

Рыбакин Иван Владиславович, студент,
Поволжский государственный университет
телекоммуникаций и информатики
Rybakin Ivan Vladislavovich, student,
Volga Region State University of
Telecommunications and Informatics

Попов Виктор Борисович, профессор,
Поволжский государственный университет
телекоммуникаций и информатики
Popov Viktor Borisovich, professor
Volga Region State University of
Telecommunications and Informatics

КОНТЕКСТНО-ЗАВИСИМЫЕ И КОНТЕКСТНО-НЕЗАВИСИМЫЕ МЕТОДЫ КОРРЕКЦИИ ОПЕЧАТОК: СРАВНИТЕЛЬНЫЙ АНАЛИЗ ЭФФЕКТИВНОСТИ CONTEXT-DEPENDENT AND CONTEXT-INDEPENDENT METHODS OF TYPO CORRECTION: A COMPARATIVE ANALYSIS OF EFFECTIVENESS

Аннотация. В статье проводится сравнительный анализ контекстно-независимых (словарные методы, редакционное расстояние) и контекстно-зависимых (языковые модели, трансформеры) методов коррекции опечаток. Рассматриваются ключевые алгоритмы каждого класса, сравниваются их точность, скорость и область применимости. Выделены условия, при которых предпочтителен каждый из подходов.

Abstract. This paper presents a comparative analysis of context-independent (dictionary methods, edit distance) and context-dependent (language models, transformers) typo correction methods. Key algorithms of each class are examined, their accuracy, speed and applicability are compared. Conditions under which each approach is preferable are identified.

Ключевые слова: Коррекция опечаток, контекстно-зависимые методы, редакционное расстояние, BERT, spell checking, NLP.

Keywords: Typo correction, context-dependent methods, edit distance, BERT, spell checking, NLP.

1. Введение

Автоматическая коррекция опечаток является фундаментальной задачей обработки естественного языка. По различным оценкам, от 1 до 3% токенов в пользовательских запросах содержат ошибки ввода [1], что негативно влияет на качество машинного перевода, анализа тональности и других NLP-приложений.

Методы коррекции делятся на два класса по признаку использования контекста. Контекстно-независимые оценивают слово изолированно; контекстно-зависимые учитывают окружение и способны исправлять real-word ошибки – случаи, когда написанное слово словарно корректно, но не соответствует смыслу. Целью данной работы является сравнительный анализ этих двух классов по ключевым показателям эффективности.

2. Характеристика методов

2.1. Контекстно-независимые методы

Базовым инструментом служат метрики редакционного расстояния. Расстояние Левенштейна определяет минимальное число операций вставки, удаления и замены символов



[2]; метрика Дамерау–Левенштейна дополнительно учитывает транспозиции. Инструмент Hunspell реализует этот подход с поддержкой морфологических правил и обеспечивает детерминированный результат при задержке порядка единиц миллисекунд. Принципиальное ограничение – неспособность обнаруживать real-word ошибки.

2.2. Контекстно-зависимые методы

Байесовская модель зашумлённого канала [3] формализует опечатку как прохождение слова через канал с шумом; языковая модель реализуется n-граммной статистикой. Нейросетевые seq2seq-модели на базе LSTM трактуют коррекцию как sequence-to-sequence преобразование [4]. Трансформерные архитектуры – BERT и T5 – используют механизм self-attention для учёта полного контекста предложения, обеспечивая наивысшее качество обработки обоих классов ошибок [5, 6].

3. Сравнительный анализ

Контекстно-независимые методы полностью неэффективны для real-word ошибок, тогда как трансформерные модели достигают F1-меры свыше 92% на этом классе. На non-word ошибках разрыв менее значителен: Hunspell демонстрирует конкурентные результаты при задержке 12 мс против 380-520 мс для BERT и T5 соответственно. Таким образом, прирост качества сопровождается нелинейным ростом вычислительных затрат.

Контекстно-независимые методы предпочтительны при жёстких ресурсных ограничениях, высококонтролируемых доменах и отсутствии обучающих данных. Контекстно-зависимые оправданы при наличии значимой доли real-word ошибок, обработке свободного текста и достаточности вычислительных ресурсов. Перспективным решением является гибридная каскадная архитектура, совмещающая быструю контекстно-независимую предфильтрацию с контекстно-зависимым анализом сложных случаев.

4. Заключение

В работе проведён сравнительный анализ контекстно-зависимых и контекстно-независимых методов коррекции опечаток. Показано, что принципиальное различие определяется типом обрабатываемых ошибок: контекстно-независимые методы ограничены non-word классом, тогда как трансформерные архитектуры обеспечивают полноценную обработку обоих классов при значительно более высокой вычислительной стоимости. Выбор подхода должен определяться требованиями целевого приложения. Результаты формируют теоретическую базу для разработки системы коррекции опечаток в рамках магистерской диссертации.

Список литературы:

1. Kukich K. Technique for automatically correcting words in text // ACM Computing Surveys. – 1992. – Vol. 24, No. 4. – P. 377–439.
2. Levenshtein V.I. Binary codes capable of correcting deletions, insertions, and reversals // Soviet Physics Doklady. – 1966. – Vol. 10, No. 8. – P. 707–710.
3. Church K., Gale W. Probability scoring for spelling correction // Statistics and Computing. – 1991. – Vol. 1, No. 2. – P. 93–103.
4. Yuan Z., Briscoe T. Grammatical error correction using neural machine translation // Proc. NAACL-HLT. – San Diego, 2016. – P. 380–386.
5. Devlin J. et al. BERT: Pre-training of deep bidirectional transformers for language understanding // Proc. NAACL-HLT. – Minneapolis, 2019. – P. 4171–4186.
6. Omelianchuk K. et al. GECToR – grammatical error correction: tag, not rewrite // Proc. Workshop BEA. – 2020. 0Ц P. 163–170.

