

# Classification

Chris Amato  
Northeastern University

Some images and slides are used from: Rob Platt,  
CS188 UC Berkeley, AIMA

# Supervised learning

Given: *Training set*  $\{(x_i, y_i) \mid i = 1 \dots N\}$ , given a labeled set of input-output pairs  $D = \{(x_i, y_i)\}_i$

Find: A good approximation to  $f: X \rightarrow Y$  **Function approximation**

Examples: what are  $X$  and  $Y$  ?

Spam Detection – Map email to {Spam, Not Spam} **Binary Classification**

Digit recognition – Map pixels to {0,1,2,3,4,5,6,7,8,9} **Multiclass Classification**

Stock Prediction – Map new, historic prices, etc. to (the real numbers) **Regression**

# Supervised learning

Goal: make predictions on novel inputs, meaning ones that we have not seen before (this is called *generalization*)

Formalize this problem as approximating a function:  $f(x)=y$

The learning problem is then to use *function approximation* to discover:  $\hat{f}(x) = \hat{y}$

# Spam example

Input: email

Output: spam/ham

Setup:

Get a large collection of example emails, each labeled "spam" or "ham"

Note: someone has to hand label all this data!

Want to learn to predict labels of new, future emails

Features: attributes used to make the ham/spam decision

Words: FREE!

Text Patterns: \$dd, CAPS

Non-text: SenderInContacts

...



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

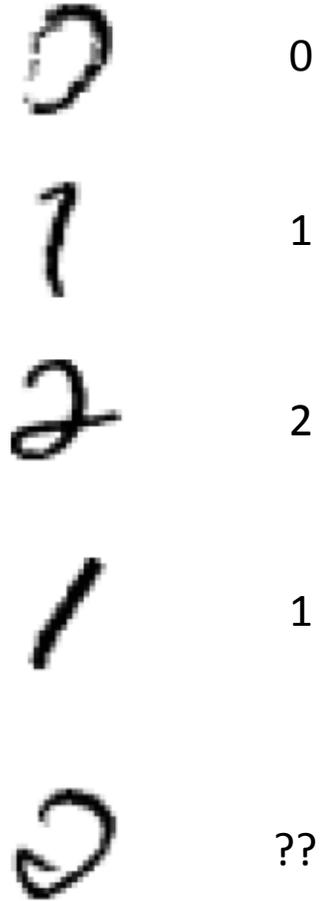
99 MILLION EMAIL ADDRESSES FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

# Digit recognition

- Input: images / pixel grids
- Output: a digit 0-9
- Setup:
  - Get a large collection of example images, each labeled with a digit
  - Note: someone has to hand label all this data!
  - Want to learn to predict labels of new, future digit images
- Features: The attributes used to make the digit decision
  - Pixels: (6,8)=ON
  - Shape Patterns: NumComponents, AspectRatio, NumLoops
  - ...



# Some applications

Document classification and email spam filtering

Image classification and handwriting recognition

Face detection and recognition

Fraud detection

Medical diagnosis

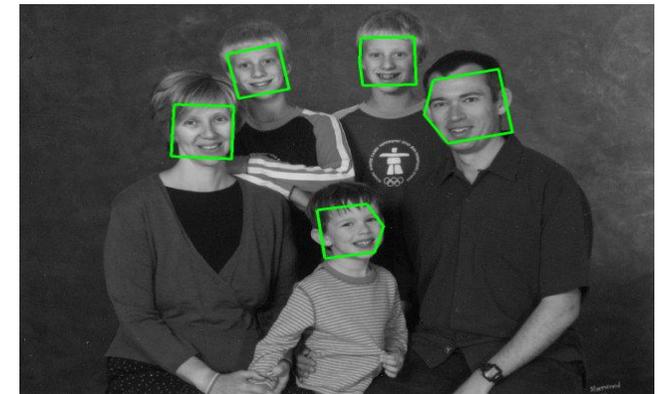
✗ Dear Sir.  
First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...

✗ TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.  
99 MILLION EMAIL ADDRESSES FOR ONLY \$99

✓ Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.



true class = 7 	true class = 2 	true class = 1 
true class = 0 	true class = 4 	true class = 1 
true class = 4 	true class = 9 	true class = 5 



# Probabilistic Classification

Want a probability distribution over possible labels, given the input vector  $x$  and training set  $D$  by  $P(y|x,D)$

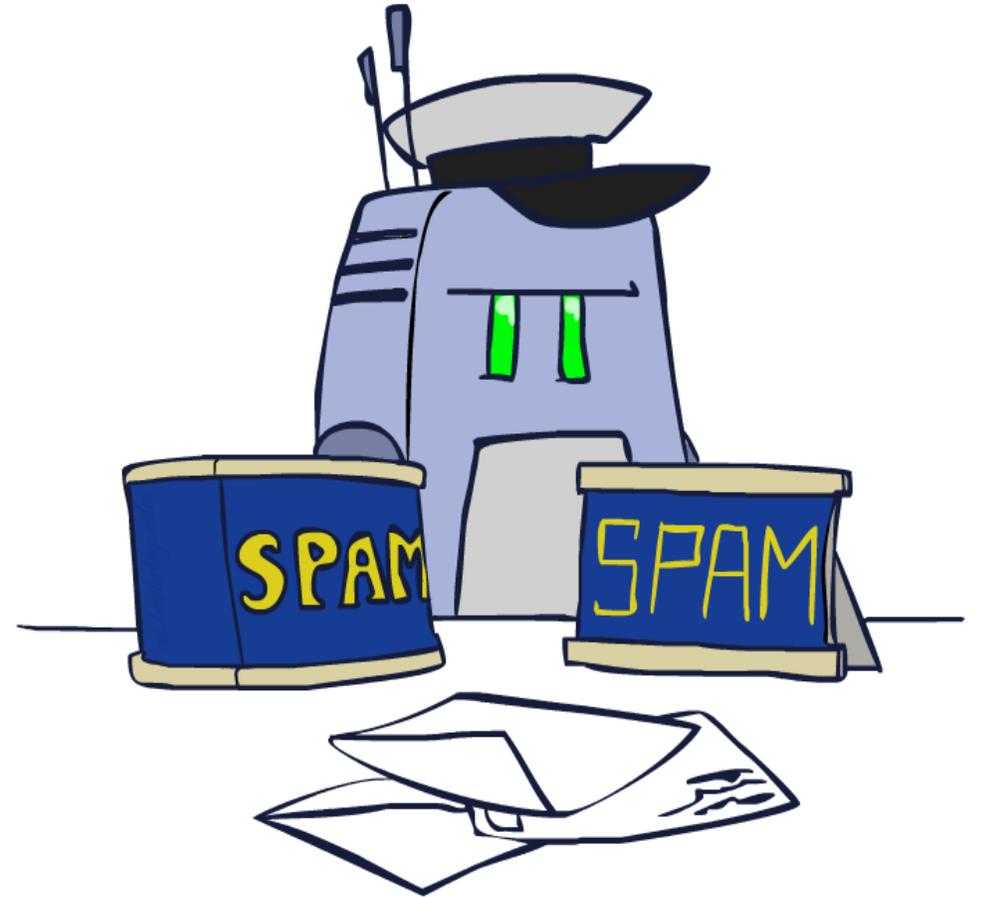
In general, this represents a vector of length  $C$

Calculate a “best guess”:  $\hat{y} = \hat{f}(x) = \operatorname{argmax}_{c=1}^C P(y = c | x, D)$

This corresponds to the most probable class label (the mode of the distribution)

# Model-Based Classification

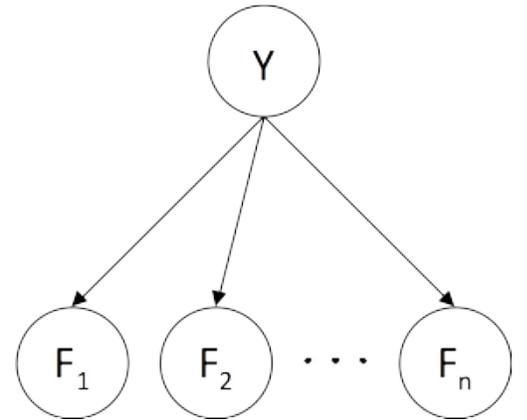
- Model-based approach
  - Build a model (e.g. Bayes' net) where both the label and features are random variables
  - Instantiate any observed features
  - Query for the distribution of the label conditioned on the features
- Challenges
  - What structure should the BN have?
  - How should we learn its parameters?



# Naïve Bayes for Digits

- Naïve Bayes: Assume all features are independent effects of the label
- Simple digit recognition version:
  - One feature (variable)  $F_{ij}$  for each grid position  $\langle i,j \rangle$
  - Feature values are on / off, based on whether intensity is more or less than 0.5 in underlying image
  - Each input maps to a feature vector, e.g.

1  $\rightarrow \langle F_{0,0} = 0 \ F_{0,1} = 0 \ F_{0,2} = 1 \ F_{0,3} = 1 \ F_{0,4} = 0 \ \dots F_{15,15} = 0 \rangle$



- Here: lots of features, each is binary valued

- Naïve Bayes model: 
$$P(Y|F_{0,0} \dots F_{15,15}) \propto P(Y) \prod_{i,j} P(F_{i,j}|Y)$$
- What do we need to learn?

# Note: Discriminative vs Generative

Can calculate  $P(y|x)$  directly, but we could also use Bayes rule to calculate  $P(y|x) = \alpha P(x|y)P(y)$

Calculating  $P(y|x)$  directly is called *discriminative*

Calculating  $P(x|y)P(y)$  is called *generative* (since it represents a way to generate the data)

# Generative models

Use Bayes rule to switch classification into a generative problem

Determine the right parameters for your model

$$P(y = c | \mathbf{x}, \theta) \propto P(x | y = c, \theta) P(y = c | \theta)$$

Unknown parameters

Class conditional density

Prior

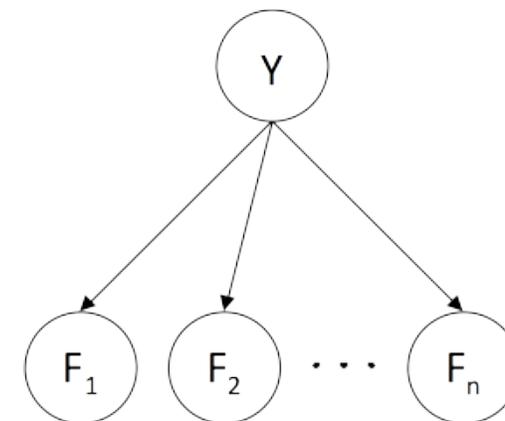
The diagram shows the equation  $P(y = c | \mathbf{x}, \theta) \propto P(x | y = c, \theta) P(y = c | \theta)$ . The two terms on the right,  $P(x | y = c, \theta)$  and  $P(y = c | \theta)$ , are circled in red. A red line points from the text 'Unknown parameters' to the  $\theta$  parameter in both terms. A red line points from 'Class conditional density' to the  $P(x | y = c, \theta)$  term, and another red line points from 'Prior' to the  $P(y = c | \theta)$  term.

# General Naïve Bayes

- A general Naive Bayes model:

$$P(Y, F_1 \dots F_n) = \underbrace{P(Y)}_{|Y| \text{ parameters}} \prod_i \underbrace{P(F_i|Y)}_{n \times |F| \times |Y| \text{ parameters}}$$

$|Y| \times |F|^n$  values



- We only have to specify how each feature depends on the class
- Total number of parameters is *linear* in  $n$
- Model is very simplistic, but often works anyway

# Inference for Naïve Bayes

- Goal: compute posterior distribution over label variable  $Y$

- Step 1: get joint probability of label and evidence for each label

$$P(Y, f_1 \dots f_n) = \begin{bmatrix} P(y_1, f_1 \dots f_n) \\ P(y_2, f_1 \dots f_n) \\ \vdots \\ P(y_k, f_1 \dots f_n) \end{bmatrix} \Rightarrow \begin{bmatrix} P(y_1) \prod_i P(f_i|y_1) \\ P(y_2) \prod_i P(f_i|y_2) \\ \vdots \\ P(y_k) \prod_i P(f_i|y_k) \end{bmatrix}$$

$\xrightarrow{+}$

- Step 2: sum to get probability of evidence

$$P(f_1 \dots f_n)$$

- Step 3: normalize by dividing Step 1 by Step 2

$$P(Y|f_1 \dots f_n)$$

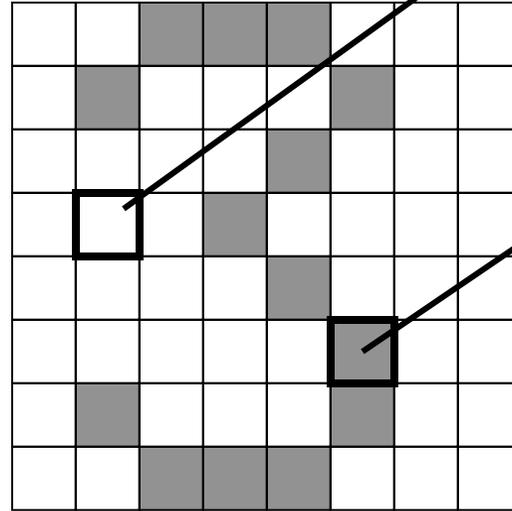
# General Naïve Bayes

- What do we need in order to use Naïve Bayes?
  - Inference method (we just saw this part)
    - Start with a bunch of probabilities:  $P(Y)$  and the  $P(F_i|Y)$  tables
    - Use standard inference to compute  $P(Y|F_1...F_n)$
    - Nothing new here
  - Estimates of local conditional probability tables
    - $P(Y)$ , the prior over labels
    - $P(F_i|Y)$  for each feature (evidence variable)
    - These probabilities are collectively called the *parameters* of the model and denoted by  $\theta$
    - Up until now, we assumed these appeared by magic, but...
    - ...they typically come from training data counts: we'll look at this soon

# Example: Conditional Probabilities

$P(Y)$

1	0.1
2	0.1
3	0.1
4	0.1
5	0.1
6	0.1
7	0.1
8	0.1
9	0.1
0	0.1



$P(F_{3,1} = on|Y)$

1	0.01
2	0.05
3	0.05
4	0.30
5	0.80
6	0.90
7	0.05
8	0.60
9	0.50
0	0.80

$P(F_{5,5} = on|Y)$

1	0.05
2	0.01
3	0.90
4	0.80
5	0.90
6	0.90
7	0.25
8	0.85
9	0.60
0	0.80

# Naïve Bayes for Text

- **Bag-of-words Naïve Bayes:**

- Features:  $W_i$  is the word at position  $i$
- As before: predict label conditioned on feature variables (spam vs. ham)
- As before: assume features are conditionally independent given label
- New: each  $W_i$  is identically distributed

- **Generative model:** 
$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i|Y)$$

*Word at position  
 $i$ , not  $i^{\text{th}}$  word in  
the dictionary!*



- **“Tied” distributions and bag-of-words**

- Usually, each variable gets its own conditional probability distribution  $P(F|Y)$
- In a bag-of-words model
  - Each position is identically distributed
  - All positions share the same conditional probs  $P(W|Y)$
  - Why make this assumption?
- Called “bag-of-words” because model is insensitive to word order or reordering

# Example: Spam Filtering

- Model:  $P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i|Y)$
- What are the parameters?

$P(Y)$

ham : 0.66
spam: 0.33

$P(W|\text{spam})$

the : 0.0156
to : 0.0153
and : 0.0115
of : 0.0095
you : 0.0093
a : 0.0086
with: 0.0080
from: 0.0075
...

$P(W|\text{ham})$

the : 0.0210
to : 0.0133
of : 0.0119
2002: 0.0110
with: 0.0108
from: 0.0107
and : 0.0105
a : 0.0100
...

- Where do these tables come from?

# Spam Example

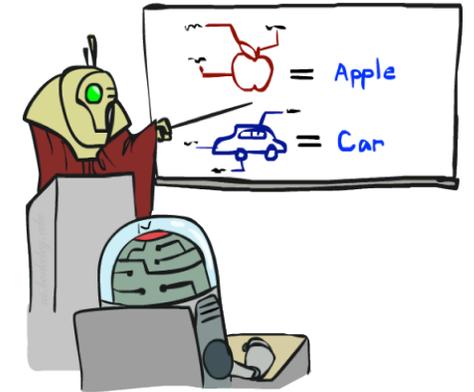
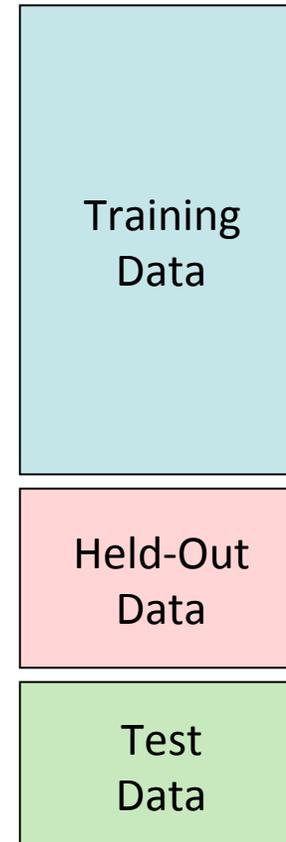
Word	P(w spam)	P(w ham)	Tot Spam	Tot Ham
(prior)	0.33333	0.66666	-1.1	-0.4
Gary	0.00002	0.00021	-11.8	-8.9
would	0.00069	0.00084	-19.1	-16.0
you	0.00881	0.00304	-23.8	-21.8
like	0.00086	0.00083	-30.9	-28.9
to	0.01517	0.01339	-35.1	-33.2
lose	0.00008	0.00002	-44.5	-44.0
weight	0.00016	0.00002	-53.3	-55.0
while	0.00027	0.00027	-61.5	-63.2
you	0.00881	0.00304	-66.2	-69.0
sleep	0.00006	0.00001	-76.0	-80.5

---

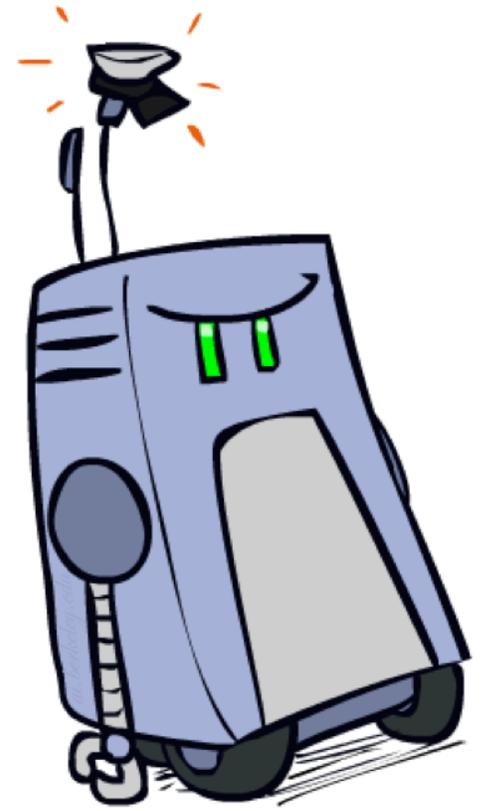
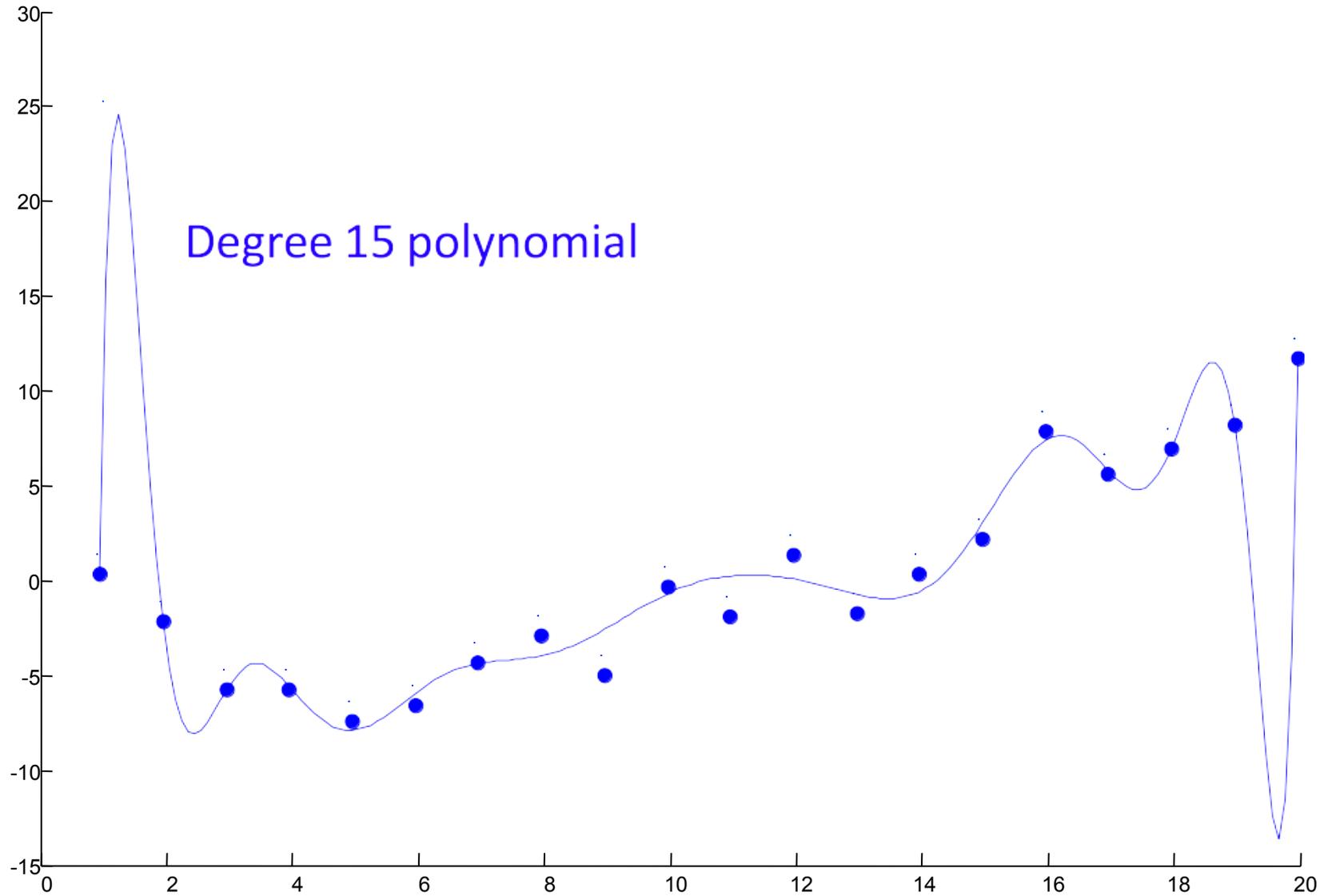
$P(\text{spam} | w) = 98.9$

# Important Concepts

- **Data:** labeled instances, e.g. emails marked spam/ham
  - Training set
  - Held out set
  - Test set
- **Features:** attribute-value pairs which characterize each  $x$
- **Experimentation cycle**
  - Learn parameters (e.g. model probabilities) on training set
  - (Tune hyperparameters on held-out set)
  - Compute accuracy of test set
  - Very important: never “peek” at the test set!
- **Evaluation**
  - Accuracy: fraction of instances predicted correctly
- **Overfitting and generalization**
  - Want a classifier which does well on *test* data
  - Overfitting: fitting the training data very closely, but not generalizing well



# Overfitting



# Example: Overfitting

$P(\text{features}, C = 2)$

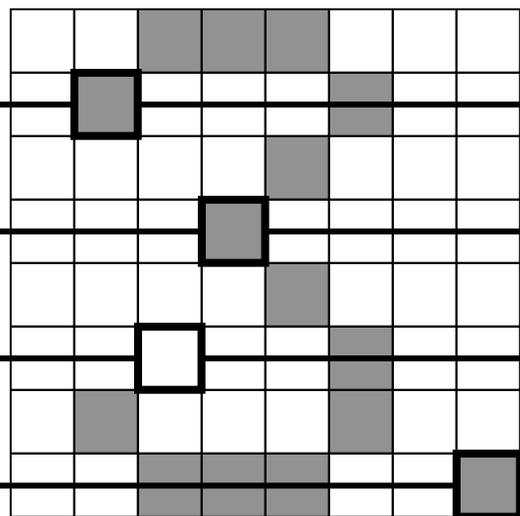
$P(C = 2) = 0.1$

$P(\text{on} | C = 2) = 0.8$

$P(\text{on} | C = 2) = 0.1$

$P(\text{off} | C = 2) = 0.1$

$P(\text{on} | C = 2) = 0.01$



$P(\text{features}, C = 3)$

$P(C = 3) = 0.1$

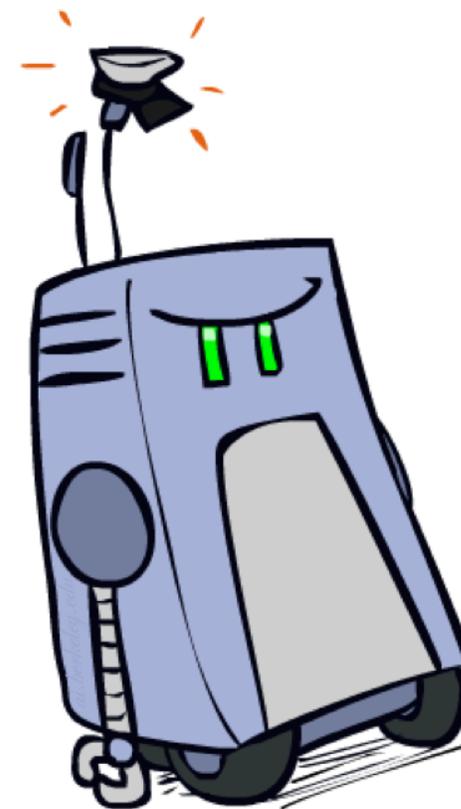
$P(\text{on} | C = 3) = 0.8$

$P(\text{on} | C = 3) = 0.9$

$P(\text{off} | C = 3) = 0.7$

$P(\text{on} | C = 3) = 0.0$

*2 wins!!*



# Example: Overfitting

- Posterior determined by *relative* probabilities (odds ratios):

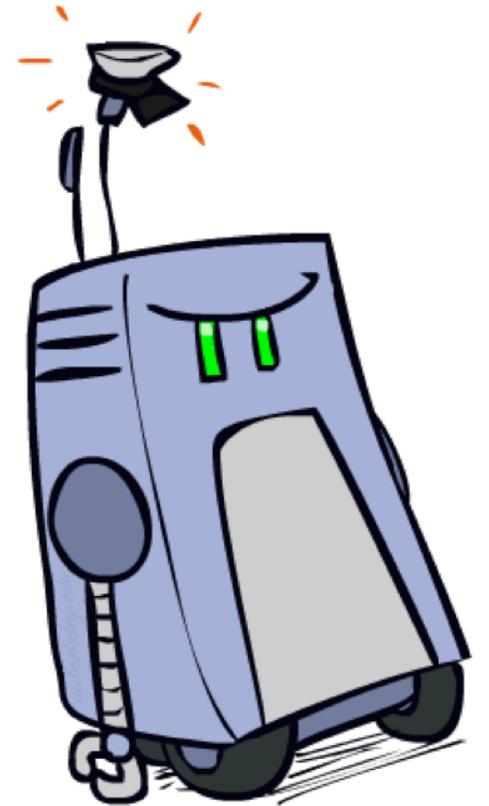
$$\frac{P(W|\text{ham})}{P(W|\text{spam})}$$

south-west	: inf
nation	: inf
morally	: inf
nicely	: inf
extent	: inf
seriously	: inf
...	

$$\frac{P(W|\text{spam})}{P(W|\text{ham})}$$

screens	: inf
minute	: inf
guaranteed	: inf
\$205.00	: inf
delivery	: inf
signature	: inf
...	

*What went wrong here?*



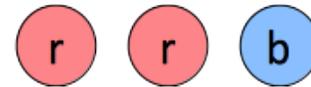
# Generalization and Overfitting

- Relative frequency parameters will **overfit** the training data!
  - Just because we never saw a 3 with pixel (15,15) on during training doesn't mean we won't see it at test time
  - Unlikely that every occurrence of "minute" is 100% spam
  - Unlikely that every occurrence of "seriously" is 100% ham
  - What about all the words that don't occur in the training set at all?
  - In general, we can't go around giving unseen events zero probability
- As an extreme case, imagine using the entire email as the only feature
  - Would get the training data perfect (if deterministic labeling)
  - Wouldn't *generalize* at all
  - Just making the bag-of-words assumption gives us some generalization, but isn't enough
- To generalize better: we need to **smooth** or **regularize** the estimates

# Parameter Estimation

- Estimating the distribution of a random variable
- *Elicitation*: ask a human (why is this hard?)
- *Empirically*: use training data (learning!)
  - E.g.: for each outcome  $x$ , look at the *empirical rate* of that value:

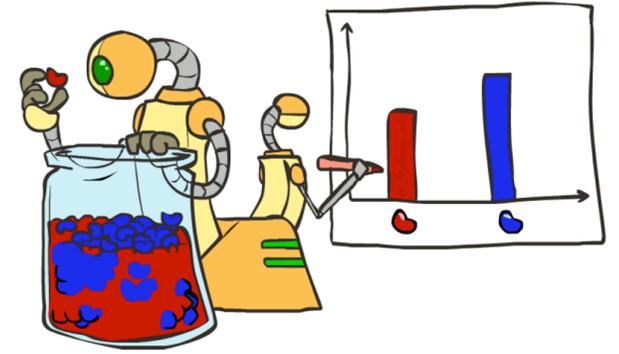
$$P_{\text{ML}}(x) = \frac{\text{count}(x)}{\text{total samples}}$$



$$P_{\text{ML}}(r) = 2/3$$

- This is the estimate that maximizes the *likelihood of the data*

$$L(x, \theta) = \prod_i P_{\theta}(x_i)$$



# Maximum Likelihood?

- Relative frequencies are the maximum likelihood estimates

$$\begin{aligned}\theta_{ML} &= \arg \max_{\theta} P(\mathbf{X}|\theta) \\ &= \arg \max_{\theta} \prod_i P_{\theta}(X_i)\end{aligned}$$

$$P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

- Another option is to consider the most likely parameter value given the data

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} P(\theta|\mathbf{X}) \\ &= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)/P(\mathbf{X}) \quad \longrightarrow \quad \text{????} \\ &= \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)\end{aligned}$$

# Maximum Likelihood Estimation

What should the parameter values be?

Calculate the *maximum-likelihood estimate (MLE)* (or really the log-likelihood)

$$\hat{\theta} \triangleq \operatorname{argmax}_{\theta} \log P(D | \theta)$$

Assume the training examples are *independent and identically distributed (iid)*, can calculate the log-likelihood

$$l(\theta) \triangleq \log P(D | \theta) = \sum_{i=1}^N \log P(y_i | x_i, \theta)$$

We can either maximize log-likelihood or minimize *negative log likelihood (NLL)*

$$NLL(\theta) \triangleq -\sum_{i=1}^N \log P(y_i | x_i, \theta)$$

# MLE vs. MAP

MAP:  $\hat{y}^{MAP} = \operatorname{argmax}_{c=1}^C P(y = c | D) = \operatorname{argmax}_{c=1}^C P(D | y = c)P(y = c)$

(or in our case)

$$\hat{\theta}^{MAP} = \operatorname{argmax}_{\theta} P(\theta | D) = \operatorname{argmax}_{\theta} P(D | \theta) P(\theta)$$

MLE:

$$\hat{\theta}^{MLE} = \operatorname{argmax}_{\theta} P(D | \theta)$$

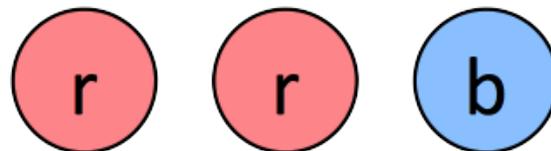
Likelihood Prior

MLE and MAP (can) converge as the dataset grows

# Unseen Events: Laplace Smoothing

- Laplace's estimate:

- Pretend you saw every outcome once more than you actually did



$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$
$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{ML}(X) = \left\langle \frac{2}{3}, \frac{1}{3} \right\rangle$$

$$P_{LAP}(X) = \left\langle \frac{3}{5}, \frac{2}{5} \right\rangle$$

- Can derive this estimate with *Dirichlet priors*

# Unseen Events: Laplace Smoothing

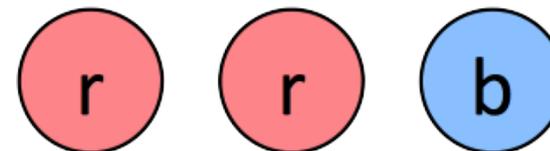
- Laplace's estimate (extended):
  - Pretend you saw every outcome  $k$  extra times

$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

- What's Laplace with  $k = 0$ ?
- $k$  is the **strength** of the prior

- Laplace for conditionals:
  - Smooth each condition independently:

$$P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$$



$$P_{LAP,0}(X) = \left\langle \frac{2}{3}, \frac{1}{3} \right\rangle$$

$$P_{LAP,1}(X) = \left\langle \frac{3}{5}, \frac{2}{5} \right\rangle$$

$$P_{LAP,100}(X) = \left\langle \frac{102}{203}, \frac{101}{203} \right\rangle$$

# Estimation: Linear Interpolation\*

- In practice, Laplace often performs poorly for  $P(X|Y)$ :
  - When  $|X|$  is very large
  - When  $|Y|$  is very large
- Another option: linear interpolation
  - Also get the empirical  $P(X)$  from the data
  - Make sure the estimate of  $P(X|Y)$  isn't too different from the empirical  $P(X)$

$$P_{LIN}(x|y) = \alpha \hat{P}(x|y) + (1.0 - \alpha) \hat{P}(x)$$

- What if  $\alpha$  is 0? 1?

# Real NB: Smoothing

- For real classification problems, smoothing is critical
- New odds ratios:

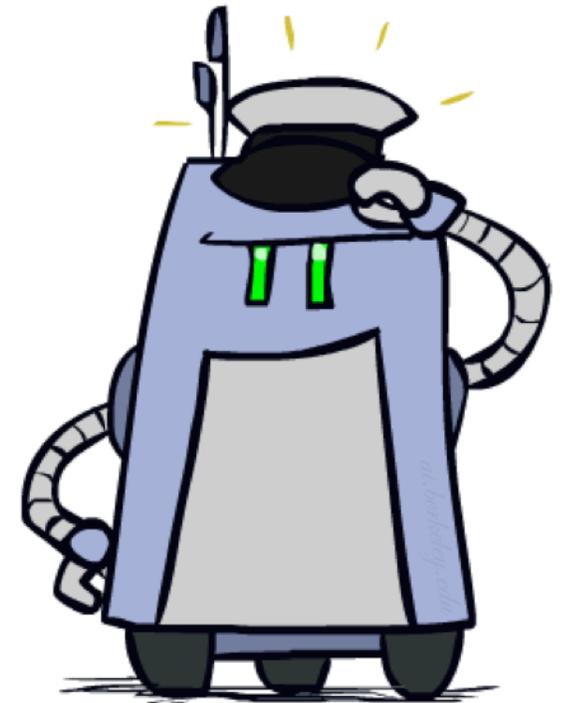
$$\frac{P(W|\text{ham})}{P(W|\text{spam})}$$

helvetica	:	11.4
seems	:	10.8
group	:	10.2
ago	:	8.4
areas	:	8.3
...		

$$\frac{P(W|\text{spam})}{P(W|\text{ham})}$$

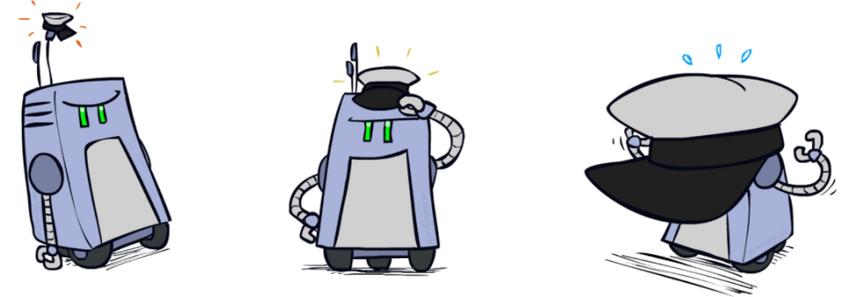
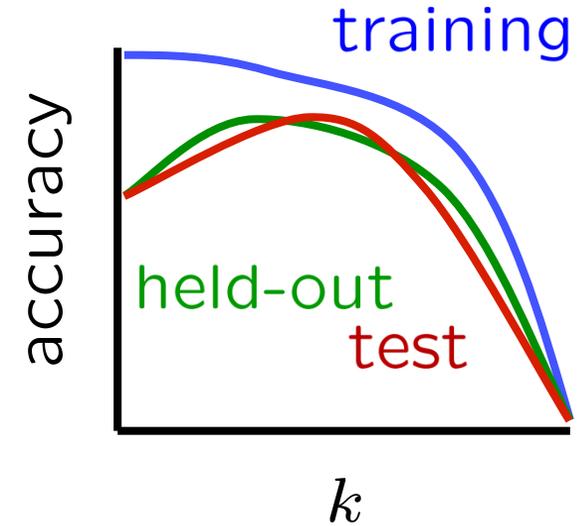
verdana	:	28.8
Credit	:	28.4
ORDER	:	27.2
<FONT>	:	26.9
money	:	26.5
...		

*Do these make more sense?*



# Tuning on Held-Out Data

- Now we've got two kinds of unknowns
  - Parameters: the probabilities  $P(X|Y)$ ,  $P(Y)$
  - Hyperparameters: e.g. the amount / type of smoothing to do,  $k$ ,  $\alpha$
- What should we learn where?
  - Learn parameters from training data
  - Tune hyperparameters on different data
    - Why?
  - For each value of the hyperparameters, train and test on the held-out data
  - Choose the best value and do a final test on the test data



# What to Do About Errors

- Problem: there's still spam in your inbox
- Need more **features** – words aren't enough!
  - Have you emailed the sender before?
  - Have 1M other people just gotten the same email?
  - Is the sending information consistent?
  - Is the email in ALL CAPS?
  - Do inline URLs point where they say they point?
  - Does the email address you by (your) name?
- Naïve Bayes models can incorporate a variety of features, but tend to do best in homogeneous cases (e.g. all features are word occurrences)



# Baselines

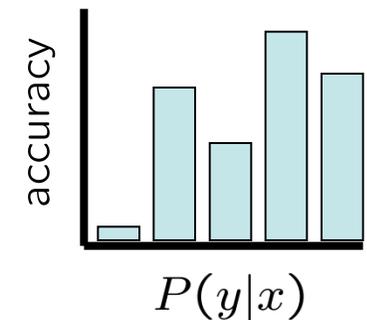
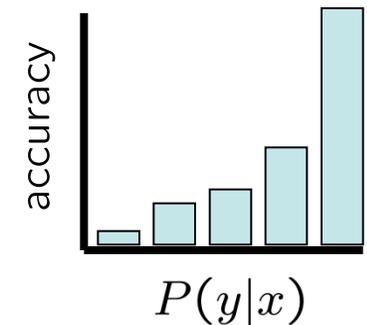
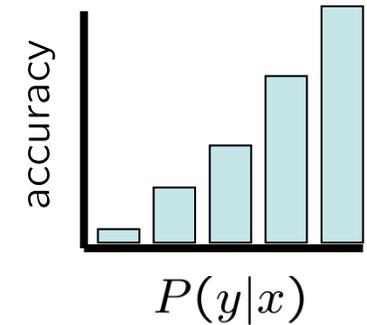
- **First step: get a baseline**
  - Baselines are very simple “straw man” procedures
  - Help determine how hard the task is
  - Help know what a “good” accuracy is
- **Weak baseline: most frequent label classifier**
  - Gives all test instances whatever label was most common in the training set
  - E.g. for spam filtering, might label everything as ham
  - Accuracy might be very high if the problem is skewed
  - E.g. calling everything “ham” gets 66%, so a classifier that gets 70% isn’t very good...
- **For real research, usually use previous work as a (strong) baseline**

# Confidences from a Classifier

- The **confidence** of a probabilistic classifier:
  - Posterior over the top label

$$\text{confidence}(x) = \max_y P(y|x)$$

- Represents how sure the classifier is of the classification
- Any probabilistic model will have confidences
- No guarantee confidence is correct



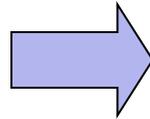
# Summary

- Bayes rule lets us do diagnostic queries with causal probabilities
- The naïve Bayes assumption takes all features to be independent given the class label
- We can build classifiers out of a naïve Bayes model using training data
- Smoothing estimates is important in real systems
- Classifier confidences are useful, when you can get them

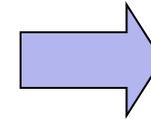
# Feature Vectors

 $x$  $f(x)$  $y$ 

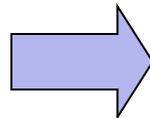
```
Hello,  
  
Do you want free print  
cartridges? Why pay more  
when you can get them  
ABSOLUTELY FREE! Just
```



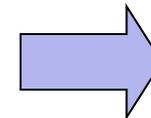
```
# free      : 2  
YOUR_NAME   : 0  
MISPELLED   : 2  
FROM_FRIEND : 0  
...
```



SPAM  
or  
+

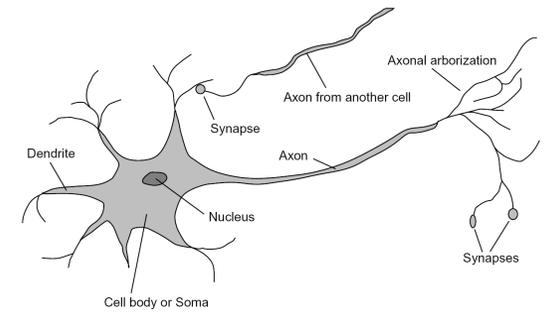


```
PIXEL-7,12  : 1  
PIXEL-7,13  : 0  
...  
NUM_LOOPS   : 1  
...
```



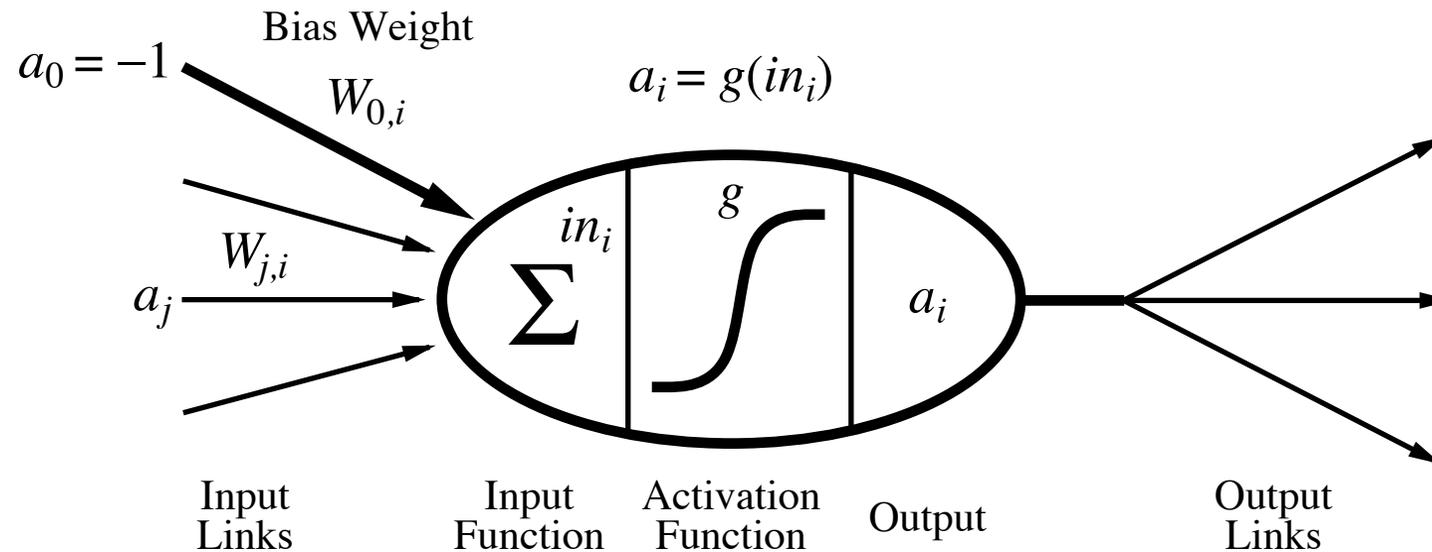
"2"

# McCulloch–Pitts “unit”



Output is a “squashed” linear function of the inputs:

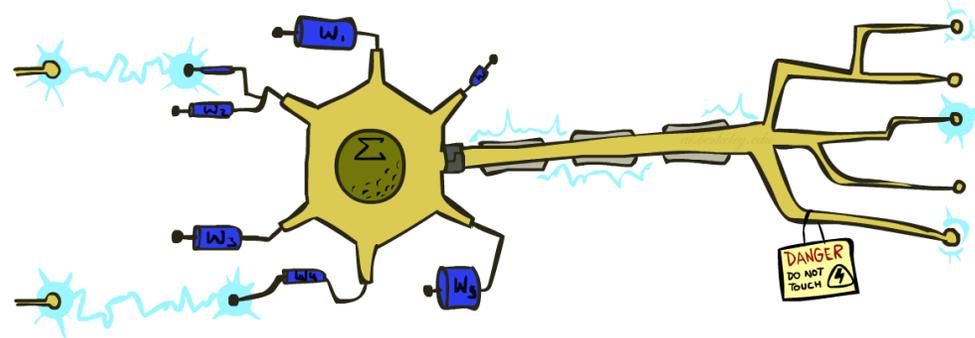
$$a_i \leftarrow g(in_i) = g\left(\sum_j W_{j,i} a_j\right)$$



A gross oversimplification of real neurons, helps to develop understanding of what networks of simple units can do

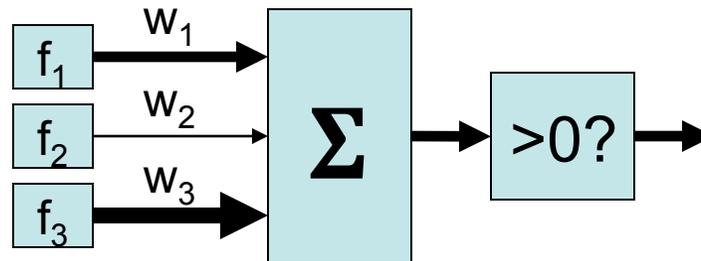
# Linear Classifiers

- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
  - Positive, output +1
  - Negative, output -1



# Weights

- Binary case: compare features to a weight vector
- Learning: figure out the weight vector from examples

```
# free      : 4  
YOUR_NAME  :-1  
MISPELLED  : 1  
FROM_FRIEND :-3  
...
```

$w$

$f(x_1)$

```
# free      : 2  
YOUR_NAME  : 0  
MISPELLED  : 2  
FROM_FRIEND : 0  
...
```

$f(x_2)$

```
# free      : 0  
YOUR_NAME  : 1  
MISPELLED  : 1  
FROM_FRIEND : 1  
...
```

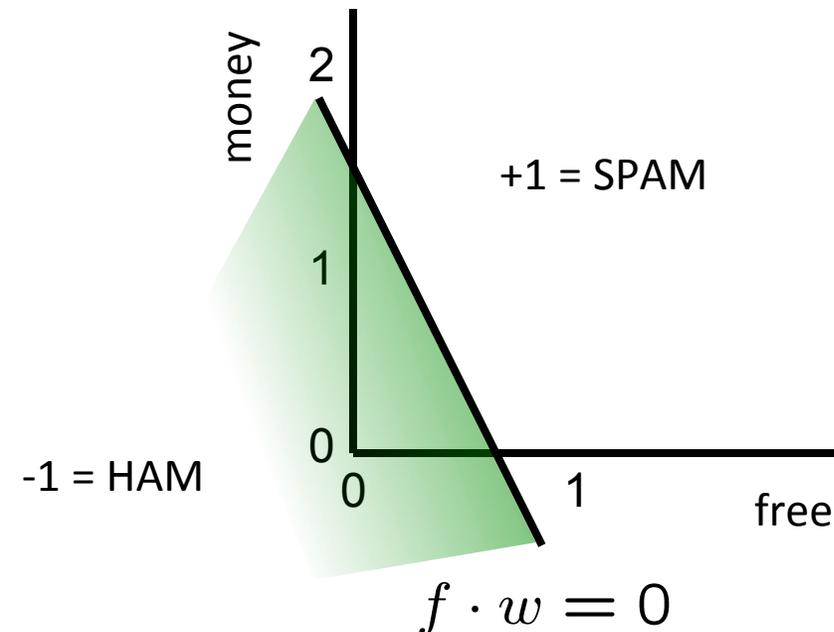
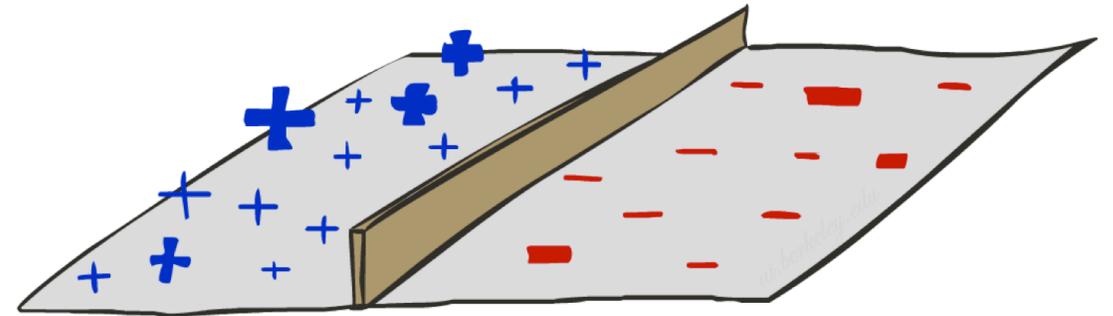
*Dot product  $w \cdot f$  positive  
means the positive class*

# Binary Decision Rule

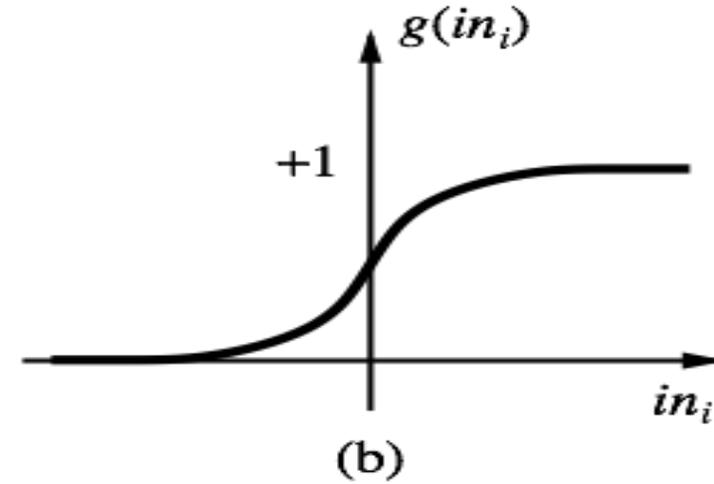
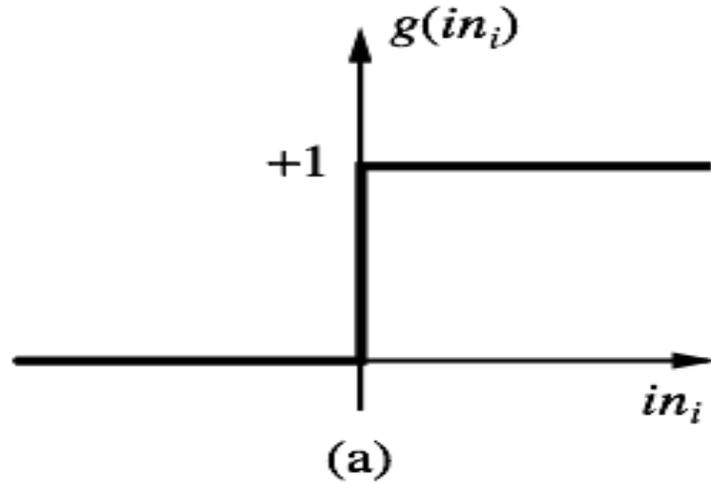
- In the space of feature vectors
  - Examples are points
  - Any weight vector is a hyperplane
  - Decision boundary:
    - One side corresponds to  $Y=+1$
    - Other corresponds to  $Y=-1$

$w$

BIAS	:	-3
free	:	4
money	:	2
...		



# Activation function

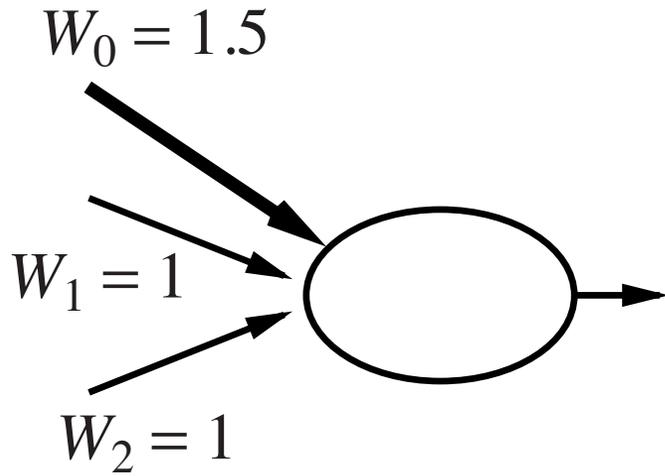


(a) is a step function or threshold function (a *perceptron*)

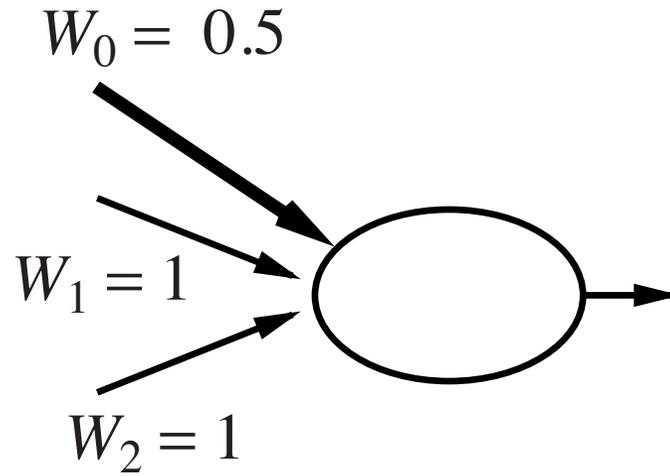
(b) is a sigmoid function  $1/(1 + e^{-x})$  (a *sigmoid perceptron*)

Changing the bias weight  $w_0$  moves the threshold location

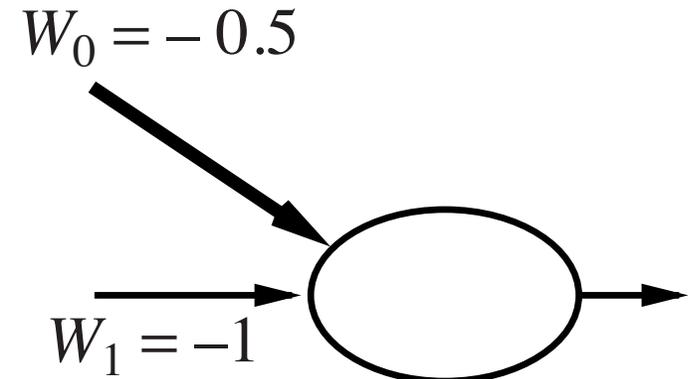
# Implementing logical functions



**AND**



**OR**

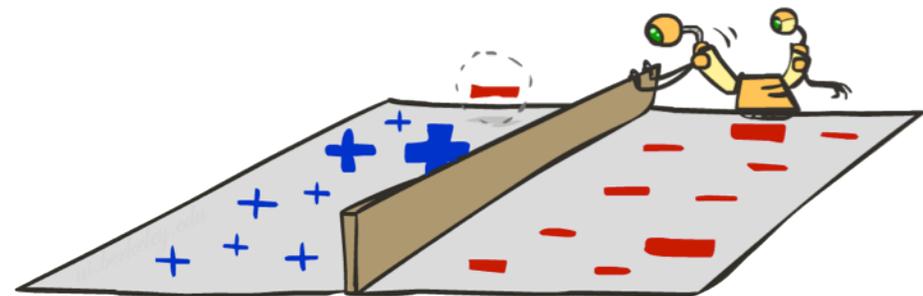
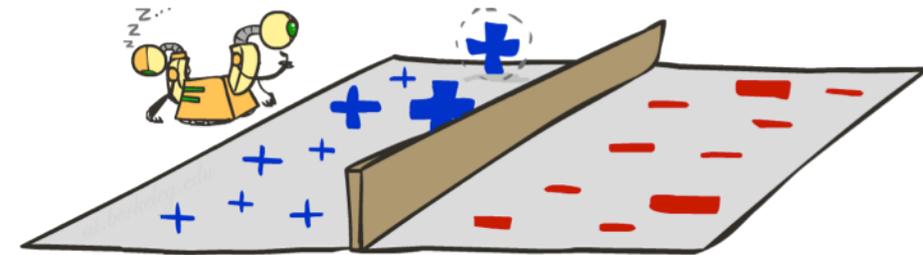
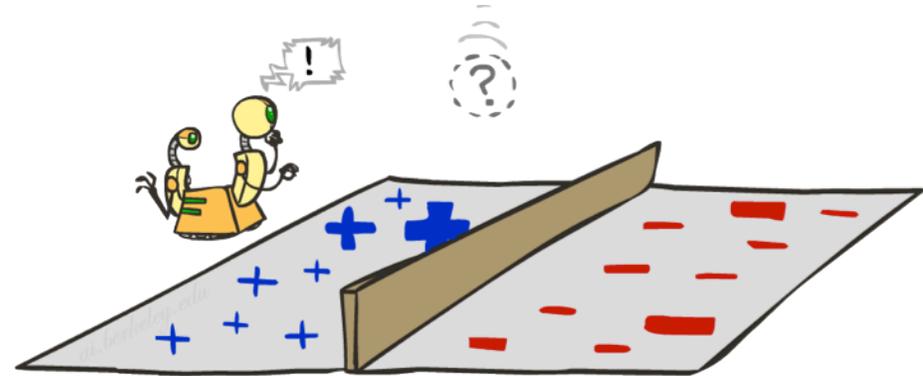


**NOT**

McCulloch and Pitts: every Boolean function can be implemented (sort of)

# Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
  - Classify with current weights
- If correct (i.e.,  $y=y^*$ ), no change!
- If wrong: adjust the weight vector

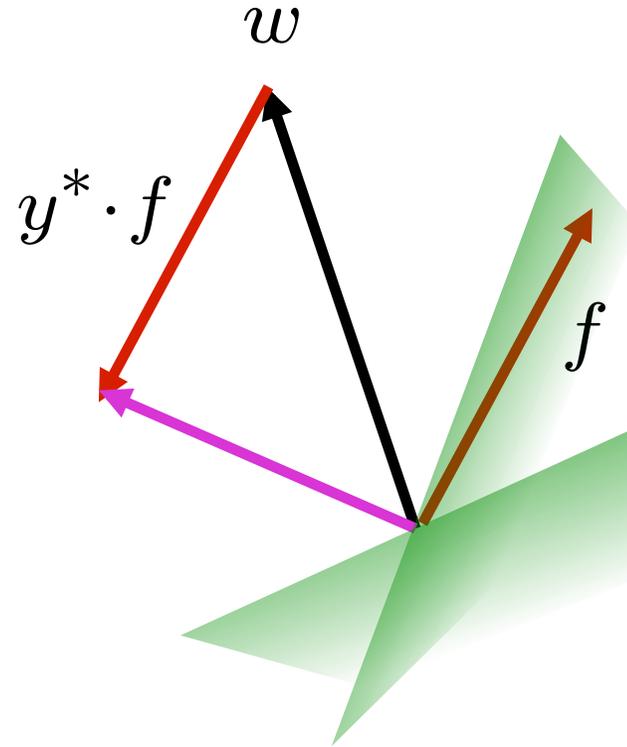


# Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
  - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

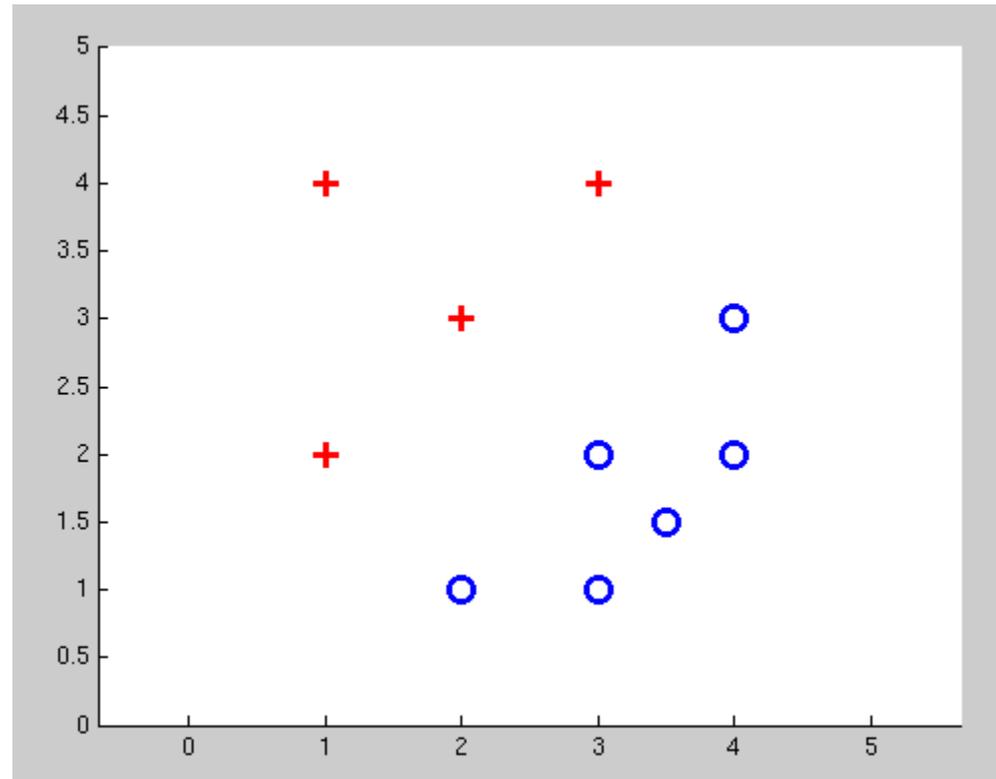
- If correct (i.e.,  $y=y^*$ ), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if  $y^*$  is -1.



$$w = w + y^* \cdot f$$

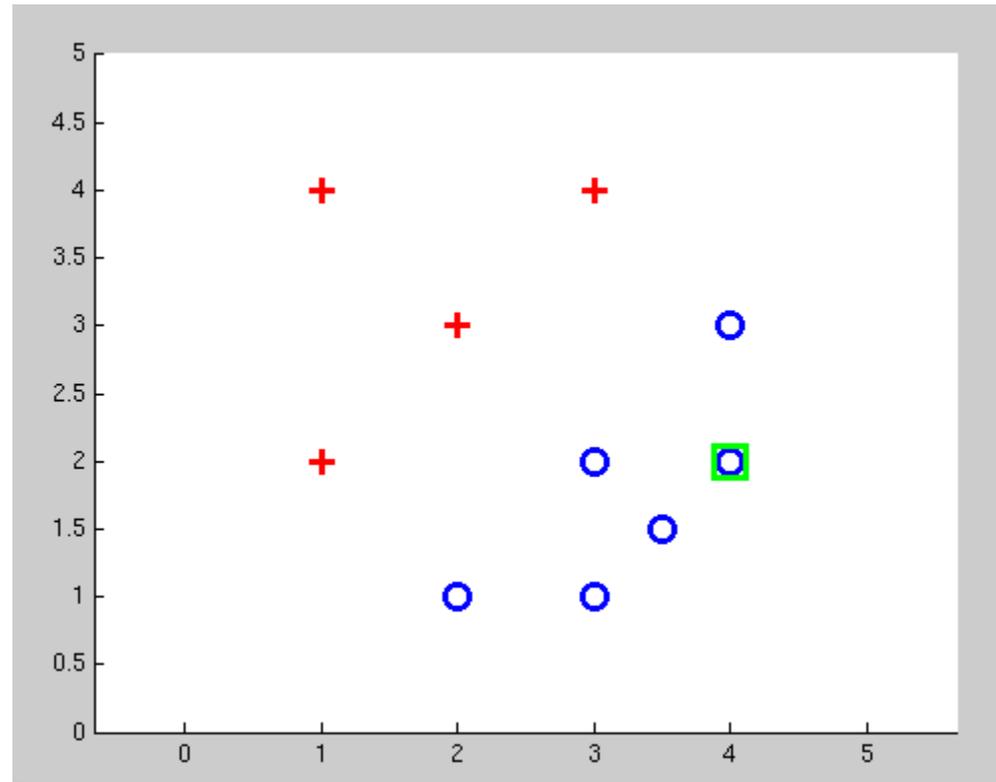
# Examples: Perceptron

- Separable Case



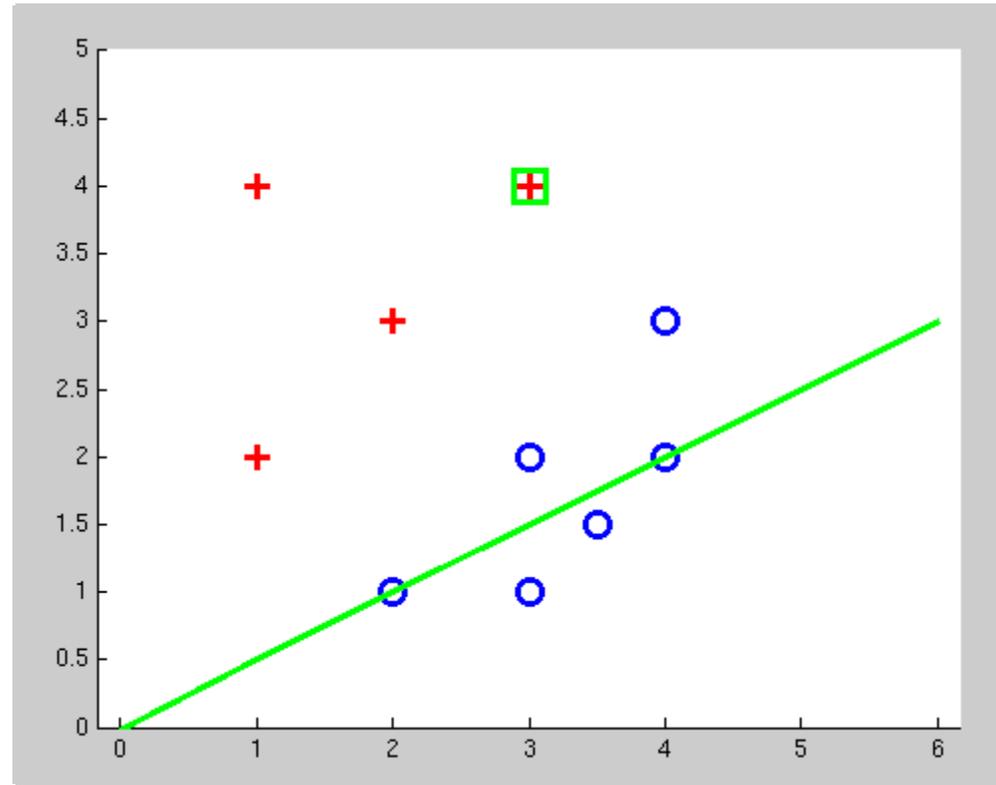
# Examples: Perceptron

- Separable Case



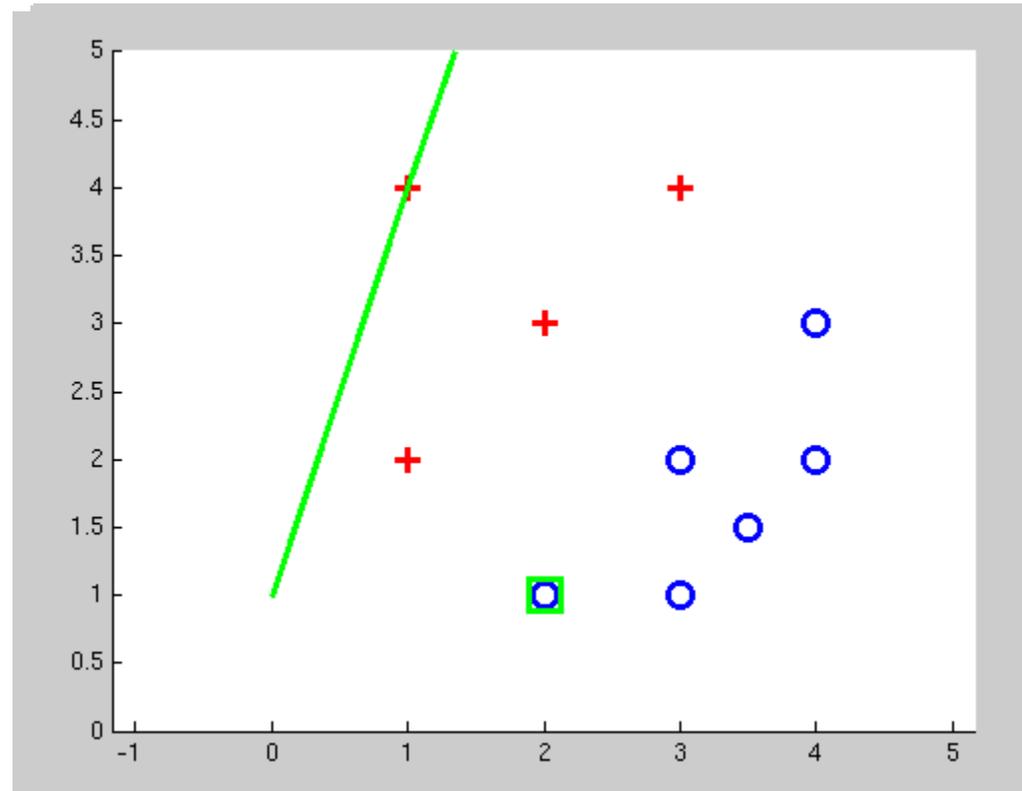
# Examples: Perceptron

- Separable Case



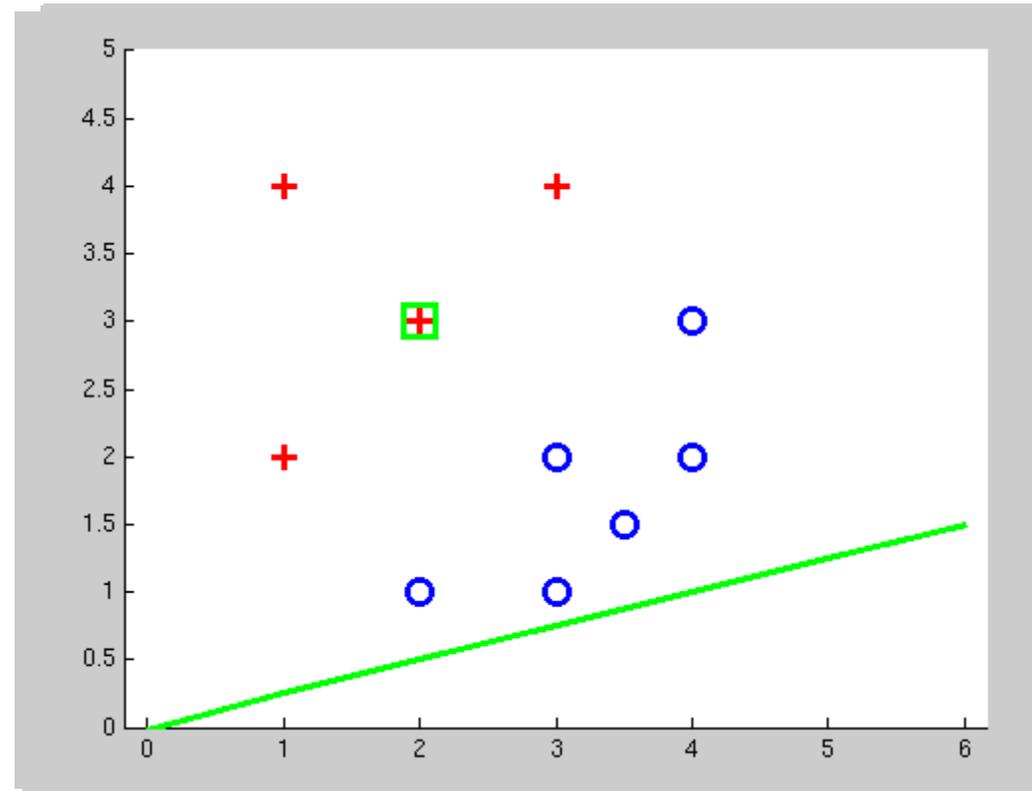
# Examples: Perceptron

- Separable Case



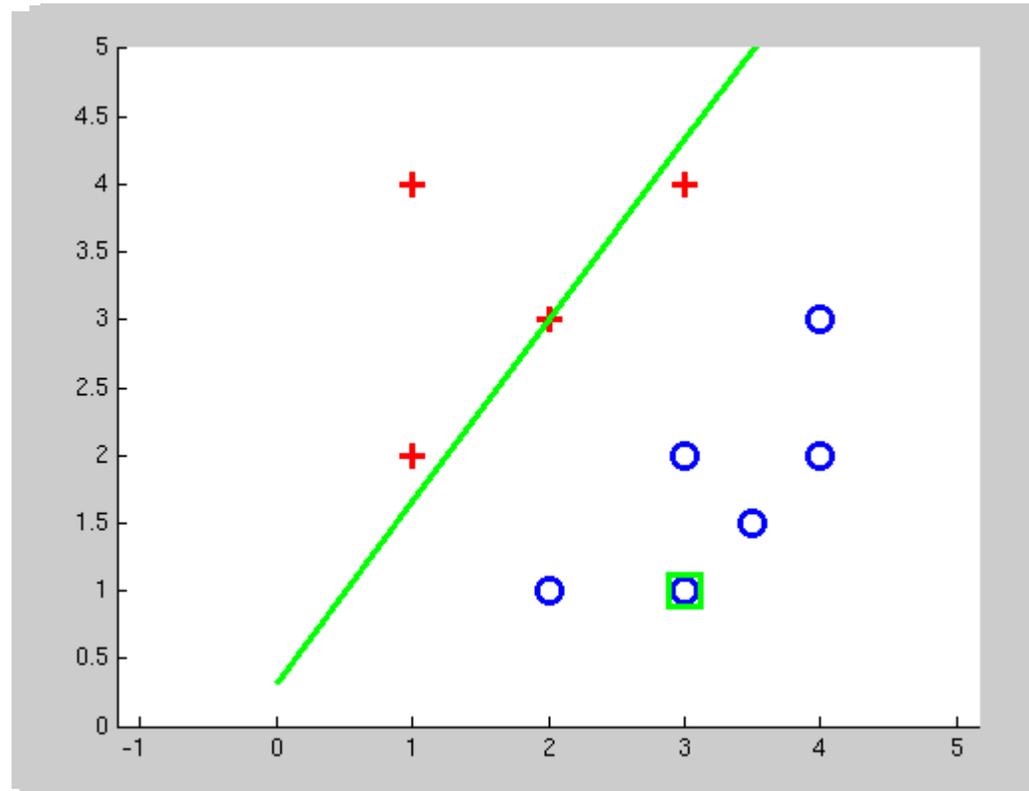
# Examples: Perceptron

- Separable Case



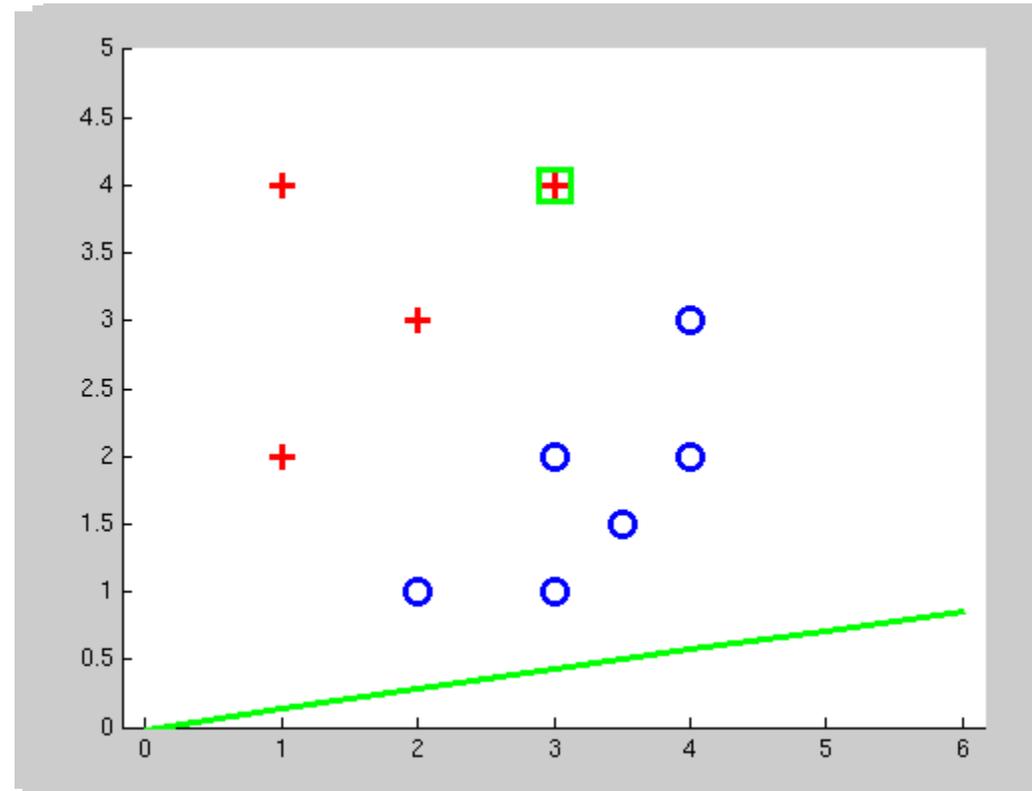
# Examples: Perceptron

- Separable Case



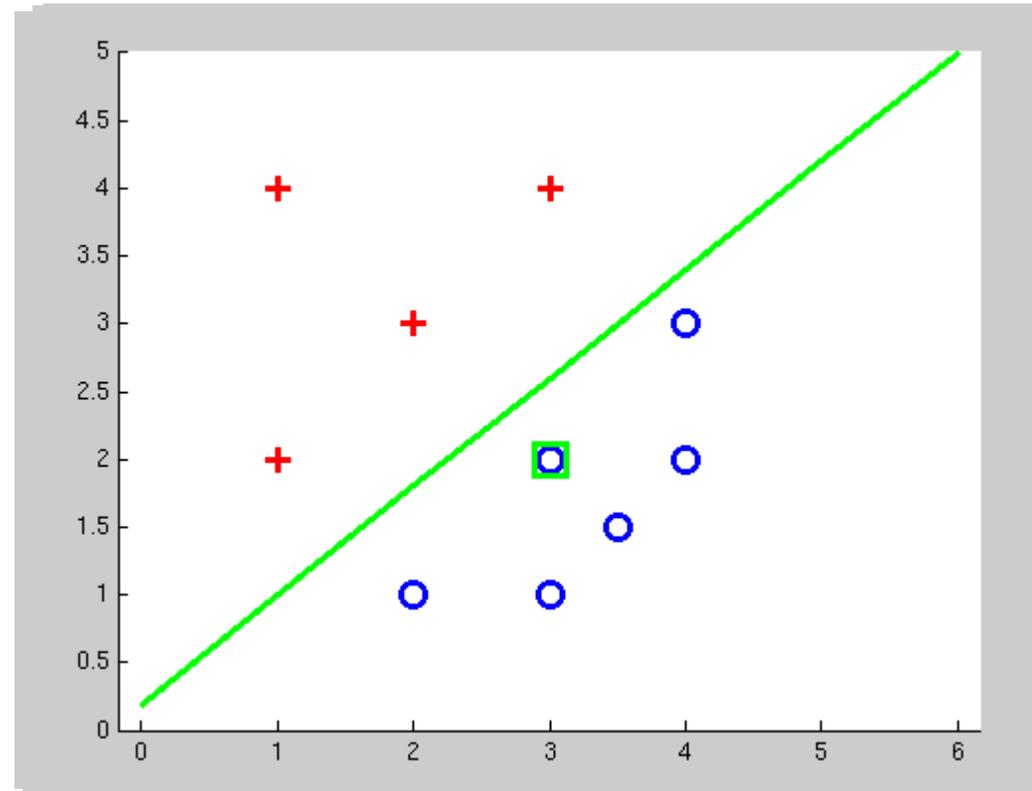
# Examples: Perceptron

- Separable Case



# Examples: Perceptron

- Separable Case



# Multiclass Decision Rule

- If we have multiple classes:
  - A weight vector for each class:

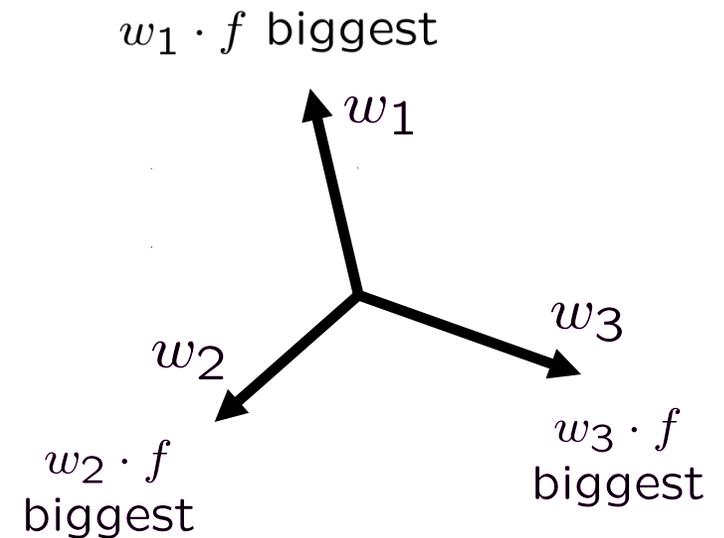
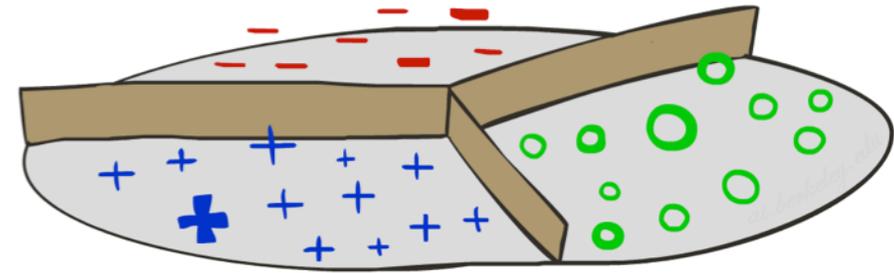
$$w_y$$

- Score (activation) of a class  $y$ :

$$w_y \cdot f(x)$$

- Prediction highest score wins

$$y = \arg \max_y w_y \cdot f(x)$$



*Binary = multiclass where the negative class has weight zero*

# Learning: Multiclass Perceptron

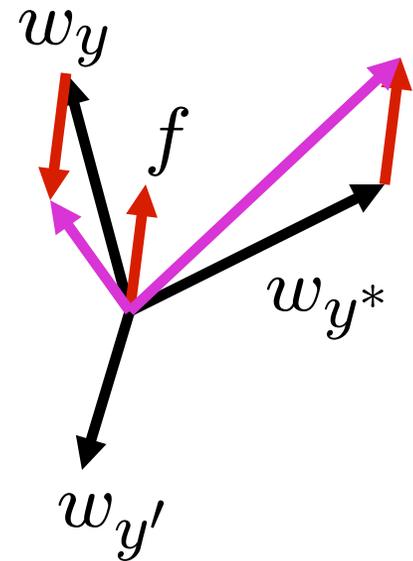
- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

$$y = \arg \max_y w_y \cdot f(x)$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$

$$w_{y^*} = w_{y^*} + f(x)$$



# Example: Multiclass Perceptron

“win the vote”

“win the election”

“win the game”

$w_{SPORTS}$

BIAS	:	1
win	:	0
game	:	0
vote	:	0
the	:	0
...		

$w_{POLITICS}$

BIAS	:	0
win	:	0
game	:	0
vote	:	0
the	:	0
...		

$w_{TECH}$

BIAS	:	0
win	:	0
game	:	0
vote	:	0
the	:	0
...		

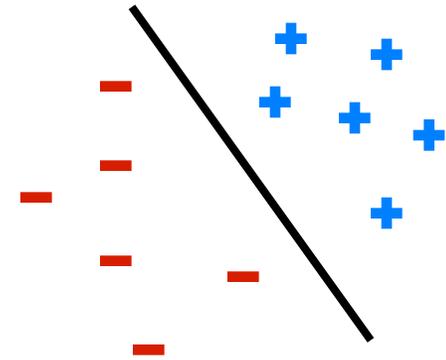
# Properties of Perceptrons

- Separability: true if some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

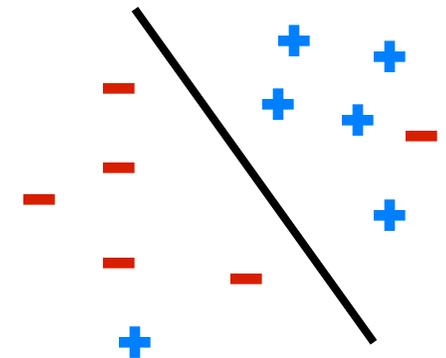
$$\text{mistakes} < \frac{k}{\delta^2}$$

k is the number of features  
 $\delta$  is the size of the margin

Separable



Non-Separable



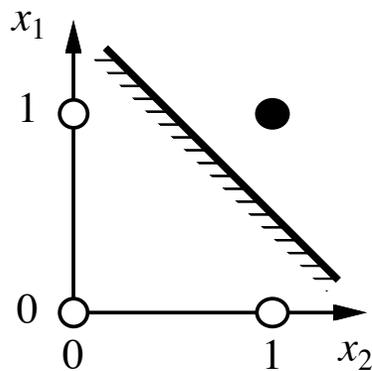
# Expressiveness of Perceptron

Consider a perceptron with  $g = \text{step function}$  (Rosenblatt, 1957, 1960)

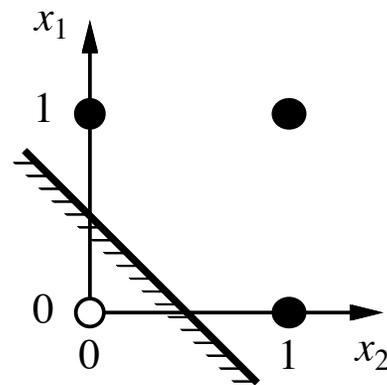
Can represent AND, OR, NOT, majority, etc., but not XOR

Represents a linear separator in input space:

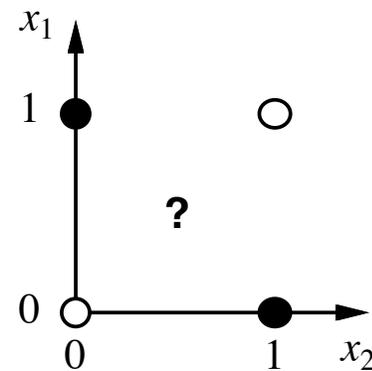
$$\sum_j W_j x_j > 0 \quad \text{or} \quad \mathbf{W} \cdot \mathbf{x} > 0$$



(a)  $x_1$  **and**  $x_2$



(b)  $x_1$  **or**  $x_2$

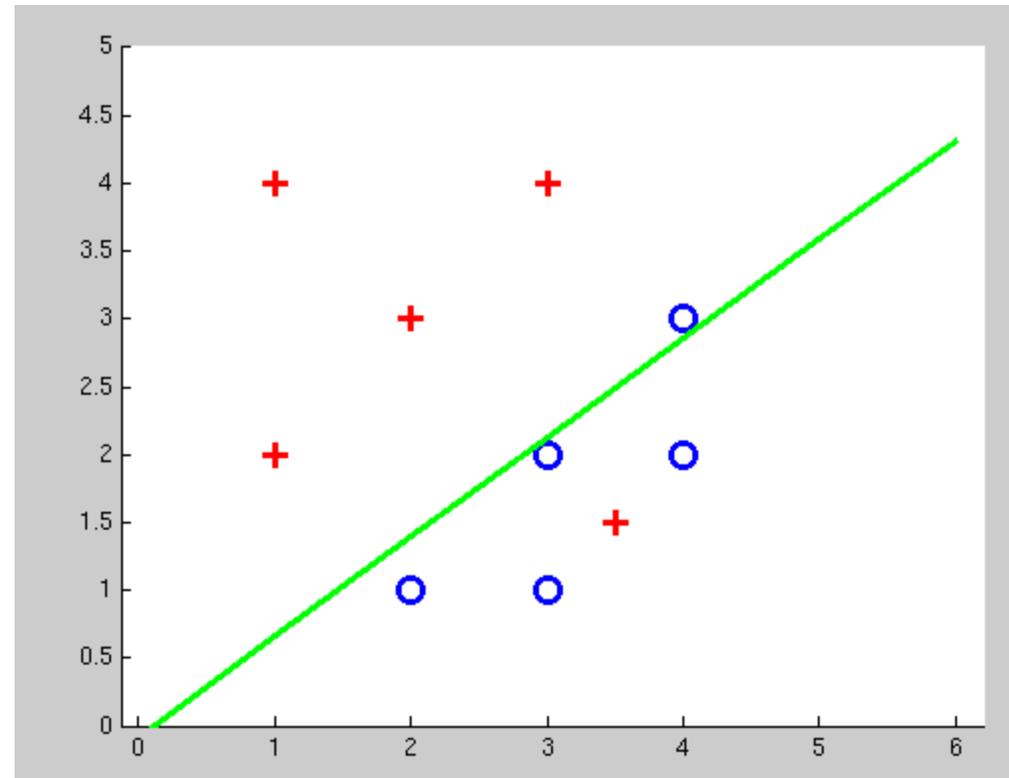


(c)  $x_1$  **xor**  $x_2$

Minsky & Papert (1969) pricked the neural network balloon

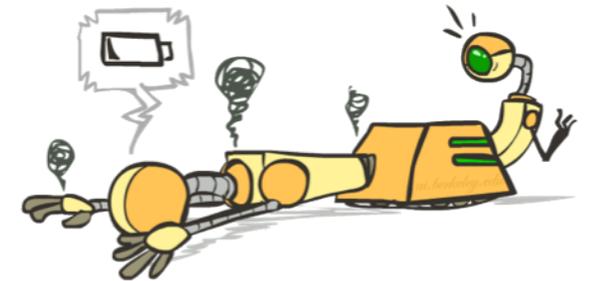
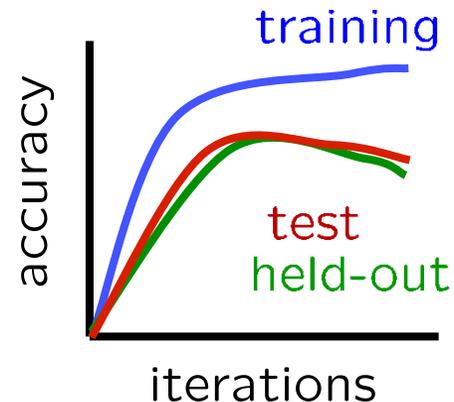
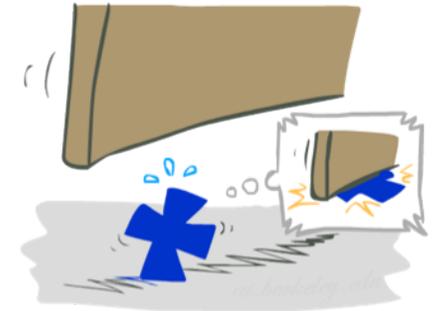
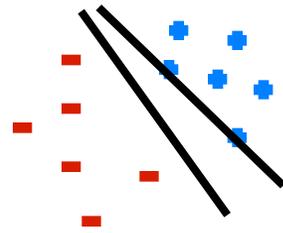
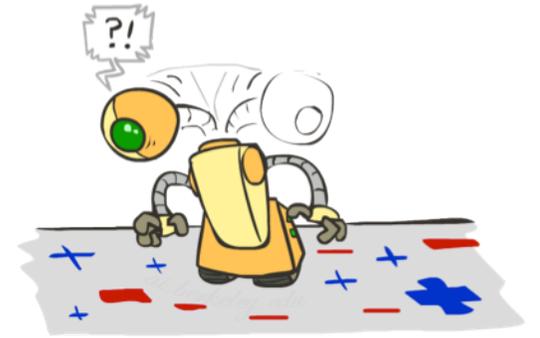
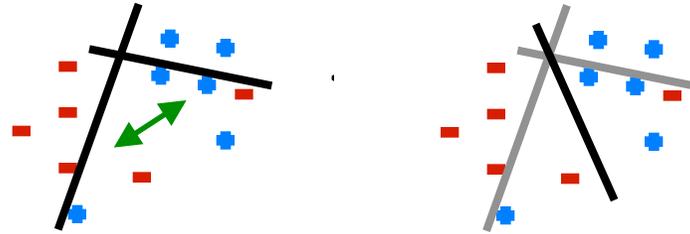
# Examples: Perceptron

- Non-Separable Case



# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)
- Mediocre generalization: finds a "barely" separating solution
- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining is a kind of overfitting



# Fixing the Perceptron

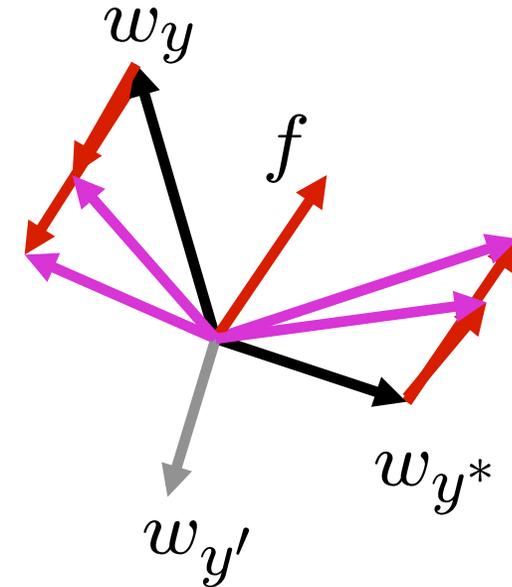
- Idea: adjust the weight update to mitigate these effects
- MIRA\*: choose an update size that fixes the current mistake...
- ... but, minimizes the change to  $w$

$$\min_w \frac{1}{2} \sum_y ||w_y - w'_y||^2$$

$$w_{y^*} \cdot f(x) \geq w_y \cdot f(x) + 1$$

- The +1 helps to generalize

\* Margin Infused Relaxed Algorithm



Guessed  $y$  instead of  $y^*$  on example  $x$  with features  $f(x)$

$$w_y = w'_y - \tau f(x)$$
$$w_{y^*} = w'_{y^*} + \tau f(x)$$

# Minimum Correcting Update

$$\min_w \frac{1}{2} \sum_y \|w_y - w'_y\|^2$$

$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$

$$\min_{\tau} \|\tau f\|^2$$

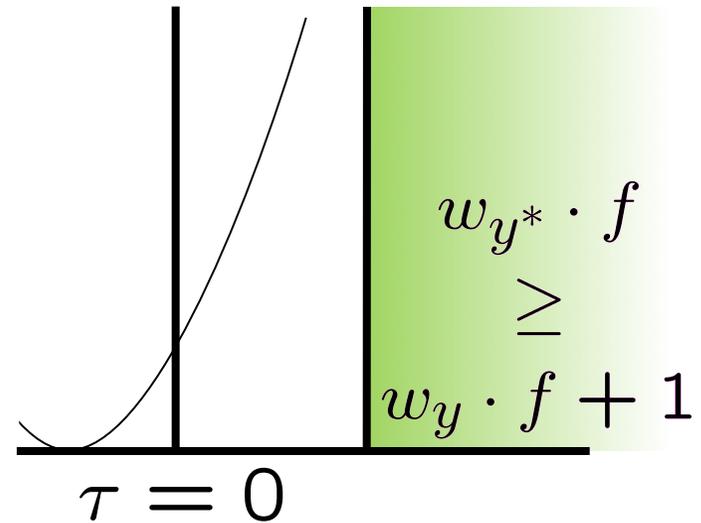
$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$

$$(w'_{y^*} + \tau f) \cdot f = (w'_y - \tau f) \cdot f + 1$$

$$\tau = \frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}$$

$$w_y = w'_y - \tau f(x)$$

$$w_{y^*} = w'_{y^*} + \tau f(x)$$



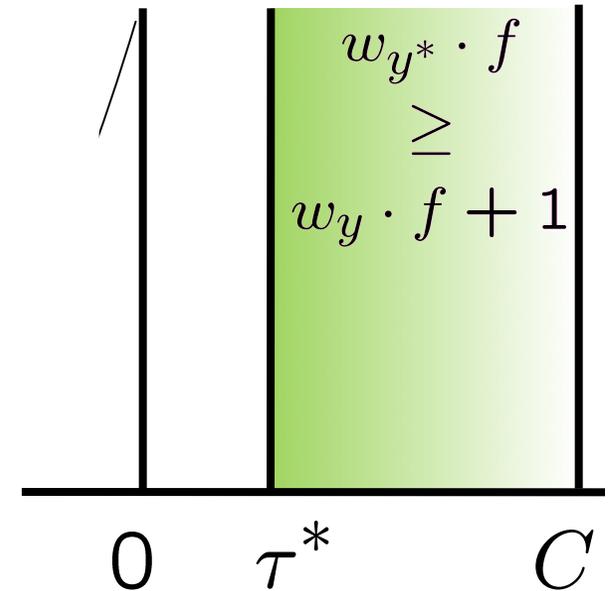
min not ?=0, or would not have made an error, so min will be where equality holds

# Maximum Step Size

- In practice, it's also bad to make updates that are too large
  - Example may be labeled incorrectly
  - You may not have enough features
  - Solution: cap the maximum possible value of  $\tau$  with some constant  $C$

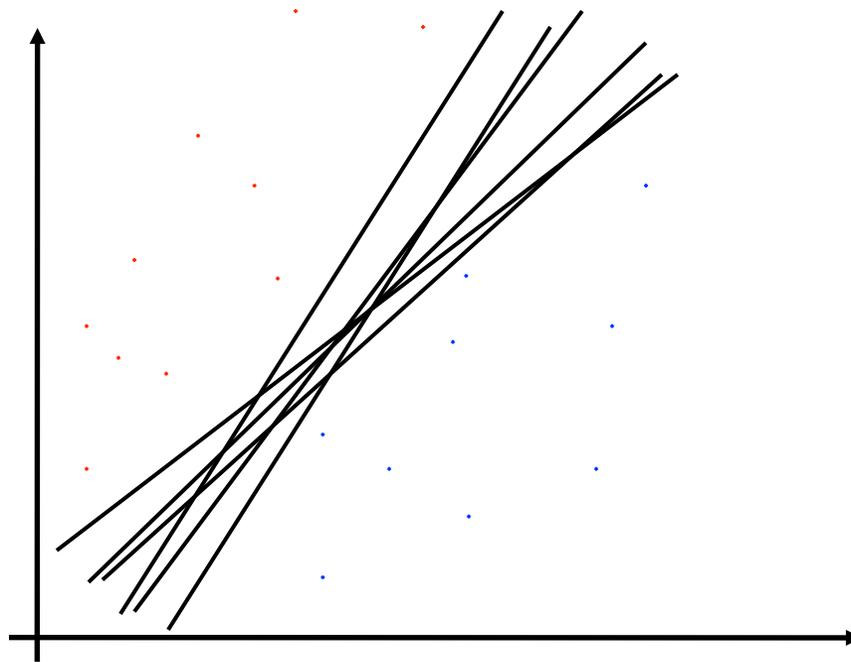
$$\tau^* = \min \left( \frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}, C \right)$$

- Corresponds to an optimization that assumes non-separable data
- Usually converges faster than perceptron
- Usually better, especially on noisy data



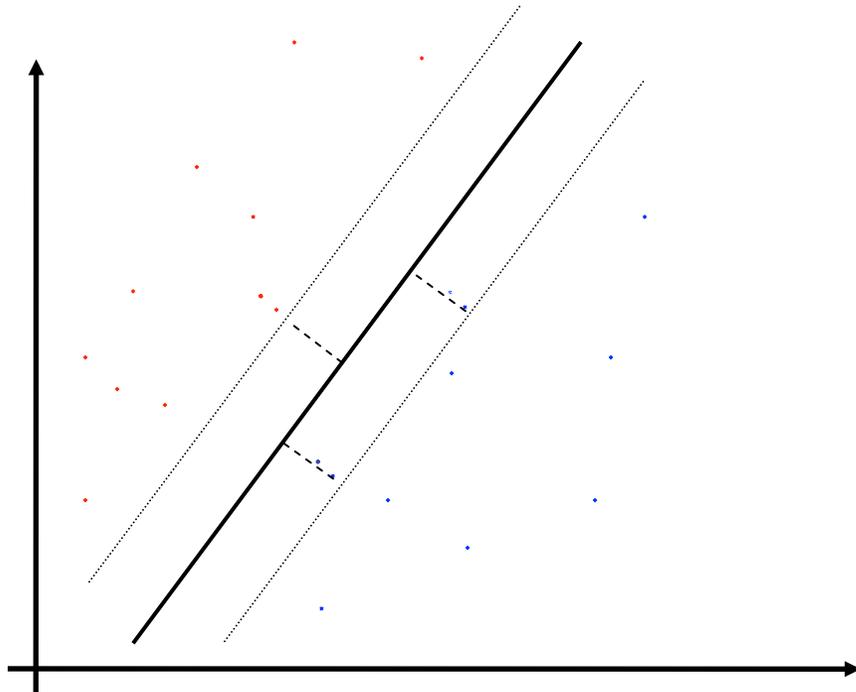
# Linear Separators

- Which of these linear separators is optimal?



# Support Vector Machines

- **Maximizing the margin:** good according to intuition, theory, practice
- Only **support vectors** matter; other training examples are ignorable
- Support vector machines (SVMs) find the separator with max margin
- Basically, SVMs are MIRA where you optimize over all examples at once



MIRA

$$\min_w \frac{1}{2} \|w - w'\|^2$$

$$w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

SVM

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

# Classification: Comparison

- Naïve Bayes

- Builds a model training data
- Gives prediction probabilities
- Strong assumptions about feature independence
- One pass through data (counting)

- Perceptrons / MIRA:

- Makes less assumptions about data
- Mistake-driven learning
- Multiple passes through data (prediction)
- Often more accurate