

CSC 580 Principles of Machine Learning

# 03 Nearest Neighbors

**Chicheng Zhang**

**Department of Computer Science**



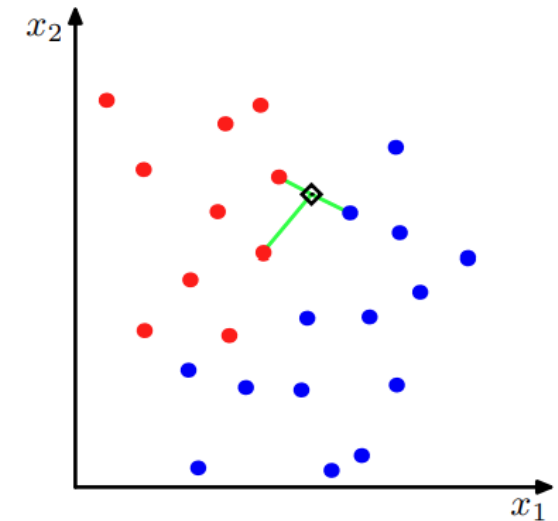
\*slides credit: built upon CSC 580 Fall 2021 lecture slides by Kwang-Sung Jun

# Motivation

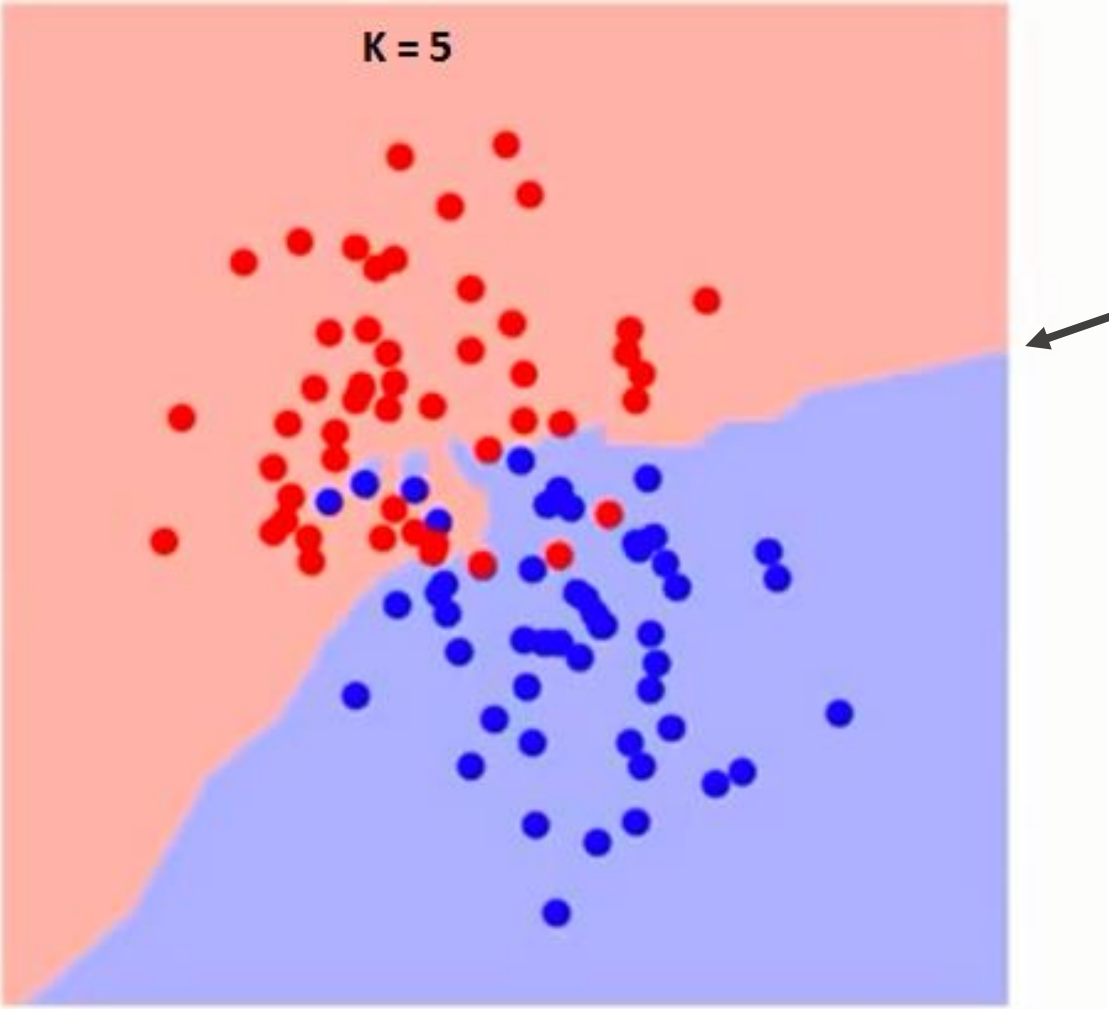
- Example: given student course survey data, predict whether Alice likes Algorithms course
- A natural idea: find a student “similar” to Alice & has taken Algorithm course before, say Jeremy
  - If Jeremy likes Algorithms, then Alice is also likely to have the same preference.
  - Or even better, find *several* similar students

# $k$ -nearest neighbors ( $k$ -NN): main concept

- Training set:  $S = \{ (x_1, y_1 ), \dots, (x_m, y_m ) \}$
- **Inductive bias**: given test example  $x$ , its label should resemble the labels of **nearby points**
- Function
  - input:  $x$
  - find the  $k$  nearest points to  $x$  from  $S$ ; call their indices  $N(x)$
  - output: the majority vote of  $\{y_i: i \in N(x)\}$ 
    - For regression, the average.



# k-NN classification example



← decision boundary

# Basics

- Oftentimes convenient to work with feature  $x \in \mathbb{R}^d$

- Distances in  $\mathbb{R}^d$ :

notation  $x(f): x = (x(1), \dots, x(d))$

- (popular) Euclidean distance  $d_2(x, x') = \sqrt{\sum_{f=1}^d (x(f) - x'(f))^2}$

- Manhattan distance  $d_1(x, x') = \sum_{f=1}^d |x(f) - x'(f)|$

- If we shift a feature, would the distance change?

- What about scaling a feature?

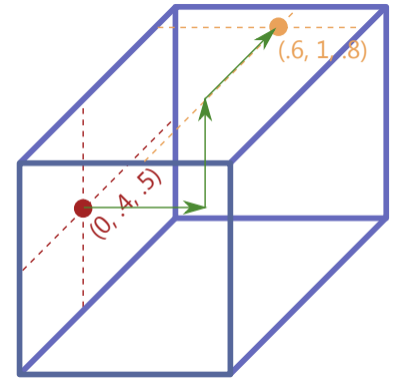
- How to extract features as **real values**?

- Boolean features:  $\{Y, N\} \rightarrow \{0, 1\}$

- Categorical features:  $\{\text{Red, Blue, Green, Black}\}$

- Convert to  $\{1, 2, 3, 4\}$ ?

- Better one-hot encoding:  $(1, 0, 0, 0), \dots, (0, 0, 0, 1)$  (IsRed?/isGreen?/isBlue?/IsBlack?)



# $k$ -NN classification: pseudocode

- Training is trivial: store the training set
- Test:

---

**Algorithm 3** KNN-PREDICT( $\mathbf{D}, K, \hat{x}$ )

---

list	→	1: $S \leftarrow []$	
		2: <b>for</b> $n = 1$ <b>to</b> $N$ <b>do</b>	
append to list	→	3: $S \leftarrow S \oplus \langle d(x_n, \hat{x}), n \rangle$	// store distance to training example $n$
		4: <b>end for</b>	
sort in first coordinate	→	5: $S \leftarrow \text{SORT}(S)$	// put lowest-distance objects first
		6: $\hat{y} \leftarrow 0$	
		7: <b>for</b> $k = 1$ <b>to</b> $K$ <b>do</b>	
		8: $\langle \text{dist}, n \rangle \leftarrow S_k$	// $n$ this is the $k$ th closest data point
		9: $\hat{y} \leftarrow \hat{y} + y_n$	// vote according to the label for the $n$ th training point
		10: <b>end for</b>	
Majority vote of $\{y_i : i \in N(x)\}$	→	11: <b>return</b> SIGN( $\hat{y}$ )	// return +1 if $\hat{y} > 0$ and -1 if $\hat{y} < 0$

---

- Time complexity (assuming distance calculation takes  $O(d)$  time)
  - $O(m d + m \log m + k) = O(m(d + \log m))$
- Faster nearest neighbor search: k-d trees, locality sensitive hashing

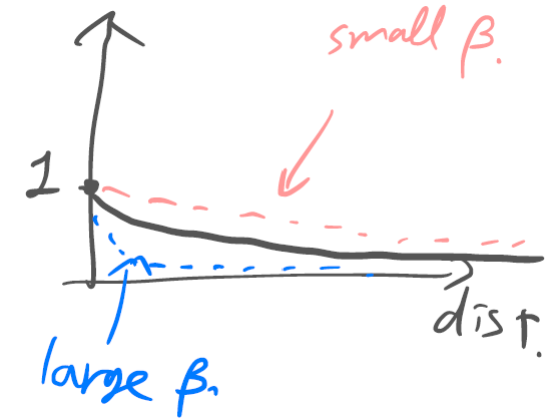
# Variations

- Classification

- Recall the majority vote rule:  $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i \in N(x)} 1\{y_i = y\}$

- Soft weighting nearest neighbors:  $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i=1}^m w_i 1\{y_i = y\},$

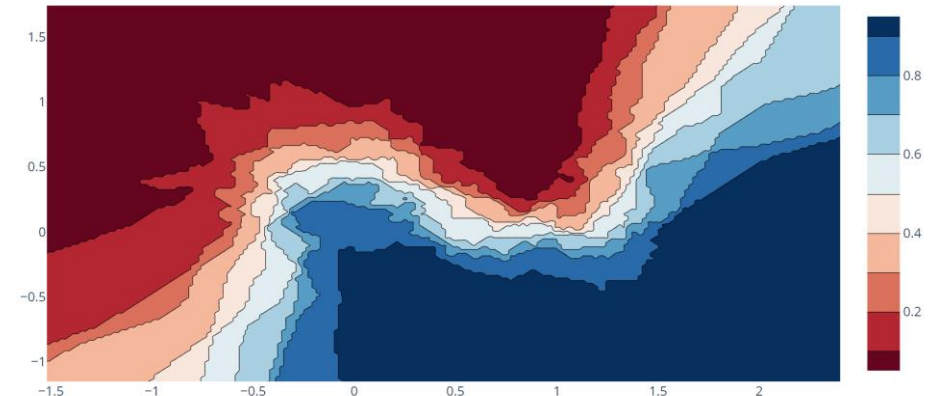
where  $w_i \propto \exp(-\beta d(x, x_i))$ , or  $\propto \frac{1}{1+d(x, x_i)^\beta}$



- Class probability estimates

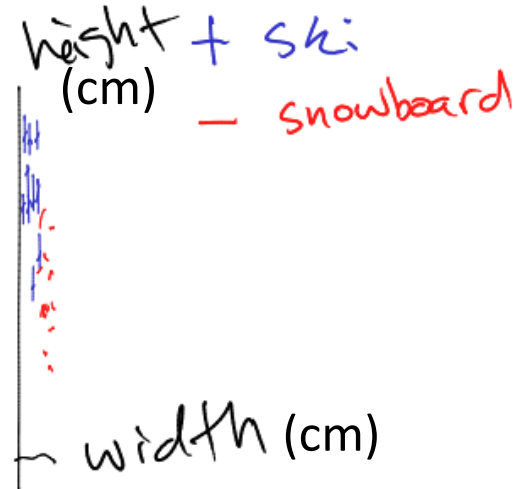
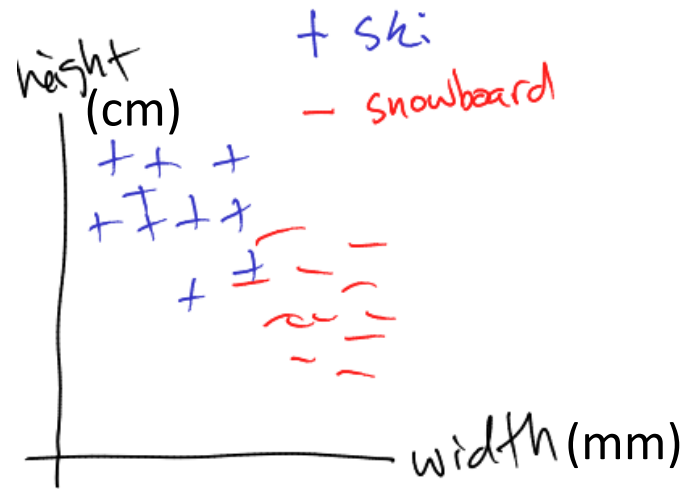
- $\hat{P}(Y = y | x) = \frac{1}{k} \sum_{i \in N(x)} 1\{y_i = y\}$

- Useful for “classification with rejection”



# Feature issue 1: scaling

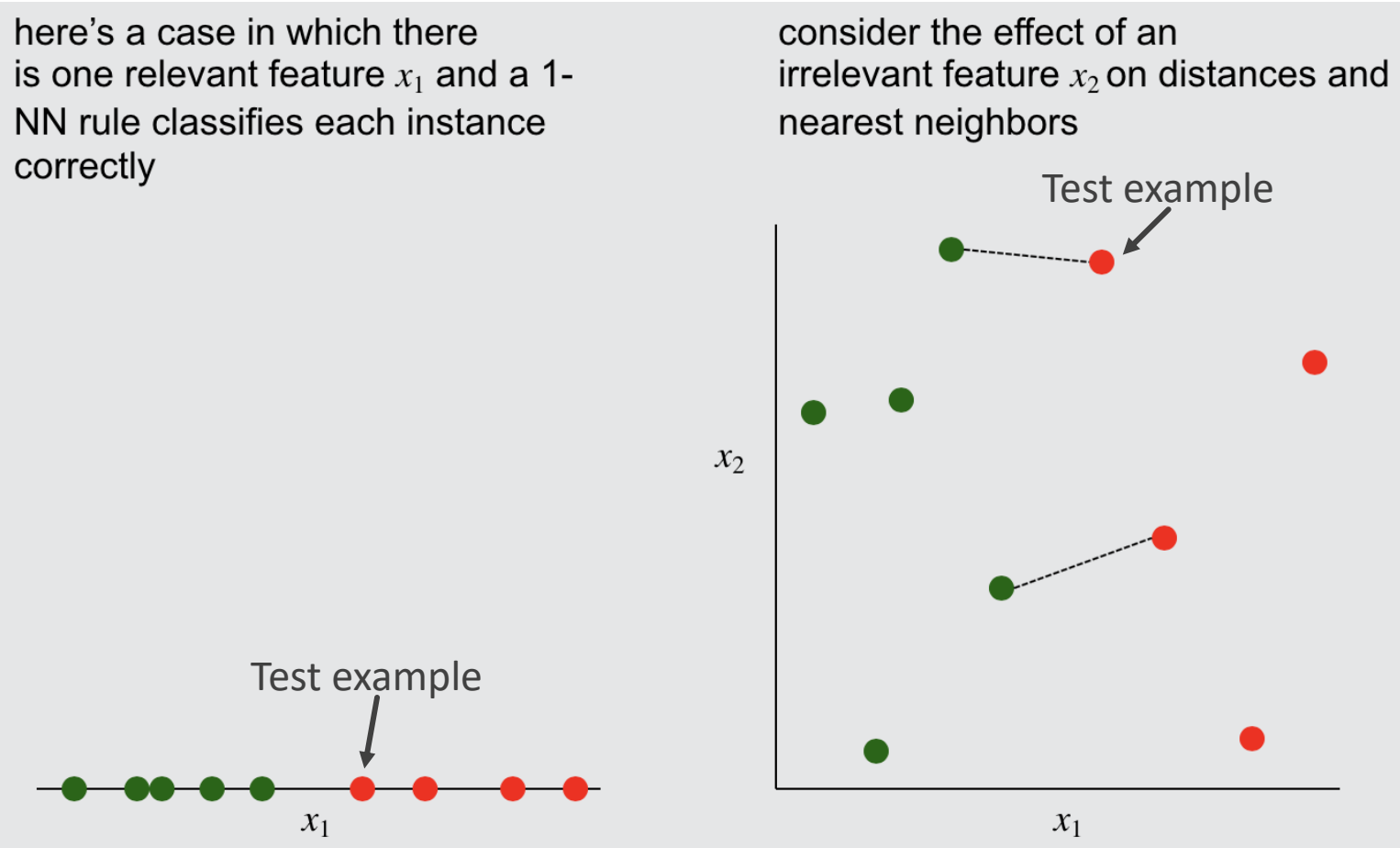
- Features having different scale can be problematic.
- Ex: ski vs. snowboard classification



- Solution: feature standardization (later in the course)



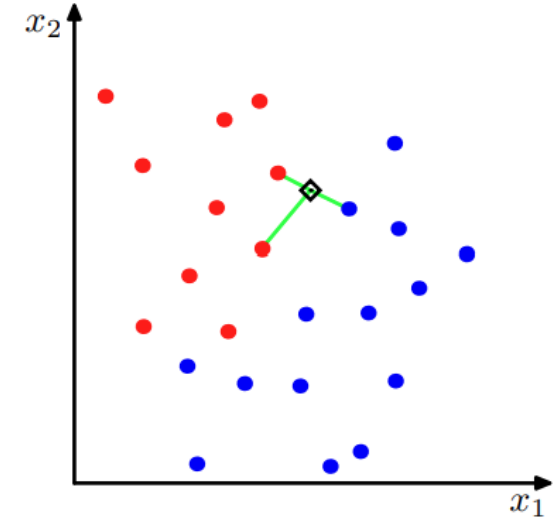
# Feature issue 2: irrelevant features



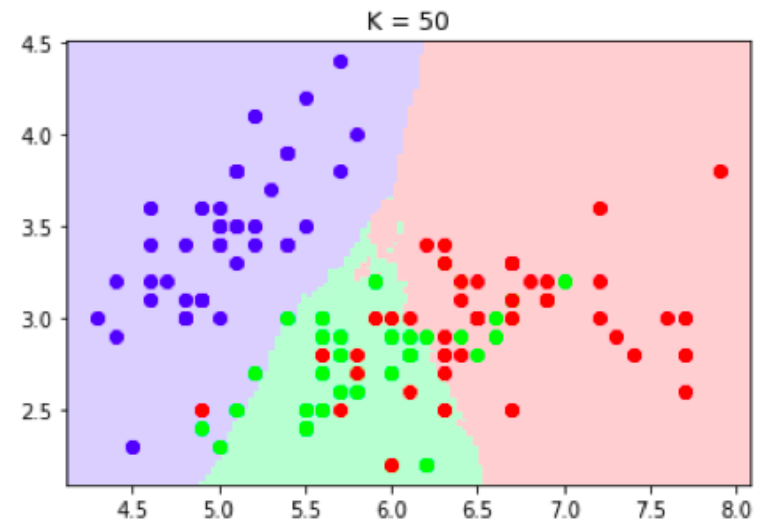
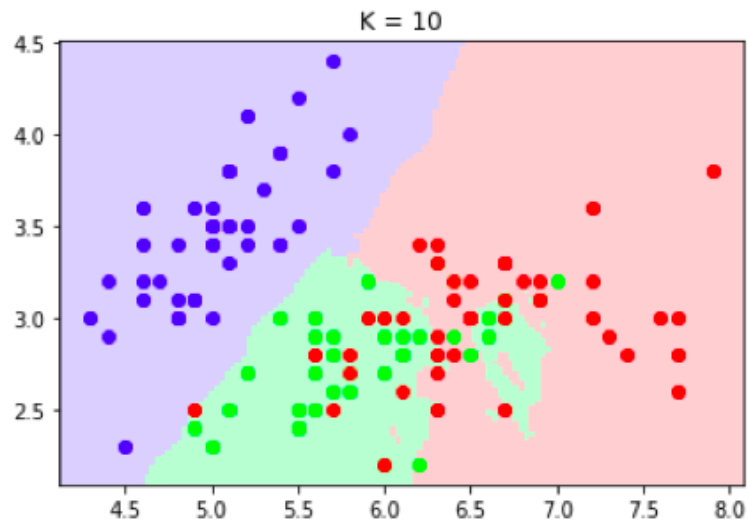
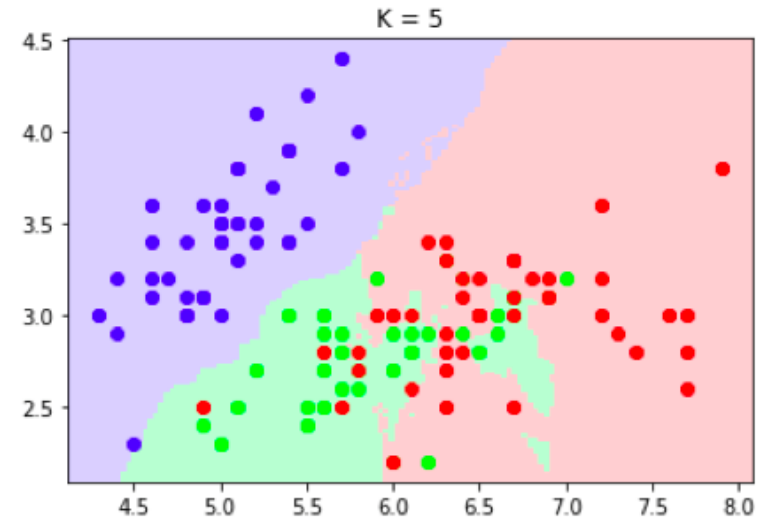
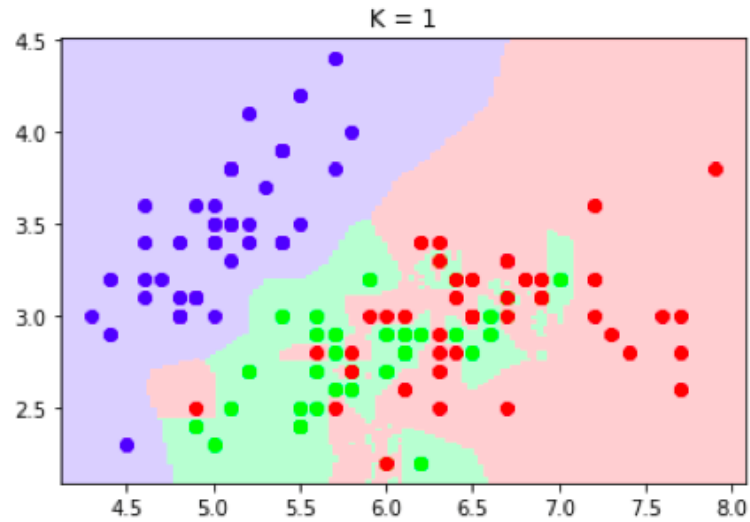
- Recall: how did we deal with these in decision trees?
- Solution: feature selection (later in the course)

# Hyperparameter tuning in $k$ -NN

- Hyperparameter:  $k$
- $k = 1$ :
  - Training error = 0, overfitting
- $k = N$ :
  - Output a constant (majority class) prediction, underfitting
- Can use hold-out validation to choose  $k$



# Hyperparameter tuning in $k$ -NN



# Comparison (feature $x \in \mathbb{R}^d$ )

- Interpretability
- Sensitivity to irrelevant features
- training time
- test time per example

## Decision Tree

## $k$ -NN

High

Medium (example-based)

Low

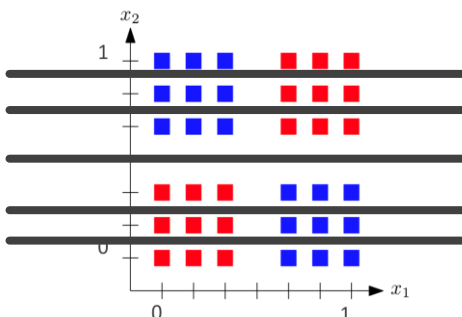
High

$O(\#nodes \cdot d \cdot (m + m \log m))$   
 $\leq \tilde{O}(d m^2)$  (when no two points have the same feature)

0

$O(\text{depth})$

$O(m(d + \log m))$



# Next lecture (9/7)

- Linear classification; the Perceptron algorithm
- Assigned reading: CIML Chap. 4