

GEOMETRIC VIEW OF DATA

David Kauchak
CS 451 – Fall 2013

Admin

Assignment 2 out

Assignment 1

Keep reading

Office hours today from: 2:35-3:20

Videos?

Proper Experimentation



u13007351 fotosearch.com

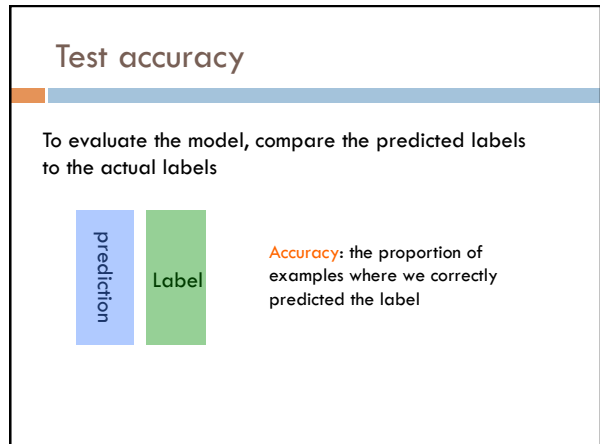
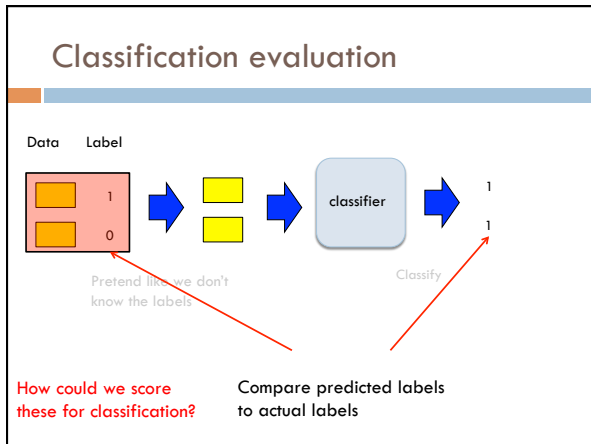
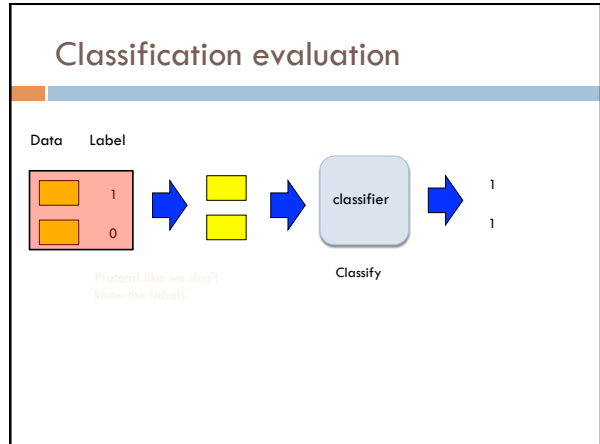
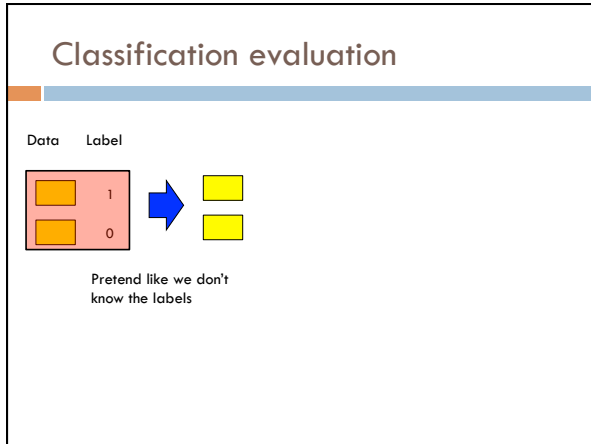
Experimental setup

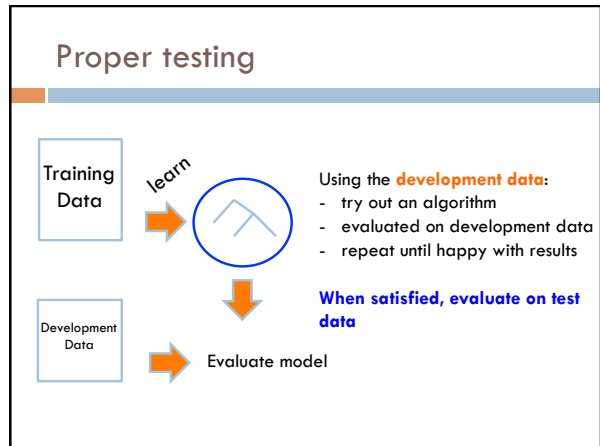
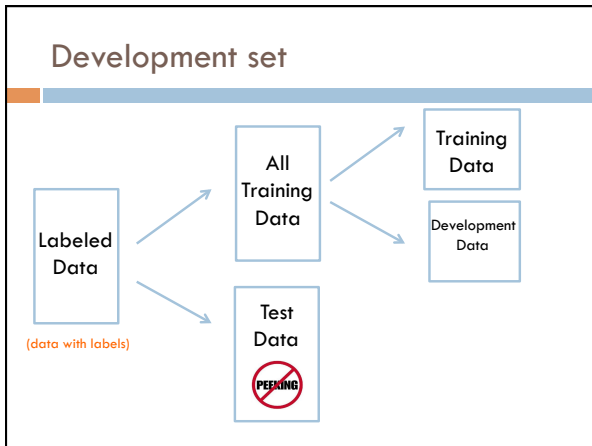
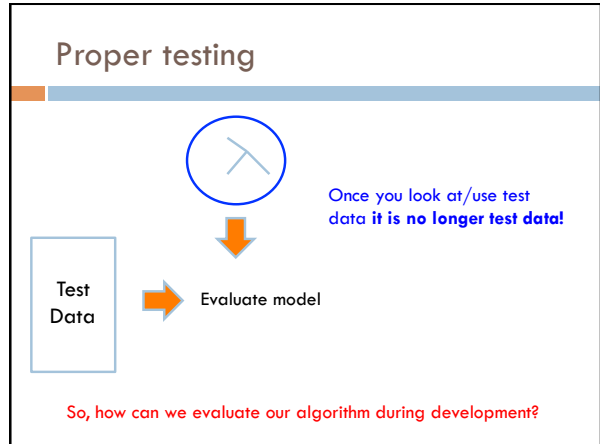
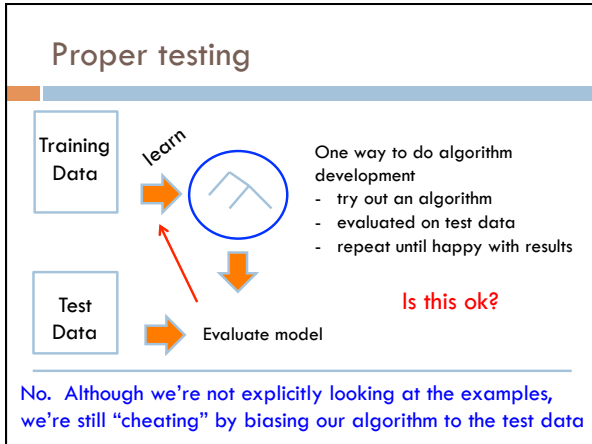
REAL WORLD USE OF ML ALGORITHMS

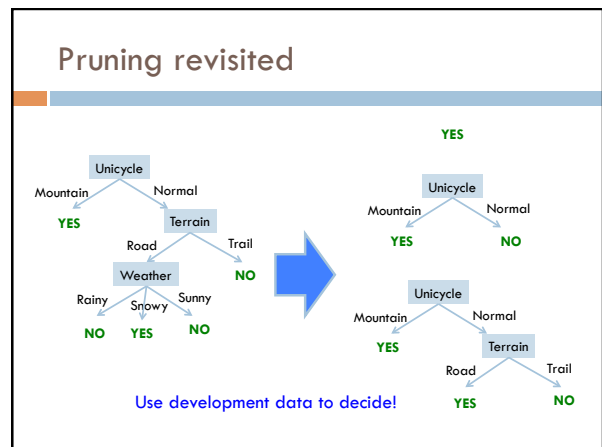
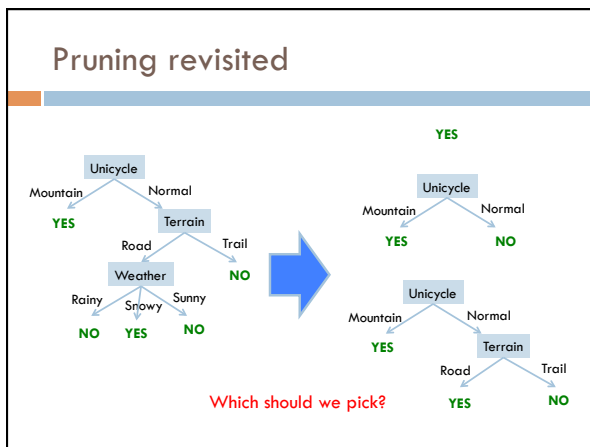
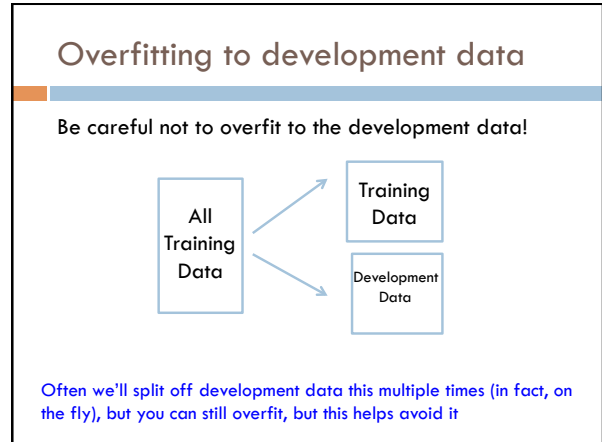
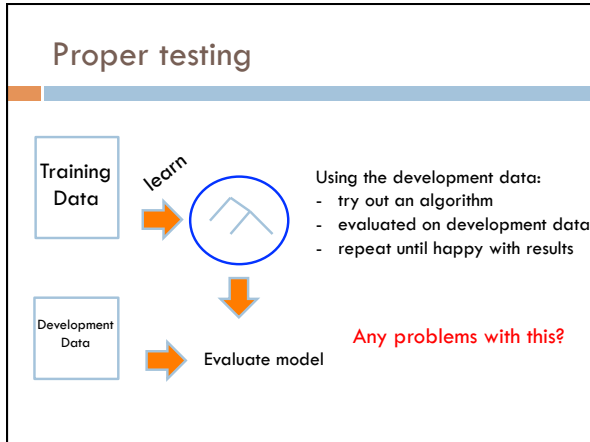
past Training Data (data with labels) → learn → (ML Model)

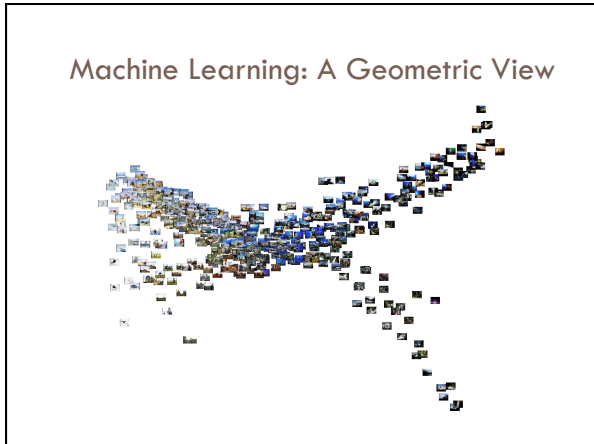
future Testing Data (data without labels) → predict → (ML Model)

How do we tell how well we're doing?









Apples vs. Bananas

Weight	Color	Label
4	Red	Apple
5	Yellow	Apple
6	Yellow	Banana
3	Red	Apple
7	Yellow	Banana
8	Yellow	Banana
6	Yellow	Apple

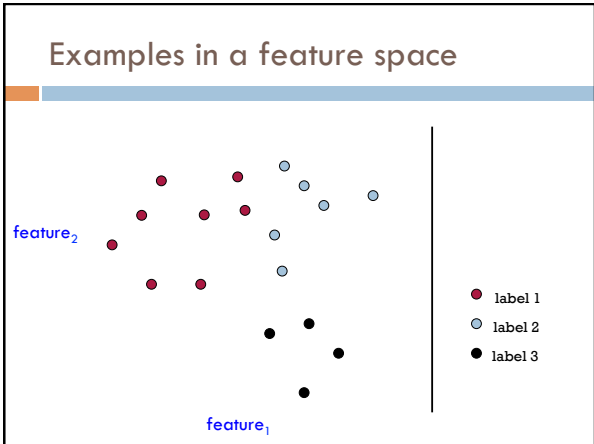
Can we visualize this data?

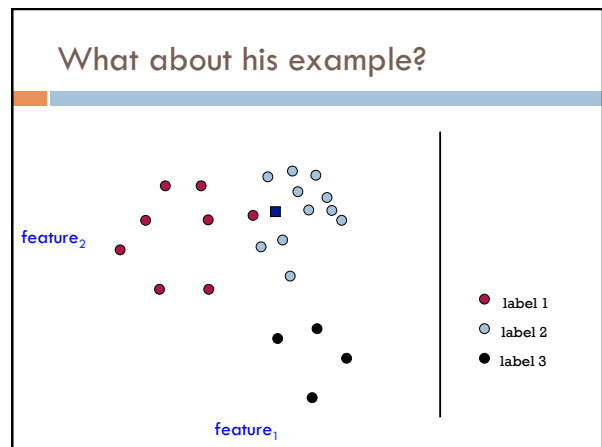
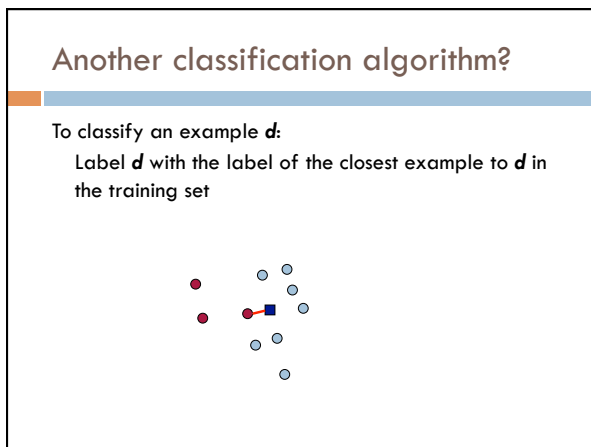
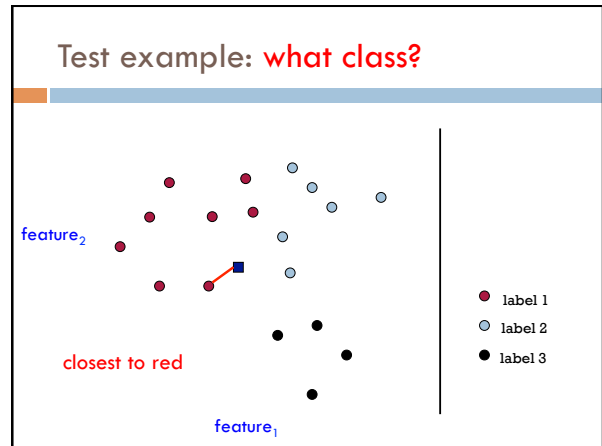
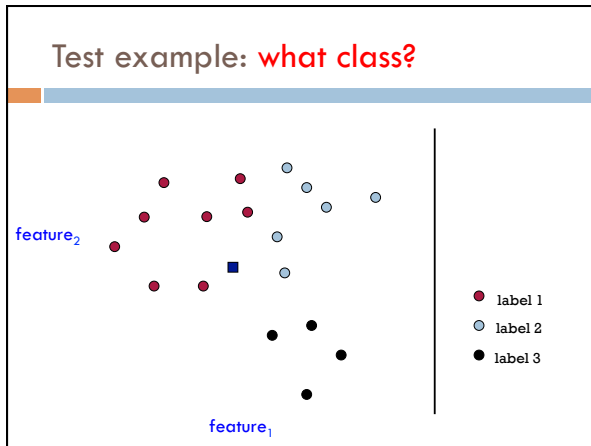
Apples vs. Bananas

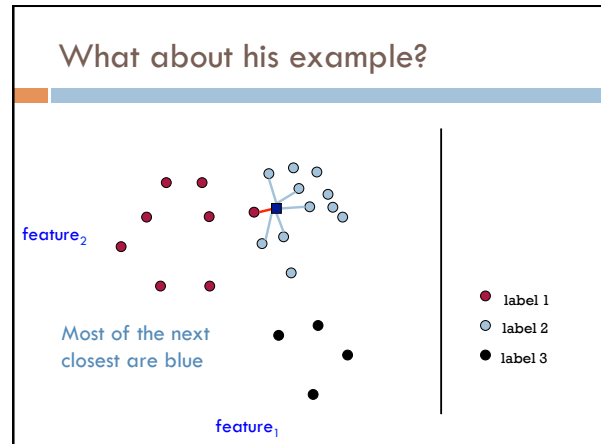
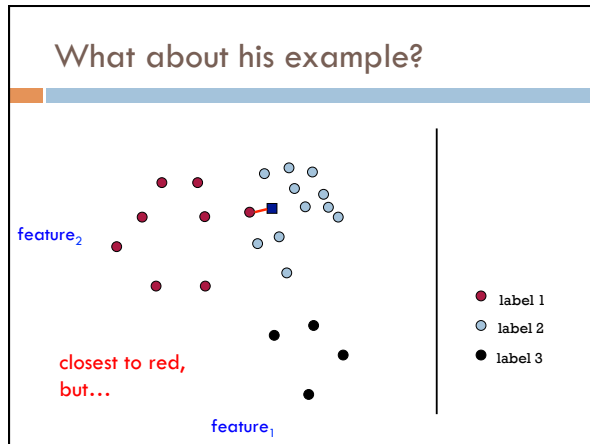
Turn features into numerical values
(read the book for a more detailed discussion of this)

Weight	Color	Label
4	0	Apple
5	1	Apple
6	1	Banana
3	0	Apple
7	1	Banana
8	1	Banana
6	1	Apple

We can view examples as points in an n -dimensional space where n is the number of features







k-Nearest Neighbor (k-NN)

To classify an example d :

- Find k nearest neighbors of d
- Choose as the label the **majority label** within the k nearest neighbors

k-Nearest Neighbor (k-NN)

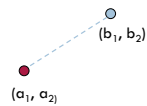
To classify an example d :

- Find k *nearest* neighbors of d
- Choose as the label the **majority label** within the k nearest neighbors

How do we measure "nearest"?

Euclidean distance

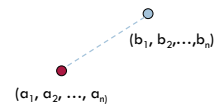
In two dimensions, how do we compute the distance?



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

Euclidean distance

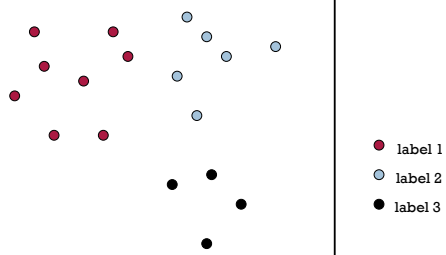
In n-dimensions, how do we compute the distance?



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

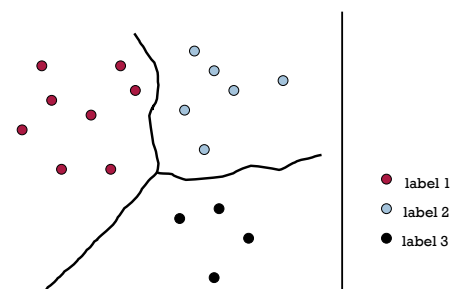
Decision boundaries

The **decision boundaries** are places in the features space where the classification of a point/example changes

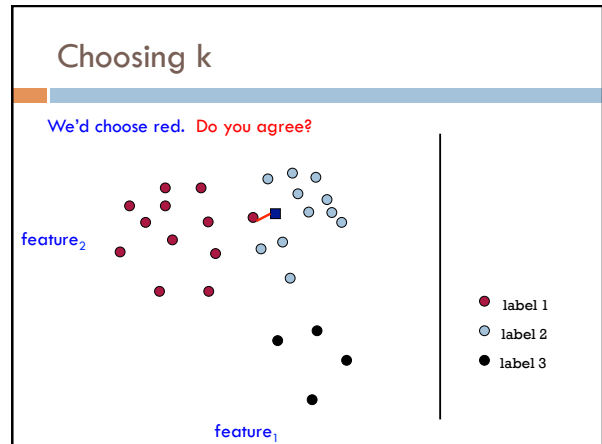
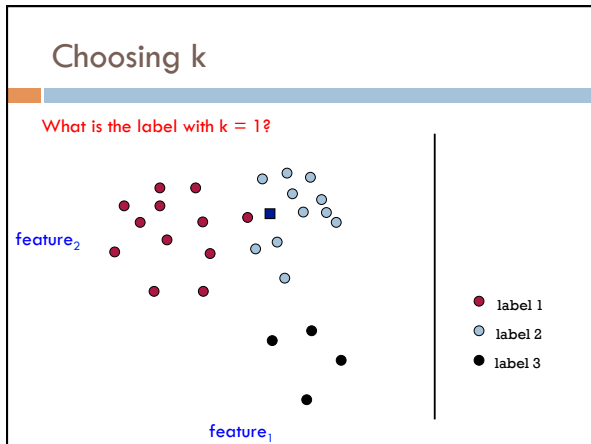
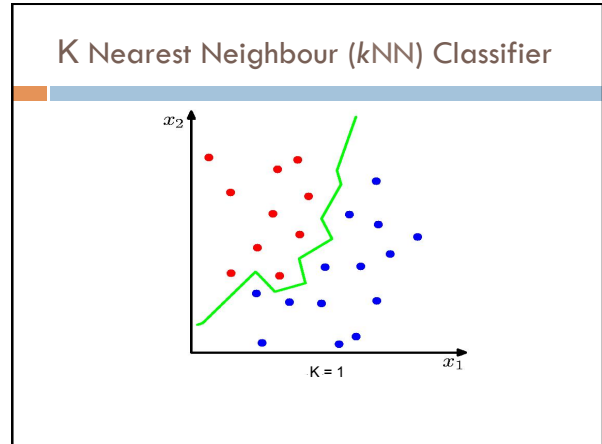
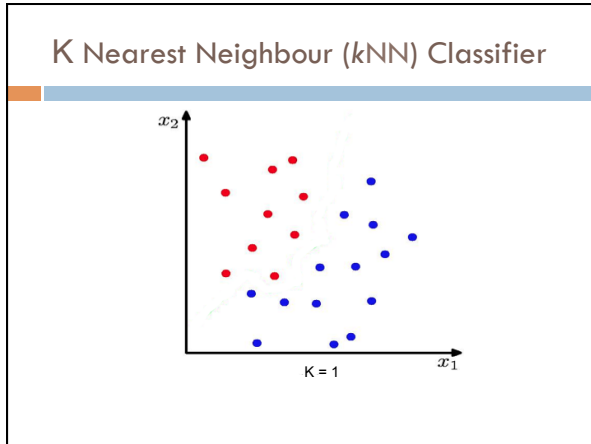


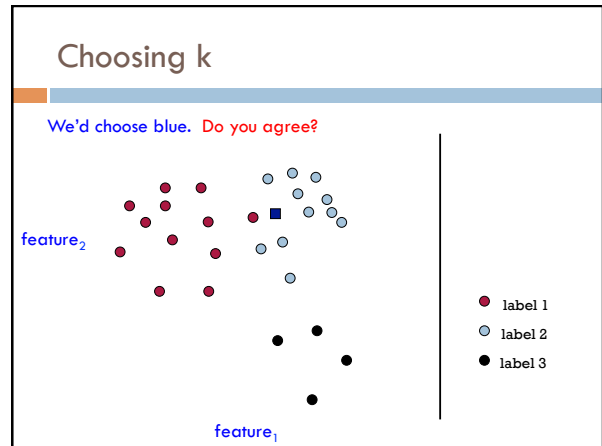
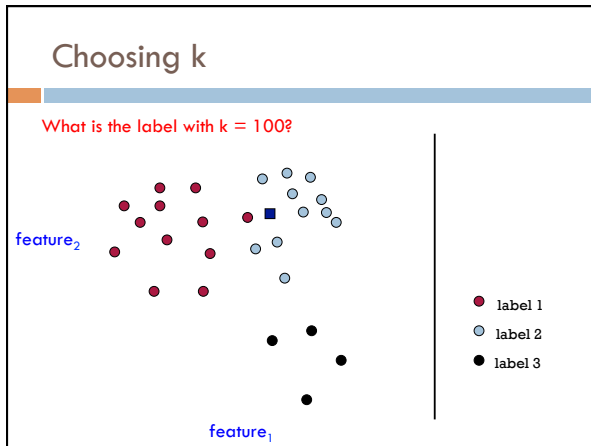
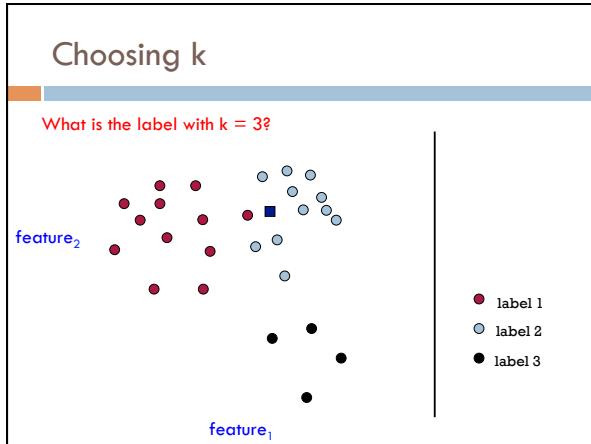
Where are the decision boundaries for k-NN?

k-NN decision boundaries

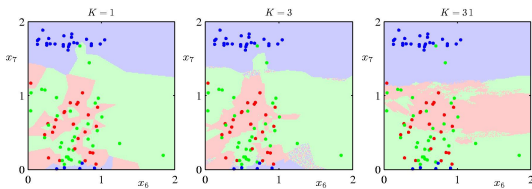


k-NN gives locally defined decision boundaries between classes





The impact of k



What is the role of k ?
 How does it relate to overfitting and underfitting?
 How did we control this for decision trees?

k-Nearest Neighbor (k-NN)

To classify an example d :

- ▣ Find k nearest neighbors of d
- ▣ Choose as the class the **majority class** within the k nearest neighbors

How do we choose k ?

How to pick k

Common heuristics:

- ▣ often 3, 5, 7
- ▣ choose an odd number to avoid ties

Use development data

k-NN variants

To classify an example d :

- ▣ Find k nearest neighbors of d
- ▣ Choose as the class the **majority class** within the k nearest neighbors

Any variation ideas?

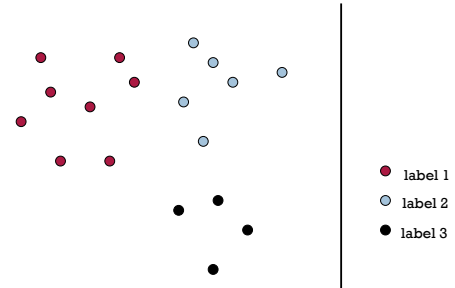
k-NN variations

Instead of k nearest neighbors, count majority from all examples within a fixed distance

Weighted k -NN:

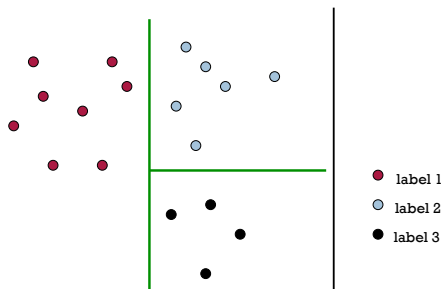
- ▣ Right now, all examples within examples are treated equally
- ▣ weight the "vote" of the examples, so that closer examples have more vote/weight
- ▣ often use some sort of exponential decay

Decision boundaries for decision trees



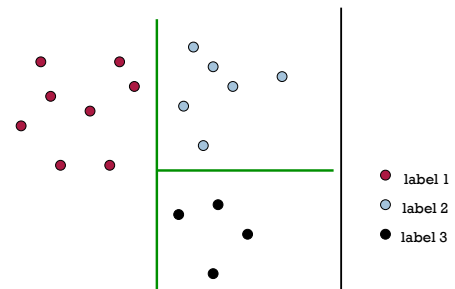
What are the decision boundaries for decision trees like?

Decision boundaries for decision trees



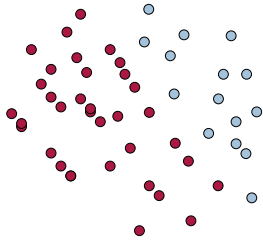
Axis-aligned splits/cuts of the data

Decision boundaries for decision trees



What types of data sets will DT work poorly on?

Problems for DT



Decision trees vs. k -NN

Which is faster to train?

Which is faster to classify?

Do they use the features in the same way to label the examples?

Decision trees vs. k -NN

Which is faster to train?

k -NN doesn't require any training!

Which is faster to classify?

For most data sets, decision trees

Do they use the features in the same way to label the examples?

k -NN treats all features equally! Decision trees "select" important features

A thought experiment

What is a 100,000-dimensional space like?

You're a 1-D creature, and you decide to buy a 2-unit apartment




2 rooms (very, skinny rooms)


Another thought experiment

What is a 100,000-dimensional space like?

Your job's going well and you're making good money. You upgrade to a 2-D apartment with 2-units per dimension



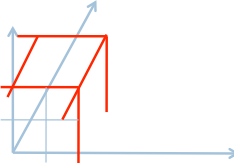
4 rooms (very, flat rooms)



Another thought experiment


What is a 100,000-dimensional space like?

You get promoted again and start having kids and decide to upgrade to another dimension.



8 rooms (very, normal rooms)

Each time you add a dimension, the amount of space you have to work with goes up exponentially




Another thought experiment

What is a 100,000-dimensional space like?

Larry Page steps down as CEO of google and they ask you if you'd like the job. You decide to upgrade to a 100,000 dimensional apartment.

**How much room do you have?
Can you have a big party?**

$2^{100,000}$ rooms (it's very quiet and lonely...) $\approx \sim 10^{30}$ rooms per person if you invited everyone on the planet



The challenge

Our intuitions about space/ distance don't scale with dimensions!

