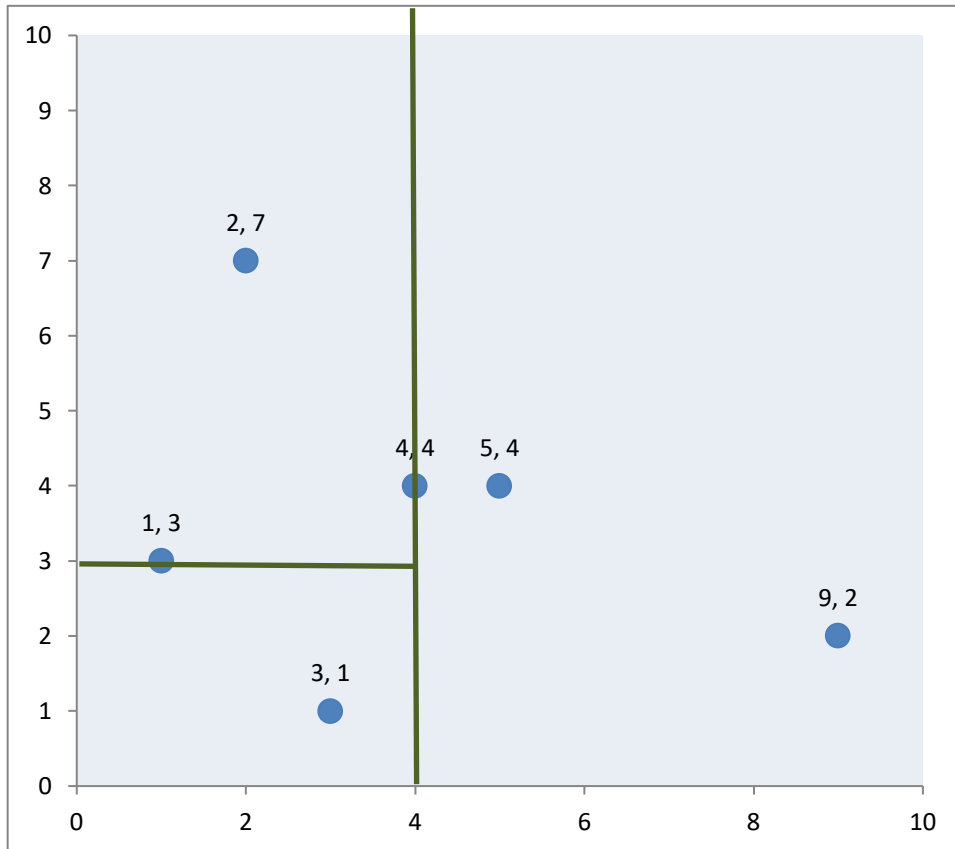


# Making a kd-tree



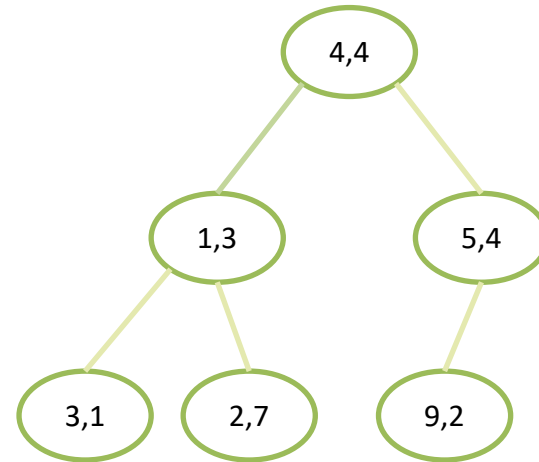
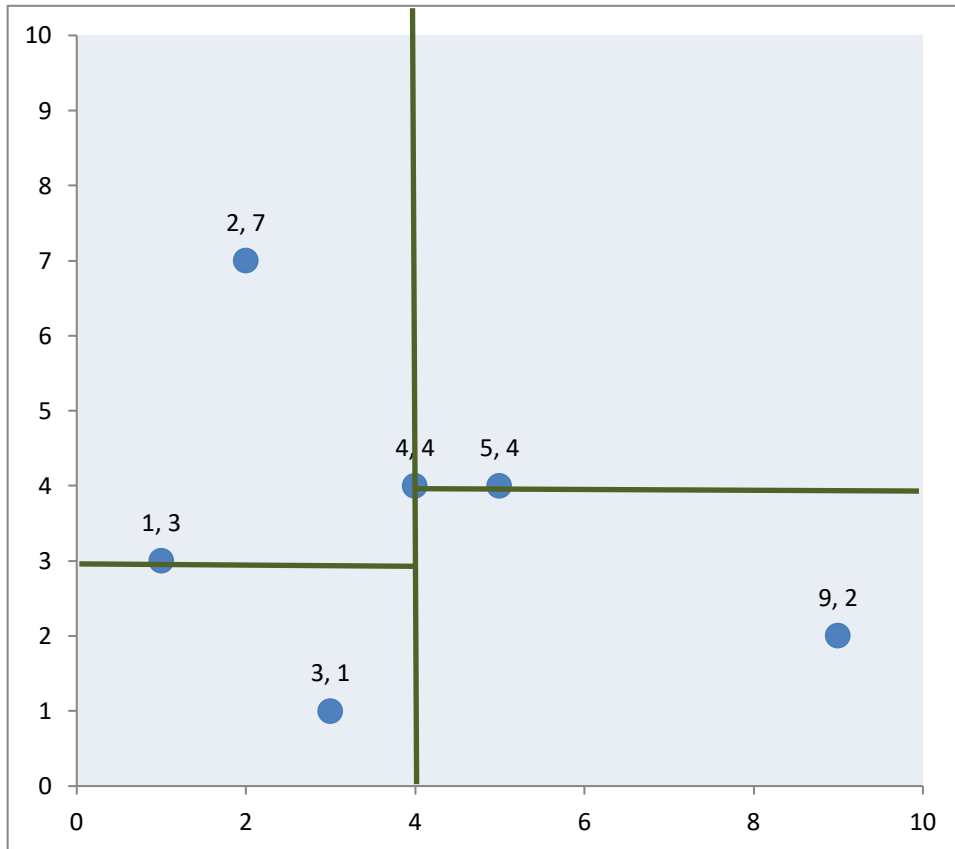
4,4

# Making a kd-tree

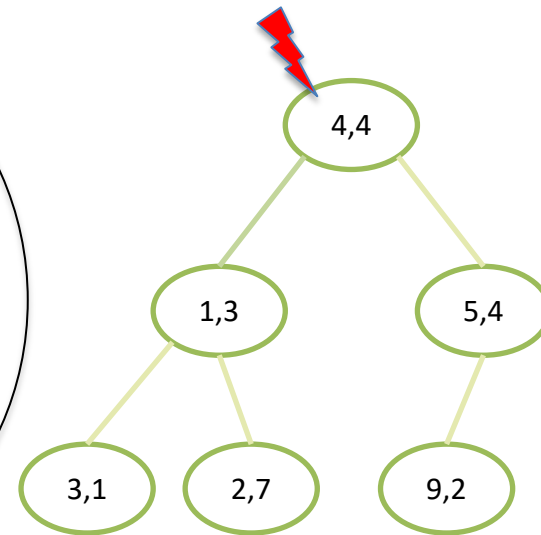
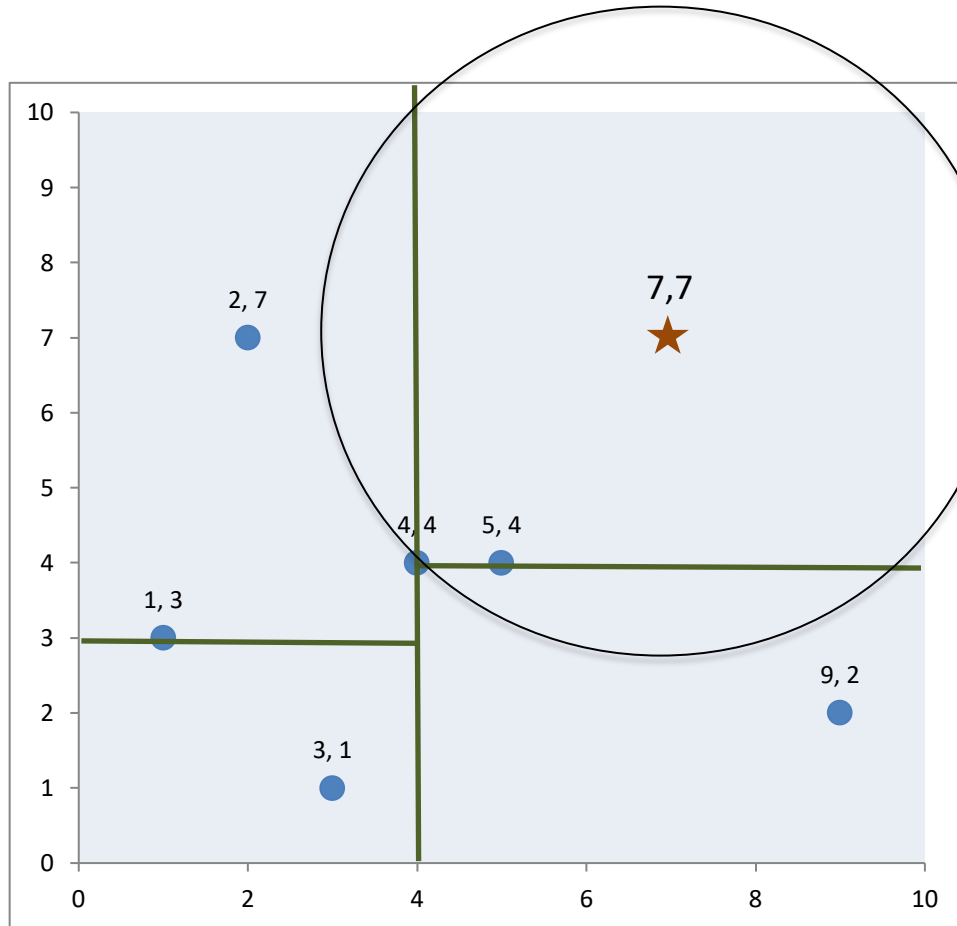




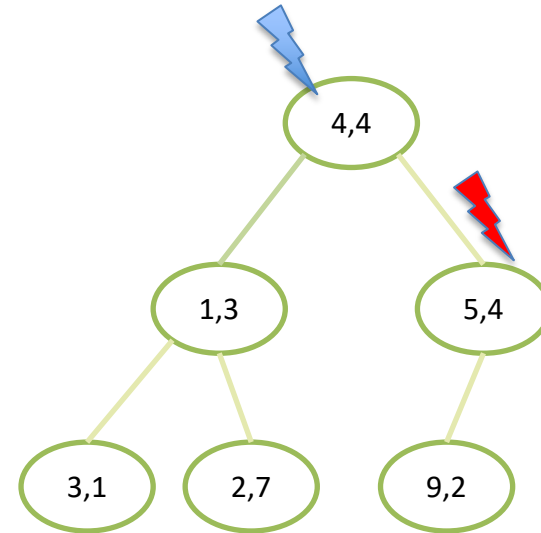
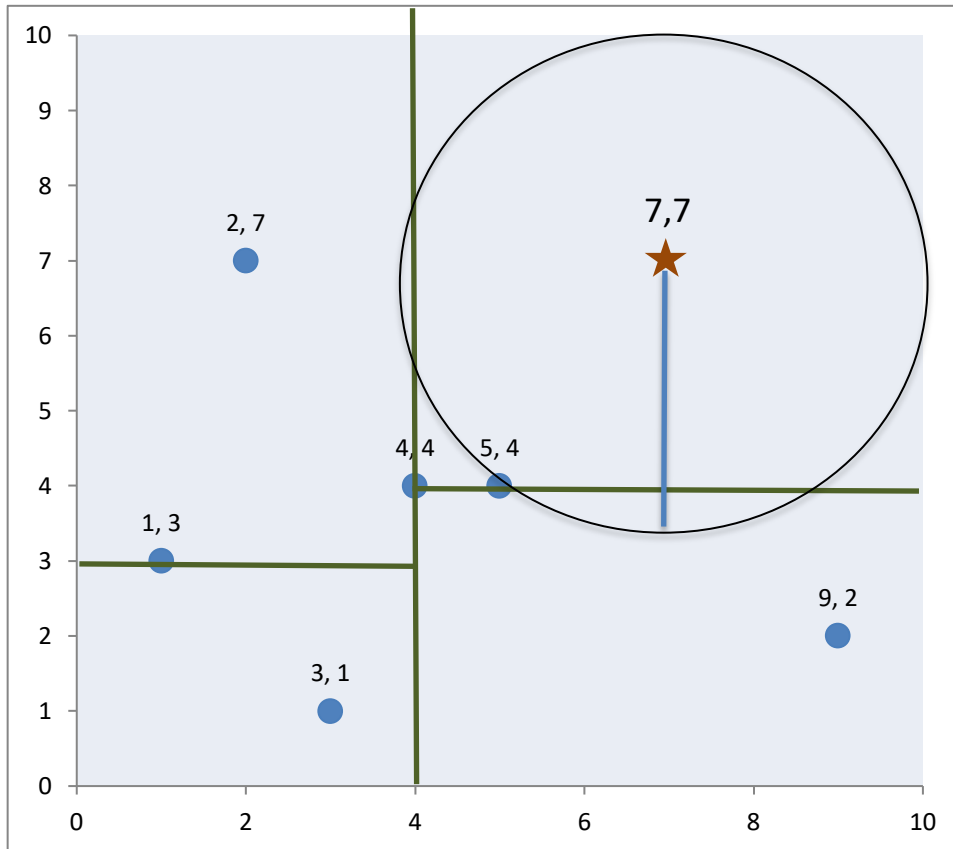
# Making a kd-tree



# Finding the nearest neighbor

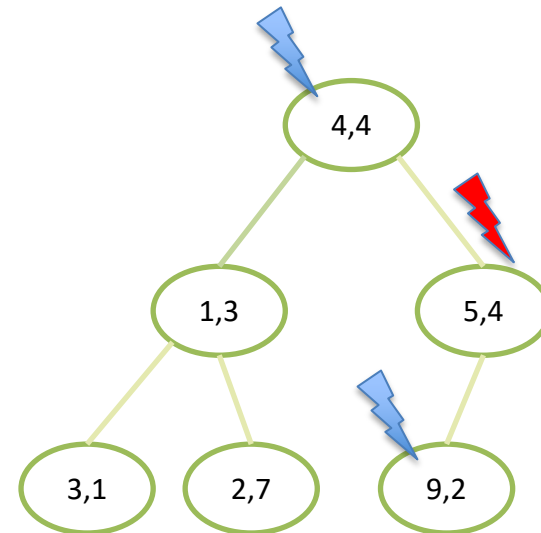
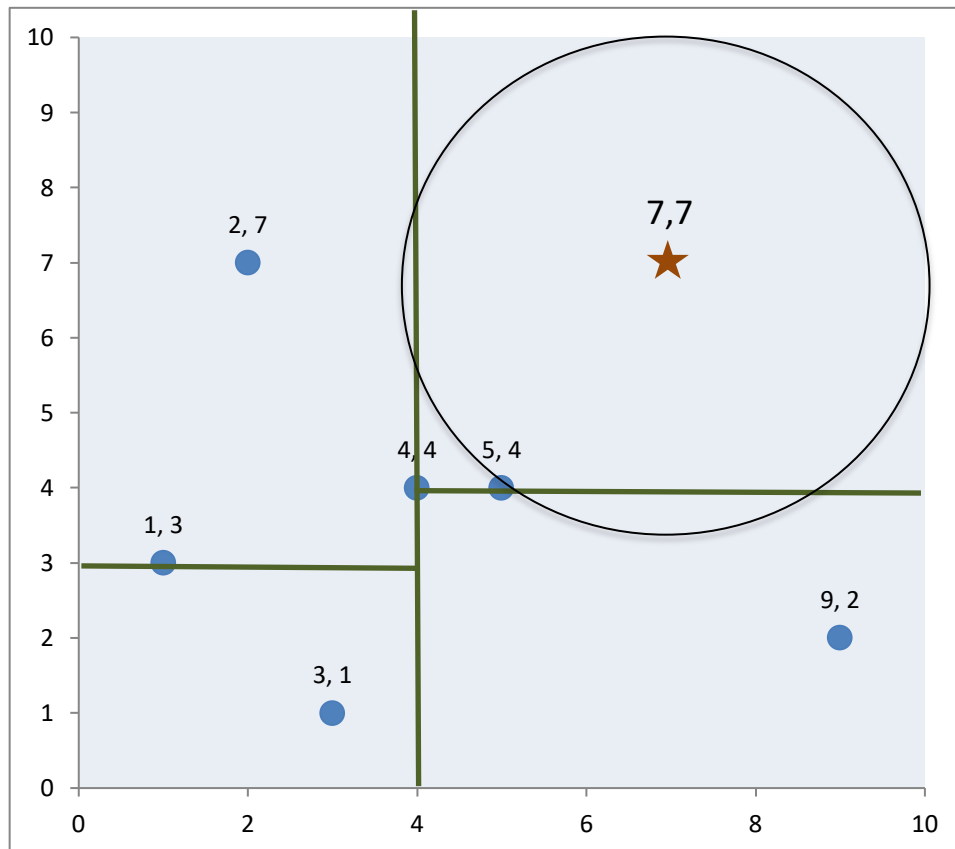


# Finding the nearest neighbor



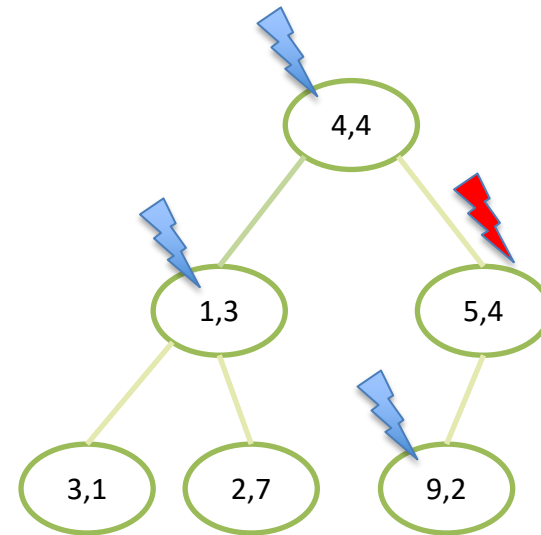
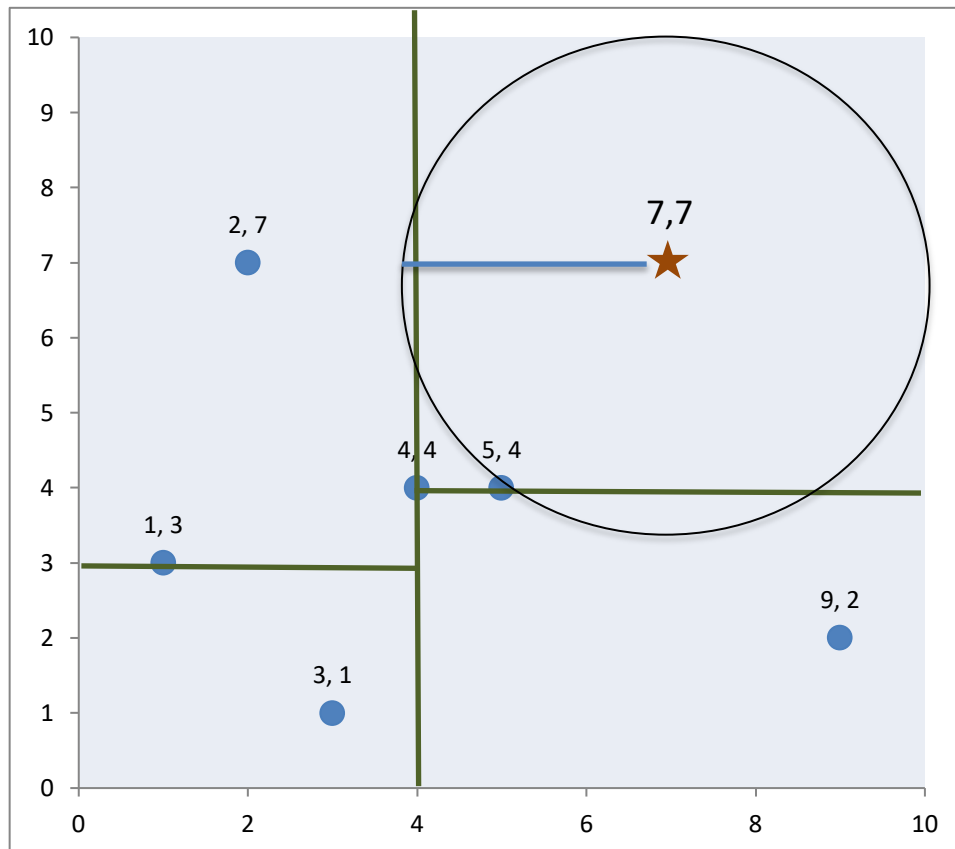
No right child, but need to check left  
Due to intersecting hypersphere

# Finding the nearest neighbor



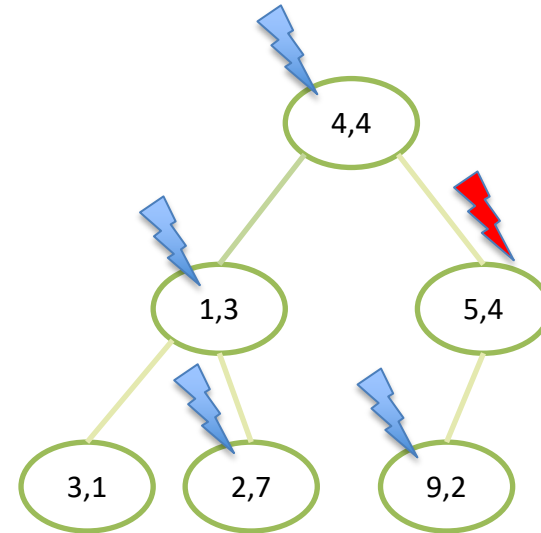
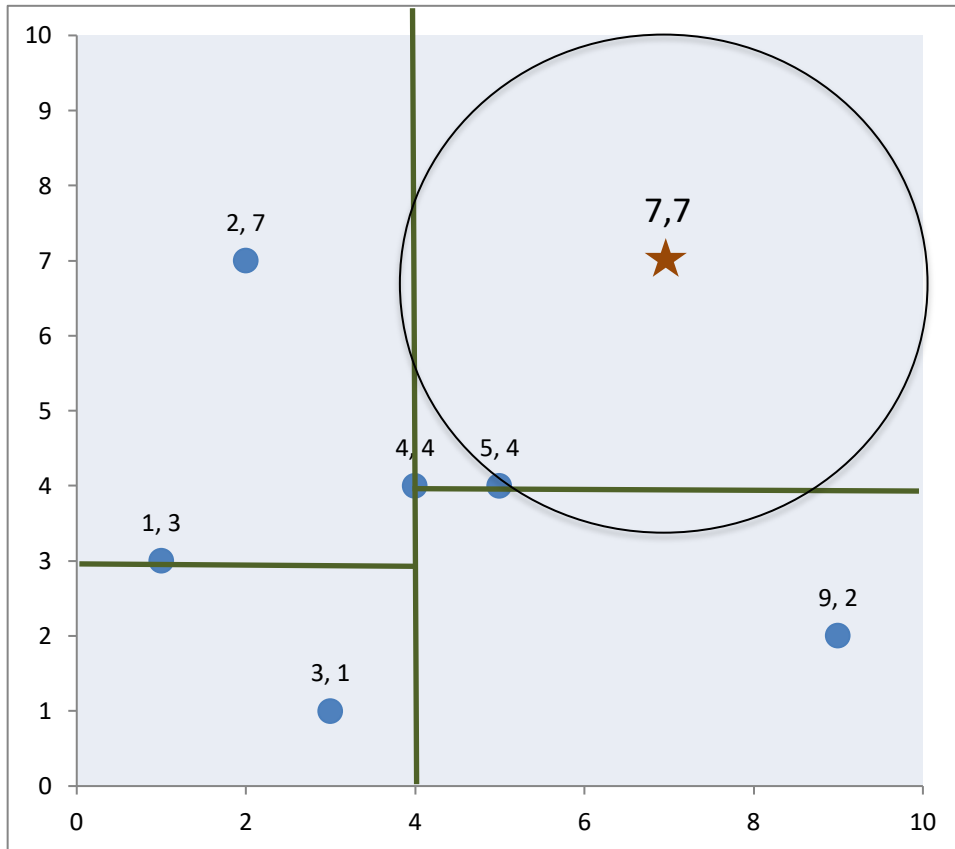
(9,2) is not closer, so keep (5,4) as nn

# Finding the nearest neighbor



Visit (1,3) due to intersecting hypersphere

# Finding the nearest neighbor



# Nearest neighbor recursive algorithm

leaf case

```
get_nn(root, q, min_d, nn)
  if root is a leaf:
    if d(root, q) < min_d:
      return d(root, q), root
    else:
      return min_d, nn
```

non-leaf

```
  else:
    if d(root, q) < min_d:
      ← min_d, nn = d(root, q), root
```

child = root of subtree  
dim { containing q  
min\_d, nn = get\_nn(child,  
q, min\_d, nn)  
# hypersphere / hyperplane  
if  $|q.\text{dim} - \text{split} - \text{root.dim}| < \text{min}_d$ :  
 return get\_nn(other\_child,  
q, min\_d, nn)  
else:  
 return min\_d, nn

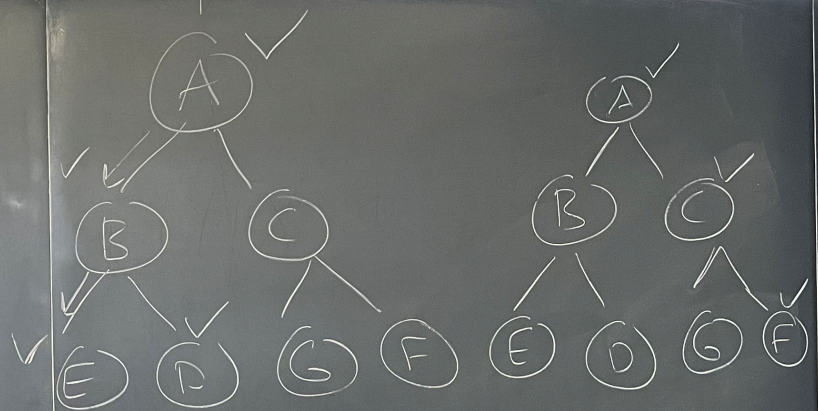
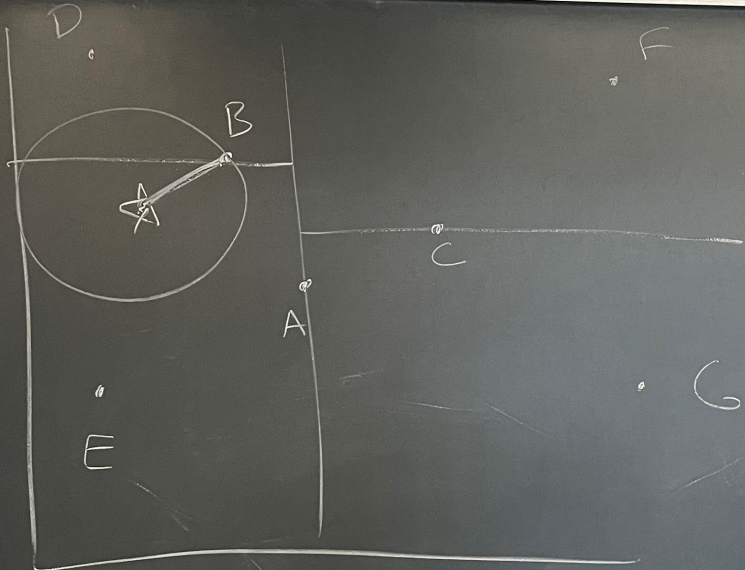
first line

dim = depth  $\propto$  d  
cycling                      increasing

$\uparrow$   
total # dimensions



# Handout 4



4 nodes \*

naive

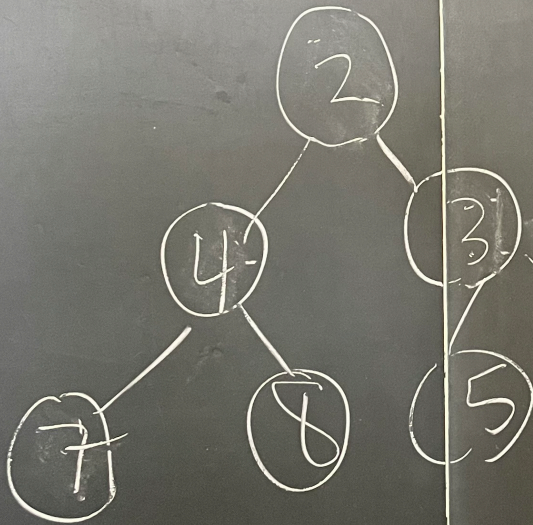
⇒ sort by distance,  
take top  $k$

distances ⇒

7, 2, 3, 4, 8, 5, 1

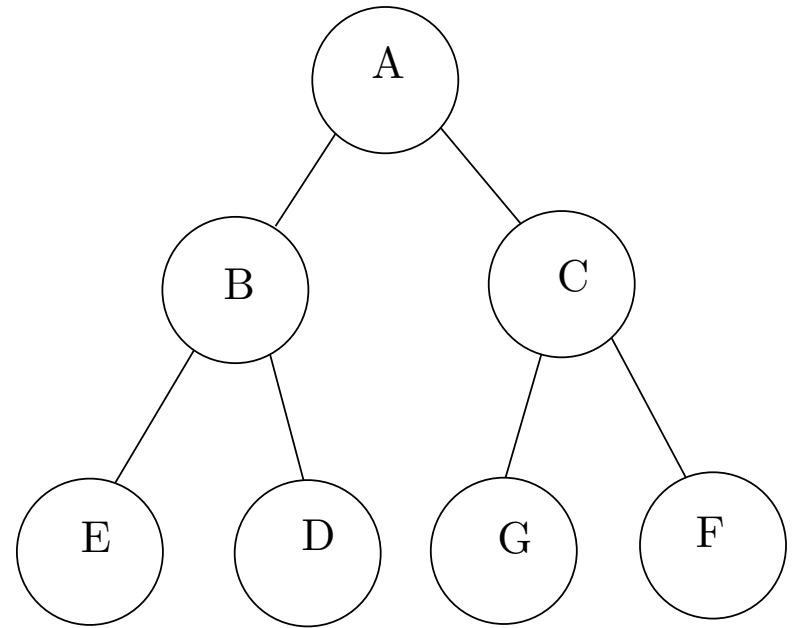
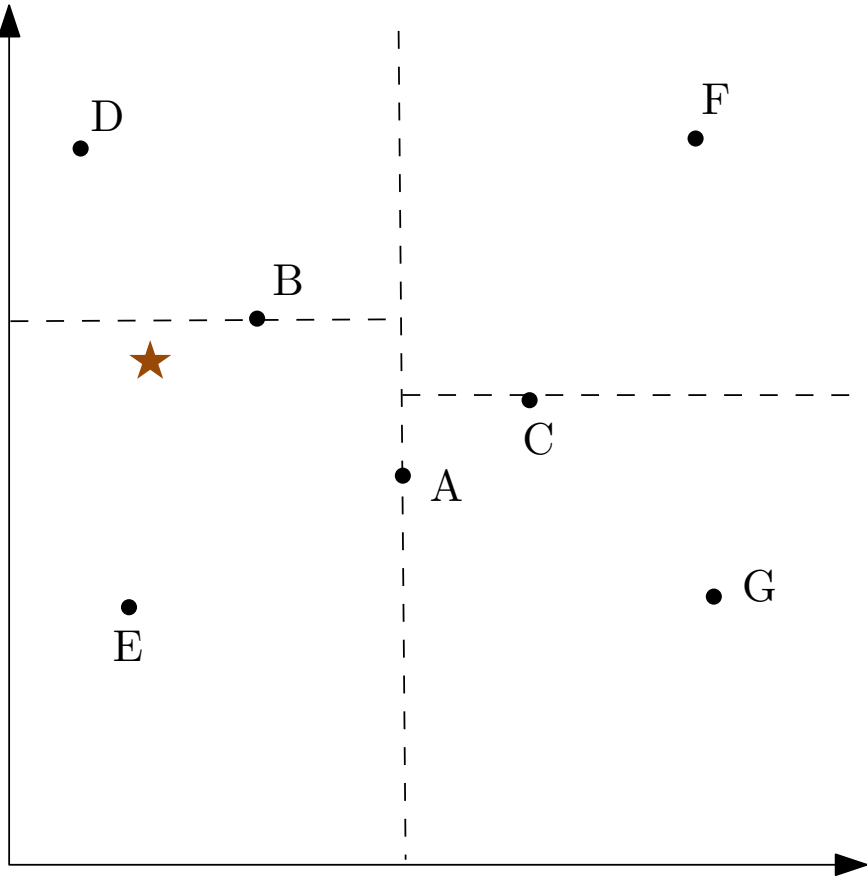
take off 1

take off 2

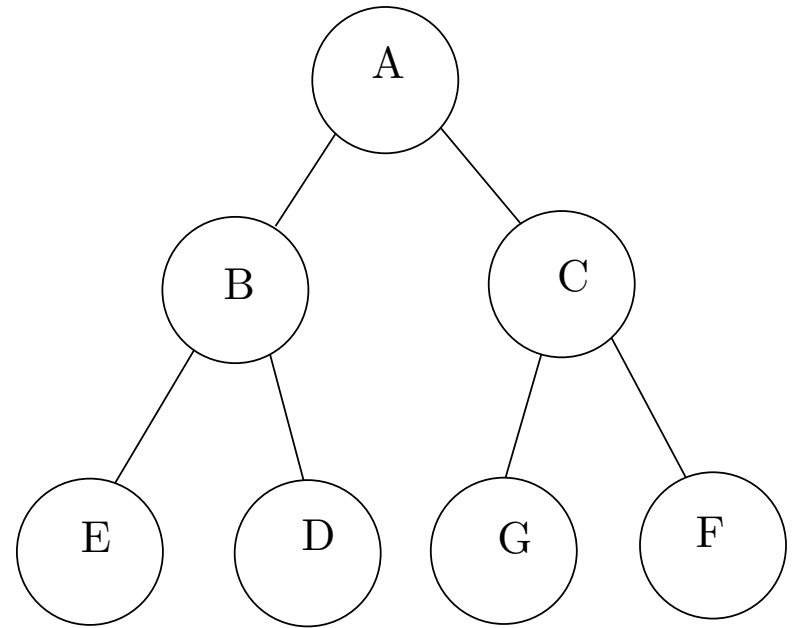
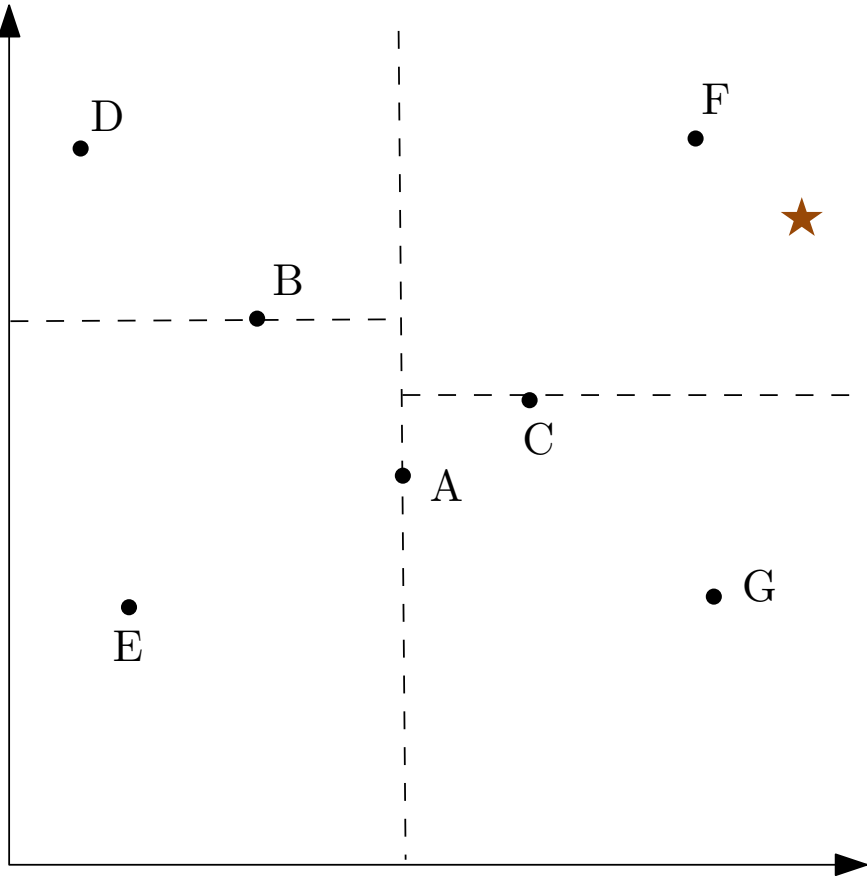


Extending to  $k > 1$

# Handout 4, example 1



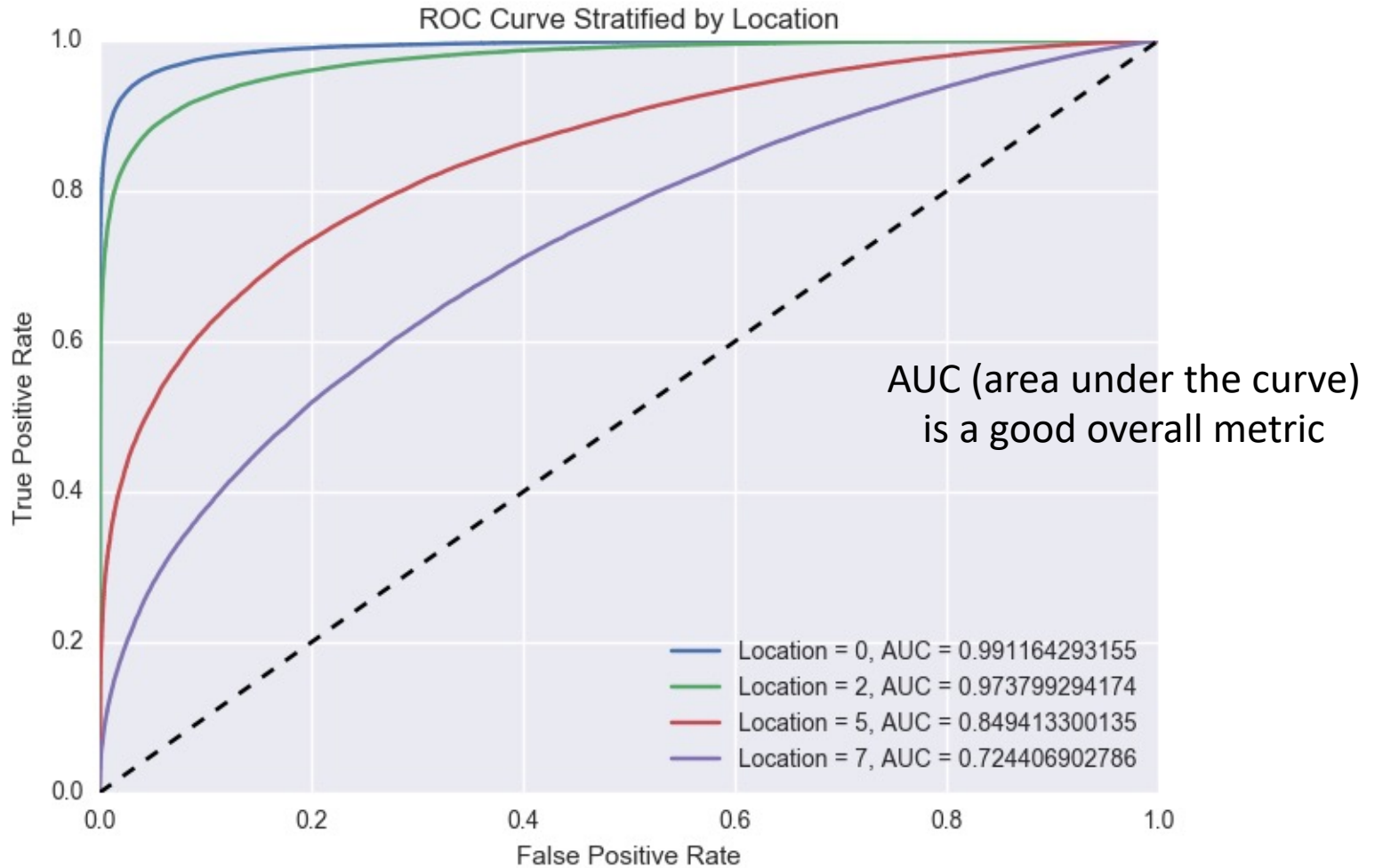
# Handout 4, example 2



# Outline for Feb 1

- Finish KD Trees
  - Nearest neighbor algorithm
  - Extending to  $k > 1$
- Evaluation metrics beyond CS260
  - AUC
  - Precision/recall curves

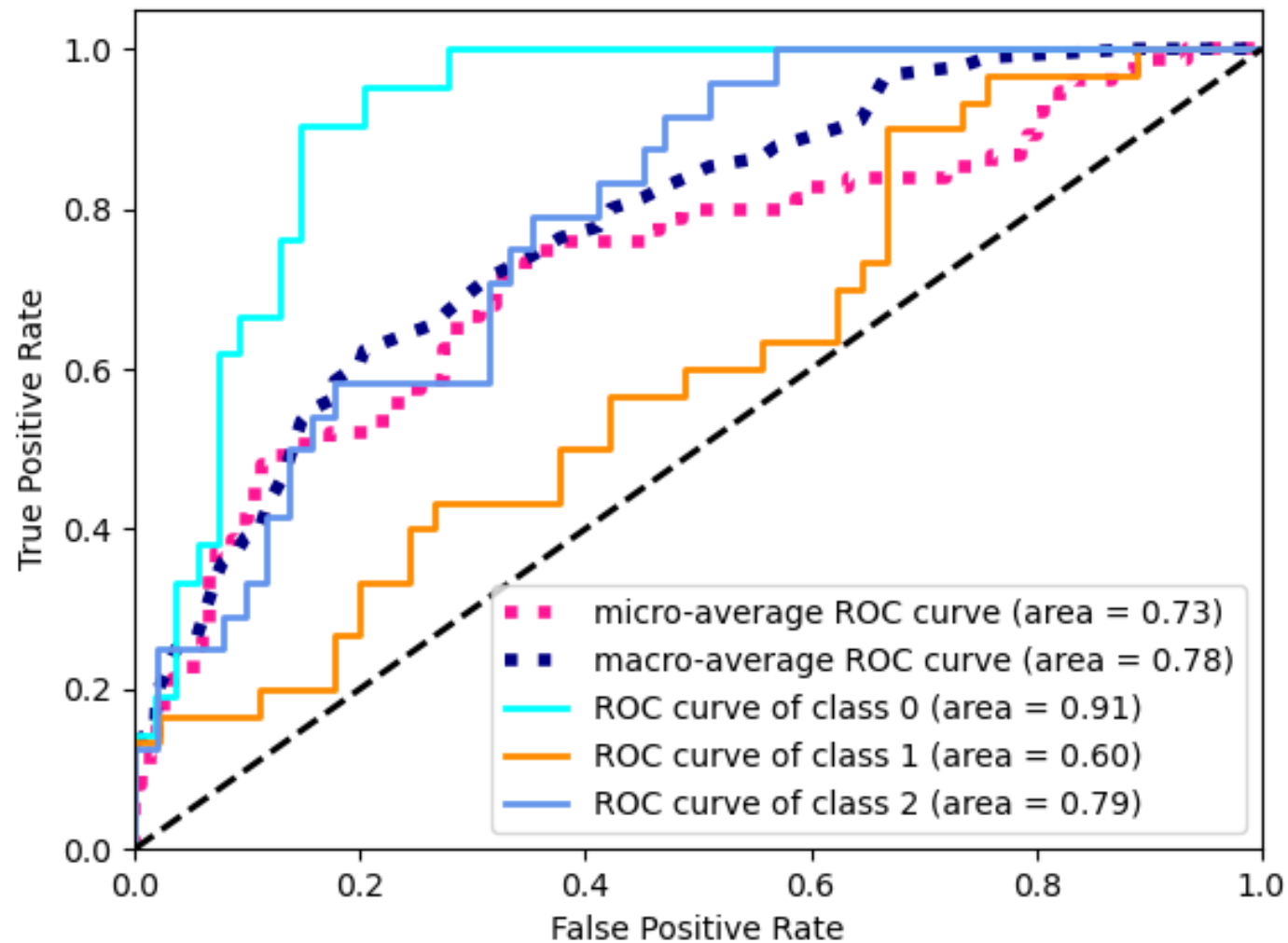
# How to compare ROC curves? AUC (area under the curve)



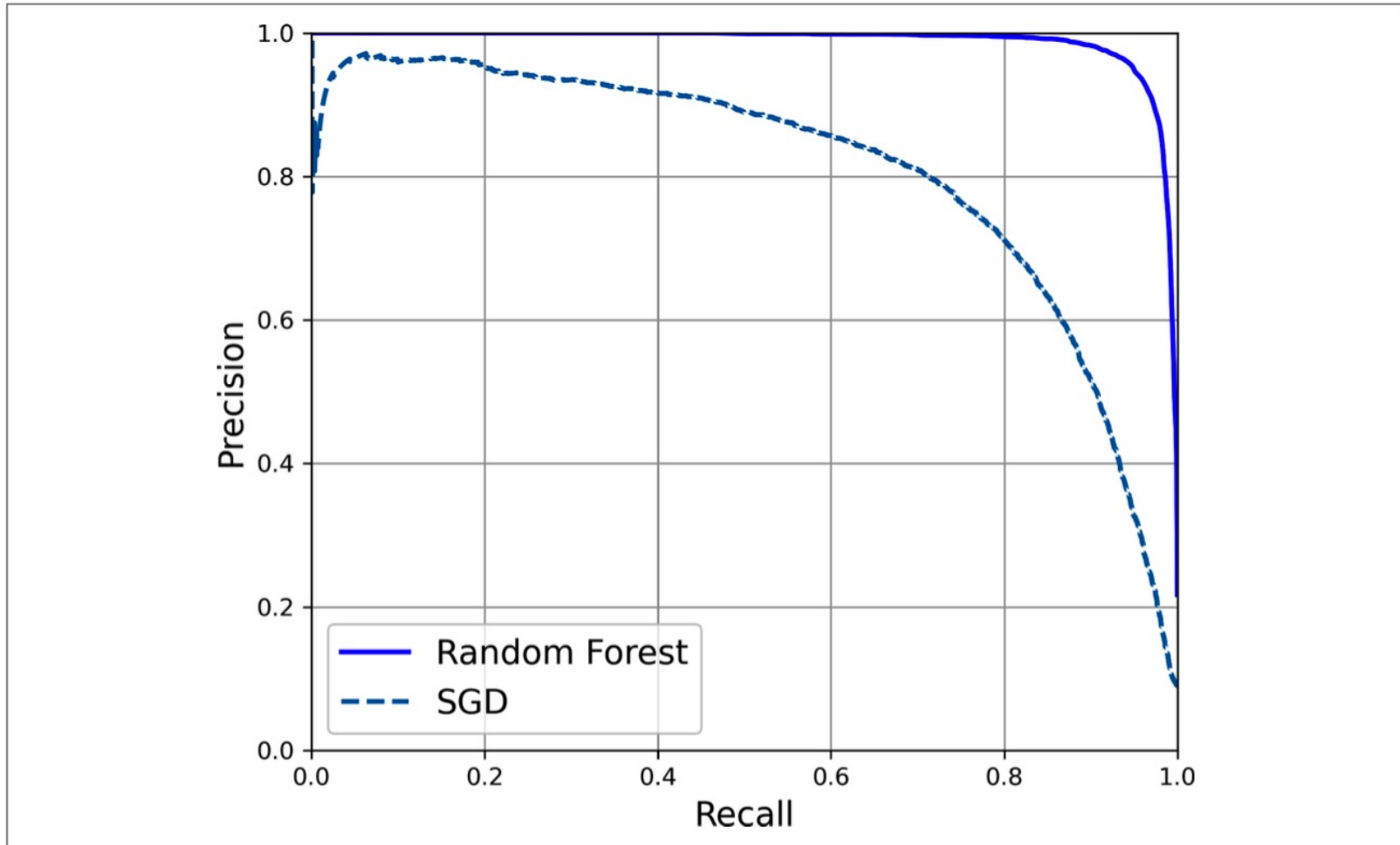
Example of a ROC curve from my research  
Chan, Perrone, Spence, Jenkins, Mathieson, Song

# AUC example

Some extension of Receiver operating characteristic to multiclass



# Precision/Recall curve



*Figure 3-8. Comparing PR curves: the random forest classifier is superior to the SGD classifier because its PR curve is much closer to the top-right corner, and it has a greater AUC*



# Precision/Recall

example from my research

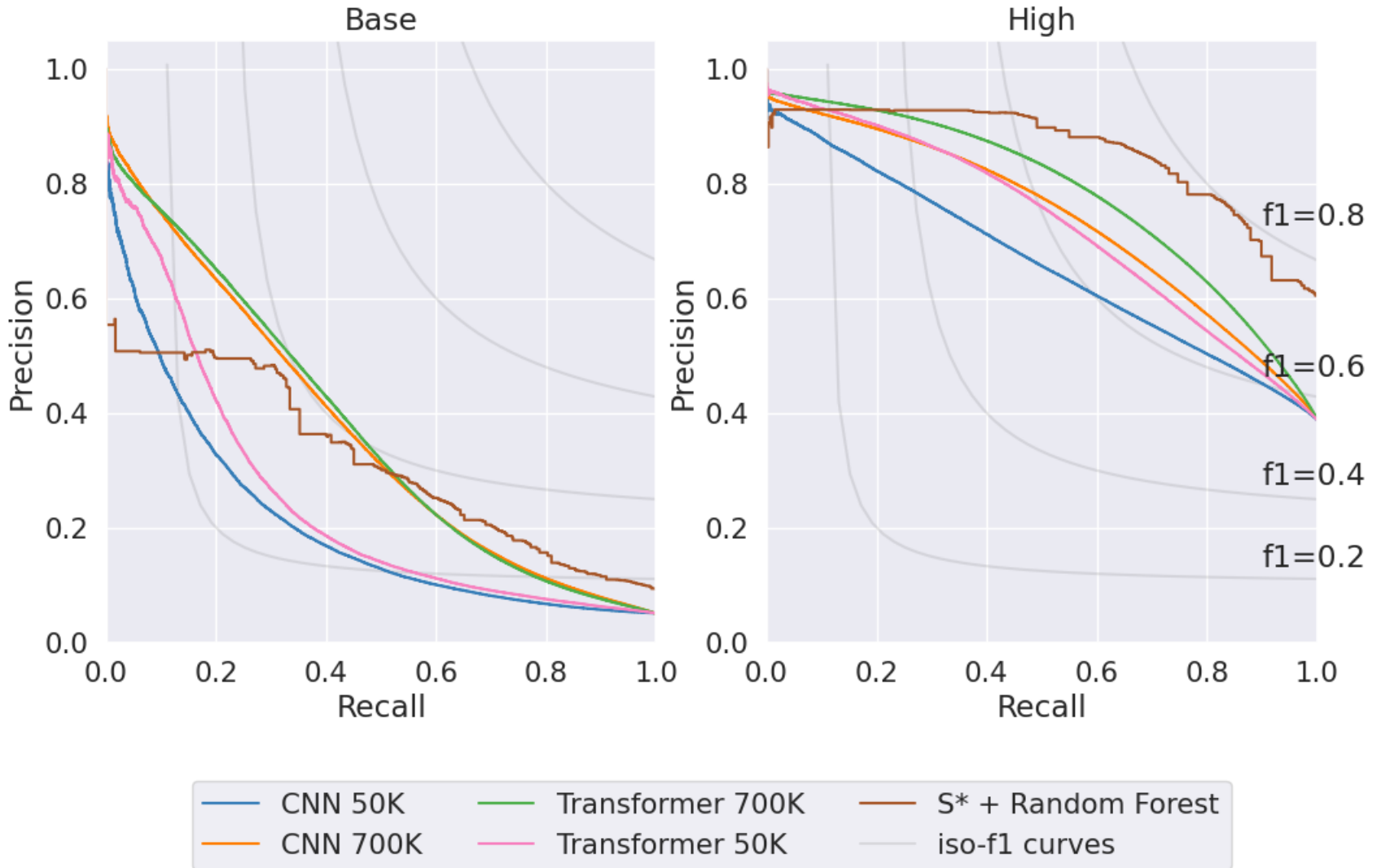
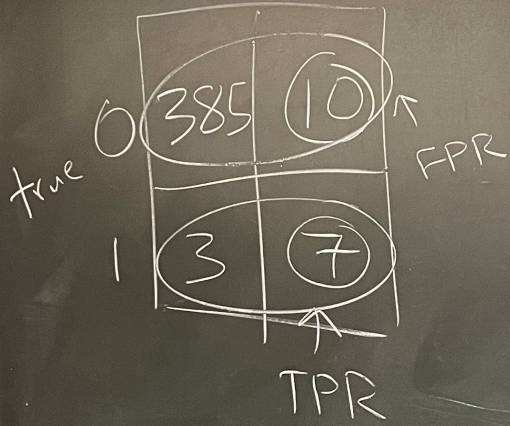
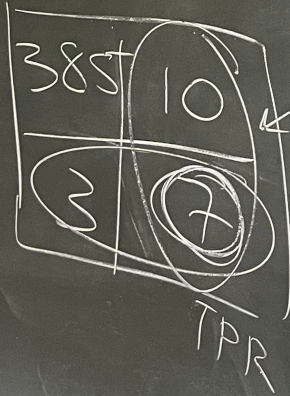
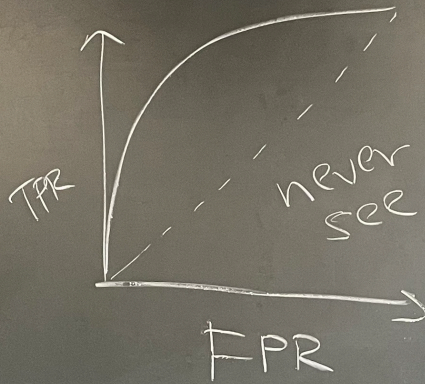


Figure: Jordan Cahoon

0 <sup>pred</sup> 1



ROC



PR curve

