

# GIS search queries correction

Vadim Fomin, NSU  
Ivan Bondarenko, MIPT

# Candidate search

String metric + dictionary lookup

Problems:

1. Too many unusual tokens:

хоум кредит банк → ухом кредит банк

2. How bad can a typo be?

ditroitigers → detroit tigers, log wood → dog food

# Candidate search

Iterative spelling correction as “unsupervised” candidate search:  
ditroitigers → ditroittigers → detroitigers → detroit tigers

Cucerzan S., Brill E. Spelling correction as an iterative process that exploits the collective knowledge of web users //Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing. – 2004.

# Candidate selection

Which suggestion is better? Why?

сабака → собака / кабака

зодорнов → задорнов / лингвистика

еда для кта → еда для кита / еда для кота

# Candidate selection

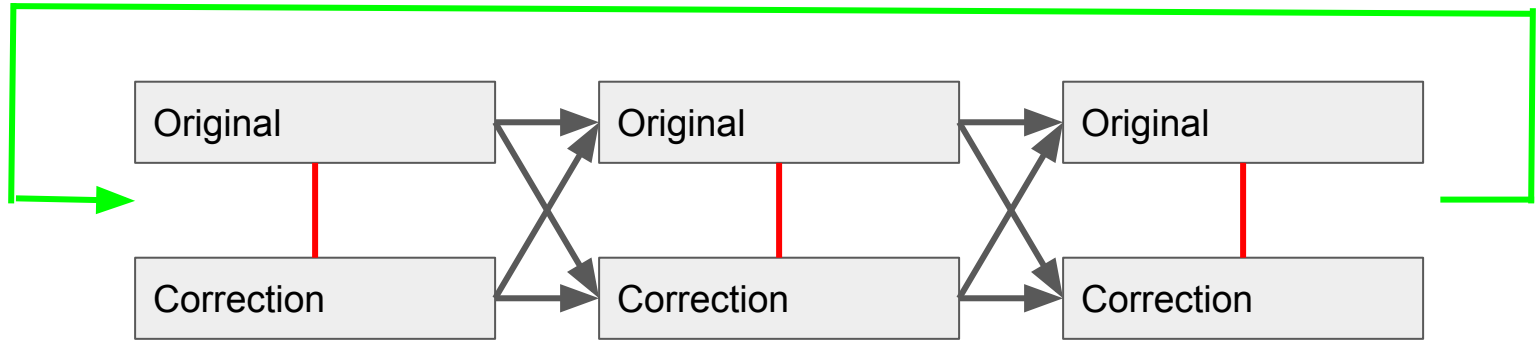
Kernighan M. D., Church K. W., and Gale W. A. (1990) A spelling correction program based on a noisy channel model. In Proceedings of the 13th conference on Computational linguistics, pp. 205–210. Association for Computational Linguistics.

# Candidate selection

Sorokin A. A., Shavrina T. O. Automatic spelling correction for Russian social media texts //Proceedings of the International Conference “Dialog”(Moscow. – 2016. – C. 688-701.

Reranking: models as features

# Overall architecture



# Training

Logistic regression

Winner – loser  $\rightarrow 1$

Loser – winner  $\rightarrow 0$



# Dataset

*"аверкеева,22", аверкеева 22, аверкиева 22,1 × 14 400 — supervised*  
*сорокина 24/б, сорокина 24 б × 949 000 — unsupervised*

1	Correction length	18	$a \rightarrow o, o \rightarrow a, e \rightarrow u, \text{ or } u \rightarrow e$
2	Simple Levenshtein distance	19	$ы \rightarrow u, ё \rightarrow o, ю \rightarrow y$ after ж, ч, ш, щ
3	Ngram-model score	20	$цы \rightarrow ци$ or $ци \rightarrow цы$
4	OOV in corrections	21	$ыва \rightarrow ова$
5	Vocabulary $\rightarrow$ OOV	22	$аро \rightarrow оро$ or $ало \rightarrow оло$
6	OOV $\rightarrow$ vocabulary	23	$э \rightarrow e$
7	Corrections that are more frequent than the originals	24	$ца \rightarrow це$
8	Simple Levenshtein distance, only OOV originals	25	$пре \rightarrow при$ or $при \rightarrow пре$
9	Simple Levenshtein distance, originals in vocabulary only	26	$э \rightarrow u$
10	1-operation corrections	27	$ё \rightarrow йо$ or $e \rightarrow йо$
11	Space deletions	28	unvoiced $\rightarrow$ voiced, or voiced $\rightarrow$ unvoiced
12	Space insertions	29	$зн \rightarrow здн, сн \rightarrow стн, сл \rightarrow стл, нст \rightarrow нтст, здн \rightarrow зн, стн \rightarrow сн, стл \rightarrow сл, \text{ or } нтст \rightarrow нст$
13	OOV that can be split into two vocabulary words	30	$хк \rightarrow гк$
14	Weighted keyboard layout Levenshtein distance	31	$н \rightarrow нн, с \rightarrow сс, м \rightarrow мм, ф \rightarrow фф$ , or vice versa
15	Weighted Levenshtein distance with insertion weight 10	32	$ь \rightarrow ъ$ or $ъ \rightarrow ь$
16	Weighted Levenshtein distance with deletion weight 10	33	insertion of ь as the fourth-to-last letter
17	Simple Levenshtein distance between phonetic codes	34	$тся \rightarrow ться$ or $ться \rightarrow тся$

# Features

семёновская 8 → семёновская <NUMBER>

# Features

Rejected features:

Word2vec as a probabilistic language model

Various settings of ngram-models

Ngram-model of morpho-tags

Ngram-model of prepositions and morpho-tags

Ngram-model of top-words and morpho-tags

Corrections that are more frequent than the originals etc.

# Feature importance



1	Correction length	18	$a \rightarrow o, o \rightarrow a, e \rightarrow u$ , or $u \rightarrow e$
2	Simple Levenshtein distance	19	$ы \rightarrow u, ё \rightarrow o, ю \rightarrow y$ after $ж, ч, ш, щ$
3	Ngram-model score	20	$цы \rightarrow ци$ or $ци \rightarrow цы$
4	OOV in corrections	21	$ыва \rightarrow ова$
5	Vocabulary $\rightarrow$ OOV	22	$аро \rightarrow оро$ or $ало \rightarrow оло$
6	OOV $\rightarrow$ vocabulary	23	$э \rightarrow e$
7	Corrections that are more frequent than the originals	24	$ца \rightarrow це$
8	Simple Levenshtein distance, only OOV originals	25	$пре \rightarrow при$ or $при \rightarrow пре$
9	Simple Levenshtein distance, originals in vocabulary only	26	$э \rightarrow u$
10	1-operation corrections	27	$ё \rightarrow йо$ or $e \rightarrow йо$
11	Space deletions	28	unvoiced $\rightarrow$ voiced, or voiced $\rightarrow$ unvoiced
12	Space insertions	29	$зн \rightarrow здн, сн \rightarrow стн, сл \rightarrow стл, нст \rightarrow нтст, здн \rightarrow зн, стн \rightarrow сн, стл \rightarrow сл$ , or $нтст \rightarrow нст$
13	OOV that can be split into two vocabulary words	30	$хк \rightarrow гк$
14	Weighted keyboard layout Levenshtein distance	31	$н \rightarrow нн, с \rightarrow сс, м \rightarrow мм, ф \rightarrow фф$ , or vice versa
15	Weighted Levenshtein distance with insertion weight 10	32	$ь \rightarrow ъ$ or $ъ \rightarrow ь$
16	Weighted Levenshtein distance with deletion weight 10	33	insertion of $ь$ as the fourth-to-last letter
17	Simple Levenshtein distance between phonetic codes	34	$тся \rightarrow ться$ or $ться \rightarrow тся$

Ngram-model score	Vocabulary → OOV
Weighted Levenshtein distance with deletion weight 10	Space insertions
Correction length	<i>ца → це</i>
Simple Levenshtein distance, originals in vocabulary only	OOV that can be split into two vocabulary words
OOV in corrections	<i>ы → и, ё → о, ю → у</i> after <i>ж, ч, ш, щ</i>
Simple Levenshtein distance, only OOV originals	insertion of <i>ь</i> as the fourth-to-last letter
Simple Levenshtein distance	<i>э → е</i>
<i>ь → ъ</i> or <i>ъ → ь</i>	unvoiced → voiced, or voiced → unvoiced
Weighted keyboard layout Levenshtein distance	<i>пре → при</i> or <i>при → пре</i>
Simple Levenshtein distance between phonetic codes	<i>цы → ци</i> or <i>ци → цы</i>
Space deletions	<i>аро → оро</i> or <i>ало → оло</i>
OOV → vocabulary	<i>зн → здн, сн → стн, сл → стл, нст → нтст, здн → зн, стн → сн, стл → сл, or нтст → нст</i>
1-operation corrections	<i>ыва → ова</i>
Corrections that are more frequent than the originals	<i>э → и</i>
<i>а → о, о → а, е → и, or и → е</i>	<i>ё → йо</i> or <i>е → йо</i>
Weighted Levenshtein distance with insertion weight 10	<i>хк → гк</i>
<i>н → нн, с → сс, м → мм, ф → фф</i> , or vice versa	<i>тся → ться</i> or <i>ться → тся</i>

# Feature importance

Ngram-model is incredibly helpful in picking a good candidate

Weighted string metrics are good at rejecting the bad ones

(Weighted Levenshtein distance with a high deletion weight is the best)

Morphological models and misspelling type features are not

(but Ъ → ъ!)



# Evaluation

original:	мастер фиш	на улленина
suggestion:	мастер фиш	над ул ленина
expected:	мастерфиш	на ул ленина
	FN	FP TP

# Evaluation

	Precision	Recall	F1
Model in question	98.4%	58.6%	73.5%
Hunspell	19.6%	22.2%	20.8%
2GIS algorithm	54.6%	31.5%	40.0%

# Evaluation: HunSpell

An extended dictionary

First candidate, if any

Result: **20.8%**

# Evaluation: 2GIS algorithm

No spelling correction unit

City-dependent data

Top-1 candidate query = output

# Evaluation: SpellRuEval

SpellRuEval	82.0%	69.3%	75.0%
-------------	-------	-------	-------



Model in question	98.4%	58.6%	73.5%
-------------------	-------	-------	-------

# Results

- An annotated corpus
- Effectiveness of ngram-models and reranking tested on a quite different corpus
- Feature engineering
- Feature importance analyzed
- An open-source system

# Acknowledgements

We are grateful to Botan Investments for supporting students' interest to machine learning.