

BERT

Грабовой Андрей Валериевич

Московский физико-технический институт

МФТИ, г. Долгопрудный

- Архитектура BERT.
- Головы BERT.
- Как обучается.
- Применение.
- Где же prior?

Разделяется условно на три части:

- ① Tokens Embedding.
- ② Self-Attention.
- ③ Pooler.

Пусть задано множество токенов:

$$\mathcal{I} = \{\mathbf{i} | \mathbf{i} = [0, \dots, 0, 1, 0, \dots, 0]^T\}$$

Задано множество предложений и множество типов токенов в предложении:

$$\mathcal{S} = \mathcal{I}^n, \quad \mathcal{T} = \{[0, 1]^T, [1, 0]^T\}^n$$

Отображения:

$$BM_1 : \mathbb{R}^{n \times L} \times \mathbb{R}^{2 \times L} \rightarrow \mathbb{R}^{n \times l}$$

$$BM_2 : \mathbb{R}^{n \times L} \times \mathbb{R}^{2 \times L} \rightarrow \mathbb{R}^{1 \times l}$$

Суперпозиция отображений:

$$BM_1 = BL_m \circ \dots \circ BL_1 \circ BSE$$

$$BM_2 = BP \circ BL_m \circ \dots \circ BL_1 \circ BSE$$

Функция BSE :

$$BSE : \mathbb{R}^{n \times L} \times \mathbb{R}^{n \times 2} \rightarrow \mathbb{R}^{n \times l}.$$

Для произвольной матрицы $\mathbf{s} \in \mathcal{S} \subset \mathbb{R}^{n \times L}$ и матрицы $\mathbf{t} \in \mathcal{T} \subset \mathbb{R}^{n \times 2}$ отображение BSE принимает следующий вид:

$$BSE(\mathbf{s}, \mathbf{t}) = \frac{\mathbf{h}_{bse} - \mathbf{E}\mathbf{h}_{bse}}{\sqrt{\mathbf{D}\mathbf{h}_{bse} + \varepsilon}} \cdot \mathbf{w}_1 + \mathbf{w}_2, \quad \mathbf{h}_{bse} = \mathbf{s}\mathbf{W}_1 + \mathbf{1}_{n \times n}\mathbf{W}_2 + \mathbf{t}\mathbf{W}_3,$$

где $\mathbf{W}_1 \in \mathbb{R}^{L \times l}$, $\mathbf{W}_2 \in \mathbb{R}^{n \times l}$, $\mathbf{W}_3 \in \mathbb{R}^{2 \times l}$.

Функция BSE имеет настраиваемые параметры: $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{w}_1, \mathbf{w}_2$.

Результат работы функции BSE обозначим:

$$\mathbf{h}_0 = BSE(\mathbf{s}, \mathbf{t}),$$

где $\mathbf{h} \in \mathbb{R}^{n \times l}$.

Функция BL :

$$BL : \mathbb{R}^{n \times l} \rightarrow \mathbb{R}^{n \times l}.$$

Для матрицы $\mathbf{h} \in \mathbb{R}^{n \times l}$ BL принимает следующий вид:

$$BL_q(\mathbf{h}) = \frac{\mathbf{u}\mathbf{W}_{3+6q} + \mathbf{a} - \mathbf{E}(\mathbf{u}\mathbf{W}_{3+6q} + \mathbf{a})}{\sqrt{D(\mathbf{u}\mathbf{W}_{3+6q} + \mathbf{a}) + \varepsilon}} \cdot \mathbf{w}_{3+4q} + \mathbf{w}_{4+4q},$$

$$\mathbf{u} = \sigma(\mathbf{a}\mathbf{W}_{4+6q}), \quad \mathbf{a} = \frac{\mathbf{c}\mathbf{W}_{5+6q} - \mathbf{E}\mathbf{c}\mathbf{W}_{5+6q}}{\sqrt{D\mathbf{c}\mathbf{W}_{5+6q} + \varepsilon}} \cdot \mathbf{w}_{5+4q} + \mathbf{w}_{6+4q}$$

$$\mathbf{c} = [\mathbf{c}_1, \dots, \mathbf{c}_{r_2}]$$

$$\mathbf{c}_j = \text{softmax}(\mathbf{h}\mathbf{W}_{6+6q}^j \odot \mathbf{h}\mathbf{W}_{7+6q}^j) \odot \mathbf{h}\mathbf{W}_{8+6q}^j$$

где для всех q матрицы $\mathbf{W}_{6+6q}^j, \mathbf{W}_{7+6q}^j, \mathbf{W}_{8+6q}^j \in \mathbb{R}^{l \times r}$, для всех j матрицы $\mathbf{c}_j \in \mathbb{R}^{n \times r}$, для всех q матрицы $\mathbf{W}_{5+6q}^j \in \mathbb{R}^{l \times l}$, $\mathbf{W}_{4+6q}^j \in \mathbb{R}^{l \times p}$, $\mathbf{W}_{3+6q}^j \in \mathbb{R}^{p \times l}$, матрица $\mathbf{c}, \mathbf{a} \in \mathbb{R}^{n \times l}$, матрица $\mathbf{u} \in \mathbb{R}^{n \times p}$.

Настраиваемые параметры:

$$\mathbf{W}_{3+6q}^j, \mathbf{W}_{4+6q}^j, \mathbf{W}_{5+6q}^j, \mathbf{W}_{6+6q}^j, \mathbf{W}_{7+6q}^j, \mathbf{W}_{8+6q}^j, \mathbf{w}_{3+4q}, \mathbf{w}_{4+4q}, \mathbf{w}_{5+4q}, \mathbf{w}_{6+4q}$$

Результат работы функции BL_q :

$$\forall q \in \{1, \dots, m\} \quad \mathbf{h}_q = BL_q(\mathbf{h}_{q-1}).$$

Функция BP :

$$BP : \mathbb{R}^{n \times l} \rightarrow \mathbb{R}^{n \times l}.$$

Для матрицы $\mathbf{h}_m \in \mathbb{R}^{n \times l}$ BP принимает следующий вид:

$$BP(\mathbf{h}_m) = \sigma(\mathbf{h}_m^1 \mathbf{W}_{9+6m}),$$

где \mathbf{h}_m^1 первая строка матрицы \mathbf{h}_m , а матрица $\mathbf{W}_{9+6m} \in \mathbb{R}^{l \times l}$

Функция BP имеет настраиваемые параметры \mathbf{W}_{9+6m}

Результат работы функции BP :

$$\mathbf{h} = BP(\mathbf{h}_m).$$

Вернемся к суперпозициям BM_1, BM_2 :

$$\mathbf{h} = BP(\mathbf{h}_m), \quad \mathbf{h}_q = BM(\mathbf{h}_{q-1}), \quad \mathbf{h}_0 = BSE(\mathbf{s}, \mathbf{t})$$

Получаем вектор эмбедингов слов:

$$BM_1(\mathbf{s}, \mathbf{t}) = \mathbf{h}_m,$$

Получаем вектор эмбединга предложения:

$$BM_2(\mathbf{s}, \mathbf{t}) = \mathbf{h}$$

Математическая модель Bert (multi task learning)

LM модель:

$$\mathbf{v} = \text{softmax}(\mathbf{h}_m \mathbf{W}_{LM}),$$

где $\mathbf{W}_{LM} \in \mathbb{R}^{l \times L}$, а \mathbf{v} это вероятность каждого токена.

NSP модель:

$$z = \sigma(\mathbf{h} \mathbf{W}_{NSP}),$$

где $\mathbf{W}_{NSP} \in \mathbb{R}^{l \times 1}$, а z это вероятность класса 1.

Функция ошибки:

$$L(\mathbf{S}, \mathbf{y}) = \sum_{\mathbf{s}_i, \mathbf{t}_i \in \mathbf{S}} \text{CrossEntropy}(\mathbf{v}_i, \mathbf{s}_i) + \sum_{\mathbf{s}_i, \mathbf{t}_i \in \mathbf{S}, \mathbf{y}_i \mathbf{y}} \text{CrossEntropyLoss}(z_i, \mathbf{y}_i)$$

Задача оптимизации:

$$L(\mathbf{S}, \mathbf{y}) \rightarrow \min_{\mathbf{W}_{all}}$$

Все параметры:

$$\mathbf{W}_{all} = [\mathbf{W}_{LM}, \mathbf{W}_{NSP}, \mathbf{W}_{9+6m}, \mathbf{W}_{3+6q}^j, \mathbf{W}_{4+6q}^j, \mathbf{W}_{5+6q}^j, \mathbf{W}_{6+6q}^j, \mathbf{W}_{7+6q}^j, \mathbf{W}_{8+6q}^j, \mathbf{w}_{3+4q}, \mathbf{w}_{4+4q}, \mathbf{w}_{5+4q}, \mathbf{w}_{6+4q}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{w}_1, \mathbf{w}_2]$$

```
BertEmbeddings(  
  (word_embeddings): Embedding(119547, 768, padding_idx=0)  
  (position_embeddings): Embedding(512, 768)  
  (token_type_embeddings): Embedding(2, 768)  
  (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
  (dropout): Dropout(p=0.1, inplace=False)  
)
```

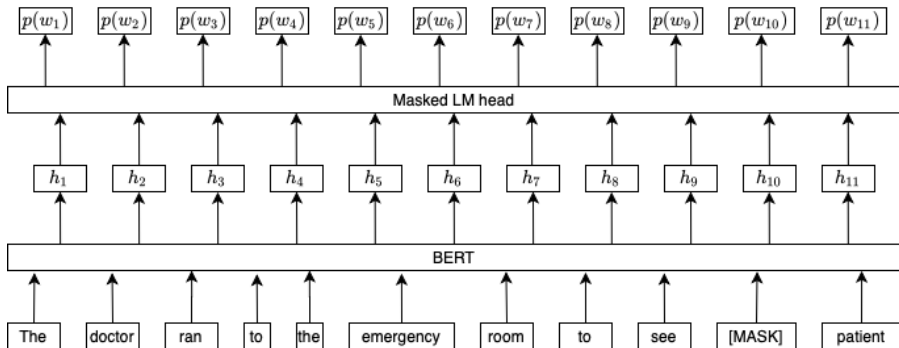
```
BertLayer(  
  (attention): BertAttention(  
    (self): BertSelfAttention(  
      (query): Linear(in_features=768, out_features=768, bias=True)  
      (key): Linear(in_features=768, out_features=768, bias=True)  
      (value): Linear(in_features=768, out_features=768, bias=True)  
      (dropout): Dropout(p=0.1, inplace=False)  
    )  
    (output): BertSelfOutput(  
      (dense): Linear(in_features=768, out_features=768, bias=True)  
      (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
      (dropout): Dropout(p=0.1, inplace=False)  
    )  
  )  
  (intermediate): BertIntermediate(  
    (dense): Linear(in_features=768, out_features=3072, bias=True)  
  )  
  (output): BertOutput(  
    (dense): Linear(in_features=3072, out_features=768, bias=True)  
    (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
    (dropout): Dropout(p=0.1, inplace=False)  
  )  
)
```

```
BertPooler(  
  (dense): Linear(in_features=768, out_features=768, bias=True)  
  (activation): Tanh()  
)
```

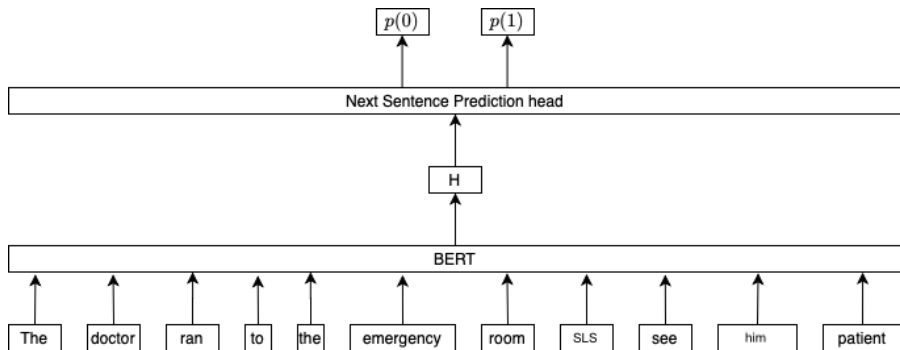
До этого мы разобрали архитектуру генерации признакового описания токенов при помощи модели BERT. Теперь перейдем к использованию сгенерированных признаков для разных задач:

- 1 Masked LM.
- 2 Next Sentence Prediction (NSP).

```
BertOnlyMLMHead(  
  (predictions): BertLMPredictionHead(  
    (transform): BertPredictionHeadTransform(  
      (dense): Linear(in_features=768, out_features=768, bias=True)  
      (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
    )  
    (decoder): Linear(in_features=768, out_features=119547, bias=True)  
  )  
)
```



```
BertOnlyNSPHead(  
  (seq_relationship): Linear(in_features=768, out_features=2, bias=True)  
)
```

Обученные параметры модели BERT применяются как генерация признакового описания токенов предложения. После чего применяются простые модели (обычно линейные), которые решают поставленную задачу машинного обучения. Примеры задач:

- 1 Words in Context.
- 2 BoolQ.
- 3 Choice of Plausible Alternatives.
- 4 и.т.д.

Полезные ссылки на датасеты:

- 1 <https://super.gluebenchmark.com/tasks>
- 2 <https://russiansuperglue.com>

- 1 Основная идея: смысл текста не зависит от решаемой задачи.
- 2 Сама модель BERT по сути является некоторым априорным заданием векторов, которые потом подстраиваются под конкретную задачу.

- ① J. Devlin, M. Chang, K. Lee, K. Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. 2018.
- ② J. Ba, J. Kiros, G. Hinton Layer Normalization. 2016.
- ③ https://huggingface.co/transformers/model_doc/bert.html#bertmodel