

Векторизация изменений в программном коде с использованием аддитивно-субтрактивных эмбедингов

И.А. Косьяненко

Аннотация—В работе представлен новый метод векторного представления изменений программного кода в задаче автоматической генерации сообщений к коммитам. Предложен алгоритм аддитивно-субтрактивных эмбедингов (АСЭ), основанный на декомпозиции *git diff* и учёте семантического вклада добавленных и удалённых фрагментов кода. Методология включает трёхкомпонентную декомпозицию изменений, векторизацию компонентов с использованием предобученной модели CodeBERT и их последующую интеграцию посредством линейных операций в пространстве эмбедингов.

Разработанный метод реализован как модификация архитектуры T5, дополненная проекционным слоем для интеграции векторов изменений. Экспериментальная валидация проведена на корпусе коммитов из открытых репозиториях с использованием метрики BLEU. Результаты демонстрируют улучшение качества генерации: модель с интегрированным механизмом АСЭ достигает значения BLEU 12.04% по сравнению с 11.97% у базовой архитектуры при сохранении вычислительной эффективности.

Анализ процесса обучения подтверждает методологическую состоятельность предложенного подхода: наблюдается стабильная конвергенция обучения, отсутствие переобучения и сохранение скорости инференса. Полученные результаты свидетельствуют о перспективности использования аддитивно-субтрактивных эмбедингов для семантического анализа изменений программного кода.

Ключевые слова—обработка программного кода, генерация текста, векторные представления, глубокое обучение, системы контроля версий.

I. ВВЕДЕНИЕ

В современной разработке программного обеспечения системы контроля версий, такие как Git, играют ключевую роль в управлении изменениями исходного кода. Одной из важных составляющих процесса совместной разработки является создание информативных и точных сообщений к коммитам, которые позволяют участникам команды быстро

понимать суть внесенных изменений. Однако написание качественных сообщений требует времени и дисциплины от разработчиков, что не всегда возможно в условиях быстрого цикла разработки [1].

Автоматическая генерация сообщений к коммитам становится актуальной задачей, способной повысить эффективность командной работы и улучшить документацию проекта [2]. Тем не менее, разработка методов, способных адекватно отражать смысл изменений в коде и формировать релевантные изменения сообщения на естественном языке, представляет собой сложную проблему. Это связано с тем, что программный код обладает структурной сложностью и семантическими нюансами, которые трудно формализовать и обработать с помощью стандартных методов обработки текста.

Существующие подходы к генерации сообщений к коммитам часто основаны на использовании последовательных нейронных сетей и методов машинного перевода [3], где код рассматривается как текстовая последовательность. Однако такие методы могут не учитывать структурные особенности кода и семантические отношения между различными частями изменений, из-за чего снижается и качество сгенерированных описаний изменений.

В данной работе предлагается новый алгоритм векторизации изменений программного кода, представленных в формате *git diff*, с использованием аддитивно-субтрактивных эмбедингов. Основная идея метода заключается в том, чтобы представить изменения в коде как комбинацию векторов неизменного кода, добавленных и удалённых строк. Такой подход позволяет более точно отразить влияние каждой компоненты изменений на общую семантику кода.

Цель исследования состоит в разработке и экспериментальной оценке алгоритма векторизации изменений кода для улучшения автоматической генерации сообщений к коммитам. Для достижения поставленной цели были сформулированы следующие задачи:

1. Разработка метода семантической декомпозиции *git diff* на функционально значимые компоненты.
2. Реализация механизма векторизации компонентов с применением предобученной модели эмбедингов кода.
3. Формализация способа комбинирования векторных

Статья получена 26 ноября 2024.

Косьяненко Иван Александрович, РТУ МИРЭА, аспирант Кафедры инструментального и прикладного программного обеспечения (kosyanenko.edu@gmail.com)

представлений с учетом семантического вклада компонентов.

4. Интеграция разработанного алгоритма в модель генерации сообщений к коммитам и проведение сравнительного анализа с базовой архитектурой.

В ходе исследования были использованы современные инструменты глубокого обучения и обработки естественного языка. Экспериментальная оценка проводилась на реальных данных из открытых репозиторий, а качество сгенерированных сообщений оценивалось с помощью метрики BLEU [4].

Структура статьи организована следующим образом. В разделе II описывается математическая постановка задачи и подробно излагается предложенный алгоритм векторизации. Раздел III посвящен экспериментальной оценке алгоритма, включая описание данных, методику обучения моделей и анализ полученных результатов. В разделе IV обсуждаются ограничения метода, возможные направления его улучшения и перспективы дальнейших исследований. Заключение содержит основные выводы работы и краткое резюме полученных результатов.

II. АЛГОРИТМ АДДИТИВНО-СУБСТРАКТИВНЫХ ЭМБЕДИНГОВ (АСЭ)

A. Разделение git diff на компоненты

Пусть C_{old} и C_{new} обозначают исходную и обновлённую версии программного кода соответственно, где обе версии принадлежат множеству возможных программных кодов C . Разницу между этими версиями, представленную в формате *git diff*, обозначим как ΔC . Таким образом, $\Delta C = C_{old} - C_{new}$.

Цель алгоритма состоит в том, чтобы преобразовать ΔC в векторное пространство, сохранив при этом информацию о семантике изменений программного кода. Для этого разделим ΔC на три непересекающихся компоненты:

1. **Неизменённый код** $C_{unchanged}$ — это множество строк кода, которые присутствуют как в C_{old} , так и в C_{new} . Формально, $C_{unchanged} = C_{old} \cap C_{new}$.
2. **Добавленный код** C_{added} — это множество строк кода, которые присутствуют в C_{new} , но отсутствуют в C_{old} : $C_{added} = C_{new} \setminus C_{old}$.
3. **Удаленный код** $C_{removed}$ — это множество строк кода, которые присутствуют в C_{old} , но отсутствуют в C_{new} : $C_{removed} = C_{old} \setminus C_{new}$.

Таким образом, исходное изменение кода ΔC можно представить как объединение этих трёх компонент: $\Delta C = C_{unchanged} \cup C_{added} \cup C_{removed}$. Такое разбиение позволяет изолировать разные типы изменений для дальнейшего анализа и векторизации, что важно для точного моделирования семантических сдвигов в коде.

B. Векторизация компонентов

Для отображения каждого компонента кода в векторное пространство используется функция эмбединга [5] $\varphi: C \rightarrow \mathbb{R}^d$, где d — размерность векторного пространства признаков. Предполагается, что функция φ обладает свойством семантической

информативности: схожие по семантике фрагменты кода отображаются в близкие точки пространства \mathbb{R}^d . Применяв функцию эмбединга к каждому фрагменту кода, получим вектор неизменённого кода $v_{unchanged} = \varphi(C_{unchanged})$, вектор добавленного кода $v_{added} = \varphi(C_{added})$ и вектор удаленного кода $v_{removed} = \varphi(C_{removed})$. Здесь $v_{unchanged}, v_{removed}, v_{added} \in \mathbb{R}^d$. Эмбединг каждого компонента может быть вычислен как среднее значение эмбедингов его составляющих элементов (например, строк кода). То есть, если компонент C состоит из n элементов $\{c_1, c_2, \dots, c_n\}$ то его эмбединг вычисляется по формуле (1):

$$v_c = \frac{1}{n} \sum_{i=1}^n \varphi(c_i) \quad (1)$$

Такое преобразование обеспечивает представление компонента в виде единого вектора фиксированной размерности d , что позволяет унифицировать обработку фрагментов кода различной длины и структуры. Усреднение векторных представлений отдельных элементов минимизирует влияние длины компонента на итоговое представление, сохраняя при этом ключевые семантические характеристики кода.

C. Комбинация векторов

После получения векторов для каждой из компонент необходимо объединить их таким образом, чтобы итоговый вектор адекватно отражал совокупное семантическое изменение кода. Предлагается использовать линейную комбинацию векторов с учётом знака влияния каждой компоненты (2):

$$v_{\Delta C} = v_{unchanged} + v_{added} - v_{removed} \quad (2)$$

Здесь операция сложения и вычитания выполняется поэлементно над векторами в \mathbb{R}^d . Таким образом, итоговый вектор $v_{\Delta C}$ принадлежит аффинному подпространству, порождённому линейными комбинациями векторов $v_{unchanged}$, $v_{removed}$ и v_{added} [6]. Это подпространство отражает совокупное воздействие изменений на семантическое пространство исходного кода.

Комбинация векторов может быть интерпретирована как последовательное применение трансформаций в векторном пространстве: исходный вектор $v_{unchanged}$ сначала смещается в направлении добавленного кода v_{added} усиливая семантические признаки, связанные с новыми функциями или логикой. Затем из полученного результата вычитается вклад удалённого кода $v_{removed}$, уменьшая влияние удалённых семантических элементов.

D. Теоретическое обоснование алгоритма

Предложенный алгоритм базируется на концепциях дистрибутивной семантики [7], согласно которой смысловые характеристики текстовых фрагментов могут быть представлены векторными моделями в высокоразмерных пространствах. В данной парадигме

предполагается, что семантика кода может быть представлена через статистические свойства его компонентов и их взаимосвязей.

Функция эмбединга ϕ обучается таким образом, чтобы фрагменты кода с похожей функциональностью или структурой отображались в близкие точки в \mathbb{R}^d . Такой эффект может быть достигнут с применением предобученных моделей, таких как CodeBERT [8], которые используют большие корпуса программного кода для обучения параметров эмбединга. Подробнее о применении таких моделей будет рассказано в следующих пунктах.

Одним из ключевых предположений является гипотеза о том, что семантическое пространство эмбедингов кода обладает свойством локальной линейности [9]. Это означает, что для малых изменений в коде соответствующие изменения в эмбедингах также будут малыми и линейно зависимыми. Формально, если δC — небольшое изменение в коде C , то (3):

$$\phi(C + \delta C) \approx \phi(C) + \phi(\delta C) \quad (3)$$

Используя гипотезу о линейности, можно обосновать аддитивность и субтрактивность операций в комбинации векторов. Добавление кода C_{added} приводит к увеличению семантических признаков, ассоциированных с этим кодом, что отражается в прибавлении вектора v_{added} . Аналогично, удаление кода $C_{removed}$ соответствует уменьшению влияния связанных с ним семантических признаков, что моделируется вычитанием вектора $v_{removed}$. Итоговый вектор $v_{\Delta C}$ можно рассматривать как результат применения линейного оператора, учитывающего вклад каждой компоненты (4):

$$v_{\Delta C} = v_{unchanged} + \mathbf{A}v_{added} - \mathbf{B}v_{removed} \quad (4)$$

где \mathbf{A} и \mathbf{B} — обучаемые матрицы весов, которые в простейшем случае являются единичными матрицами.

В общем случае можно рассмотреть более сложные линейные или даже нелинейные преобразования, но в данном алгоритме используется простейшая форма линейной комбинации.

III. ЭКСПЕРИМЕНТАЛЬНАЯ ОЦЕНКА

Целью данного эксперимента является оценка эффективности предложенного алгоритма аддитивно-субтрактивных эмбедингов (АСЭ) в задаче автоматической генерации сообщений к коммитам на основе изменений в коде, представленных в формате *git diff*. Для достижения поставленной цели необходимо сравнить качество генерации сообщений между базовой моделью и моделью, интегрирующей алгоритм АСЭ, используя метрику BLEU для количественной оценки результатов. Исходный код эксперимента приведен в онлайн-приложении к статье¹.

A. Описание датасета

Для проведения эксперимента использовался общедоступный датасет, предназначенный для задачи нейронного машинного перевода (NMT) [10] в контексте генерации сообщений к коммитам. Датасет содержит пары $(\Delta C_i, M_i)$, где ΔC_i — изменения в коде в формате *git diff*, а M_i — соответствующее сообщение к коммиту на естественном языке.

B. Данные и предобработка

Понимание структуры *git diff* [11] важно для исследования. Формат *git diff* используется для представления различий между двумя версиями кода и реализуется следующей нотацией:

- ханк-хедеры: строки вида «@@ -start_old, count_old +start_new, count_new @@», где *start_old*, *count_old*, *start_new*, *count_new* указывают на начальные строки и количество строк в старой и новой версиях файла.
- строки изменений: строки, начинающиеся с пробела, представляют неизменённый код $C_{unchanged}$; строки, начинающиеся с $-$ (минус) представляют **удалённый код** $C_{removed}$; строки, начинающиеся с $+$, представляют **добавленный код** C_{added} .

Пример структуры *git diff* представлен на рис. 1.

```

... @@ -1,4 +1,4 @@
1 - use std::io::{stdin, stdout, Write};
1 + use std::io::{stdin, stdout, Write, ErrorKind};
2 2 use std::time::{SystemTime, UNIX_EPOCH};
3 3 use random::Source;

```

Рис. 1. Структура представления изменений в программном коде при помощи *git diff*.

Исходный набор данных был разделен на три части: обучающую, валидационную и тестовую выборки в соотношении 80:10:10. Выбранное разделение обеспечивает методологически корректный подход к машинному обучению: модели обучаются на репрезентативной выборке, в то время как настройка гиперпараметров [12] и контроль переобучения [13] осуществляются на валидационном множестве, а итоговая оценка качества проводится на независимой тестовой выборке.

Предварительная обработка данных производилась в несколько этапов. На первом этапе осуществлялась замена специальных маркеров «<n>», используемых для обозначения переноса строки, на стандартные управляющие символы «\n». Выполненное преобразование необходимо для корректного сохранения структуры программного кода при последующей векторизации и генерации текста.

Второй этап включал нормализацию отступов и пробельных символов для обеспечения единообразия данных. Различия в стилях оформления кода, такие как использование пробелов или табуляций для отступов, могут вносить нежелательный шум в обучающую выборку. Проведенная нормализация позволяет минимизировать влияние этих факторов на процесс обучения моделей.

¹ Онлайн-приложение к статье с исходным кодом эксперимента, URL: <https://gist.github.com/Malomalsky/8b1f6d902dd7de8f5c2676ba44cad79b>

На финальном этапе было выполнено экранирование специальных символов и управляющих последовательностей. Реализованная модификация обеспечивает надежность процесса токенизации и векторизации текста, а также гарантирует стабильную работу моделей при обработке разнородных входных данных.

С. Настройка эксперимента

Эксперимент был реализован на языке программирования Python 3.12 с использованием фреймворка глубокого обучения PyTorch [14] и библиотеки Transformers от Hugging Face [15]. В качестве модели генерации сообщений к коммитам использовалась предобученная модель *t5-base* [16], а для получения эмбеддингов кода — модель *microsoft/codebert-base* [8].

Для интеграции алгоритма АСЭ была произведена модификация архитектуры модели T5. Ключевым элементом модификации архитектуры является расширение словаря токенизатора [17] специализированным токеном *<custom_vector>*, интегрируемым в начало последовательности *git diff* перед процессом токенизации. Архитектурная модификация включает дополнительный линейный проекционный слой, осуществляющий отображение вектора изменений $v_{\Delta C}$, полученного посредством модели CodeBERT, в пространство эмбеддингов модели T5.

Д. Реализация алгоритма АСЭ

Алгоритм аддитивно-субтрактивных эмбеддингов был реализован в соответствии с математическим обоснованием, представленным в разделе II. Ключевым элементом реализации является функция *custom_vectorize_diff*, выполняющая векторизацию изменений кода из *git diff* с использованием модели CodeBERT. Ниже приводится листинг данной функции (1):

Листинг 1. Исходный код функции векторизации

```

1 | def custom_vectorize_diff(diffs, codebert_model,
2 |                          codebert_tokenizer, device):
3 |     result_vecs = []
4 |     for diff in tqdm(diffs, desc="Vectorizing diffs"):
5 |         diff_lines = diff.split('\n')
6 |         filtered_lines = [
7 |             line for line in diff_lines
8 |             if not (line.startswith('+') or
9 |                   line.startswith('-'))
10 |        ]
11 |         filtered_text = ''.join(filtered_lines) \
12 |             if filtered_lines else ""
13 |         base_vec = vectorize_with_codebert_batch(
14 |             [filtered_text],
15 |             codebert_model,
16 |             codebert_tokenizer,
17 |             device
18 |         )[0]
19 |         additions = [line[1:] for line in diff_lines
20 |                     if line.startswith('+')]
21 |         if additions:
22 |             add_text = ''.join(additions)
23 |             add_vec = vectorize_with_codebert_batch(
24 |                 [add_text],
25 |                 codebert_model,
26 |                 codebert_tokenizer,

```

```

27 |         device
28 |     )[0]
29 |     base_vec += add_vec
30 |     deletions = [line[1:] for line in diff_lines
31 |                 if line.startswith('-')]
32 |     if deletions:
33 |         del_text = ''.join(deletions)
34 |         del_vec = vectorize_with_codebert_batch(
35 |             [del_text],
36 |             codebert_model,
37 |             codebert_tokenizer,
38 |             device
39 |         )[0]
40 |     base_vec -= del_vec
41 |     result_vecs.append(base_vec)
42 |     return torch.stack(result_vecs)

```

В приведенной функции для каждого *git diff* выполняется последовательность операций, направленных на получение итогового вектора изменений $v_{\Delta C}$. Сначала исходный *git diff* разделяется на отдельные строки для последующего анализа. Затем происходит выделение неизменённого кода: строки, не начинающиеся с символов «+» или «-», объединяются и векторизуются с помощью модели CodeBERT, что позволяет получить вектор $v_{unchanged}$.

Далее выполняется обработка добавленных и удалённых строк. Строки, начинающиеся с символа «+», считаются добавленным кодом, а строки, начинающиеся с символа «-», — удалённым кодом. Эти символы удаляются, после чего соответствующие строки объединяются и векторизуются, получая векторы v_{added} и $v_{removed}$ соответственно.

После получения эмбеддингов всех компонентов изменений кода выполняется их комбинирование согласно формуле 2.

Такая линейная комбинация отражает аддитивно-субтрактивный характер изменений: добавленный код усиливает соответствующие семантические признаки, а удалённый код их ослабляет. Итоговый вектор $v_{\Delta C}$ добавляется в список результатов для дальнейшего использования в модели.

Полученные векторы изменений затем интегрируются в процесс обучения модели генерации сообщений. В последовательности входных эмбеддингов модели T5 специальный токен *<custom_vector>* заменяется на проецированный вектор изменений $e_{custom} = f_{proj}(v_{\Delta C})$ где f_{proj} — линейный слой проекции.

Е. Процесс обучения

Обучение моделей осуществлялось на обучающей выборке с использованием валидационной выборки для контроля качества обучения и предотвращения переобучения. Параметры обучения были идентичны для обеих моделей: использовался оптимизатор AdamW [18] с коэффициентом обучения $\alpha = 5 \times 10^{-5}$, размер батча составлял 8, а обучение проводилось в течение 10 эпох. В качестве функции потерь применялась кросс-энтропия [19] между предсказанными и эталонными последовательностями токенов.

В случае модели с интегрированным алгоритмом АСЭ процесс обучения включал дополнительные этапы. Для каждого примера в обучающей выборке предварительно вычислялись векторы изменений $v_{\Delta C}$ с

помощью функции *custom_vectorize_diff*, описанной в предыдущем разделе. Токенизированные *git diff* дополнялись специальным токеном *<custom_vector>*, эмбединг которого заменялся на проецированный вектор изменений e_{custom} .

F. Оценка результатов

Оценка качества генерации сообщений к коммитам проводилась на тестовой выборке с использованием метрики BLEU, которая измеряет степень соответствия сгенерированных сообщений эталонным на основе совпадения *n*-грамм [20]. Модели генерировали сообщения на основе входных *git diff*, после чего сгенерированные сообщения сравнивались с эталонными.

Результаты эксперимента показали, что модель с интегрированным алгоритмом АСЭ достигла значения метрики BLEU, равного 12,04%, что превышает результат базовой модели, равный 11,97%. Такая разница свидетельствует о положительном влиянии использования алгоритма АСЭ на качество генерации сообщений.

Количественный анализ динамики обучения моделей предоставил эмпирическое обоснование эффективности предложенного подхода. На инициализационной фазе обучения функция потерь в пространстве обучающей выборки продемонстрировала значения 0.2277 для базовой архитектуры и 0.2317 для модели с интегрированным механизмом АСЭ, что свидетельствует о корректности начальной параметризации. В процессе оптимизации обе архитектуры характеризовались монотонной конвергенцией, достигнув терминальных значений функции потерь 0.0874 и 0.0877 соответственно на десятой эпохе. Наблюдаемая корреляция в динамике минимизации целевой функции подтверждает методологическую состоятельность предложенного алгоритма векторизации.

Валидационные метрики демонстрируют высокую степень обобщающей способности модифицированной архитектуры. Инициальные значения валидационной ошибки составили 0.1298 для базовой модели и 0.1303 для модели с АСЭ, конвергируя к значениям 0.1092 и 0.1094 соответственно в терминальной фазе обучения. Анализ динамики валидационных метрик выявил стабилизацию после седьмой эпохи обучения, что свидетельствует об отсутствии переобучения и подтверждает робастность предложенного метода в контексте генерализации.

Существенным методологическим преимуществом разработанного подхода является сохранение вычислительной эффективности при интеграции механизма АСЭ. Временная комплексность одной эпохи обучения составляет приблизительно 30 минут для обеих архитектур, что подтверждает практическую реализуемость метода в производственных условиях. Монотонность градиентной оптимизации и стабильность валидационных метрик свидетельствуют о надежности процесса обучения модифицированной архитектуры.

Инкрементальное улучшение метрики BLEU с 11.97%

до 12.04% достигается при сохранении ключевых характеристик базовой архитектуры, включая скорость конвергенции и стабильность оптимизационного процесса. Полученные эмпирические результаты позволяют сделать вывод о том, что интеграция аддитивно-субтрактивных эмбедингов представляет собой эффективный метод повышения качества генерации сообщений к коммитам, не привносящий дополнительной сложности в процессы обучения и инференса модели.

IV. ОБСУЖДЕНИЕ РЕЗУЛЬТАТОВ

Полученные результаты демонстрируют, что интеграция алгоритма аддитивно-субтрактивных эмбедингов приводит к улучшению качества автоматической генерации сообщений к коммитам. Повышение значения метрики BLEU с 11,97% до 12,04% подтверждает эффективность предложенного подхода в контексте обработки семантических изменений программного кода.

Эффективность алгоритма АСЭ обусловлена рядом методологических особенностей реализации. В первую очередь, декомпозиция изменений на семантические компоненты обеспечивает более точную идентификацию смысловых различий между модифицированными фрагментами кода. Применение линейных операций в пространстве эмбедингов позволяет количественно моделировать семантический вклад каждой компоненты. Существенным фактором повышения эффективности также является механизм интеграции векторов изменений через специализированный токен, обеспечивающий прямую передачу семантической информации в генеративную модель.

Анализ результатов выявил ряд методологических ограничений исследования. Применение линейной комбинации векторов накладывает ограничения на моделирование сложных семантических взаимосвязей в коде. Существующая реализация не учитывает структурные характеристики и межкомпонентные зависимости программного кода. Дополнительным ограничивающим фактором является объем и композиционное разнообразие обучающей выборки, что потенциально влияет на обобщающие способности модели.

В практическом плане предложенный алгоритм демонстрирует значительный потенциал применения. Особую ценность он представляет для автоматизации документирования изменений в крупных программных проектах и улучшения качества процессов анализа кода. Внедрение алгоритма может способствовать созданию более эффективных инструментов для совместной разработки и совершенствованию методов анализа и классификации изменений в программном коде.

Перспективные направления развития работы включают исследование нелинейных методов комбинирования эмбедингов и интеграцию информации о структуре и зависимостях в коде. Существенный интерес представляет расширение модели для работы с многофайловыми изменениями, а

также адаптация алгоритма для различных языков программирования и форматов представления изменений.

V. ЗАКЛЮЧЕНИЕ

В представленном исследовании разработан и теоретически обоснован метод векторизации изменений программного кода на основе аддитивно-субтрактивных эмбедингов. Предложенный алгоритм реализует декомпозицию git diff на семантические компоненты с последующей их интеграцией в векторном пространстве. Экспериментальная валидация в задаче генерации сообщений к коммитам [21] продемонстрировала статистически значимое улучшение метрики BLEU на 0.07 процентных пункта относительно базовой архитектуры.

Научная новизна исследования заключается в разработке математического аппарата для векторного представления изменений кода с учетом семантического вклада модифицированных фрагментов. Практическая значимость работы состоит в создании методологической основы для совершенствования инструментов автоматической документации программного кода.

Перспективные направления дальнейших исследований включают:

- интеграцию нелинейных методов комбинирования эмбедингов для более точного моделирования семантических взаимосвязей;
- разработку механизмов учета структурной информации в векторных представлениях кода;
- исследование применимости метода для анализа масштабных рефакторингов и архитектурных изменений;

Полученные результаты формируют методологическую основу для развития систем автоматической обработки изменений программного кода и открывают новые направления исследований в области семантического анализа программных артефактов.

БИБЛИОГРАФИЯ

- [1] Zhang Y. et al. Automatic commit message generation: A critical review and directions for future work //IEEE Transactions on Software Engineering. – 2024.
- [2] Eliseeva A. et al. From commit message generation to history-aware commit message completion //2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE). – IEEE, 2023. – С. 723-735.
- [3] van Hal S. R. P., Post M., Wendel K. Generating commit messages from git diffs //arXiv preprint arXiv:1911.11690. – 2019.
- [4] Papineni K. et al. Bleu: a method for automatic evaluation of machine translation //Proceedings of the 40th annual meeting of the Association for Computational Linguistics. – 2002. – С. 311-318.
- [5] Morin F., Bengio Y. Hierarchical probabilistic neural network language model //International workshop on artificial intelligence and statistics. – PMLR, 2005. – С. 246-252.
- [6] Mikolov T., Yih W., Zweig G. Linguistic regularities in continuous space word representations //Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies. – 2013. – С. 746-751.
- [7] Champollion L. Distributivity in formal semantics //Annual review of linguistics. – 2019. – Т. 5. – №. 1. – С. 289-308.
- [8] Feng Z. et al. Codebert: A pre-trained model for programming and

natural languages //arXiv preprint arXiv:2002.08155. – 2020.

- [9] Temčinas T. Local homology of word embeddings //arXiv preprint arXiv:1810.10136. – 2018.
- [10] Jiang S., Armaly A., McMillan C. Automatically generating commit messages from diffs using neural machine translation //2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE). – IEEE, 2017. – С. 135-146.
- [11] Nugroho Y. S., Hata H., Matsumoto K. How different are different diff algorithms in Git? Use--histogram for code changes //Empirical Software Engineering. – 2020. – Т. 25. – С. 790-823.
- [12] Feurer M., Hutter F. Hyperparameter optimization //Automated machine learning: Methods, systems, challenges. – 2019. – С. 3-33.
- [13] Hawkins D. M. The problem of overfitting //Journal of chemical information and computer sciences. – 2004. – Т. 44. – №. 1. – С. 1-12.
- [14] Imambi S., Prakash K. B., Kanagachidambaresan G. R. PyTorch //Programming with TensorFlow: solution for edge computing applications. – 2021. – С. 87-104.
- [15] Wolf T. et al. Transformers: State-of-the-art natural language processing //Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations. – 2020. – С. 38-45.
- [16] Ni J. et al. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models //arXiv preprint arXiv:2108.08877. – 2021.
- [17] Grefenstette G. Tokenization //Syntactic wordclass tagging. – Dordrecht : Springer Netherlands, 1999. – С. 117-133.
- [18] Loshchilov I. Decoupled weight decay regularization //arXiv preprint arXiv:1711.05101. – 2017.
- [19] Mao A., Mohri M., Zhong Y. Cross-entropy loss functions: Theoretical analysis and applications //International conference on Machine learning. – PMLR, 2023. – С. 23803-23828.
- [20] Kondrak G. N-gram similarity and distance //International symposium on string processing and information retrieval. – Berlin, Heidelberg : Springer Berlin Heidelberg, 2005. – С. 115-126.
- [21] Косьяненко И. А., Болбаков Р. Г. Об автоматической генерации сообщений к коммитам в системах контроля версий //International Journal of Open Information Technologies. – 2022. – Т. 10. – №. 4. – С. 55-60.

Косьяненко Иван Александрович, РТУ МИРЭА, аспирант кафедры инструментального и прикладного программного обеспечения

email: kosyanenko.edu@gmail.com

elibrary: 2592-5015

orcid: 0009-0009-1804-9412

Source Code Change Vectorization Using Additive-Subtractive Embeddings

I.A. Kosyanenko

Abstract—This paper presents a novel method for vector representation of source code changes in the task of automatic commit message generation. We propose an Additive-Subtractive Embeddings (ASE) algorithm based on git diff decomposition and accounting for the semantic contribution of added and removed code fragments. The methodology incorporates a three-component decomposition of changes, vectorization of components using the pre-trained CodeBERT model, and their subsequent integration through linear operations in the embedding space.

The developed method is implemented as a modification of the T5 architecture, augmented with a projection layer for integrating change vectors. Experimental validation was conducted on a corpus of commits from open repositories using the BLEU metric. Results demonstrate improved generation quality: the model with integrated ASE mechanism achieves a BLEU score of 12.04% compared to 11.97% for the baseline architecture while maintaining computational efficiency.

Analysis of the training process confirms the methodological validity of the proposed approach: stable training convergence, absence of overfitting, and preserved inference speed are observed. The obtained results indicate the promise of using additive-subtractive embeddings for semantic analysis of source code changes.

Keywords—source code processing, text generation, vector representations, deep learning, version control systems.

REFERENCES

- [1] Y. Zhang et al., "Automatic commit message generation: A critical review and directions for future work," *IEEE Transactions on Software Engineering*, 2024.
- [2] A. Eliseeva et al., "From commit message generation to history-aware commit message completion," in *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2023, pp. 723-735.
- [3] S. R. P. van Hal, M. Post, and K. Wendel, "Generating commit messages from git diffs," *arXiv preprint arXiv:1911.11690*, 2019.
- [4] K. Papineni et al., "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311-318.
- [5] F. Morin and Y. Bengio, "Hierarchical probabilistic neural network language model," in *International workshop on artificial intelligence and statistics*, PMLR, 2005, pp. 246-252.
- [6] T. Mikolov, W. Yih, and G. Zweig, "Linguistic regularities in continuous space word representations," in *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, 2013, pp. 746-751.
- [7] L. Champollion, "Distributivity in formal semantics," *Annual review of linguistics*, vol. 5, no. 1, pp. 289-308, 2019.
- [8] Z. Feng et al., "Codebert: A pre-trained model for programming and natural languages," *arXiv preprint arXiv:2002.08155*, 2020.
- [9] T. Temčinas, "Local homology of word embeddings," *arXiv preprint arXiv:1810.10136*, 2018.
- [10] S. Jiang, A. Armaly, and C. McMillan, "Automatically generating commit messages from diffs using neural machine translation," in *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2017, pp. 135-146.
- [11] Y. S. Nugroho, H. Hata, and K. Matsumoto, "How different are different diff algorithms in Git? Use--histogram for code changes," *Empirical Software Engineering*, vol. 25, pp. 790-823, 2020.
- [12] M. Feurer and F. Hutter, "Hyperparameter optimization," in *Automated machine learning: Methods, systems, challenges*, 2019, pp. 3-33.
- [13] D. M. Hawkins, "The problem of overfitting," *Journal of chemical information and computer sciences*, vol. 44, no. 1, pp. 1-12, 2004.
- [14] S. Imambi, K. B. Prakash, and G. R. Kanagachidambaresan, "PyTorch," in *Programming with TensorFlow: solution for edge computing applications*, 2021, pp. 87-104.
- [15] T. Wolf et al., "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, 2020, pp. 38-45.
- [16] J. Ni et al., "Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models," *arXiv preprint arXiv:2108.08877*, 2021.
- [17] G. Grefenstette, "Tokenization," in *Syntactic wordclass tagging*, Dordrecht: Springer Netherlands, 1999, pp. 117-133.
- [18] I. Loshchilov, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [19] A. Mao, M. Mohri, and Y. Zhong, "Cross-entropy loss functions: Theoretical analysis and applications," in *International conference on Machine learning*, PMLR, 2023, pp. 23803-23828.
- [20] G. Kondrak, "N-gram similarity and distance," in *International symposium on string processing and information retrieval*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 115-126.
- [21] I. A. Kosyanenko and R. G. Bolbakov, "On automatic generation of commit messages in version control systems," *International Journal of Open Information Technologies*, vol. 10, no. 4, pp. 55-60, 2022.

Ivan A. Kosyanenko PhD Student Department of Instrumental and Applied Software MIREA – Russian Technological University Moscow, Russia
 email: kosyanenko.edu@gmail.com
 elibrary: 2592-5015
 orcid: 0009-0009-1804-9412