

Text Classification – Naïve Bayes

June 17, 2016

Why Text Classification?

- Users may have ongoing information needs
 - Might want to track developments in a particular topic such as “multicore computer chips”
- The classification of documents by topic capture the generality and scope of the problem space.

Classification Problems

- Email filtering: spam / non spam
- Email foldering / tagging: Work, Friends, Family, Hobby
- Research articles by topics: Machine Learning, Data Mining, Algorithms
- Sentiment Analysis: positive / negative
- Emotion Detection: anger, happiness, joy, sadness, etc.
- Tumor: malignant / benign
- Medical diagnosis: Not ill, Cold, Flu

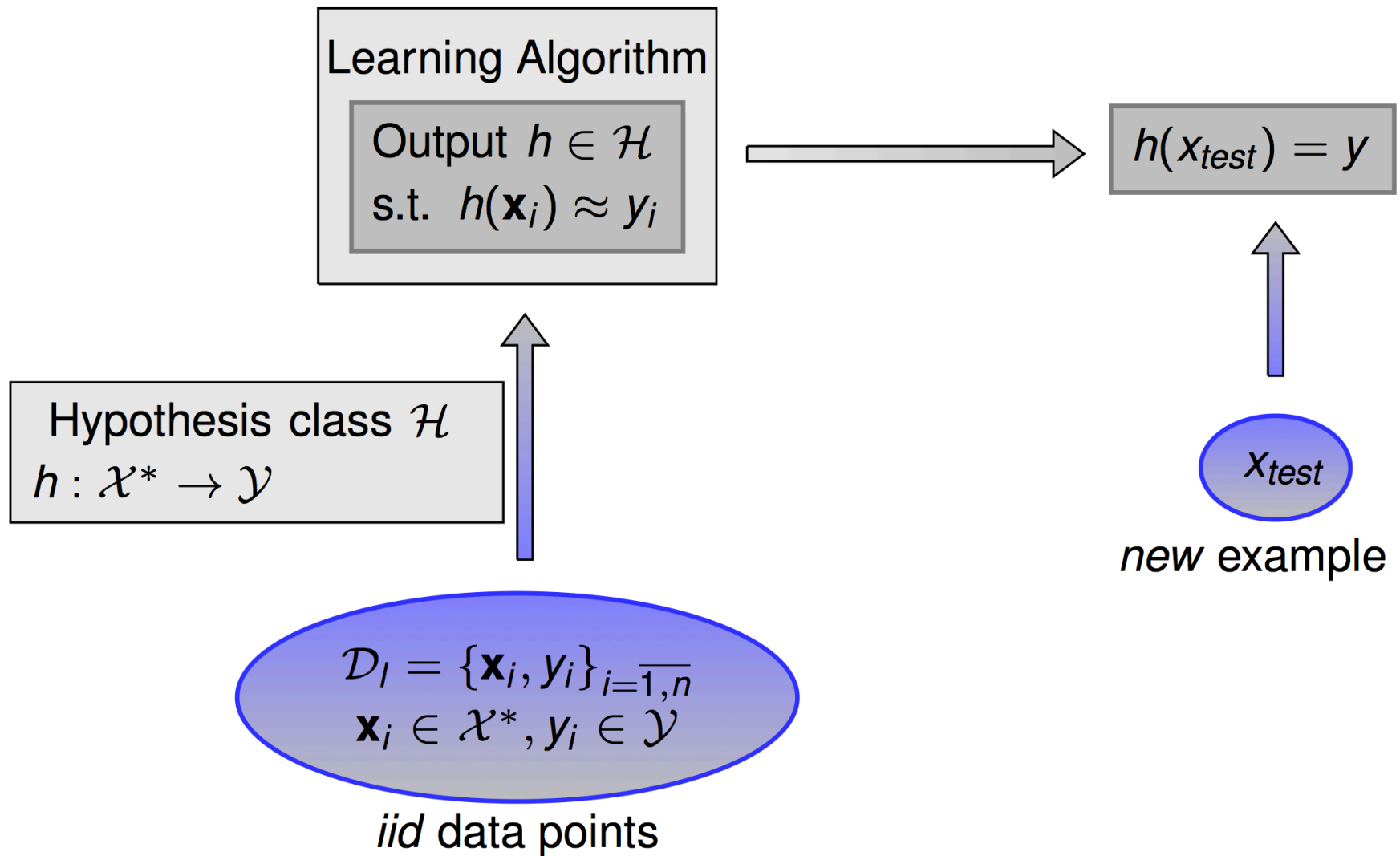
Data Representation

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no

Data Representation

- N = number of training examples
- x 's = “input” variable / features
- y 's = “output” variable / “target” variable
- (x, y) – one training example
- $(x^{(i)}, y^{(i)})$ – the i^{th} training example

Training and Classification



Summary of Basic Probability Formulas

- Product rule: probability of a conjunction of two events A and B

$$P(A \wedge B) = P(A | B)P(B) = P(B | A)P(A)$$

- Sum rule: probability of a disjunction of two events A and B

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

- Bayes theorem: the posterior probability of A given B

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

- Theorem of total probability: if events A_1, \dots, A_n are mutually exclusive with $\sum_{i=1}^n P(A_i) = 1$

$$P(B) = \sum_{i=1}^n P(B | A_i)P(A_i)$$

Bayes Classifiers for Categorical Data

Task: Classify a new instance x based on a tuple of attribute values $x = \langle x_1, x_2, \dots, x_n \rangle$ into one of the classes $c_j \in C$

$$\begin{aligned} c_{MAP} &= \operatorname{argmax}_{c_j \in C} P(c_j \mid x_1, x_2, \dots, x_n) \\ &= \operatorname{argmax}_{c_j \in C} \frac{P(x_1, x_2, \dots, x_n \mid c_j) P(c_j)}{P(x_1, x_2, \dots, x_n)} \\ &= \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_n \mid c_j) P(c_j) \end{aligned}$$

Example	Color	Shape	Class
1	red	circle	positive
2	red	circle	positive
3	red	square	negative
4	blue	circle	negative

← attributes

← values

Joint Distribution

- The joint probability distribution for a set of random variables, X_1, \dots, X_n gives the probability of every combination of values: $P(X_1, \dots, X_n)$

positive			negative		
	circle	square		circle	square
red	0.20	0.02	red	0.05	0.30
blue	0.02	0.01	blue	0.20	0.20

- The probability of all possible conjunctions can be calculated by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = ?$$

$$P(\text{red}) = ?$$

Joint Distribution

- The joint probability distribution for a set of random variables, X_1, \dots, X_n gives the probability of every combination of values: $P(X_1, \dots, X_n)$

positive			negative		
	circle	square		circle	square
red	0.20	0.02	red	0.05	0.30
blue	0.02	0.01	blue	0.20	0.20

- The probability of all possible conjunctions can be calculated by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = 0.20 + 0.05 = 0.25$$

$$P(\text{red}) = 0.20 + 0.02 + 0.05 + 0.3 = 0.57$$

Joint Distribution

- The joint probability distribution for a set of random variables, X_1, \dots, X_n gives the probability of every combination of values: $P(X_1, \dots, X_n)$

	positive	negative
	circle	square
red	0.20	0.02
blue	0.02	0.01

	circle	square
red	0.05	0.30
blue	0.20	0.20

- The probability of all possible conjunctions can be calculated by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = 0.20 + 0.05 = 0.25$$

$$P(\text{red}) = 0.20 + 0.02 + 0.05 + 0.3 = 0.57$$

- Therefore, all conditional probabilities can also be calculated.

Joint Distribution

- The joint probability distribution for a set of random variables, X_1, \dots, X_n gives the probability of every combination of values: $P(X_1, \dots, X_n)$

	positive			negative	
	circle	square		circle	square
red	0.20	0.02	red	0.05	0.30
blue	0.02	0.01	blue	0.20	0.20

- The probability of all possible conjunctions can be calculated by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = 0.20 + 0.05 = 0.25$$

$$P(\text{red}) = 0.20 + 0.02 + 0.05 + 0.3 = 0.57$$

- Therefore, all conditional probabilities can also be calculated.

$$P(\text{positive} \mid \text{red} \wedge \text{circle}) = ?$$

Joint Distribution

- The joint probability distribution for a set of random variables, X_1, \dots, X_n gives the probability of every combination of values: $P(X_1, \dots, X_n)$

	positive		negative	
	circle	square	circle	square
red	0.20	0.02	0.05	0.30
blue	0.02	0.01	0.20	0.20

- The probability of all possible conjunctions can be calculated by summing the appropriate subset of values from the joint distribution.

$$P(\text{red} \wedge \text{circle}) = 0.20 + 0.05 = 0.25$$

$$P(\text{red}) = 0.20 + 0.02 + 0.05 + 0.3 = 0.57$$

- Therefore, all conditional probabilities can also be calculated.

$$P(\text{positive} | \text{red} \wedge \text{circle}) = \frac{P(\text{positive} \wedge \text{red} \wedge \text{circle})}{P(\text{red} \wedge \text{circle})} = \frac{0.20}{0.25} = 0.80$$

Bayes Classifiers

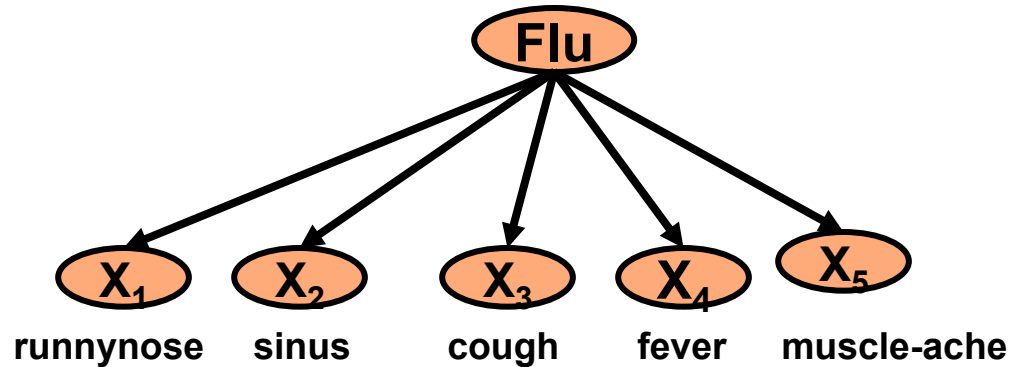
$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$

Bayes Classifiers

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_n | c_j)$
 - $O(|X|^n | C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.
 - Need to make some sort of independence assumptions about the features to make learning tractable.

The Naïve Bayes Classifier

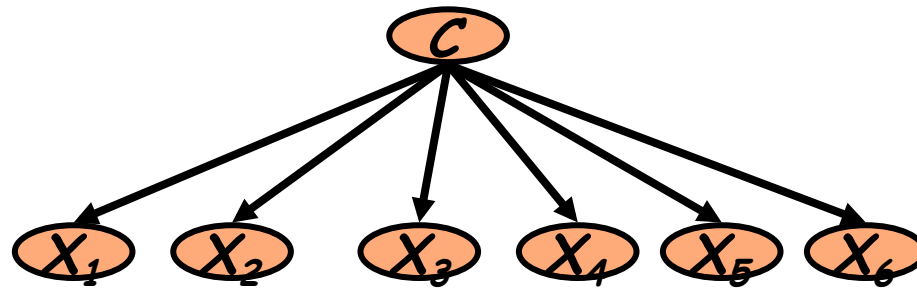


- **Conditional Independence Assumption:** attributes are independent of each other given the class:

$$P(X_1, \dots, X_5 | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_5 | C)$$

- Multi-valued variables: multivariate model
- Binary variables: multivariate Bernoulli model

Learning the Model

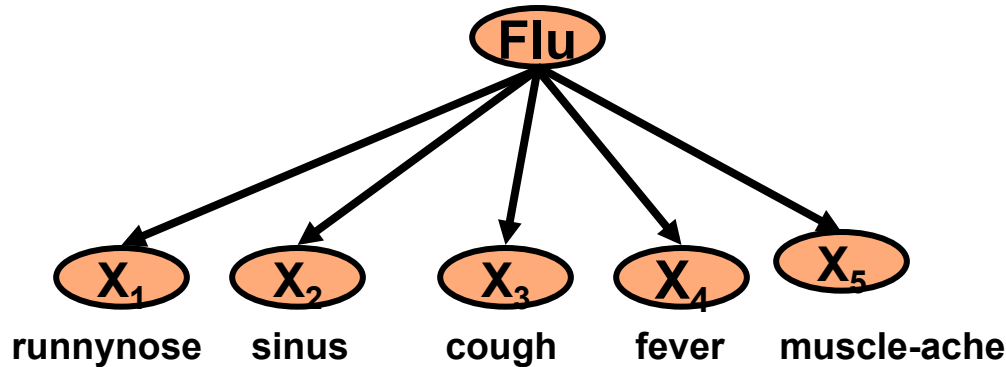


- First attempt: maximum likelihood estimates
 - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{N(C = c_j)}{N}$$

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j)}{N(C = c_j)}$$

Problem with Max Likelihood



$$P(X_1, \dots, X_5 \mid C) = P(X_1 \mid C) \cdot P(X_2 \mid C) \cdot \dots \cdot P(X_5 \mid C)$$

- What if we have seen no training cases where patient had no flu and muscle aches?

$$\hat{P}(X_5 = t \mid C = nf) = \frac{N(X_5 = t, C = nf)}{N(C = nf)} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$\ell = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

Smoothing to Improve Generalization on Test Data

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k}$$

of values of X_i

Underflow Prevention

- Multiplying lots of probabilities, which are between 0 and 1 by definition, can result in floating-point underflow.
- Since $\log(xy) = \log(x) + \log(y)$, it is better to perform all computations by summing logs of probabilities rather than multiplying probabilities.
- Class with highest final un-normalized log probability score is still the most probable.

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j)$$

Probability Estimation Example

Ex	Size	Color	Shape	Class
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

Probability	positive	negative
$P(Y)$		
$P(\text{small} \mid Y)$		
$P(\text{medium} \mid Y)$		
$P(\text{large} \mid Y)$		
$P(\text{red} \mid Y)$		
$P(\text{blue} \mid Y)$		
$P(\text{green} \mid Y)$		
$P(\text{square} \mid Y)$		
$P(\text{triangle} \mid Y)$		
$P(\text{circle} \mid Y)$		

Probability Estimation Example

Ex	Size	Color	Shape	Class
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

Probability	positive	negative
$P(Y)$	0.5	0.5
$P(\text{small} \mid Y)$	0.5	0.5
$P(\text{medium} \mid Y)$	0.0	0.0
$P(\text{large} \mid Y)$	0.5	0.5
$P(\text{red} \mid Y)$	1.0	0.5
$P(\text{blue} \mid Y)$	0.0	0.5
$P(\text{green} \mid Y)$	0.0	0.0
$P(\text{square} \mid Y)$	0.0	0.0
$P(\text{triangle} \mid Y)$	0.0	0.5
$P(\text{circle} \mid Y)$	1.0	0.5

Naïve Bayes Example

Probability	positive	negative
P(Y)	0.5	0.5
P(small Y)	0.4	0.4
P(medium Y)	0.1	0.2
P(large Y)	0.5	0.4
P(red Y)	0.9	0.3
P(blue Y)	0.05	0.3
P(green Y)	0.05	0.4
P(square Y)	0.05	0.4
P(triangle Y)	0.05	0.3
P(circle Y)	0.9	0.3

Test Instance:
<medium ,red, circle>

$$c_{MAP} = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

Naïve Bayes Example

Probability	positive	negative
$P(Y)$	0.5	0.5
$P(\text{medium} \mid Y)$	0.1	0.2
$P(\text{red} \mid Y)$	0.9	0.3
$P(\text{circle} \mid Y)$	0.9	0.3

Test Instance:
<medium ,red, circle>

$P(\text{positive} \mid X) = ?$

$P(\text{negative} \mid X) = ?$

$$c_{MAP} = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

Naïve Bayes Example

Probability	positive	negative
P(Y)	0.5	0.5
P(medium Y)	0.1	0.2
P(red Y)	0.9	0.3
P(circle Y)	0.9	0.3

$$c_{MAP} = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

Test Instance:
<medium ,red, circle>

$$\begin{aligned}
 P(\text{positive} | X) &= P(\text{positive}) * P(\text{medium} | \text{positive}) * P(\text{red} | \text{positive}) * P(\text{circle} | \text{positive}) / P(X) \\
 &\quad 0.5 \quad * \quad 0.1 \quad * \quad 0.9 \quad * \quad 0.9 \\
 &= 0.0405 / P(X) = 0.0405 / 0.0495 = 0.8181
 \end{aligned}$$

$$\begin{aligned}
 P(\text{negative} | X) &= P(\text{negative}) * P(\text{medium} | \text{negative}) * P(\text{red} | \text{negative}) * P(\text{circle} | \text{negative}) / P(X) \\
 &\quad 0.5 \quad * \quad 0.2 \quad * \quad 0.3 \quad * \quad 0.3 \\
 &= 0.009 / P(X) = 0.009 / 0.0495 = 0.1818
 \end{aligned}$$

$$P(\text{positive} | X) + P(\text{negative} | X) = 0.0405 / P(X) + 0.009 / P(X) = 1$$

$$P(X) = (0.0405 + 0.009) = 0.0495$$

Naïve Bayes for Text Classification

Two models:

- Multivariate Bernoulli Model
- Multinomial Model

Model 1: Multivariate Bernoulli

- One feature X_w for each word in dictionary
- $X_w = \text{true (1)}$ in document d if w appears in d
- Naive Bayes assumption:
 - Given the document's topic, appearance of one word in the document tells us nothing about chances that another word appears
- Parameter estimation

$$\hat{P}(X_w = 1 | c_j) = ?$$

Model 1: Multivariate Bernoulli

- One feature X_w for each word in dictionary
- $X_w = \text{true (1)}$ in document d if w appears in d
- Naive Bayes assumption:
 - Given the document's topic, appearance of one word in the document tells us nothing about chances that another word appears
- Parameter estimation

$$\hat{P}(X_w = 1 | c_j) = \text{fraction of documents of topic } c_j \text{ in which word } w \text{ appears}$$

Multinomial Naïve Bayes

- **Class conditional unigram language**
 - Attributes are text positions, values are words.
 - One feature X_i for each word position in document
 - feature's values are all words in dictionary
 - Value of X_i is the word in position i
 - **Naïve Bayes assumption:**
 - Given the document's topic, word in one position in the document tells us nothing about words in other positions

$$\begin{aligned}c_{NB} &= \operatorname{argmax}_{c_j \in C} P(c_j) \prod_i P(x_i | c_j) \\ &= \operatorname{argmax}_{c_j \in C} P(c_j) P(x_1 = \text{"our"} | c_j) \cdots P(x_n = \text{"text"} | c_j)\end{aligned}$$

- Too many possibilities!

Multinomial Naive Bayes Classifiers

- Second assumption:
 - Classification is *independent* of the positions of the words (word appearance does not depend on position)

$$P(X_i = w | c) = P(X_j = w | c)$$

for all positions i, j , word w , and class c

- Use same parameters for each position
- Result is bag of words model (over tokens)

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_i P(w_i | c_j)$$

Multinomial Naïve Bayes for Text

- Modeled as generating a bag of words for a document in a given category by repeatedly sampling with replacement from a vocabulary $V = \{w_1, w_2, \dots, w_m\}$ based on the probabilities $P(w_j | c_i)$.
- Smooth probability estimates with Laplace m -estimates assuming a uniform distribution over all words ($p = 1/|V|$) and $m = |V|$

Naïve Bayes Classification

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_i P(x_i | c_j)$$

Parameter Estimation

- Multivariate Bernoulli model:

$$\hat{P}(X_w = 1 | c_j) = \begin{array}{l} \text{fraction of documents of topic } c_j \\ \text{in which word } w \text{ appears} \end{array}$$

- Multinomial model:

$$\hat{P}(X_i = w | c_j) = \begin{array}{l} \text{fraction of times in which} \\ \text{word } w \text{ appears} \\ \text{across all documents of topic } c_j \end{array}$$

- Can create a mega-document for topic j by concatenating all documents in this topic
- Use frequency of w in mega-document

Classification

- Multinomial vs Multivariate Bernoulli?
- Multinomial model is almost always more effective in text applications!

Naïve Bayes - Spam Assassin

- Naïve Bayes has found a home in spam filtering
 - Paul Graham's *A Plan for Spam*
 - A mutant with more mutant offspring...
 - Widely used in spam filters
 - Classic Naive Bayes superior when appropriately used
 - According to David D. Lewis
 - But also many other things: black hole lists, etc.
- Many email topic filters also use NB classifiers

Naive Bayes is Not So Naive

- Naïve Bayes: First and Second place in KDD-CUP 97 competition, among 16 (then) state of the art algorithms

Goal: Financial services industry direct mail response prediction model: Predict if the recipient of mail will actually respond to the advertisement – 750,000 records.

- Robust to Irrelevant Features

Irrelevant Features cancel each other without affecting results

- Very good in domains with many equally important features
- A good baseline for text classification!
- Very Fast: Learning with one pass of counting over the data; testing linear in the number of attributes, and document collection size
- Low Storage requirements