

## CS 188: Artificial Intelligence

### Naïve Bayes: RECAP



Instructor: Marco Alvarez --- University of Rhode Island

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley.  
All CS188 materials are available at <http://ai.berkeley.edu>.]

1

## Model-Based Classification



2

## 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

- Challenge
  - How should we learn its parameters?



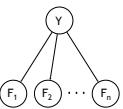
3

## General Naïve Bayes

- A general Naïve Bayes model:

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

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

4

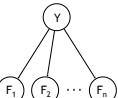
## Naïve Bayes for Digits

- Naïve Bayes: Assume all features are independent effects of the label
- Simple digit recognition version:
  - One feature (variable)  $F_{i,j}$  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.

$$\begin{bmatrix} 1 & \dots & 0 & 0,0 = 0 & F_{0,1} = 0 & F_{0,2} = 1 & F_{0,3} = 1 & F_{0,4} = 0 & \dots & F_{15,15} = 0 \end{bmatrix}$$

▪ 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?



## Example: Conditional Probabilities

$P(Y)$	$P(F_{3,1} = on Y)$	$P(F_{5,5} = on Y)$
1 0.1	1 0.01	1 0.05
2 0.1	2 0.05	2 0.01
3 0.1	3 0.05	3 0.90
4 0.1	4 0.30	4 0.80
5 0.1	5 0.80	5 0.90
6 0.1	6 0.90	6 0.90
7 0.1	7 0.05	7 0.25
8 0.1	8 0.60	8 0.85
9 0.1	9 0.50	9 0.60
0 0.1	0 0.80	0 0.80

## Example: Spam Filtering

$$\text{Model: } P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i | Y)$$

- What are the parameters?

$P(Y)$	$P(W spam)$	$P(W ham)$
ham : 0.66	the : 0.0156	the : 0.0210
spam: 0.33	to : 0.0153	to : 0.0133
	and : 0.0115	of : 0.0119
	of : 0.0095	2002: 0.0110
	you : 0.0093	with: 0.0108
	a : 0.0086	from: 0.0107
	with: 0.0080	and : 0.0105
	from: 0.0075	a : 0.0100
	...	...

- Where do these tables come from?

7

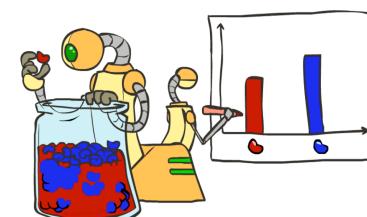
## Spam Example

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

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

8

## Parameter Estimation



9

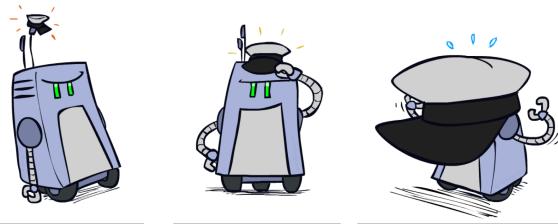
## Parameter Estimation

- 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:
- $P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$
- $P_{ML}(r) = 2/3$
- This is the estimate that maximizes the *likelihood of the data*
- $$L(x, \theta) = \prod_i P_\theta(x_i)$$
- $$P(r) = \theta = \max_\theta P(r)P(\theta)$$
- $$P(b) = 1 - \theta = \max_\theta \theta^2(1 - \theta)$$
- $$2\theta - 3\theta = 0$$
- $$2\theta - 3\theta^2 = 0$$
- $$\theta = 2/3$$



10

## Generalization and Overfitting



11

## Example: Overfitting

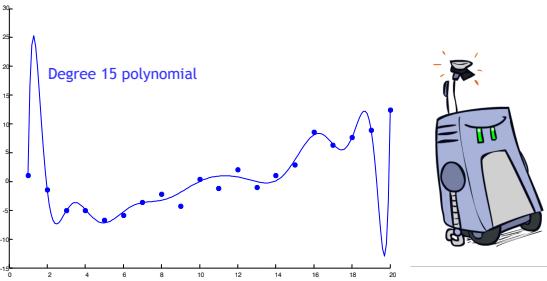
$P(\text{features}, C = 2)$	$P(\text{features}, C = 3)$
$P(C = 2) = 0.1$	$P(C = 3) = 0.1$
$P(\text{on} C = 2) = 0.8$	$P(\text{on} C = 3) = 0.8$
$P(\text{on} C = 2) = 0.1$	$P(\text{on} C = 3) = 0.9$
$P(\text{off} C = 2) = 0.1$	$P(\text{off} C = 3) = 0.7$
$P(\text{on} C = 2) = 0.01$	$P(\text{on} C = 3) = 0.0$

2 wins!!



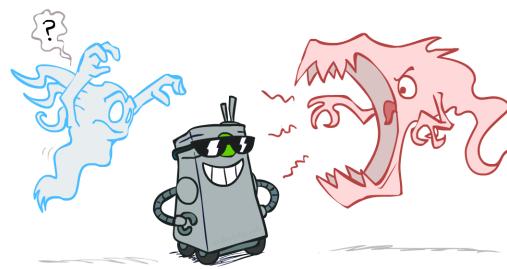
12

## Overfitting



13

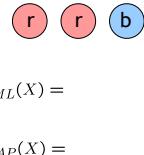
## Smoothing



14

## Laplace Smoothing

- Laplace's estimate:
    - Pretend you saw every outcome once more than you actually did
- $$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$
- $$= \frac{c(x) + 1}{N + |X|}$$
- Can derive this estimate with *Dirichlet priors* (see cs281a)



15

## Laplace Smoothing

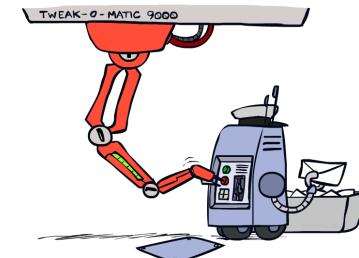
- Laplace's estimate (extended):
    - Pretend you saw every outcome  $k$  extra times
- $$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$
- $P_{LAP,0}(X) =$
- What's Laplace with  $k = 0$ ?  
k is the **strength** of the prior
- $P_{LAP,1}(X) =$
- Laplace for conditionals:
    - Smooth each condition independently:
- $$P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$$



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

## 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?
- For even better ways to estimate parameters, as well as details of the math, see cs281a, cs288



16

17

18

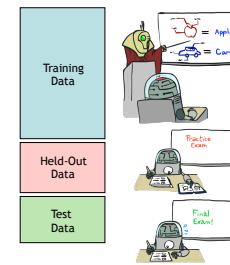
## Training and Testing



19

## Important Concepts

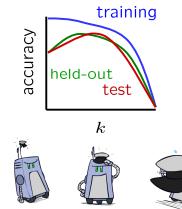
- Data: labeled instances, e.g. emails marked spam/ham
  - Training set
  - Held out set
  - Test set
- Features: attribute-value pairs which characterize each  $x$
- Experimentation cycle
  - Learn parameters (e.g. model probabilities) on training set
  - Tune hyperparameters on held-out set
  - Compute 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
  - We'll investigate overfitting and generalization formally in a few lectures



20

## Tuning on Held-Out Data

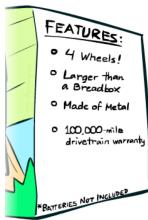
- Now we've got two kinds of unknowns
  - Parameters: the probabilities  $P(X|Y)$ ,  $P(Y)$
  - Hyperparameters: e.g. the amount / type of smoothing to do,  $k$ ,  $\alpha$
- What should we learn where?
  - Learn parameters from training data
  - Tune hyperparameters on different data
    - Why?
  -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



21

## 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
- Next class we'll talk about classifiers which let you easily add arbitrary features more easily



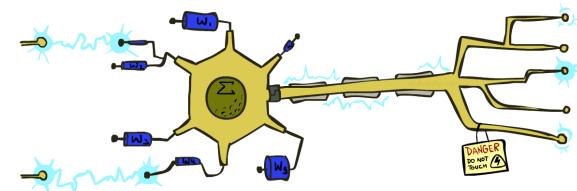
22

## 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

23

## CS 188: Artificial Intelligence Perceptrons



Instructor: Marco Alvarez --- University of Rhode Island  
 [These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley.  
 All CS188 materials are available at <http://ai.berkeley.edu>.]

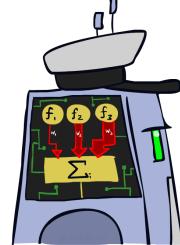
24

## Error-Driven Classification



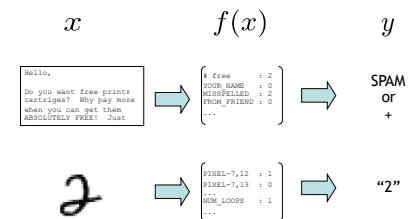
25

## Linear Classifiers



26

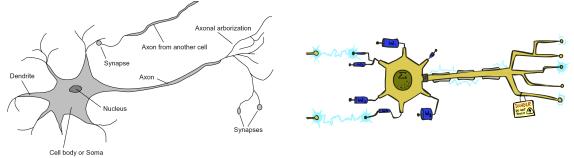
## Feature Vectors



27

## Some (Simplified) Biology

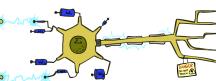
- Very loose inspiration: human neurons



28

## Linear Classifiers

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



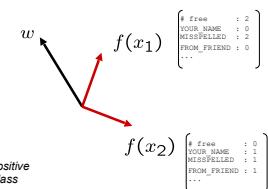
- If the activation is:
  - Positive, output +1
  - Negative, output -1

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

29

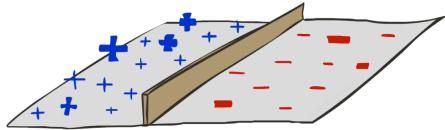
## Weights

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



30

## Decision Rules



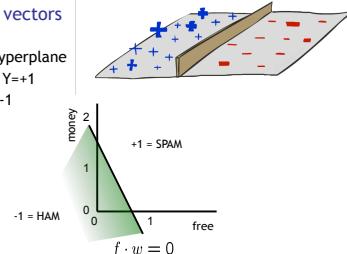
31

## 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
...



32

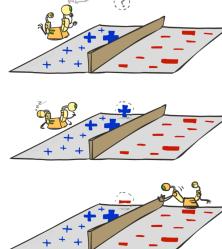
## Weight Updates



33

## 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



34

## Learning: Binary Perceptron

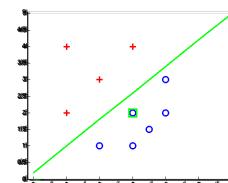
- Start with weights = 0
- For each training instance:
  - Classify with current weights
  - $y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$
  - If correct (i.e.,  $y=y^*$ ), no change!
  - If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if  $y^*$  is -1.

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

35

## Examples: Perceptron

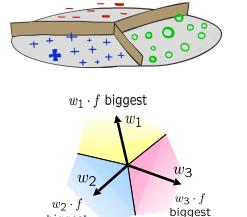
- Separable Case



36

## Multiclass Decision Rule

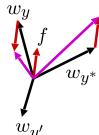
- If we have multiple classes:
  - A weight vector for each class:
  - Score (activation) of a class  $y$ :  
 $w_y \cdot f(x)$
  - Prediction highest score wins:  
 $y = \arg \max_y w_y \cdot f(x)$



37

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

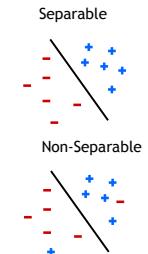


38

## 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

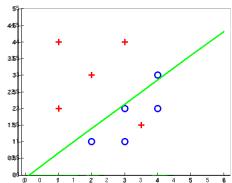
$$\text{mistakes} < \frac{k}{\delta^2}$$



39

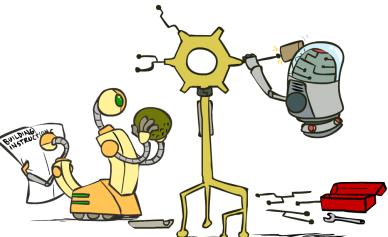
## Examples: Perceptron

- Non-Separable Case



40

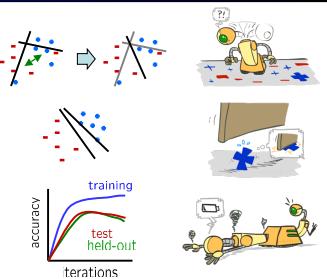
## Improving the Perceptron



41

## 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



42

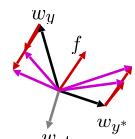
## Fixing the Perceptron

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

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

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

\*Margin Infused Relaxed Algorithm



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

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

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

43

## Minimum Correcting Update

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

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

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

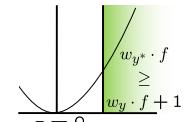
$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$

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

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

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

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



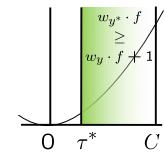
$\min_{\tau=0}$ , or would not have made an error, so  $\min$  will be where equality holds

44

## 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  $C$

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

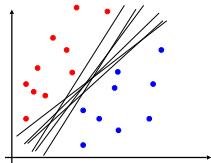


- Corresponds to an optimization that assumes non-separable data
- Usually converges faster than perceptron
- Usually better, especially on noisy data

45

## Linear Separators

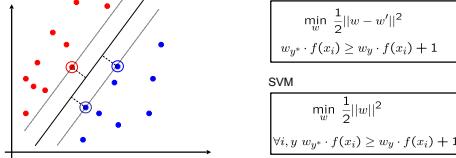
- Which of these linear separators is optimal?



46

## 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



47

## 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

48