

Lecture 4: Spam Filtering

Siddharth Garg
sg175@nyu.edu

Spam Detection: Features

- Recall features used in the UCI Spam database

48 continuous real [0,100] attributes of type word_freq_WORD

- Even easier way to encode features:

- $x_i = 1$ if term i appears in a document; 0 otherwise

- Boolean features

- Assume M Boolean features, $x = (x_1, x_2, \dots, x_M)$

- We want to map this M -dimensional Boolean input to a Boolean output y

- *Thoughts?*

- Instead of using LR (or SVM) we will start with an even simpler

approach referred to as "Naive Bayes"

Ref: Metzis, Vangelis, John Andrioutsos, and Georgios Paliouras. "Spam filtering with naive bayes-which naive bayes?." In CEAS, vol. 17, pp. 28-69. 2006.

Naiive Bayes for Spam Filtering

- Assume M Boolean feature, $x = (x_1, x_2, \dots, x_M)$
- Each email is either $\{s=\text{spam}, l=\text{legit}\}$
- We begin by computing:

“Bernoulli Naiive Bayes”

$$P\{spam \mid x\} = \frac{P\{x \mid spam\} * P\{spam\}}{P\{x\}}$$

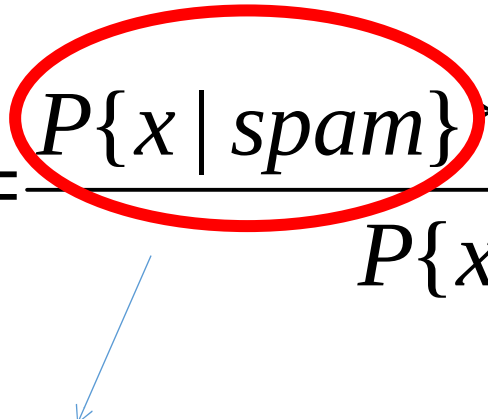
Bayes Rule

$$P\{A \cap B\} = P\{A \mid B\} * P\{B\}$$

Ref: Metsis, Vangelis, Ion Androutsopoulos, and Georgios Paliouras. "Spam filtering with naive bayes- which naive bayes?." In *CEAS*, vol. 17, pp. 28-69. 2006.

Naiive Bayes for Spam Filtering

- We begin by computing:

$$P\{spam | x\} = \frac{P\{x | spam\} * P\{spam\}}{P\{x\}}$$


$$P\{x_1, x_2, \dots, x_M | spam\} = P\{x_1 | spam\} * P\{x_2 | spam\} * \dots * P\{x_M | spam\}$$

Assuming that term occurrences are independent (given class)

Is this a reasonable assumption?

Naiive Bayes for Spam Filtering

$$P\{x_1 | spam\} * P\{x_2 | spam\} * .. * P\{x_M | spam\}$$



How do we estimate this from the training dataset?

$$P\{x_1 = 1 | spam\} = p_{i,s}$$

$$= (\text{\#Spam emails that contain term } i) / (\text{\#spam emails})$$



What happens if term i never occurred in any spam email in the

Laplacian Smoothing

$$p_{i,s} (\cancel{\text{\#Spam emails that contain term } i}) / (\cancel{\text{\#spam emails}})$$

$$= (\text{\#Spam emails that contain term } i+1) / (\text{\#spam emails}+2)$$

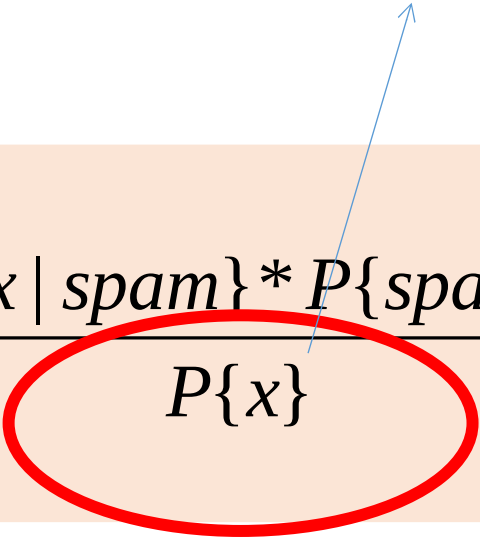
Equivalent to assuming two addition spam emails in the training dataset, of which one contains all terms and the other is empty

$$P\{x_i = 0 \mid \text{spam}\} = 1 - p_{i,s}$$

$$P\{x_1, x_2, \dots, x_M \mid \text{spam}\} = \prod_{i=1}^M p_{i,s}^{x_i} (1 - p_{i,s})^{1-x_i}$$

Naiive Bayes for Spam Filtering

$$P\{x\} = P\{spam\} * P\{x | spam\} + P\{legit\} * P\{x | legit\}$$


$$P\{spam | x\} = \frac{P\{x | spam\} * P\{spam\}}{P\{x\}} \quad \text{Vs.} \quad P\{legit | x\} = \frac{P\{x | legit\} * P\{legit\}}{P\{x\}}$$

Or:

$$P\{spam | x\} \geq threshold$$


Spam Detection: Occurences


- Recall features used in the UCI Spam database

48 continuous real [0,100] attributes of type word_freq_WORD

- Let's consider a different representation that is closer to the UCI spambase features: **Term Frequencies** (TF)
 - x_i # times term i appears in a document $x \in \mathbb{N}^M$
 - Each document is represented by $x = (x_1, x_2, \dots, x_M)$, a vector of term frequencies
 - We will again use a Naïve Bayes approach to classify documents as either spam or legit
 - **"Multinomial Naïve Bayes"**

Applying Bayes Rule

$$P\{spam | x\} = \frac{P\{x | spam\} * P\{spam\}}{P\{x\}}$$


$$P\{x_1 | spam\} * P\{x_2 | spam\} * .. * P\{x_M | spam\}$$


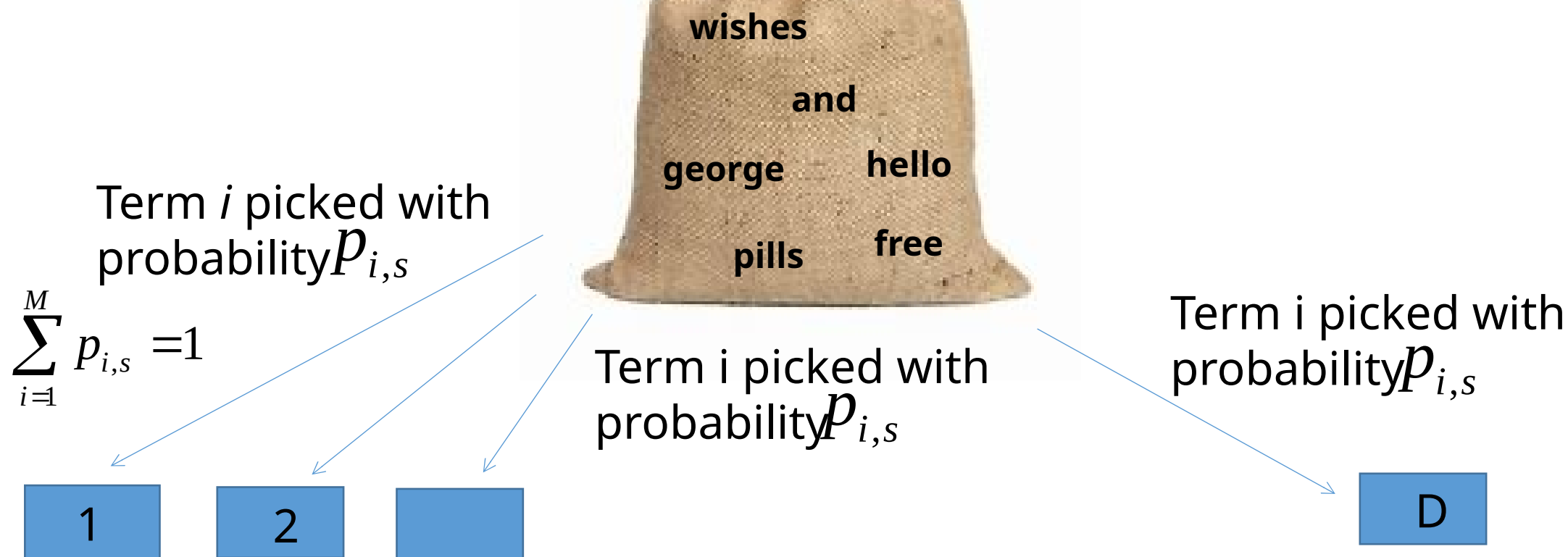
Independence
assumption shows up
again!

But how do we estimate the probability $P\{x_1 = t | spam\}$

What if there is no document in the training dataset where term 1 occurs t time

“Bag of Words” Model

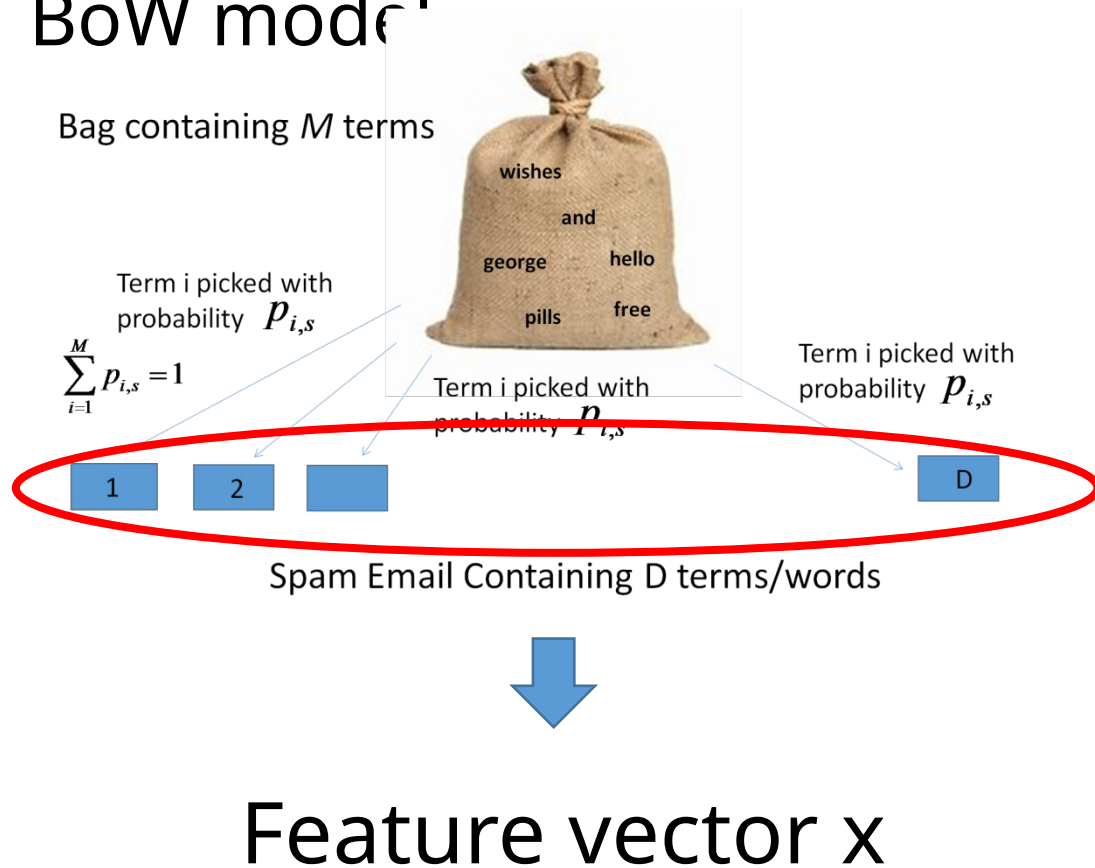
Bag containing M terms



Spam Email Containing D terms/words

Likelihood Estimation

- Say you have a spam e-mail of length D generated using the BoW model'



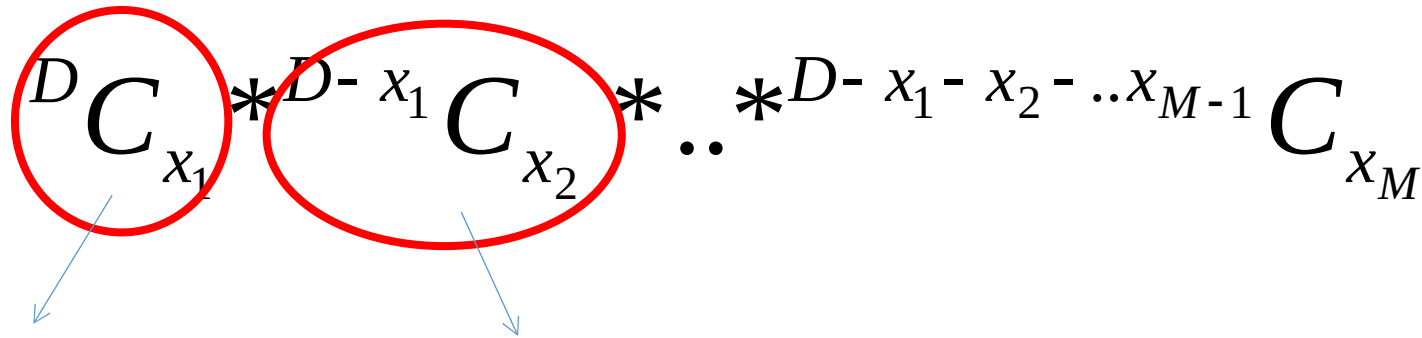
$$P\{x \mid \text{spam}, D\} = \prod_{i=1}^M (p_{i,s})^{x_i}$$

Assuming term and positional independence

Are we done?

Likelihood Estimation

- Recall that the BoW model does not keep track of the positions in which terms appear
 - We must account for all possible ways of arranging
 - x_1 instances of term 1 and
 - x_2 instances of term 2 and
 - ... x_M instances of term M into D locations

$$\binom{D}{x_1} * \binom{D - x_1}{x_2} * \dots * \binom{D - x_1 - x_2 - \dots - x_{M-1}}{x_M}$$


Choose x_1 locations
from a total of D
locations

Choose x_2 locations from
remaining $D - x_1$ locations

Likelihood Estimation

$${}^D C_{x_1} * {}^{D-x_1} C_{x_2} * \dots * {}^{D-x_1-x_2-\dots-x_{M-1}} C_{x_M} = \frac{D!}{x_1! (D-x_1)!} * \frac{(D-x_1)!}{x_2! (D-x_1-x_2)!} \dots 1$$

$$= \frac{D!}{x_1! x_2! \dots x_M!}$$

$$P\{x \mid spam, D\} = D! \prod_{i=1}^M \frac{(p_{i,s})^{x_i}}{x_i!}$$

[Typo: this should be $x_{\{i\}}$]

Note that this expression is conditioned on the length of the e-mail D . In practice, emails can be of varying lengths.

Accounting for Document Length

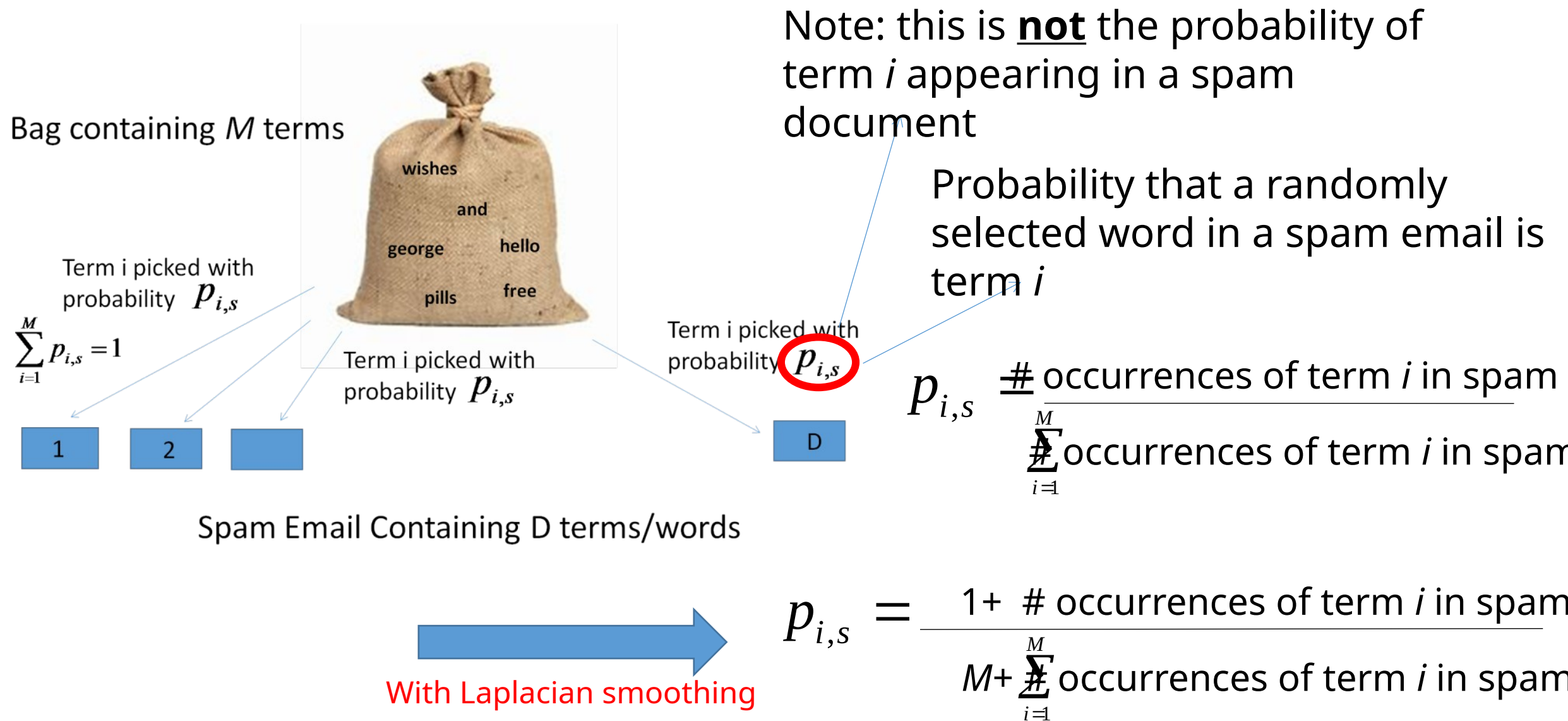
$$P\{x \mid \text{spam}\} = P\{x \mid \text{spam}, D\} P\{D \mid \text{spam}\} = P\{x \mid \text{spam}, D\} P\{D\}$$

Assume email length is independent of whether email is spam or legit.

Putting it all together:

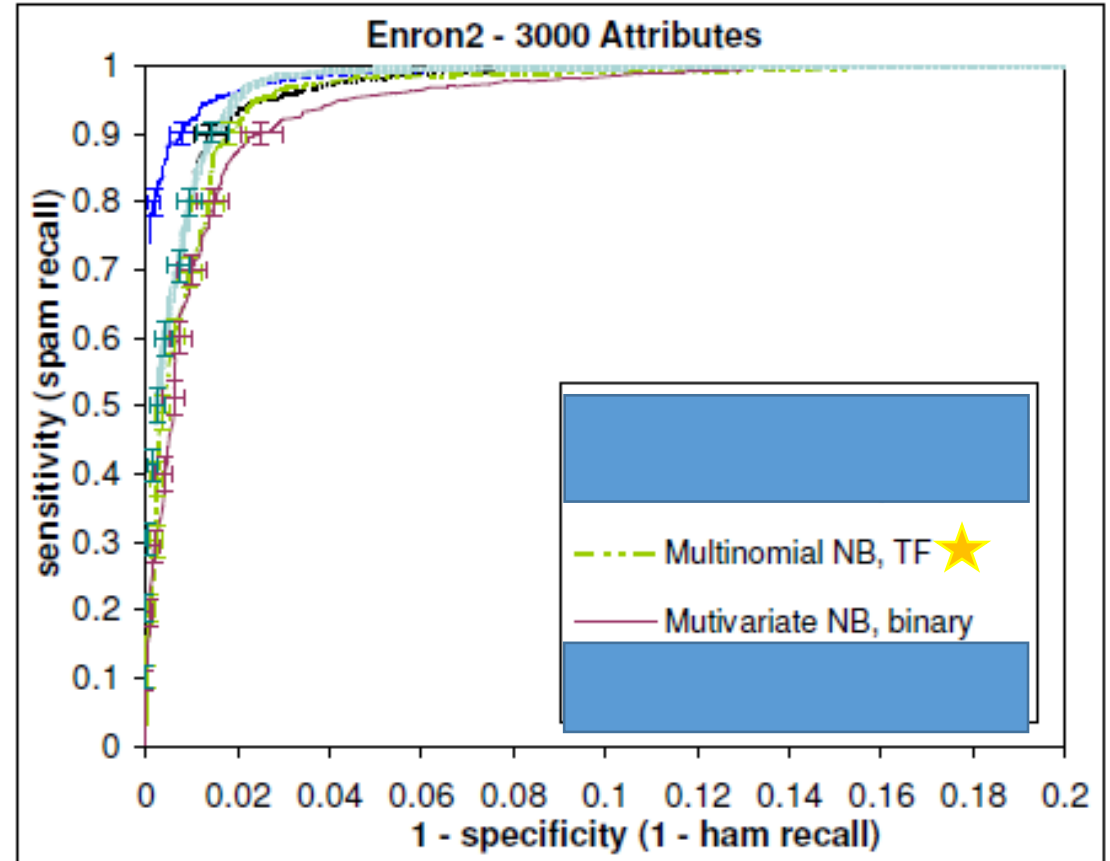
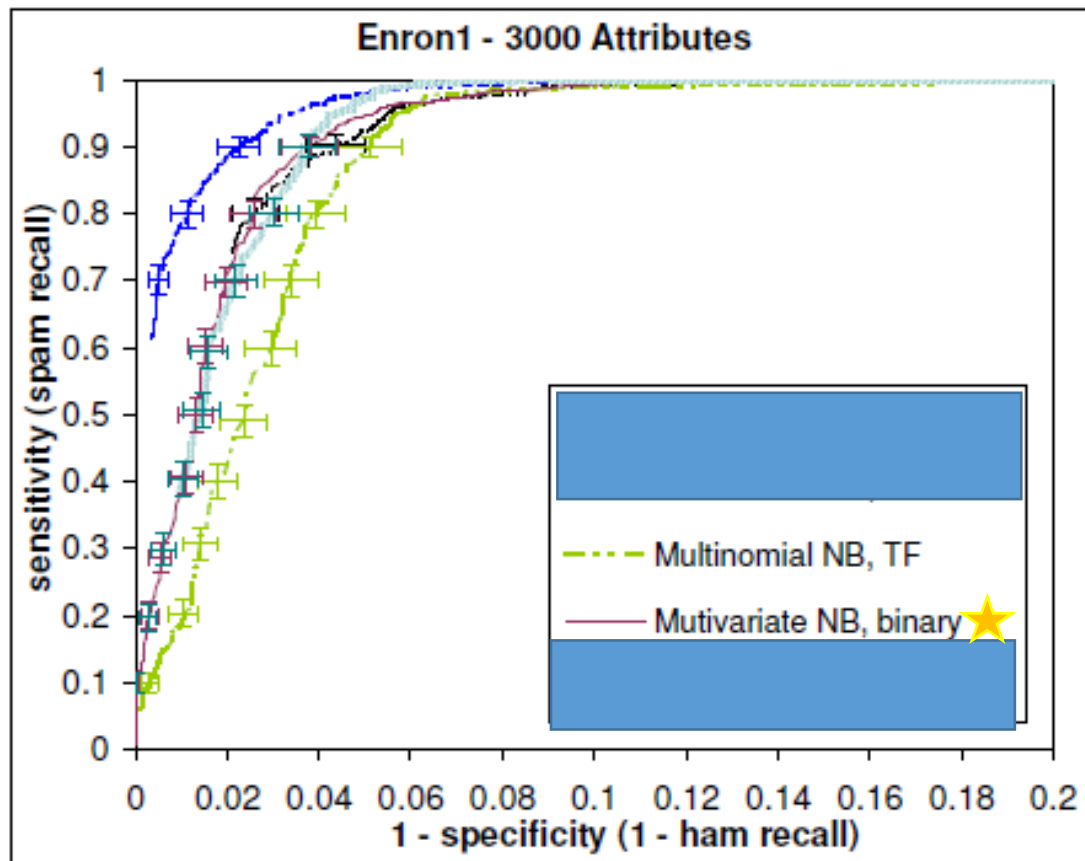
$$P\{\text{spam} \mid x\} = \frac{P\{x \mid \text{spam}, D\} \cancel{P\{D\}} P\{\text{spam}\}}{\cancel{P\{x\}}} \quad \textbf{Vs.} \quad P\{\text{legit} \mid x\} = \frac{P\{x \mid \text{legit}, D\} \cancel{P\{D\}} P\{\text{legit}\}}{\cancel{P\{x\}}}$$

Estimating Model Parameters



Bernoulli NB Vs. Multinomial NB with TF

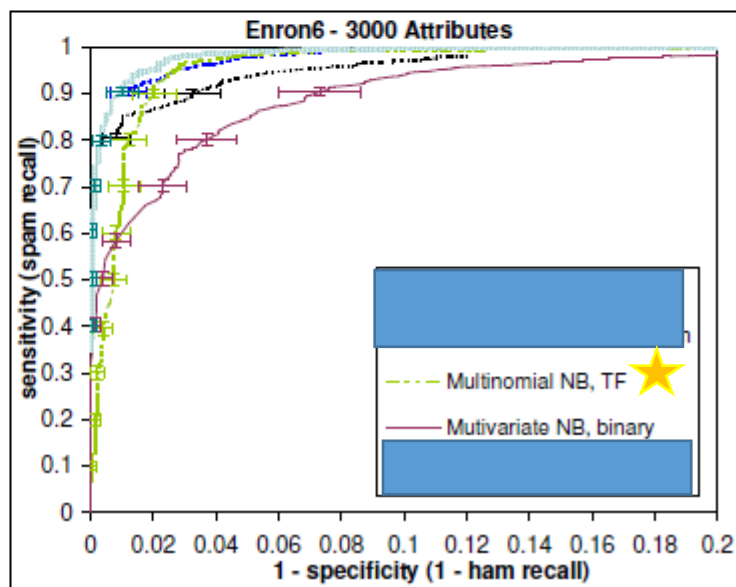
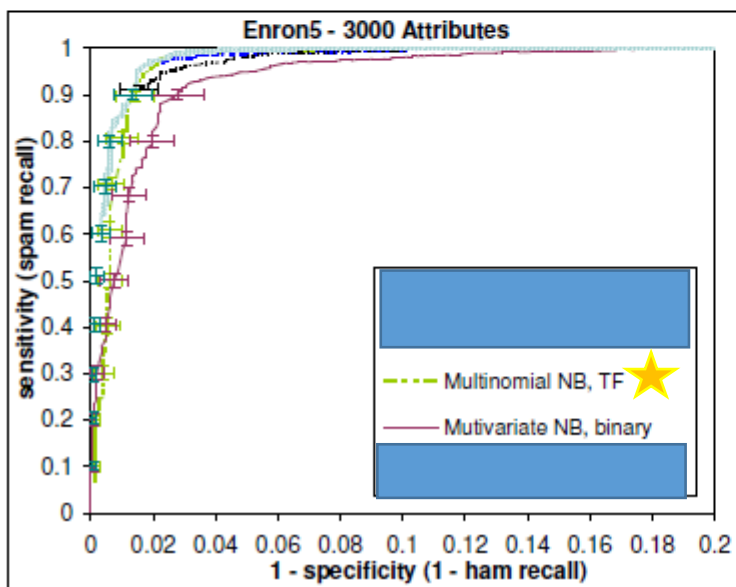
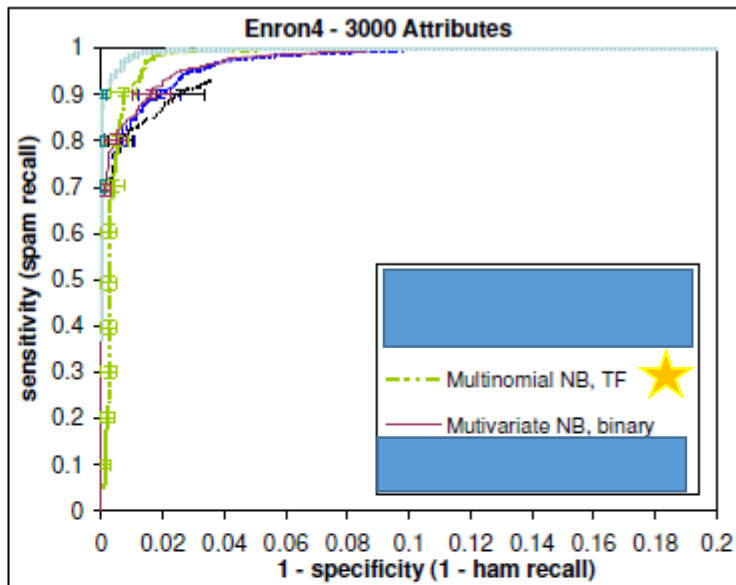
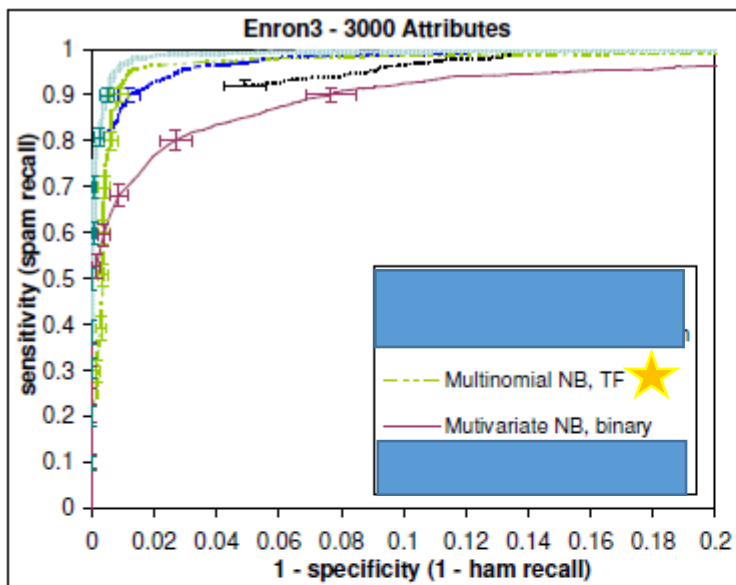
- Data for 6 different users from ENRON dataset
 - Augmented with spam emails from various sources (legit = "ham")
 - Top-3000 features selected (we will discuss feature selection soon)



% of legit emails classified as spam

Bernoulli NB Vs. Multinomial NB with

TF
% of spam emails predicted as spam



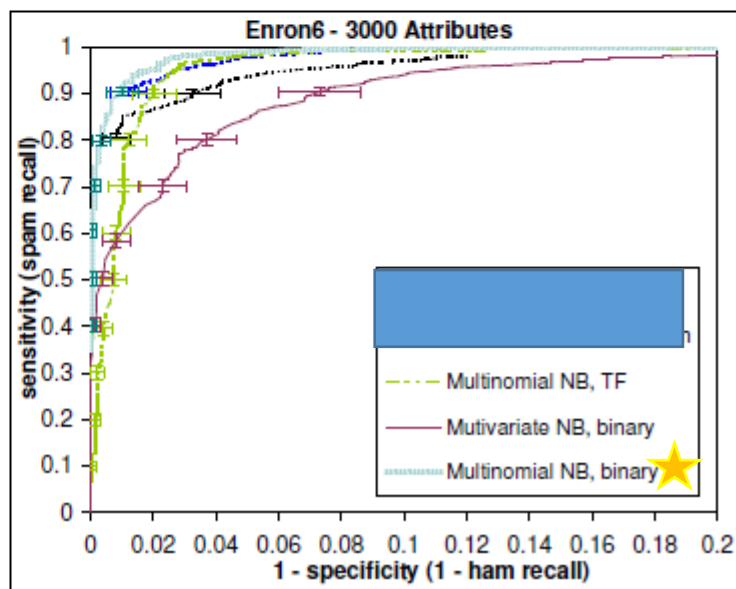
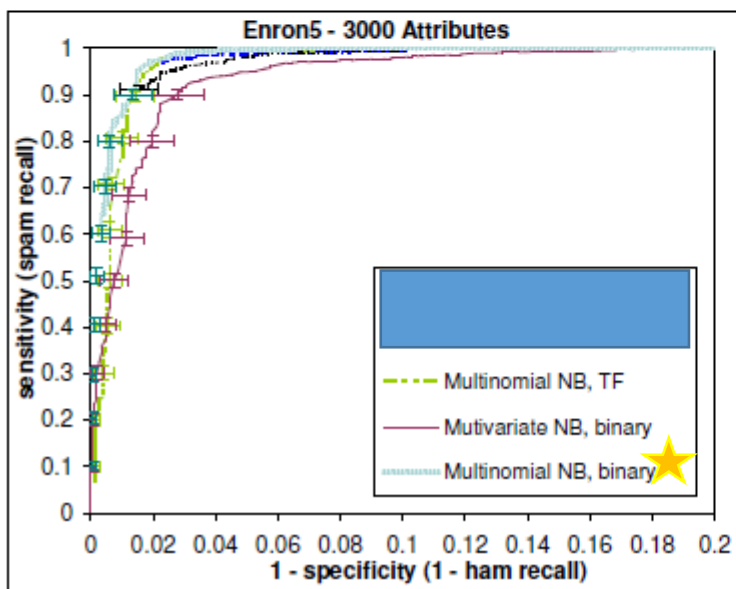
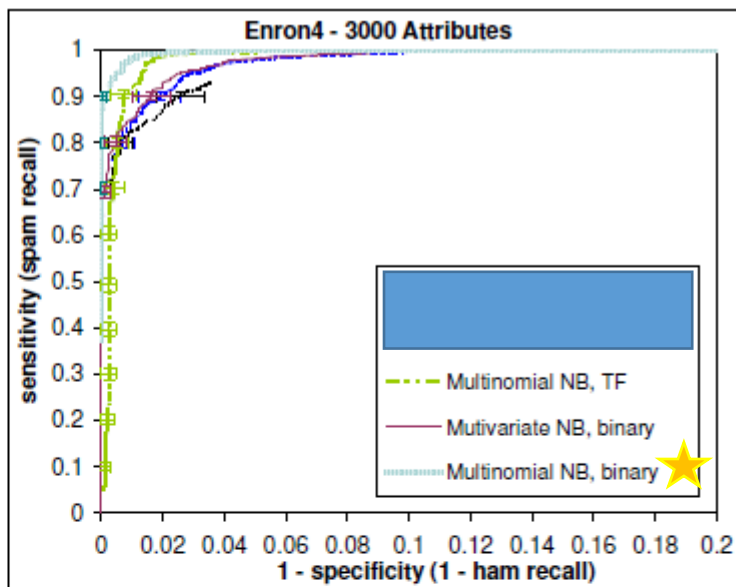
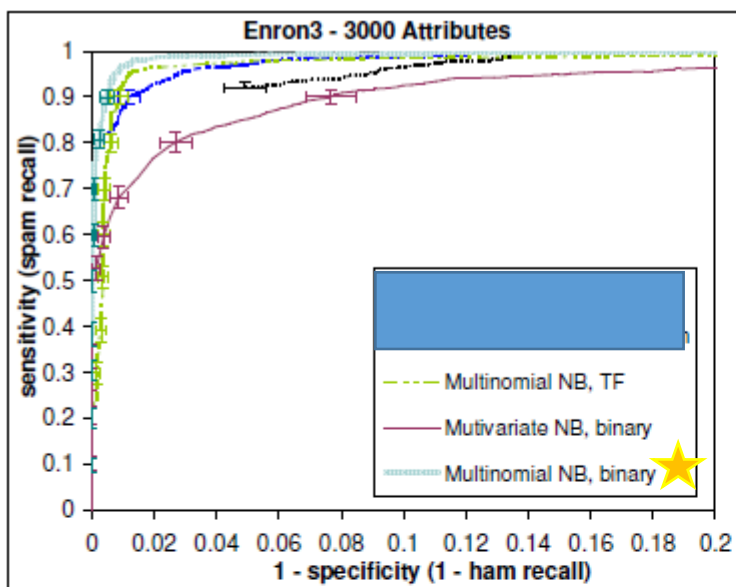
% of legit emails classified as spam

Tempted to conclude that using term frequencies instead of binary occurrences helped in spam filtering.

Is there any other reason multinomial NB with TF might have outperformed Bernoulli NB?

Bernoulli NB Vs. Multinomial NB with

TF
% of spam emails predicted as spam



Multinomial NB with Binary instead of TF features performs the best!

% of legit emails classified as spam

Multinomial NB with Binary Features

- Let $x = (x_1, x_2, \dots, x_M)$ be the TF features. Binary features are derived from the TF features as follows:

$$\bar{x} = (\bar{x}_1 = \min(1, x_1), \bar{x}_2 = \min(1, x_2), \dots, \bar{x}_M = \min(1, x_M))$$

- Transformation is applied to both the training and test data and the multinomial model is used for prediction, i.e.,

$$P\{\bar{x} | spam\} = p(D) D! \prod_{i=1}^M \frac{(p_{i,s})^{\bar{x}_i}}{\bar{x}_i!} \quad \left\{ \begin{array}{ll} p_{i,s} & \text{if } \bar{x}_i = 1 \\ 1 & \text{if } \bar{x}_i = 0 \end{array} \right.$$

[Typo: this should be $\bar{x}_{\{i\}}$]

Multinomial Vs. Bernoulli NB

Multinomial

$$P\{\bar{x} \mid spam\} = \cancel{p(D)} D! \prod_{i=1}^M (p_{i,s})^{\bar{x}_i}$$

Bernoulli

$$P\{x \mid spam\} = \prod_{i=1}^M p_{i,s}^{x_i} (1 - p_{i,s})^{1-x_i}$$

How are the two different?

1. Multinomial model ignores negative evidence
2. $p_{i,s}$ is estimated differently

$\frac{1 + \text{\# occurrences of term } i \text{ in spam}}{M + \sum_{i=1}^M \text{\# occurrences of term } i \text{ in spam}}$

$\frac{1 + \text{\#Spam emails that contain term } i}{2 + \text{\#spam emails}}$

Why Ignore Negative Evidence?

Schneider, Karl-Michael. "On word frequency information and negative evidence in Naive Bayes text classification." *Advances in Natural Language Processing*. Springer, Berlin, Heidelberg, 2004. 474-485.

Table 1. Statistics of the ling-spam corpus

	Total	Ling	Spam
Documents	2893	2412 (83.4%)	481 (16.6%)
Vocabulary	59,829	54,860 (91.7%)	11,250 (18.8%)

Vocabulary	Total		Ling		Spam	
	Words	Documents	Words	Documents	Words	Documents
Full	226.5	11.0	226.9	9.1	224.5	1.8
MI 5000	138.5	80.2	133.8	64.5	162.5	15.6
MI 500	44.0	254.5	39.6	190.9	66.2	63.7

Observation 1: >80% of words never occur in spam documents, while only 10% of words never occur in legit documents



Observation 2: On average, documents only contain a very small fraction of words from the vocabulary

For Bernoulli NB, probability of a document is mostly determined by words that do not appear in the document!

Why is Multinomial Binary Features better than Term Frequencies?

Multinomial TF assumes repeated instances of the same word occur independently, but that is not the case -> for example, if a word appears once it is more likely to appear multiple times. Therefore multinomial TF is a poor model for the underlying data.