

## Sentiment Analysis



## Sentiment Analysis

this movie was **great!** would **watch again** **+**

the movie was **gross** and **overwrought**, but I **liked** it **+**

this movie was **not** really very **enjoyable** **-**

- Bag-of-words doesn't seem sufficient (discourse structure, negation)
- There are some ways around this: extract bigram feature for "not X" for all X following the not



## Pang et al. (2002)

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

▸ Simple feature sets can do pretty well!

▸ Learning alg. doesn't matter too much

- ME = "Maximum Entropy" = what we call Logistic Regression

Bo Pang, Lillian Lee, Shivakumar Vaithyanathan (2002)



## Wang and Manning (2012)

- 10 years later — revisited basic BoW classifiers vs. other methods

Method	RT-s	MPQA
MNB-uni	77.9	85.3
MNB-bi	<b>79.0</b>	<b>86.3</b>
SVM-uni	76.2	86.1
SVM-bi	77.7	<b>86.7</b>
NBSVM-uni	<b>78.1</b>	85.3
NBSVM-bi	<b>79.4</b>	<b>86.3</b>
RAE	76.8	85.7
RAE-pretrain	77.7	<b>86.4</b>
Voting-w/Rev.	63.1	81.7
Rule	62.9	81.8
BoF-noDic.	75.7	81.8
BoF-w/Rev.	76.4	84.1
Tree-CRF	77.3	86.1

Before neural nets had taken off — results weren't that great

Kim (2014) CNNs **81.5** **89.5**

Wang and Manning (2012)

## Multiclass Examples



## Entailment

- Three-class task over sentence pairs

A soccer game with multiple males playing.

ENTAILS

Some men are playing a sport.

- Not clear how to do this with simple bag-of-words features

A black race car starts up in front of a crowd of people.

CONTRADICTS

A man is driving down a lonely road

A smiling costumed woman is holding an umbrella.

NEUTRAL

A happy woman in a fairy costume holds an umbrella.

Bowman et al. (2015)



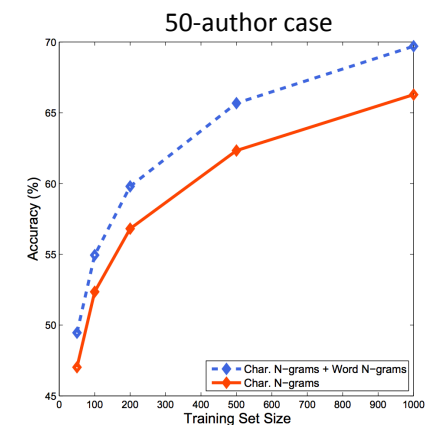
## Authorship Attribution

- Statistical methods date back to 1930s and 1940s
  - Based on handcrafted heuristics like stopword frequencies
  - Early work: Shakespeare's plays, Federalist papers (Hamilton v. Madison)
- Twitter: given a bunch of tweets, can we figure out who wrote them?
  - Schwartz et al. EMNLP 2013: 500M tweets, take 1000 users with at least 1000 tweets each
- Task: given a held-out tweet by one of the 1000 authors, who wrote it?



## Authorship Attribution

- SVM with character 4-grams, words 2-grams through 5-grams
- 1000 authors, 200 tweets per author => 30% accuracy
- 50 authors, 200 tweets per author => 71.2% accuracy



Schwartz et al. (2013)



## Authorship Attribution

- k-signature: n-gram that appears in k% of the authors tweets but not appearing for anyone else — suggests why these are so effective

Signature Type	10%-signature	Examples
Character n-grams	' ^ _ '	REF oh ok ^ _ ^ Glad you found it!
		Hope everyone is having a good afternoon ^ _ ^
		REF Smirnoff lol keeping the goose in the freezer ^ _ ^
	'yew '	gurl <b>yew</b> serving me tea nooch
		REF about wen <b>yew</b> and ronnie see each other
		REF lol so <b>yew</b> goin to check out tini's tonight huh???

Schwartz et al. (2013)