

Thumbs Up? Sentiment Classification using Machine Learning Techniques

Pang, Lee, Vaithyanathan - EMNLP 2002

Slides by: Dor Cohen, Itai Gat

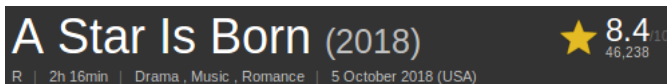
IE&M @ Technion

October 23, 2018

Problem definition

Sentiment analysis

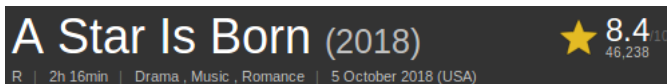
Motivation



- Should we watch this movie?

Sentiment analysis

Motivation



- Should we watch this movie?
- Ideally: read each review and decide



Amazing!

[gastraube](#) 31 August 2018

This movie is perfect, it made me cry from the beginning to the end, it deserves several Oscar's nominations. Lady Gaga is a true artist.

Sentiment analysis

Motivation

A Star Is Born (2018) ★ 8.4¹⁰
46,238

R | 2h 16min | Drama , Music , Romance | 5 October 2018 (USA)

- Should we watch this movie?
- Ideally: read each review and decide



10/10

Amazing!

[gastraube](#) 31 August 2018

This movie is perfect, it made me cry from the beginning to the end, it deserves several Oscar's nominations. Lady Gaga is a true artist.



2/10

Directing and acting miscarriage

[dailynewsandinformations](#) 5 October 2018

Maybe only the sound on this film is up to par...everything else is down in mediocrity...

Problem definition

- Find mapping from text to binary label
 - ▶ Supervised learning

Problem definition

- Find mapping from text to binary label
 - ▶ Supervised learning
- For m numeric features (extracted from text) we define:

Definition (Binary classifier)

$$f : X \rightarrow y$$

where $X \in \mathbb{R}^m, y \in \{0, 1\}$

Problem definition

- Find mapping from text to binary label
 - ▶ Supervised learning
- For m numeric features (extracted from text) we define:

Definition (Binary classifier)

$$f : X \rightarrow y$$

where $X \in \mathbb{R}^m, y \in \{0, 1\}$

- Evaluate by loss function
- e.g., Zero-one loss: $L(x, y, f_w) = \mathbf{1}\{f_w(x) \neq y\}$
 - ▶ w denotes learned parameters

Data: IMDB Movie Reviews

- Lucky for us: user rating provides **supervised** learning
- Converted into 3 categories:
 - ▶ *Positive, negative, (neutral - not used)*
- Avoid bias issues:
 - ▶ 20 reviews per author per sentiment
 - ▶ 752 negative vs 1301 positive
 - ▶ total of 144 reviewers

Human based sentiment classifiers

- **Hypothesis:** certain words indicate on sentiment type
- **Test:** count positive vs. negative words

Human	Proposed words	Accuracy	Ties ¹
1	positive (5): dazzling, brilliant.. negative (5): suck, terrible..	58%	75%
2	positive (11): gripping, spectacular.. negative (6): cliched, boring..	64%	39%

Table: Baseline results for human word lists, data is balanced (700 vs. 700)

¹Documents percentage where sentiments rated equally

Human based sentiment classifiers

Should we worry about high rate of ties?

- Proposed list is relatively short (usually effect is 0 vs. 0)
 - ▶ Not necessarily the reason for low accuracy!
- Authors propose their list
 - ▶ Backed up with statistics

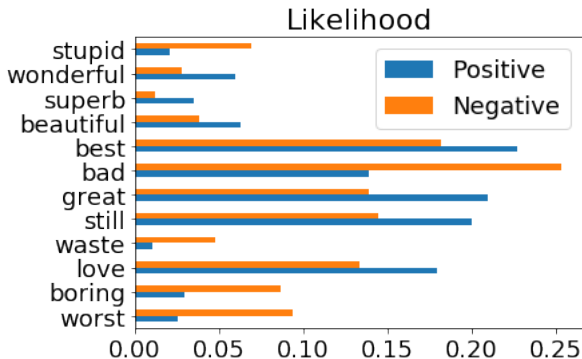
Human	Proposed words	Accuracy	Ties
3+Stats	positive (7): love, wonderful.. negative (7): bad, worst, '?', '!',...	69%	16%

Table: Results where words (total 14) were chosen based on data statistics

Human based sentiment classifiers

(2018) Data analysis

We reproduced the analysis, following are example estimates



- Note words occurrences were binarized

Methods

Bag of words

Example

- d_1 : "Audio rocks"
- d_2 : "Act boring"

	act	audio	boring	rocks
mapping	0	1	2	3
d1	0	1	0	1
d2	1	0	1	0

- d_3 : "Boring effects", $d_3^{bow} = ?$

Bag of words

Example

- d_1 : "Audio rocks"
- d_2 : "Act boring"

	act	audio	boring	rocks
mapping	0	1	2	3
d1	0	1	0	1
d2	1	0	1	0

- d_3 : "Boring effects", $d_3^{bow} = ?$
- $d_3^{bow} = (0, 0, 1, 0)$

Naive Bayes classifier

- Assign class which maximizes: $c^* = \operatorname{argmax}_c P(c|d)$

Naive Bayes classifier

- Assign class which maximizes: $c^* = \operatorname{argmax}_c P(c|d)$
- Recap:

Definition (Bayes theorem)

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

Naive Bayes classifier

- Assign class which maximizes: $c^* = \operatorname{argmax}_c P(c|d)$
- Recap:

Definition (Bayes theorem)

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

- To estimate $P(d|c)$ we **naively** assume f_i 's are independent

Naive Bayes classifier

- Assign class which maximizes: $c^* = \operatorname{argmax}_c P(c|d)$
- Recap:

Definition (Bayes theorem)

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

- To estimate $P(d|c)$ we **naively** assume f_i 's are independent
 - ▶ Hence $\widehat{P(d|c)} = \prod_{i=1}^m P(f_i|c)^{n_i(d)}$

Naive bayes example

Assume the following Bow model with 4 documents for each class:

	act	audio	boring	rocks
mapping	0	1	2	3
positive	1	3	1	3
negative	3	1	3	1

Naive bayes example

Assume the following Bow model with 4 documents for each class:

	act	audio	boring	rocks
mapping	0	1	2	3
positive	1	3	1	3
negative	3	1	3	1

- $d_3 = \text{"Boring effects"} , d_3^{bow} = (0, 0, 1, 0)$

Naive bayes example

Assume the following Bow model with 4 documents for each class:

	act	audio	boring	rocks
mapping	0	1	2	3
positive	1	3	1	3
negative	3	1	3	1

- $d_3 = \text{"Boring effects"} , d_3^{bow} = (0, 0, 1, 0)$
- Recall $\widehat{P}(d|c) = \prod_{i=1}^m P(f_i|c)^{n_i(d)}$

Naive bayes example

Assume the following Bow model with 4 documents for each class:

	act	audio	boring	rocks
mapping	0	1	2	3
positive	1	3	1	3
negative	3	1	3	1

- $d_3 = \text{"Boring effects"} , d_3^{bow} = (0, 0, 1, 0)$
- Recall $\widehat{P(d|c)} = \prod_{i=1}^m P(f_i|c)^{n_i(d)}$
- $P(d|pos) = \frac{1}{4}^0 * \frac{1}{4}^0 * \frac{1}{4}^1 * \frac{1}{4}^0 = \frac{1}{4}$

Naive bayes example

Assume the following Bow model with 4 documents for each class:

	act	audio	boring	rocks
mapping	0	1	2	3
positive	1	3	1	3
negative	3	1	3	1

- $d_3 = \text{"Boring effects"} , d_3^{bow} = (0, 0, 1, 0)$
- Recall $\widehat{P(d|c)} = \prod_{i=1}^m P(f_i|c)^{n_i(d)}$
- $P(d|pos) = \frac{1}{4}^0 * \frac{1}{4}^0 * \frac{1}{4}^1 * \frac{1}{4}^0 = \frac{1}{4}$
- $P(d|neg) = \frac{1}{4}^0 * \frac{1}{4}^0 * \frac{3}{4}^1 * \frac{1}{4}^0 = \frac{3}{4}$

SVM vs. Logistic regression

- Can view both parametrically: $f(x_i, W) = Wx_i + b$
- **Train** with gradient descent:
 - ▶ Initialize parameters
 - ▶ 1. Update parameters following loss gradient
 - ▶ 2. Repeat until convergence

SVM vs. Logistic regression

- Can view both parametrically: $f(x_i, W) = Wx_i + b$
- **Train** with gradient descent:
 - ▶ Initialize parameters
 - ▶ 1. Update parameters following loss gradient
 - ▶ 2. Repeat until convergence

Methods differ in their loss functions:

Definition (SVM loss)

$$L_i = \sum_{j \neq i} \max(0, s_j - s_{y_i} + \delta), \text{ where } s_j = f(x_i, W)_j$$

Definition (Softmax loss)

$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right)$$

SVM vs. Logistic regression

Evaluate loss example

$$\begin{array}{|c|c|c|c|} \hline 0.04 & 0.5 & -0.1 & 0.2 \\ \hline 0.7 & 0.2 & +0.03 & -0.05 \\ \hline \end{array}
 \quad
 \begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 1 \\ \hline 0 \\ \hline \end{array}
 \quad
 +
 \quad
 \begin{array}{|c|} \hline 0.2 \\ \hline -0.3 \\ \hline \end{array}$$

W X_i b

$y_i = 0$

$$L_i = \sum_{j \neq i} \max(0, s_i - s_j - \Delta)$$

SVM loss

$$L_i = \max(0, -0.37 - 0.1 + 1) = 0.63$$

$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right)$$

Softmax loss

$$\begin{array}{|c|} \hline 0.1 \\ \hline -0.37 \\ \hline \end{array}
 \xrightarrow{\text{exp}}
 \begin{array}{|c|} \hline 1.1 \\ \hline 0.76 \\ \hline \end{array}
 \xrightarrow{\text{normalize}}
 \begin{array}{|c|} \hline 0.6 \\ \hline 0.4 \\ \hline \end{array}
 \rightarrow L_i = -\log 0.6 = 0.22$$

"Probabilities"

Results

Experimental setup

- Features:

- ▶ d : "Audio quality rocks"
- ▶ Unigram: {'audio', 'rocks',... }
- ▶ Bigram: {'audio quality',... }
- ▶ N-gram !

Setup

- Unigram/bigram appear at least 4/7 times
- Uniform class distributions ("balanced" dataset)
- 3-fold average accuracies reported
- Expressions with negation handled with unified NOT tag
- Punctuation treated as separate lexicon, no stemming used

Results and discussion

Using feature presence

ID	Features	count	NB	ME	SVM
2	unigrams	16165	81.0%	80.4%	82.9%
3	uni+bigrams	32330	80.6%	80.8	82.7%
4	bigrams	16165	77.3%	77.4%	77.1%
5	unigrams+POS	16695	81.5%	80.4%	81.9%
6	adjectives	2633	77.0%	77.7%	75.1%
7	top 2633 unigrams	2633	80.3%	81.0%	81.4%
8	unigram+position	22430	81.0%	80.1%	81.6%

Results and discussion

Using feature presence

ID	Features	count	NB	ME	SVM
2	unigrams	16165	81.0%	80.4%	82.9%
3	uni+bigrams	32330	80.6%	80.8	82.7%
4	bigrams	16165	77.3%	77.4%	77.1%
5	unigrams+POS	16695	81.5%	80.4%	81.9%
6	adjectives	2633	77.0%	77.7%	75.1%
7	top 2633 unigrams	2633	80.3%	81.0%	81.4%
8	unigram+position	22430	81.0%	80.1%	81.6%

- Adding bigrams doesn't improve results; Bigrams alone is worse
- Part-of-speech: "I love this movie" vs. "This is a love story"
- Position based on dividing text into quarters

Conclusions

Conclusions

- Unigrams presence setting achieves the best performance
 - ▶ Apply feature selection algorithms
- Contrarily, performance isn't comparable to topic classification

Review example

"This film should be brilliant. It sounds like a great plot, the actors are first grade, and the supporting cast is good as well, and Stallone is attempting to deliver a good performance. However, it can't hold up."

Conclusions

- Unigrams presence setting achieves the best performance
 - ▶ Apply feature selection algorithms
- Contrarily, performance isn't comparable to topic classification

Review example

"This film should be brilliant. It sounds like a great plot, the actors are first grade, and the supporting cast is good as well, and Stallone is attempting to deliver a good performance. However, it can't hold up."

- Difficult for bag-of-words classifiers
- Authors suggest determining the **focus** of each sentence

Thank you for participating!
Questions?