

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ»  
(СПбГУ)

Образовательная программа бакалавриата «Науки о данных»



**Отчёт по практике**  
**Учебная практика (научно-исследовательская работа)**

Выполнил студент 4 курса бакалавриата  
(группа 22.Б05-мкн)  
Михеев Артём Антонович

Научный руководитель:  
Кандидат физико-математических наук  
Ершов Василий Алексеевич

Санкт-Петербург

Текстовые эмбединги финансовых новостей, обратимые по смыслу.

## Содержание

<b>Введение. Постановка задачи. Актуальность.</b>	<b>3</b>
Полезные статьи.	4
Формулировка задачи	4
<b>Основные результаты практики</b>	<b>5</b>
<i>Dataset.</i>	5
<i>Метрики.</i>	5
<i>Архитектура моделей.</i>	7
<i>Обучение</i>	10
<i>Графики</i>	11
<i>Результаты</i>	13
<i>Генерация. GREEDY vs BEAM-SEARCH</i>	14
LOSS на уровне предложений	15
Изучение полученных эмбедингов	17
<i>Примеры восстановления.</i>	17
<b>Заключение</b>	<b>18</b>

# Введение. Постановка задачи.

## Актуальность.

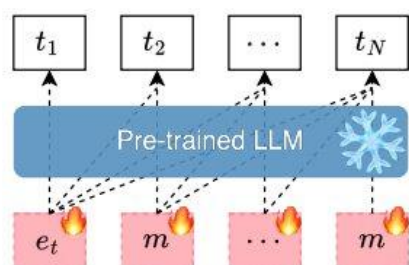
Обращать эмбединги токенов довольно просто. Находим эмбединг в матрице эмбедингов и номер строчки является порядковым номером токена. Используя токенизатор, декодируем и получаем символы.

Предложений бывает очень много, и никакой памяти не хватит для хранения всех эмбедингов. Задача заключается в построении текстовых эмбедингов с возможностью восстановления по смыслу. Обратить полностью не получится, так как какая-то часть информации теряется, но получить схожее по смыслу хотелось бы.

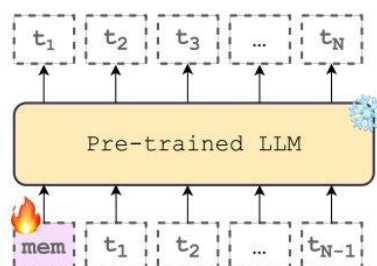
Предположим, мы проектируем мультимодальную модель для обработки данных разной природы. К примеру, такая модель может на основе числовых финансовых данных (временных рядов), новостного текста, вспомогательных таблиц предсказывать различные показатели: индексы фондового рынка, курс валют, возможен ли дефолт? Такие модели плохо интерпретируемые. Встает задача совместной генерации этих данных и текста-интерпретации – объяснения, почему произошло изменение этих финансовых показателей. Мы не можем генерировать токен за токеном, так как loss будет сосредотачиваться только на тексте, плохо генерируя показатели. Нам важно уловить общий смысл, получить одно представление (вектор) этого объяснения. Если получить такое представление, из которого можно сгенерировать текст, тогда для новой задачи нужно только научить получать это представление, а декодер уже есть.

## *Полезные статьи.*

Авторы изучали возможность получения представления конкретного текста, то есть без encoder части. Вводились токены памяти (на схеме - mem) перед текстом и изучалось сжатие в эти токены. Основной акцент был на изучении локальности полученных представлений при разной инициализации. Авторегрессионный подход к генерации. И прямой, менее качественный, но быстрый, только на основе эмбединга текста, без предыдущих токенов. Был сделан вывод о возможности построения encoder части.



*прямой*



*авторегрессионный*

## Формулировка задачи

*Проверка гипотезы о возможности (и способах) конструирования такого эмбединга и изучения насколько можно сжимать тексты сохраняя смысл.*

# Основные результаты практики

Наша задача спроектировать и обучить как энкодер для получения одного представления, так более важная часть модели - decoder для смыслового восстановления.

## ***Dataset.***

Будем обучаться на финансовых новостях [Financial News Headlines Data](#) с kaggle.

Заголовки этих наборов данных, взятые с официальных сайтов CNBC, The Guardian и Reuters, отражают обзор экономики и фондового рынка США за каждый день за 2018-2020.

Пример:

*U.S. economy faces significant risks, long road to recovery: IMF staff.*

*UBS, Morgan Stanley expected to lead Vodafone Tower IPO: sources.*

*As big U.S. banks let customers delay payments, loan losses remain unclear.*

*China says it will act to protect its interests after UK Huawei ban.*

Всего 53370 заголовков. Делим данные на train и validation, validation size = 0.1.

## ***Метрики.***

Надо оценивать на сколько декодируемый текст похож на оригинал. Будем оценивать следующие метрики:

1) Совпадение слов.

*ACCURACY*. Сопоставляем токены по позициям.

2) Семантическая близость.

*BERTScore*.

Модель: `'microsoft/deberta-xlarge-mnli'`

- Получает эмбединги для каждого токена из двух предложений.
- Вычисляет семантическое соответствие токенов, сопоставляя каждый токен из output с ближайшим по смыслу токеном в input (и наоборот).  $O(T^2)$ .
- Итоговые P, R, F1 — агрегированные значения сходства между токенами.

`bert_score("The car was damaged in the accident.", "The vehicle got wrecked during the crash.") = 0.91`

`bert_score("The economy is growing rapidly.", "My dog loves running in the park.") = 0.35`

*PARAPH-SIM*.

Модель: `'sentence-transformers/paraphrase-mpnet-base-v2'`

- Каждое предложение сжимается в один вектор фиксированной размерности 768.
- Вычисляется косинусное сходство между двумя векторами:

$$\text{cos\_sim}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

`paraph_sim("The car was damaged in the accident.", "The vehicle got wrecked during the crash.") = 0.821`

paraph\_sim("The economy is growing rapidly.", 'My dog loves running in the park.') =  
**0.17**

## ***Размерности.***

B - batch\_size

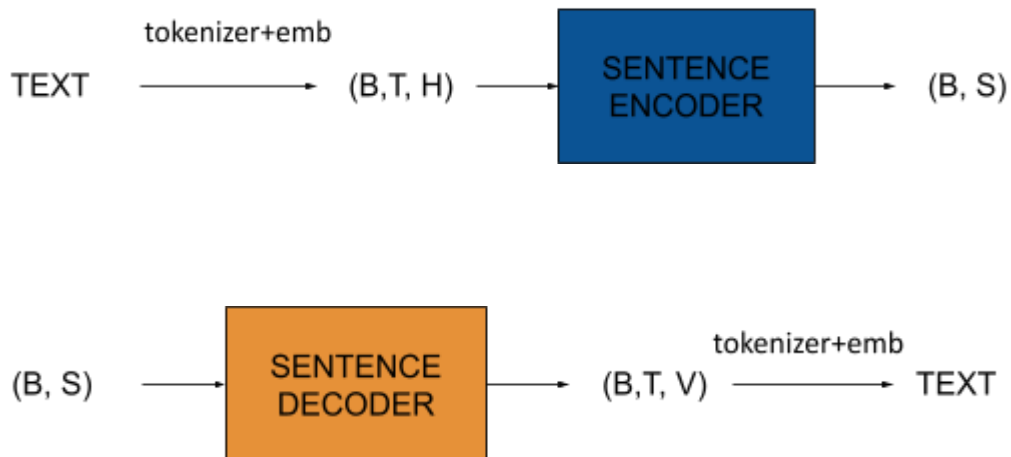
T - длина предложения в токенах с padding для обучения по бачам.

H - размерность эмбедингов токенов.

S - размерность эмбедингов предложений.

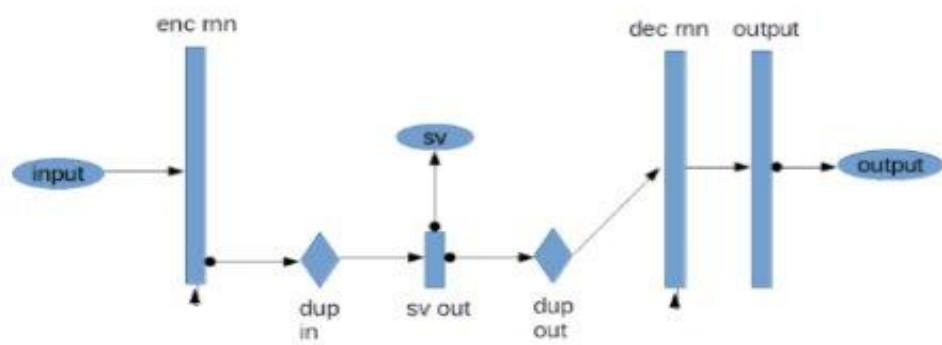
V - размер словаря из токенов.

На вход модели подаётся тензор размера (B, T, H). На выходе logits (B,T,V)



## ***Архитектура моделей.***

Прошлые подходы на основе RNN показали способность учиться, однако качество по метрикам было не очень.



## Новые подходы на основе transformer архитектуры.

За основу encoder был взят **FinBERT** — это предварительно обученная модель для анализа настроений (сентимента) финансовых текстов.

На выходе **FinBERT** эмбединги для каждого токена и **cls** токен.

Агрегируем с обучаемыми весами выход **FinBERT** с помощью **Attention pooling** слоя в эмбединг предложения. Также для эксперимента подаем только **cls**.

Decoder обучаем с нуля двумя способами.

(1) За один проход. Decoder на основе transformer\_encoder.

**sent\_emb**: эмбединг предложения. Расширяем до T, учитывая позиции (прибавляем позиционные эмбединги) и подаем в decoder для self-attention.

Предсказываем все токены на основе **sent\_emb**.

$P(\text{tgt}[1:T] \mid \text{sent\_emb})$

(2) Авторегрессионно. Decoder на основе transformer\_decoder.

**sent\_emb**: эмбединг предложения. Подаем в decoder как вход для cross-attention.

**tgt**: токены входа. Подаем в decoder для masked self-attention и cross-attention.

Предсказываем следующий токен на основе предыдущих и **sent\_emb**.

$P(\text{tgt}[N+1] \mid \text{tgt}[1:N], \text{sent\_emb})$

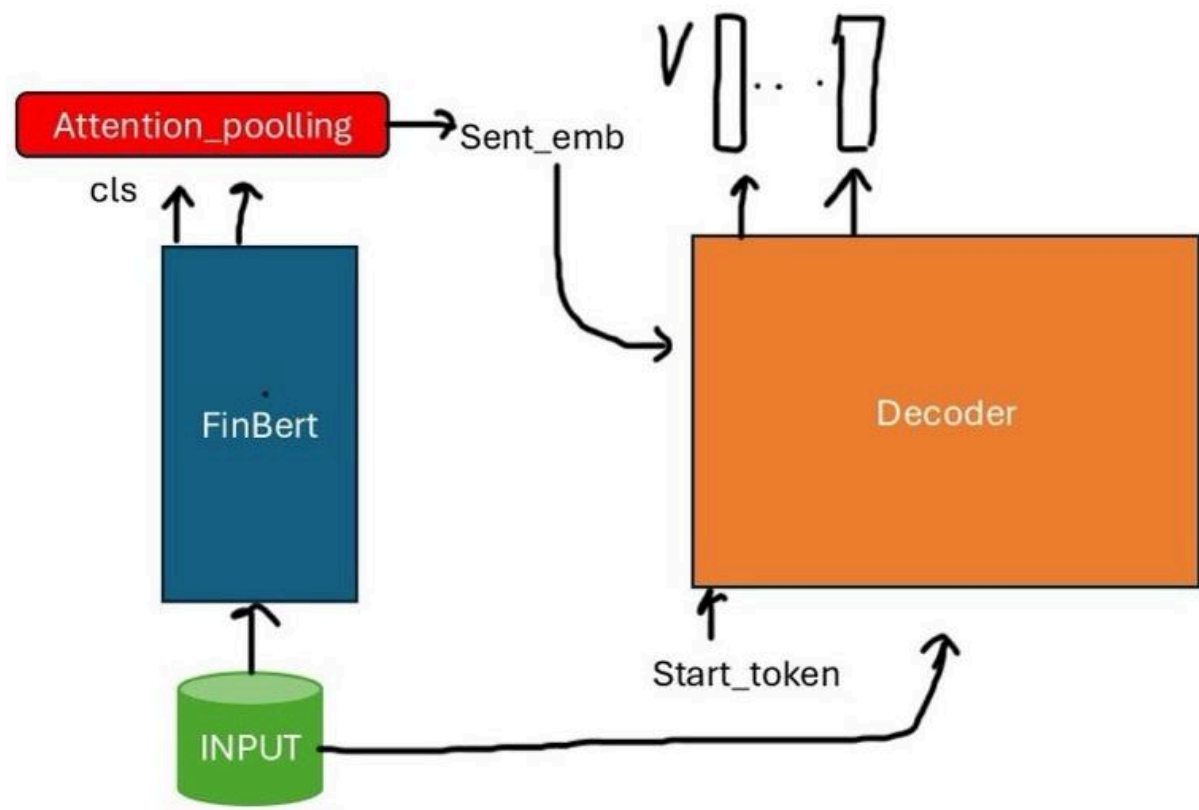


## Attention Pooling

```
class AttentionPooling(nn.Module):
    def __init__(self, hidden_dim, output_dim):
        super().__init__()
        self.attn = nn.Linear(hidden_dim, 1)
        self.proj = nn.Sequential(
            nn.Linear(hidden_dim, output_dim),
            nn.GELU(),
            nn.LayerNorm(output_dim),
        )
    def forward(self, x, mask=None):
        attn_scores = self.attn(x).squeeze(-1) # (B, T)
        if mask is not None:
            attn_scores = attn_scores.masked_fill(mask == 0, -1e9)
        attn_weights = torch.softmax(attn_scores, dim=1) # (B, T)
        return torch.sum(attn_weights.unsqueeze(-1) * self.proj(x),
dim=1) # (B, S)
```

Последний слой для получения *logits*.

```
self.fc_out = nn.Sequential(
    nn.Dropout(0.5),
    nn.Linear(sent_dim, vocab_size)
)
```



## Обучение

Loss: *cross\_entropy\_loss*

+ fine\_tuning with *sent\_sim\_loss*= $\cos(\text{sent\_emb}(\text{input}), \text{sent\_emb}(\text{decoded}))$

optimizer: AdamW (learning\_rate=0.0001, weight\_decay=0.01)

lr\_scheduler: ReduceLROnPlateau (mode='max', factor=0.1,  
min\_lr=1e-5, patience=5)

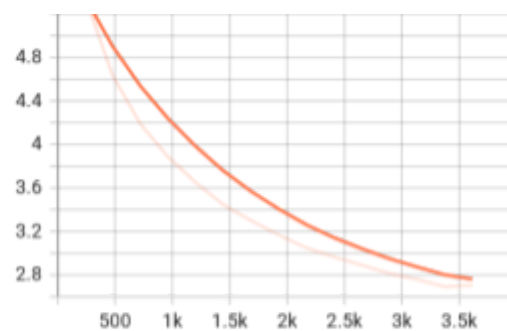
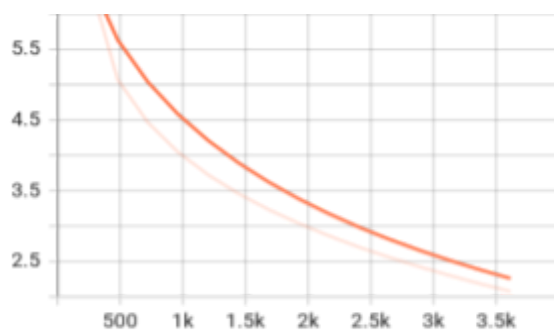
batch\_size: 128-200

## Графики

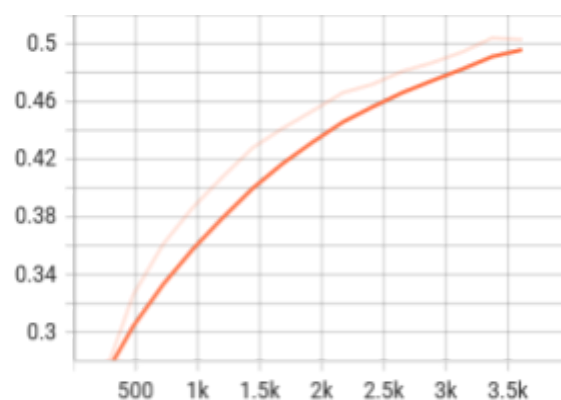
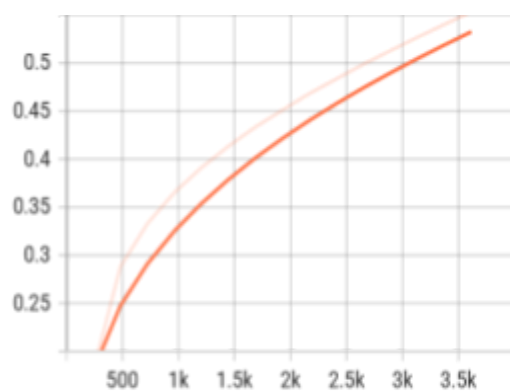
*Train*

*VALIDATION*

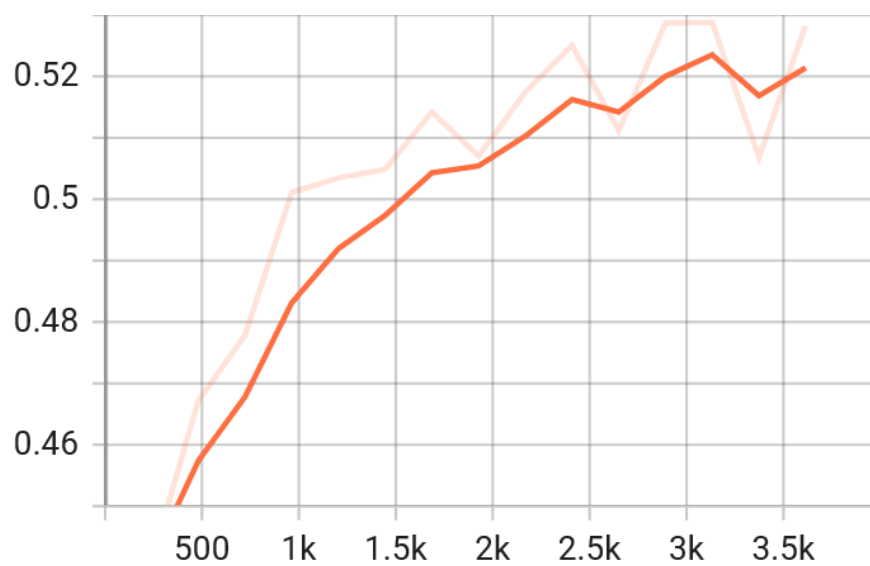
*Loss*



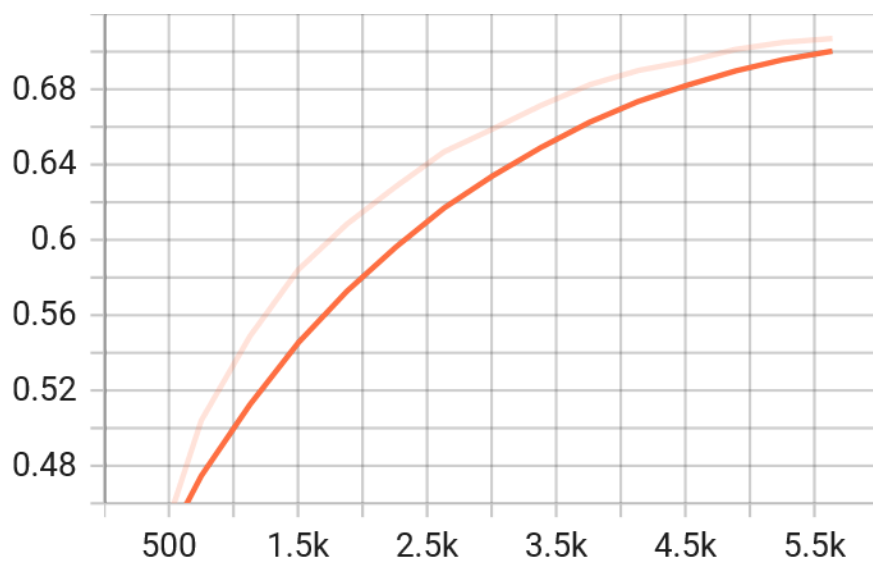
*Accuracy*



*BERT-score*



*Top-5*



## *Разные размерности*

Зависимость метрик от архитектуры и размера эмбединга предложения.

<i>model</i>	<i>dim</i>	<i>Valid_accuracy</i>	<i>BERT_score</i>
<i>За один проход</i>	<i>1568</i>	<i>0.357</i>	
<i>Аutoreгрессионно CLS</i>	<i>768</i>	<i>0.406</i>	<i>0.58</i>
<i>Аutoreгрессионно attention_pooling</i>	<i>768</i>	<i>0.502</i>	<i>0.65</i>
	<i>2560</i>	<i>0.552</i>	<i>0.675</i>

Аторегрессионные модели по качеству сильно лучше, однако медленнее.

*Attention pooling* увеличивает качеству по сравнению с *cls* эмбеделлингом. Так же этот слой позволяет менять размерность *sent\_emb*.

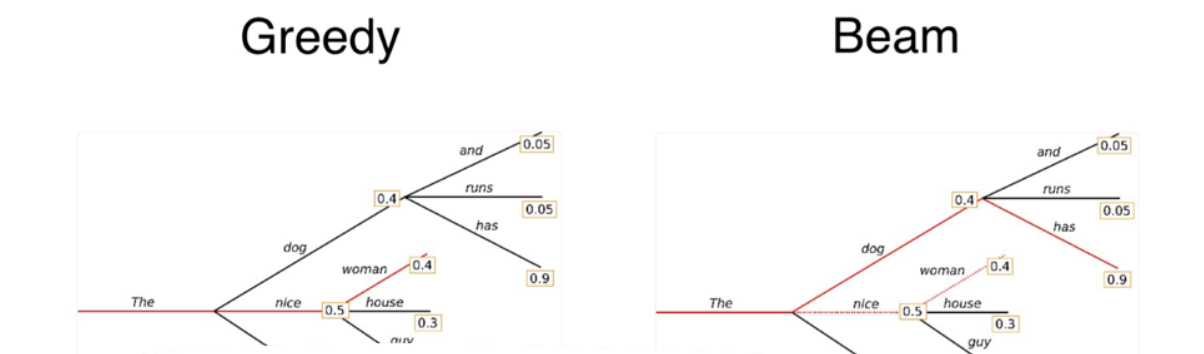
Увеличивая размерность вектора, метрики растут.

Как часто входные токены в первых k токенах по вероятности.

<i>Top-k</i>	<i>1</i>	<i>3</i>	<i>5</i>	<i>10</i>	<i>15</i>	<i>25</i>	<i>50</i>	<i>100</i>
<i>Model_768</i>	<i>0.50</i>	<i>0.63</i>	<i>0.68</i>	<i>0.73</i>	<i>0.75</i>	<i>0.78</i>	<i>0.81</i>	<i>0.84</i>
<i>Model-2560</i>	<i>0.55</i>	<i>0.68</i>	<i>0.72</i>	<i>0.77</i>	<i>0.79</i>	<i>0.81</i>	<i>0.84</i>	<i>0.86</i>

Заметно, что top-5 сильно лучше, чем top-1. С увеличением k скорость роста метрик сильно уменьшается. Возможно, что top-k для маленького k содержит синонимы или перефразы.

## ***Генерация. GREEDY vs BEAM-SEARCH***



Изменив алгоритм генерации с жадного на поиск наиболее вероятной последовательности, модель стала реже заикаться (повторять слова).

bert\_score:                      0.66 —————> 0.7

paraph\_sim:                      0.64 —————> 0.69

## *LOSS на уровне предложений*

Было замечено, что cross\_entropy loss переобучался на train с 20 эпохи. Поэтому были придуманы и реализованы loss на уровне предложений, а не токенов.

### *Cosine\_similarity loss*

Для исходного и восстановленного текста можно получить специальные эмбединги, и считать cosine\_similarity между ними.

paragraph\_model:'sentence-transformers/paraphrase-mpnet-base-v2'

Выбор пал на paragraph\_model, используемую в вычислении метрики paragraph\_sim, так как она обучена получать эмбединги всего текста близкие у похожих по смыслу. Проще использовать свою же модель, но для разных по смыслу предложений cos\_sim большой.

Сравнение.

	paragraph_model	my_model
I enjoy reading books in my spare time.		
I like reading novels during my free time.	0.918	0.792
Python is a popular programming language.	0.162	0.544

Генерировать их распределения над токенами напрямую нельзя, не будет течь градиент.

Сэмплируем из logits с помощью gumbel trick.

```
distribution = F.gumbel_softmax(logits, tau=1.0, hard=True)
```

Straight-Through Estimator (STE) `distribution=one_hot + (soft - soft.detach())`

```
gen = torch.matmul(distribution, emb_matrix) - эмбединги токенов
```

Подаем в paragraph\_model gen вместо токенов.

```
emb_inp = paragraph_model(input_text)
```

```
emb_gen = paragraph_model(inputs_embeds=gen)
```

$\text{Loss} = \text{cosine\_similarity}(\text{emb\_inp}, \text{emb\_gen})$

Этот loss оказался очень не стабильным. Предлагается использовать его как добавку к базовому loss с небольшим коэффициентом в качестве регуляризатора.

## *Loss на перефразах*

Наша задача восстановить по смыслу, можем это учитывать в loss. И напрямую обучать восстанавливать перефразы.

`gen_paraphrase_model-"ramsrigouthamg/t5-large-paraphraser-diverse-high-quality"`

С помощью `gen_paraphrase_model` создаем Paraphrase Dataset, по 5 перефраз на пример из train.

Для разнообразной генерации можно изменять параметры.

Параметры генерации: `top_k=120`, `top_p=0.9`, `temperature=1.7`

Пример из Paraphrase Dataset

*base:*

Weak sales growth hits American Eagle outlook; shares slide

*paraphrase1*

The American Eagle's outlook is marred by poor sales growth; the stock plummets.

*paraphrase2:*

The American Eagle's prospects are harmed by poor sales growth; the stock plummets.

*paraphrase3:*

The American Eagle's prospects suffer as a result of poor sales growth; the stock plummets.

$\text{loss} = \text{ce\_loss} + \min(\text{ce\_loss1}, \dots, \text{ce\_loss5}),$

$\text{ce\_loss}\{i\}$  - cross\_entropy loss для i-ой перефразы

Это поможет восстанавливать именно перефразы.

Этот loss не хуже базового, менее склонен к переобучению.



## *Изучение полученных эмбедингов*

Задача изучить наши эмбединги на локальность, чтобы при неточной их генерации, восстановление ухудшалось несильно.

Добавить к эмбедингам нормальный шум  $N(0, \sigma)$ , измерить качество восстановления.

$Bert\_score(\sigma = 0.1) \sim Bert\_score(\sigma = 0)$   $Bert\_score(\sigma = 1) < Bert\_score(\sigma = 0) * 0.9$

## *Примеры восстановления.*

<i>INPUT</i>	<i>DECODED</i>
<i>consumer giants spurn risks to chase online subscribers</i>	<i>consumer giants spurn consumer risks to online subscribers</i>
<i>united technologies wins \$ 2. 2 billion u. s. defense contract : pentagon</i>	<i>united technologies wins \$ 2. 2 billion u. s. defense contract : pentagon</i>
<i>australia watchdog says vodafone misled customers over digital purchases</i>	<i>australia watchdog says hsbc over wireless firms misled customers payments</i>
<i>bayer says monsanto likely kept files on influential people across europe</i>	<i>bayer says monsanto likely put on some countries sick people in europe</i>
<i>wall street rises on economic data, easing trade worries</i>	<i>wall street rises on trade worries, fed data data</i>

*malaysia regulator to probe if airasia  
broke rules in airbus deals*

*malaysia regulator to probe airasia in  
regulators if boeing deals*

(1-2) целиком восстанавливаются. В (3) компания Vodofone восстановлена как фирма, занимающаяся беспроводной связью. Payments почти синоним purchases. (6) boeing - airbus. Остальные примеры имеют похожие слова.

## Заключение

Основные цели проекта были достигнуты. Была показана возможность построения текстового эмбединга с возможностью частичного восстановления по смыслу с BERTScore=0.7. Были исследованы и реализованы различные архитектуры, а также придуманы специальные функции потерь, напрямую отражающие идею восстановления текста похожего по смыслу.

Дальнейшее изучение:

Изучить эмбединги на локальность.

Использовать предобученные LLM(Llama, Pythia) как более мощные decoder-архитектуры.

Попробовать для нашей задачи rnn-подобные архитектуры(Mamba, Gated Delta Net). Они основаны на состоянии всего текста, и, инициализируя модель состоянием нашего текста, можно попросить повторить, что было до. Проблема в том, что это состояние имеет размерность  $\text{dim} \times \text{dim} \gg \text{dim}$ . Это предстоит изучить.

Данный подход открывает новые подходы к совместной генерации мультимодальных моделей. Так для объяснения изменения финансовых показателей можно предсказывать один вектор. А наш декодер сгенерирует текст, сохранивший общий смысл.

Список использованных литературных источников и информационных материалов.

[\[1\] Cramming 1568 Tokens into a Single Vector and Back Again: Exploring the Limits of Embedding Space Capacity](#)

[\[2\] Exploring the Latent Capacity of LLMs for One-Step Text Generation](#)

Репозиторий проекта.

***Github:*** *sber\_practice*