Using Probabilities in Understanding Sentiment

Aida Nematzadeh nematzadeh@berkeley.edu

Language Learning and Probabilities

Learning word meanings:



dax means dog



Look at the dax!

A dax!

Language Learning and Probabilities

Anticipating the next words:

She is drinking coffee lemonade

a chair doogh

In a cold winter day, she is drinking hot cocoa

Using knowledge given the situation: hot drink is better on a cold day.

lemonade

Combining Evidence & Knowledge

In a cold winter day, she is drinking hot cocoa h₁ evidence

lemonade h₂

The Bayes rule:



posterior $P(h|e) = \frac{P(e|h)P(h)}{\sum_{h'}(e|h')P(h')}$ P(e)

Evidence & Knowledge in Action

In a cold winter day, she is drinking hot cocoa h₁ evidence lemonade h₂

Let's calculate the probability of each hypothesis!

When Hypotheses Are Classes









H₂: the dog class

Many language processing tasks are classification!

Sentiment Analysis

Positive/negative orientation (sentiment) of text:

- Book, restaurant, movie and product reviews
- Political text



It's just <u>gorgeous</u>, like a flipbook made of dreamy vintage postcards that are somehow about contemporary life in Los Angeles.



The cinematography and special effects are fantastic, but don't actually compensate for a weak storyline, and forgettable musical numbers.

Spam Detection

Classify email to spam and non-spam

	Margaret Linda Hogan	My Greetings, - My Greetings, I am Ms Margaret Linda Hog
*	Mailbox Validation	Your mailbox Will be closed if has exceed storage limit.
*	E-mail Verification Port.	Verify Your Account to avoid closure Dear nematzadeh
*	S.Mani ICCAIRO 20.	International Conference on Control, Artificial Intelligence
*	Admin Notification	Re-validate your account to avoid termination - Dear nen

"Dear winner"; "Click this link"; "Urgent: send me your credit card information"; ...

Authorship Attribution

Find the text author and author's characteristics:

• Gender, age, etc

Study claims Agatha Christie had Alzheimer's

Textual analysis detects signs of sharply declining faculties towards the end of beloved mystery writer's life

An in-depth analysis of Agatha Christie's novels has suggested that the muchloved author of more than 80 mysteries was suffering from Alzheimer's disease.

Academics at the University of Toronto studied a selection of Christie's novels written between the ages of 28 and 82, counting the numbers of different words, indefinite nouns and phrases used in each.

Learning Word Representations

Word2Vec models:



Class 1: positive examples



Class 2: negative examples

Classification: Formalism

Given an input & fixed classes C=c₁, c₂, ..., c_M, find:

- the predicted class c_i
- The probability of each class

Supervised training: uses data points & their gold-standard labels, (d_1, c_1) , (d_2, c_2) , ..., (d_N, c_N)

Goal: Find the correct class for the new data point

Classification Algorithms

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

posterior

Full distribution

Can use prior info

Naive Bayes

Discriminative models

P(d|c).

likelihood

Easier to train

Used more in practice

Support Vector Machine & Word2Vec

Naive Bayes Classifier: Input

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!



bag-of-word representation

Find the best/correct class for the document:

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c|d) \operatorname{posterior}$$

i.e., the class with the maximum posterior prob. **M.A.P estimation**

What are the disadvantages of this?

How do we calculate P(c | d)? Let's expand it.

Find the best/correct class for the document:

$$\begin{split} \hat{c} &= \operatorname*{argmax}_{c \in C} P(c|d) \text{ posterior} \\ \hat{c} &= \operatorname*{argmax}_{c \in C} \frac{\underset{P(d|c)P(c)}{\text{prior}}}{P(d|c)P(c)} \\ \hline P(d)_{\text{Is P(d) of important?}} \end{split}$$

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(d|c)P(c)$$

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(d|c) P(c)$$

How do we calculate P(d | c)?

Represent the document as a set of features.

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} \begin{array}{c} \overbrace{P(f_1, f_2, \dots, f_n | c)}^{\text{likelihood}} \begin{array}{c} \text{prior} \\ \overbrace{P(c)}^{\text{prior}} \end{array}$$

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} \begin{array}{c} \overbrace{P(f_1, f_2, \dots, f_n | c)}^{\text{likelihood}} \begin{array}{c} \text{prior} \\ \overbrace{P(c)}^{\text{prior}} \end{array}$$

How do we calculate $P(f_1, f_2, ..., f_n | c)$?

Naive Bayes assumption: Independent features.

$$P(f_1, f_2, \dots, f_n | c) = P(f_1 | c) \cdot P(f_2 | c) \cdot \dots \cdot P(f_n | c)$$
$$c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{f \in F} P(f | c)$$

Naive Bayes on Documents

Given a new document:

positions \leftarrow all word positions in test document $c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{i \in positions} P(w_i | c)$

Practical trick:

 $c_{NB} = \underset{c \in C}{\operatorname{argmax}} \log P(c) + \sum_{i \in positions} \log P(w_i|c)$

Let's Look an Example!

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

$$c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{i \in positions} P(w_i|c)$$

Let's Look an Example!

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

$$\hat{P}(c) = \frac{N_c}{N_{doc}}$$

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

Add-one Smoothing

What is P("predictable" | +)?

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

$$= \frac{count(w_i, c) + 1}{(\sum_{w \in V} count(w, c)) + |V|}$$
Adding prior probability to unseen words.

Unknown words (words not in V) are removed.

Naive Bayes Demo

Affective Text: Data Annotated for Emotions and Polarity

- Affective Text is a data set consisting of 1000 test headlines and 200 development headlines, each of them annotated with the six Eckman emotions and the polarity orientation. [download] (July 13, 2007).
 - Carlo Strapparava and Rada Mihalcea, SemEval-2007 Task 14: Affective Text, in Proceedings of the 4th International Workshop on the Semantic Evaluations (SemEval 2007), Prague, Czech Republic, June 2007. [pdf]
 - Read more about the task <u>here</u>.

Improving Sentiment Analysis

Binarize the frequency of a word in a document.

Four original documents:	NB Counts + -		Binary Counts		
_ it was pathetic the worst part was the	and	2	0	1	0
boxing scenes	boxing	0	1	0	1
boxing scenes	film	1	0	1	0
 no plot twists or great scenes 	great	3	1	2	1
+ and satire and great plot twists	it	0	1	0	1
+ great scenes great film	no	0	1	0	1
	or	0	1	0	1
After per-document binarization:	part	0	1	0	1
- it was pathetic the worst part boxing	pathetic	0	1	0	1
it was pathetic the worst part boxing	plot	1	1	1	1
scenes	satire	1	0	1	0
 no plot twists or great scenes 	scenes	1	2	1	2
+ and satire great plot twists	the	0	2	0	1
+ great scenes film	twists	1	1	1	1
0	was	0	2	0	1

0

1

0

1

worst

What are other possible scores?

Improving Sentiment Analysis

Dealing with negation:

I didn't love the food **vs** I loved the food.

Add a prefix after negation (n't, not, no, never) *I didn't NOT-love* NOT-the NOT-food:

Improving Sentiment Analysis

Use **sentiment lexicon**s when there is no enough labeled training data.

- Linguistic Inquiry and Word Count (LIWC)
- NRC Word-Emotion Association Lexicon
- + : admirable, beautiful, confident, dazzling, ecstatic, favor, glee, great
- : awful, bad, bias, catastrophe, cheat, deny, envious, foul, harsh, hate

Use Two features -- positive vs negative words

Naive Bayes as a Language Model

Model prob. of generating sentences in the lang. NB: set of class-specific unigram lang. models.

Likelihood of a sentence: $P(s|c) = \prod_{i \in positions} P(w_i|c)$

How can we calculate P(s)?

Evaluating Classifier Performance

Consider a binary detection task.

• Label text spam (*positive*) or ~spam (*negative*).

Gold labels: human labels used as ground truth.

• Gold label is either spam (*true*) or ~spam (*false*).

Need *metrics* to quantify classifier's performance.

Evaluation Metrics: Precision/Recall

Accuracy is a natural metric, but rarely used.

• Problem with unbalanced classes. (Why?)

Precision: correctly labeled + / all labeled +.

Recall: correctly labeled + / all truly +.

F-measure combines both: harmonic mean.

- Weighs min of two more heavily.
- Who to invite to your valentine's day party?

Evaluation Metrics

