# Twitter as a Corpus for Sentiment Analysis and Opinion Mining

Alexander Pak, Patrick Paroubek

Université Paris-Sud 11, LIMSI-CNRS

# Microblogging

Microblogging = posting small blog entries

Eg.: "@alex: I'm presenting now my paper at LREC'10"

Platforms:
- Twitter
- Tumblr
- Plurk

# Twitter

# Twitter

Twitter – social network for publishing short messages (tweets)

1 tweet contains:
 maximum: 140 character
 in average: 1 sentence

More than 1 billion tweets per month

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Twitter for opinion mining

## People are expressing their opinions in tweets

Eg.: "CelineBG: @itsRyanButler u should come to Malta (europe) it's below Italy..we have sun nearly all year round =) we have amazing beaches =) follow me"

## Twitter is multilingual

## More than 14 billion tweets

## Twitter API for data retrieval

Content of Tweets

72   75

174        117

751

811

- ■ News (green)
- ■ Spam (magenta)
- ■ Self-Promotion (orange)
- ■ Pointless Babble (dark red)
- ■ Conversational (blue)
- ■ Pass Along Value (olive)

Kelly, Ryan, ed. (2009-08-12), "Twitter Study - August 2009" (PDF), Twitter Study Reveals Interesting ResultsAbout Usage, San Antonio, Texas: Pear Analytics. http://www.pearanalytics.com/wp-content/uploads/2009/08/Twitter-Study-August-2009.pdf

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Corpus collection

Use emoticons as noisy sentiment labels

## Positive tweets with :) =) :D
@mia_jones oh lovely! I'm heading to Malta & Italy next week!! Can't wait :)

## Negative tweets with :( :'( ;(
Supposed to be flying tonight, now stuck in Malta until Thursday. Homesick :(

## Use newspapers' tweets for neutral texts
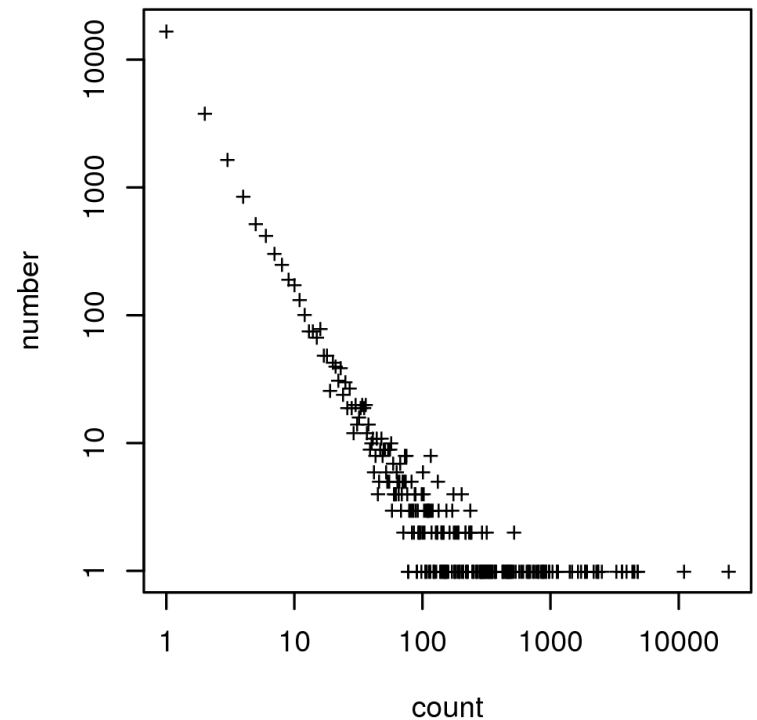@nytimes: Iron Man Defeats Robin Hood at North American Box Office

# Corpus analysis
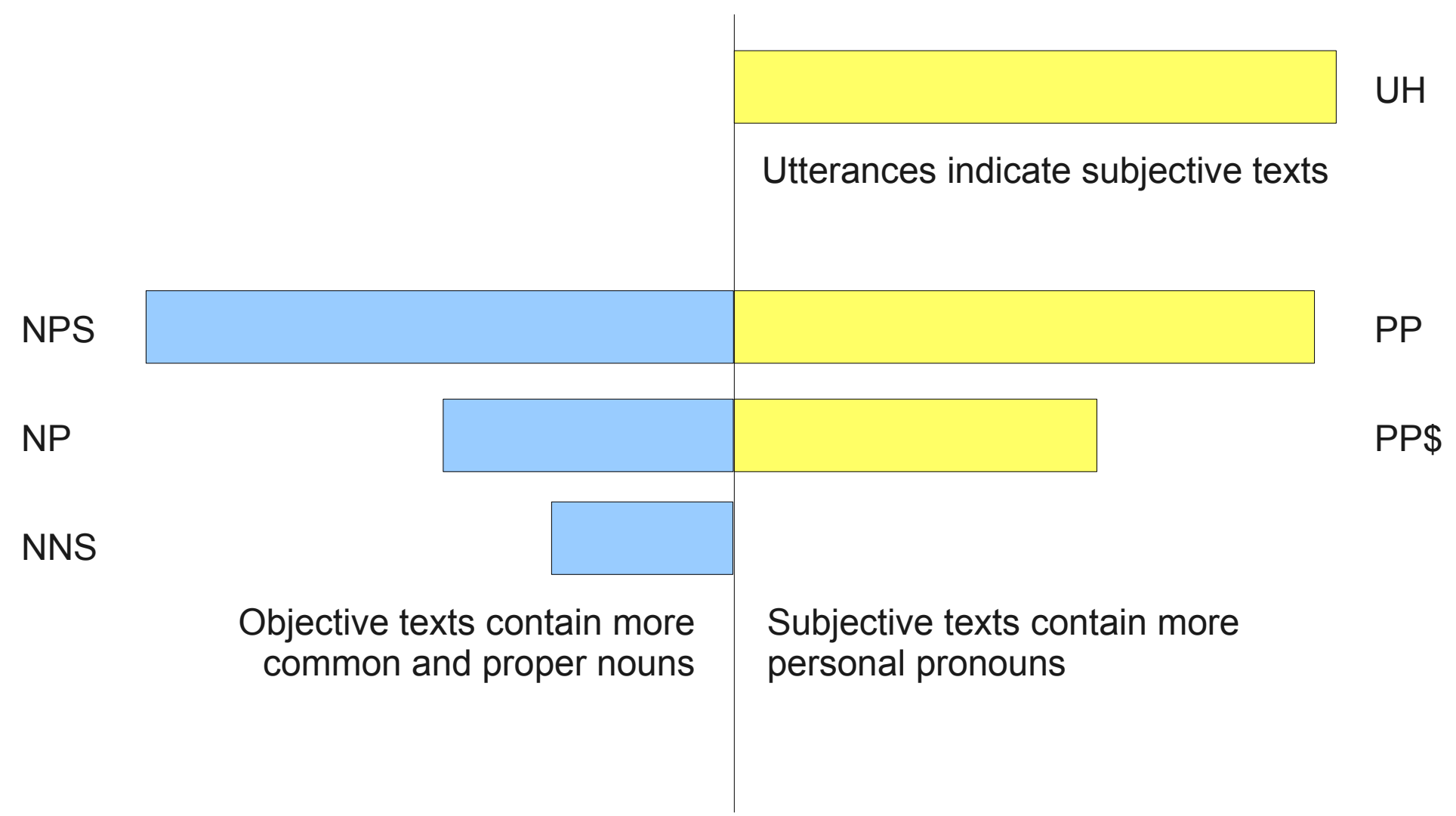
Collected 300'000 positive, negative and neutral tweets

Distribution of word frequencies is Zipfian

Use TreeTagger for POS tagging



Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Objective vs. subjective tweets



UH

Utterances indicate subjective texts

NPS — PP

NP — PP$

NNS

Objective texts contain more common and proper nouns

Subjective texts contain more personal pronouns

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Objective vs. subjective tweets

VBN ▮▮  VBP

VBD

Verbs in objective texts are usually in the 3d person

Authors write about themselves or address the audience

VBZ ▮▮  MD

VB

Past participle is used for stating facts

Modal verbs are used to express emotions

JJR ▮▮  JJS

Comparative adjectives state facts

Superlative adjectives express emotions

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Positive vs. negative tweets

VBN

VBD

Negative tweets often contain past tense

RBS

Superlative adjectives may indicate
positive tweets

POS

Positive tweets more often contain
possessive endings

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Building a classifier

Use the corpus to train a sentiment classifier

Use Naïve Bayes classifier

2 types of features: n-grams and POS

Bigrams showed the best performance

Handle negations by attaching negation particle
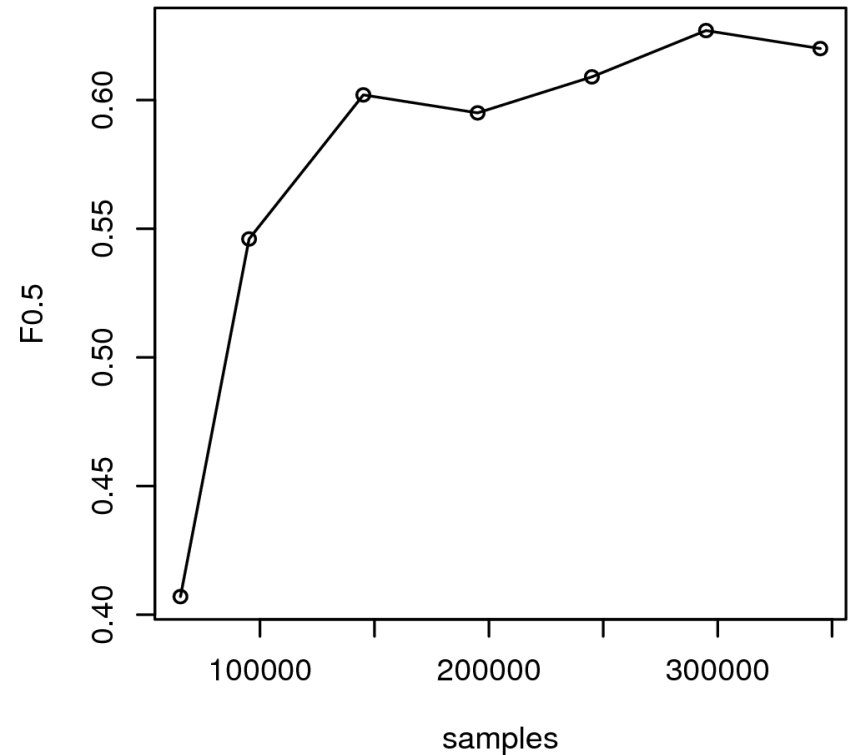Eg.: I do not like fish: I do+not, do+not like, not+like fish

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Building a classifier

Use hand annotated
tweets for evaluation:

Positive: 108
Negative: 75
Neutral: 33
Total: 216



Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Increasing accuracy

Classify tweets with high confidence of precision

Other tweets are left as "undecided"

"decision" = ratio of classified tweets

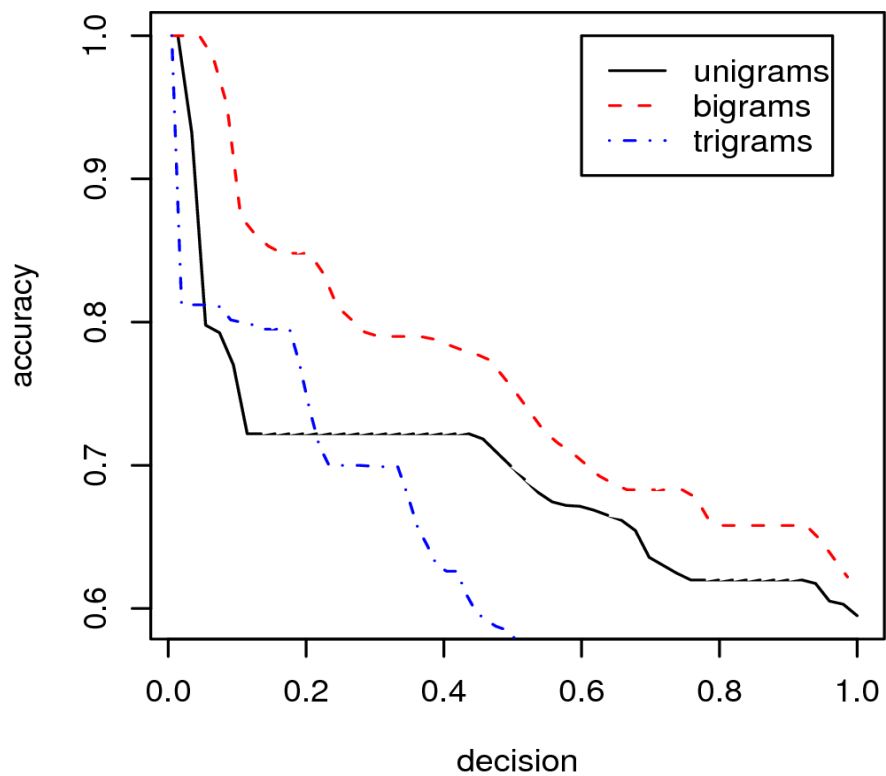Select n-grams with high salience (ignore n-grams with same frequency in all three sets)

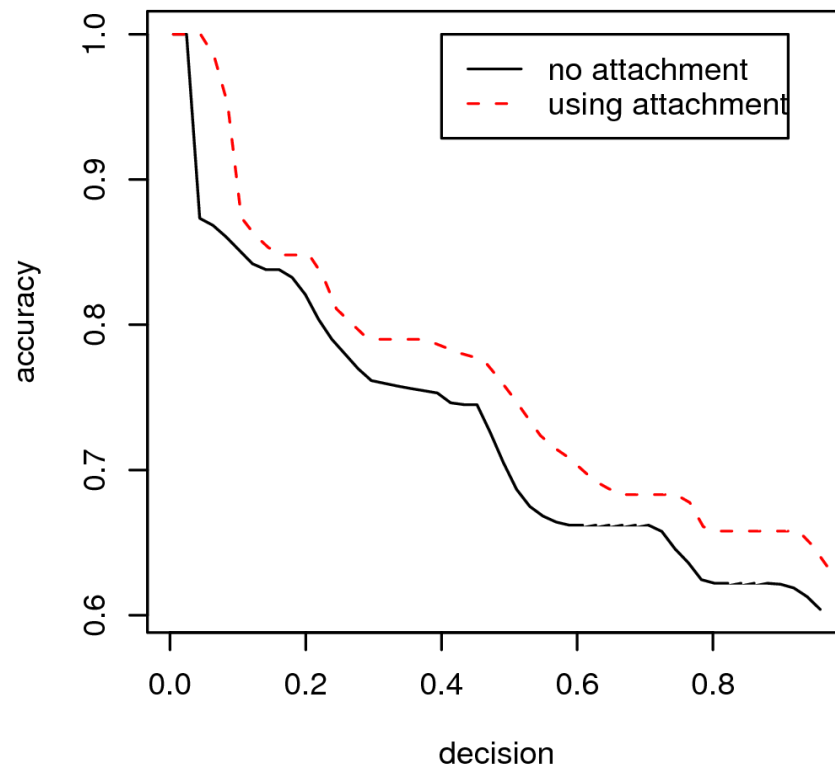Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Increasing accuracy

| N-gram | Salience |
|--------|----------|
| So sad | 0.975 |
| Miss my | 0.972 |
| So sorry | 0.962 |
| Love your | 0.961 |
| I'm sorry | 0.96 |
| Sad I | 0.959 |
| I hate | 0.959 |
| Lost my | 0.959 |
| Have great | 0.958 |

# Results



Comparison of n-gram order

Impact of negation attachment

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Prototype



Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Conclusion

Twitter can be used as a sentiment-labeled corpus

Naive-Bayes with bigram and POS features can perform a precise sentiment classification

Future work: collect more tweets, form a multilingual corpus

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)

# Thank you!

Twitter as a Corpus for Sentiment Analysis and Opinion Mining (A. Pak and P. Paroubek)