

ADVANCED CLASSIFICATION
TECHNIQUES

David Kauchak
CS 159 – Fall 2014


Admin

ML lab next Monday

Project proposals: Sunday at 11:59pm

Project proposal presentations

Machine Learning: A Geometric View



Apples vs. Bananas

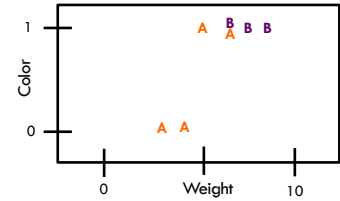
Weight	Color	Label
4	Red	Apple
5	Yellow	Apple
6	Yellow	Banana
3	Red	Apple
7	Yellow	Banana
8	Yellow	Banana
6	Yellow	Apple

Can we visualize this data?

Apples vs. Bananas

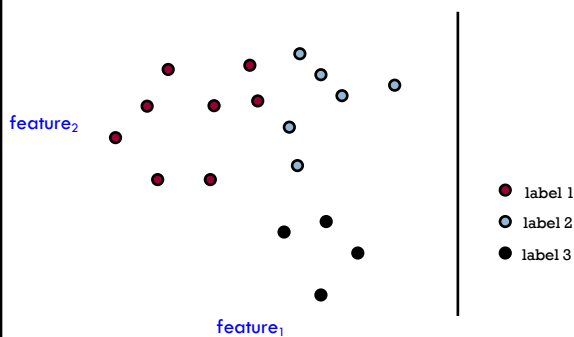
Turn features into numerical values

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

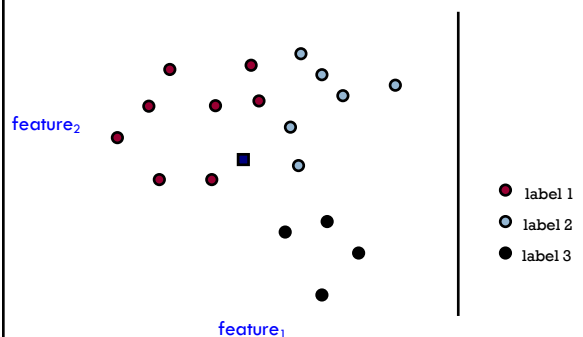


We can view examples as points in an n -dimensional space where n is the number of features called the **feature space**

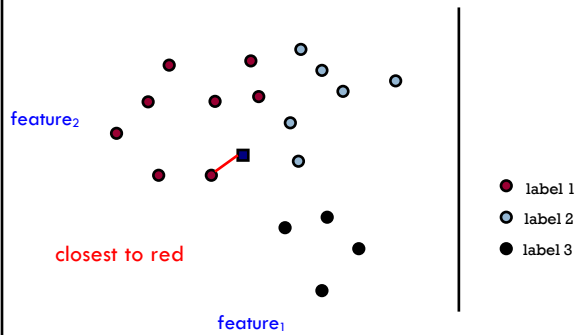
Examples in a feature space



Test example: what class?



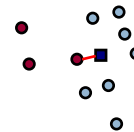
Test example: what class?



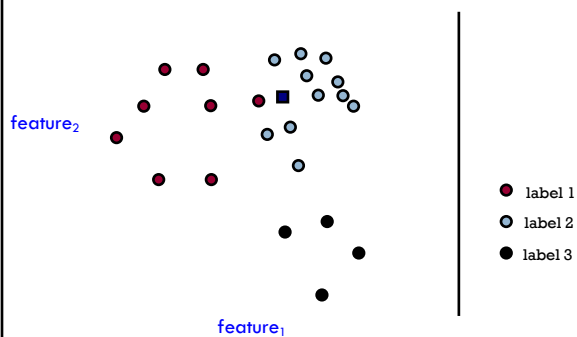
Another classification algorithm?

To classify an example d :

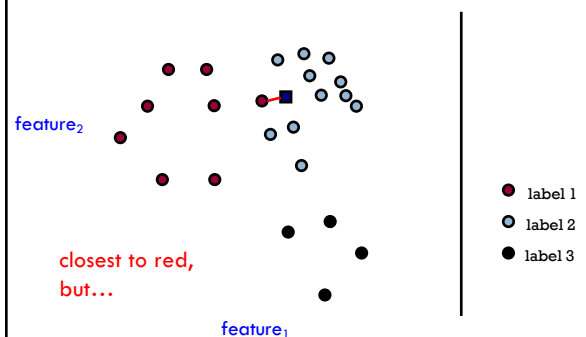
Label d with the label of the closest example to d in the training set



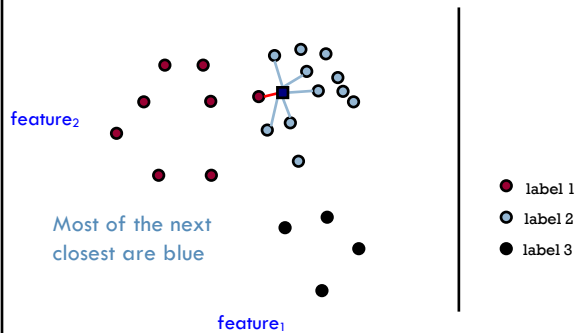
What about this example?



What about this example?



What about this example?



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)

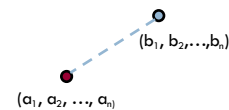
To classify an example d :

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

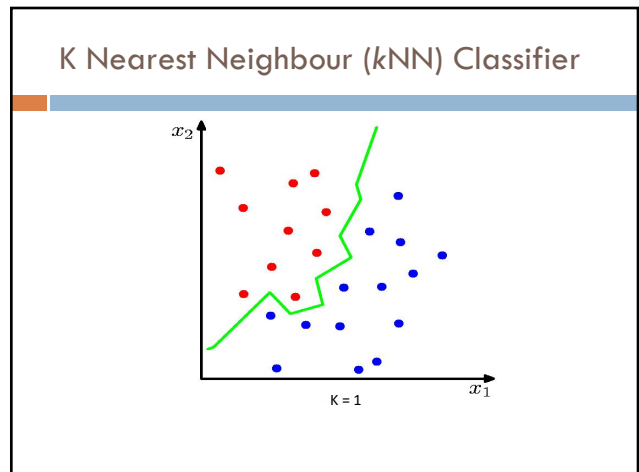
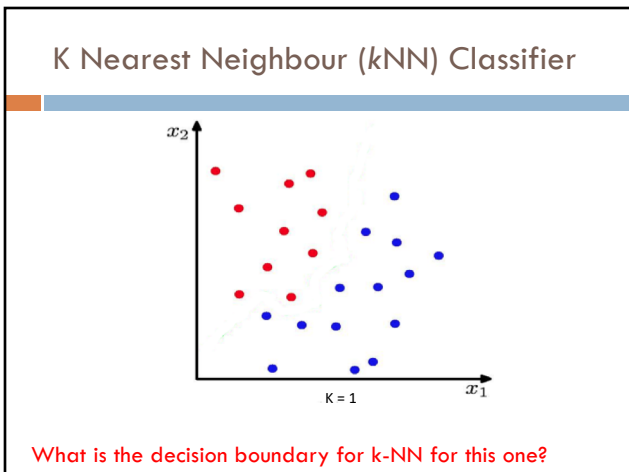
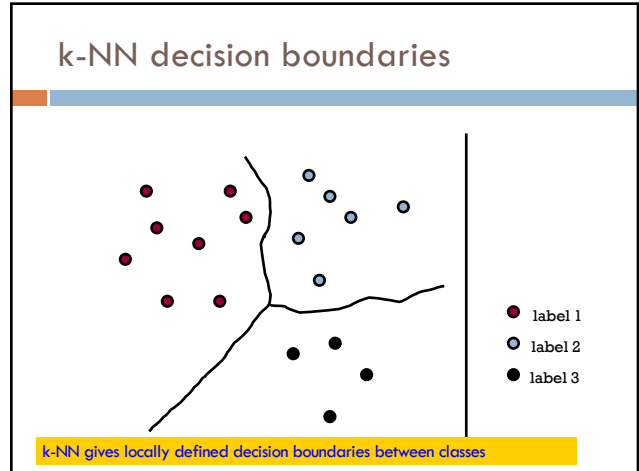
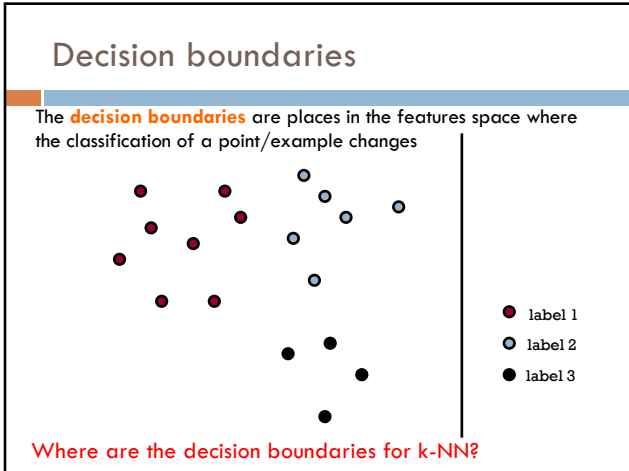
How do we measure “nearest”?

Euclidean distance

Euclidean distance! (or L1 or ...)



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



Machine learning models

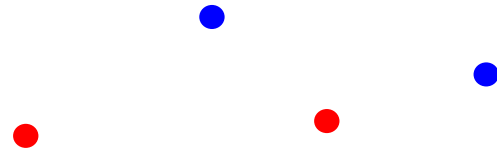
Some machine learning approaches make strong assumptions about the data

- ▣ If the assumptions are true this can often lead to better performance
- ▣ If the assumptions aren't true, they can fail miserably

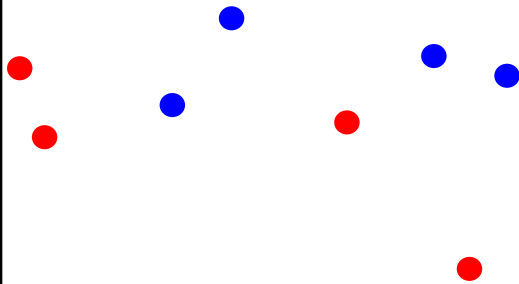
Other approaches don't make many assumptions about the data

- ▣ This can allow us to learn from more varied data
- ▣ But, they are more prone to overfitting
- ▣ and generally require more training data

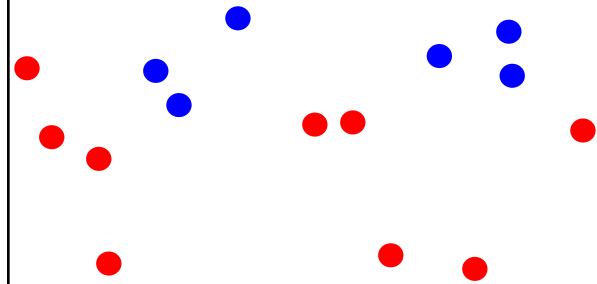
What is the data generating distribution?

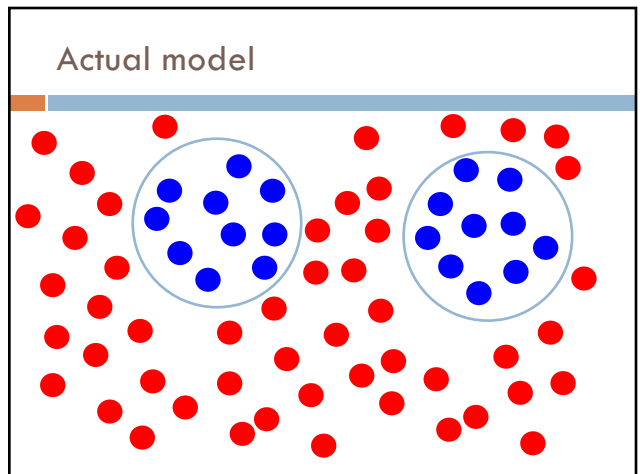
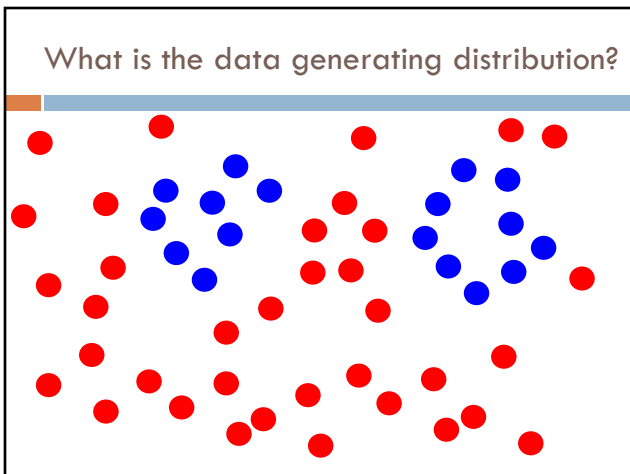
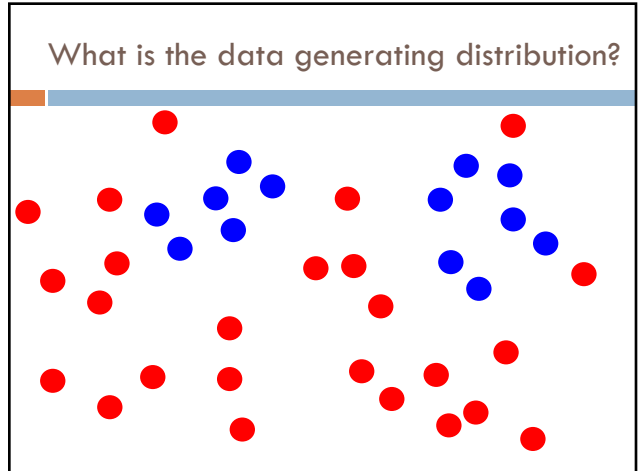
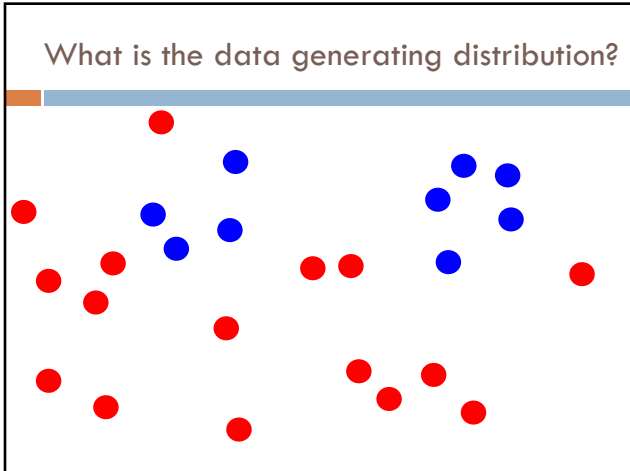


What is the data generating distribution?



What is the data generating distribution?



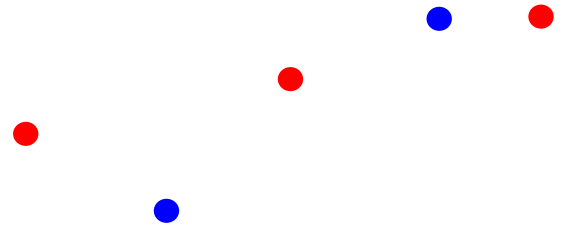


Model assumptions

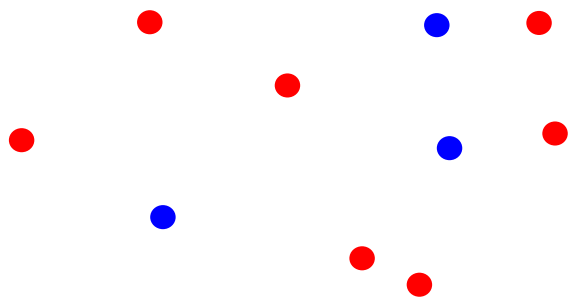
If you don't have strong assumptions about the model, it can take you a longer to learn

Assume now that our model of the blue class is two circles

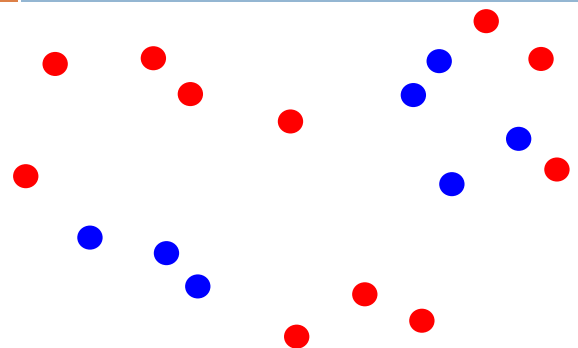
What is the data generating distribution?

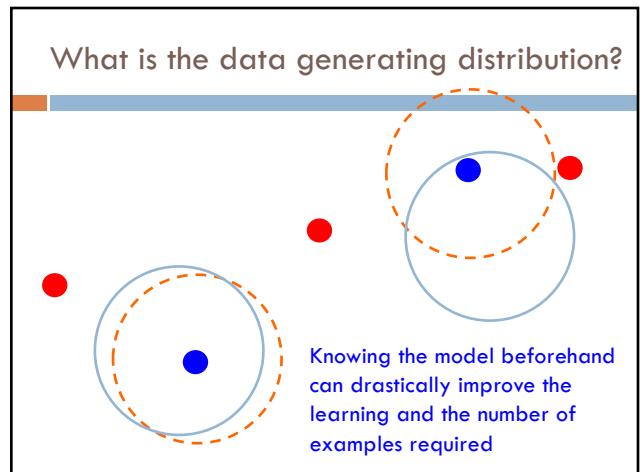
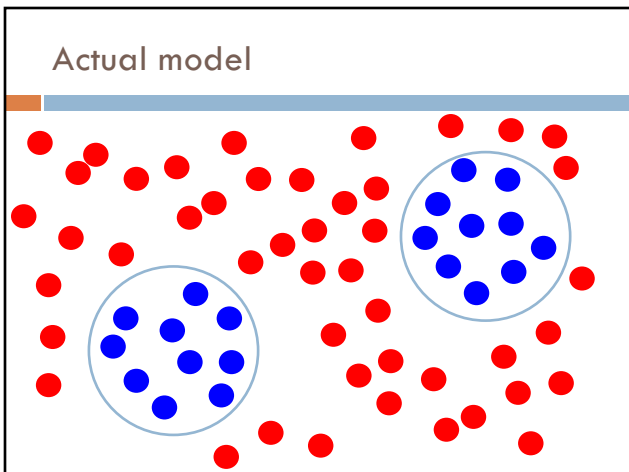
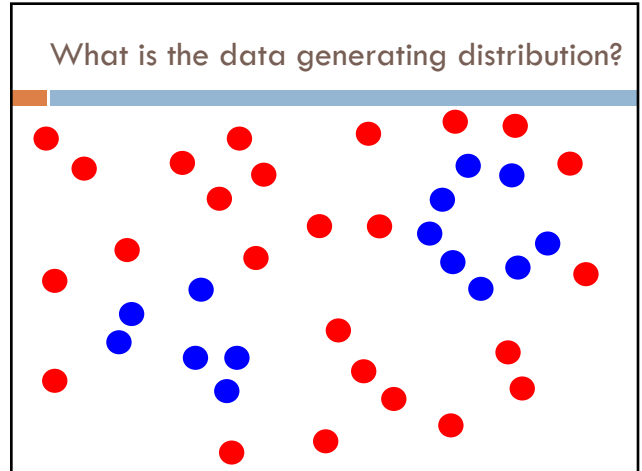
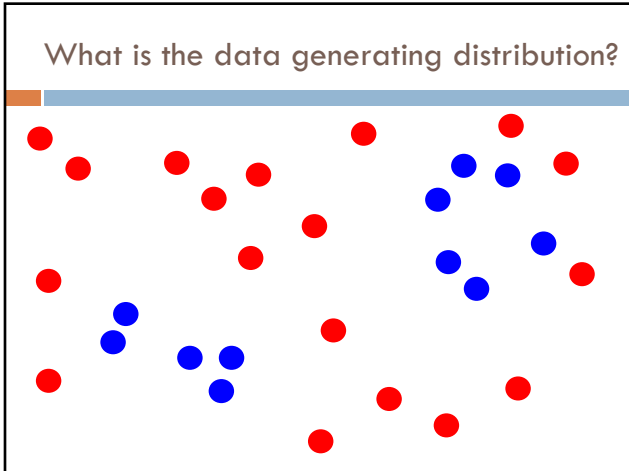


What is the data generating distribution?

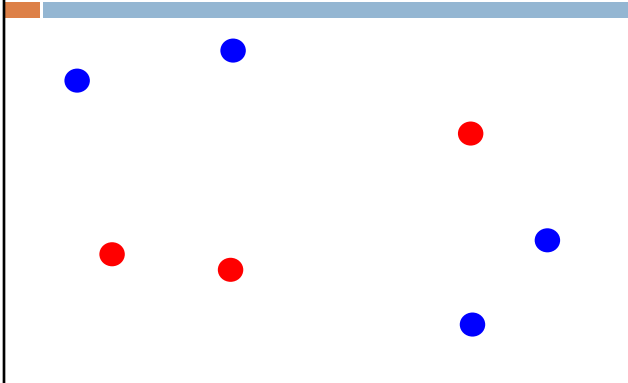


What is the data generating distribution?

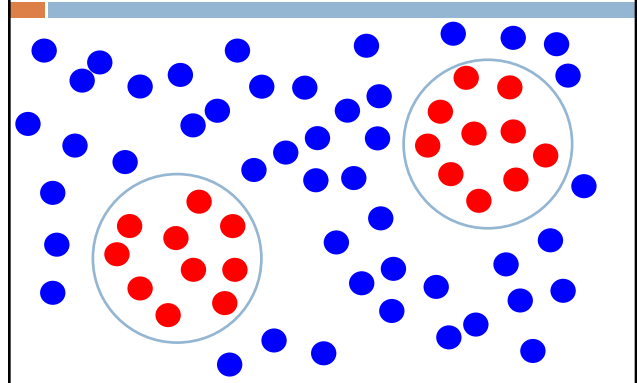




What is the data generating distribution?



Make sure your assumption is correct, though!

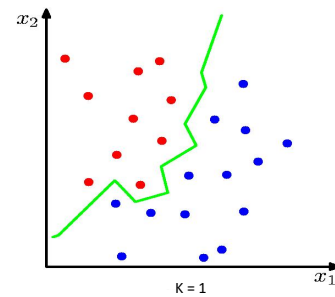


Machine learning models

What were the *model* assumptions (if any) that *k*-NN and NB made about the data?

Are there training data sets that could never be learned correctly by these algorithms?

k-NN model

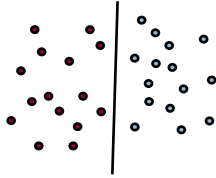


Linear models

A strong assumption is *linear separability*:

- in 2 dimensions, you can separate labels/classes by a line
- in higher dimensions, need hyperplanes

A *linear model* is a model that assumes the data is linearly separable



Hyperplanes

A hyperplane is line/plane in a high dimensional space



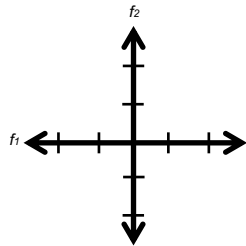
What defines a line?

What defines a hyperplane?

Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$



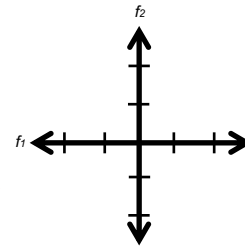
Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

What does this line look like?



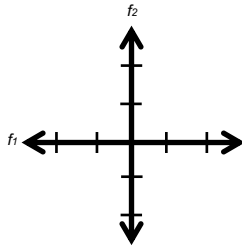
Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

-2	1
-1	0.5
0	0
1	-0.5
2	-1



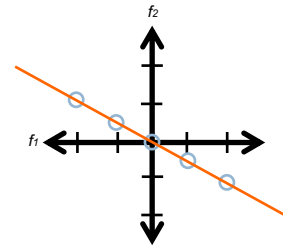
Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

-2	1
-1	0.5
0	0
1	-0.5
2	-1



Defining a line

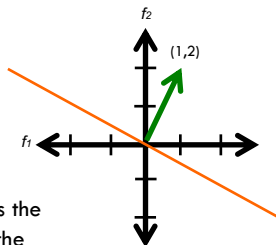
Any pair of values (w_1, w_2) defines a line through the origin:

$$0 = w_1 f_1 + w_2 f_2$$

$$0 = 1f_1 + 2f_2$$

$$w = (1, 2)$$

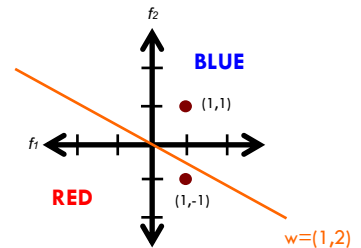
We can also view it as the line perpendicular to the weight vector



Classifying with a line

Mathematically, how can we classify points based on a line?

$$0 = 1f_1 + 2f_2$$



Classifying with a line

Mathematically, how can we classify points based on a line?

$0 = 1f_1 + 2f_2$

$(1,1): 1*1 + 2*1 = 3$

$(1,-1): 1*1 + 2*(-1) = -1$

The sign indicates which side of the line

Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$0 = w_1f_1 + w_2f_2$

$0 = 1f_1 + 2f_2$

How do we move the line off of the origin?

Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$a = w_1f_1 + w_2f_2$

$-1 = 1f_1 + 2f_2$

-2	
-1	
0	
1	
2	

Defining a line

Any pair of values (w_1, w_2) defines a line through the origin:

$a = w_1f_1 + w_2f_2$

$-1 = 1f_1 + 2f_2$

-2	0.5
-1	0
0	-0.5
1	-1
2	-1.5

Linear models

A linear model in n -dimensional space (i.e. n features) is defined by $n+1$ weights:

In two dimensions, a line:

$$0 = w_1 f_1 + w_2 f_2 + b \quad (\text{where } b = -a)$$

In three dimensions, a plane:

$$0 = w_1 f_1 + w_2 f_2 + w_3 f_3 + b$$

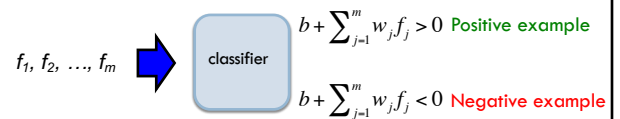
In n -dimensions, a *hyperplane*

$$0 = b + \sum_{i=1}^n w_i f_i$$



Classifying with a linear model

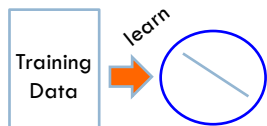
We can classify with a linear model by checking the sign:



Learning a linear model

Geometrically, we know what a linear model represents

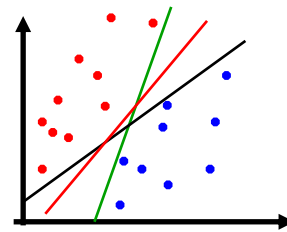
Given a linear model (i.e. a set of weights and b) we can classify examples

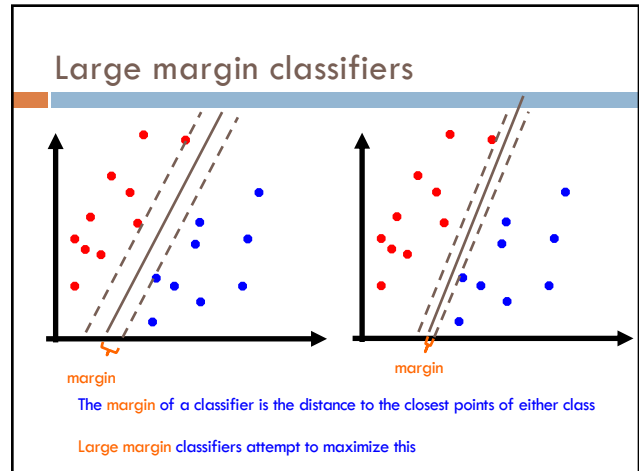
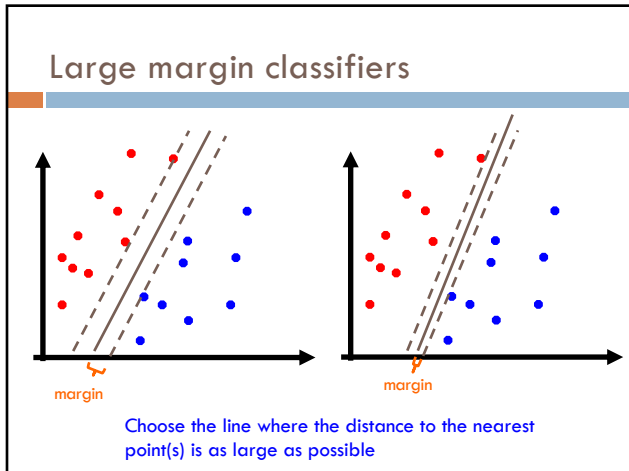


(data with labels)

How do we learn a linear model?

Which hyperplane would you choose?





Large margin classifier setup

Select the hyperplane with the largest margin where the points are classified correctly!

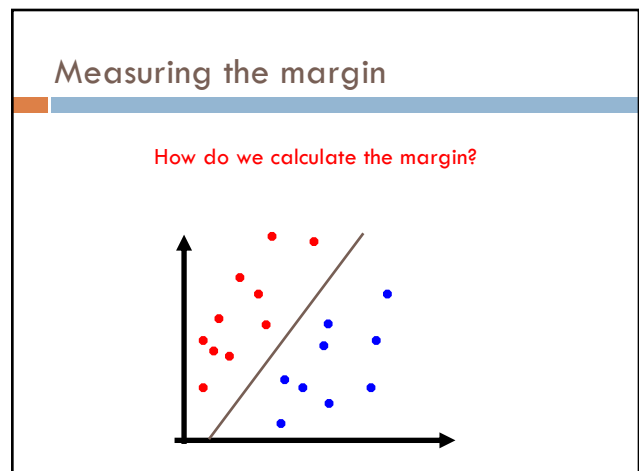
Setup as a **constrained optimization problem**:

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

$$y_i(w \cdot x_i + b) > 0 \quad \forall i \quad \text{what does this say?}$$

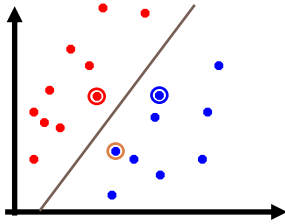
y_i : label for example i , either 1 (positive) or -1 (negative)
 x_i : our feature **vector** for example i



Support vectors

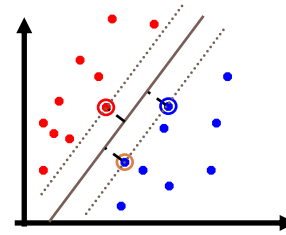
For any separating hyperplane, there exist some set of "closest points"

These are called the support vectors



Measuring the margin

The margin is the distance to the support vectors, i.e. the "closest points", on either side of the hyperplane



Support vector machine problem

Posted as a [quadratic optimization problem](#)

Maximize/minimize a quadratic function

Subject to a set of linear constraints

Many, many variants of solving this problem

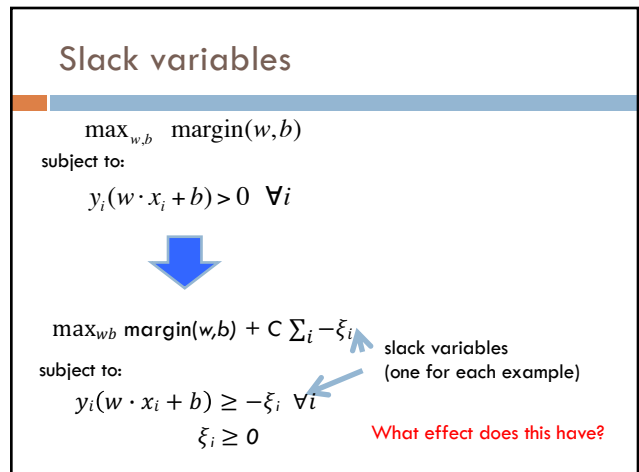
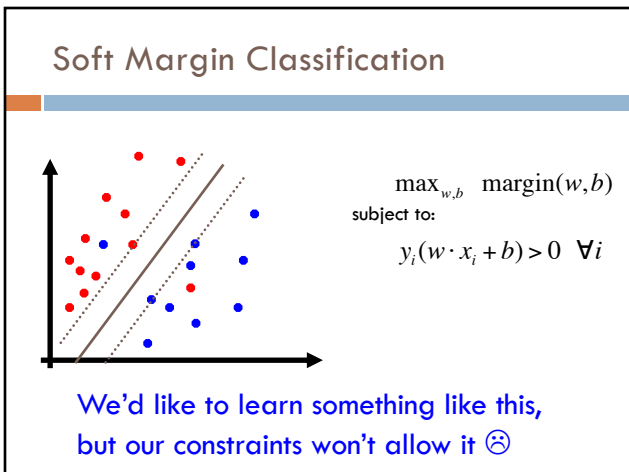
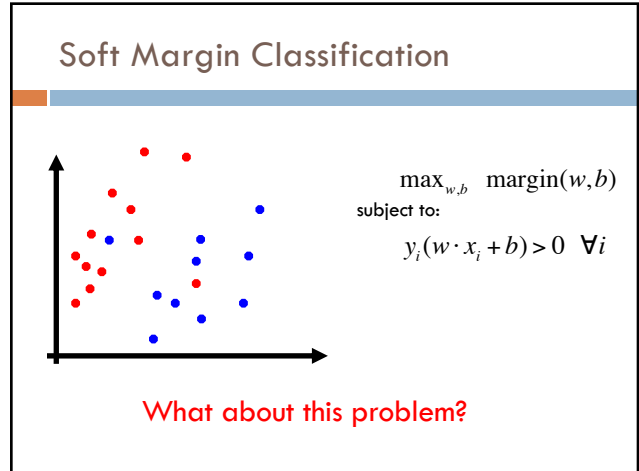
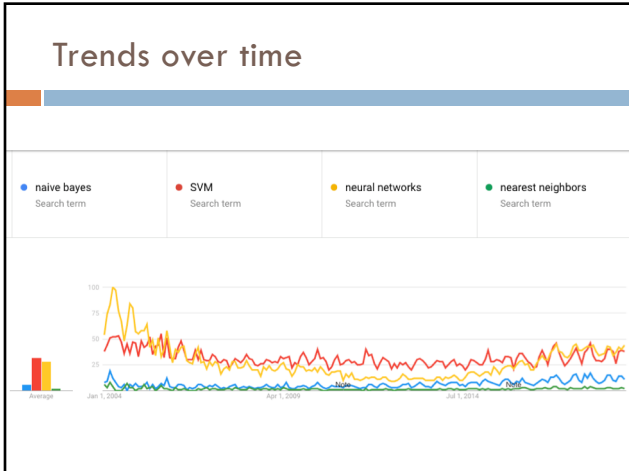
One of the most successful classification approaches

Support vector machines

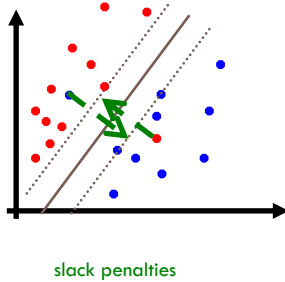
One of the most successful (if not the most successful) classification approach:

decision tree	About 2,240,000 results (0.32 sec)
Support vector machine	About 2,180,000 results (0.36 sec)
k nearest neighbor	About 844,000 results (0.33 sec)
Naïve Bayes	About 71,300 results (0.32 sec)

Google
scholar



Slack variables



$$\max_{w,b} \text{margin}(w,b) + C \sum_i \xi_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq -\xi_i \quad \forall i$$

$$\xi_i \geq 0$$

Slack variables

margin

trade-off between margin maximization and penalization

$$\max_{w,b} \text{margin}(w,b) + C \sum_i \xi_i$$

penalized by how far from "correct"

subject to:

$$y_i(w \cdot x_i + b) \geq -\xi_i \quad \forall i$$

allowed to make a mistake

$$\xi_i \geq 0$$

Soft margin SVM

$$\max_{w,b} \text{margin}(w,b) + C \sum_i \xi_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq -\xi_i \quad \forall i$$

$$\xi_i \geq 0$$

Still a **quadratic optimization problem!**

Other successful classifiers in NLP

Perceptron algorithm

- Linear classifier
- Trains "online"
- Fast and easy to implement
- Often used for tuning parameters (not necessarily for classifying)

Logistic regression classifier (aka Maximum entropy classifier)

- Probabilistic classifier
- Doesn't have the NB constraints
- Performs very well
- More computationally intensive to train than NB

Resources

SVM

- ▣ SVM light: <http://svmlight.joachims.org/>
- ▣ Others, but this one is awesome!

Maximum Entropy classifier

- ▣ <http://nlp.stanford.edu/software/classifier.shtml>

General ML frameworks:

- ▣ Python: scikit-learn, MLpy
- ▣ Java: Weka (<http://www.cs.waikato.ac.nz/ml/weka/>)
- ▣ Many others...

Quiz 3

Mean 23: 80%

Median: 23.5 (81%)