

Normalizing Social Media Texts by Combining Word Embeddings and Edit Distances in a Random Forest Regressor

Rob van der Goot
r.van.der.goot@rug.nl

28-05-2016

- 1 Problem
- 2 Error Detection
- 3 Generation
- 4 Ranking
- 5 Conclusion
- 6 Future Work

Outline

- 1 Problem
- 2 Error Detection
- 3 Generation
- 4 Ranking
- 5 Conclusion
- 6 Future Work

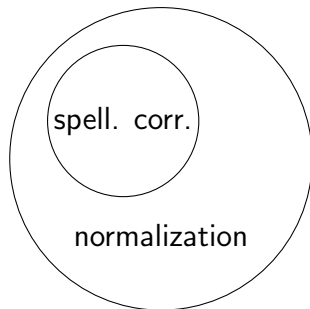
Problem

- Adapt Natural Language Processing pipelines to noisy (web) data

Problem

- Adapt Natural Language Processing pipelines to noisy (web) data
- Normalize

Spelling Correction vs. Normalization



Problem

Spelling Correction

abilites	abilities
teh	the
kingdon	kingdom

Problem

Normalization

abilites

teh

kingdon

doin

Bham

2

The

ggggrrrrreeeeeeaaaaaaattttttt

ur

abilities

the

kingdom

doing

Birmingham

to

There

great

your

Problem

Traditional spelling correction framework:

- Error detection
- Candidate generation
- Ranking of candidates

Problem

- Train set: 2,577 tweets from (Li and Liu 2014)
- Test set: LexNorm (Han and Baldwin 2011) 549 tweets

Outline

- 1 Problem
- 2 Error Detection**
- 3 Generation
- 4 Ranking
- 5 Conclusion
- 6 Future Work

Error Detection

Spelling correction:

- Dictionary lookup

Error Detection

- Often skipped in normalization methods
- Here as well, because the goal is to be used in a pipeline
- All tokens are considered to be a possible error/disfluency
- Recall = 100%
- But the original word is always kept!

Outline

- 1 Problem
- 2 Error Detection
- 3 Generation**
- 4 Ranking
- 5 Conclusion
- 6 Future Work

Spelling correction:

- Lexical edit distance
- Phonetic edit distance (Double Metaphone)

Spelling correction:

- Lexical edit distance
- Phonetic edit distance (Double Metaphone)
- Good results

Spelling correction:

- Lexical edit distance
- Phonetic edit distance (Double Metaphone)
- Good results
- So we use an existing system (Aspell)

Other disfluencies

- A more data aware model is necessary

Other disfluencies

- A more data aware model is necessary
- Semi-supervised

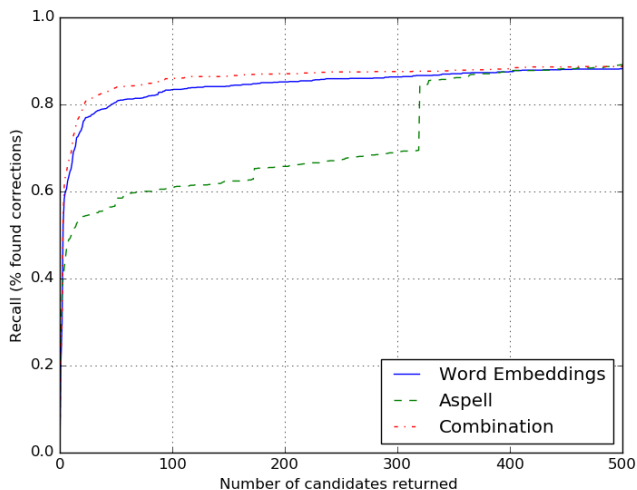
Other disfluencies

- A more data aware model is necessary
- Semi-supervised
- Word Embeddings

Word Embeddings

- Model taken from (Godin et al. 2015)
- Trained on 400 million Tweets
- 3,039,345 words
- Use cosine distance to find top-n words in vector-space

Generation



Outline

- 1 Problem
- 2 Error Detection
- 3 Generation
- 4 Ranking**
- 5 Conclusion
- 6 Future Work

Spelling correction:

- Combination of edit distances

Previous approaches:

- Ngram based approaches
- Combine Ranking with generation

Ranking

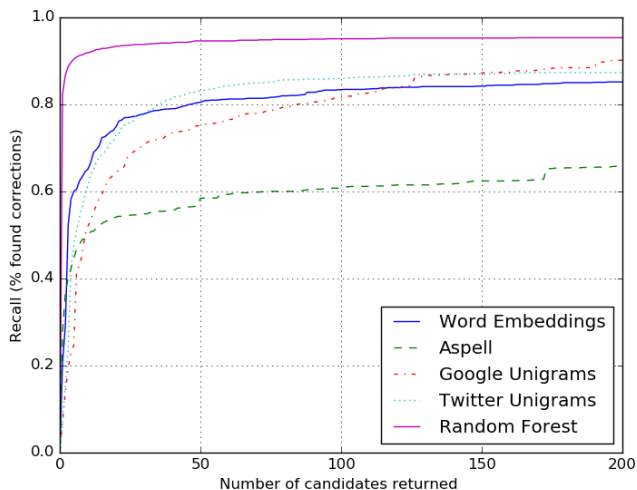
My approach:

- Use features from generation
- Supplement these features with N-Gram features
- Google Ngrams ¹ & Twitter Ngrams ²
- Combine all features in a Random Forest Classifier
- Default parameters Scikit Learn, except for the number of trees = 100

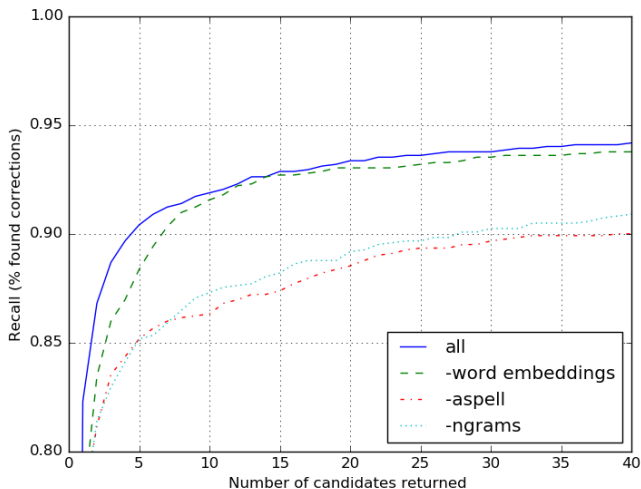
¹Brants and Franz 2006

²Herdağdelen 2013

Ranking



Ranking (ablation)



Ranking

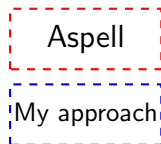
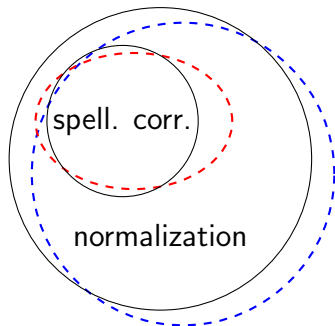
Ranking System	top1	top3	top10	top20	upper bound
(Li and Liu 2012)	73.0	81.9	86.7	89.2	94.2
(Li and Liu 2014)	77.14	86.96	93.04	94.82	95.90
(Li and Liu 2015)	87.58				
Our system	82.31	88.70	91.89	93.37	93.37

Outline

- 1 Problem
- 2 Error Detection
- 3 Generation
- 4 Ranking
- 5 Conclusion**
- 6 Future Work

Conclusion

Overview



Conclusion

For the normalization task:

- Word embeddings complement edit distances well
- A random forest classifier works very well for ranking
- This is a simple system, with a reasonable performance

Outline

- 1 Problem
- 2 Error Detection
- 3 Generation
- 4 Ranking
- 5 Conclusion
- 6 Future Work**

Future Work

- Multilingual/multiword embeddings
- Generation (build own language models)

Future Work

- Multilingual/multiword embeddings
- Generation (build own language models)
- Parameter tuning, add domain specific information
- Find candidate with: "word.*"

Future Work

- This system was created for use in a pipeline system

Future Work

- This system was created for use in a pipeline system
- Parse a word graph based on the output of this normalization

Future Work

<https://bitbucket.org/robvandergerg/errcor>