

# CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



**HVERFORD**  
COLLEGE

Sit somewhere new!

# Admin

- **EVERYONE**: Sign in again
- Sorelle office hours **Thursday: 4-5pm in H110**
- **Lab 1** was due last night
- **Lab 2** due Thursday Feb 8
  - Don't wait til the last minute!
- **TA hour schedule** on Piazza

# Outline for Jan 30

- Python style and implementation notes
- Overfitting
- K-nearest neighbors
- KD Trees

Logistic Regression and Gradient Descent Review:  
Moved to discussion of softmax



# Outline for Jan 30

- Python style and implementation notes
- Overfitting
- K-nearest neighbors
- KD Trees

# Python style

- Decompose code into natural functions
- Avoid global variables (sometimes useful)
- Include a file header with purpose, author, and date
- Include headers for each function
- No lines over 80 chars
- Variable names implicitly show type
- Include line breaks and comments!

# Python style

- “Snake-case” not “camel-case”
  - ~~linearSearch~~
  - linear\_search
- Alphabetize imports and don't use “\*”
  - ~~from numpy import \*~~
  - import numpy as np

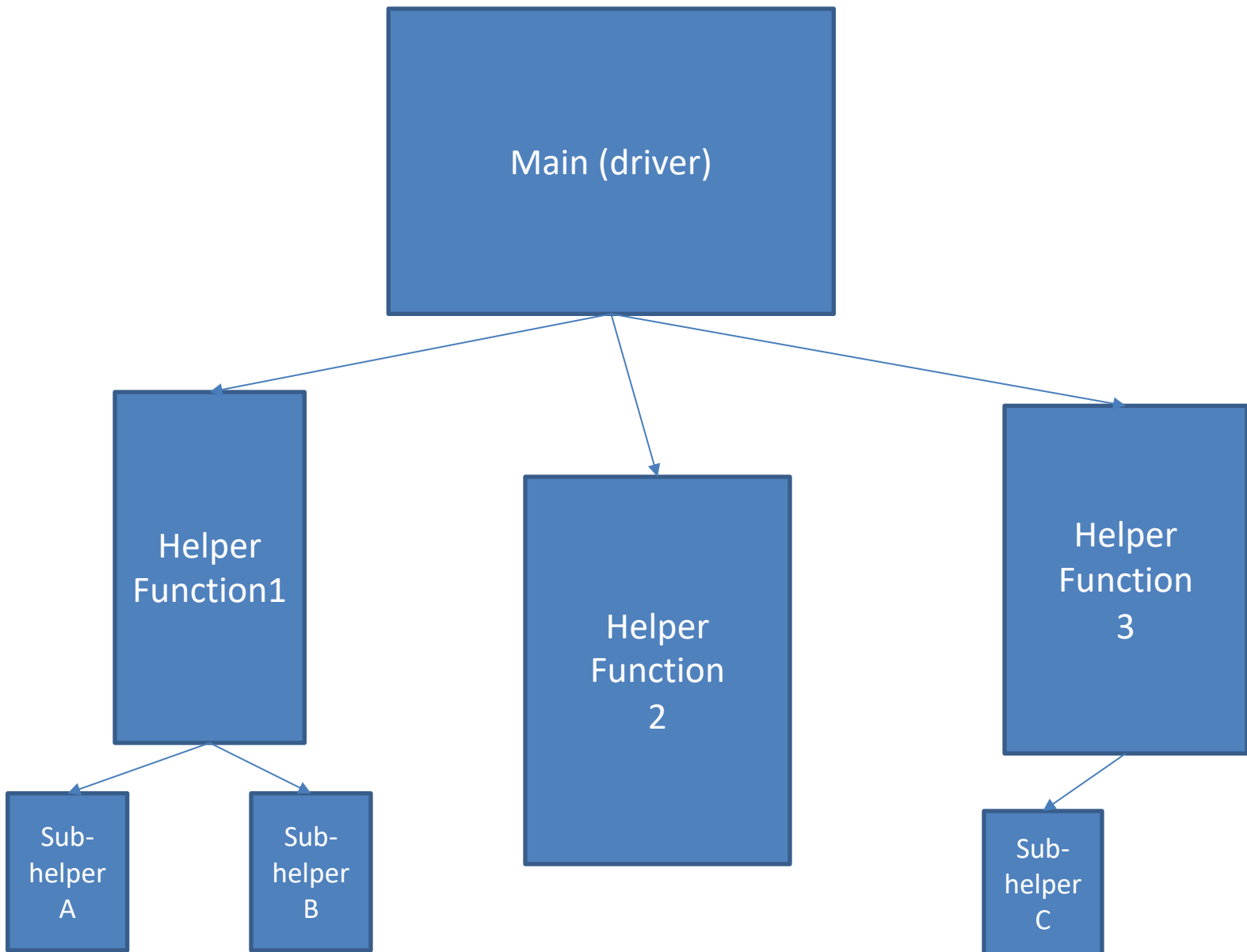
# Python style examples

```
'''
Author      : Allison Gong
Class       : CS260
Date        : 7/30/21
Description : This is the main driver program for the Introduction to Python
portion of Lab 1. This intro consists of a series of coding exercises that cover
some of the core syntax, data structures, and functionality of Python.
'''
```

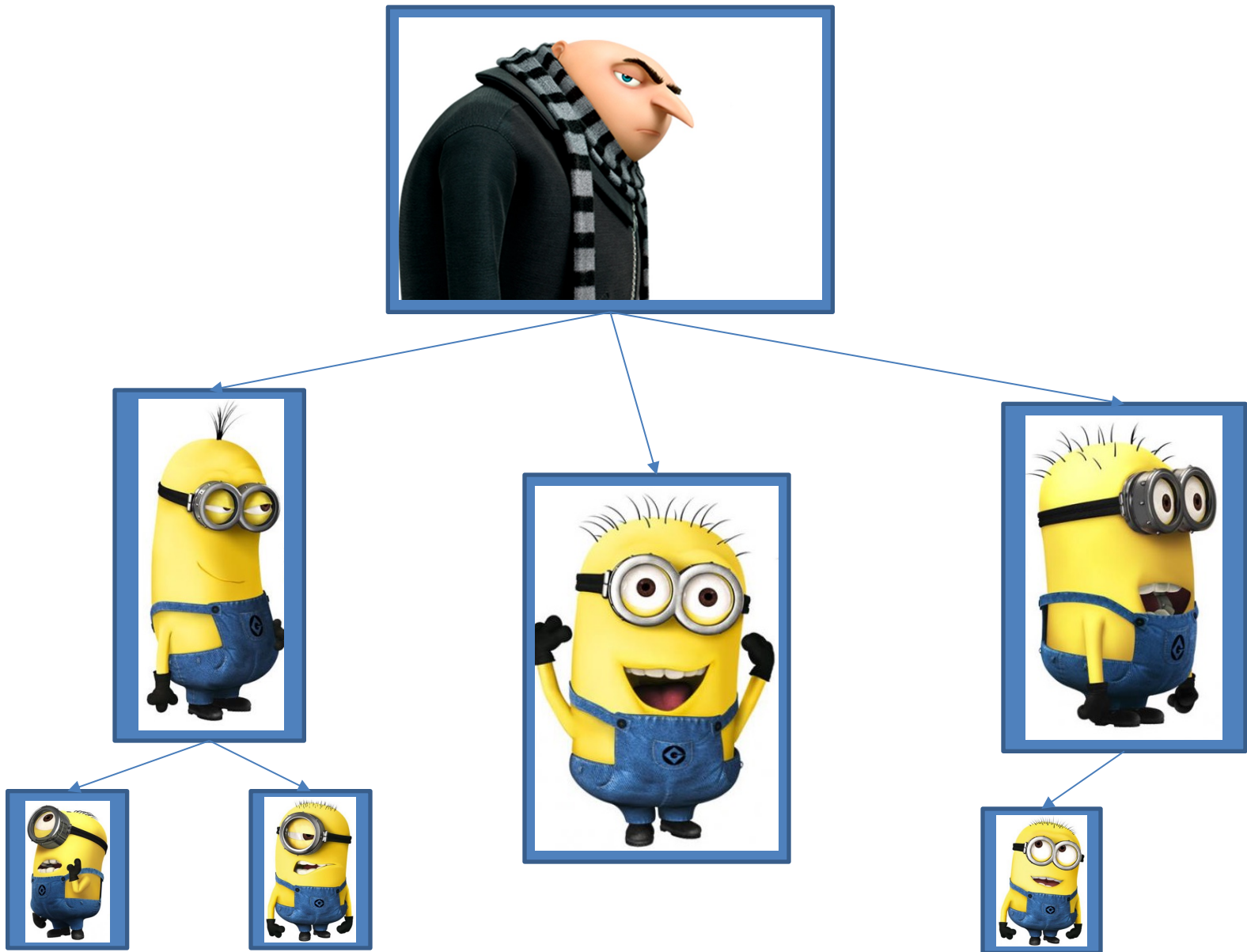
```
# add up 10 random numbers
total = 0
for i in range(10):
    total += random.randrange(6) # includes 0, excludes 6
print("sum:", total)
```

```
def fib(n):
    '''
    Compute and return the nth Fibonacci number.
    n: non-negative integer
    return: nth Fibonacci number
    '''
    # code here
```

# Structure of main and “helper” functions



# Structure of main and “helper” functions



```
"""
```

Given an input phrase and a letter, count how many times that letter appears in the phrase. For example:

```
phrase: creative code
```

```
letter: e
```

```
Number of e's: 3
```

```
Author: Jeff Knerr & Sara Mathieson
```

```
Date: 9/21/18
```

```
"""
```

```
def main():
```

```
    # ask the user for a phrase and a letter
```

```
    phrase = input("phrase: ")
```

```
    letter = input("letter: ")
```

```
    num_chars = len(phrase)
```

```
    # set up accumulator variable count
```

```
    count = 0
```

```
    for i in range(num_chars):
```

```
        # add on 1 each time we see the desired letter
```

```
        if phrase[i] == letter:
```

```
            count = count + 1
```

```
    # example of string formatting (%s for str, %i for int)
```

```
    print("Number of %s's: %i" % (letter, count))
```

```
main()
```

# Reminder: steps of top-down-design (TDD)



# Reminder: steps of top-down-design (TDD)

- 1) Design a **high-level main function** that captures the basic idea of the program.

# Reminder: steps of top-down-design (TDD)

- 1) Design a **high-level main function** that captures the basic idea of the program.
- 2) As you're writing/designing main, think about which details can be **abstracted into small tasks**. Make names for these functions and write their signatures below main.

# Reminder: steps of top-down-design (TDD)

- 1) Design a **high-level main function** that captures the basic idea of the program.
- 2) As you're writing/designing main, think about which details can be **abstracted into small tasks**. Make names for these functions and write their signatures below main.
- 3) **“Stub” out the functions**. This means that they should work and return the correct type so that your code runs, but they don't do the correct task yet. For example, if a function should return a list, you can return []. Or if it returns a boolean, you can return False.

# Reminder: steps of top-down-design (TDD)

- 1) Design a **high-level main function** that captures the basic idea of the program.
- 2) As you're writing/designing main, think about which details can be **abstracted into small tasks**. Make names for these functions and write their signatures below main.
- 3) **“Stub” out the functions**. This means that they should work and return the correct type so that your code runs, but they don't do the correct task yet. For example, if a function should return a list, you can return []. Or if it returns a boolean, you can return False.
- 4) Iterate on your design until you have a working main and stubbed out functions. Then start **implementing** the functions, starting from the “bottom up”.

# Outline for Jan 30

- Python style and implementation notes
- **Overfitting**
- K-nearest neighbors
- KD Trees

# Overfitting with a high-degree polynomial

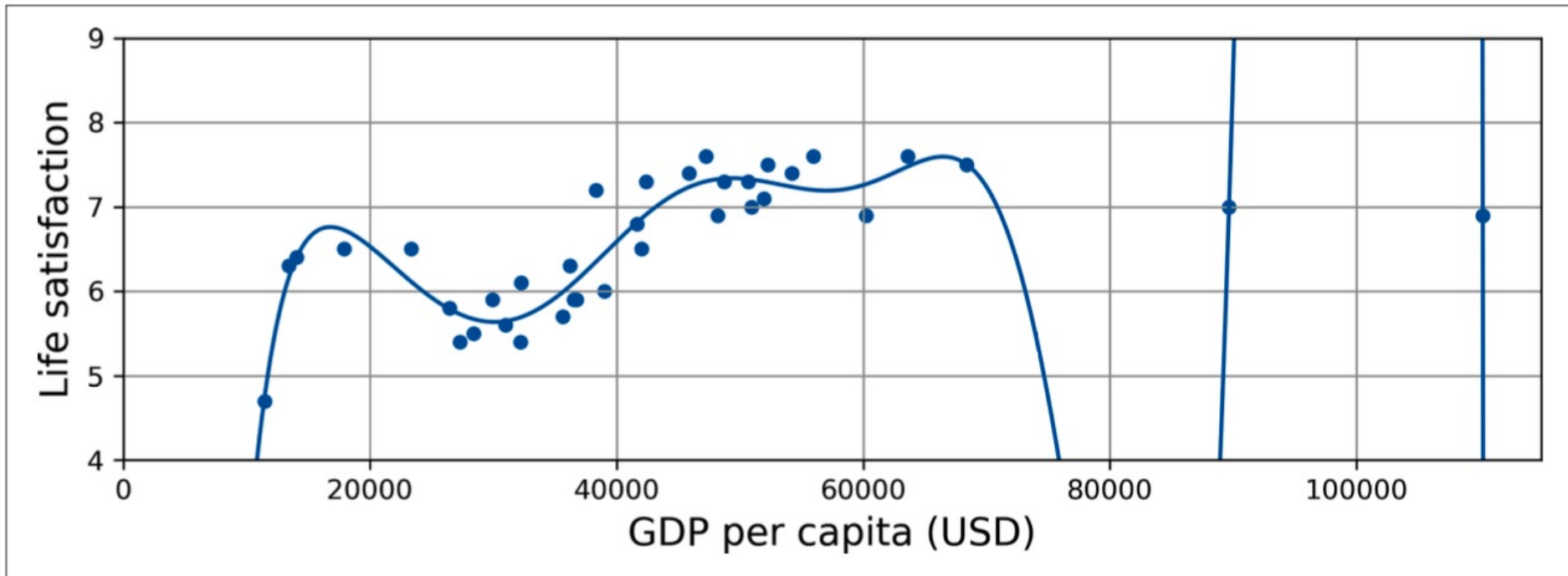


Figure 1-23. Overfitting the training data

# Overfitting

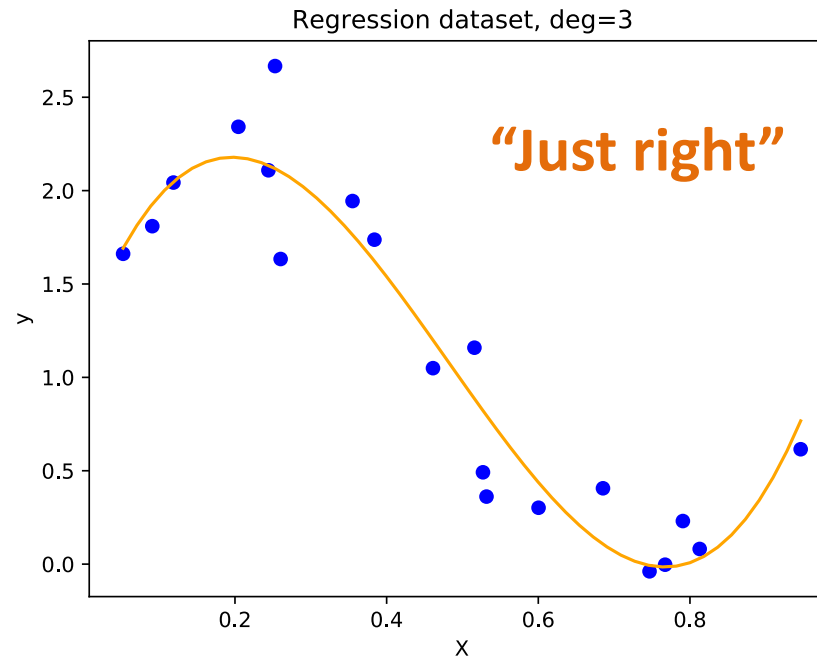
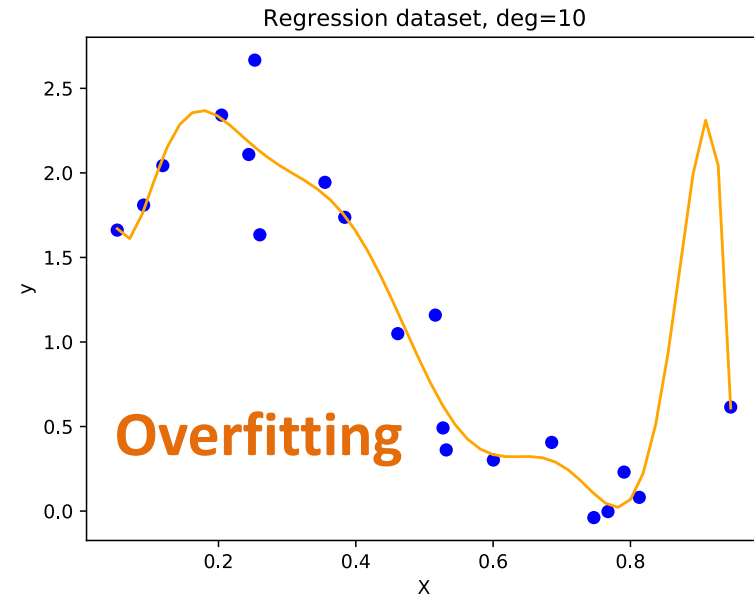
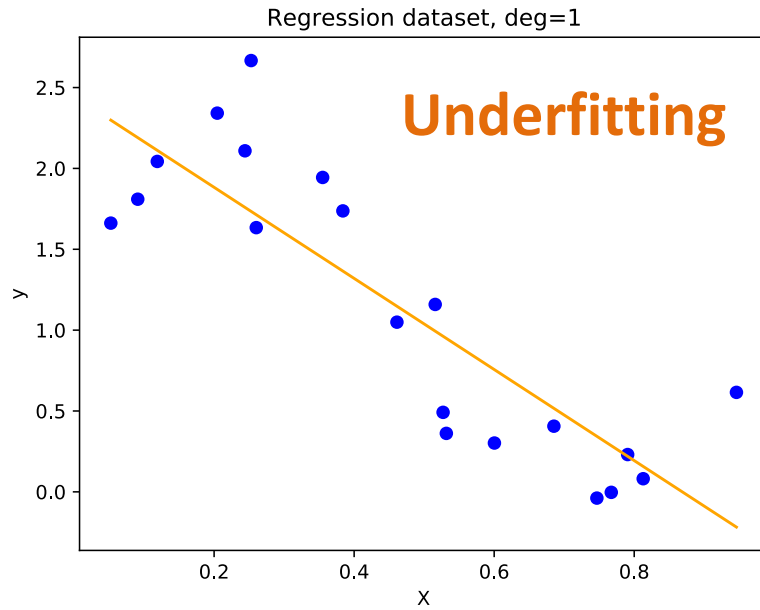
- “**Overfitting** happens when the model is too complex relative to the amount and noisiness of the training data.” (Geron Chap 1)
- **Solutions**
  - Reduce the complexity of the model
  - Get more training data
  - Reduce noise in the training data

# Terminology

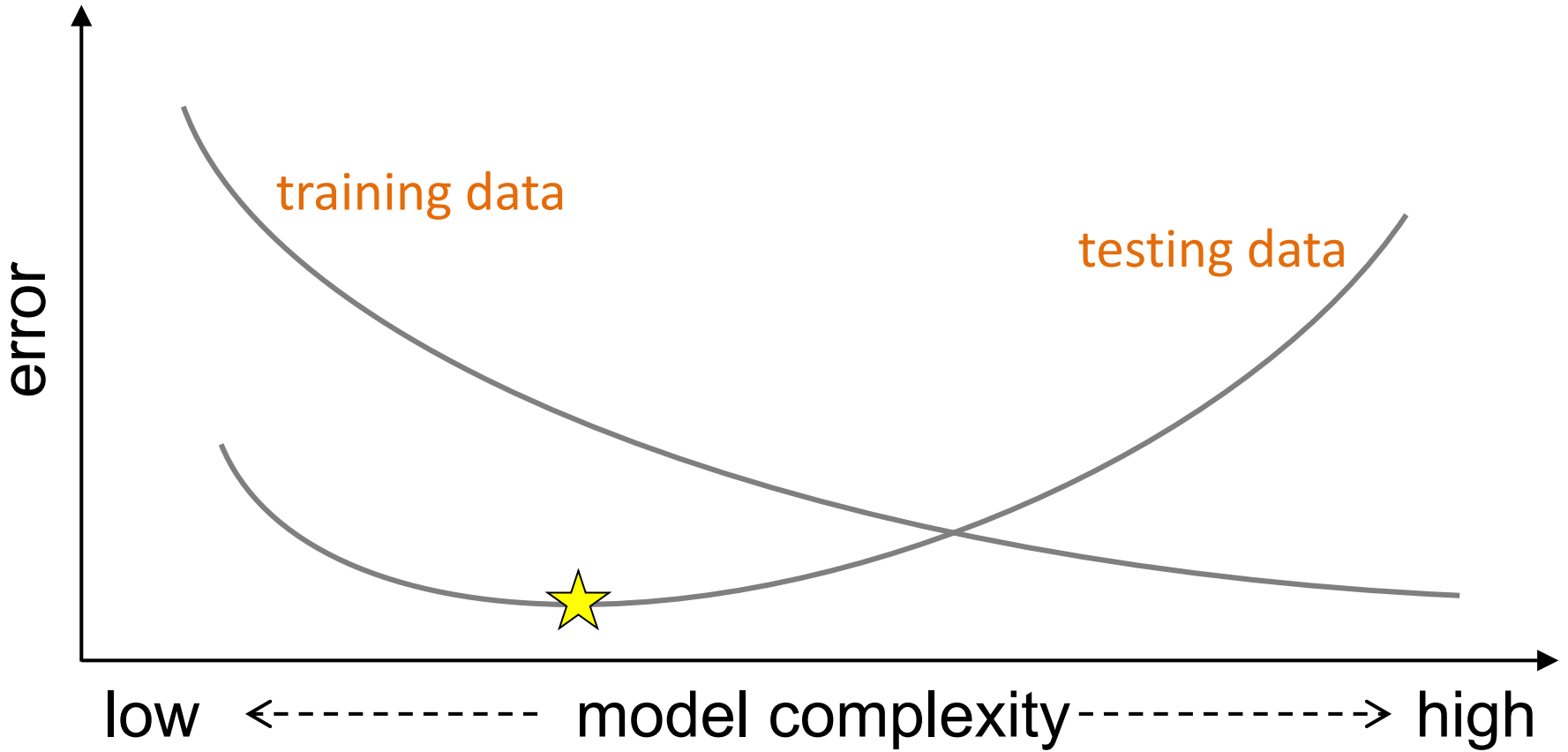
- *Underfitting*: “had the opportunity to learn something but didn’t” (Duame)
- *Overfitting*: memorized individual training examples (fit to noise) and can’t generalize



# Under and over-fitting (CS260 example)



# Common pattern observed in machine learning



# Validation data

- Is it wrong to use the test data to determine the model complexity?
- **Yes!**

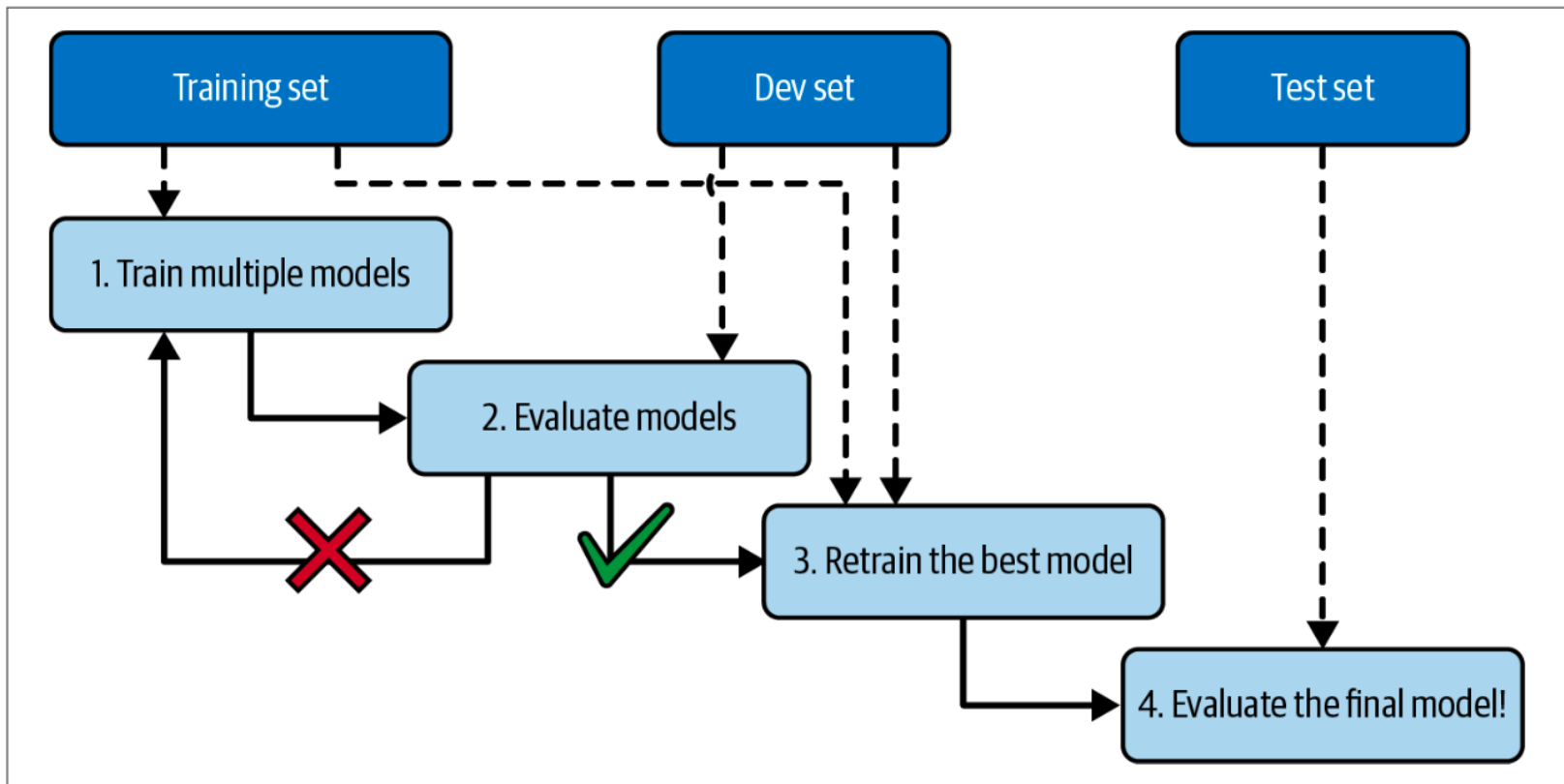


Figure 1-25. Model selection using holdout validation

# Outline for Jan 30

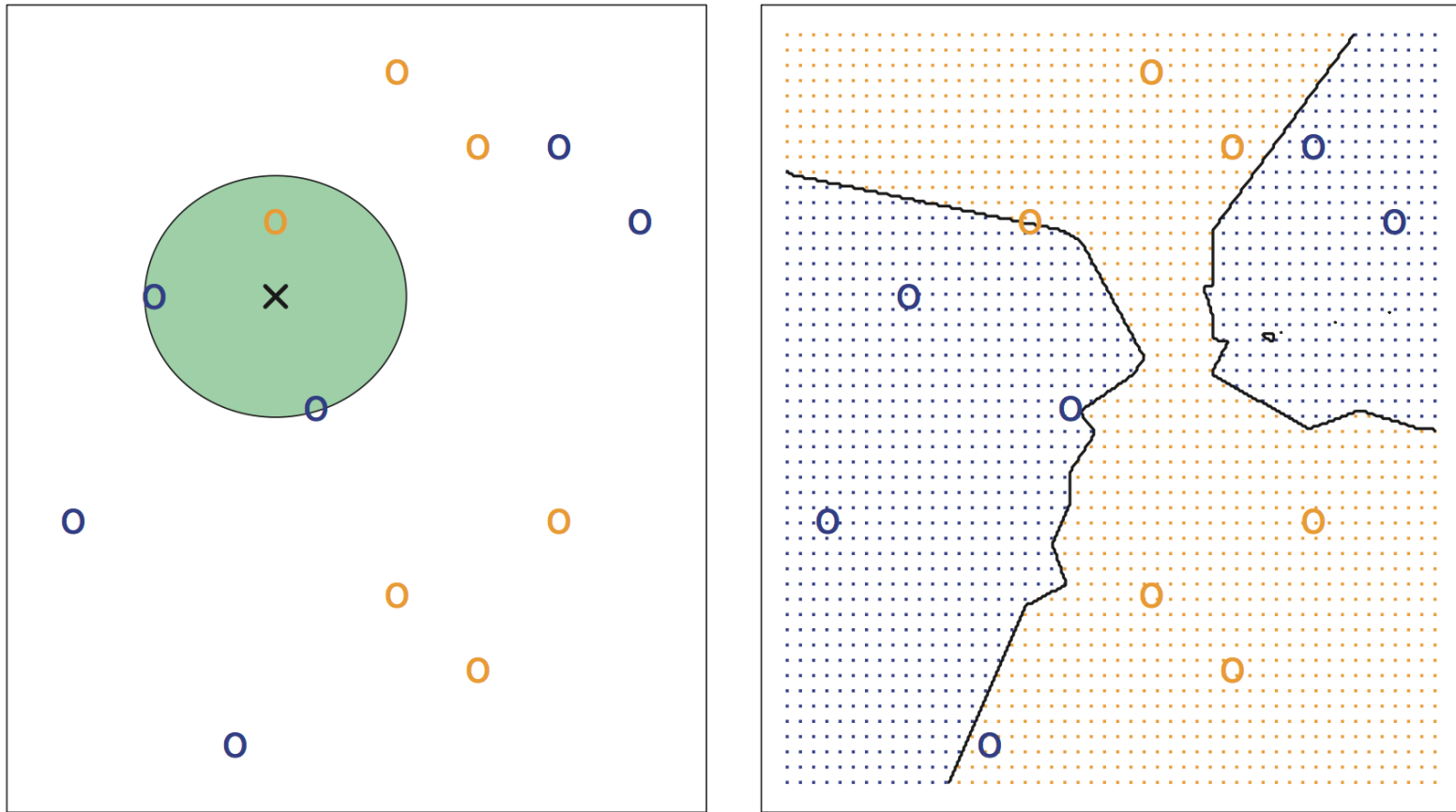
- Python style and implementation notes
- Overfitting
- **K-nearest neighbors**
- KD Trees

# Nearest Neighbors

Why would we be interested in finding a point's nearest neighbor in a set of points?

- Fill in missing data
- Prediction unknown values (labels, output, etc)
- subroutine for some clustering methods

# K-nearest neighbors creates implicit decision boundaries



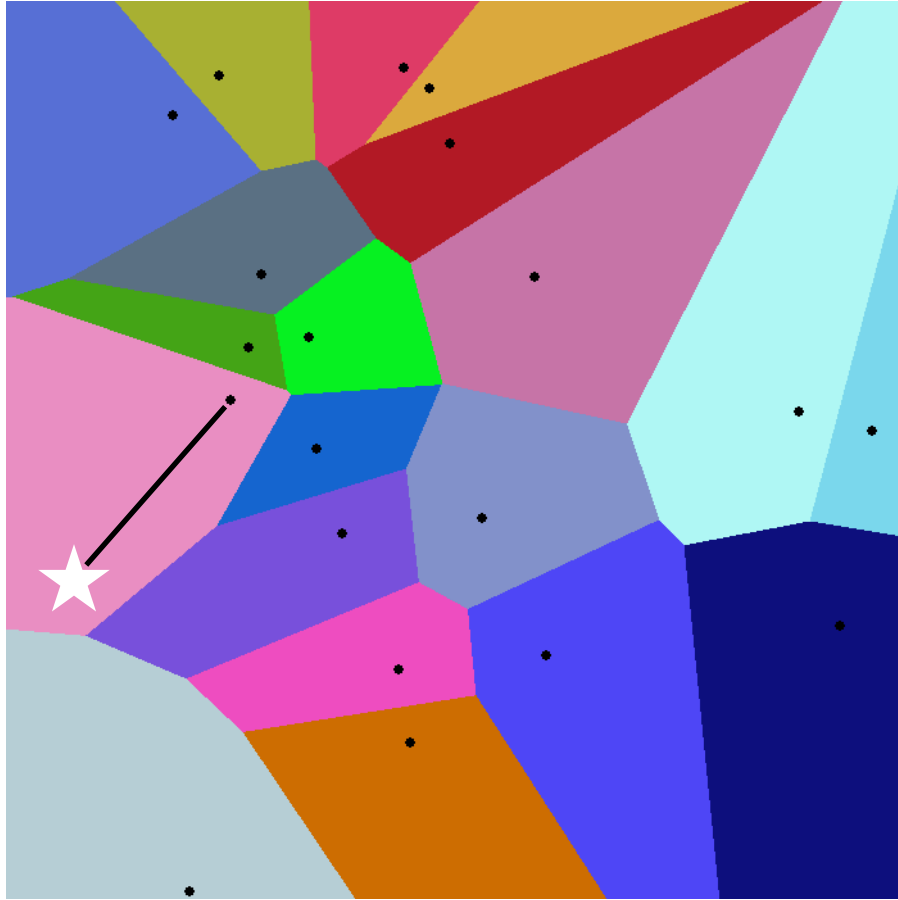
*Decision boundary:* separates regions of the feature space that would be classified as positive or negative (or multiclass)

Figure 2.14 from ISL book, KNN with two classes ( $C=2$ ), and  $K=3$

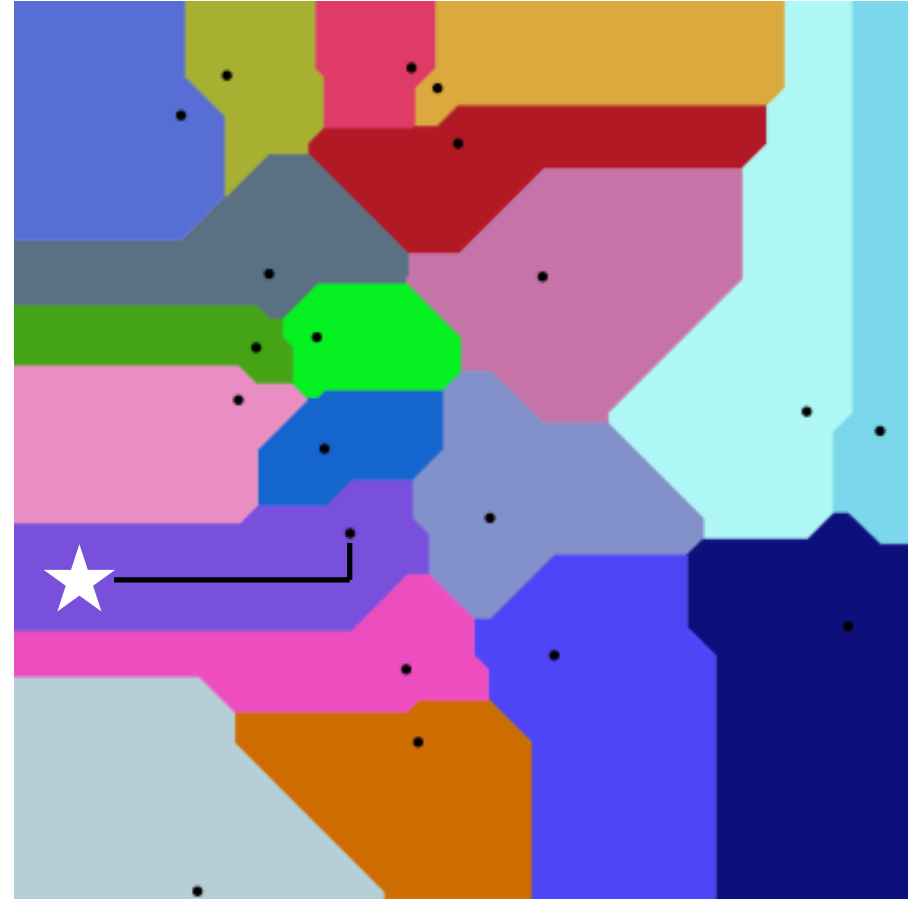
# Voronoi Diagrams

Nearest neighbor queries in 2D

☆ query point



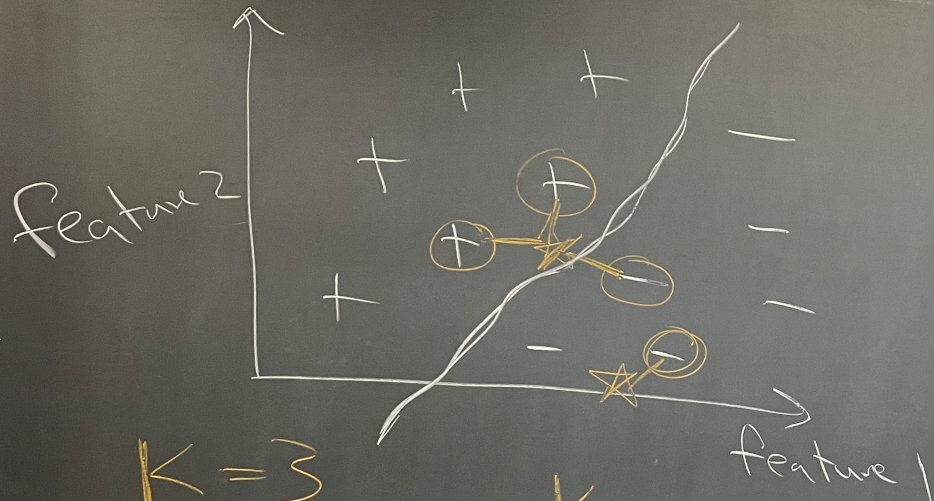
Euclidean distance



Manhattan distance



# KNN for binary classification



$K=3$   
label =  $\oplus$

$K=1$   
label =>  $\ominus$

Naive KNN  
 $knn(\text{points}, q)$  — training data, query

for  $p$  in points:  
if  $d(p, q) < \min$   
 $\min = d(p, q)$   
 $nn = p$



① runtime if  $\text{len}(\text{points}) = n$   
 $\Rightarrow O(n)$  ignoring # dims

② what about higher  $K$ ?

---

Euclidean dist  $a = (a_1, a_2, \dots, a_d)$

$b = (b_1, b_2, \dots, b_d)$

$$d_e(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots}$$

$$d_m(a, b) = |a_1 - b_1| + |a_2 - b_2| + \dots$$

↙ Manhattan



## Comparison of decision boundaries

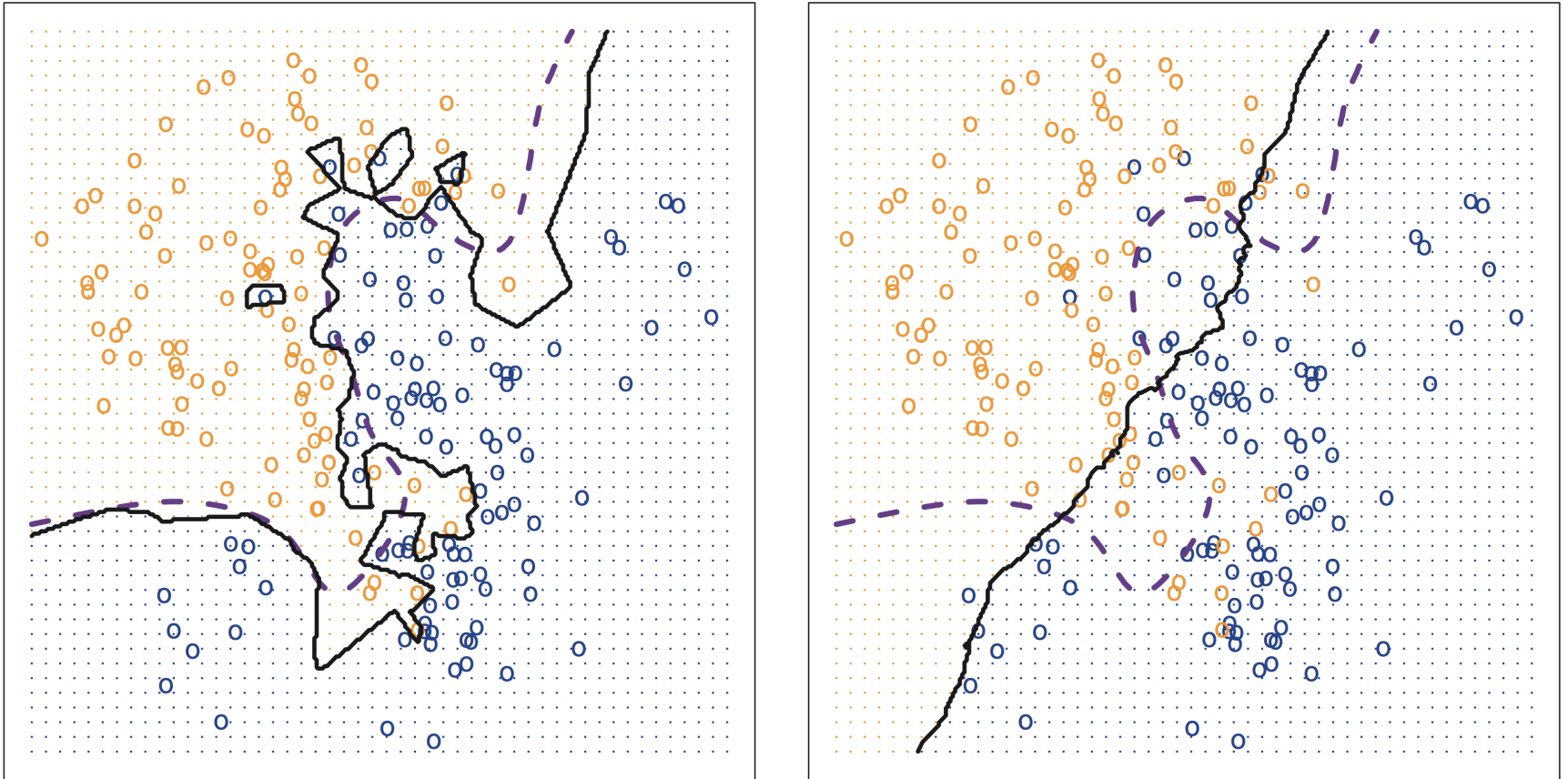
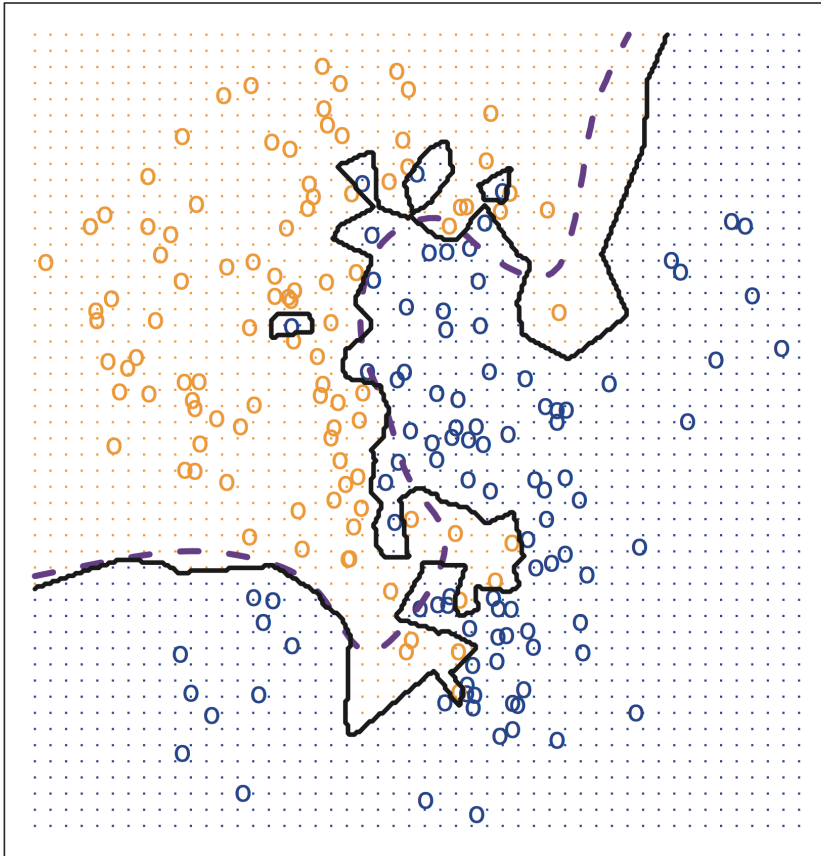


Figure 2.16 from ISL book (dashed line is “ideal” boundary)

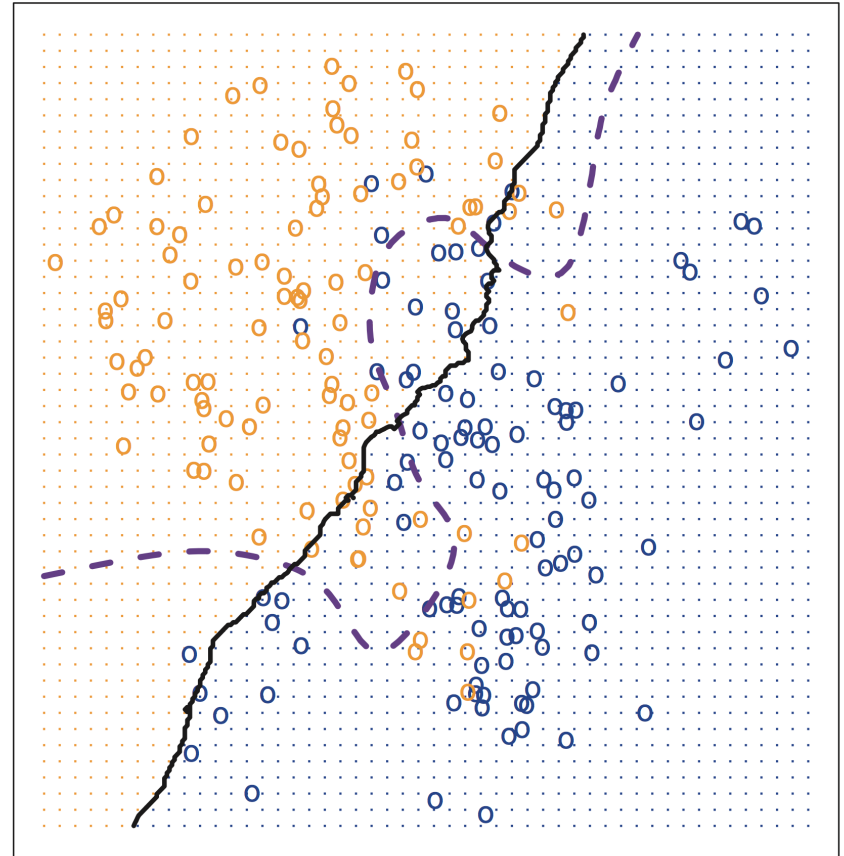
# Comparison of decision boundaries

KNN:  $K=1$



Overfitting

KNN:  $K=100$



Underfitting

Figure 2.16 from ISL book (dashed line is "ideal" boundary)

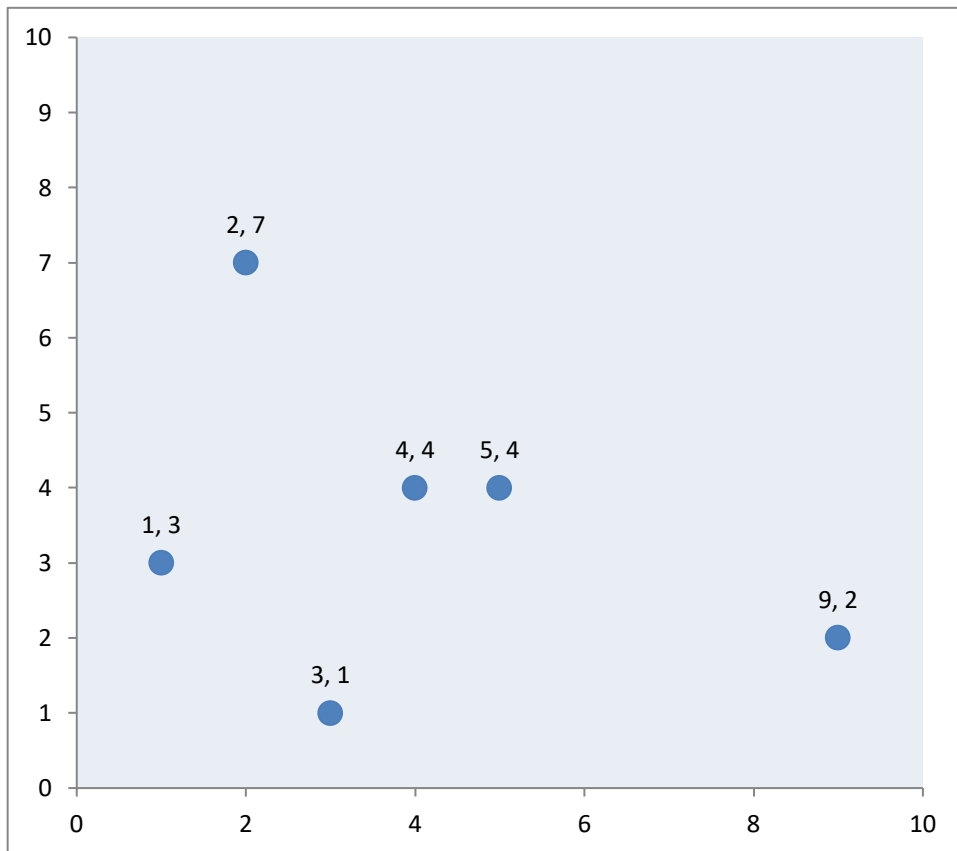
# Calculating the nearest neighbor

- What is the "naïve" approach?
- How long does it take to find the nearest neighbor of a point? In 2D? In  $d$ -dimensions?
- How could we do better?

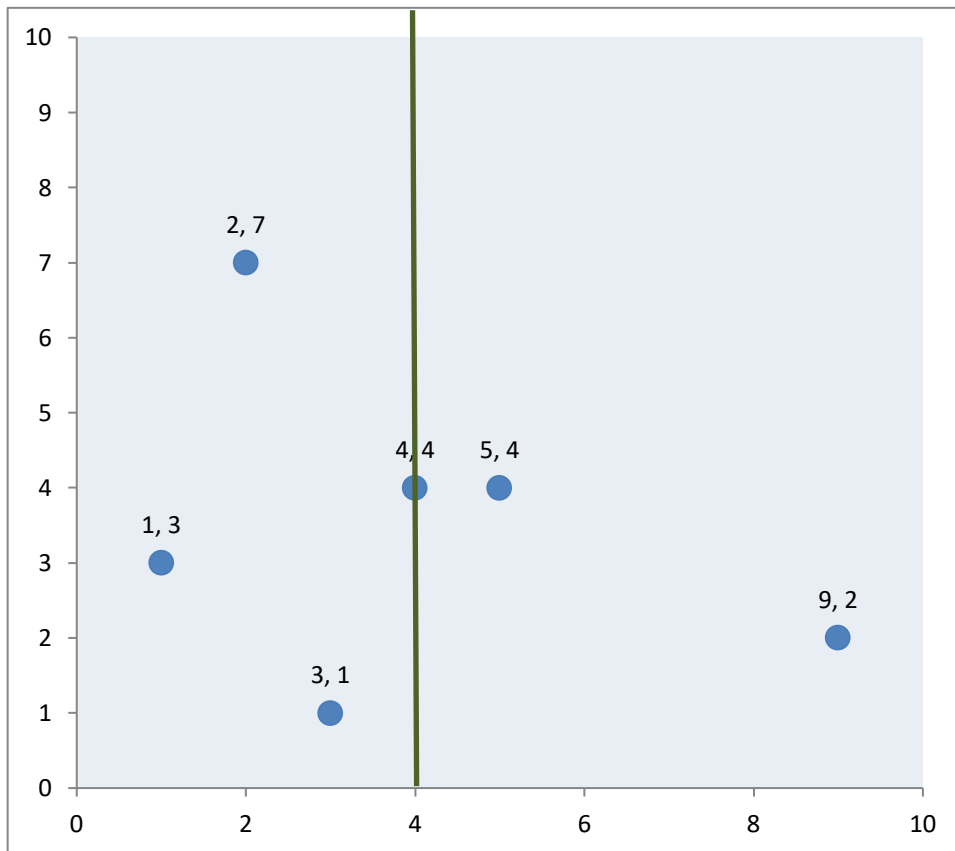
# Outline for Jan 30

- Python style and implementation notes
- Overfitting
- K-nearest neighbors
- **KD Trees**

# Making a kd-tree

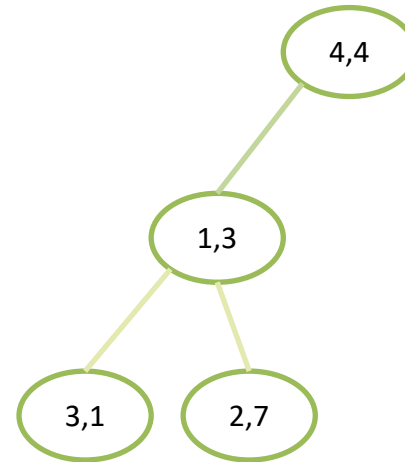
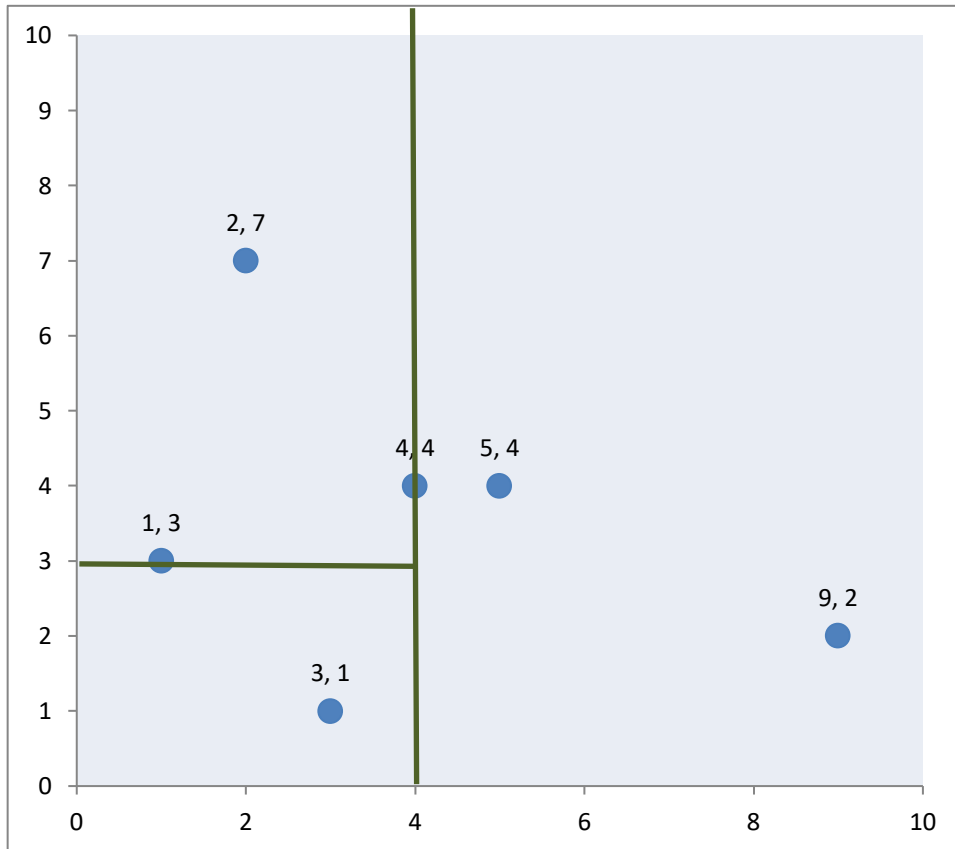


# Making a kd-tree



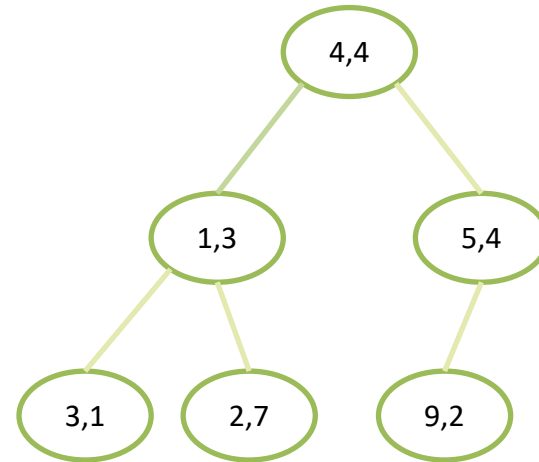
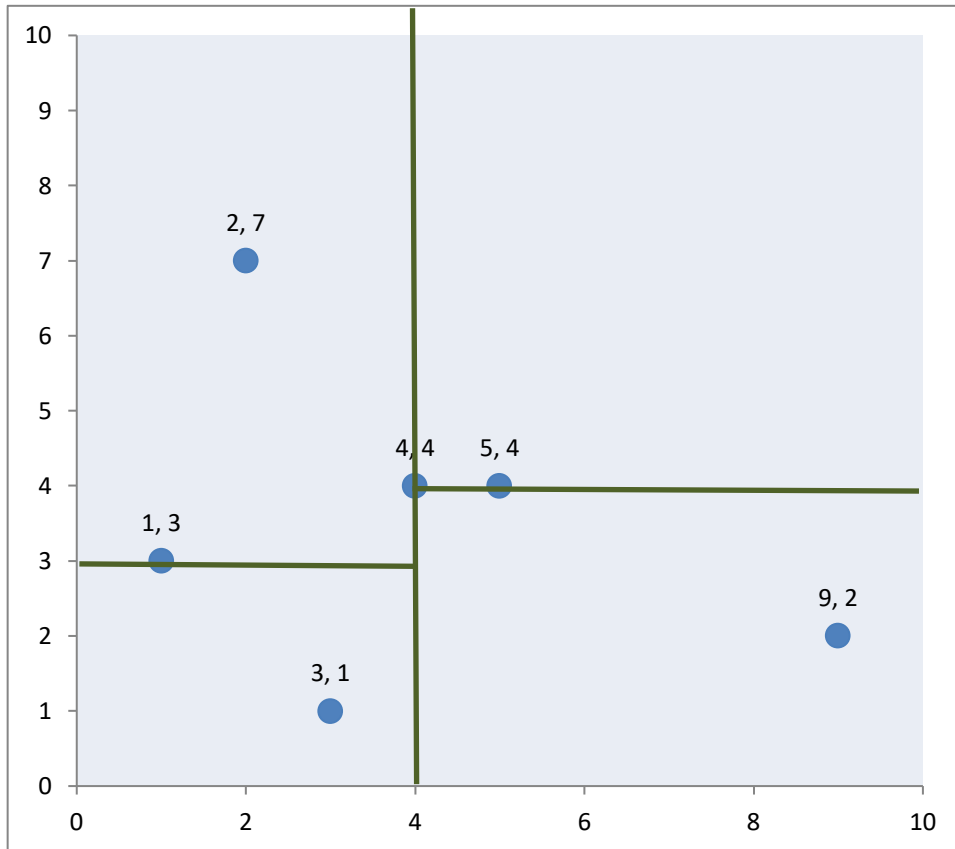
4,4

# Making a kd-tree





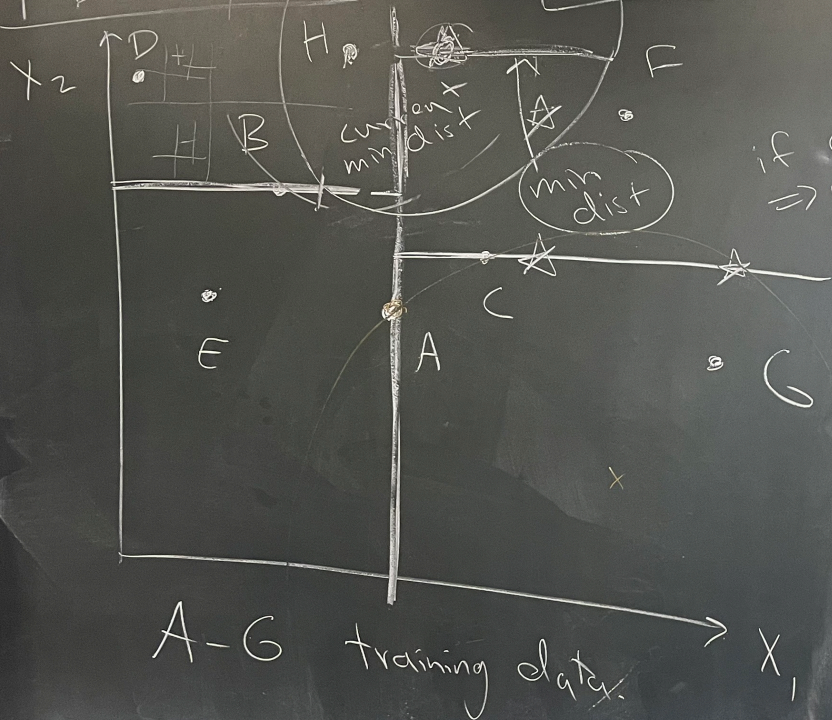
# Making a kd-tree





# KD-trees

$d=2$



if even  
=> more on  
the left

$kdtree(\text{points}, \text{depth})$   
 $\text{dim} = \text{depth} \% d$

Sort points & take median  
 make node at median

$\text{node.point} = \text{median}$

$\text{node.left} = kdtree(\text{points on the left}, \text{depth} + 1)$

$\text{node.right} = kdtree(\text{points on the right}, \text{depth} + 1)$

return ?



$\boxed{\text{dim}=0}$

Kdtree()

D

E

B

$\textcircled{A}$

C

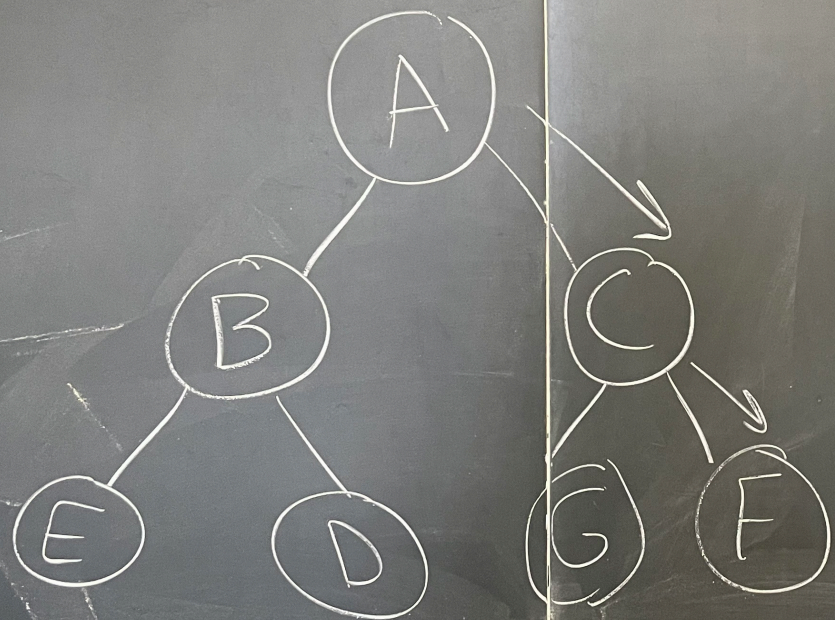
F

G

(Sort in  $\text{dim}=0$ )

Split

Kdtree



$\boxed{\text{dim}=1}$

E  $\textcircled{B}$  D

$\boxed{\text{dim}=1}$

G  $\textcircled{C}$  F