

Is this spam?

----- Forwarded message -----
From: E-Lotto <infoalert@wp.pl>
Date: Fri, Sep 16, 2016 at 3:24 PM
Subject: Ref#.EL16BDXAUL16-03
To:

Hi,

E-Lotto congratulates you as the winner of \$2,500,000.00. Email Rep (Aaron Martins) at
aaronmarts@excite.com
with Ref#.EL16BDXAUL20

UMass Amherst

CS 585

Fall 2016

9/20 + 9/22

Classification Lecture
NB + LogReg + Eval

Slides from J&M 3rd ed. website

Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods



James Madison



Alexander Hamilton

Male or female author?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochin-China; the central area with its imperial capital at Hue was the protectorate of Annam...
2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, volume 23, number 3, pp. 321-346

#1: Saigon. Grey, Anthony

#2: Jerusalem the Golden. Drabble, Margaret
From shlomo argamon

Positive or negative movie review?



- unbelievably disappointing



- Full of zany characters and richly applied satire, and some great plot twists



- this is the greatest screwball comedy ever filmed



- It was pathetic. The worst part about it was the boxing scenes.

What is the subject of this article?

MEDLINE Article



MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...



Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- ...

Text Classification: definition

Prediction for a document

- *Input:*
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
- *Output:* a predicted class $c \in C$

Classification Methods: Hand-coded rules

- Rules based on combinations of words or other features
 - spam: black-list-address OR ("dollars" AND "have been selected")
- Accuracy can be high
 - If rules carefully refined by expert
- But building and maintaining these rules is expensive

~~Class~~
Problems

- Less Explicit Lang.
- Language Evolves
- False Positive
- Multiple word SENSES
- Lexical Coverage

Classification Methods: Supervised Machine Learning

① Train

- **Input:**
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$
 - A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$
- **Output:**
 - a learned classifier $f: d \rightarrow c$

11

② Predict on new data

Models for Sup. learning

⇒ Naive Bayes

- k -Nearest Neighbors

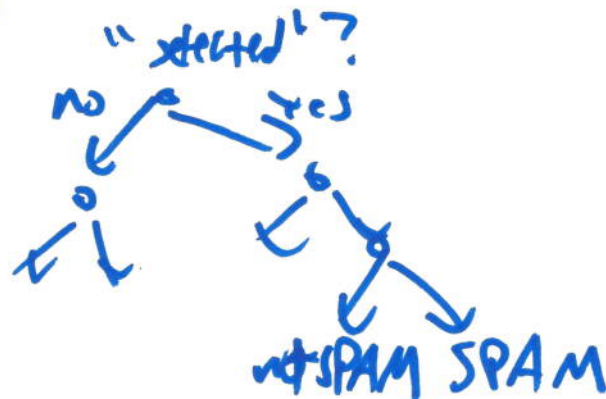
$$f(d) = c$$

⇒ Neural Networks

- Decision Trees

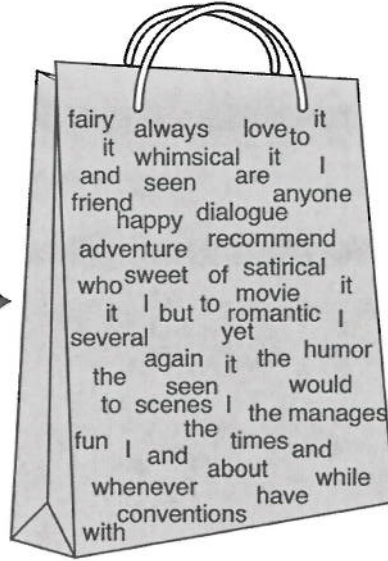
⇒ Support V. M.

⇒ Logistic Regression



The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



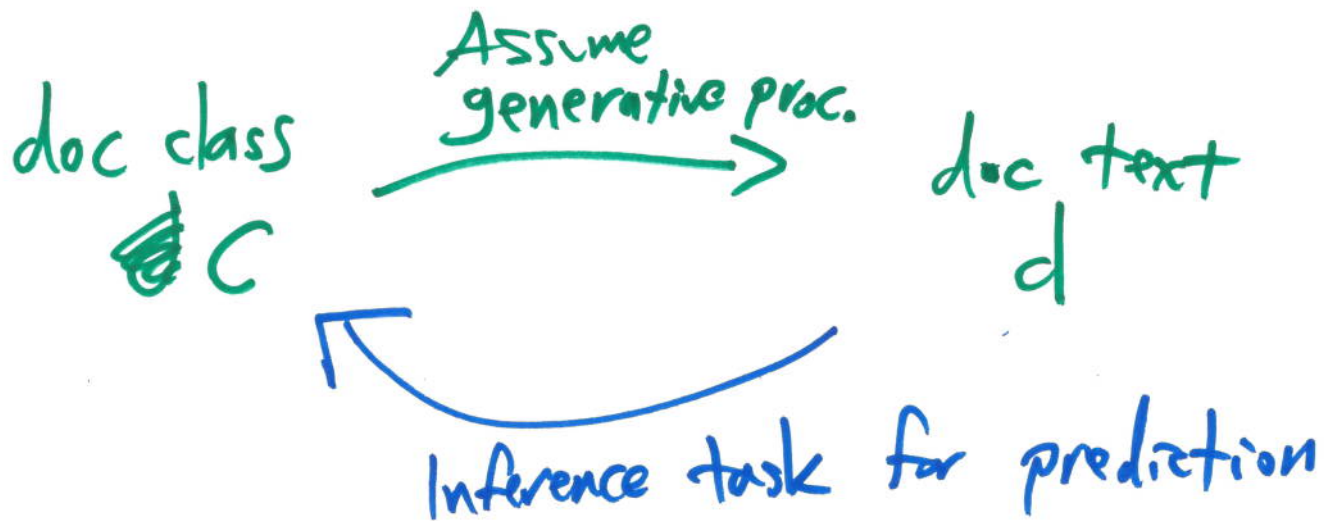
it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

15

Classifier with BOW

$$f(\dots) = c$$

Bayes Rule for Doc Classif.



$$P(c|d) = \frac{P(d|c) P(c)}{P(d)}$$

← Likelihood

→ Prior

↓ Ignore ???

$$= \sum_{c \in C} P(d|c) P(c)$$

$$\max_{x \in \{-1, 0, 1, 2\}} x^2 = 4$$

$$\operatorname{argmax}_{x \in \{-1, 0, 1, 2\}} x^2 = 2$$

↑↑

"argmax" notation

NB Classifier

$$C_{MAP} = \arg \max_{c \in C} P(c|d)$$

"MAP"
max.
posterior

$$= \arg \max_{c \in C} \frac{P(d|c) P(c)}{P(d)}$$

BR

← Constant

$$= \arg \max_{c \in C} P(d|c) P(c)$$

$$= \arg \max_{c \in C} \underbrace{P(x_1, x_2, \dots, x_N | c)}_{\text{Doc Features}} P(c)$$

Learning?

⇓

$$\approx O(|X|^N C)$$

Too many Parameters

⇓

Just count!

Easy!

Multinomial Naïve Bayes Independence Assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$



Can estimate

Multinomial Naïve Bayes Classifier

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

⇓ Conditional independence assumption

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{x \in X} P(x | c)$$

⇓ Unigram LM for class c

Naïve Bayes as a Language Model

- Which class assigns the higher probability to s?

Model pos	
0.1	I
0.1	love
0.01	this
0.05	fun
0.1	film

Model neg	
0.2	I
0.001	love
0.01	this
0.005	fun
0.1	film

$P(I \text{ love this film} | c)$

$$P(\text{text} | c = \text{pos}) = P(I | c = p) P(\text{love} | c = p) \dots$$

$$0.1 \times 0.1 \times 0.01 \times 0.05 \times 0.1$$

I love this fun film

-vs-

$$P(\text{text} | c = \text{neg}) = 0.2 \times 0.001 \times 0.01 \times 0.005 \times 0.1$$

Learning the Multinomial Naïve Bayes Model

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

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

↓
Total # tokens
in c_j

Parameter estimation

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

fraction of times word w_i appears
among all words in documents of topic c_j

- Create mega-document for topic j by concatenating all docs in this topic
 - Use frequency of w in mega-document

Problem with Maximum Likelihood

- What if we have seen no training documents with the word *fantastic* and classified in the topic **positive** (*thumbs-up*)?

$$\hat{P}(\text{"fantastic"} \mid \text{positive}) = \frac{\text{count}(\text{"fantastic"}, \text{positive})}{\sum_{w \in V} \text{count}(w, \text{positive})} = 0$$

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

$$c_{MAP} = \operatorname{argmax}_c \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

Laplace (add-1) smoothing for Naïve Bayes

$$\hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)}$$
$$= \frac{\text{count}(w_i, c) + 1}{\left(\sum_{w \in V} \text{count}(w, c) \right) + |V|}$$

$$\frac{\text{count}(w_i, c) + \alpha}{\left(\sum_w \dots \right) + \alpha |V|}$$

Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms

- For each c_j in C do

$docs_j \leftarrow$ all docs with class = c_j

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k | c_j)$ terms

- $Text_j \leftarrow$ single doc containing all $docs_j$

- For each word w_k in *Vocabulary*

$n_k \leftarrow$ # of occurrences of w_k in $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

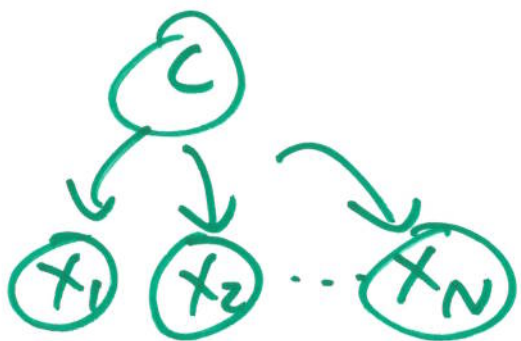
size? $\Rightarrow |V| \cdot |C|$

100k

Summary: Naive Bayes is Not So Naive

- Very Fast, low storage requirements
- Robust to Irrelevant Features
 - Irrelevant Features cancel each other without affecting results
- Very good in domains with many equally important features
 - Decision Trees suffer from *fragmentation* in such cases – especially if little data
- Optimal if the independence assumptions hold: If assumed independence is correct, then it is the Bayes Optimal Classifier for problem
- A good dependable baseline for text classification
 - **But we will see other classifiers that give better accuracy**

Multinom. Naive Bayes



To classify

$$\text{argmax}_c P(c, d)$$

$$= \dots P(c) P(d|c)$$

$$= \dots P(c) \prod_{i=1}^N P(x_i|c)$$

Cond. Indep. Assump.

NB \Leftrightarrow Unigram LM only if features = words

Other Features

- Is the word an emotion?

Happy emo?

Sad emo?

IS_HAPPY_EMU

Lists

$$P(:) | c)$$

$$P(:] | c)$$

$$P(:-)) | c) = ? ?$$

Might generalize better

$$P(\text{IS-HAPPY-EMU} | c)$$

Naïve Bayes in Spam Filtering

- SpamAssassin Features:
 - Mentions Generic Viagra
 - Online Pharmacy
 - Mentions millions of (dollar) ((dollar) NN,NNN,NNN.NN)
 - Phrase: impress ... girl
 - From: starts with many numbers
 - Subject is all capitals
 - HTML has a low ratio of text to image area
 - One hundred percent guaranteed
 - Claims you can be removed from the list
 - 'Prestigious Non-Accredited Universities'
 - http://spamassassin.apache.org/tests_3_3_x.html

Features: Domain Knowledge

Feature Engineering

NBS as a Log-Linear Model

$$\text{Pred Time } P(c|d) = \frac{1}{P(d)} P(c) \prod_{i=1}^N P(x_i|c)$$

$$\log P(c|d) = \underbrace{-\log P(d)}_{\text{Constant w.r.t. } c} + \underbrace{\log P(c)}_{w_{c,0}} + \sum_{i=1}^N \underbrace{\log P(x_i|c)}_{w_{x_i,c}}$$

$$\text{Log-Lin Model } P(c|d) = \frac{e^{\sum_{j=1}^n w_{j,c} f_j}}{\sum_{c'} e^{\sum_j w_{j,c'} f_j}}$$

Multinom. Logistic Regression

$$\log P(\text{dog}|\text{pos}) \equiv w_{\text{dog}, \text{pos}}$$

For MNB equivalence

$$f_j = \text{num. times word } j \text{ occurs}$$

NB vs Log. Reg

NB: Generative Model $P(c, d) = P(c) P(d|c)$

Train: maximize likelihood $\uparrow \Rightarrow$ Counting

Log. Reg: Discriminative Model $P(c|d)$ [Has no opinion about $P(d|c)$]

Train: maximize likelihood $\rightarrow \uparrow$

Algo: iterative gradient descent
 \Rightarrow Returns \vec{w}

Why? No CI assumptions in LR

Evaluation

Data Splits : All annotated data



↓↓
Learn Model



↑ Predict here
Check an
accuracy metric

Tuning hyperparameters? \Rightarrow Dev
- Pseudocount α

Metrics

		doc	doc				
<u>Gold:</u>	⊕	⊕	⊕	⊖	⊖	⊖	⊖
<u>Pred:</u>	⊕	⊕	⊖	⊖	⊖	⊕	⊕
			↑			↑	↑
			error			error	error
			False Neg			False Pos	

Acc Rate: % correct = $\frac{4}{7}$

Confusion Matrix

		Gold	
		⊕	⊖
Pred	⊕	2 TP	2 FP
	⊖	1 FN	2 TN

Prec & Recall

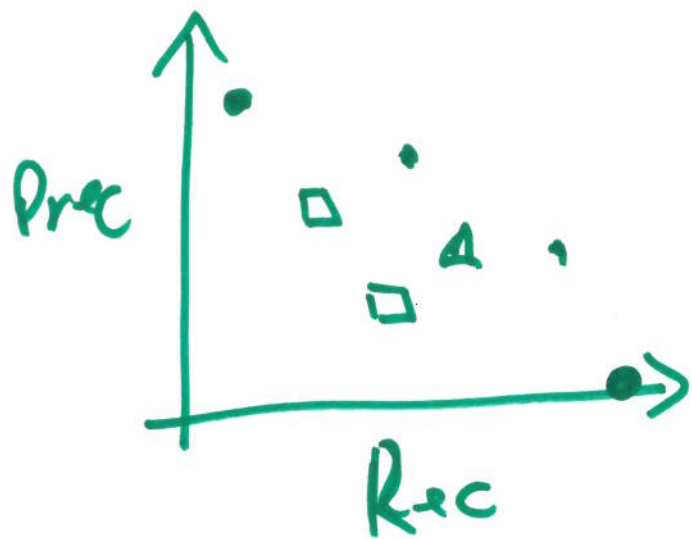
Prec = % correct of POS preds

$$= \frac{TP}{TP+FP} = \frac{2}{4} = \frac{1}{2}$$

Recall = % of gold POS we found

$$= \frac{TP}{TP+FN} = \frac{2}{2+1} = \frac{2}{3}$$

Prec & Recall Trade off



Single Number from P, R

$$F\text{-Score} = \frac{2 PR}{P+R}$$

"Harmonic Mean"

Development Loop - Feature Eng

