

# Embedding-based Query Spelling Correction

---

**Ines Zelch**<sup>1,2</sup>

Gustav Lahmann<sup>1</sup>

Matthias Hagen<sup>1</sup>

<sup>1</sup>Friedrich-Schiller-Universität Jena

<sup>2</sup>Leipzig University

WOWS, March 28 2024

`webis.de`

# Query Spelling Correction

I can't beleeve this is really happening. I just can't beleev it

believe  
Bellevue  
belief

[<https://i.stack.imgur.com/DCUcD.png>]

# Query Spelling Correction

Type	Misspelling		Correction
Deletion	corosion	→	corrosion
Insertion	occurr	→	occur
Space	abouta	→	about a
Special character	isnt	→	isn't
Substitution	grammer	→	grammar
Transposition	rewriet	→	rewrite

[Hagen et al. 2017, SIGIR]

# Query Spelling Correction

Type	Misspelling		Correction
Deletion	corosion	→	corrosion
Insertion	occurr	→	occur
Space	abouta	→	about a
Special character	isnt	→	isn't
Substitution	grammer	→	grammar
Transposition	rewriet	→	rewrite

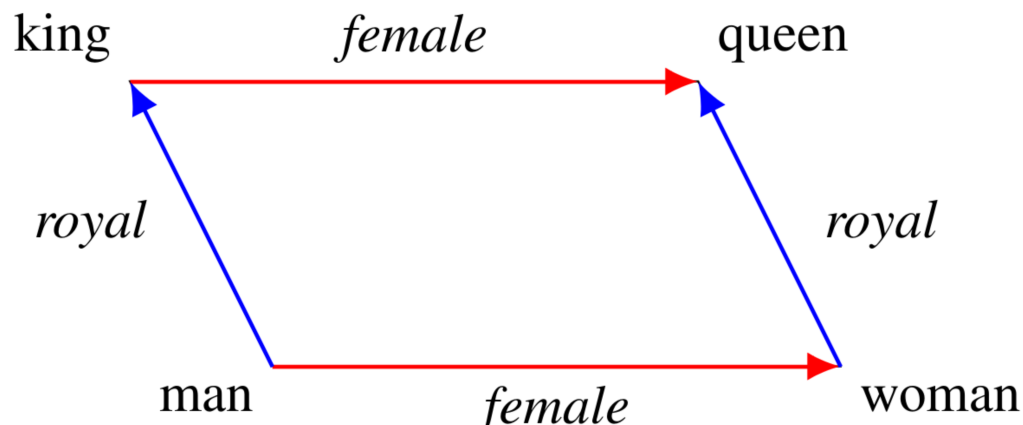
- real-word errors (e.g.: `their` instead of `there`)
- non-word errors (e.g.: `grammer`)

[Hagen et al. 2017, SIGIR]

# Embeddings

- Blog post: embedding-based correction [Rushton, 2018]

<https://edrushton.medium.com/a-simple-spell-checker-built-from-word-vectors-9f28452b6f26>

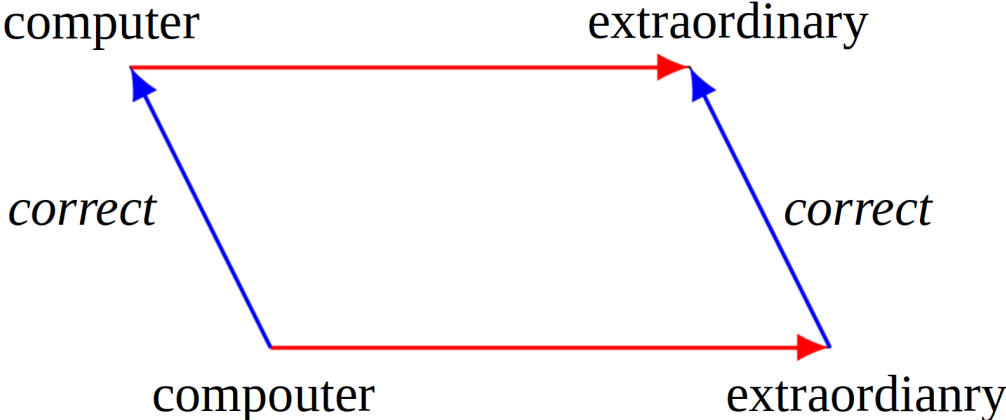


$$\vec{king} - \vec{man} + \vec{woman} \approx \vec{queen}$$

[Ethayarajh et al. 2019, ACL]

# Query Spelling Correction with Embeddings

- Misspelling property encoded in embeddings



$$\vec{computer} - \vec{compouter} + \vec{extraordianry} \approx \vec{extraordinary}$$

[Rushton, 2018]

# Query Spelling Correction with Embeddings

- ❑ 300-dim GloVe embeddings [nlp.stanford.edu/projects/glove/](http://nlp.stanford.edu/projects/glove/)
- ❑ Compute correction vector  $\vec{v}_{correct}$  (Oxford dict of common misspellings)
- ❑ Add  $\vec{v}_{correct}$  to each GloVe embedding
- ❑ Save nearest neighbor
- ❑ Look up correction for each query word

# Evaluation – Retrieval Effectiveness

None

Rushton

GPT-3.5

Hunspell

pyspellchecker



# Evaluation – Retrieval Effectiveness

---

	ANTIQUÉ	TREC DL '19	TREC DL '20
Spell Correction	BM25 PL2	BM25 PL2	BM25 PL2

---

None

Rushton

GPT-3.5

Hunspell

pyspellchecker

---

nDCG@10

# Evaluation – Retrieval Effectiveness

	ANTIQUUE		TREC DL '19		TREC DL '20	
Spell Correction	BM25	PL2	BM25	PL2	BM25	PL2
None	0.51	0.49	0.48	0.47	0.49	0.48
Rushton	0.51	0.49	0.48	0.47	0.49	0.48
GPT-3.5						
Hunspell						
pyspellchecker						

nDCG@10

# Evaluation – Retrieval Effectiveness

---

	ANTIQUUE		TREC DL '19		TREC DL '20	
Spell Correction	BM25 PL2		BM25	PL2	BM25	PL2
None	0.51	0.49	0.48	0.47	0.49	0.48
Rushton	0.51	0.49	0.48	0.47	0.49	0.48
GPT-3.5	0.49	0.47	0.48	0.47	0.48	0.47
Hunspell	0.48	0.46	0.36	0.34	0.42	0.41
pyspellchecker	0.43	0.41	0.31	0.30	0.33	0.32

---

nDCG@10

# Evaluation – Corrections

---

	ANTIQUUE				
	#err.	corrected			
Error Type		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

---

- R – Rushton
- H – Hunspell
- P – pyspellchecker
- G – GPT3.5

# Evaluation – Corrections Examples

Error Type	ANTIQUE				
	#err.	corrected			
		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
<b>Incorr. changes</b>		<b>0</b>	<b>52</b>	<b>75</b>	<b>42</b>

Original query:

why is gordon ramsey so popular

Rushton: —

Hunspell: —

pyspellchecker:

why is cordon raise so popular

GPT-3.5:

why is gordon ramsay so popular?

# Evaluation – Corrections Examples

Error Type	ANTIQUE				
	#err.	corrected			
		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

Original query:

why is gordon ramsey so popular

Rushton: —

Hunspell: —

pyspellchecker:

why is cordon raise so popular

GPT-3.5:

why is gordon ramsay so popular?

Limitation Rushton: GloVe vocabulary

# Evaluation – Corrections Examples

Error Type	ANTIQUE				
	#err.	corrected			
		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

Original query:

[...] copies of letters of  
 commendation **pr**sented [...]

Rushton, Hunspell, pypell., GPT-3.5:

[...] copies of letters of  
 commendation **pre**sented [...]

# Evaluation – Corrections Examples

Error Type	ANTIQUE				
	#err.	corrected			
		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

Original query:

how can i keep my **rabit** indoors

Rushton:

how can i keep my **rabbit** indoors

Hunspell:

how can i keep my **habit** indoors

pyspellchecker:

how can i keep my **barit** indoors

GPT-3.5:

how can i keep my **rabbit** indoors



# Evaluation – Corrections Examples

---

	ANTIQUE				
	#err.	corrected			
Error Type		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

---

Original query:

what is wifi vs bluetooth

Rushton, GPT-3.5: —

Hunspell, pyspellchecker:

what is wife vs bluetooth

# Evaluation – Corrections Examples

---

	ANTIQUE				
	#err.	corrected			
Error Type		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

---

Original query:

why hot air raise up

GPT-3.5:

why does hot air raise up?

Original query:

is cdg airport in main paris

GPT-3.5:

[...] corrected search query:  
is cdg airport in main paris?

# Evaluation – Corrections Examples

ANTIQUE					
Error Type	#err.	corrected			
		R	H	P	G
Deletion	14	4	6	5	12
Insertion	5	1	1	1	4
Space	13	1	1	2	11
Special char.	19	-	2	6	9
Substitution	8	1	2	2	8
Transposition	4	2	2	1	4
Incorr. changes		0	52	75	42

Original query:

how can we get concentration  
 on something

Rushton: —

Hunspell:

how can we get concentration  
 something

pyspellchecker:

how can we get concentration

GPT-3.5:

how can we get concentration  
 on something?

Limitation Rushton: Spaces, real-word errors

# Evaluation – Corrections Examples

Error Type	ANTIQUÉ					TREC DL '19					TREC DL '20				
	#err.	corrected				#err.	corrected				#err.	corrected			
		R	H	P	G		R	H	P	G		R	H	P	G
Deletion	14	4	6	5	12	-	-	-	-	-	-	-	-	-	-
Insertion	5	1	1	1	4	1	-	1	1	1	-	-	-	-	-
Space	13	1	1	2	11	-	-	-	-	-	-	-	-	-	-
Special char.	19	-	2	6	9	1	-	-	-	-	3	-	-	-	-
Substitution	8	1	2	2	8	1	-	1	1	1	1	-	-	-	-
Transposition	4	2	2	1	4	-	-	-	-	-	-	-	-	-	-
Incorr. changes		0	52	75	42		1	12	15	5		0	11	23	15

# Wrap-Up

- ❑ Considered spell correctors worsen retrieval results (on average)
- ❑ Embedding-based approach better than common spell correctors

# Wrap-Up

- ❑ Considered spell correctors worsen retrieval results (on average)
- ❑ Embedding-based approach better than common spell correctors

## Limitations

- ❑ GloVe vocabulary
- ❑ Single words
- ❑ No context

# Wrap-Up

- ❑ Considered spell correctors worsen retrieval results (on average)
- ❑ Embedding-based approach better than common spell correctors

## Limitations

- ❑ GloVe vocabulary
- ❑ Single words
- ❑ No context

## Future Work

- ❑ Other correction approaches
- ❑ Other embeddings
- ❑ Create misspelling query variants