

ПРАВИТЕЛЬСТВО РОССИЙСКОЙ ФЕДЕРАЦИИ
ФГАОУ ВО НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»

Факультет компьютерных наук
Образовательная программа «Прикладная математика и информатика»

УДК XXXXX

Отчет об исследовательском проекте на тему:
Активное обучение в задаче вопросно-ответных систем

Выполнила:

студентка группы ММОВС21_О_1
Ведерникова Мария Алексеевна

(подпись)

(дата)

Принял руководитель проекта:

Цвигун Аким Олегович
Преподаватель факультета компьютерных наук НИУ ВШЭ

(подпись)

(дата)

Москва 2023

Содержание

Аннотация	3
1 Введение	4
2 Обзор существующих решений	5
3 Постановка задачи	6
3.1 Активное обучение	6
3.2 Стратегии активного обучения	7
4 Вопросо-ответная система	8
4.1 Описание датасета	8
4.2 Двухэтапная схема генерации ответа	9
4.3 Обучение моделей для классификации и генерации	10
4.4 Валидационные метрики	12
5 Активное обучение	12
5.1 Схема активного обучения	12
5.2 Эксперименты	13
6 Результаты	14
6.1 Серия экспериментов №1	14
6.2 Серия экспериментов №2	15
6.3 Распределение скоров кандидатов во время отбора	16
7 Выводы	18
8 Заключение	19
Список литературы	20

Аннотация

Разметка данных для создания вопросо-ответных система требует значительных затрат времени и ресурсов, поскольку каждый экземпляр должен быть вручную размечен человеком. Эта сложность особенно возрастет в случае длинных текстов, из которых необходимо извлекать ответ, и свободной формы ответа на естественном языке, ведь в этом случае как содержание ответов, так и их формулировки могут варьироваться. Активное обучение – хорошо известный метод, который может сократить трудозатраты по извлечению информации за счет выбора наиболее перспективных экземпляров для разметки на основании исследования уже размеченного датасета. Однако, до сих пор не было эффективных стратегий использования активного обучения для задач вопросо-ответных систем. Это связано с тем, что многие стратегии активного обучения полагаются на оценку неопределенности, которая может ухудшить производительность модели за счет обучения на экземплярах из областей низкой плотности распределения. В данной работе производится сравнение разных стратегий активного обучения, в том числе основанных на принципе разнообразия в противовес неопределенности, для задачи генеративного построения ответа на естественном языке в условиях очень длинного контекста, превышающего максимальную длину входной последовательности для многих моделей глубокого обучения. Показано, что рассмотренные стратегии активного обучения не дают прироста качества в сравнении со случайным отбором документов для разметки.

Ключевые слова

Глубокое обучение, обработка естественного языка, активное обучение, вопросо-ответные системы, генерация текста

1 Введение

Вопросо-ответные системы предполагают сопоставление паре (контекст, вопрос) ответа, который создается с учетом информации, полученной из контекста и содержит верный ответ на вопрос. При этом контекст, вопрос и ответ являются текстами на естественном языке. Существует два основных подхода к решению данной задачи: экстрактивный, при котором ответ является последовательной частью контекста [2], и генеративный, при котором текст ответа не ограничен контекстом и создается, как правило, с помощью нейронных моделей класса seq-to-seq, способных генерировать текст [29].

Составление аннотации к тексту с ответом на вопрос вручную – сложная задача, которая требует внимательного прочтения длинного исходного документа, выделения важной информации и, наконец, написания небольшого текста. Каждый из этих этапов может занять много времени, в результате чего добавление каждого экземпляра в размеченный датасет для дальнейшего обучения на нем вопросо-ответной системы обходится очень дорого.

Активное обучение [16] – хорошо известная техника, которая помогает существенно сократить объем работы по разметке данных, необходимой для достижения определенного уровня производительности модели машинного обучения. Например, в задаче маркировки последовательности исследователи сообщают о возможности достижения 99% от качества, получаемого при обучении на полном датасете, при использовании 16% данных с помощью активного обучения [6]. Это делает активное обучение особенно важным, когда разметка является дорогостоящей, что характерно для вопросо-ответных задач.

Активное обучение работает итеративно: на каждой итерации (1) модель обучается на размеченном на текущий момент наборе данных; (2) модель используется для выбора нескольких наиболее информативных экземпляров из большого пула неразмеченных данных; (3) информативные экземпляры предоставляются экспертам на разметку; (4) наконец, полученные на предыдущем шаге экземпляры с разметкой добавляются в размеченный набор данных, и начинается новая итерация.

Стратегии активного обучения могут быть основаны на методах оценки неопределенности [25], [23]. Такие методы предполагают, что экземпляры, на которых предсказания модели, обученной на текущей итерации, наиболее неопределенны, являются самыми информативными для обучения модели на следующей итерации. Однако, такие методы работают не всегда и в некоторых случаях не более полезны, чем добавление случайно выбранных экземпляров [3]. Более того, в некоторых случаях они вносят больше шума и ухудшают производительность моделей обобщения за счет смещения распределения [15], [34] или предпо-

читательного выбора выбросов [38], [27]. Поэтому подходы, основанные на неопределенности, могут адаптироваться под задачу с помощью учета плотности распределения экземпляров или добавления в разметку определенной части случайно выбранных документов. Еще одним вариантом является выбор документов для разметки, полностью основанный на распределении экземпляров.

В данной работе рассмотрены различные стратегии активного обучения в задаче вопросо-ответных систем с применением генеративного подхода. Код для воспроизведения экспериментов доступен онлайн¹.

2 Обзор существующих решений

Генеративные вопросо-ответные системы. Первые модели вида seq2seq [35] появились в 2014 году. Механизм внимания [30] и архитектура transformer [31], основанная на нем, легла в основу многих языковых моделей, ставших базовыми для таких задач, как машинный перевод и суммаризации текстов. Генеративные вопросо-ответные системы рассматривались в основном для датасетов, требующих генерации ответов, таких как NarrativeQA [37], CoQA [14] или ELI5 [18]. Эти датасеты были сгенерированы таким образом, что ответ не является частью контекста, что делает затруднительным применение экстрактивных моделей. В работе [20] показано, что генеративные модели конкурентоспособны и в задачах понимания прочитанного, таких как SQuAD [36], где ответы представляют собой часть контекста. Идея сочетания предварительного извлечения части контекста, содержащей ответ, с последующим применением генеративной модели для формирования ответа, была рассмотрена в работах [10], [32], [24].

Активное обучение при генерации естественного языка. В то время как многие работы используют активное обучение для задач классификации текстов или разметки последовательностей [13], [12], [7], задачам генерации естественного языка уделяется не так много внимания. Среди работ в этой области стоит отметить [8], [9], [11]. Эти работы посвящены нейросетевому машинному переводу текста и предлагают несколько возможных стратегий запросов на основе неопределенности для активного обучения. В работе [4] успешно применяется активное обучение в интерактивном машинном переводе. В работе [3] рассмотрены различные стратегии активного обучения в задаче генеративной суммаризации текста, в том числе предложен новый метод выбора экземпляров для активного обучения, основанный на

¹Репозиторий с кодом для воспроизведения экспериментов: <https://github.com/masha239/ActiveLearningInQuestionAnswering>

внутридоменном разнообразии. Сразу несколько задач, в которых требуется генерация текста, например, генерация ответа, перефразирование и суммаризация, и несколько различных стратегий активного обучения рассмотрены в статье [5], авторы которой не находят значимых отличий в результативности этих стратегий от случайного выбора экземпляров.

3 Постановка задачи

3.1 Активное обучение

В этом разделе представлено формальное определение процедуры активного обучения для генерации текста. Здесь и во всей остальной части работы входная последовательность обозначается как $\mathbf{x} = (x_1, \dots, x_m)$, и выходную последовательность как $\mathbf{y} = (y_1, \dots, y_n)$, где m и n – длины \mathbf{x} и \mathbf{y} соответственно. При этом x_i составлена из текста вопроса и контекста, объединенных в одну последовательность (такую последовательность в дальнейшем мы будем называть вопросом с контекстом).

Пусть $\mathcal{D} = (\mathbf{x}^{(k)}, \mathbf{y}^{(k)})_{k=1}^K$ – набор данных пар (вопрос с контекстом, ответ). Рассмотрим вероятностную модель $p_{\mathbf{w}}(\mathbf{y}|\mathbf{x})$, параметризованную вектором \mathbf{w} . Обычно $p_{\mathbf{w}}(\mathbf{y}|\mathbf{x})$ представляет собой нейронную сеть, а оценка параметров производится с помощью метода максимального правдоподобия:

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}), \quad (1)$$

где $L(\mathcal{D}, \mathbf{w}) = \sum_{k=1}^K \log p_{\mathbf{w}}(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$ – логарифмическое правдоподобие.

Многие методы активного обучения основаны на стратегии жадных запросов, которые выбирают экземпляры для аннотации, оптимизируя определенный критерий $\mathcal{A}(\mathbf{x} | \mathcal{D}, \hat{\mathbf{w}})$, называемый функцией приобретения:

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{x}} \mathcal{A}(\mathbf{x} | \mathcal{D}, \hat{\mathbf{w}}) \quad (2)$$

Выбранный экземпляр \mathbf{w}^* затем аннотируется целевым значением \mathbf{y} (ответ на вопрос) и добавляется в обучающий набор данных: $\mathcal{D} := \mathcal{D} \cup (\mathbf{x}^*, \mathbf{y}^*)$. Далее параметры модели \mathbf{w} обновляются и процесс выбора новых экземпляров продолжается до тех пор, пока не будет достигнуто желаемое качество модели или не будет получено максимально возможное количество аннотаций в рамках решения задачи.

Выбор функции приобретения имеет решающее значение для успеха активного обуче-

ния. Распространенным подходом для приобретения является выбор экземпляров с высокой неопределенностью модели. Ниже будут рассмотрены несколько мер неопределенности, используемых при решении задачи генерации ответа.

3.2 Стратегии активного обучения

Нормализованная вероятность последовательности впервые была предложена в 2007 году [40] и широко использовалась в дальнейшем [22], [39], [41]. Нормализованная вероятность последовательности (Normalized Sequence Probability, NSP) обобщает понятие вероятности на сгенерированную последовательность, нормализуя ее по длине:

$$\text{NSP}(\mathbf{x}) = 1 - \bar{p}_{\hat{\mathbf{w}}}(\mathbf{y} | \mathbf{x}), \quad (3)$$

где геометрическое среднее вероятностей токенов, предсказанных моделью, определяется как $\bar{p}_{\hat{\mathbf{w}}}(\mathbf{y} | \mathbf{x}) = \exp\{\frac{1}{n} \log p_{\mathbf{w}}(\mathbf{x} | \mathbf{y})\}$.

Внутридоменное разнообразие. Стратегия отбора кандидатов по критерию внутридоменного разнообразия (In-Domain Diversity Sampling, IDDS) [3] не использует неопределенность в предсказаниях модели. Авторы статьи рассмотрели различные стратегии, основанные на неопределенности, и пришли к выводу, что в задаче суммаризации текста данные стратегии не дают прироста в качества в сравнении со случайным выбором кандидатов; аналогичные выводы для ряда стратегий делаются и в [5].

Стратегия внутридоменного разнообразия заключается в выборе кандидатов, максимально репрезентативных по отношению к объектам неразмеченного пула и при этом "непохожих" на объекты, присутствующие в уже размеченной выборке.

$$\text{IDDS}(\mathbf{x}) = \lambda \frac{\sum_{j=1}^{|U|} s(\mathbf{x}, \mathbf{x}_j)}{|U|} - (1 - \lambda) \frac{\sum_{j=1}^{|L|} s(\mathbf{x}, \mathbf{x}_j)}{|L|}, \quad (4)$$

где $s(\mathbf{x}, \mathbf{x}_j)$ – функция близости между двумя текстами, U – множество неразмеченных объектов, L – множество размеченных объектов, λ – гиперпараметр.

Близость пары объектов s определяется на основе скалярного произведения нейросетевых эмбедингов входных последовательностей. В данной работе использовалось два подхода к вычислению близости: в первом эмбединги строились на основании исключительно текста вопроса, во втором – вопроса с контекстом. В обоих случаях в качестве модели для эмбединга была взята предобученная модель distilbert-base-uncased [17], эмбединг получался путем извлечения последнего скрытого состояния, соответствующего токenu [CLS]. В

случае вычисления эмбединга по вопросу с контекстом также было рассмотрено использование дообученной на тренировочном датасете модели, использующейся в качестве бинарного классификатора (подробнее об этом см. в разделе 4.2).

Минимальная уверенность. Данная стратегия [1] не относится к модели, которая генерирует текст, и представляет собой классическую стратегию минимальной уверенности для классификатора.

$$\text{IDDS}(\mathbf{x}) = 1 - p_{\mathbf{w}}(\mathbf{y} \mid \mathbf{x}), \quad (5)$$

В настоящей работе применяется двухэтапная схема генерации ответа, и на первом этапе применяется бинарный классификатор текста, который оценивает возможность наличия ответа в выделенной части контекста. От качества предсказаний бинарного классификатора напрямую зависит качество генерации ответа по всему контексту, кроме того, предположительно, неопределенность на первом этапе предсказания может быть связана и со сложностью выделения ответа в данном тексте. По эти соображениям данная стратегия также была включена в рассмотрение.

4 Вопросо-ответная система

В данном разделе описана схема решения задачи в применении к конкретному датасету, а также метрики, по которым оценивалось качество генерации ответа.

4.1 Описание датасета

Датасет Google Natural Questions [28] является одним из популярных бенчмарков для валидации моделей в задаче построения ответа на вопрос на английском языке. Вопросы, которые используются в датасете, заданы реальными пользователями поисковой системы, а ответы были найдены на сайте-энциклопедии www.wikipedia.org. Таким образом, контекстом в данном случае является соответствующая страница википедии. В основной версии датасета контекст сохранен как html-страница, включая специальные символы разметки. Однако, авторами также опубликована и версия без html-разметки, которая представляет текст на естественном языке и лучше подходит для дообучения языковых моделей, которые изначально также обучались на естественном языке.

Ответы на вопросы вручную выделены ассессорами из контекста. В датасете представлены три вида ответов: бинарные (да/нет), короткие (дата какого-либо события, имя

человека и т.д.; как правило, длина короткого ответа не превосходит 10 слов) и длинные (развернутый ответ, который, как правило, состоит из одного или нескольких предложений). Короткие и длинные ответы всегда являются частью контекста и заданы с помощью указания номеров первого и последнего символов. В то время как одному документу соответствует либо один длинный ответ (либо длинные ответы для данного документа отсутствуют), коротких ответов может быть больше одного.

В данной работе рассматривались только задачи генерации коротких ответов, и, соответственно, подмножество исходного датасета, содержащее хотя бы один короткий ответ. Кроме того, в датасет для исследования вошли только те экземпляры, длина контекста которых не превышает 20000 символов.

После описанной выше фильтрации тренировочный датасет насчитывает 40783 документов документов, тестовый – 500 документов, валидационный – 1693 документа.

4.2 Двухэтапная схема генерации ответа

Особенность датасета Google Natural Questions заключается в очень длинном контексте, который невозможно подать на вход большинству стандартных генеративных моделей в виду ограничения на длину входной последовательности. Существует два распространенных подхода к решению данной задачи: использование моделей, механизм внимания которых позволяет обработку длинной входной последовательности, например, longformer [26], и предварительный поиск наиболее релевантной вопросу части контекста [21]. В настоящей работе применен второй подход.

Как показано на см. Рисунке 4.1, для получения ответа контекст предварительно разбивается на части, каждая из которых, кроме, возможно, последней, содержит 10 предложений. Соседние части перекрываются между собой, начало каждой следующей части отстоит на 5 предложений от предыдущего. Входная последовательность модели составляется из строки "question: текста вопроса, строки "context: и, наконец, части контекста. Каждая из полученных входных последовательностей оценивается бинарным классификатором, в качестве которого выступает модель distilbert-base-uncased [17]. Последовательность с максимальным скором подается на вход генеративной модели t5-small [19], которая, в свою очередь, генерирует выходную последовательность, которая будет рассматриваться как ответ на вопрос.

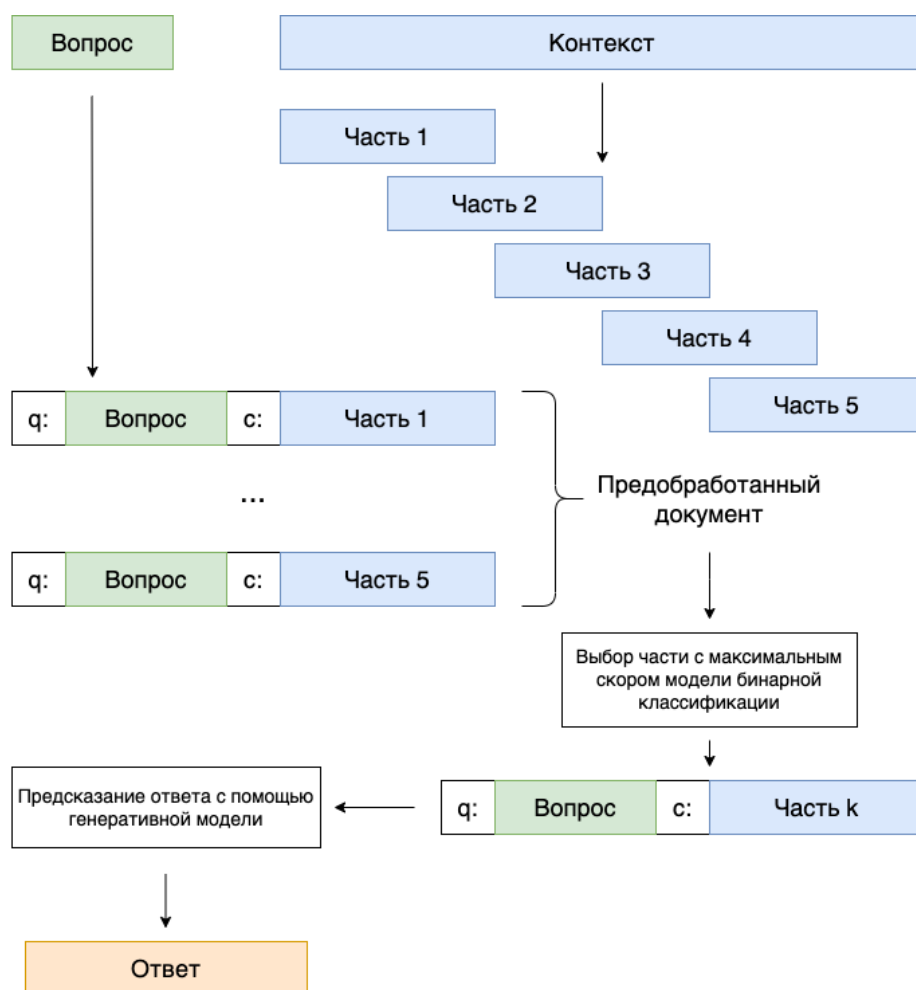


Рис. 4.1: Схема генерации ответа

4.3 Обучение моделей для классификации и генерации

Предобученные модели для классификации² и для генерации³ тренируются на одном и том же датасете отдельно.

Обучение модели классификации. Каждой части контекста сопоставляется метка **1**, если хотя бы один из ответов попадает в данную часть, и метка **0** в противном случае. На этапе обучения мы продолжаем рассматривать в качестве экземпляров датасета полные документы, а не их части, чтобы избежать дисбаланса встречаемости документов в зависимости от их длины (и, следовательно, количества частей). Поэтому во время каждой эпохи обучения часть контекста, которая будет передана модели, выбирается случайным образом.

Поскольку части без ответа встречаются приблизительно в 8 раз чаще, и при этом именно точное выделение частей контекста с ответом наиболее важно для поставленной задачи, вероятности выбора частей с меткой **1** и частей с меткой **0** подобраны следующим

²Модель `distillbert-base-uncased`: <https://huggingface.co/distilbert-base-uncased>, дата обращения 01.06.2023

³Модель `t5-small`: <https://huggingface.co/t5-small>, дата обращения 01.06.2023

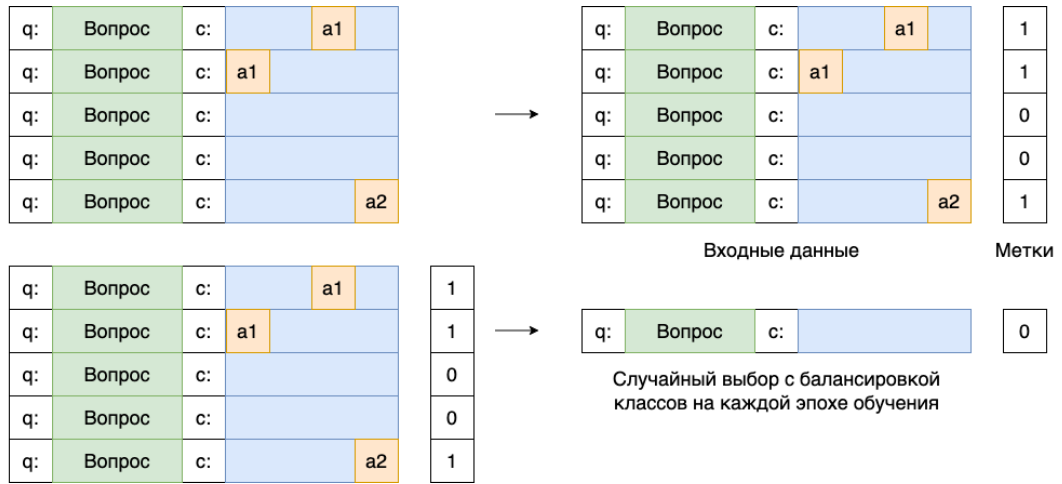


Рис. 4.2: Выбор части контекста для классификации на каждой эпохе обучения

образом: $p_0 = \frac{0.5}{n}$, $p_1 = \frac{0.5}{m}$, где p_0 – вероятность выбрать часть с меткой **0**, p_1 – вероятность выбрать часть с меткой **1**, n и m – количество частей в документе с метками **0** и **1** соответственно. Таким образом, две части одного документа с одинаковыми метками выбираются равновероятно и вероятность выбрать часть с меткой **1** равна 0.5.

Обучение генеративной модели. В обучении модели, которая применяется на втором этапе, участвуют только те части контекста, которые содержат хотя бы один ответ (и, соответственно, имеют метку **1** в задаче классификации). Поскольку контекст может содержать несколько различных ответов, в случае, если в части контекста есть хотя бы два не совпадающих ответа, для обучения равновероятно выбирается любой из них.

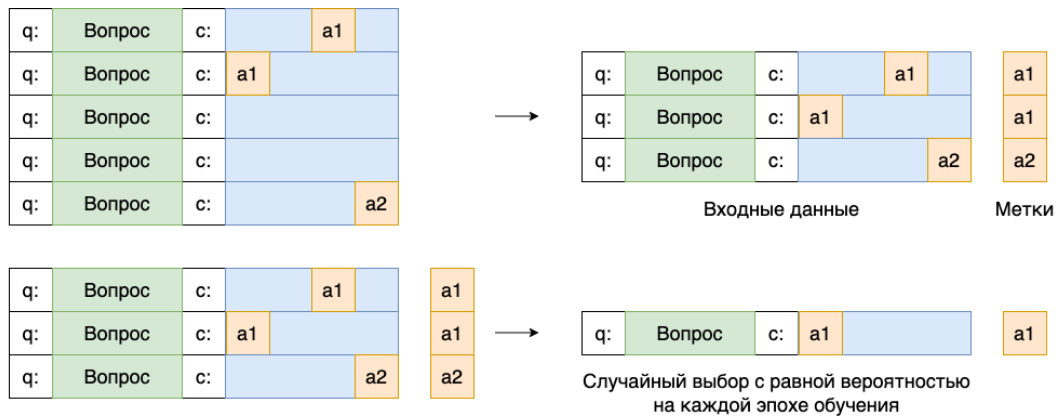


Рис. 4.3: Выбор части контекста для генерации ответа на каждой эпохе обучения

На каждой эпохе обучения выбор входной последовательности происходит равновероятно среди всех частей с ответом, которые содержатся в данном экземпляре.

4.4 Валидационные метрики

Для оценки качества системы вычисляются широко используемые метрики ROUGE [33], а именно ROUGE-1 и ROUGE-2, а также метрика точного совпадения предсказания и ответа Exact Match (в случае нескольких возможных ответов она будет равна 1, если предсказание совпало хотя бы с одним из них).

На этапе валидации метрики вычисляются для обеих подзадач, причем, в отличие от этапа обучения, используются все части каждого экземпляра валидационного датасета в случае подзадачи классификации и все части, содержащие ответ, в случае задачи генерации ответа. Для задачи классификации вычисляется ROC AUC – метрика, отражающая способность модели правильно ранжировать экземпляры, принадлежащие к различным классам. Для задачи генерации вычисляются те же метрики, что и для общей задачи генерации ответа по вопросу и контексту. Их можно интерпретировать как качество вопросо-ответной системы на задаче с "коротким" контекстом, которая не требует двухэтапного подхода.

5 Активное обучение

5.1 Схема активного обучения

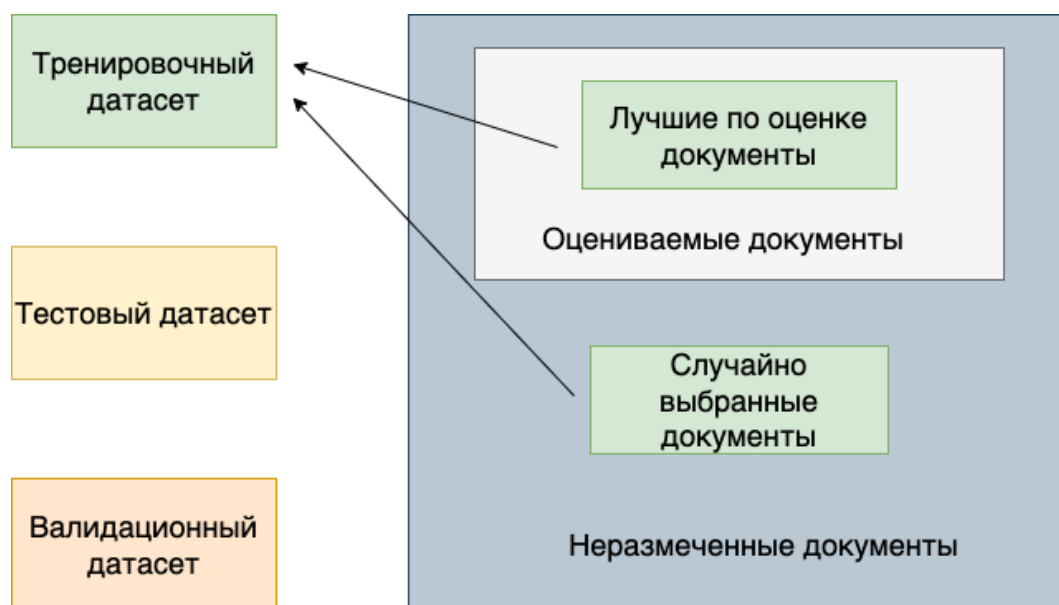


Рис. 5.1: Схема активного обучения

Схема активного обучения определяется стратегией отбора кандидатов, количеством документов, размеченных на первом этапе N , количеством документов, размечаемых на каждом этапе отбора M , количеством документов, среди которых проводится отбор M_POOL

(случайное сэмплирование из пула неразмеченных документов перед отбором требуется для сокращения времени вычислений) и долей случайно выбранных документов α . При $\alpha = 1$ отбор становится случайным.

В качестве базовых рассматриваются модели, обученные на случайно отобранных документах. С ними сравнивается качество моделей, обученных на документах, отобранных при помощи различных стратегий активного обучения.

При отборе документов стратегии минимальной нормализованной вероятности последовательности и внутридоменного разнообразия применяются к той части документа-кандидата, которая получает максимальный скор модели-классификатора, обученной на предыдущем этапе активного обучения.

5.2 Эксперименты

В рамках данной работы были проведены две серии экспериментов.

Серия экспериментов №1 с $N = 500$, $M = 500$ и $M_POOL = 5000$. В данной серии рассматривались стратегии отбора на основании:

- минимальной уверенности
- минимальной нормализованной вероятности последовательности
- внутридоменного разнообразия с $\lambda = 0.67$, расчетом близости между частями контекста с вопросами и использованием предобученной модели классификатора (без дообучения на текущую задачу).

Серия экспериментов №2 с $N = 50$, $M = 50$ и $M_POOL = 1000$. В данной серии рассматривались различные стратегии отбора на основании внутридоменного разнообразия с $\lambda = 0.67$:

- с расчетом близости между частями контекста с вопросами и использованием предобученной модели классификатора (без дообучения на текущую задачу)
- с расчетом близости между вопросами и использованием предобученной модели классификатора (без дообучения на текущую задачу),
- с расчетом близости между частями контекста с вопросами и использованием дообученной на текущую задачу модели классификатора.

6 Результаты

6.1 Серия экспериментов №1

На графиках ключевых метрик в решении полной задачи 6.1, а также метрик двух подзадач вопросо-ответной системы 6.2, 6.3 отражены следующие стратегии: (1) обучение на всем датасете, (2) случайный отбор, (3) стратегия минимальной уверенности, (4) стратегия минимальной нормализованной вероятности последовательности, (5) стратегия внутримоделного разнообразия с $\lambda = 0.67$.

Можно видеть, что никакая из рассмотренных стратегий активного обучения не выигрывает в качестве у случайного отбора ни в решении полной задачи, ни в обеих ее подзадачах.

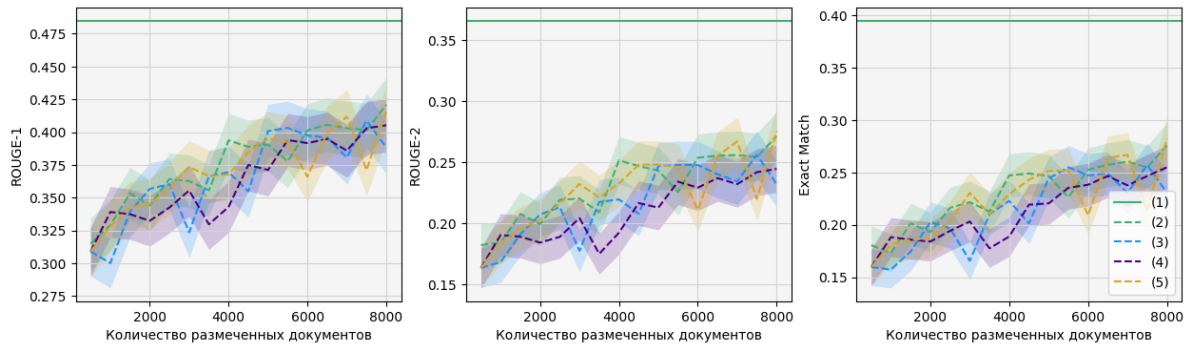


Рис. 6.1: Изменение метрик ROUGE-1, ROUGE-2 и Exact Match на валидационном датасете в задаче генерации ответа на вопрос по полному документу во время активного обучения.

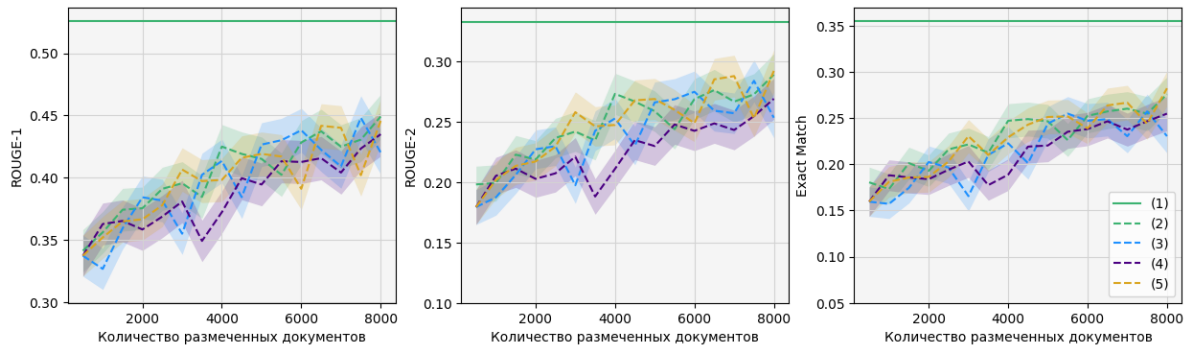


Рис. 6.2: Изменение метрик ROUGE-1, ROUGE-2 и Exact Match на валидационном датасете в задаче генерации ответа по заданной части контекста во время активного обучения.

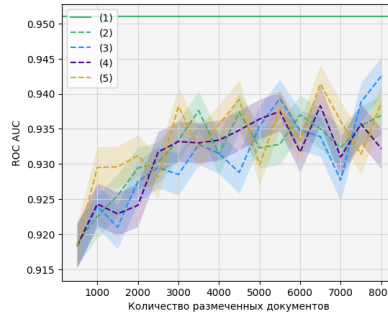


Рис. 6.3: Изменение метрики ROC AUC на валидационном датасете в задаче предсказания наличия ответа в заданной части контекста во время активного обучения.

6.2 Серия экспериментов №2

В этой серии шаг активного обучения выбран в 10 раз меньше, а стратегии активного являются вариациями стратегии внутридоменного разнообразия с $\lambda = 0.67$: (1) случайный отбор, (2) расчет близости между частями контекста с вопросами и использованием дообученной на текущую задачу модели классификатора, (3) расчет близости между частями контекста с вопросами и использованием предобученной модели классификатора (без дообучения на текущую задачу), (4) расчет близости между частями контекста с вопросами и использованием дообученной на текущую задачу модели классификатора.

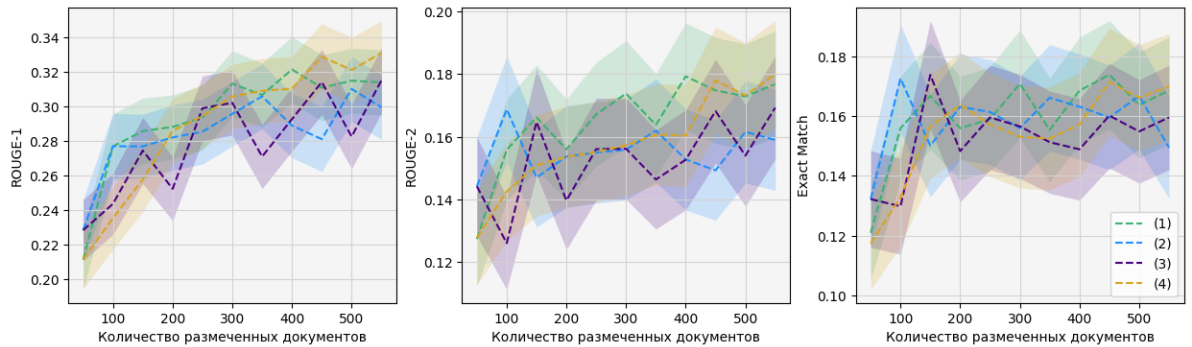


Рис. 6.4: Изменение метрик ROUGE-1, ROUGE-2 и Exact Match на валидационном датасете в задаче генерации ответа на вопрос по полному документу во время активного обучения.

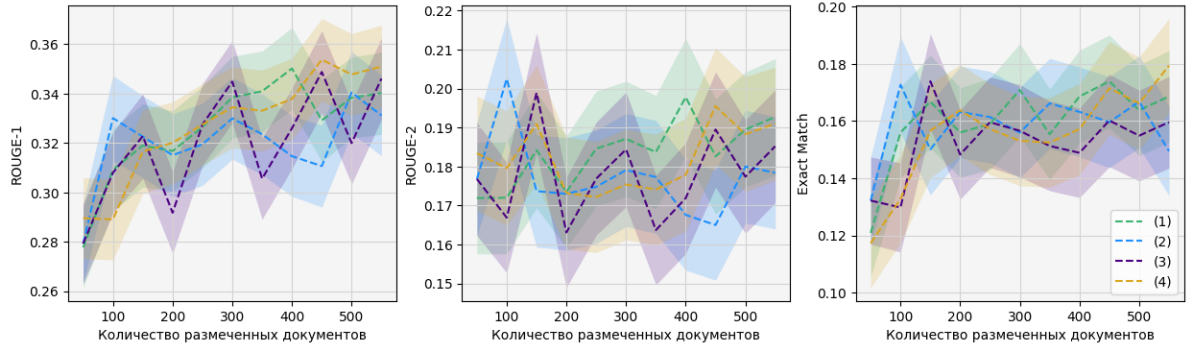


Рис. 6.5: Изменение метрик ROUGE-1, ROUGE-2 и Exact Match на валидационном датасете в задаче генерации ответа по заданной части контекста во время активного обучения.

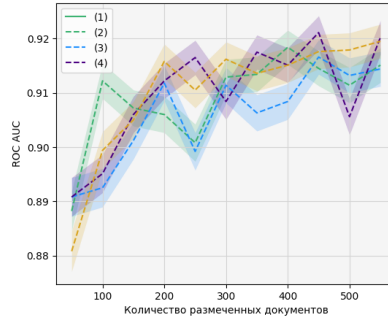


Рис. 6.6: Изменение метрики ROC AUC на валидационном датасете в задаче предсказания наличия ответа в заданной части контекста во время активного обучения.

В этой серии экспериментов наблюдается та же картина, что и в предыдущей: случайный выбор кандидатов не уступает представленным стратегиям активного обучения.

6.3 Распределение скоров кандидатов во время отбора

Для метода минимальной нормализованной вероятности последовательности (аналога метода минимальной уверенности для последовательностей) и метода внутридоменного разнообразия построены графики, отражающие распределение сора, по которому отбираются кандидаты, среди всех M_POOL кандидатов:

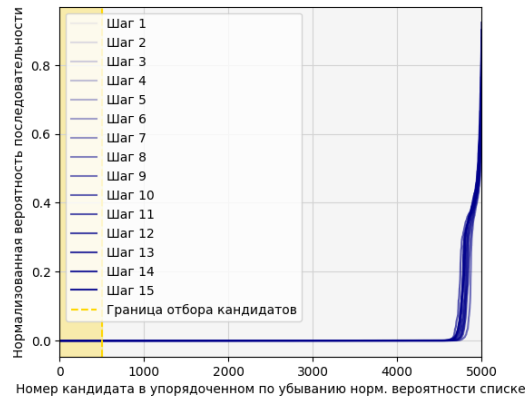


Рис. 6.7: Распределение нормированной вероятности последовательности среди всех кандидатов для каждого шага активного обучения.

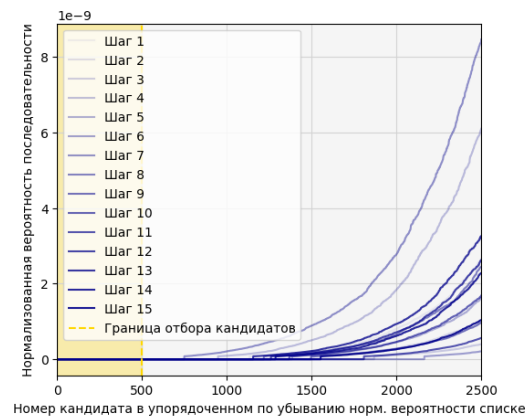


Рис. 6.8: Распределение нормированной вероятности последовательности среди 2500 кандидатов с ее наименьшим значением для каждого шага активного обучения.

Как можно видеть в случае нормированной вероятности последовательности, для подавляющего большинства кандидатов данный показатель почти равен 0, и только у приблизительно 6% кандидатов данная вероятность превышает 0.1. При этом оцениваемые последовательности предварительно отбираются моделью-классификатором, которая даже на первых шагах обучения достигает метрики ROC AUC около 0.9 (см. Рисунок 6.3), и, соответственно, в большинстве случаев верно выделяет часть контекста, содержащую в себе ответ. Метрика Exact Match при любой стратегии обучения уже на первом шаге превышает 0.15 (см. Рисунок 6.1). Метрики ROUGE-1, ROUGE-2 и Exact Match при использовании данной стратегии почти во всех точках проигрывают случайному отбору кандидатов, что заставляет предположить, что стратегия нормированной вероятности последовательности выбирает нерелевантные документы для обучающего датасета. Таким образом, отбор по стратегии минимальной нормированной вероятности последовательности представляется сомнительным для данной задачи.

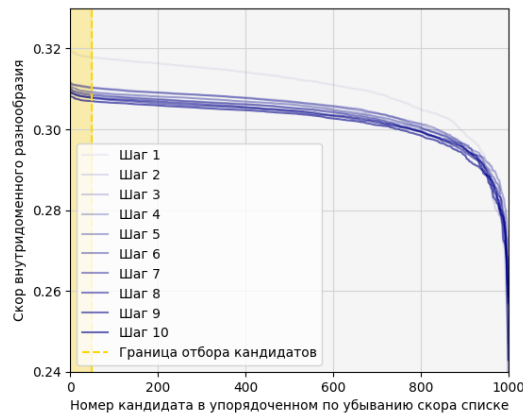


Рис. 6.9: Распределение IDS среди всех кандидатов для каждого шага активного обучения.

Распределение IDS (см. формулу 4) среди кандидатов сложнее поддается интерпретации. В районе минимальных значений данного показателя на Рисунке 6.9 можно видеть резкое снижение графика – вероятно, выбросы из распределения. В остальном же значения варьируются достаточно слабо, изменяясь у первых 80% кандидатов на графике менее чем на 4%.

7 Выводы

Отсутствие наблюдаемого преимущества в качестве моделей, обученных на датасетах, которые созданы при помощи активного обучения, по крайней мере отчасти может объясняться тем, что, в отличие от классической схемы активного обучения, при отборе кандидатов по отдельным частям документов в обучающий датасет добавляются документы целиком, и каждая из моделей вместе с выбранной частью-кандидатом получает и все остальные части данного документа.

По всей видимости, стратегия нормированной вероятности последовательности не подходит в качестве релевантной меры неопределенности, сопоставляя большинству документов скоры, крайне близкие к 0 (даже для тех документов, в которых бинарный классификатор верно выбрал часть, содержащую ответ, а генеративная модель верно его предсказала). Метрики качества при такой стратегии ниже, чем при случайном выборе.

В случае стратегии минимальной уверенности, которая опирается исключительно на предсказания бинарного классификатора, значимых отличий от случайного выбора нельзя увидеть даже в метрике ROC AUC, вычисляемой во время первого этапа предсказания для валидационного датасета. Однако, качество предсказаний бинарного классификатора сам по себе достаточно высоко – уже при обучении на 1000 документов метрика ROC AUC превы-

шает значение 0.925, тогда как при обучении на всем датасете, состоящем из 40783 документов, значение ROC AUC достигает 0.951. Можно предположить, что достаточно быстрое достижение близких к максимальным показателей качества затрудняет поиск наиболее информативных примеров для активного обучения.

Стратегия внутридоменного разнообразия, успешно примененная в задаче суммаризации, также не показала значимых отличий от случайной. Возможно, причиной тому является неоптимальный выбор параметра λ , формат эксперимента или особенности отбора кандидатов: в проведенных экспериментах отбирался сразу топ- M кандидатов, что напрямую влияет значения IDDS по сравнению с жадным отбором кандидатов по одному. Также в данном эксперименте не проводилось предварительного дообучения модели для получения эмбедингов на неразмеченном датасете для более явного "знакомства" модели с доменом [3].

8 Заключение

Выбор стратегии активного обучения в задаче генерации естественного языка остается открытой проблемой, решение которой может существенно сократить затраты человеческих ресурсов при обучении языковых моделей. В данной работе были рассмотрены несколько стратегий применительно к задаче генеративного построения ответа на вопрос с длинным контекстом. Показано, что ни одна из этих стратегий не дает прироста в качестве предсказания в сравнении со случайным выбором документов для разметки.

Список литературы

- [1] “A Sequential Algorithm for Training Text Classifiers”. B: *Proceedings of Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*. 1994, с. 3—12.
- [2] “A Survey of Extractive Question Answering”. B: *2022 International Conference on High Performance Big Data and Intelligent Systems (HDIS)*. 2022.
- [3] “Active Learning for Abstractive Text Summarization”. B: *Findings of the Association for Computational Linguistics: EMNLP 2022*. 2022, с. 5128—5152.
- [4] “Active learning for interactive neural machine translation of data streams”. B: *Proceedings of the 22nd Conference on Computational Natural Language Learning, CoNLL 2018, Brussels, Belgium, October 31 - November 1, 2018*. 2018, с. 151—160.
- [5] “Active Learning for Natural Language Generation”. B: *arXiv:2305.15040*. 2023.
- [6] “Active Learning for Sequence Tagging with Deep Pre-trained Models and Bayesian Uncertainty Estimates”. B: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*. 2021, с. 1698—1712.
- [7] “Active learning for sequence tagging with deep pre-trained models and Bayesian uncertainty estimates”. B: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 2021, с. 1698—1712.
- [8] “Active learning for statistical phrase-based machine translation”. B: *Human Language Technologies Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, May 31 - June 5, 2009, Boulder, Colorado, USA*. 2009, с. 415—423.
- [9] Vamshi Ambati. “Active Learning and Crowdsourcing for Machine Translation in Low Resource Scenarios”. B: *Ph.D. thesis, Language Technologies Institute School of Computer Science Carnegie Mellon University, USA*. 2012.
- [10] “AmbigQA: Answering ambiguous open-domain questions”. B: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*. 2020, с. 5783—5797.
- [11] “Batchmode semi-supervised active learning for statistical machine translation”. B: *Computer Speech Language*, 27(2). 2013, с. 397—406.
- [12] “Bayesian active learning with pretrained language models”. B: *arXiv preprint arXiv:2104.08320*. 2021.

- [13] “Cartography active learning”. B: *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*. 2021, c. 395—406.
- [14] “CoQA: A conversational question answering challenge”. B: *Transactions of the Association for Computational Linguistics*. 7. 2019, c. 249—266.
- [15] Sanjoy Dasgupta. “Two faces of active learning”. B: *Theoretical computer science* 412(19) (2011), c. 1767—1781.
- [16] Zoubin Ghahramani David A Cohn и Michael I Jordan. “Active learning with statistical models”. B: *Journal of artificial intelligence research* 4 (1996), c. 129—145.
- [17] “DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter”. B: *arXiv:1910.01108* 2020.
- [18] “ELI5: Long form question answering.” B: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2019, c. 3558—3567.
- [19] “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. B: *arXiv preprint arXiv:1910.10683*. 2020.
- [20] “Exploring the limits of transfer learning with a unified text-totext transformer”. B: *Journal of Machine Learning Research* 21(140) (2020), c. 1—61.
- [21] “HopRetriever: Retrieve Hops over Wikipedia to Answer Complex Questions”. B: *arXiv preprint arXiv:2012.15534*. 2020.
- [22] “Improving back-translation with uncertainty-based confidence estimation.” B: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019, c. 791—802.
- [23] “Learning hidden markov models for information extraction actively from partially labeled text”. B: *Künstliche Intell.* 16(2) (2002), c. 17—22.
- [24] “Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering”. B: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*. 2021, c. 874—880.

- [25] David D. Lewis и William A. Gale. “A Sequential Algorithm for Training Text Classifiers”. B: *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. Dublin, Ireland, 3-6 July 1994 (Special Issue of the SIGIR Forum)*. 1994, c. 3—12.
- [26] “Longformer: The long-document transformer”. B: *arXiv preprint arXiv:2004.05150*. 2020.
- [27] “Mind your outliers! investigating the negative impact of outliers on active learning for visual question answering.” B: *59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*. 2021, c. 7265—7281.
- [28] “Natural Questions: A Benchmark for Question Answering Research”. B: *Transactions of the Association for Computational Linguistics* 7(15) (2019), c. 453—466.
- [29] “Neural Generative Question Answering”. B: *Proceedings of 2016 NAACL Human-Computer Question Answering Workshop*. 2016, c. 36—42.
- [30] “Neural machine translation by jointly learning to align and translate”. B: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. 2015.
- [31] “Neural machine translation by jointly learning to align and translate”. B: *3rd International Conference on Learning Representations, ICLR 2015*. 2015.
- [32] “Retrieval-augmented generation for knowledge-intensive nlp tasks”. B: *arXiv preprint arXiv:2005.11462*. 2021.
- [33] “ROUGE: A package for automatic evaluation of summaries”. B: *Text Summarization Branches Out*. 2004, c. 74—81.
- [34] “Sampling bias in deep active classification: An empirical study”. B: *Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP/IJCNLP)*. 2019, c. 4058—4068.
- [35] “Sequence to sequence learning with neural networks”. B: *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014*. 2014, c. 3104—3112.
- [36] “SQuAD: 100,000+ questions for machine comprehension of text”. B: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. 2016, c. 2383—2392.

- [37] “The NarrativeQA reading comprehension challenge”. B: *Transactions of the Association for Computational Linguistics*. 6. 2018, c. 317–328.
- [38] “Toward optimal active learning through sampling estimation of error reduction”. B: *Eighteenth International Conference on Machine Learning*. 2001, c. 4058–4068.
- [39] “Wat zei je? detecting out-of-distribution translations with variational transformers”. B: *CoRR*, *abs/2006.08344*. 2020.
- [40] “Word-level confidence estimation for machine translation”. B: *Comput. Linguistics* 33(1) (2007), c. 9–40.
- [41] “You need only uncertain answers: Data efficient multilingual question answering”. B: *Workshop on Uncertainty and Ro-Bustness in Deep Learning*. 2020.