

TEXT CLASSIFICATION

Feng Gao, Feb. 26th, 2015

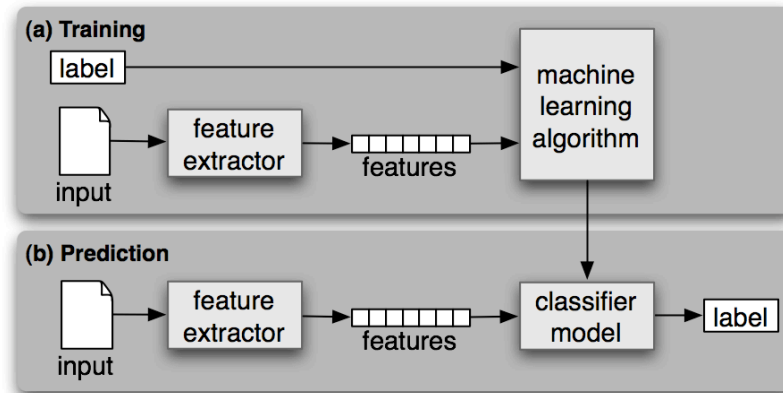
What is classification?



- ▣ Classification or categorization is the task of assigning objects from a universe to two or more classes or categories.
- ▣ Classification is the task of choosing the correct class label for a given input.

What is supervised Classification?

- A classifier is called supervised if it is built based on training corpora containing the correct label for each input.



What is text classification?

▣ The classifier:

- Input: a document x
- Output: a predicted class y from some fixed set of labels y_1, \dots, y_k

▣ The learner:

- Input: a set of m hand-labeled documents $(x_1, y_1), \dots, (x_m, y_m)$
- Output: a learned classifier $f: X \rightarrow y$

Application

- ▣ Personal email sorting
- ▣ Automatic detection of *spam* pages
- ▣ Automatic detection of sexually explicit content
- ▣ Automatic classification of a review as positive or negative
- ▣ Topic-specific or *vertical* search



Representation of text



Text Representation

$$f(\text{document}) = y$$

ARGENTINE 1986/87 GRAIN/OILSEED REGISTRATIONS
BUENOS AIRES, Feb 26
Argentine grain board figures show crop registrations of grains, oilseeds and their products to February 11, in thousands of tonnes, showing those for future shipments month, 1986/87 total and 1985/86 total to February 12, 1986, in brackets:

- Bread wheat prev 1,655.8, Feb 872.0, March 164.6, total 2,692.4 (4,161.0).
- Maize Mar 48.0, total 48.0 (nil).
- Sorghum nil (nil)
- Oilseed export registrations were:
- Sunflowerseed total 15.0 (7.9)
- Soybean May 20.0, total 20.0 (nil)

The board also detailed export registrations for subproducts, as follows....



?

What is the ~~best~~ **simplest useful** representation for the document x being classified?

Text Representation

- Document is represented as a vector of attribute values
- **Attributes:**
“Bag of words” method: Use a set of words as attributes

Text Representation

▣ Attribute values:

Method 1:

use 0 or 1 as attribute value

Method 2:

use the absolute or relative frequency of each word

Method 3:

use TF-IDF weight as the attribute value

Method 1

Training data sets:

▣ Method 1:

	word₁	word₂	...	word_m	Class
document₁	0	1	...	1	C1
document₂	1	0	...	1	C2
...
document_n	1	0		0	C2

Method 2

Training data sets:

- Method 2 with absolute frequency:

	word₁	word₂	...	word_m	Class
document₁	0	3	...	1	C1
document₂	2	0	...	3	C2
...
document_n	5	0		0	C2

Method 3

□ TF: term frequency

- Definition: $TF = t_{ij}$
- frequency of term i in document j
- Purpose: makes the frequent words *for the document* more important

□ IDF: inverted document frequency

- Definition: $IDF = \log(N/n_i)$
- n_i : number of documents containing term i
- N : total number of documents

□ TF-IDF value of a term i in document j

- Definition: $TF \times IDF = t_{ij} * \log(N/n_i)$

Text Processing

- ▣ Word (token) extraction
- ▣ Stop words removal
- ▣ Stemming
- ▣ Feature Selection



Text Processing

f(ARGENTINE 1986/87 GRAIN/OILSEED REGISTRATIONS
BUENOS AIRES, Feb 26
Argentine grain board figures show crop registrations of grains, oilseeds and their products to February 11, in thousands of tonnes, showing those for future shipments month, 1986/87 total and 1985/86 total to February 12, 1986, in brackets:
• Bread wheat prev 1,655.8, Feb 872.0, March 164.6, total 2,692.4 (4,161.0).
• Maize Mar 48.0, total 48.0 (nil).
• Sorghum nil (nil)
• Oilseed export registrations were:
• Sunflowerseed total 15.0 (7.9)
• Soybean May 20.0, total 20.0 (nil)
The board also detailed export registrations for subproducts, as follows....)=y

f((argentine, 1986, 1987, grain, oilseed,
registration, buenos, aires, feb, 26,
argentine, grain, board, figures, show, crop,
registrations, of, grains, oilseeds, and, their,
products, to, february, 11, in, ...)=y

Common refinements: **remove stopwords**, **stemming**, collapsing multiple occurrences of words into one....

Word (token) extraction

- ❑ Extract all the words in a document
- ❑ Convert them into lower cases

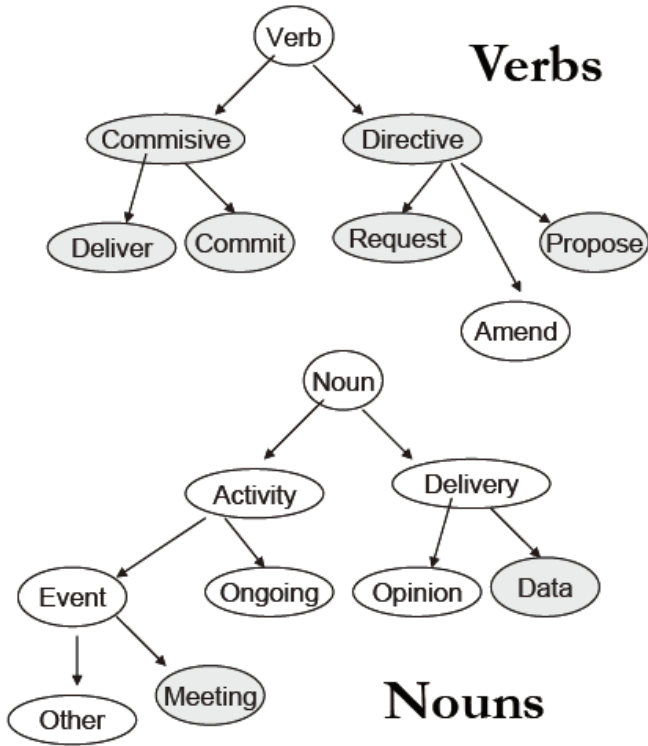


Classifying Email into Acts



- From EMNLP-04, Learning to Classify Email into Speech Acts, Cohen-Carvalho-Mitchell
- An Act is described as a verb-noun pair (e.g., propose meeting, request information) - Not all pairs make sense. One single email message may contain multiple acts.
- Try to describe commonly observed behaviors, rather than all possible speech acts in English. Also include non-linguistic usage of email (e.g. delivery of files)

Classifying Email into Acts



Classifying Email into Acts

Symbol	Pattern
[number]	any sequence of numbers
[hour]	[number]:[number]
[wvhh]	“why, where, who, what, or when”
[day]	the strings “Monday, Tuesday, ..., or Sunday”
[day]	the strings “Mon, Tue, Wed, ..., or Sun”
[pm]	the strings “P.M., PM, A.M. or AM”
[me]	the pronouns “me, her, him, us or them”
[person]	the pronouns “I, we, you, he, she or they”
[aafter]	the strings “after, before or during”
[filetype]	the strings “.doc, .pdf, .ppt, .txt, or .xls”

Table 1: Some PreProcessing Substitution Patterns

1-gram	3-gram
?	[person] need to
please	[wvhh] do [person]
[wvhh]	let [me] know
could	would [person]
do	do [person] think
can	are [person] meeting
of	could [person] please
[me]	do [person] need

5-gram
[wvhh] do [person] think ?
let [me] know [wvhh] [person]
a call [number]-[number]
give [me] a call [number]
please give give [me] a call
[person] would be able to
take a look at it
[person] think [person] need to

Word (token) extraction for Email

Request	Commit	Meeting
[wwhh] do [person] think do [person] need to and let [me] know call [number]-[number] would be able to [person] think [person] need let [me] know [wwhh] do [person] think ? [person] need to get ? [person] need to a copy of our do [person] have any [person] get a chance [me] know [wwhh] that would be great	is good for [me] is fine with [me] i will see [person] i think i can i will put the i will try to i will be there will look for [person] \$[number] per person am done with the at [hour] i will [day] is fine with each of us will i will bring copies i will do the	[day] at [hour] [pm] on [day] at [hour] [person] can meet at [person] meet at [hour] will be in the is good for [me] to meet at [hour] at [hour] in the [person] will see [person] meet at [hour] in [number] at [hour] [pm] to go over the [person] will be in let's plan to meet meet at [hour] [pm]
dData	Propose	Deliver
- forwarded message begins forwarded message begins here is in my public in my public directory [person] have placed the please take a look [day] [hour] [number] [number] [number] [day] [number] [hour] [date] [day] [number] [day] in our game directory in the etc directory the file name is is in our game fyi - forwarded message just put the file my public directory under	[person] would like to would like to meet please let [me] know to meet with [person] [person] meet at [hour] would [person] like to [person] can meet tomorrow an hour or so meet at [hour] in like to get together [hour] [pm] in the [after] [hour] or [after] [person] will be available think [person] can meet was hoping [person] could do [person] want to	forwarded message begins here [number] [number] [number] [number] is good for [me] if [person] have any if fine with me in my public directory [person] will try to is in my public will be able to just wanted to let [pm] in the lobby [person] will be able please take a look can meet in the [day] at [hour] is in the commons at

Stop words removal

- ▣ The most frequently used words in English
- ▣ Examples of stop words
 - ▣ the, of, and, to, a, ...
- ▣ Typically about 400 to 500 such words
- ▣ Additional domain-specific stop words
- ▣ Stop words are usually removed



Stemming

- ❑ find the root/stem of a word
- ❑ Reduce the number of words
- ❑ Improve effectiveness of text classification
- ❑ For example:
 - ❑ discussed
 - ❑ discusses
 - ❑ discussing
 - ❑ Discuss
 - ❑ Stem: discuss



Example Stemming Rules

▣ Remove ending

- ▣ If a word ends with *s*, preceded by a consonant other than an *s*, then delete the *s*.

▣ Transform words

- ▣ If a word ends with “ies” but not “eies” or “aies”, then “ies” is replaced with “y”.



Feature Selection



- ▣ Selecting the “bag of words” to represent documents
- ▣ Why do we need to select?
 - ▣ Learning program may not be able to handle all possible features
 - ▣ Good features can result in higher accuracy

Feature Selection Methods

▣ Class independent methods (Unsupervised)

- ▣ Document Frequency (DF)
- ▣ Term Strength (TS)

▣ Class-dependent methods (Supervised)

- ▣ Information Gain (IG)
- ▣ Mutual Information (MI)
- ▣ χ^2 statistic (CHI)

Document Frequency (DF)

Document frequency of a word

- $DF(w) = \text{number of documents containing } w$

Advantages

- Can remove rare words (hence noise)
- Easy to compute

Disadvantages

- Class independent
- Some infrequent terms can be good discriminators, which cannot be selected by this method.

Information Gain

- ▣ A measure of importance of the feature
- ▣ The number of “bits of information” gained by knowing the word is present or absent

$$\text{Gain}(\omega) = - \sum_{i=1}^k P(C_i) \log P(C_i) + P(\omega) \sum_{i=1}^k P(C_i|\omega) \log P(C_i|\omega) \\ + P(\bar{\omega}) \sum_{i=1}^k P(C_i|\bar{\omega}) \log P(C_i|\bar{\omega})$$

- ▣ Rank the words according to their information gain value
- ▣ Select the first m words with high gain values

Information Gain



▣ Advantage

- ▣ Consider the classes

▣ Disadvantage

- ▣ computationally expensive

▣ Remove rare words (appears 1 or 2 times)

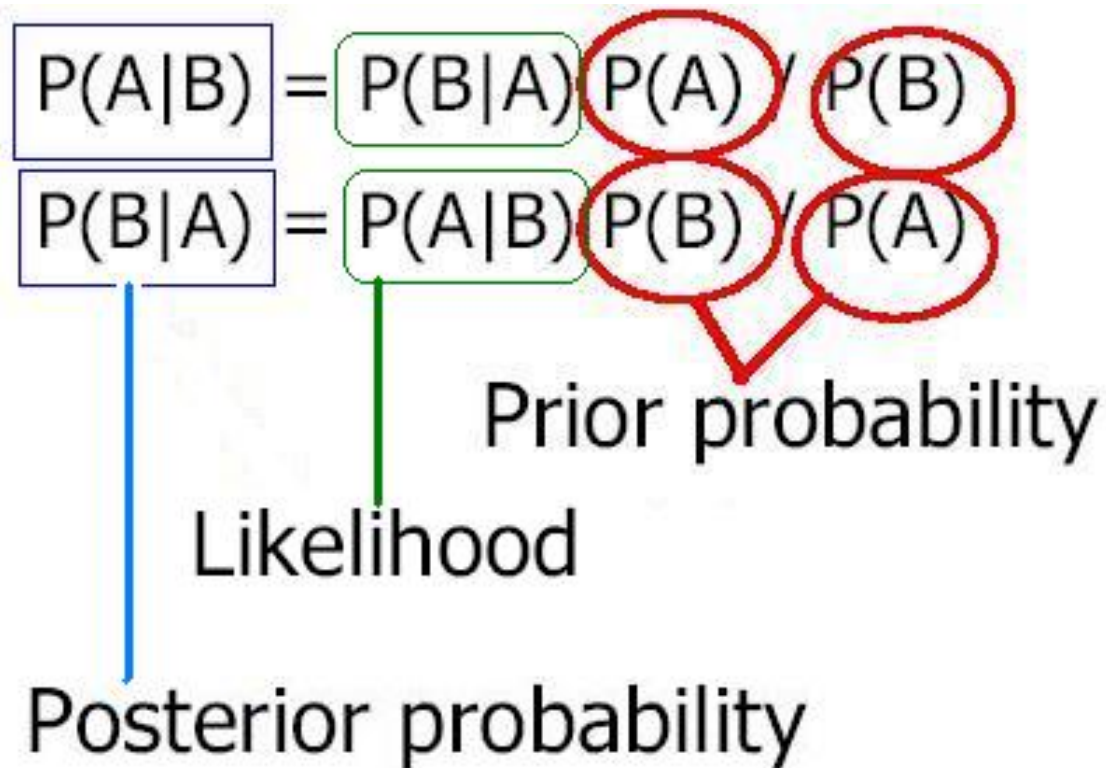
- ▣ reduce the amount of computation, and
- ▣ remove noisy words that have by-chance correlations with the classes.

What Do People Do In Practice?



- ▣ Infrequent term removal
 - ▣ infrequent across the whole collection (i.e. DF)
 - ▣ met in a single document
- ▣ Most frequent term removal (i.e. removing stop words)
- ▣ Stemming. (*often*)
- ▣ Use a class-dependent method (e.g., the information gain method) to select features.

Naive Bayes



Text Classification with Naive Bayes

- ▣ Represent document x as list of words w_1, w_2, \dots
- ▣ For each y , build a probabilistic model $\Pr(X | Y=y)$ of “documents” in class y
- ▣ To classify, find the y which was most likely to generate x —i.e., which gives x the best score according to $\Pr(x | y)$

$$f(x) = \operatorname{argmax}_y \Pr(x|y) * \Pr(y)$$

Text Classification with Naive Bayes

- How to estimate $\Pr(X | Y)$?
- Simplest useful process to generate a bag of words:
 - pick word 1 according to $\Pr(W | Y)$
 - repeat for word 2, 3, ...
 - each word is generated independently of the others (which is clearly not true) but means

$$\Pr(w_1, \dots, w_n | Y = y) = \prod_{i=1}^n \underbrace{\Pr(w_i | Y = y)}$$

How to estimate $\Pr(W|Y)$?

Text Classification with Naive Bayes

- How to estimate $\Pr(X | Y)$?

$$\Pr(w_1, \dots, w_n | Y = y) = \prod_{i=1}^n \Pr(w_i | Y = y)$$

Estimate $\Pr(w|y)$ by looking at the data...

$$\Pr(W = w | Y = y) = \frac{\text{count}(W = w \text{ and } Y = y)}{\text{count}(Y = y)}$$

- This gives score of zero if x contains a brand-new word w_{new}

Text Classification with Naive Bayes

- How to estimate $\Pr(X | Y)$?

$$\Pr(w_1, \dots, w_n | Y = y) = \prod_{i=1}^n \underbrace{\Pr(w_i | Y = y)}$$

... and also **imagine** m
examples with $\Pr(w|y)=p$

$$\Pr(W = w | Y = y) = \frac{\text{count}(W = w \text{ and } Y = y) + mp}{\text{count}(Y = y) + m}$$

- This $\Pr(W | Y)$ is a *multinomial distribution*
- This use of m and p is a *Dirichlet prior for the multinomial*

Text Classification with Naive Bayes

- ▣ Putting this together:
- ▣ for each document x_i with label y_i
 - for each word w_{ij} in x_i
 - $\text{count}[w_{ij}][y_i]++$
 - $\text{count}[y_i]++$
 - $\text{count}++$

Text Classification with Naive Bayes

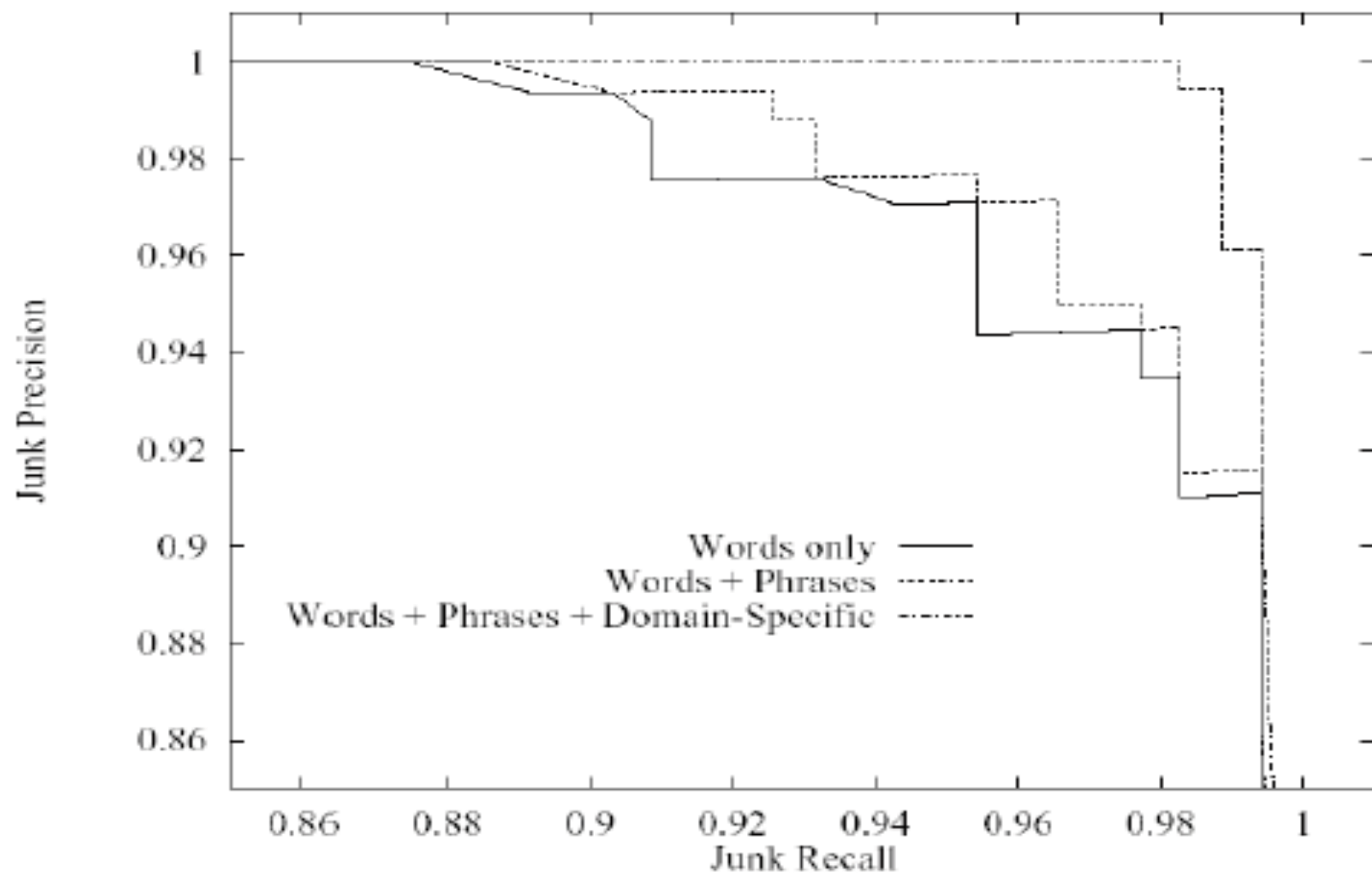
- to classify a new $x = w_1 \dots w_n$, pick y with top score:

$$score(y, w_1 \dots w_k) = \lg \frac{\text{count}[y]}{\text{count}} + \sum_{i=1}^n \lg \frac{\text{count}[w_i][y] + 0.5}{\text{count}[y] + 1}$$

key point: we only need counts for words that actually appear in x

Naïve Bayes for SPAM filtering

- ▣ Sahami et al, 1998
- ▣ Used bag of words, + special phrases (“FREE!”) and + special features (“from *.edu” , ...)



Naïve Bayes for SPAM filtering

	Classified Junk	Classified Legitimate	Total
Actually Junk	36 (92.0% precision)	9	45
Actually Legitimate	3	174 (95.0% precision)	177
Total	39	183	222

Naive Bayes Summary

□ Pros:

- Very fast and easy-to-implement
- Well-understood formally & experimentally

□ Cons:

- Seldom gives the very best performance
- “Probabilities” $\Pr(y | x)$ are not accurate

The Voted Perceptron

Training

Input: a labeled training set $\langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \rangle$
number of epochs T

Output: a list of weighted perceptrons $\langle (\mathbf{v}_1, c_1), \dots, (\mathbf{v}_k, c_k) \rangle$

- Initialize: $k := 0$, $\mathbf{v}_1 := \mathbf{0}$, $c_1 := 0$.
- Repeat T times:
 - For $i = 1, \dots, m$:
 - * Compute prediction: $\hat{y} := \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$
 - * If $\hat{y} = y$ then $c_k := c_k + 1$.
 - else $\mathbf{v}_{k+1} := \mathbf{v}_k + y_i \mathbf{x}_i$;
 - $c_{k+1} := 1$;
 - $k := k + 1$.

The Voted Perceptron

Prediction

Given: the list of weighted perceptrons: $\langle (\mathbf{v}_1, c_1), \dots, (\mathbf{v}_k, c_k) \rangle$
an unlabeled instance: \mathbf{x}

compute a predicted label \hat{y} as follows:

$$s = \sum_{i=1}^k c_i \text{sign}(\mathbf{v}_i \cdot \mathbf{x}); \quad \hat{y} = \text{sign}(s) .$$

Classifying Reviews as Favorable or Not



- ▣ Turney, ACL 2002
- ▣ Dataset: 410 reviews from Epinions
 - ▢ Autos, Banks, Movies, Travel Destinations
- ▣ Learning method:
 - ▢ Extract 2-word phrases containing an adverb or adjective (eg “unpredictable plot”)
 - ▢ Classify reviews based on average Semantic Orientation

Classifying Reviews as Favorable or Not

$$SO(\textit{phrase}) = \text{PMI}(\textit{phrase}, \text{“excellent”}) \\ - \text{PMI}(\textit{phrase}, \text{“poor”})$$

$$\text{PMI}(\textit{word}_1, \textit{word}_2) = \log_2 \left[\frac{p(\textit{word}_1 \ \& \ \textit{word}_2)}{p(\textit{word}_1) p(\textit{word}_2)} \right]$$

**Computed using
queries to web
search engine**

Classifying Reviews as Favorable or Not

Extracted Phrase	Part-of-Speech Tags	Semantic Orientation
online experience	JJ NN	2.253
low fees	JJ NNS	0.333
local branch	JJ NN	0.421
small part	JJ NN	0.053
online service	JJ NN	2.780
printable version	JJ NN	-0.705
direct deposit	JJ NN	1.288
well other	RB JJ	0.237
inconveniently located	RB VBN	-1.541
other bank	JJ NN	-0.850
true service	JJ NN	-0.732
Average Semantic Orientation		0.322

Classifying Reviews as Favorable or Not

Table 5. The accuracy of the classification and the correlation of the semantic orientation with the star rating.

Domain of Review	Accuracy	Correlation
Automobiles	84.00 %	0.4618
Honda Accord	83.78 %	0.2721
Volkswagen Jetta	84.21 %	0.6299
Banks	80.00 %	0.6167
Bank of America	78.33 %	0.6423
Washington Mutual	81.67 %	0.5896
Movies	65.83 %	0.3608
The Matrix	66.67 %	0.3811
Pearl Harbor	65.00 %	0.2907
Travel Destinations	70.53 %	0.4155
Cancun	64.41 %	0.4194
Puerto Vallarta	80.56 %	0.1447
All	74.39 %	0.5174

Classifying Reviews as Favorable or Not

Table 5. The accuracy of the classification and the correlation of the semantic orientation with the star rating.

Domain of Review	Accuracy	Correlation
Automobiles	84.00 %	0.4618
Honda Accord	83.78 %	0.2721
Volkswagen Jetta	84.21 %	0.6299
Banks	80.00 %	0.6167
Bank of America	78.33 %	0.6423
Washington Mutual	81.67 %	0.5896
Movies	65.83 %	0.3608
The Matrix	66.67 %	0.3811
Pearl Harbor	65.00 %	0.2907
Travel Destinations	70.53 %	0.4155
Cancun	64.41 %	0.4194
Puerto Vallarta	80.56 %	0.1447
All	74.39 %	0.5174

Classifying Reviews as Favorable or Not

- ▣ Pang et al, EMNLP 2002
- ▣ 700 movie reviews (ie all in same domain); Naïve Bayes, MaxEnt, and linear SVMs; accuracy with different representations x for a document
- ▣ Interestingly, the off-the-shelf methods work well... perhaps better than Turney's method.

Classifying Movie Reviews

	Features	# of features	frequency or presence?	NB	ME	SVM
(1)	unigrams	16165	freq.	78.7	N/A	72.8
(2)	unigrams	"	pres.	81.0	80.4	82.9
(3)	unigrams+bigrams	32330	pres.	80.6	80.8	82.7
(4)	bigrams	16165	pres.	77.3	77.4	77.1
(5)	unigrams+POS	16695	pres.	81.5	80.4	81.9
(6)	adjectives	2633	pres.	77.0	77.7	75.1
(7)	top 2633 unigrams	2633	pres.	80.3	81.0	81.4
(8)	unigrams+position	22430	pres.	81.0	80.1	81.6

Classifying Movie Reviews

- Assume the classifier is same form as Naive Bayes, which can be written:

$$\Pr(y | w_1, w_2, \dots, w_N) = \frac{1}{Z} \prod_i \lambda_i f(y, w_i)$$

- Set weights (λ 's) to maximize probability of the training data:

$$\prod_{(x_j, y_j) \in D} \Pr(y_j | x_j) + \underbrace{\Pr(\lambda | Q)}_{\text{prior on parameters}}$$

prior on parameters

MaxEnt classification

