

CS 6501 Natural Language Processing

Text Classification and Naive Bayes

Yangfeng Ji

September 10, 2019

Department of Computer Science
University of Virginia



ENGINEERING

Overview

1. Problem Definition
2. Bag-of-Words Representation
3. Naive Bayes Classifiers
4. Classification Evaluation

Problem Definition

Case I: Sentiment Analysis



The screenshot shows the Yelp interface with a red header containing the Yelp logo. Below the header, the section is titled "Recommended Reviews". On the right, there is a search bar labeled "Search reviews" and a magnifying glass icon. Below the search bar, there are tabs for "Yelp Sort", "Date", "Rating", and "Elites". On the far right, there is a language selector showing "English" and a count of "16". The main review is by "Jenn P." from "San Francisco, CA". It shows a profile picture of a woman, a red star icon, "1 friend", and "22 reviews". The review is dated "10/17/2013" and has a 5-star rating. The text of the review is: "Absolutely Outstanding! The Grounds at Grace Vineyards are stunning...there are SO many photo ops. I must give 5 stars for Steve the owner he is simply wonderful. He was so organized, flexible and prompt I never was stressed. The food was great and the vino was delicious! If your looking for a beautiful venue with many things included this is the place."

Recommended Reviews

Search reviews

Yelp Sort Date Rating Elites English 16

Jenn P.
San Francisco, CA
1 friend
22 reviews

★★★★★ 10/17/2013

Absolutely Outstanding! The Grounds at Grace Vineyards are stunning...there are SO many photo ops. I must give 5 stars for Steve the owner he is simply wonderful. He was so organized, flexible and prompt I never was stressed. The food was great and the vino was delicious! If your looking for a beautiful venue with many things included this is the place.

[Pang et al., 2002]

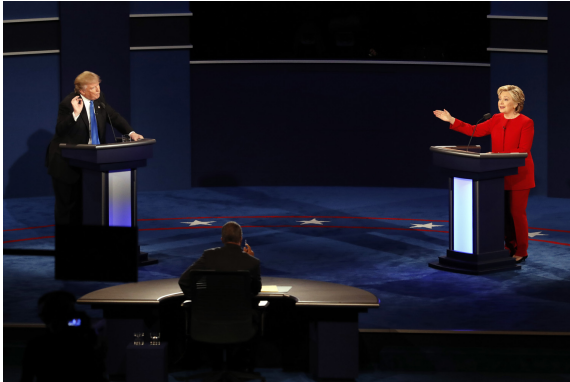
Case II: Topic Classification



Example topics

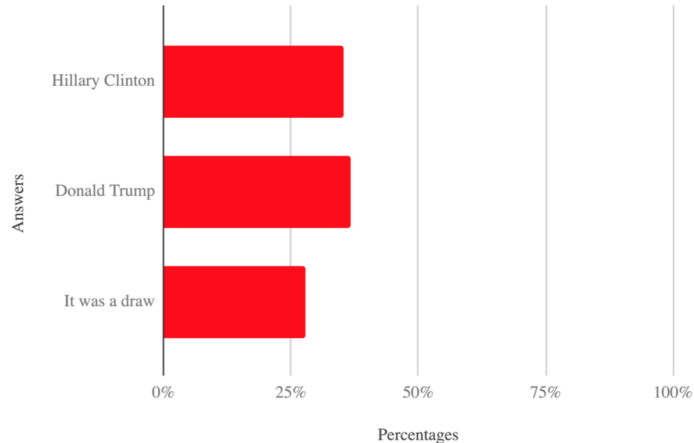
- ▶ Business
- ▶ Arts
- ▶ Technology
- ▶ Sports
- ▶ ...

Case III: Presidential Election Debates



Case III: Presidential Election Debates (II)

Republicans - Who do you think was the biggest winner of the debate?



Classification

- ▶ Input: a text x
- ▶ Output: $y \in \mathcal{Y}$, where \mathcal{Y} is the predefined category set (sample space)
 - ▶ Example: $\mathcal{Y} = \{\text{POSITIVE}, \text{NEGATIVE}\}$



Probabilistic Formulation

With the conditional probability $P(Y | X)$, the prediction on Y for a given text $X = x$ is

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} P(Y = y | X = x) \quad (1)$$

Or, for simplicity

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} p(y | x) \quad (2)$$

Key Questions

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} P(Y = y \mid X = x) \quad (3)$$

1. How to represent a text as x ?
 - ▶ Bag-of-words representation
2. How to estimate $P(y \mid x)$?
 - ▶ Naive Bayes classifier

Bag-of-Words Representation

Bag-of-Words Representation

Tokenization: convert a text into a collection of tokens:

*Super quick and really
friendly **staff**. I'd like starting
off my mornings at this **store**!!* \Rightarrow super quick and really
friendly staff . I 'd ... store
!!

NLTK function

```
nltk.tokenize.wordpunct_tokenize
```

Vocab

Collect tokens to build a vocab:

super quick and really friendly
staff . I 'd ... store ! !

⇒

$$\left\{ \begin{array}{c} \text{SUPER} \\ \dots \\ \text{QUICK} \\ \text{FOOD} \\ \text{FRIENDLY} \\ \text{EAT} \\ \dots \\ \text{STAFF} \end{array} \right\}$$

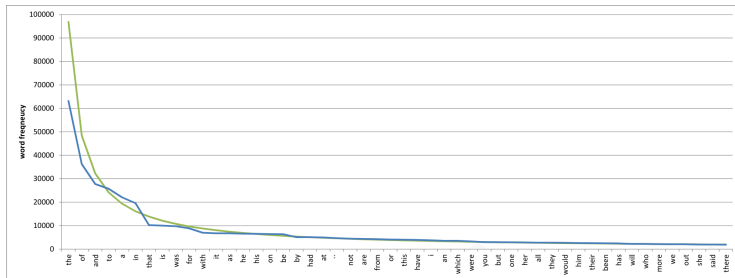
Should we collect all the tokens?

Preprocessing for Building Vocab

1. convert all characters to lowercase

$$\text{UVa}, \text{UVA} \rightarrow \text{uVa}$$

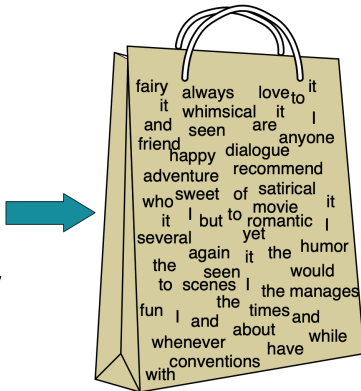
2. map low frequency words to a special token UNK



$$\text{Zipf's law: } f(w_t) \propto 1/r_t$$

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
...	...

Information Embedded in BoW Representation

- ▶ Lose:
 - ▶ word order
 - ▶ sentence boundary
 - ▶ paragraph boundary
 - ▶ ...
- ▶ Keep: words in texts

Naive Bayes Classifiers

Decision Rule

Given a document x , the classification can be conducted as

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} p(y \mid x; \theta). \quad (4)$$

where θ is the parameter of the distribution.

Bayes' Theorem

In Bayes' rule, the conditional probability $p(y | x)$ can be computed via the joint probability $p(x, y)$ as follow

$$\begin{aligned} p(y | x) &= \frac{p(x, y)}{p(x)} \\ &= \frac{p(x | y)p(y)}{p(x)} \end{aligned} \tag{5}$$

- ▶ $p(y, \theta)$: prior probability of $Y = y$
- ▶ $p(x | y; \theta)$:
 - ▶ conditional probability of $X = x$ given y
 - ▶ likelihood function of $Y = y$ given x

Simplification

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y \in \mathcal{Y}} p(y \mid x; \theta) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \underbrace{p(x \mid y)}_{\text{likelihood}} \underbrace{p(y)}_{\text{prior}}\end{aligned}\tag{6}$$

Benefit of equation 6: no need to compute $p(x)$ explicitly.

Naive Bayes Assumption

To model $p(x | y)$: we assume words within a text are independent with each other

- ▶ The distribution of $p(x | y)$ is similar to the one of modeling tossing a dice with V faces for n times
- ▶ This is a naive assumption

Conditional probability $p(x | y; \theta)$ can be written as a multinomial distribution

$$p(x | y; \theta) \propto \prod_{i=1}^V \theta_{i,y}^{x_i} \quad (7)$$

Parameters

- ▶ Prior probability (categorical distribution):

$$p(y) = \theta_y \quad (8)$$

- ▶ Likelihood (multinomial distribution):

$$p(x \mid y; \theta) \propto \prod_{i=1}^V \theta_{i,y}^{x_i} \quad (9)$$

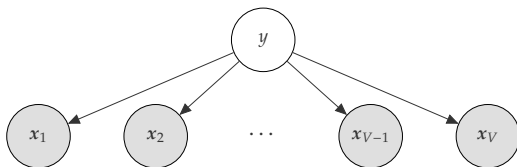
- ▶ Number of the parameters $K + KV = (K + 1)V$

Probabilistic Graphical Model

Joint probability:

$$p(\mathbf{x}, y; \boldsymbol{\theta}) \propto \theta_y \cdot \prod_{i=1}^V \theta_{i,y}^{x_i} \quad (10)$$

The conditional independence can be represented in the following graphical model:



How to estimate $\boldsymbol{\theta} = \{\theta_y, \theta_{i,y}\}, \forall i, y$

Example

Training examples

Label	Text	Words
NEG	just plain boring	3
NEG	entirely predictable and lacks energy	5
NEG	no surprises and very few laughs	6
Pos	very powerful	2
Pos	the most fun film of the summer	7

[Jurafsky and Martin, 2019]

Parameter Estimation (I)

For prior probability

$$\begin{aligned}\theta_y &= \frac{\text{Number of texts with label } y}{\text{Total number of texts}} \\ &= \frac{N_y}{N_{\text{total}}}\end{aligned}\tag{11}$$

Example

There are three documents: two positive and three negative

$$p(\text{POSITIVE}) = \frac{2}{5} \quad p(\text{NEGATIVE}) = \frac{3}{5}\tag{12}$$

Parameter Estimation (II)

For $\theta_{i,y} = P(x_i \mid y)$

$$\theta_{i,y} = \frac{\text{count}(x_i, y)}{\sum_{i=1}^V \text{count}(x_i, y)} \quad (13)$$

$$= \frac{\text{count}(x_i, y)}{\text{count}(y)} \quad (14)$$

Therefore

$$\sum_{i=1}^V \theta_{i,y} = 1 \quad (15)$$

Example: Build a vocabulary

Label	Text	Words
NEG	just plain boring	3
NEG	entirely predictable and lacks energy	5
NEG	no surprises and very few laughs	6
Pos	very powerful	2
Pos	the most fun film of the summer	7

Vocab:

- ▶ $\mathcal{V} = \{\text{just, plain, boring, entirely, } \dots, \text{summer}\}$
- ▶ $V = 20$

Example: Estimate $p(x_i | y)$

Label	Text	Words
NEG	just plain boring	3
NEG	entirely predictable and lacks energy	5
NEG	no surprises and very few laughs	6
Pos	very powerful	2
Pos	the most fun film of the summer	7

For example

$$p(\text{just} | \text{NEG}) = \frac{1}{3 + 5 + 6} = \frac{1}{14} \quad p(\text{just} | \text{POS}) = \frac{0}{2 + 7} = \frac{0}{9}$$

Example: Prediction

For a given sentence “predictable and boring”

$$p(\text{predictable} \mid \text{NEG}) = \frac{1}{14} \quad p(\text{predictable} \mid \text{POS}) = \frac{0}{9}$$

$$p(\text{and} \mid \text{NEG}) = \frac{2}{14} \quad p(\text{and} \mid \text{POS}) = \frac{0}{9}$$

$$p(\text{boring} \mid \text{NEG}) = \frac{1}{14} \quad p(\text{boring} \mid \text{POS}) = \frac{0}{9}$$

Therefore,

$$p(S, \text{neg}) = \frac{3}{5} \cdot \frac{1}{14} \cdot \frac{2}{14} \cdot \frac{1}{14}$$

$$p(S, \text{pos}) = \frac{2}{5} \cdot \frac{0}{9} \cdot \frac{0}{9} \cdot \frac{0}{9}$$

Example: Prediction (II)

For a given sentence “predictable with no fun”

$$p(\text{predictable} \mid \text{NEG}) = \frac{1}{14} \quad p(\text{predictable} \mid \text{POS}) = \frac{0}{9}$$

$$p(\text{no} \mid \text{NEG}) = \frac{1}{14} \quad p(\text{and} \mid \text{POS}) = \frac{0}{9}$$

$$p(\text{fun} \mid \text{NEG}) = \frac{0}{14} \quad p(\text{fun} \mid \text{POS}) = \frac{1}{9}$$

Therefore,

$$p(S, \text{neg}) = \frac{3}{5} \cdot \frac{1}{14} \cdot \frac{1}{14} \cdot \frac{0}{14}$$

$$P(S, \text{pos}) = \frac{2}{5} \cdot \frac{0}{9} \cdot \frac{0}{9} \cdot \frac{1}{9}$$

Smoothing

To eliminate zero probability, adding a small number α to the counts:

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \alpha}{\sum_{i=1}^V \{\text{count}(x_i, y) + \alpha\}} \quad (16)$$

$$= \frac{\text{count}(x_i, y) + \alpha}{\sum_{i=1}^V \text{count}(x_i, y) + \alpha V} \quad (17)$$

with $\alpha > 0$ as a **hyper-parameter**.

For example, with $\alpha = 1$ and $V = 20$

$$p(\text{fun} | \text{NEG}) = \frac{0 + 1}{14 + 20} \quad p(\text{fun} | \text{POS}) = \frac{1 + 1}{9 + 20}$$

An Alternative View

- ▶ Likelihood: $p(\mathbf{x} \mid y; \boldsymbol{\theta}) \propto \prod_{i=1}^V \theta_{i,y}^{x_i}$
- ▶ Prior: $p(y; \boldsymbol{\theta}) = \theta_y$

$$p(\mathbf{x}, y; \boldsymbol{\theta}) \propto \theta_y \cdot \prod_{i=1}^V \theta_{i,y}^{x_i} \quad (18)$$

Or

$$\log p(\mathbf{x}, y; \boldsymbol{\theta}) \propto \log \theta_y + \sum_{i=1}^V (x_i \cdot \log \theta_{i,y}) \quad (19)$$

An Alternative View (II)

$$\log p(\mathbf{x}, y) \propto \log \theta_y + \sum_{i=1}^V (x_i \cdot \log \theta_{i,y}) \quad (20)$$

$$= b_y + \mathbf{w}_y^\top \mathbf{x} \quad (21)$$

where

$$b_y = \log \theta_y$$

$$\mathbf{w}_y^\top = [\log \theta_{1,y}, \dots, \log \theta_{V,y}]$$

$$\mathbf{x}^\top = [x_1, \dots, x_V]$$

It's a linear classifier (with respect to \mathbf{x})

Classification Evaluation

A Development Set

- ▶ Training set $\mathcal{T} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- ▶ Development set $\mathcal{D} = \{(\mathbf{x}^{(j)}, y^{(j)})\}_{j=1}^M$
- ▶ Test set $\mathcal{U} = \{(\mathbf{x}^{(l)}, y^{(l)})\}_{l=1}^L$

Cross-validation

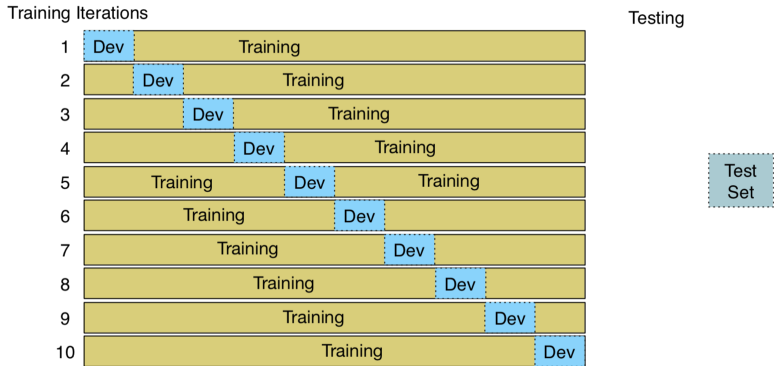


Figure: 10-fold cross validation.

Evaluation Measurements

- ▶ Accuracy
- ▶ Precision, recall, and F-measure

[Eisenstein, 2018, Sec 4.4]

Accuracy

Given N examples, with $y^{(i)}$ is the ground-truth label of the i -th example, and $\hat{y}^{(i)}$ is the predicted label

$$\text{ACC}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_{i=1}^N \delta(y^{(i)}, \hat{y}^{(i)}) \quad (22)$$

δ function is defined as

$$\delta(y, \hat{y}) = \begin{cases} 1 & y = \hat{y} \\ 0 & y \neq \hat{y} \end{cases} \quad (23)$$

Confusion Matrix

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	True Positive (TP)	False Positive (FP)
	NEGATIVE	False Negative (FN)	True Negative (TN)

Accuracy:

$$\text{acc}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (24)$$

Recall, Precision and F Measure

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	True Positive (TP)	False Positive (FP)
	NEGATIVE	False Negative (FN)	True Negative (TN)

Recall:

$$r(\mathbf{y}, \hat{\mathbf{y}}) = \frac{TP}{TP + FN} \quad (25)$$

Precision:

$$p(\mathbf{y}, \hat{\mathbf{y}}) = \frac{TP}{TP + FP} \quad (26)$$

F measure:

$$F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2 \cdot p \cdot r}{p + r} \quad (27)$$

Example: Balanced case

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	480	30
	NEGATIVE	20	470

- ▶ Accuracy: $\text{acc} = \frac{480+470}{480+20+30+470} = 0.95$
- ▶ Precision: $p = \frac{480}{480+30} \approx 0.94$
- ▶ Recall: $r = \frac{480}{480+20} \approx 0.96$
- ▶ F-measure: $F = \frac{2 \times 0.94 \times 0.96}{0.94 + 0.96} \approx 0.95$

Example: Unbalanced case

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	80	30
	NEGATIVE	20	870

- ▶ Accuracy: $\text{acc} = \frac{80+870}{80+20+30+870} = 0.95$
- ▶ Precision: $p = \frac{80}{80+30} \approx 0.73$
- ▶ Recall: $r = \frac{80}{80+20} = 0.80$
- ▶ F-measure: $F = \frac{2 \times 0.94 \times 0.96}{0.94 + 0.96} \approx 0.76$

Reference



Eisenstein, J. (2018).
Natural Language Processing.
MIT Press.



Jurafsky, D. and Martin, J. (2019).
Speech and language processing.



Pang, B., Lee, L., and Vaithyanathan, S. (2002).
Thumbs up?: sentiment classification using machine learning techniques.
In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages
79–86. Association for Computational Linguistics.