

# CS 112: Modeling Uncertainty in Information Systems

Prof. Jenn Wortman Vaughan

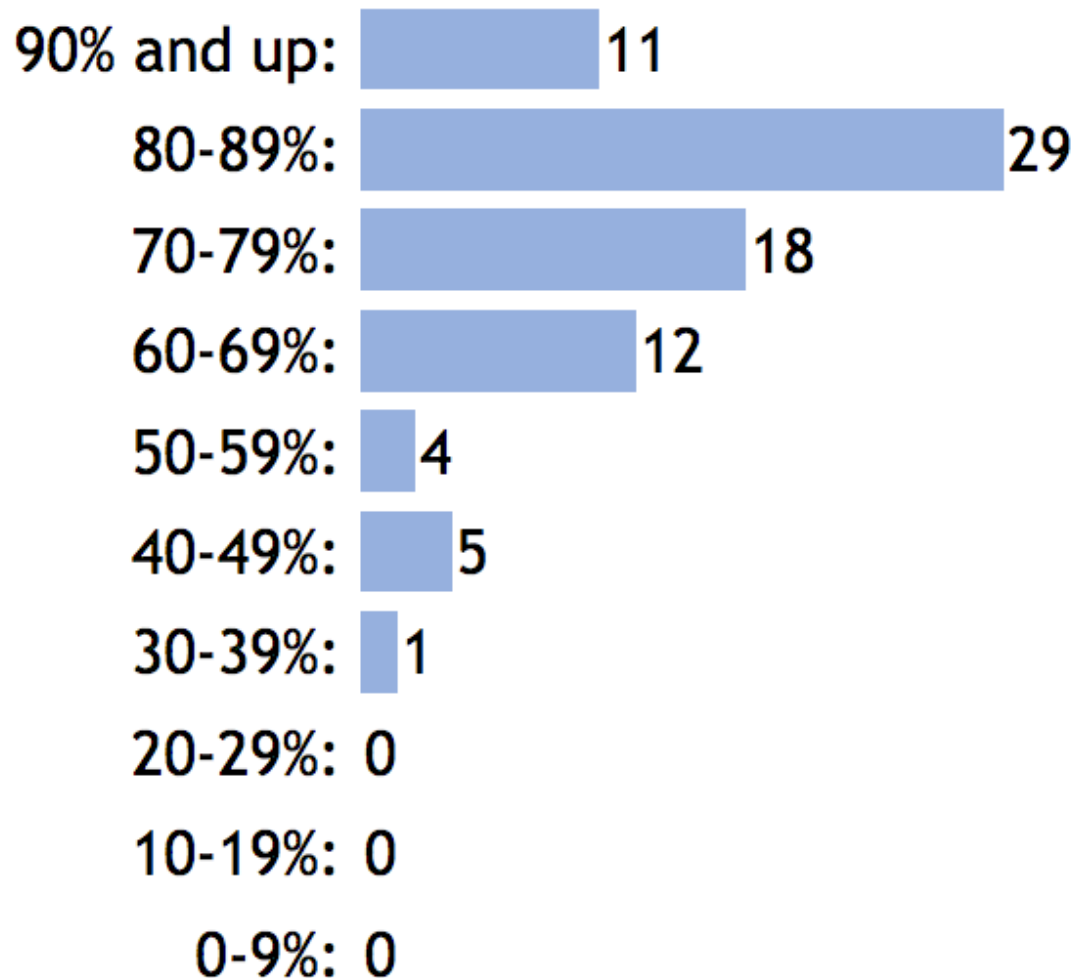
May 21, 2012

Lecture 14

# Reminders & Announcements

- Homework 4 has been posted on the website and is due on **Wednesday, May 30**
- The assignment is more open-ended than Homework 3

# Midterms



Mean: 76.4

Median: 79.5

Max: 98

Min: 37

Regrade requests should be submitted within one week

# Today

- A “Naive Bayes” classifier for spam filtering  
(or, everything you need to know for homework 4)

# Hypothesis Testing

- The **maximum likelihood (ML) hypothesis** is the hypothesis that makes the data most likely

$$H^{\text{ML}} = \operatorname{argmax}_i P(D | H_i)$$

- The **maximum a posteriori (MAP) hypothesis** is the hypothesis with the maximum posterior probability

$$H^{\text{MAP}} = \operatorname{argmax}_i P(H_i | D) = \operatorname{argmax}_i P(D | H_i) P(H_i)$$

# Parameter Estimation

Suppose that we would like to estimate the **unknown bias**  $p$  of a coin based on observations of the outcomes  $X_1, \dots, X_n$  of  $n$  independent tosses of the coin

What is the maximum likelihood estimate of  $p$ ?

# Parameter Estimation

Suppose that we would like to estimate the **unknown bias**  $p$  of a coin based on observations of the outcomes  $X_1, \dots, X_n$  of  $n$  independent tosses of the coin

What is the maximum likelihood estimate of  $p$ ?

$$\frac{1}{n} \sum_{i=1}^n X_i$$

# Parameter Estimation

Suppose that we would like to estimate the **unknown bias**  $p$  of a coin based on observations of the outcomes  $X_1, \dots, X_n$  of  $n$  independent tosses of the coin

What is the maximum likelihood estimate of  $p$ ?

$$\frac{1}{n} \sum_{i=1}^n X_i = \frac{\# \text{ times we observed heads}}{n}$$



# Classifying Spam

- Suppose that we would like to classify a new email message as either spam or not spam

# Classifying Spam

- Suppose that we would like to classify a new email message as either spam or not spam
- We can represent the email message as a **vector of features**, e.g., presence or absence of the word “cash”, or presence or absence of the recipient’s name

# Classifying Spam

- Suppose that we would like to classify a new email message as either spam or not spam
- We can represent the email message as a **vector of features**, e.g., presence or absence of the word “cash”, or presence or absence of the recipient’s name
- We can use **previously labeled emails** (also represented as feature vectors) to build a probabilistic model

# Classifying Spam

- Suppose that we would like to classify a new email message as either spam or not spam
- We can represent the email message as a **vector of features**, e.g., presence or absence of the word “cash”, or presence or absence of the recipient’s name
- We can use **previously labeled emails** (also represented as feature vectors) to build a probabilistic model
- Using this model, we can calculate a MAP hypothesis to classify the new email

# What Do We Know?

- When classifying a single email, what are our hypotheses?

# What Do We Know?

- When classifying a single email, what are our hypotheses?
- What is our observed data?

# What Do We Know?

- When classifying a single email, what are our hypotheses?
- What is our observed data?
- The MAP hypothesis is the one that maximizes

$$P(F_1 = f_1, \dots, F_d = f_d | H_i) P(H_i)$$

# What Do We Know?

- When classifying a single email, what are our hypotheses?
- What is our observed data?
- The MAP hypothesis is the one that maximizes

$$P(F_1 = f_1, \dots, F_d = f_d | H_i) P(H_i)$$

Our goal: Use the labeled emails to estimate this value for each hypothesis  $H_i$  so that we can find the MAP hypothesis



# Step 1: Estimate the Prior

How can we estimate  $P(H_i)$ ?

# Step 1: Estimate the Prior

How can we estimate  $P(H_i)$ ?

- For each previously labeled email  $j$ , let

$$X_j = \begin{cases} 1, & \text{if email is spam} \\ 0, & \text{otherwise} \end{cases}$$

# Step 1: Estimate the Prior

How can we estimate  $P(H_i)$ ?

- For each previously labeled email  $j$ , let

$$X_j = \begin{cases} 1, & \text{if email is spam} \\ 0, & \text{otherwise} \end{cases}$$

- If our emails are i.i.d., these are Bernoulli random variables with **unknown parameter  $P(H_1)$**  – can estimate this unknown parameter using **maximum likelihood**

# Step 1: Estimate the Prior

How can we estimate  $P(H_i)$ ?

- For each previously labeled email  $j$ , let

$$X_j = \begin{cases} 1, & \text{if email is spam} \\ 0, & \text{otherwise} \end{cases}$$

- If our emails are i.i.d., these are Bernoulli random variables with **unknown parameter  $P(H_1)$**  – can estimate this unknown parameter using **maximum likelihood**

$$P(H_1) = \frac{1}{n} \sum_{j=1}^n X_j$$

## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

- We could use maximum likelihood here too...

# ML for Multinomials

Suppose we would like to estimate the unknown parameters  $p_1, \dots, p_k$  of a **multinomial** (e.g., rolls of a die) based on  $n$  independent observations

What is the maximum likelihood estimate of each  $p_j$ ?

# ML for Multinomials

Suppose we would like to estimate the unknown parameters  $p_1, \dots, p_k$  of a **multinomial** (e.g., rolls of a die) based on  $n$  independent observations

What is the maximum likelihood estimate of each  $p_j$ ?

$$\frac{\text{\# times we observed outcome } j}{n}$$



## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

- We could use maximum likelihood here too... Why is this a bad idea?

## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

- We could use maximum likelihood here too... Why is this a bad idea?
- Instead, we make the **Naive Bayes assumption** that all feature values are conditionally independent given  $H_i$

## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

- We could use maximum likelihood here too... Why is this a bad idea?
- Instead, we make the **Naive Bayes assumption** that all feature values are conditionally independent given  $H_i$

$$P(F_1 = f_1, \dots, F_d = f_d | H_i) = \prod_{j=1}^d P(F_j = f_j | H_i)$$

## Step 2: Make Some Assumptions

How can we estimate  $P(F_1 = f_1, \dots, F_d = f_d | H_i)$  from data?

- We could use maximum likelihood here too... Why is this a bad idea?
- Instead, we make the **Naive Bayes assumption** that all feature values are conditionally independent given  $H_i$

$$P(F_1 = f_1, \dots, F_d = f_d | H_i) = \prod_{j=1}^d P(F_j = f_j | H_i)$$

- For each  $i$ , have to estimate  $d$  parameters instead of  $2^d - 1$

## Step 3: Estimate the Feature Probabilities

How can we estimate  $P(F_j = f_j | H_i)$  from data?

## Step 3: Estimate the Feature Probabilities

How can we estimate  $P(F_j = f_j | H_i)$  from data?

- We can use **maximum likelihood** again

## Step 3: Estimate the Feature Probabilities

How can we estimate  $P(F_j = f_j | H_i)$  from data?

- We can use **maximum likelihood** again

$$P(F_j = f_j | H_i) = \frac{\text{\# examples w/ feature } j = f_j \text{ and label } = H_i}{\text{\# examples w/ label } = H_i}$$

## Step 3: Estimate the Feature Probabilities

How can we estimate  $P(F_j = f_j | H_i)$  from data?

- We can use **maximum likelihood** again

$$P(F_j = f_j | H_i) = \frac{\text{\# examples w/ feature } j = f_j \text{ and label } = H_i}{\text{\# examples w/ label } = H_i}$$

Problem: What happens if we don't observe an example with a particular feature value and label together??



## Step 3: Estimate the Feature Probabilities

How can we estimate  $P(F_j = f_j | H_i)$  from data?

- We can use **maximum likelihood with smoothing**

$$P(F_j = f_j | H_i) = \frac{(\# \text{ examples w/ feature } j = f_j \text{ and label } = H_i) + 1}{(\# \text{ examples w/ label } = H_i) + 2}$$

# The Naive Bayes Classifier

- For each  $i$ , calculate

$$P(F_1 = f_1, \dots, F_d = f_d \mid H_i)P(H_i)$$

# The Naive Bayes Classifier

- For each  $i$ , calculate

$$P(F_1 = f_1, \dots, F_d = f_d \mid H_i)P(H_i)$$

Naive Bayes Independence Assumption

# The Naive Bayes Classifier

- For each  $i$ , calculate

$$\prod_{j=1}^d P(F_j = f_j | H_i) P(H_i)$$

# The Naive Bayes Classifier

- For each  $i$ , calculate

$$\prod_{j=1}^d P(F_j = f_j | H_i) P(H_i)$$

ML estimate



# The Naive Bayes Classifier

- For each  $i$ , calculate

$$\prod_{j=1}^d P(F_j = f_j | H_i) P(H_i)$$

ML estimate using  
smoothing

ML estimate

# The Naive Bayes Classifier

- For each  $i$ , calculate

$$\prod_{j=1}^d P(F_j = f_j | H_i) P(H_i)$$

ML estimate using smoothing

ML estimate

- The MAP hypothesis is the one that maximizes this

# Exercise: Classifying Email

“Jenn”	“cash”	“viagra”	spam
1	0	0	0
1	1	0	0
0	0	0	0
0	0	1	1
0	1	0	1
1	0	0	???