

CS 188: Artificial Intelligence

Machine Learning

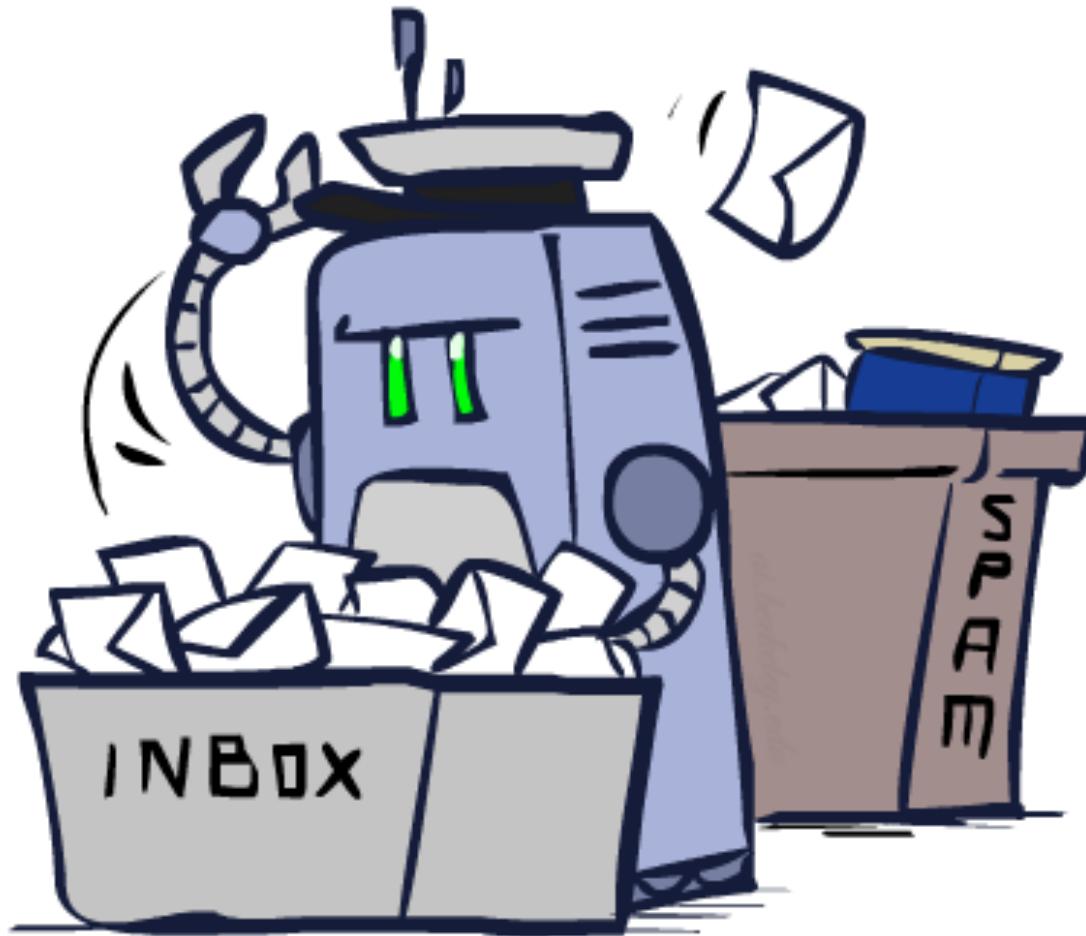


Summer 2024: Eve Fleisig & Evgeny Pobachienko

Machine Learning

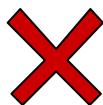
- Up until now: how use a model to make optimal decisions
- Machine learning: how to acquire a model from data / experience
 - Learning parameters (e.g. probabilities)
 - Learning structure (e.g. BN graphs)
 - Learning hidden concepts (e.g. clustering, neural nets)
- Today: model-based classification with Naive Bayes

Classification



Example: Spam Filter

- Input: an 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: The attributes used to make the ham / spam decision
 - Words: FREE!
 - Text Patterns: \$dd, CAPS
 - Non-text: SenderInContacts, WidelyBroadcast
 - ...



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.

Example: 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: $(6,8)=\text{ON}$
 - Shape Patterns: NumComponents, AspectRatio, NumLoops
 - ...
 - Features are increasingly induced rather than crafted

A handwritten digit '0'.

0

A handwritten digit '1'.

1

A handwritten digit '2'.

2

A handwritten digit '1'.

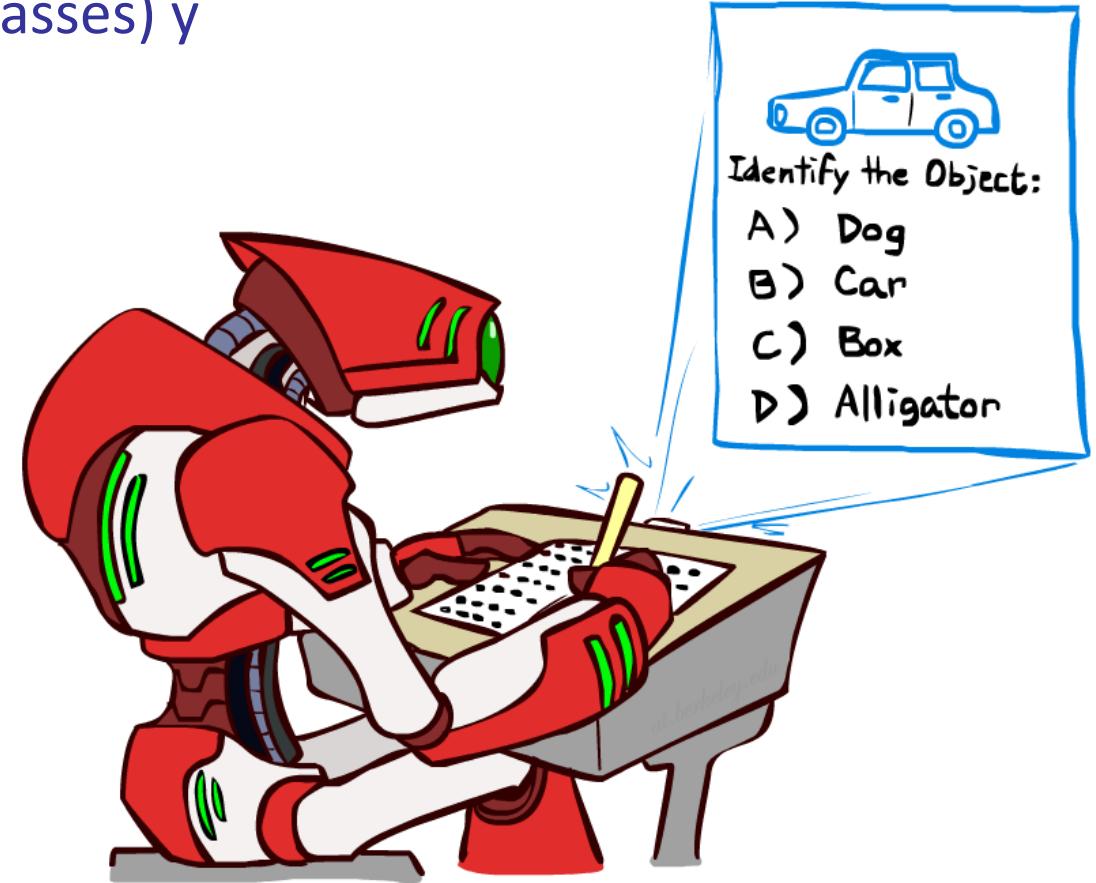
1

A handwritten digit that looks like a '4'.

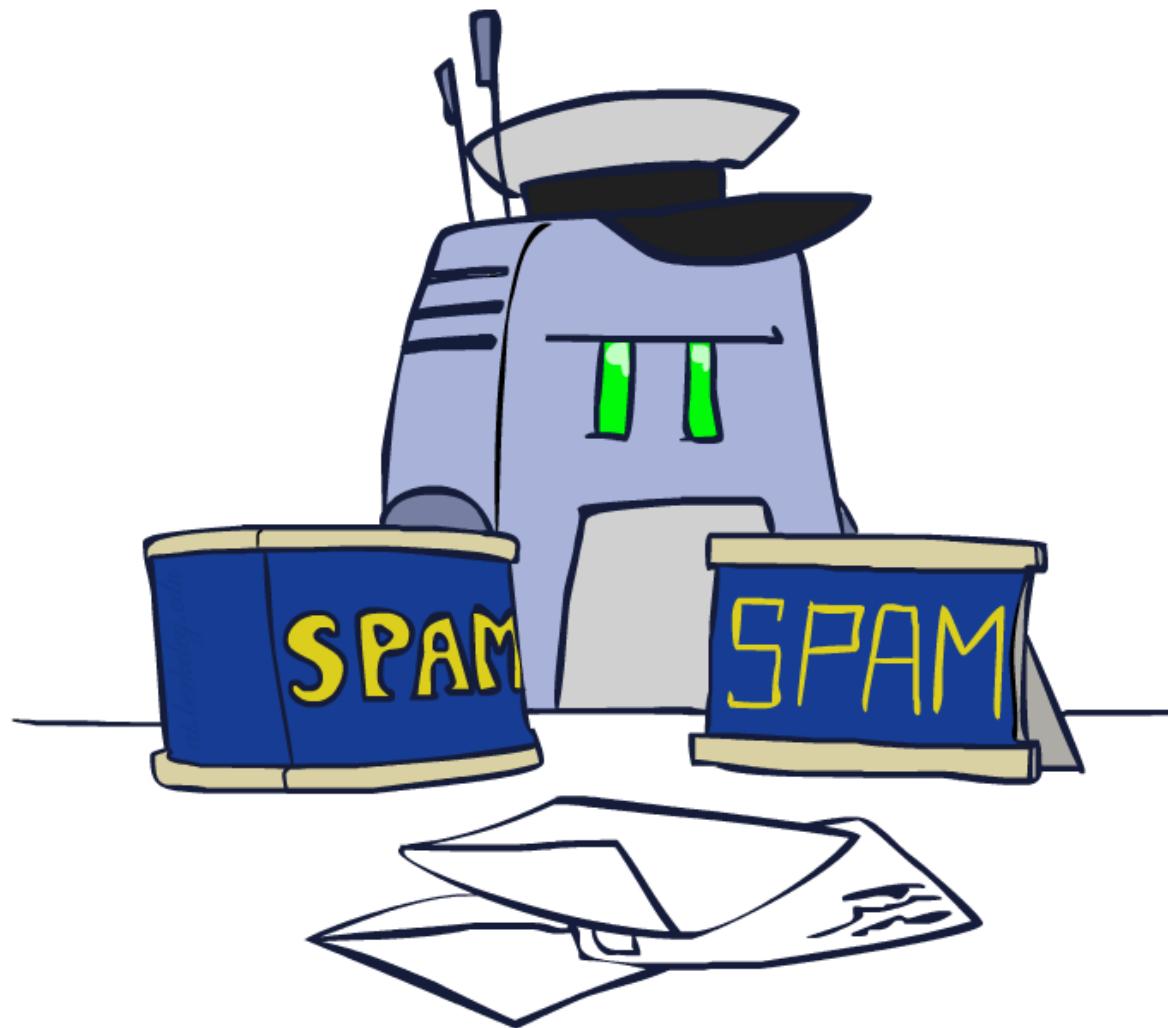
??

Other Classification Tasks

- Classification: given inputs x , predict labels (classes) y
 - Binary classification (2 classes)
 - Multiclass classification (>2 classes)
- Examples:
 - Medical diagnosis (input: symptoms, classes: diseases)
 - Fraud detection (input: account activity, classes: fraud / no fraud)
 - Automatic essay grading (input: document, classes: grades)
 - Customer service email routing
 - Review sentiment
 - Language ID
 - ... many more
- Classification is an important commercial technology!



Model-Based Classification



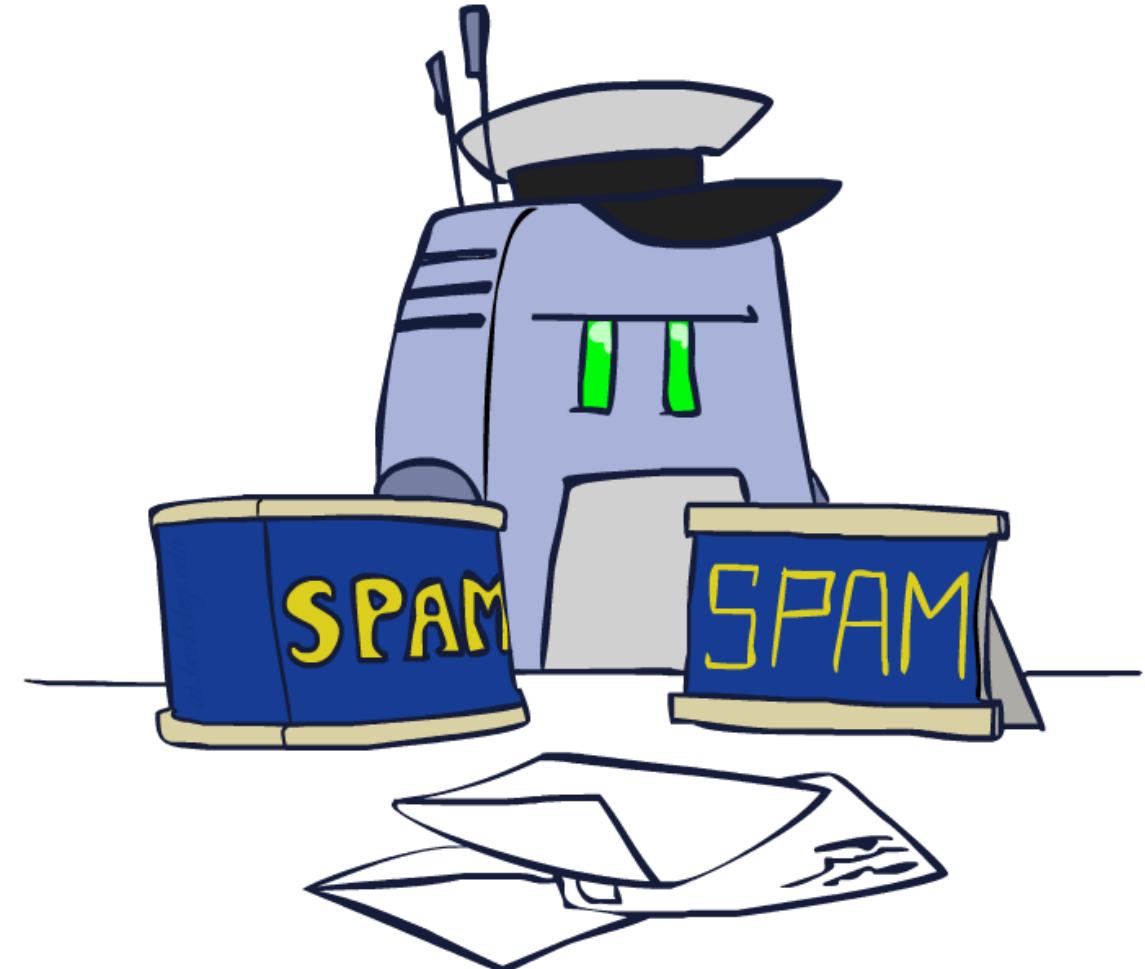
Model-Based Classification

- Model-based approach

- Build a model (e.g. Bayes' net) where both the output label and input 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

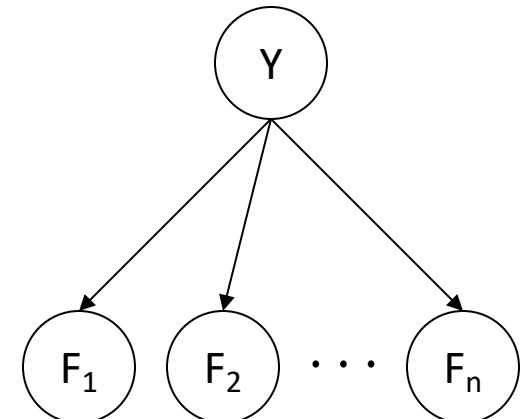
- Simple digit recognition version:

- 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.



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

- 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?



General Naïve Bayes

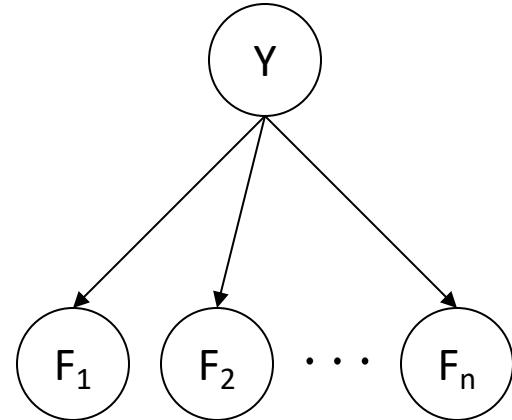
- A general Naive Bayes model:

$|Y|$ parameters

$$P(Y, F_1 \dots F_n) = P(Y) \prod_i P(F_i | Y)$$

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

$n \times |F| \times |Y|$
parameters



- 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

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

- Step 2: sum to get probability of evidence
- Step 3: normalize by dividing Step 1 by Step 2

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

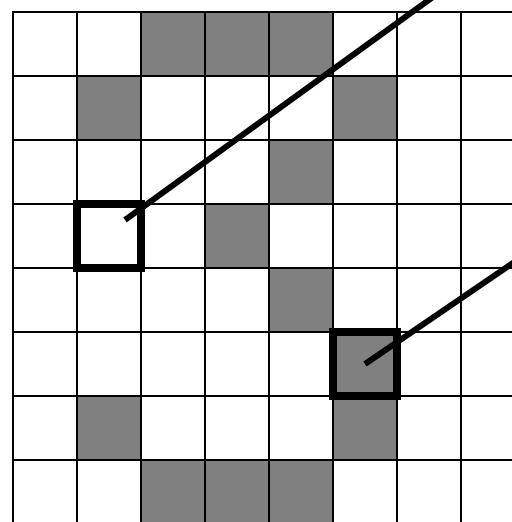
General Naïve Bayes

- 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 $P(Y|F_1 \dots F_n)$
 - 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 θ
 - 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



$P(F_{3,1} = on|Y) \quad P(F_{5,5} = on|Y)$

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

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^{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
...

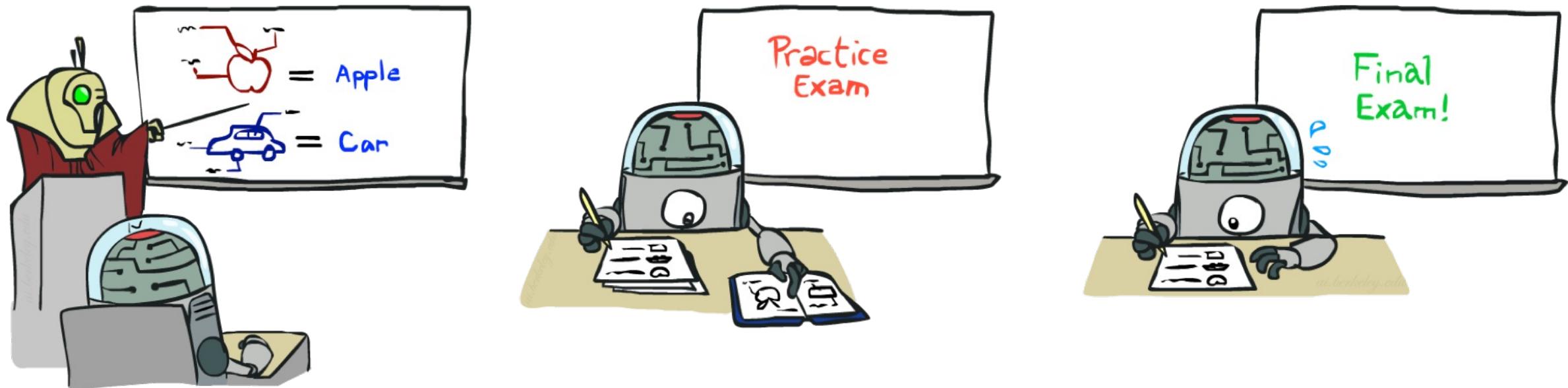
- Where do these tables come from?

Spam Example

Word	P(w spam)	P(w ham)	Tot Spam	Tot Ham
(prior)	0.33333	0.66666	-1.1	-0.4

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

Training and Testing

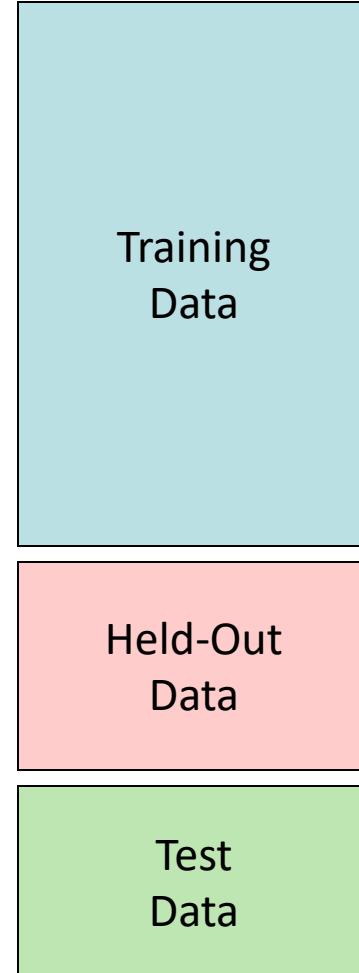


Empirical Risk Minimization

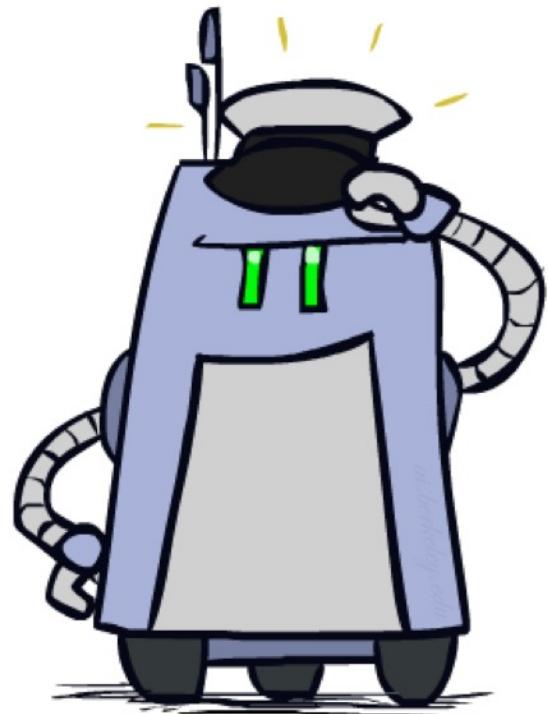
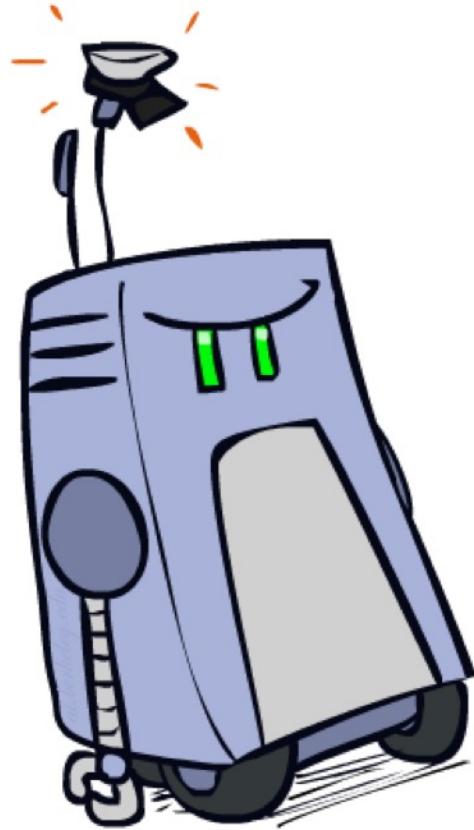
- Empirical risk minimization
 - Basic principle of machine learning
 - We want the model (classifier, etc) that does best on the true test distribution
 - Don't know the true distribution so pick the best model on our actual training set
 - Finding "the best" model on the training set is phrased as an optimization problem
- Main worry: overfitting to the training set
 - Better with more training data (less sampling variance, training more like test)
 - Better if we limit the complexity of our hypotheses (regularization and/or small hypothesis spaces)

Important Concepts

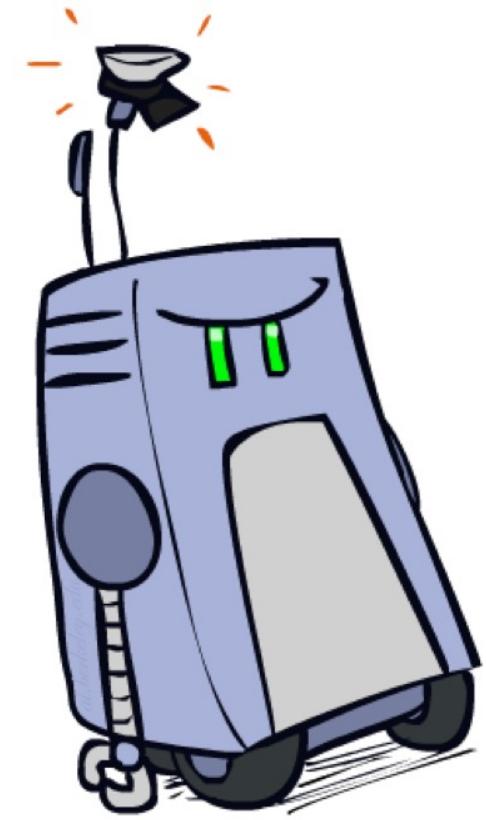
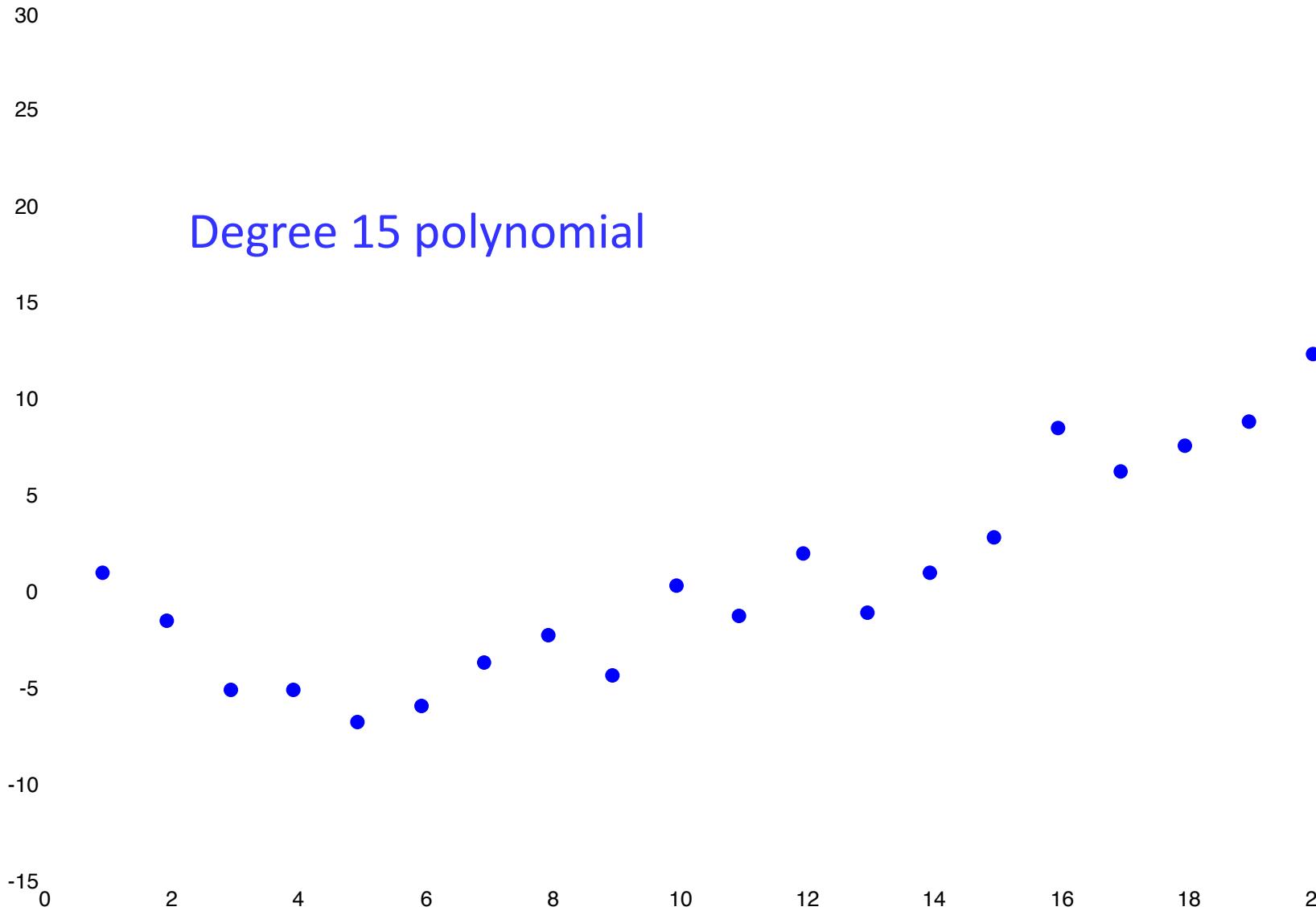
- Data: labeled instances (e.g. emails marked spam/ham)
 - Training set
 - Held out set (“development” or “validation” set)
 - Test set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
 - Learn parameters (e.g. model probabilities) on training set
 - (Tune hyperparameters on held-out set)
 - Compute accuracy of test set
 - Very important: never “peek” at the test set!
- Evaluation (many metrics possible, e.g. accuracy)
 - 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
 - We’ll investigate overfitting and generalization formally in a few lectures



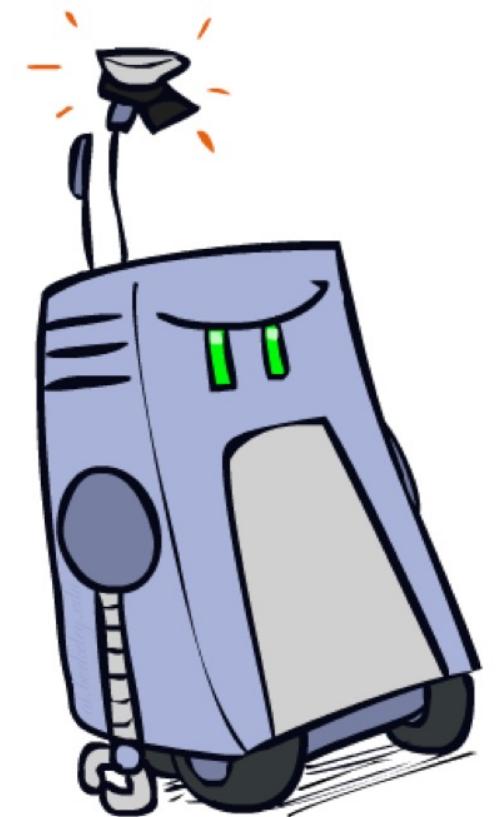
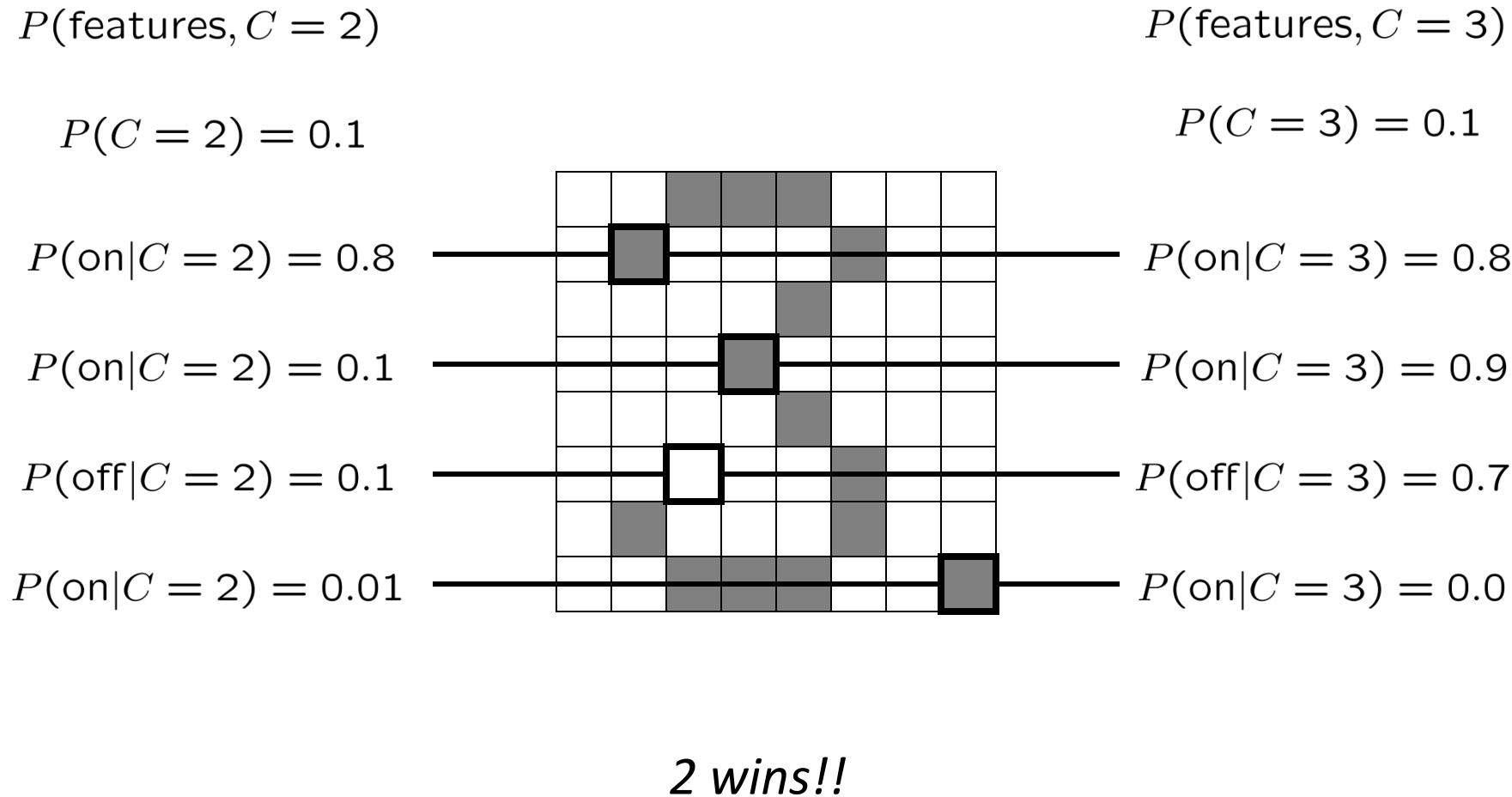
Generalization and Overfitting



Overfitting



Example: Overfitting



Example: Overfitting

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

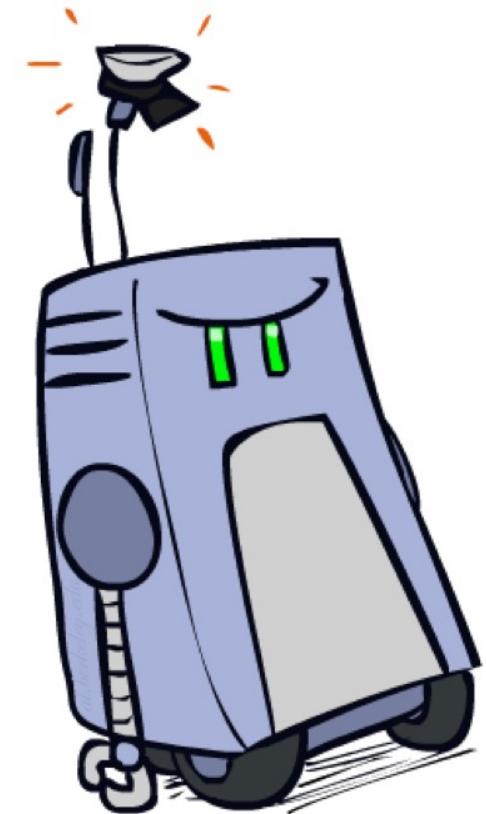
$$\frac{P(W|\text{ham})}{P(W|\text{spam})}$$

```
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  
...  
...
```

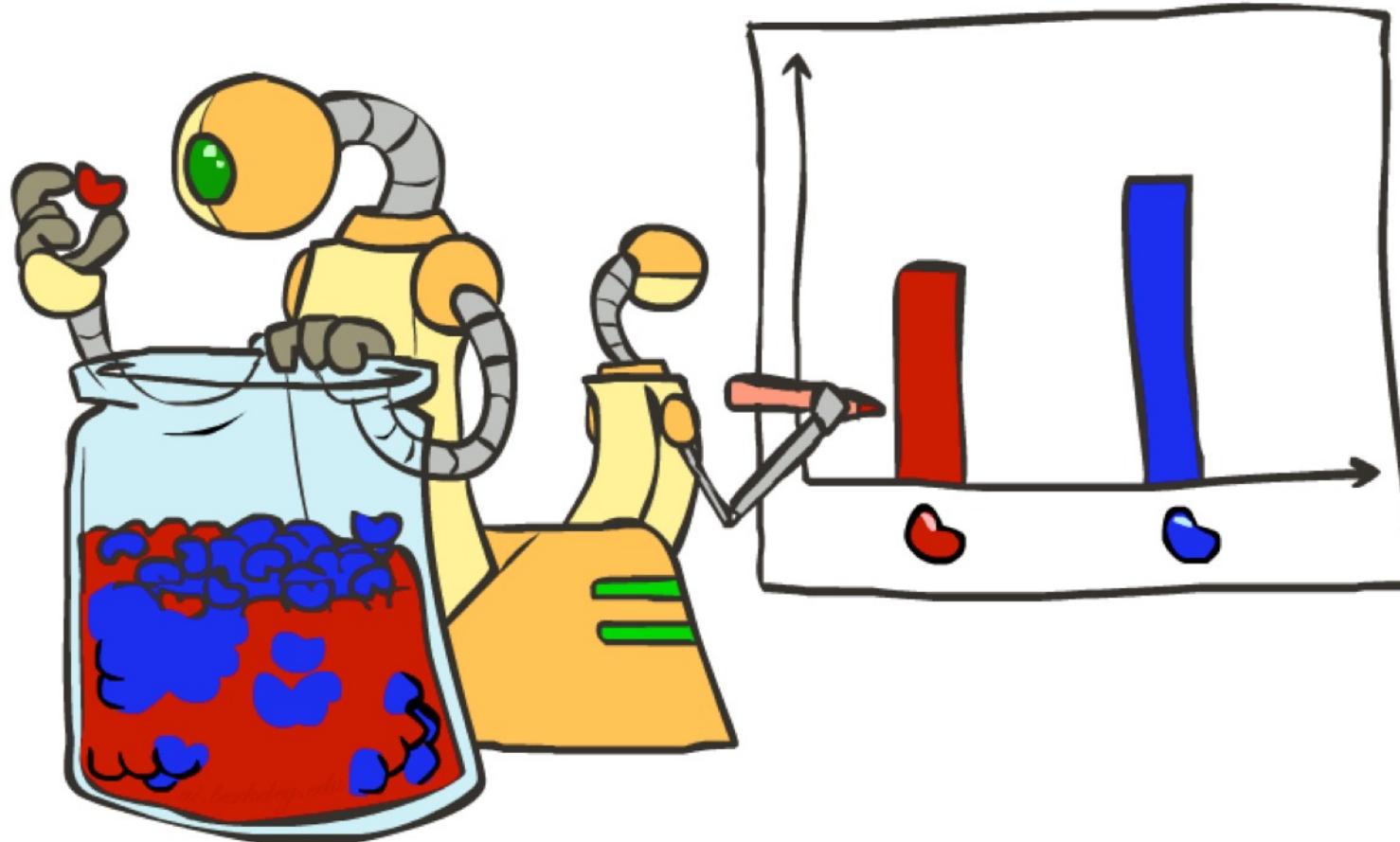
What went wrong here?



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 (e.g. document ID)
 - 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



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:

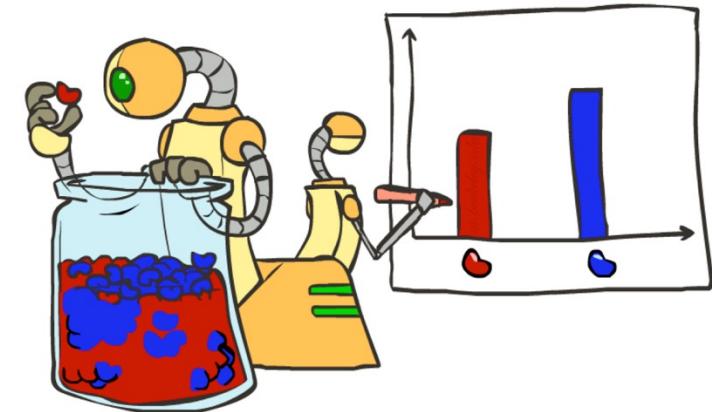
$$P_{\text{ML}}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

r r b

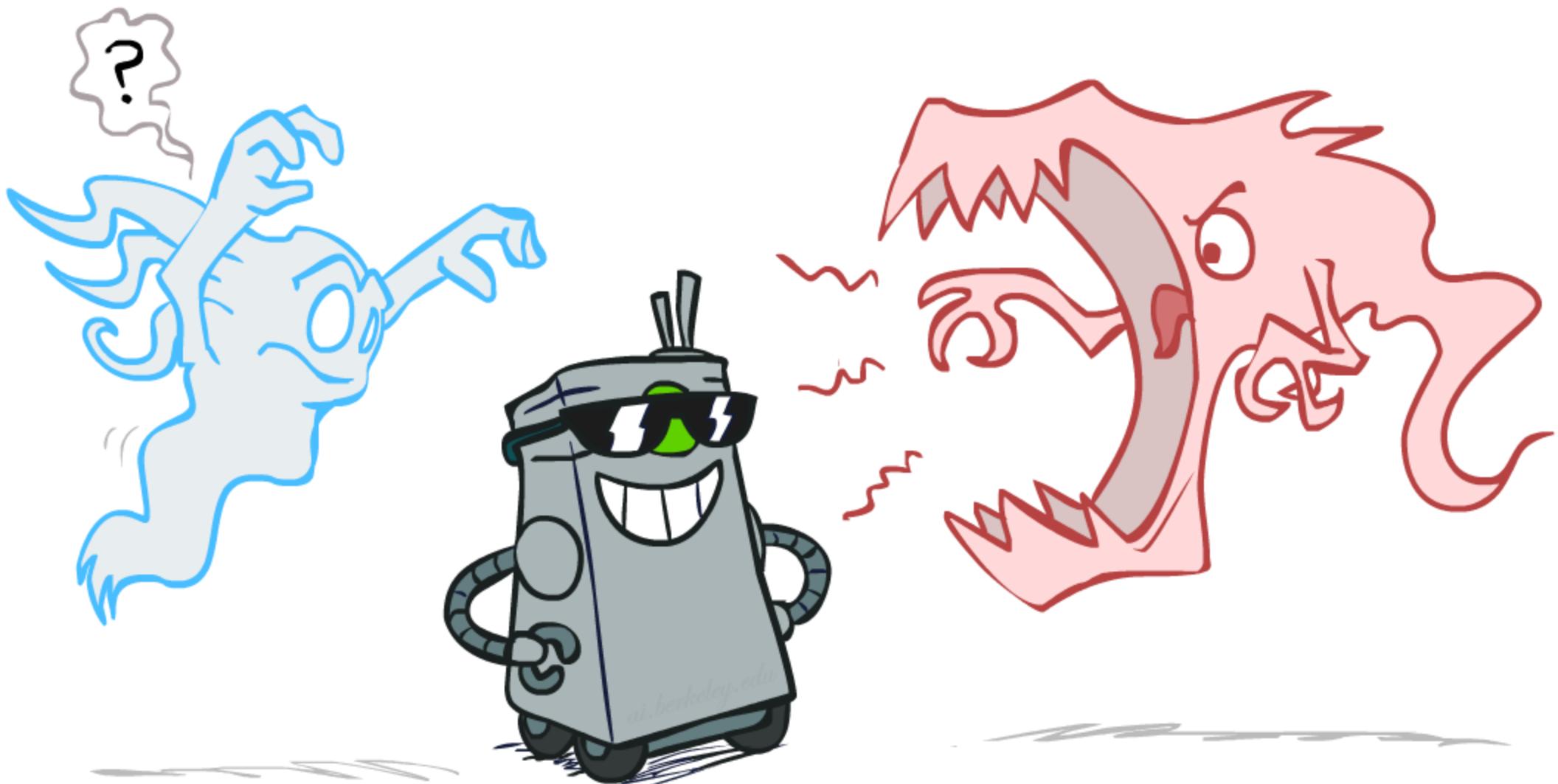
$$P_{\text{ML}}(\text{r}) = 2/3$$

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

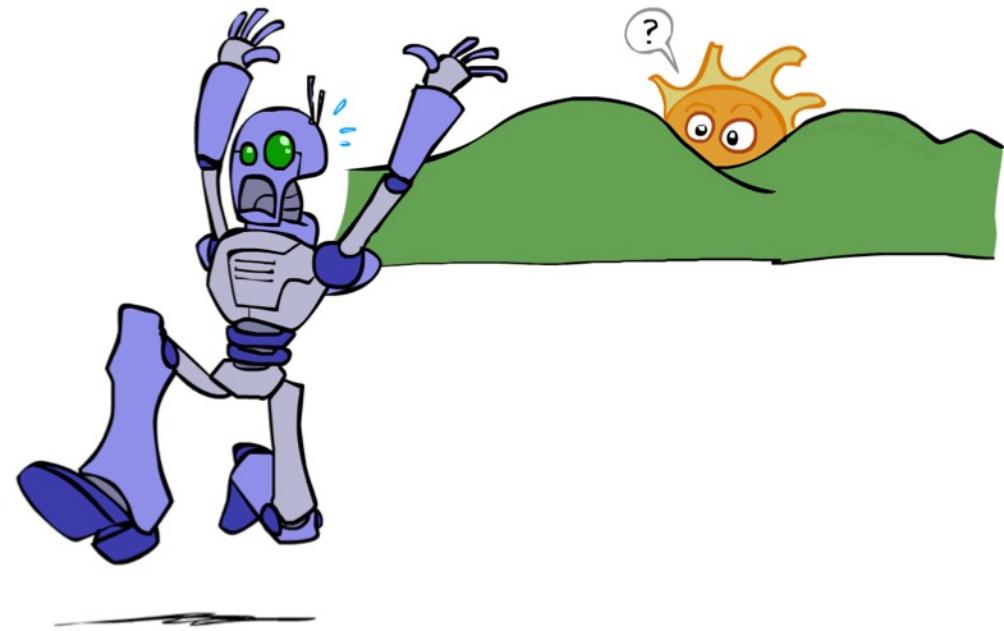
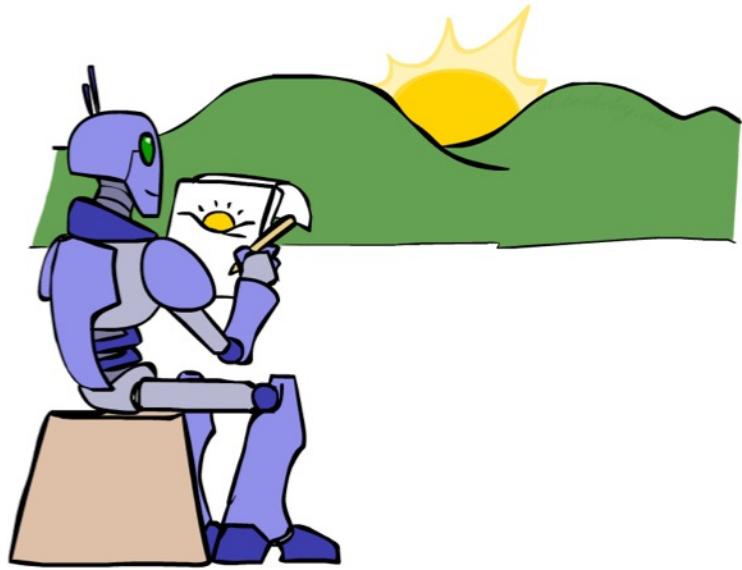
$$L(x, \theta) = \prod_i P_\theta(x_i)$$



Smoothing



Unseen Events



Laplace Smoothing

- Laplace's estimate:

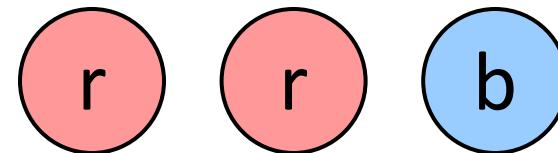
- Pretend you saw every outcome once more than you actually did

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

$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

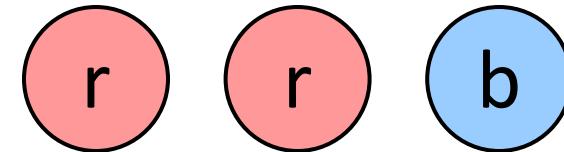


- Can derive this estimate with *Dirichlet priors* (see cs281a)

Laplace Smoothing

- Laplace's estimate (extended)– “Add-k smoothing”:
 - Pretend you saw every outcome k extra times

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



$$P_{LAP,0}(X) =$$

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

$$P_{LAP,1}(X) =$$

- Laplace for conditionals:
 - Smooth each condition independently:

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

$$P_{LAP,100}(X) =$$

[Bonus] 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 α is 0? 1?
- For even better ways to estimate parameters, as well as details of the math, see cs281a, cs288

Real NB: Smoothing

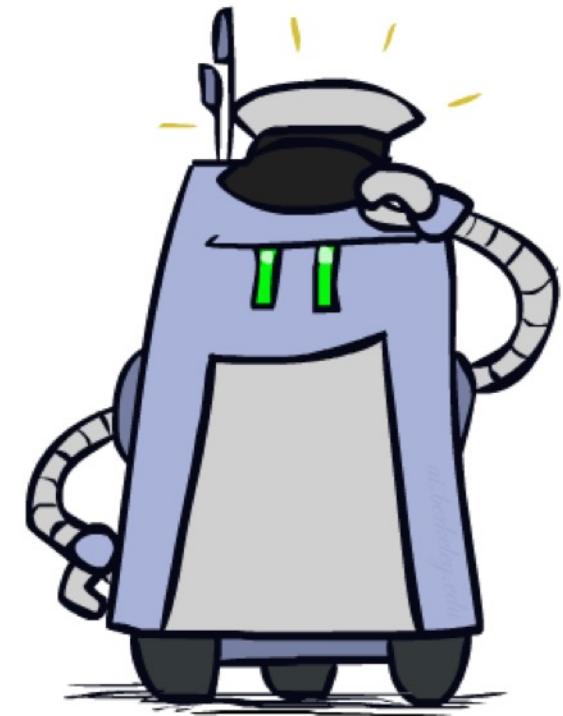
- For real classification problems, smoothing is critical
- New odds ratios:

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

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

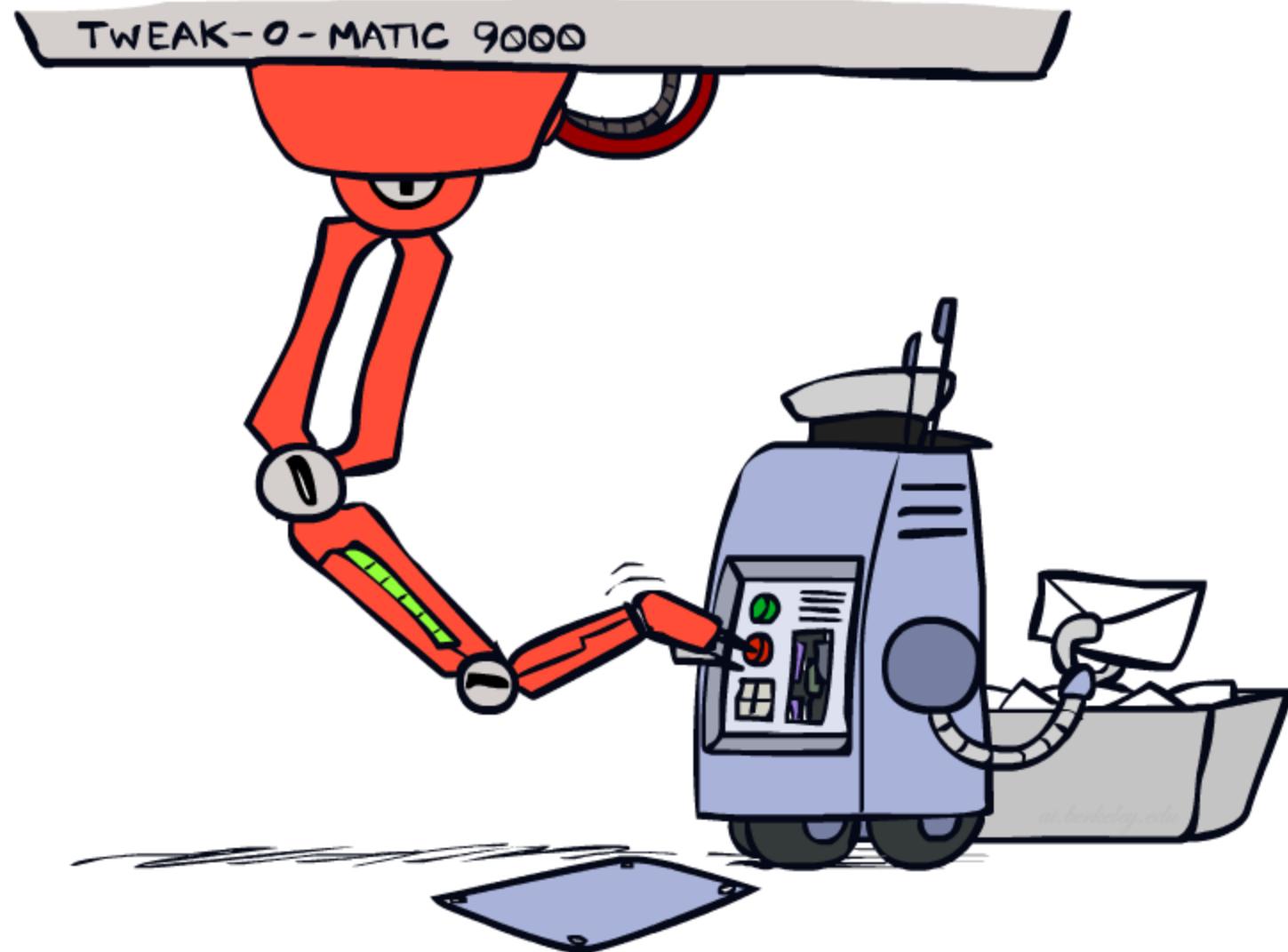
$$\frac{P(W|\text{spam})}{P(W|\text{ham})}$$

verdana	:	28.8
Credit	:	28.4
ORDER	:	27.2
	:	26.9
money	:	26.5
...		



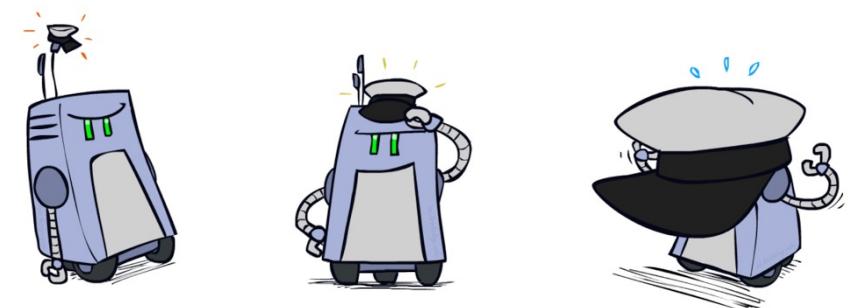
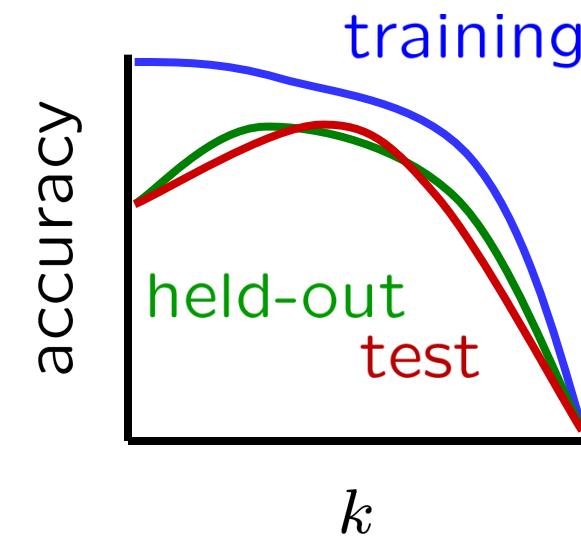
Do these make more sense?

Tuning

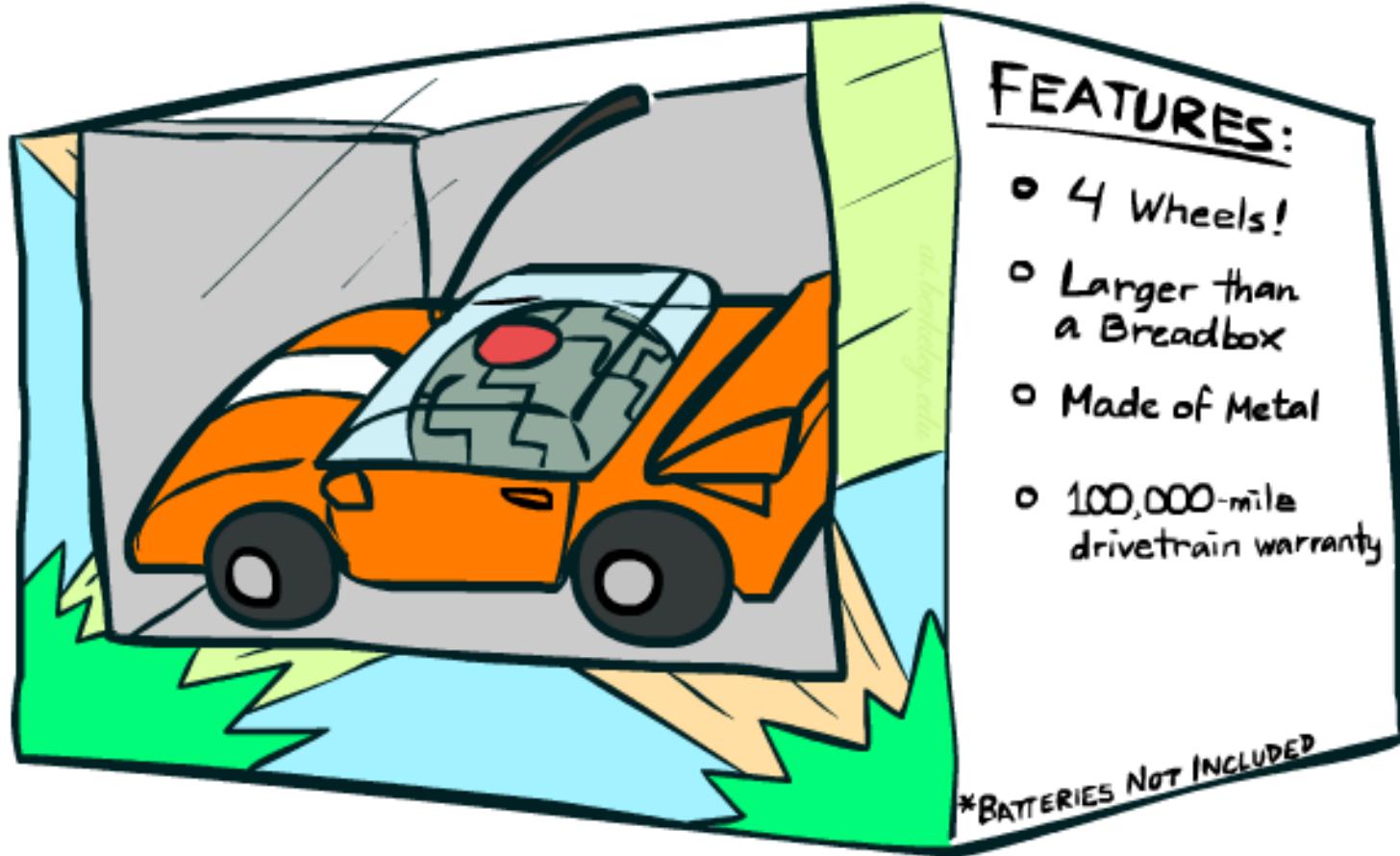


Tuning on Held-Out Data

- 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 , α
- 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



Features



Errors, and What to Do

- Examples of errors

Dear GlobalSCAPE Customer,

GlobalSCAPE has partnered with ScanSoft to offer you the latest version of OmniPage Pro, for just \$99.99* - the regular list price is \$499! The most common question we've received about this offer is - Is this genuine? We would like to assure you that this offer is authorized by ScanSoft, is genuine and valid. You can get the . . .

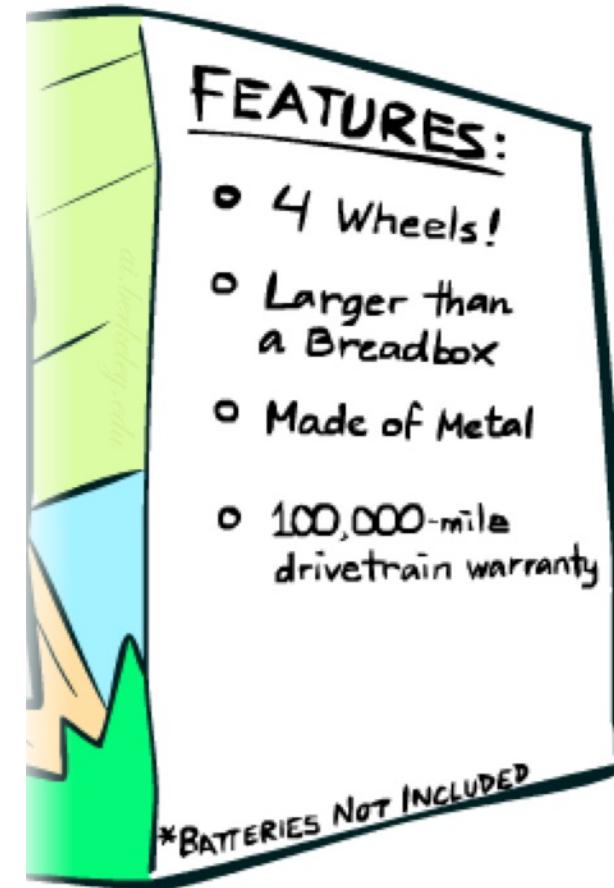
. . . To receive your \$30 Amazon.com promotional certificate, click through to

<http://www.amazon.com/apparel>

and see the prominent link for the \$30 offer. All details are there. We hope you enjoyed receiving this message. However, if you'd rather not receive future e-mails announcing new store launches, please click . . .

What to Do About Errors?

- Need more features—words aren't enough!
 - Have you emailed the sender before?
 - Have 1K 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?
- Can add these information sources as new variables in the NB model
- Later this week we'll talk about classifiers which let you easily add arbitrary features more easily, and, later, how to induce new features



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:

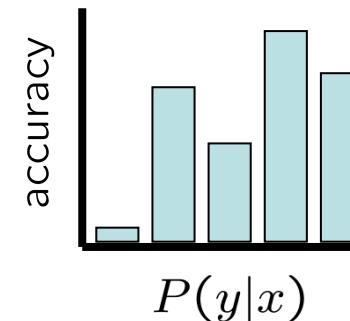
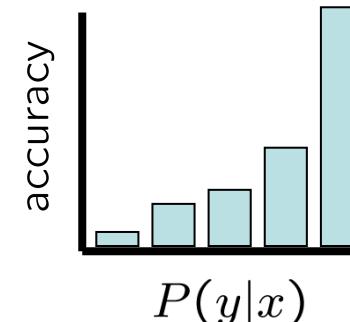
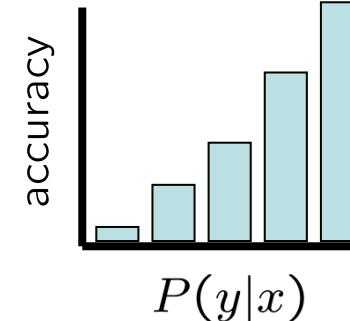
- Posterior probability of the top label

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

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

- Calibration

- Weak calibration: higher confidences mean higher accuracy
- Strong calibration: confidence predicts accuracy rate
- What's the value of calibration?



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

CS 188: Artificial Intelligence

Machine Learning



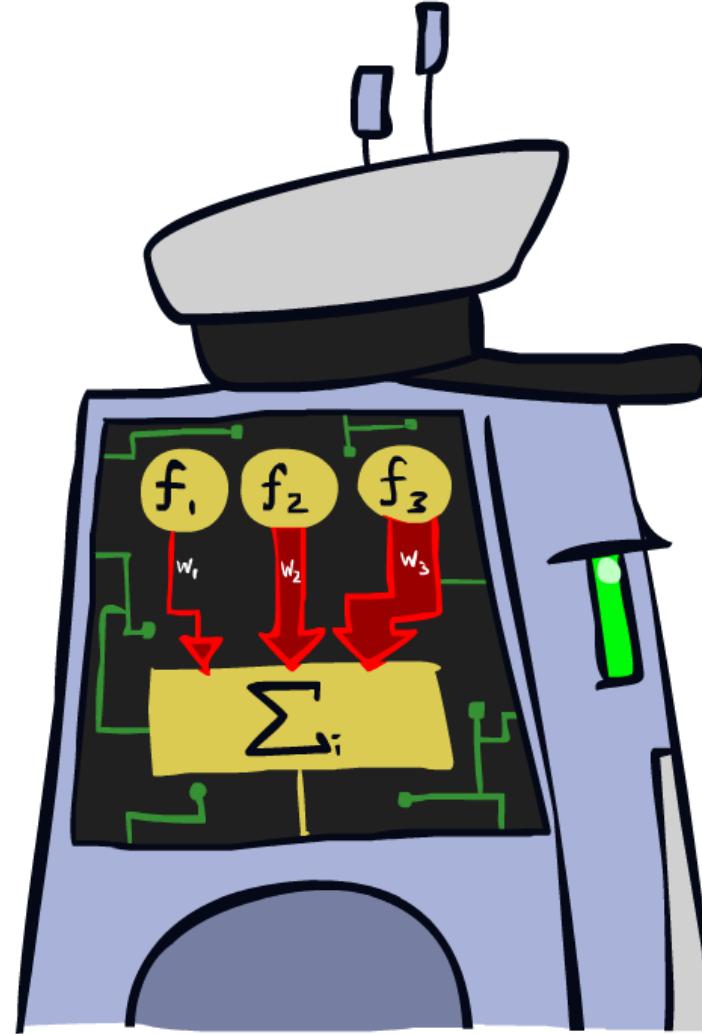
Summer 2024: Eve Fleisig & Evgeny Pobachienko

Demo: Catching AI-Generated Text

- Feature design
 - Complexity in feature design vs. model design
- Evaluation
 - Accuracy, precision & recall, F1 score
- Generalization
- Calibration
- Robustness

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Linear Classifiers



Feature Vectors

x

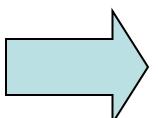
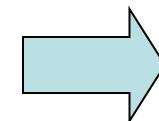
```
Hello,  
  
Do you want free printr  
cartridges? Why pay more  
when you can get them  
ABSOLUTELY FREE! Just
```

$f(x)$

$$\begin{Bmatrix} \# \text{ free} & : 2 \\ \text{YOUR_NAME} & : 0 \\ \text{MISSPELLED} & : 2 \\ \text{FROM_FRIEND} & : 0 \\ \dots \end{Bmatrix}$$

y

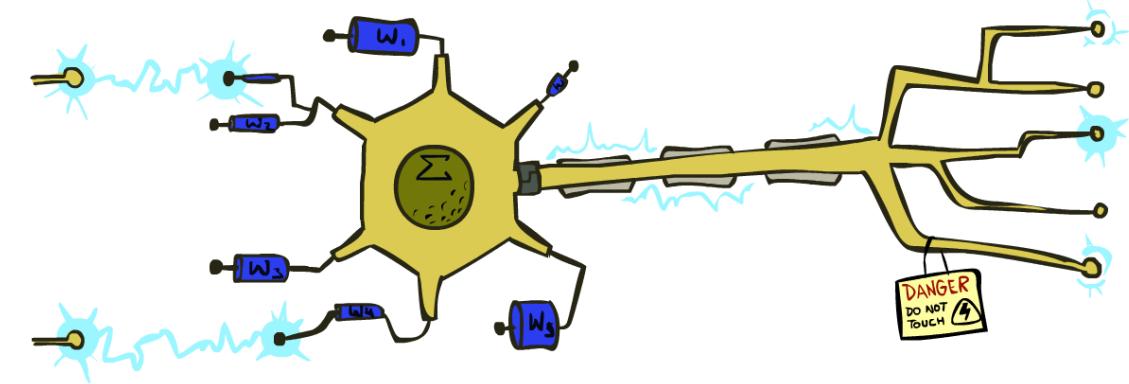
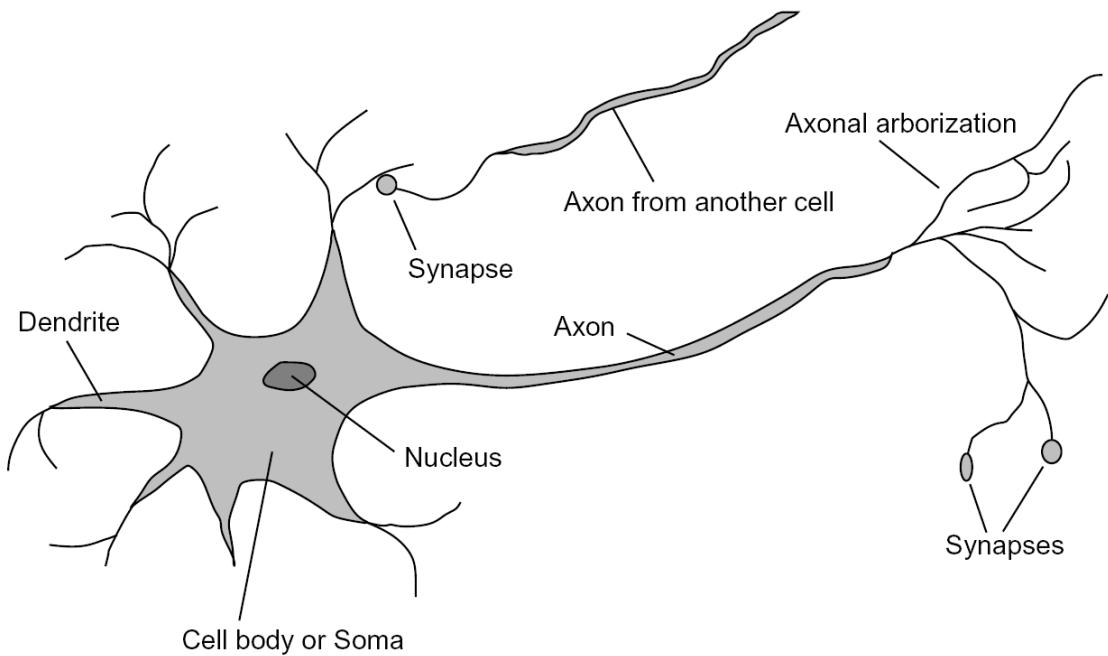
SPAM
or
+


$$\begin{Bmatrix} \text{PIXEL-7,12} & : 1 \\ \text{PIXEL-7,13} & : 0 \\ \dots \\ \text{NUM_LOOPS} & : 1 \\ \dots \end{Bmatrix}$$


“2”

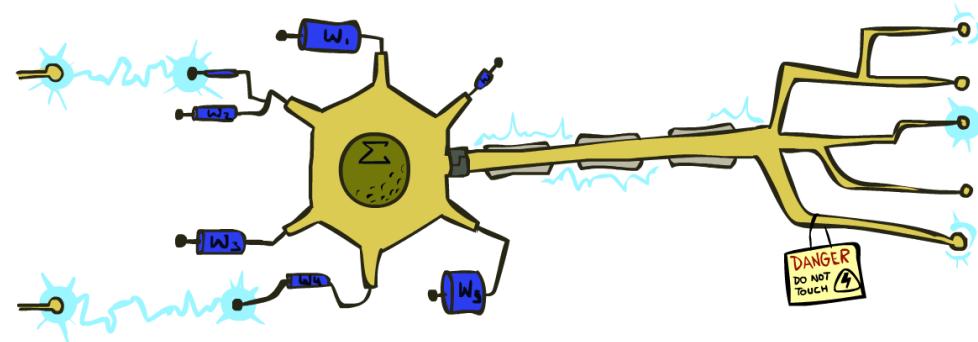
Some (Simplified) Biology

- Very loose inspiration: human neurons



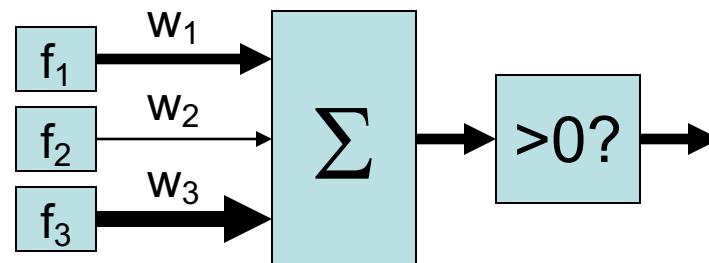
Linear Classifiers

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



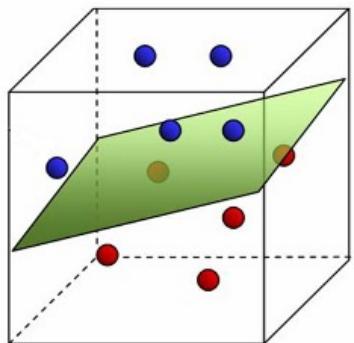
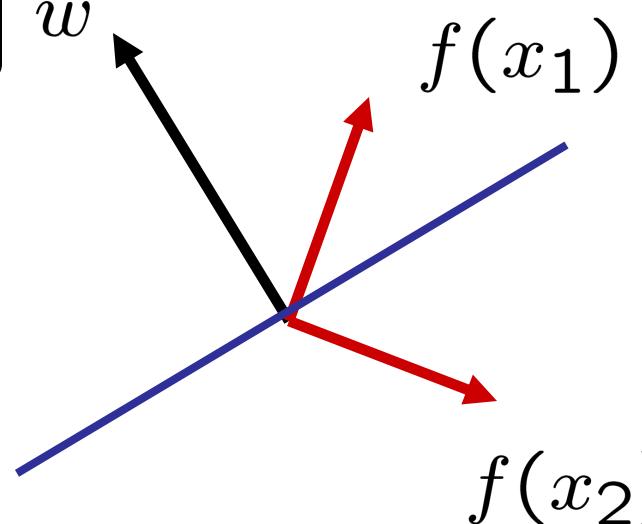
$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

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



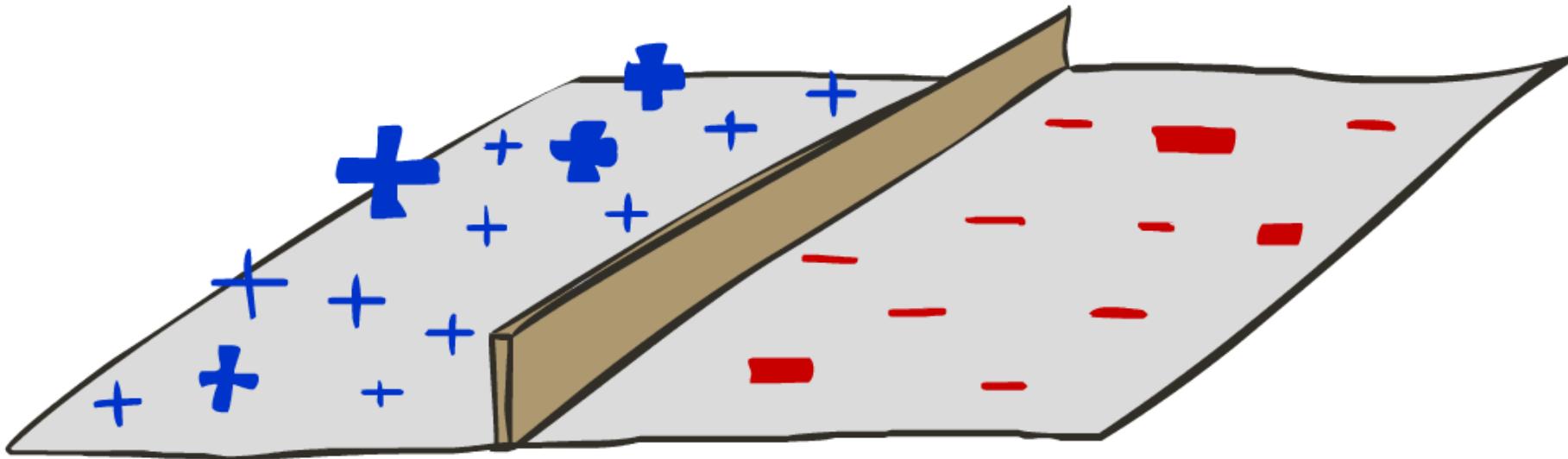
Weights

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


$$\begin{cases} \# \text{ free} & : 4 \\ \text{YOUR_NAME} & : -1 \\ \text{MISSPELLED} & : 1 \\ \text{FROM_FRIEND} & : -3 \\ \dots \end{cases}$$
 w 
$$\begin{cases} \# \text{ free} & : 2 \\ \text{YOUR_NAME} & : 0 \\ \text{MISSPELLED} & : 2 \\ \text{FROM_FRIEND} & : 0 \\ \dots \end{cases}$$
$$\begin{cases} \# \text{ free} & : 0 \\ \text{YOUR_NAME} & : 1 \\ \text{MISSPELLED} & : 1 \\ \text{FROM_FRIEND} & : 1 \\ \dots \end{cases}$$

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

Decision Rules

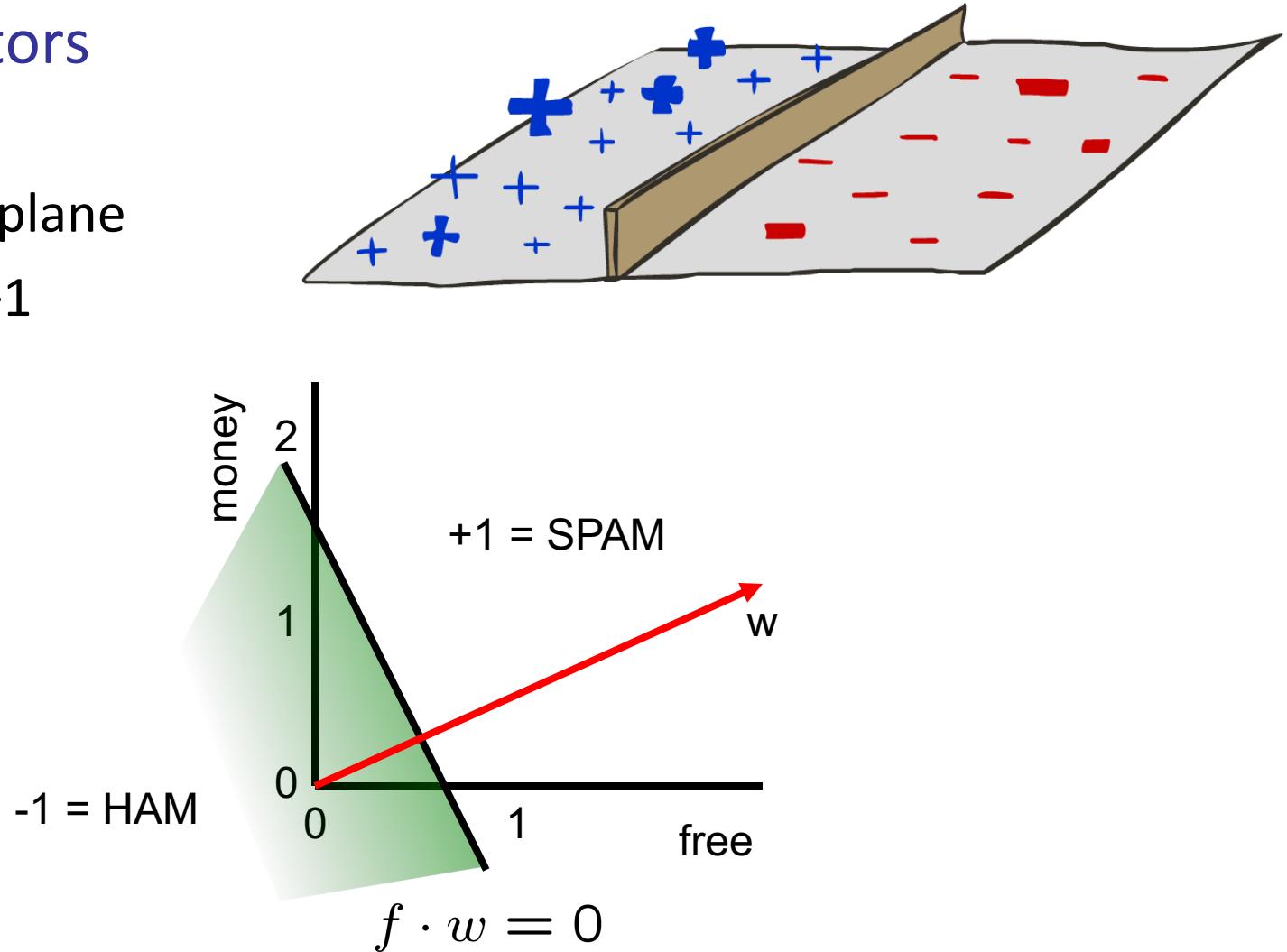


Binary Decision Rule

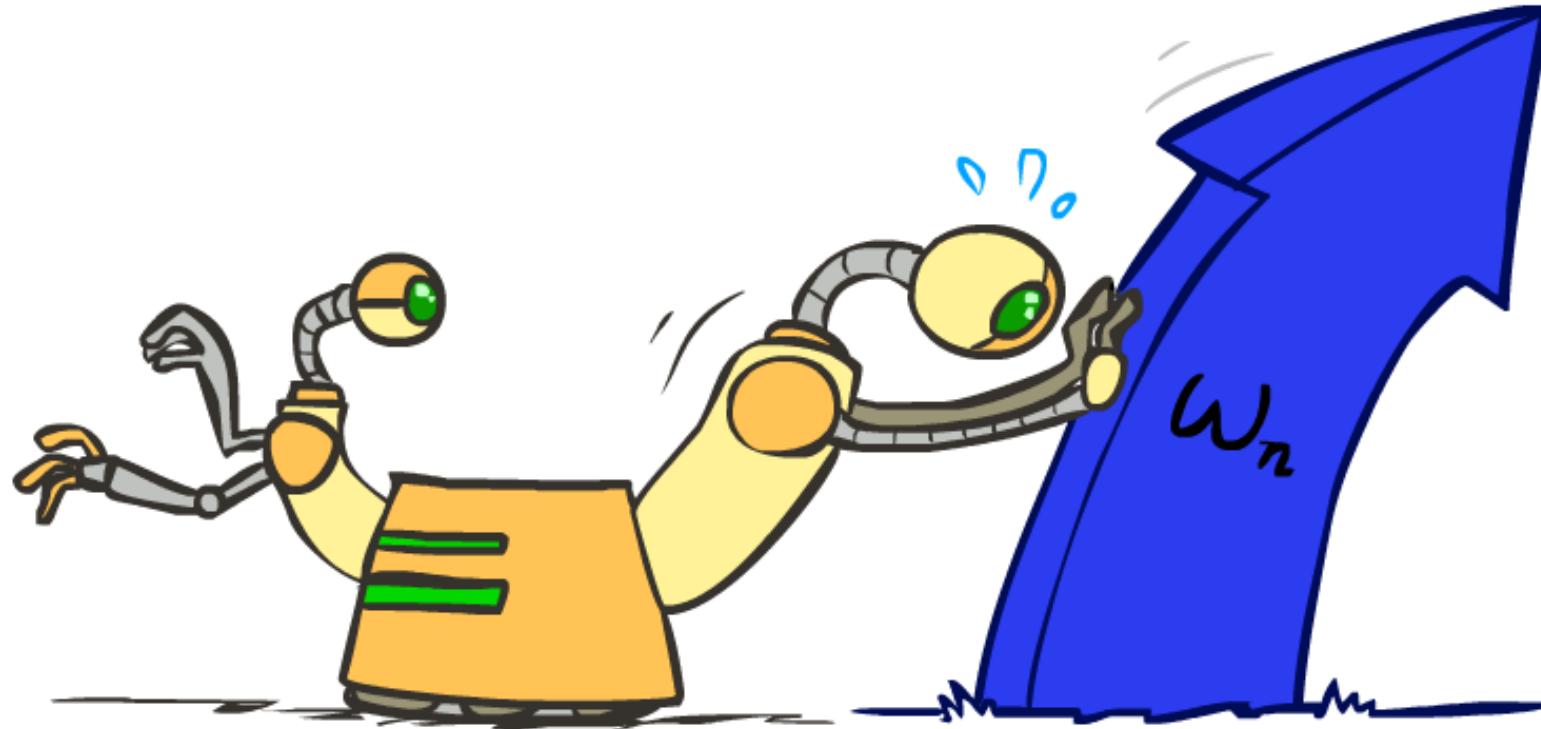
- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to $Y=+1$
 - Other corresponds to $Y=-1$

w

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

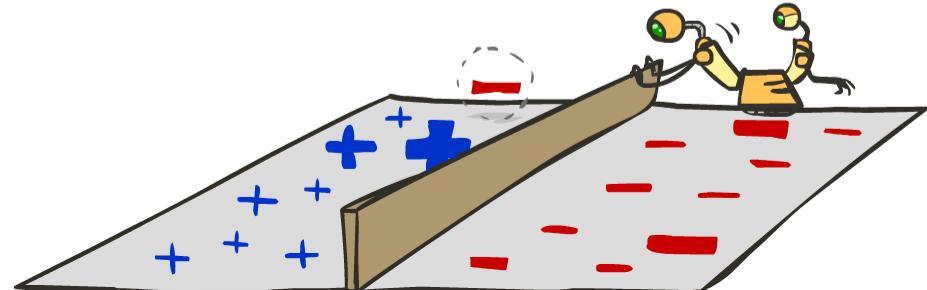
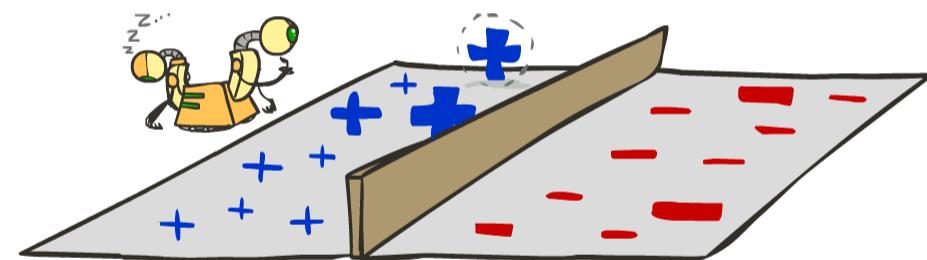
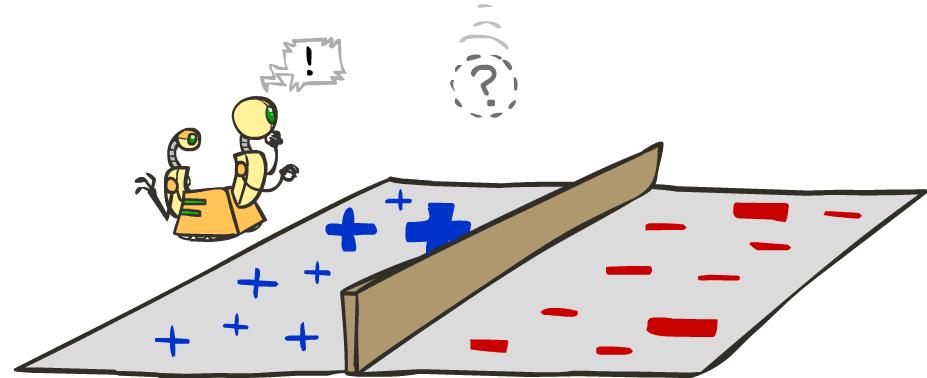


Weight Updates



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



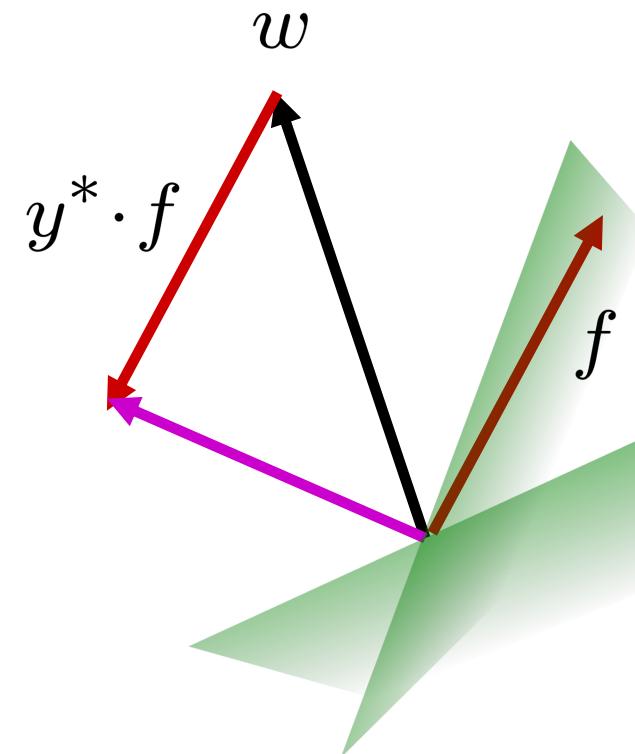
Learning: Binary Perceptron

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

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

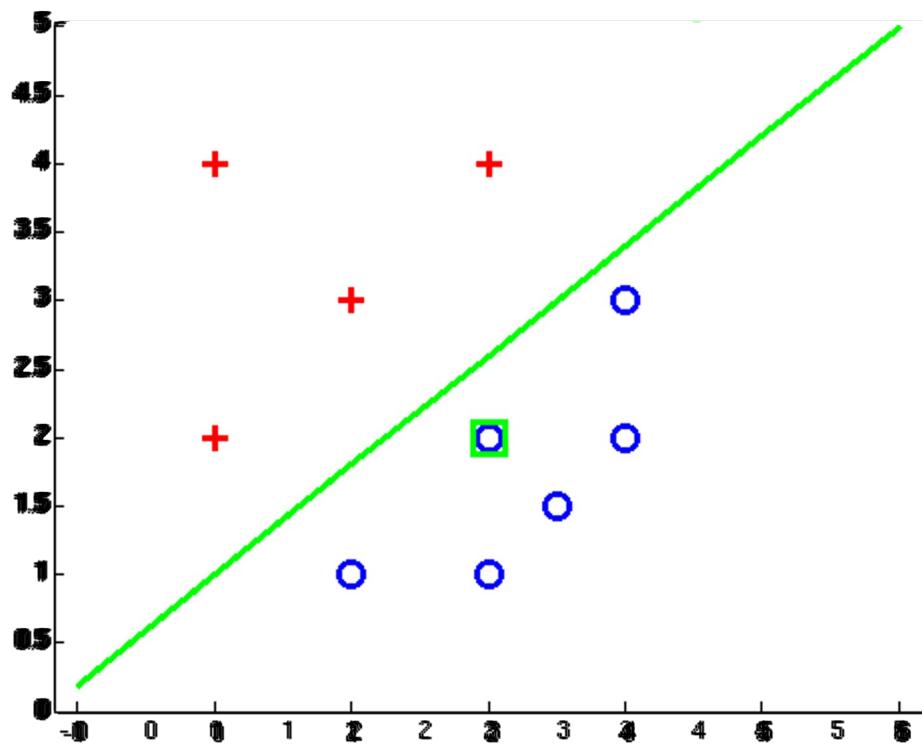
- If correct (i.e., $y=y^*$), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y^* is -1.

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



Examples: Perceptron

- Separable Case



Multiclass Decision Rule

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

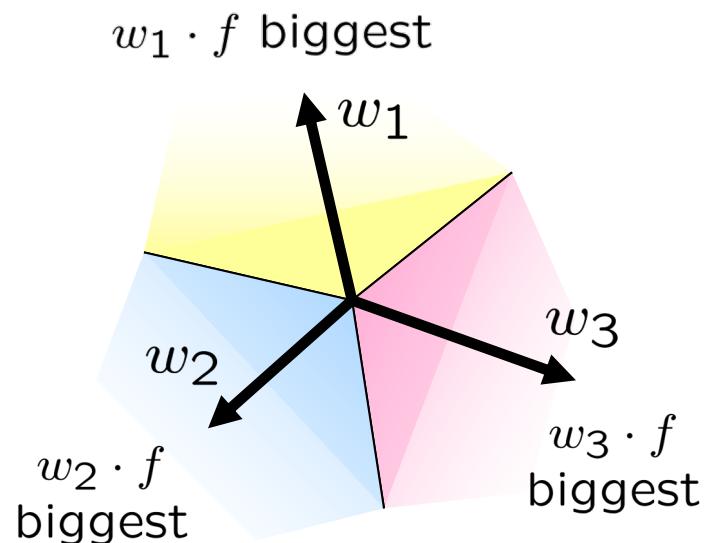
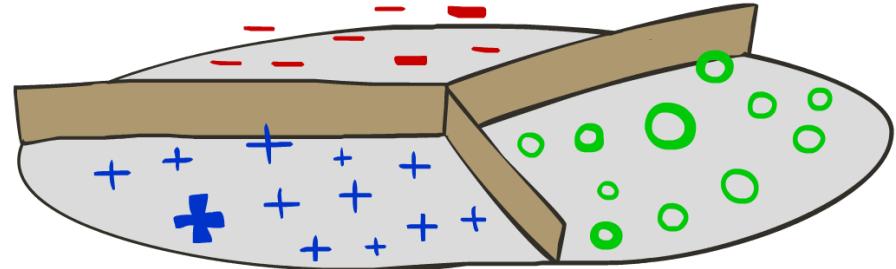
$$w_y$$

- 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

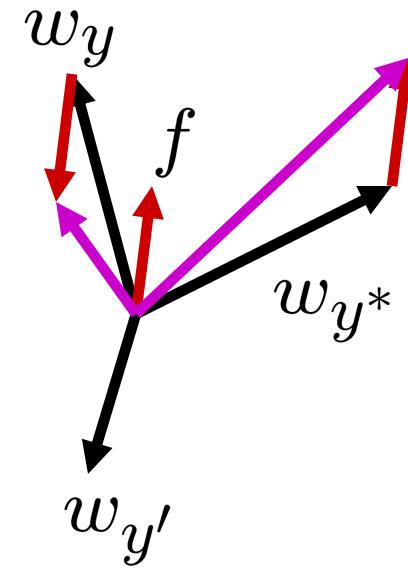
- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

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

- 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}

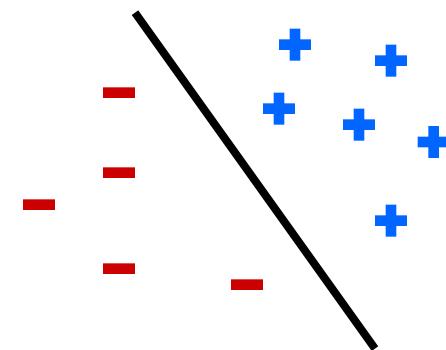
BIAS	:	0
win	:	0
game	:	0
vote	:	0
the	:	0
...		

Properties of Perceptrons

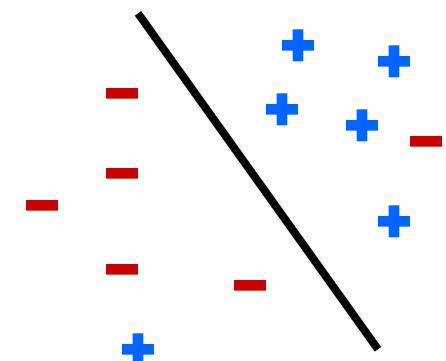
- Separability: true if some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

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

Separable

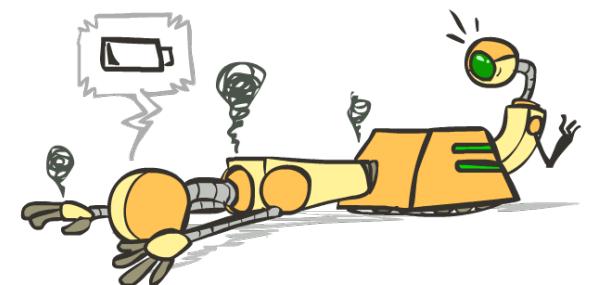
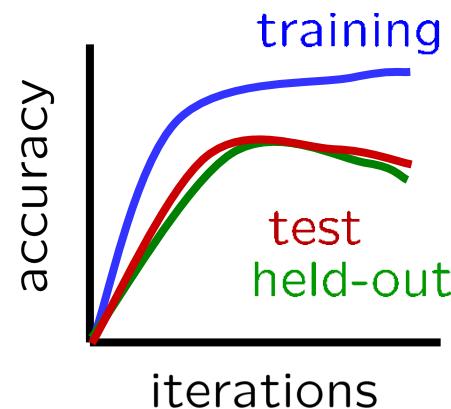
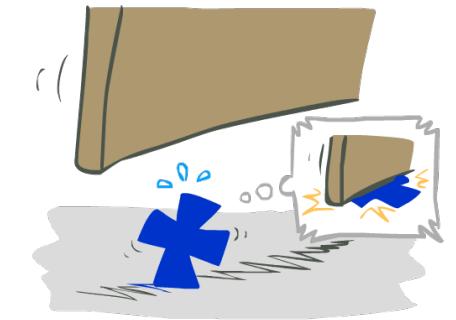
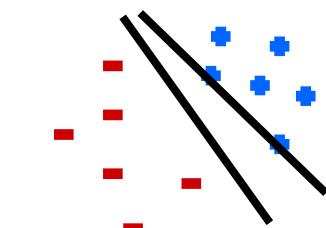
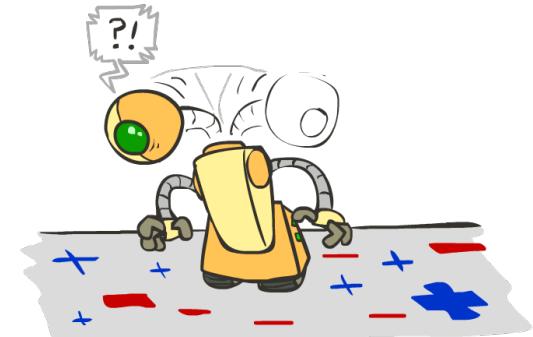
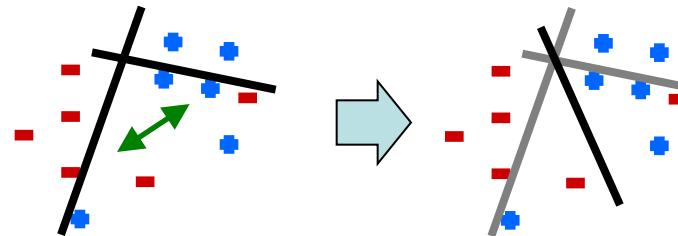


Non-Separable

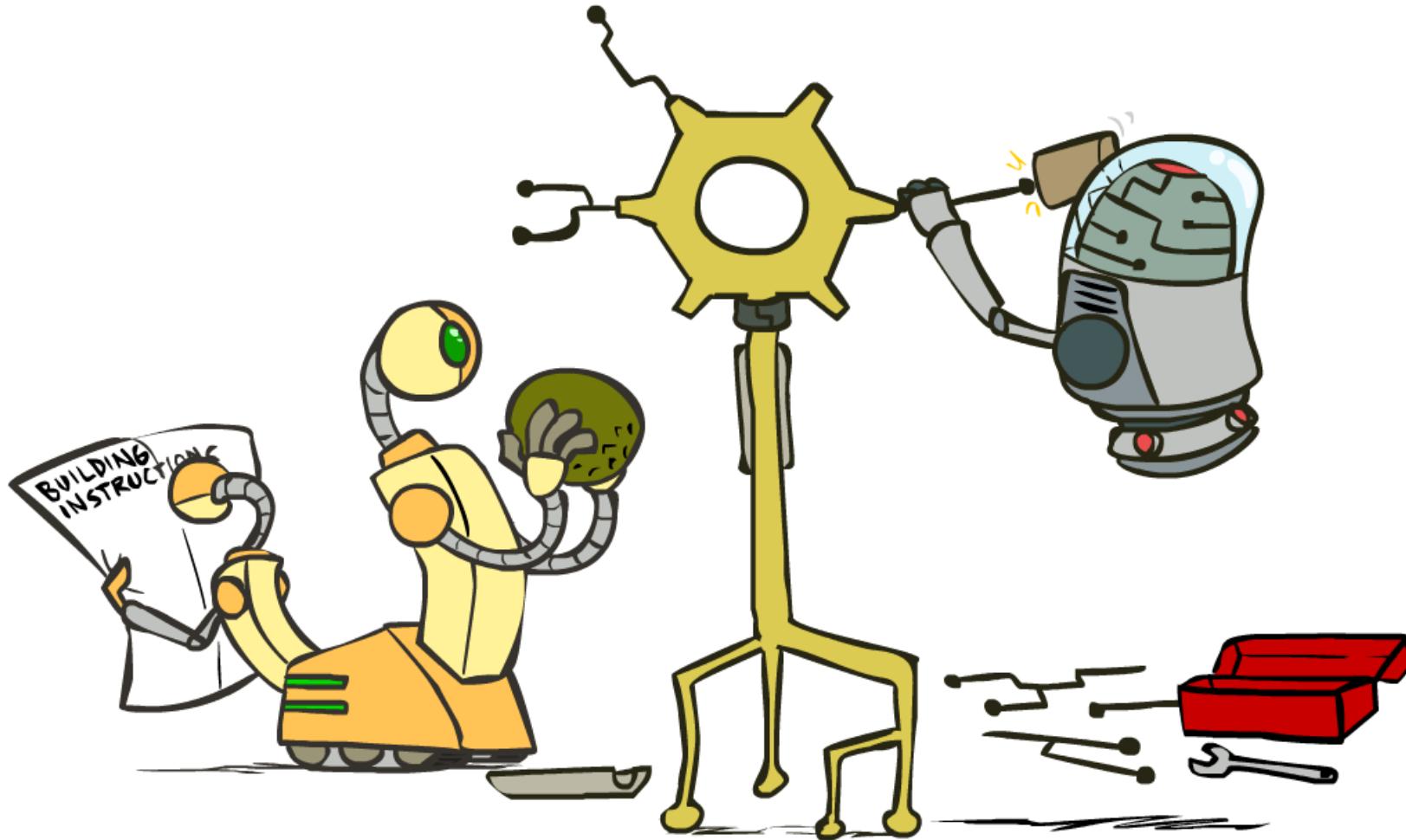


Problems with the Perceptron

- 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

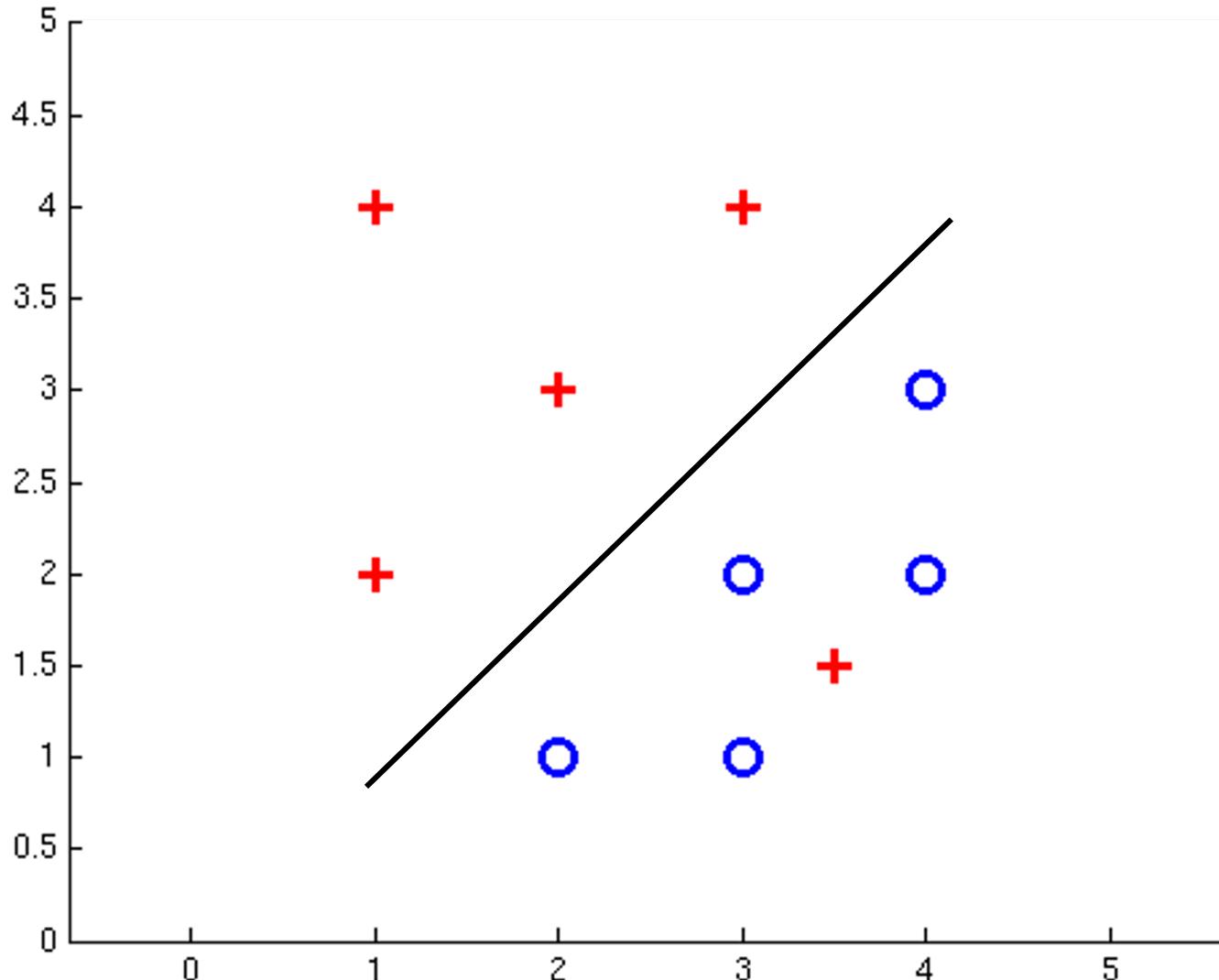


Improving the Perceptron

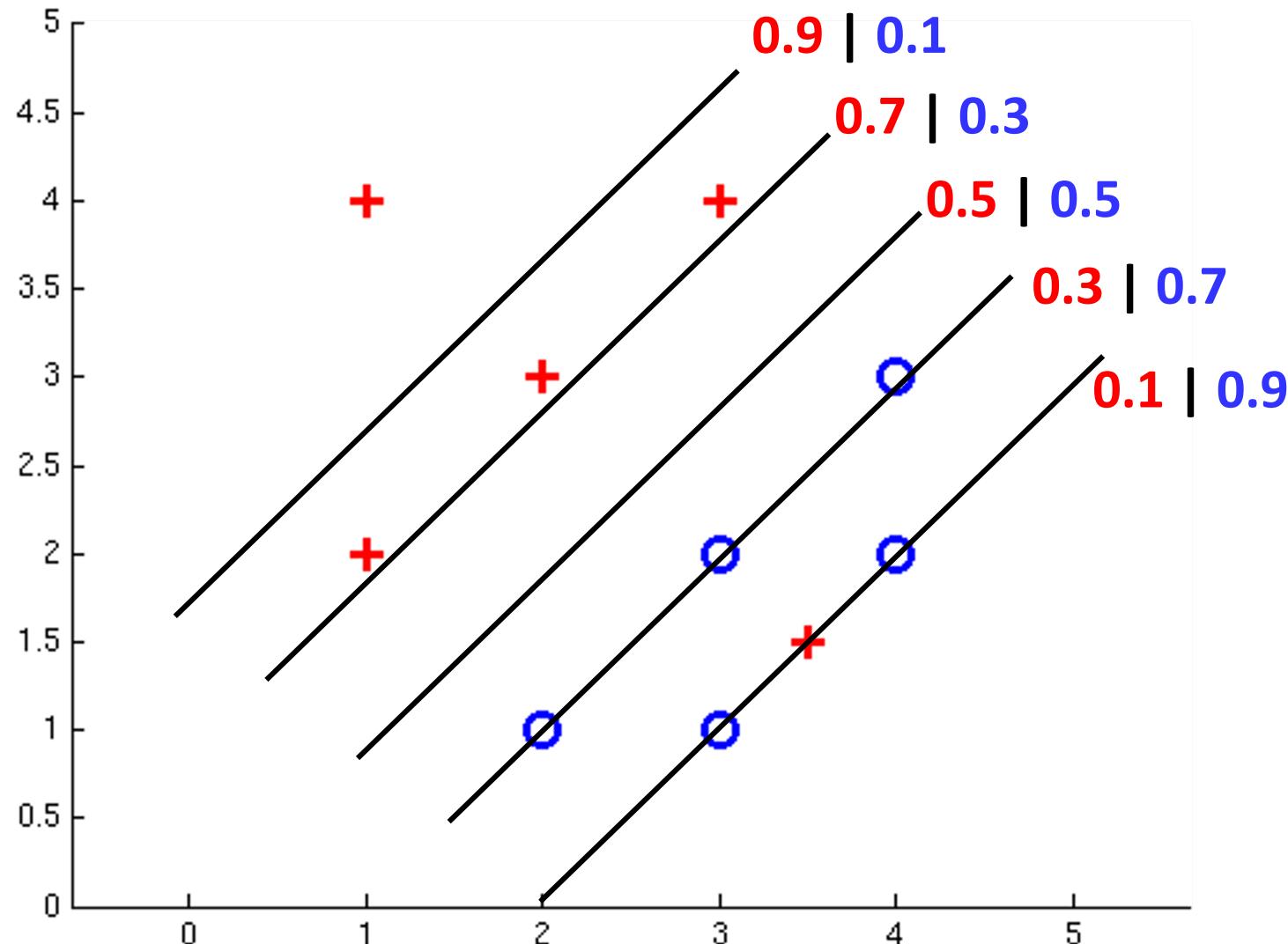


Non-Separable Case: Deterministic Decision

Even the best linear boundary makes at least one mistake



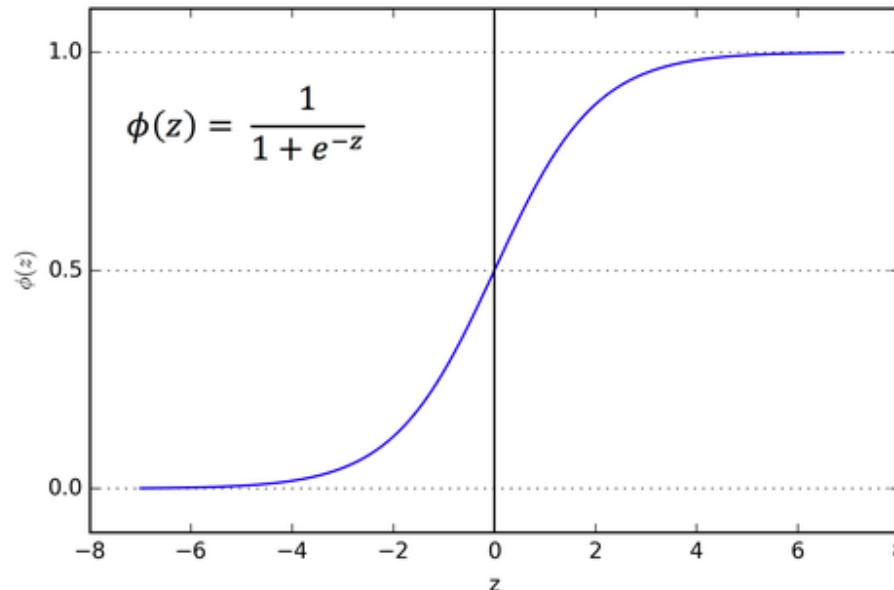
Non-Separable Case: Probabilistic Decision



How to get probabilistic decisions?

- Perceptron scoring: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive \rightarrow want probability going to 1
- If $z = w \cdot f(x)$ very negative \rightarrow want probability going to 0
- Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



Best w?

- Maximum likelihood estimation:

$$\max_w \text{ll}(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

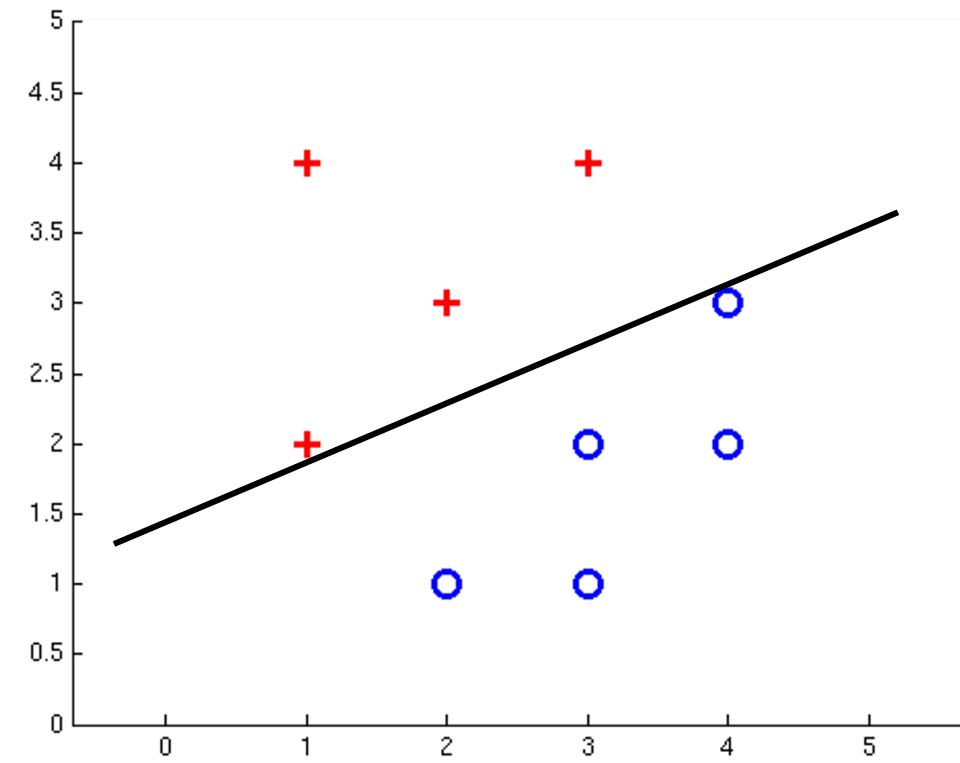
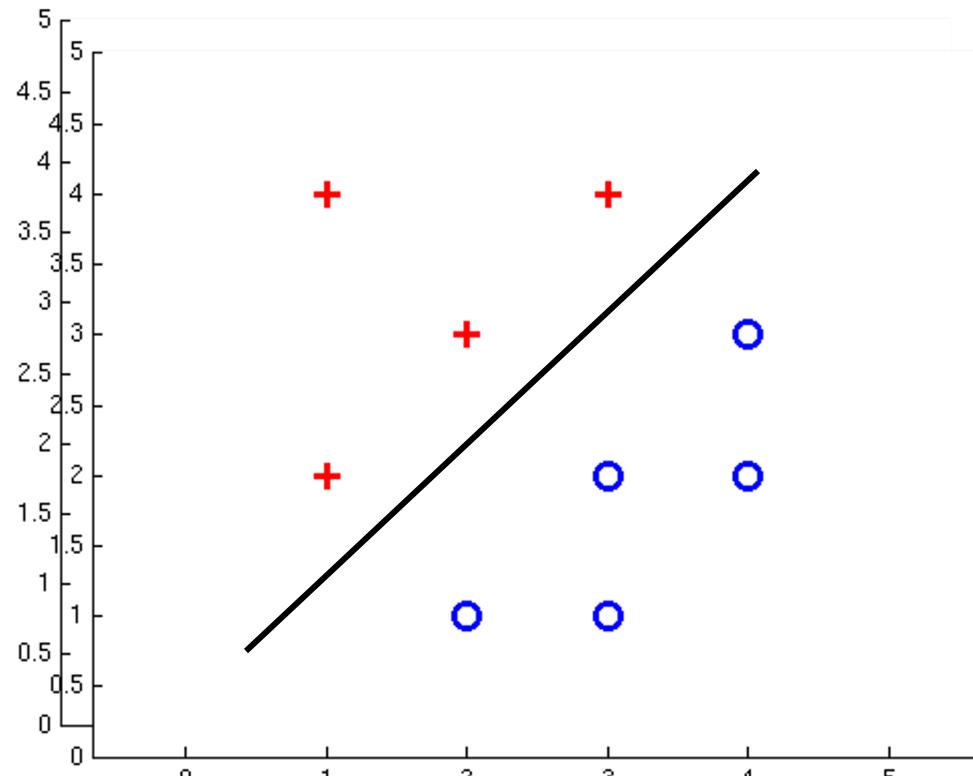
with:

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

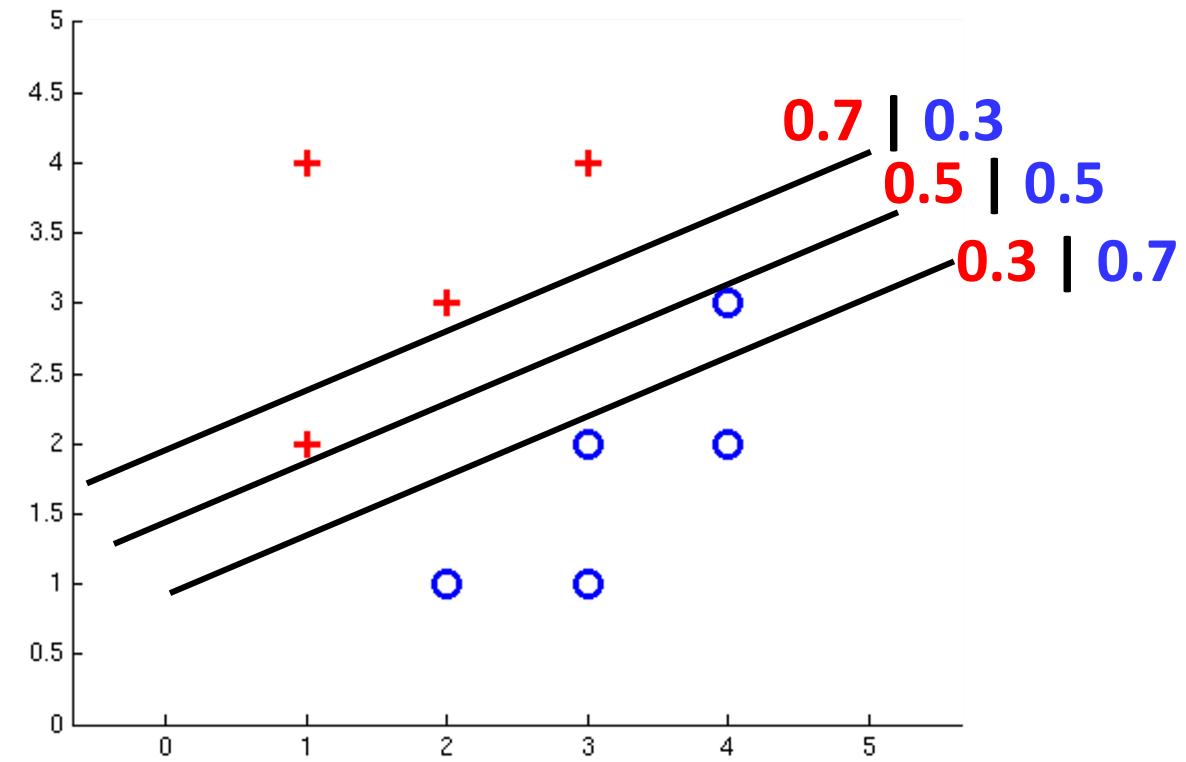
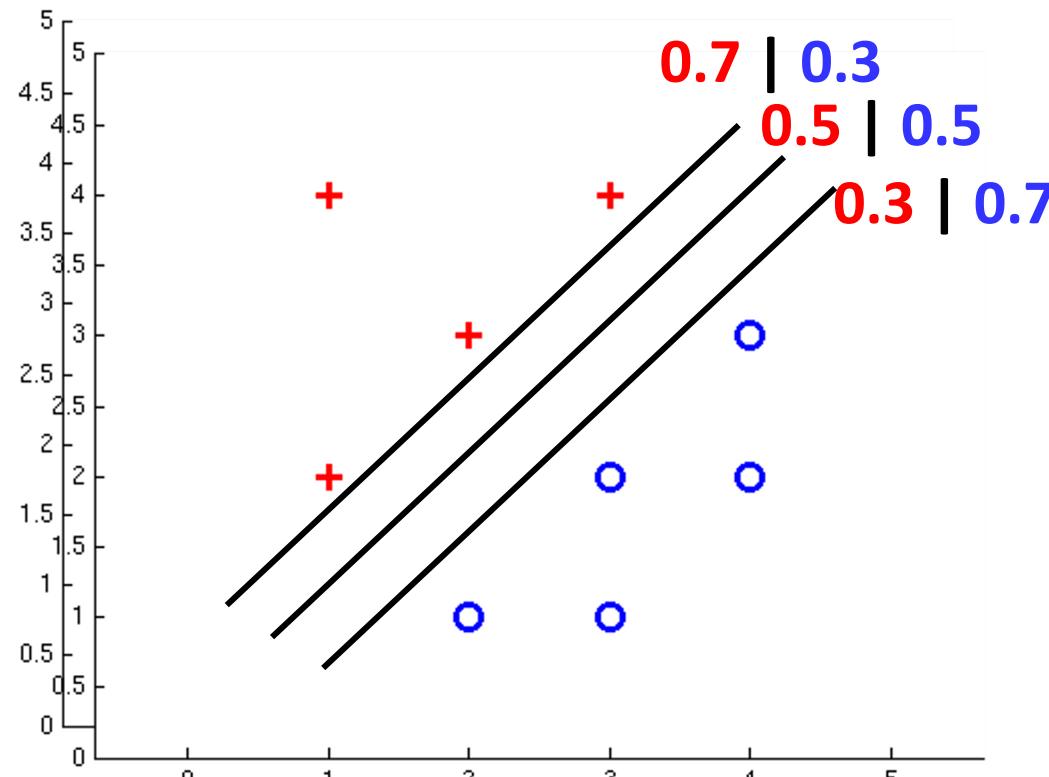
$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

= Logistic Regression

Separable Case: Deterministic Decision – Many Options



Separable Case: Probabilistic Decision – Clear Preference



Multiclass Logistic Regression

- Recall Perceptron:

- A weight vector for each class:

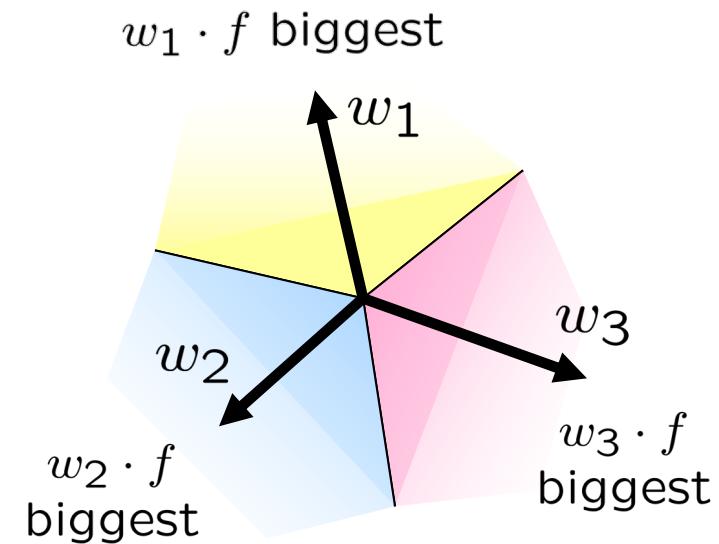
$$w_y$$

- Score (activation) of a class y :

$$w_y \cdot f(x)$$

- Prediction highest score wins

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



- How to make the scores into probabilities?

$$z_1, z_2, z_3 \rightarrow \underbrace{\frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}}_{\text{softmax activations}}$$

original activations

Best w?

- Maximum likelihood estimation:

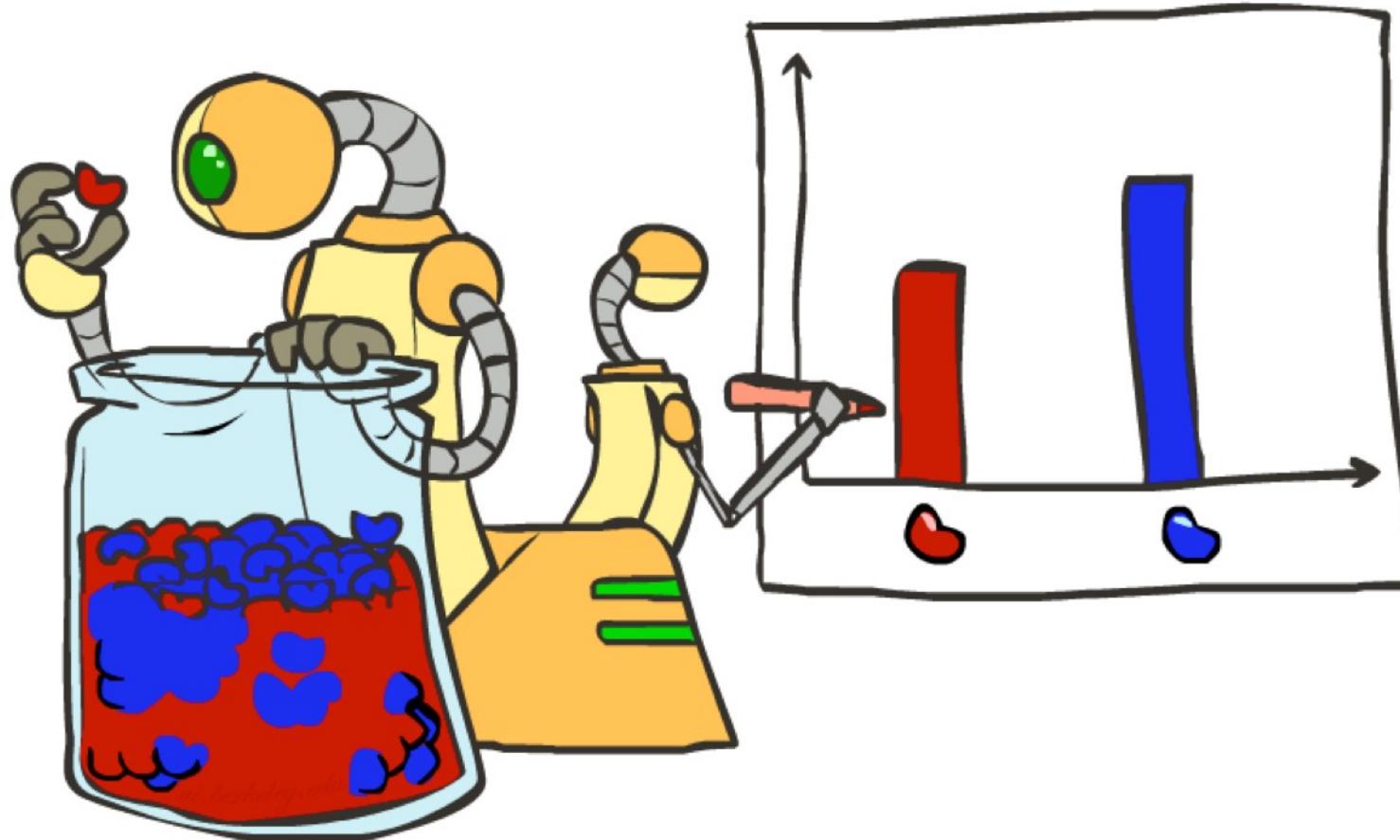
$$\max_w \text{ll}(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

with:

$$P(y^{(i)} | x^{(i)}; w) = \frac{e^{w_y \cdot f(x^{(i)})}}{\sum_y e^{w_y \cdot f(x^{(i)})}}$$

= Multi-Class Logistic Regression

Maximum Likelihood Estimation



Parameter Estimation with Maximum Likelihood

- Estimating the distribution of a random variable
- Use training data (learning!)
 - For each outcome x , look at the **empirical rate** of that value:

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

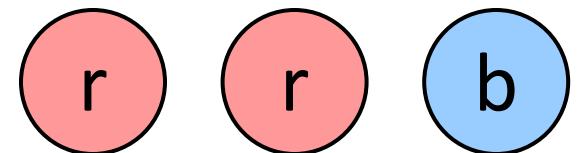
- Example: probability of $x=\text{red}$ given the training data:

$$P_{ML}(\text{red}) = \frac{2}{3}$$

- This estimate maximizes the **likelihood of the data** for the parametric model:

$$\begin{aligned}L(\theta) &= P(\text{red, red, blue} \mid \theta) = P_\theta(\text{red}) \cdot P_\theta(\text{red}) \cdot P_\theta(\text{blue}) \\&= \theta^2 \cdot (1 - \theta)\end{aligned}$$

x	red	blue
$P_\theta(x)$	θ	$1 - \theta$



Parameter Estimation with Maximum Likelihood

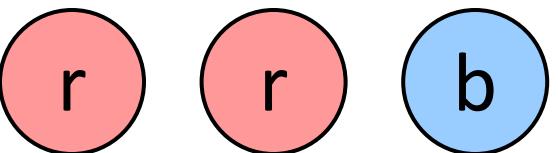
- Likelihood function:

$$\begin{aligned}L(\theta) &= P(\text{r}, \text{r}, \text{b} \mid \theta) = P_\theta(\text{r}) \cdot P_\theta(\text{r}) \cdot P_\theta(\text{b}) \\&= \theta^2 \cdot (1 - \theta) \\&= \theta^2 - \theta^3\end{aligned}$$

X	red	blue
$P_\theta(x)$	θ	$1 - \theta$

- MLE: find the θ that maximizes data likelihood

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$$



- Approach: take derivatives and set to 0

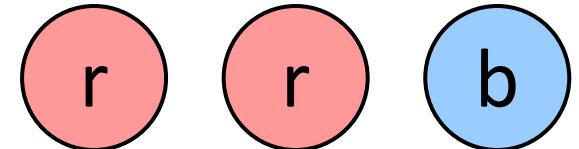
$$\begin{aligned}\frac{\partial L(\theta)}{\partial \theta} &= 2\theta - 3\theta^2 \\&= \theta(2 - 3\theta)\end{aligned}$$

- Find the maximum at $\theta = \frac{2}{3}$

Parameter Estimation (General Case)

- **Model:**

x	red	blue
$P_\theta(x)$	θ	$1 - \theta$



- **Data:** draw N balls. N_r come up red, N_b come up blue

- Dataset: $D = \{x_1, \dots, x_n\}$

- Ball draws are independent and identically distributed (i.i.d.):

$$P(D | \theta) = \prod_i P(x_i | \theta) = \prod_i P_\theta(x_i) = \theta^{N_r} \cdot (1 - \theta)^{N_b}$$

- **Maximum likelihood estimation:** find θ that maximizes $P(D | \theta)$

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

- Approach: take derivative and set to 0

Parameter Estimation (General Case)

- Maximum likelihood estimation: find θ that maximizes $P(D | \theta)$

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

$$\begin{aligned}\frac{\partial}{\partial \theta} \log P(D | \theta) &= \frac{\partial}{\partial \theta} [N_r \log(\theta) + N_b \log(1 - \theta)] \\ &= N_r \frac{\partial}{\partial \theta} \log(\theta) + N_b \frac{\partial}{\partial \theta} \log(1 - \theta) \\ &= N_r \frac{1}{\theta} - N_b \frac{1}{1-\theta} \\ &= 0\end{aligned}$$

Multiply by $\theta(1 - \theta)$:

$$\begin{aligned}N_r(1 - \theta) - N_b\theta &= 0 \\ N_r - \theta(N_r + N_b) &= 0\end{aligned}$$

$$\hat{\theta} = \frac{N_r}{N_r + N_b}$$

Example from Discussion 6B

1 Maximum Likelihood Estimation

Recall that a Geometric distribution is defined as the number of Bernoulli trials needed to get one success. $P(X = k) = p(1 - p)^{k-1}$.

We observe the following samples from a Geometric distribution:

$$x_1 = 5, x_2 = 8, x_3 = 3, x_4 = 5, x_5 = 7$$

What is the maximum likelihood estimate for p ?

$$L(p) = P(X = x_1)P(X = x_2)P(X = x_3)P(X = x_4)P(X = x_5) \quad (1)$$

$$= P(X = 5)P(X = 8)P(X = 3)P(X = 5)P(X = 7) \quad (2)$$

$$= p^5(1 - p)^{23} \quad (3)$$

$$\log(L(p)) = 5 \log(p) + 23 \log(1 - p) \quad (4)$$

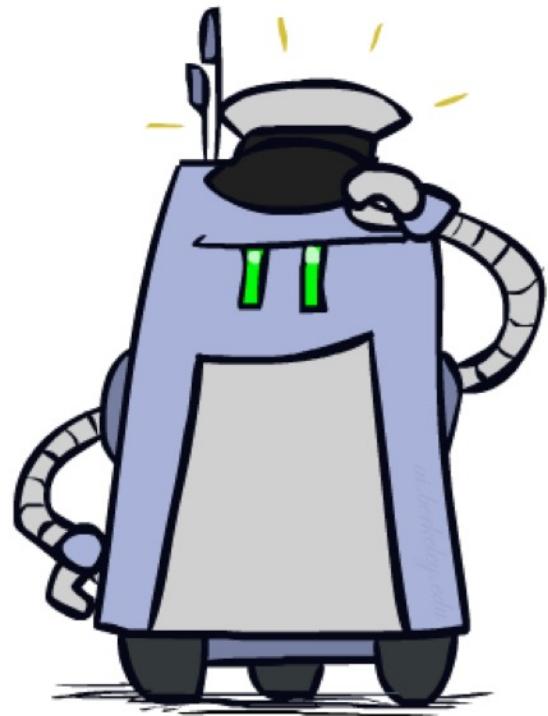
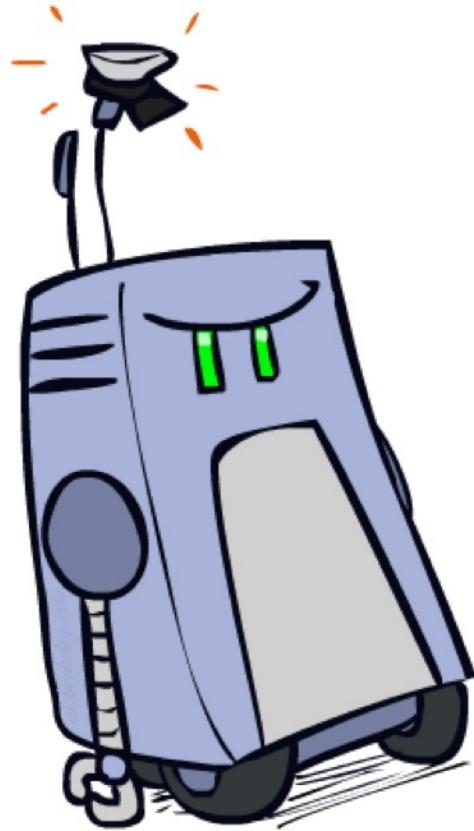
(5)

We must maximize the log-likelihood of p , so we will take the derivative, and set it to 0.

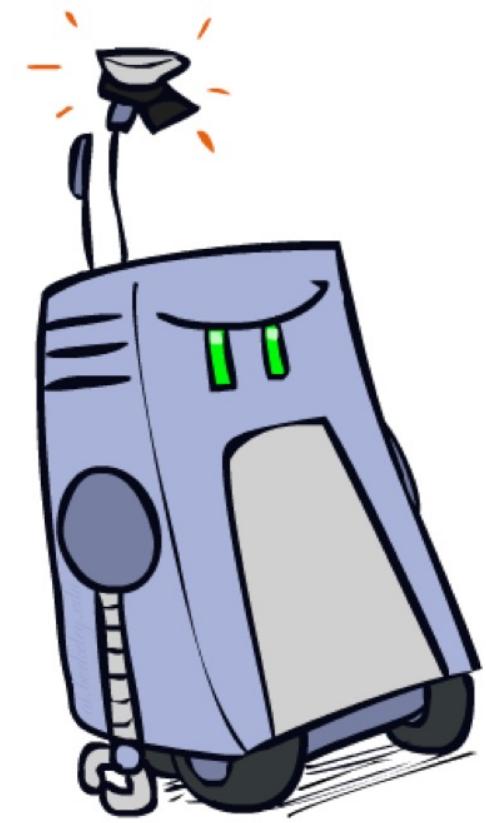
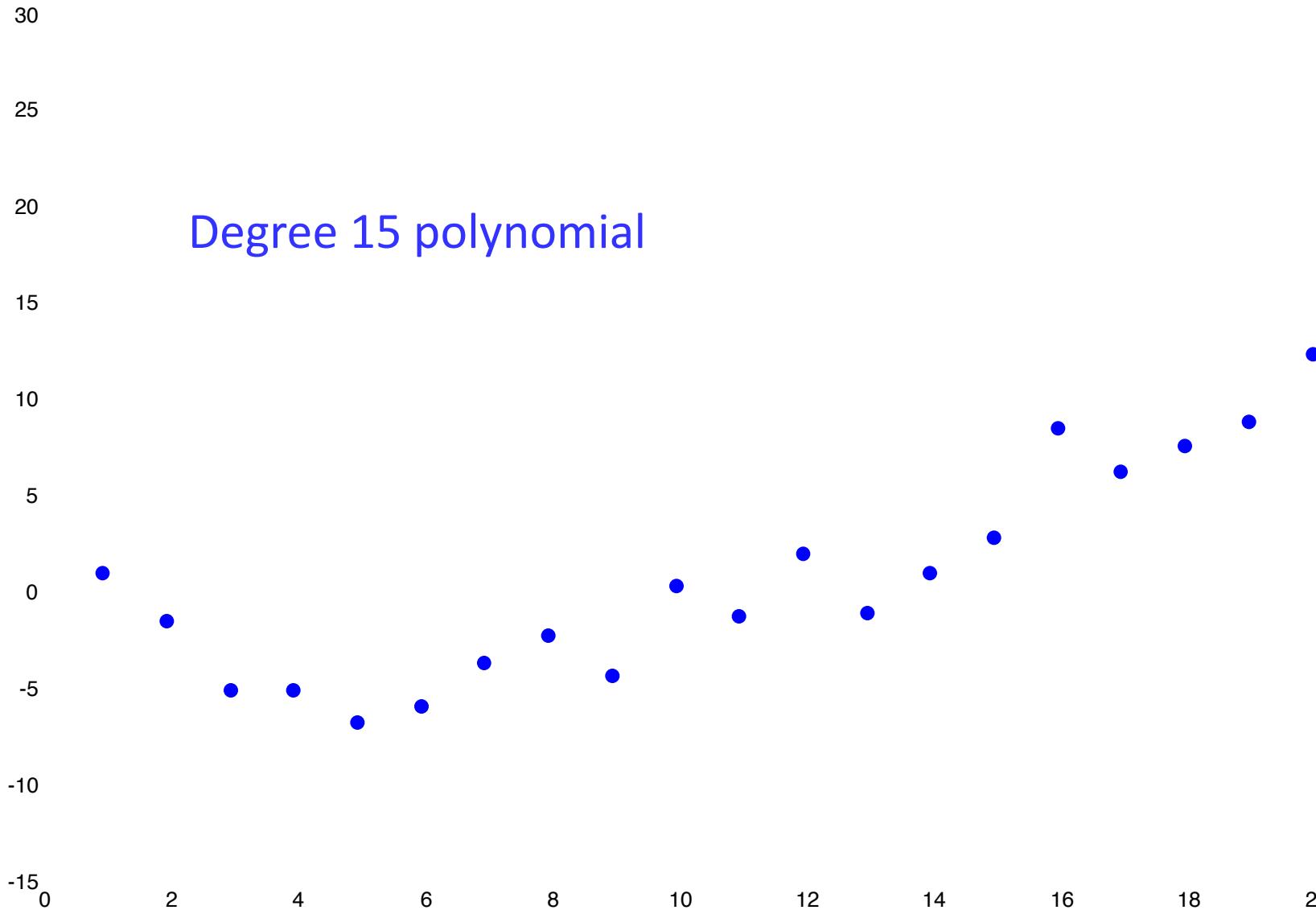
$$0 = \frac{5}{p} - \frac{23}{1 - p} \quad (6)$$

$$p = 5/28 \quad (7)$$

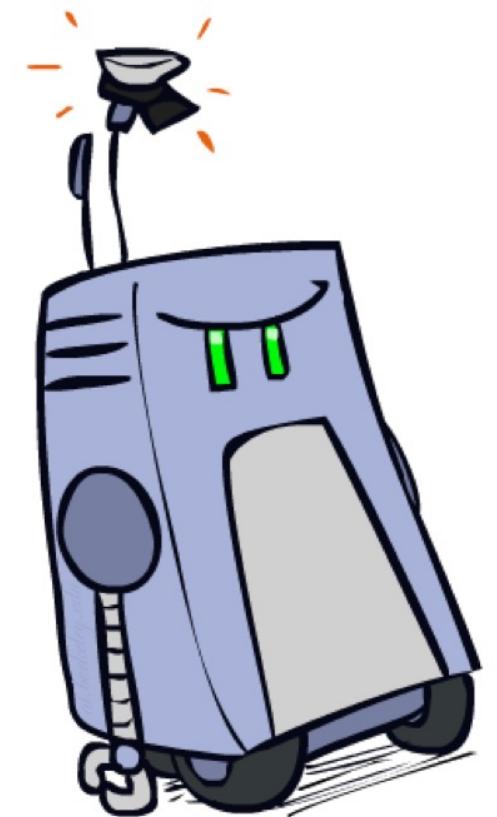
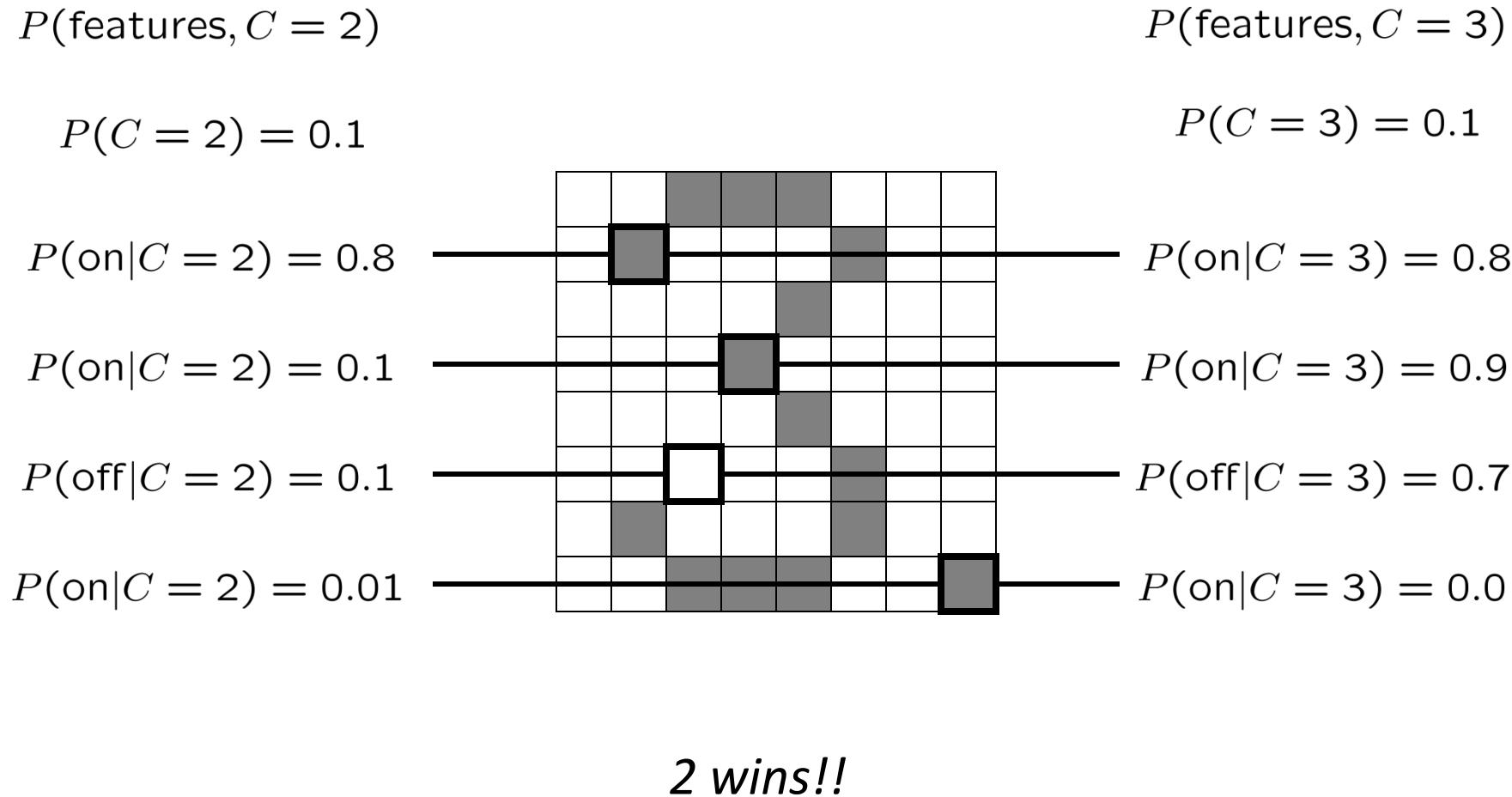
Regularization



Recall: Overfitting

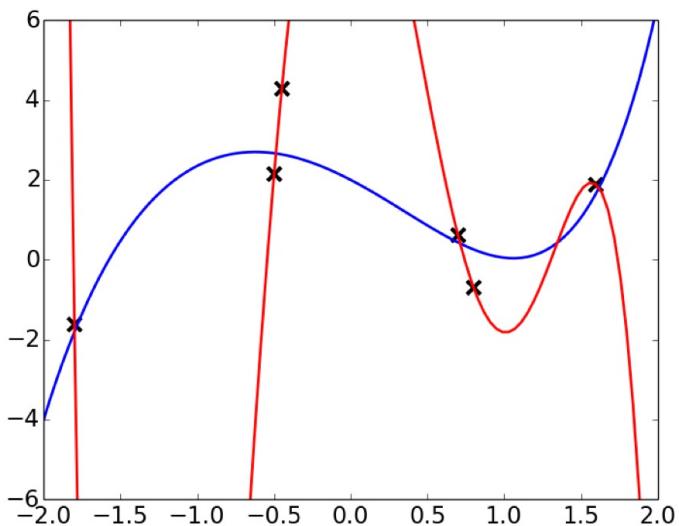


Example: Overfitting

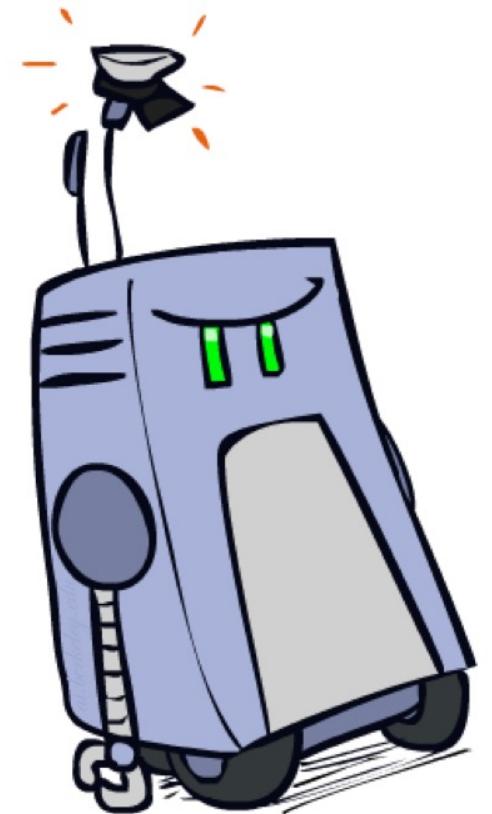


Recall: Overfitting

- Observation: polynomials that overfit tend to have large coefficients



$$y = 0.1x^5 + 0.2x^4 + 0.75x^3 - x^2 - 2x + 2$$
$$y = -7.2x^5 + 10.4x^4 + 24.5x^3 - 37.9x^2 - 3.6x + 12$$



- Let's try to keep coefficients small!

L1 and L2 Regularization

- Previously:

$$\hat{w} = \arg \max_w \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; w)$$

- Now: add a penalty term to keep the weight vector small

L1
(aka lasso regression)

$$\hat{w} = \arg \max_w \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; w) - \alpha \sum_{i=1}^n |w_i|$$

L2
(aka ridge regression)

$$\hat{w} = \arg \max_w \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; w) - \alpha \sum_{i=1}^n w_i^2$$