

# Naïve Bayes Classifiers and Logistic Regression

Doug Downey

Northwestern EECS 349

Winter 2014

# Naïve Bayes Classifiers

- Combines all ideas we've covered
  - Conditional Independence
  - Bayes' Rule
  - Statistical Estimation
  - Bayes Nets
- ...in a simple, yet accurate classifier
  - Classifier: Function  $f(\mathbf{x})$  from  $\mathbf{X} = \{<x_1, \dots, x_d>\}$  to *Class*
  - E.g.,  $\mathbf{X} = \{<\text{GRE, GPA, Letters}>\}$ , *Class* = {yes, no, wait}

# Probability => Classification (1 of 2)

- Classification Task:
  - Learn function  $f(\mathbf{x})$  from  $\mathbf{X} = \{<x_1, \dots, x_d>\}$  to *Class*
  - Given: Examples  $D = \{(\mathbf{x}, y)\}$
- Probabilistic Approach
  - Learn  $P(\text{Class} = y \mid \mathbf{X} = \mathbf{x})$  from  $D$
  - Given  $\mathbf{x}$ , pick the maximally probable  $y$

# Probability => Classification (2 of 2)

- More formally
  - $f(\mathbf{x}) = \arg \max_y P(\text{Class} = y \mid \mathbf{X} = \mathbf{x}, \theta_{\text{MAP}})$
  - $\theta_{\text{MAP}}$  : MAP parameters, learned from data
    - That is, parameters of  $P(\text{Class} = y \mid \mathbf{X} = \mathbf{x})$
  - ...we'll focus on using MAP estimate, but can also use ML or Bayesian
- Predict next coin flip? Instance of this problem
  - $X = \text{null}$
  - Given  $D = \text{hhht...tth}$ , estimate  $P(\theta \mid D)$ , find MAP
  - Predict  $\text{Class} = \text{heads}$  iff  $\theta_{\text{MAP}} > \frac{1}{2}$

# Example: Text Classification

Dear Sir/Madam,  
We are pleased to inform you of the result of the Lottery Winners International programs held on the 30/8/2004. Your e-mail address attached to ticket number: EL-23133 with serial Number: EL-123542, batch number: 8/163/EL-35, lottery Ref number: EL-9318 and drew lucky numbers 7-1-8-36-4-22 which consequently won in the 1st category, you have therefore been approved for a lump sum pay out of US\$1,500,000.00 (One Million, Five Hundred Thousand United States dollars)

• SPAM

NOT SPAM?

# Representation

- $X$  = document
- Estimate  $P(\text{Class} = \{\text{spam, non-spam}\} \mid X)$
- Question: how to represent  $X$ ?
  - One dimension for each possible e-mail, i.e. possible permutation of words?
    - No.
  - Lots of possibilities, common choice: “bag of words”

Dear Sir/Madam,  
We are pleased to inform you of the result of the Lottery Winners International programs held on the 30/8/2004. Your e-mail address attached to ticket number: EL-23133 with serial Number: EL-123542, batch number: 8/163/EL-35, lottery Ref number: EL-9318 and drew lucky numbers 7-1-8-36-4-22 which consequently won in the 1st category, you have therefore been approved for a lump sum pay out of US\$1,500,000.00 (One Million, Five Hundred Thousand United States dollars)

...

Sir	1
Lottery	10
Dollars	7
With	38
...	

# Bag of Words

- Ignores Word Order, i.e.
  - No emphasis on title
  - No compositional meaning (“Cold War” -> “cold” and “war”)
  - Etc.
  - But, massively reduces dimensionality/complexity
- Still and all...
  - Recording presence or absence of a 100,000-word vocab entails  $2^{100,000}$  distinct vectors

# Naïve Bayes Classifiers



- $P(\text{Class} | \mathbf{X})$  for  $|\text{Val}(\mathbf{X})| = 2^{100,000}$  requires  $2^{100,000}$  parameters
  - Problematic.
- Bayes' Rule:
$$P(\text{Class} | \mathbf{X}) = P(\mathbf{X} | \text{Class}) P(\text{Class}) / P(\mathbf{X})$$
- Assume presence of word  $i$  is independent of all other words given  $\text{Class}$ :
$$P(\text{Class} | \mathbf{X}) = \prod_i P(w_i | \text{Class}) P(\text{Class}) / P(\mathbf{X})$$
- Now only 200,001 parameters for  $P(\text{Class} | \mathbf{X})$



# Naïve Bayes Assumption

- Features are conditionally independent given class
  - *Not*  $P(\text{“Republican”}, \text{“Democrat”}) = P(\text{“Republican”})P(\text{“Democrat”})$   
but instead  
$$P(\text{“Republican”}, \text{“Democrat”} \mid \text{Class} = \text{Politics}) =$$
$$P(\text{“Republican”} \mid \text{Class} = \text{Politics})P(\text{“Democrat”} \mid \text{Class} = \text{Politics})$$
- Still, an absurd assumption
  - (“Lottery”  $\perp$  “Winner” | SPAM)? (“lunch”  $\perp$  “noon” | Not SPAM)?
- But: offers massive tractability advantages and works quite well in practice
  - Lesson: Overly strong independence assumptions sometimes allow you to build an accurate model where you otherwise couldn’t

# Getting the parameters from data

- Parameters  $\theta = \langle \theta_{ij} = P(w_i | \text{Class} = j) \rangle$
- Maximum Likelihood: Estimate  $P(w_i | \text{Class} = j)$  from  $D$  by counting
  - Fraction of documents in class  $j$  containing word  $i$
  - But if word  $i$  never occurs in class  $j$ ?
- Commonly used MAP estimate:
  - $$\frac{(\# \text{ docs in class } j \text{ with word } i) + 1}{(\# \text{ docs in class } j) + |V|}$$

# Caveats

- Naïve Bayes effective as a *classifier*
- **Not** as effective in producing probability estimates
  - $\prod_i P(w_i | Class)$  pushes estimates toward 0 or 1
- In practice, numerical underflow is typical at classification time
  - Compare sum of logs instead of product

# Discriminative vs. Generative training

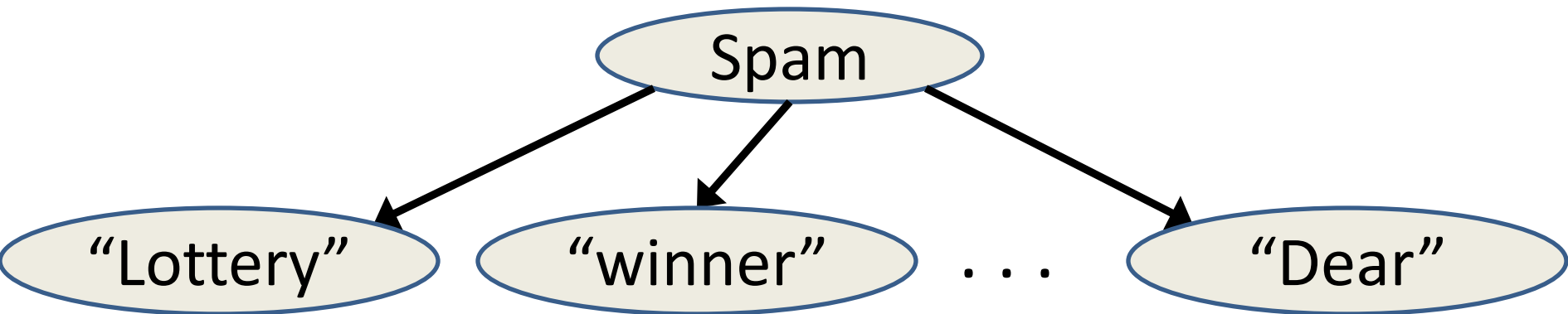
- Say our graph  $G$  has variables  $\mathbf{X}$ ,  $\mathbf{Y}$
- Previous method learns  $P(\mathbf{X}, \mathbf{Y})$
- But often, the only inferences we care about are of form  $P(\mathbf{Y} | \mathbf{X})$ 
  - $P(\textit{Disease} | \textit{Symptoms} = \mathbf{e})$
  - $P(\textit{StockMarketCrash} | \textit{RecentPriceActivity} = \mathbf{e})$

# Discriminative vs. Generative training

- Learning  $P(\mathbf{X}, \mathbf{Y})$ : **generative** training
  - Learned model can “generate” the full data  $\mathbf{X}, \mathbf{Y}$
- Learning only  $P(\mathbf{Y} | \mathbf{X})$ : **discriminative** training
  - Model **can't** assign probs. to  $\mathbf{X}$  – only  $\mathbf{Y}$  given  $\mathbf{X}$
- Idea: Only model what we care about
  - Don't “waste data” on params irrelevant to task
  - Side-step false independence assumptions in training (example to follow)

# Generative Model Example

- Naïve Bayes model
  - $Y$  binary {1=spam, 0=not spam}
  - $\mathbf{X}$  an  $n$ -vector: message has word (1) or not (0)
  - Re-write  $P(Y | \mathbf{X})$  using Bayes Rule, apply Naïve Bayes assumption
  - $2n + 1$  parameters, for  $n$  observed variables

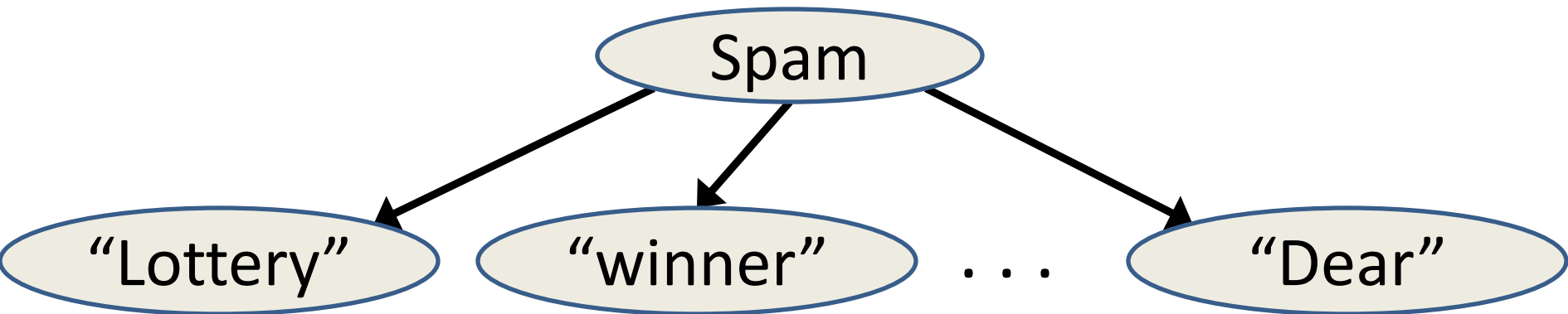


# Generative => Discriminative (1 of 3)

- But  $P(Y | \mathbf{X})$  can be written more compactly

$$P(Y | \mathbf{X}) = \frac{1}{1 + \exp(w_0 + w_1 x_1 + \dots + w_n x_n)}$$

- Total of  $n + 1$  parameters  $w_i$



# Generative => Discriminative (2 of 3)

- One way to do conversion (vars binary):

$$\exp(w_0) = \frac{P(Y = 0) P(X_1=0 | Y=0) P(X_2=0 | Y=0) \dots}{P(Y = 1) P(X_1=0 | Y=1) P(X_2=0 | Y=1) \dots}$$

for  $i > 0$ :

$$\exp(w_i) = \frac{P(X_i=0 | Y=1) P(X_i=1 | Y=0)}{P(X_i=0 | Y=0) P(X_i=1 | Y=1)}$$



# Generative => Discriminative (3 of 3)

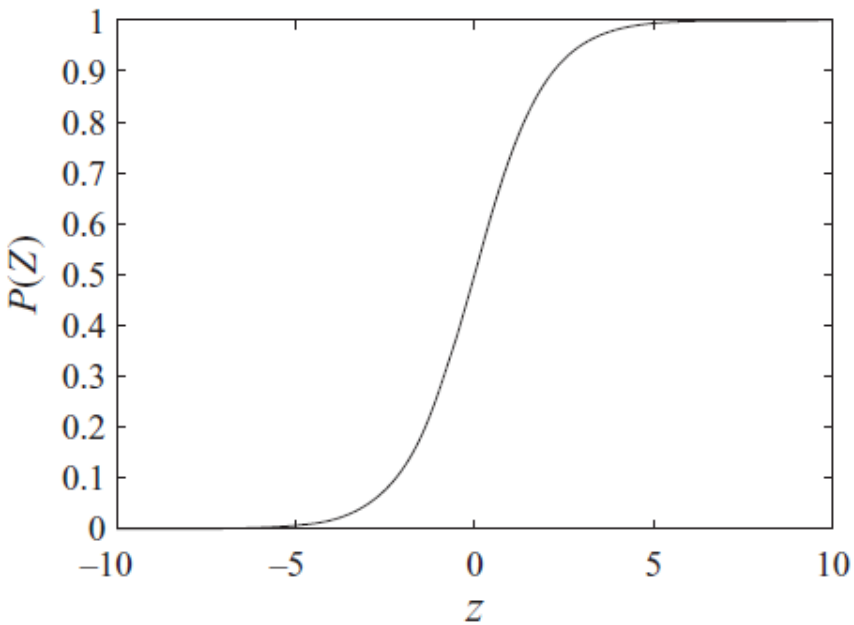
- We reduced  $2n + 1$  parameters to  $n + 1$ 
  - Bias vs. Variance arguments says this must be better, right?
- Not exactly. If we construct  $P(Y | \mathbf{X})$  to be equivalent to Naïve Bayes (as before)
  - then it's...equivalent to Naïve Bayes
- Idea: optimize the  $n + 1$  parameters directly, using training data

# Discriminative Training

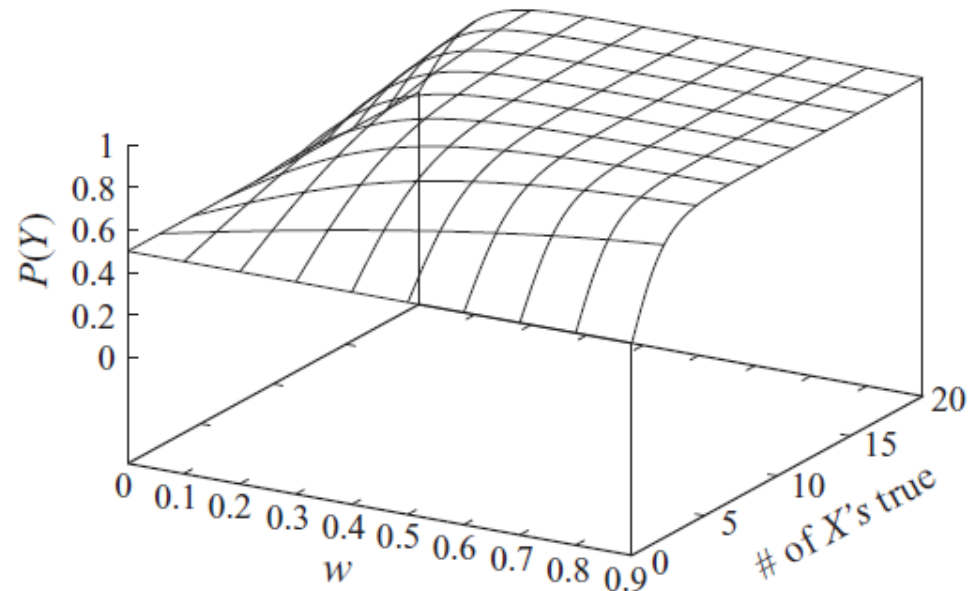
- In our example:

$$P(Y | \mathbf{X}) = \frac{1}{1 + \exp(w_0 + w_1 x_1 + \dots + w_n x_n)}$$

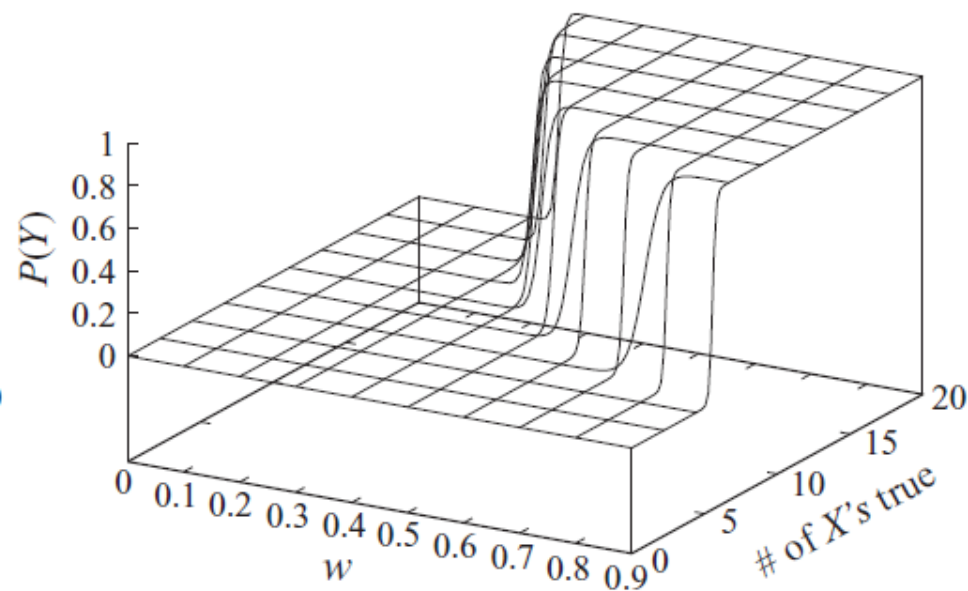
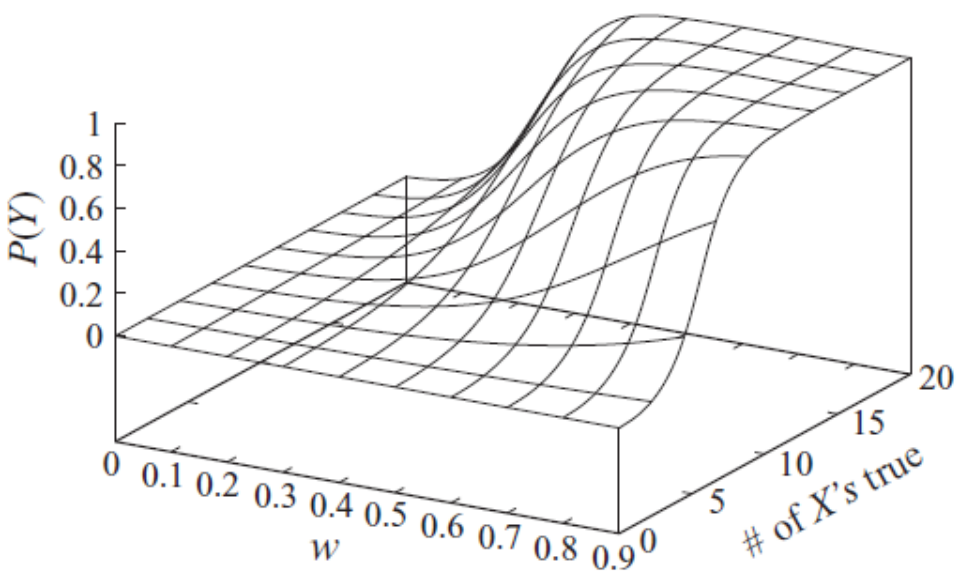
- Goal: find  $w_i$  that maximize likelihood of training data  $Y$ s given training data  $\mathbf{X}$ s
  - Known as “logistic regression”
  - Solved with gradient ascent techniques
  - A convex (actually concave) optimization problem



(a)



(b)



# Naïve Bayes vs. LR

- Naïve Bayes “trusts its assumptions” in training
- Logistic Regression doesn’t – recovers better when assumptions violated

# NB vs. LR: Example

Training Data

SPAM	Lottery	Winner	Lunch	Noon
1	1	1	0	0
1	1	1	1	1
0	0	0	1	1
0	1	1	0	1

- Naïve Bayes will classify the last example incorrectly, even after training on it!
- Whereas Logistic Regression is perfect with e.g.,  
 $w_0 = 0.1$   
 $w_{\text{lottery}} = w_{\text{winner}} = w_{\text{lunch}} = -0.2$   
 $w_{\text{noon}} = 0.4$

# Logistic Regression in practice

- Can be employed for any numeric variables  $X_i$ 
  - or for other variable types, by converting to numeric (e.g. indicator) functions
- “Regularization” plays the role of priors in Naïve Bayes
- Optimization tractable, but (way) more expensive than counting (as in Naïve Bayes)

# Discriminative Training

- Naïve Bayes vs. Logistic Regression one illustrative case
- Applicable more broadly, whenever queries  $P(\mathbf{Y} | \mathbf{X})$  known *a priori*

Data Set	MNB-FM	SFE	MNB	NBEM	LProp	Logist
Apte (10)	0.306	0.271	<b>0.336</b>	0.306	0.245	0.208
Apte (100)	<b>0.554</b>	0.389	0.222	0.203	0.263	0.330
Apte (1k)	<b>0.729</b>	0.614	0.452	0.321	0.267	0.702
Amzn (10)	<b>0.542</b>	0.524	0.508	0.475	0.470*	0.499
Amzn (100)	<b>0.587</b>	0.559	0.456	0.456	0.498*	0.542
Amzn (1k)	0.687	0.611	0.465	0.455	0.539*	<b>0.713</b>
RCV1 (10)	<b>0.494</b>	0.477	0.387	0.485	-	0.272
RCV1 (100)	<b>0.677</b>	0.613	0.337	0.470	-	0.518
RCV1 (1k)	0.772	0.735	0.408	0.491	-	<b>0.774</b>

\* Limited to 5 of 10 Amazon categories