

Natural Language Processing with Deep Learning

Lecture 10 – Text classification 6: BERT part two

Prof. Dr. Ivan Habernal

December 15, 2023

Natural Language Processing Group
Paderborn University

We focus on Trustworthy Human Language Technologies



www.trusthlt.org

Motivation

Last time we started with the transformer and BERT, but we

- skipped some important architectural details
- were unprecise with some graphical representation
- did not talk about pre-training and fine-tuning

Let's fix that today!

After this lecture you should be able to build BERT

BERT — Encoder architecture in detail

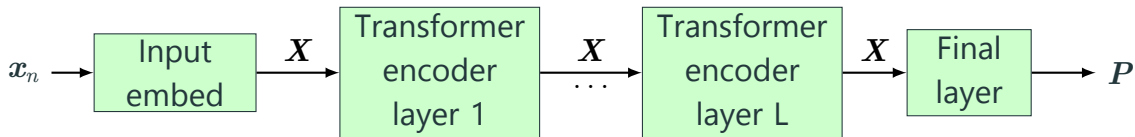
BERT — Encoder architecture in detail

Input and pre-training

Pre-training

Downstream tasks and fine-tuning

Transformer encoder (BERT)



As usual, green boxes are functions with trainable parameters

BERT (encoding-only transformer, forward pass)

1: **function** ETRANSFORMER(x ; \mathcal{W})

2: ...

Input:

x — $x \in V^*$, a sequence of token IDs

\mathcal{W} — all trainable parameters

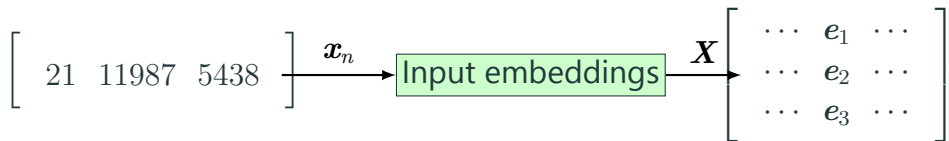
Output:

Typically an embedding vector for each input token

Or: $\mathbf{P} \in (0, 1)^{\ell_x \times N_V}$, where each row of \mathbf{P} is a distribution over the vocabulary

Input embeddings

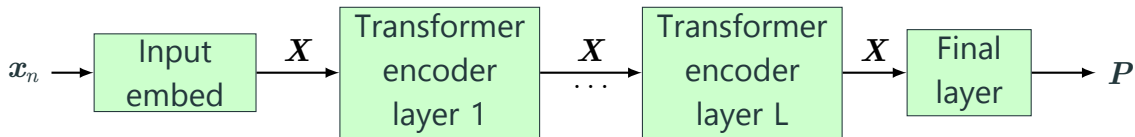
The cat sat $\mathbf{x}_n = (21 \ 11987 \ 5438)$



BERT (encoding-only transformer, forward pass)

```
1: function ETRANSFORMER( $x$ ;  $\mathcal{W}$ )  
2:    $\ell \leftarrow \text{length}(x)$   
3:   for  $t \in [\ell]$  :  $e_t \leftarrow \mathbf{W}_e[x[t], :] + \mathbf{W}_p[t, :]$        $\triangleright$  Token emb. + positional emb.  
4:    $\mathbf{X} \leftarrow \text{Stack row-wise}[e_1, e_2, \dots e_\ell]$   
5:   ...
```

Transformer encoder (BERT)

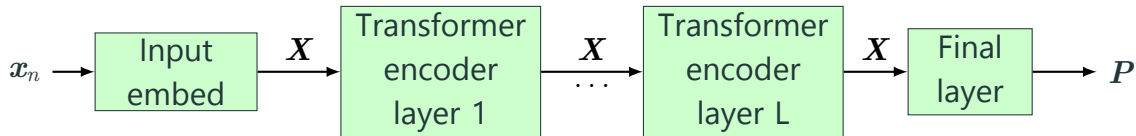


The transformer encoder layer is repeated L-times (each with different parameters)

BERT (encoding-only transformer, forward pass)

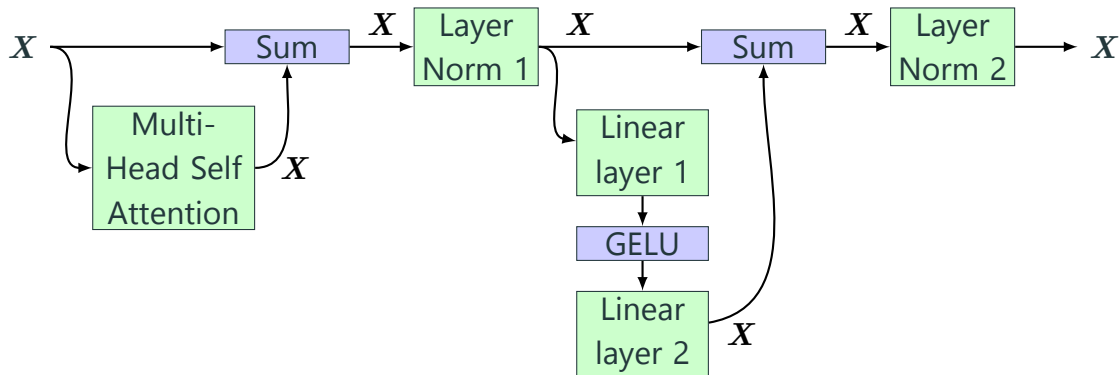
```
1: function ETRANSFORMER( $x$ ;  $\mathcal{W}$ )  
2:    $\ell \leftarrow \text{length}(x)$   
3:   for  $t \in [\ell]$  :  $e_t \leftarrow \mathbf{W}_e[x[t], :] + \mathbf{W}_p[t, :]$        $\triangleright$  Token emb. + positional emb.  
4:    $\mathbf{X} \leftarrow \text{Stack row-wise}[e_1, e_2, \dots, e_\ell]$   
5:   for  $l = 1, 2, \dots, L$  do  
6:     ...
```

Transformer encoder (BERT)



Let's look at a single transformer encoder layer

Transformer encoder layer (BERT)

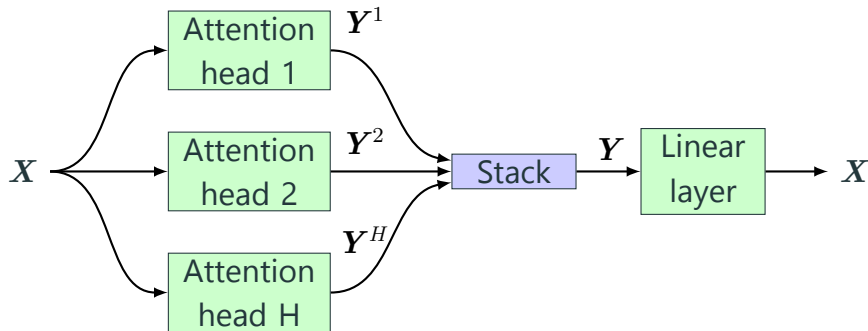


Let's focus on Multi-Head Self Attention

BERT (encoding-only transformer, forward pass)

```
1: function ETRANSFORMER( $x$ ;  $\mathcal{W}$ )
2:    $\ell \leftarrow \text{length}(x)$ 
3:   for  $t \in [\ell]$  :  $e_t \leftarrow \mathbf{W}_e[x[t], :] + \mathbf{W}_p[t, :]$        $\triangleright$  Token emb. + positional emb.
4:    $\mathbf{X} \leftarrow \text{Stack row-wise}[e_1, e_2, \dots, e_\ell]$ 
5:   for  $l = 1, 2, \dots, L$  do
6:      $\mathbf{X} \leftarrow \mathbf{X} + \text{MHATTENTION}(\mathbf{X} | \mathcal{W}_l)$        $\triangleright$  Multi-head att., residual conn
7:   ...
```

Multi-head unmasked self-attention (BERT)



Multi-head bidirectional / unmasked self-attention

Input: $\mathbf{X} \in \mathbb{R}^{\ell_x \times d_x}$, vector representations of the sequence of length ℓ_x

Output: $\tilde{\mathbf{V}} \in \mathbb{R}^{\ell_x \times d_{\text{out}}}$, updated vector representations of tokens in \mathbf{X}

Hyper-param: H , number of attention heads

Params for each $h \in [H]$: \mathcal{W}_{qkv}^h :

- $\mathbf{W}_q^h, \mathbf{W}_k^h \in \mathbb{R}^{d_x \times d_{\text{attn}}}$, $\mathbf{b}_q^h, \mathbf{b}_k^h \in \mathbb{R}^{d_{\text{attn}}}$, $\mathbf{W}_v \in \mathbb{R}^{d_x \times d_{\text{mid}}}$, $\mathbf{b}_v \in \mathbb{R}^{d_{\text{mid}}}$
- $\mathbf{W}_o \in \mathbb{R}^{H \cdot d_{\text{mid}} \times d_{\text{out}}}$, $\mathbf{b}_o \in \mathbb{R}^{d_{\text{out}}}$

1: **function** MHATTENTION($\mathbf{X}; \mathcal{W}$)

2: **for** $h \in [H]$ **do**

3: $\mathbf{Y}^h \leftarrow \text{ATTENTION}(\mathbf{X}; \mathcal{W}_{qkv}^h)$

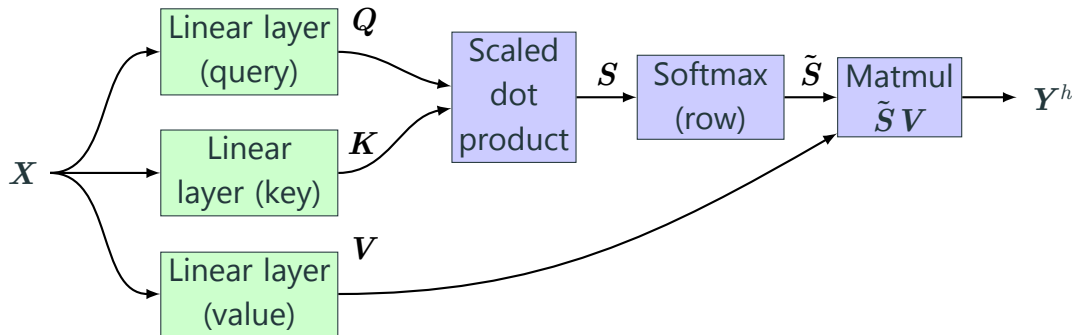
▷ $\mathbf{Y}^h \in \mathbb{R}^{\ell_x \times d_{\text{mid}}}$

4: $\mathbf{Y} \leftarrow [\mathbf{Y}^1; \mathbf{Y}^2; \dots; \mathbf{Y}^H]$

▷ $\mathbf{Y} \in \mathbb{R}^{\ell_x \times H \cdot d_{\text{mid}}}$

5: **return** $\tilde{\mathbf{V}} = \mathbf{Y}\mathbf{W}_o + \mathbf{b}_o$

Single unmasked self-attention head (BERT)



Bidirectional / unmasked self-attention (recap from last lecture)

Input: $\mathbf{X} \in \mathbb{R}^{\ell_x \times d_x}$, vector representations of the sequence of length ℓ_x

Output: $\tilde{\mathbf{V}} \in \mathbb{R}^{\ell_x \times d_{\text{out}}}$, updated vector representations of tokens in \mathbf{X}

Params \mathcal{W}_{qkv} : $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{d_x \times d_{\text{attn}}}$, $\mathbf{b}_q, \mathbf{b}_k \in \mathbb{R}^{d_{\text{attn}}}$, $\mathbf{W}_v \in \mathbb{R}^{d_x \times d_{\text{out}}}$, $\mathbf{b}_v \in \mathbb{R}^{d_{\text{out}}}$

1: **function** ATTENTION($\mathbf{X}; \mathcal{W}_{qkv}$)

2: $\mathbf{Q} \leftarrow \mathbf{X} \mathbf{W}_q +_{(\text{rows})} \mathbf{b}_q$

▷ Query $\in \mathbb{R}^{\ell_x \times d_{\text{attn}}}$

3: $\mathbf{K} \leftarrow \mathbf{X} \mathbf{W}_k +_{(\text{rows})} \mathbf{b}_k$

▷ Key $\in \mathbb{R}^{\ell_x \times d_{\text{attn}}}$

4: $\mathbf{V} \leftarrow \mathbf{X} \mathbf{W}_v +_{(\text{rows})} \mathbf{b}_v$

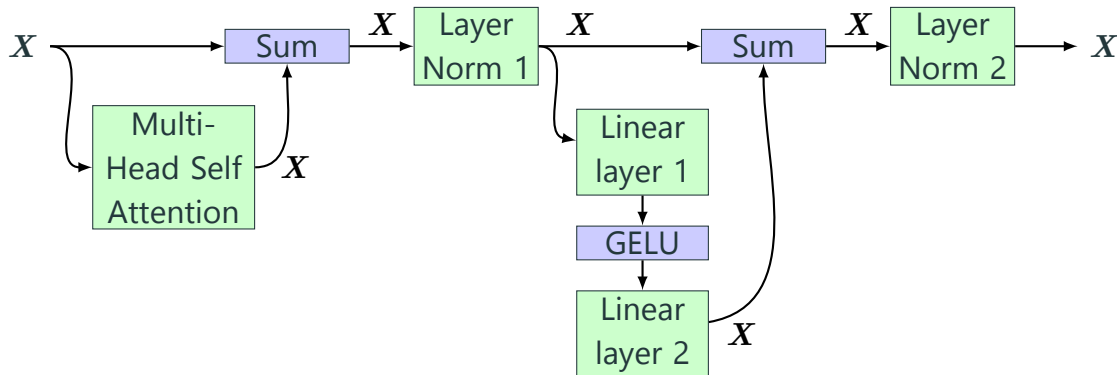
▷ Value $\in \mathbb{R}^{\ell_x \times d_{\text{out}}}$

5: $\mathbf{S} \leftarrow \frac{1}{\sqrt{d_{\text{attn}}}} (\mathbf{Q} \mathbf{K}^\top)$

▷ Scaled score $\in \mathbb{R}^{\ell_x \times \ell_x}$

6: **return** $\tilde{\mathbf{V}} = \text{softmax}_{\text{row}}(\mathbf{S}) \mathbf{V}$

Transformer encoder layer (BERT)



Let's add Layer Normalization and GELU

Simplifying notation: Perform LAYERNORM on each row

Recall: LayerNorm

Input: $e \in \mathbb{R}^d$ (output of a layer), Output: $\hat{e} \in \mathbb{R}^d$

Params: $\gamma, \beta \in \mathbb{R}^d$, trainable element-wise scale and offset

- 1: **function** LAYERNORM($e | \gamma, \beta$)
- 2: $m \leftarrow \frac{1}{d} \sum_{i=1}^d e[i]$ ▷ 'Sample mean' of e
- 3: $v \leftarrow \frac{1}{d} \sum_{i=1}^d (e[i] - m)^2$ ▷ 'Sample variance' of e
- 4: **return** $\hat{e} = \frac{e-m}{\sqrt{v}} \odot \gamma + \beta$ ▷ Offset and scale

- 1: **function** LAYERNORMEACHROW($X \in \mathbb{R}^{m \times n} | \gamma, \beta$)
- 2: **for** $t \in [m]$ **do**
- 3: $X[t, :] \leftarrow \text{LAYERNORM}(X[t, :] | \gamma, \beta)$
- 4: **return** X

GELU — Gaussian Error Linear Units

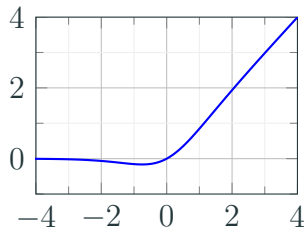
Recall: CDF $\Phi(x)$ of standard normal $X \sim \mathcal{N}(0; 1)$

$$\Phi(x) = \Pr(X \leq x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{t^2}{2}\right) dt$$

D. Hendrycks and K. Gimpel (2016).
Gaussian Error Linear Units (GELUs).
arXiv: 1606.08415

For vectors $x \in \mathbb{R}^n$, $\text{GELU}(x)$ is applied element-wise

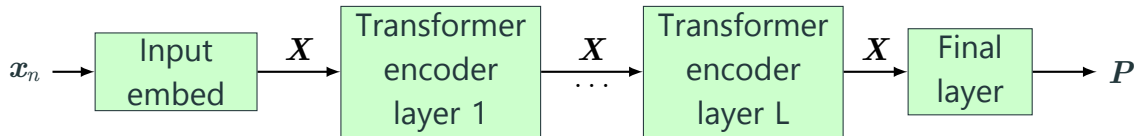
$$\begin{aligned}\text{GELU}(x) &= x \cdot \Phi(x) \\ &\approx x \cdot \sigma(1.702x) \quad (\text{if speed} > \text{exactness})\end{aligned}$$



BERT (encoding-only transformer, forward pass)

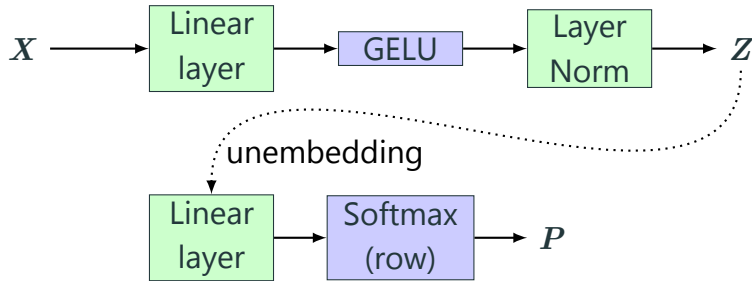
```
1: function ETRANSFORMER( $x$ ;  $\mathcal{W}$ )
2:    $\ell \leftarrow \text{length}(x)$ 
3:   for  $t \in [\ell]$  :  $e_t \leftarrow \mathbf{W}_e[x[t], :] + \mathbf{W}_p[t, :]$  ▷ Token emb. + positional emb.
4:    $\mathbf{X} \leftarrow \text{Stack row-wise}[e_1, e_2, \dots, e_\ell]$ 
5:   for  $l = 1, 2, \dots, L$  do
6:      $\mathbf{X} \leftarrow \mathbf{X} + \text{MHATTENTION}(\mathbf{X} | \mathcal{W}_l)$  ▷ Multi-head att., residual conn
7:      $\mathbf{X} \leftarrow \text{LAYERNORMPERROW}(\mathbf{X} | \gamma_l^1, \beta_l^1)$ 
8:      $\mathbf{X} \leftarrow \mathbf{X} + \left( \text{GELU}(\mathbf{X} \mathbf{W}_l^{\text{mlp1}} +_{(\text{row})} \mathbf{b}_l^{\text{mlp1}}) \mathbf{W}_l^{\text{mlp2}} +_{(\text{row})} \mathbf{b}_l^{\text{mlp2}} \right)$  ▷ MLP
9:      $\mathbf{X} \leftarrow \text{LAYERNORMPERROW}(\mathbf{X} | \gamma_l^2, \beta_l^2)$ 
10:  ...
```

Transformer encoder (BERT)



Let's look at the final layers

Final layer (BERT)



BERT (encoding-only transformer, forward pass)

```
1: function ETRANSFORMER( $x; \mathcal{W}$ )
2:    $\ell \leftarrow \text{length}(x)$ 
3:   for  $t \in [\ell] : e_t \leftarrow \mathbf{W}_e[x[t], :] + \mathbf{W}_p[t, :]$  ▷ Token emb. + positional emb.
4:    $\mathbf{X} \leftarrow \text{Stack row-wise}[e_1, e_2, \dots, e_\ell]$ 
5:   for  $l = 1, 2, \dots, L$  do
6:      $\mathbf{X} \leftarrow \mathbf{X} + \text{MHATTENTION}(\mathbf{X} | \mathcal{W}_l)$  ▷ Multi-head att., residual conn
7:      $\mathbf{X} \leftarrow \text{LAYERNORMPERROW}(\mathbf{X} | \gamma_l^1, \beta_l^1)$ 
8:      $\mathbf{X} \leftarrow \mathbf{X} + \left( \text{GELU}(\mathbf{X} \mathbf{W}_l^{\text{mlp1}} +_{(\text{row})} \mathbf{b}_l^{\text{mlp1}}) \mathbf{W}_l^{\text{mlp2}} +_{(\text{row})} \mathbf{b}_l^{\text{mlp2}} \right)$  ▷ MLP
9:      $\mathbf{X} \leftarrow \text{LAYERNORMPERROW}(\mathbf{X} | \gamma_l^2, \beta_l^2)$ 
10:     $\mathbf{X} \leftarrow \text{GELU}(\mathbf{X} \mathbf{W}_f +_{(\text{row})} \mathbf{b}_f)$ 
11:     $\mathbf{X} \leftarrow \text{LAYERNORMPERROW}(\mathbf{X} | \gamma_l, \beta_l)$ 
12:    return  $\mathbf{P} = \text{softmax}(\mathbf{X} \mathbf{W}_u)$  ▷ Project to vocab., probabilities
```

BERT parameters and hyperparameters

Hyperparameters: $\ell_{\max}, L, H, d_e, d_{\text{mlp}}, d_f \in \mathbb{N}$

Parameters:

$\mathbf{W}_e \in \mathbb{R}^{N_V \times d_e}$, $\mathbf{W}_p \in \mathbb{R}^{\ell_{\max} \times d_e}$, the token and positional embedding matrices

For $l \in [L]$: \mathcal{W}_l , multi-head attention parameters for layer l :

- $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2$, two sets of layer-norm parameters
- $\mathbf{W}_l^{\text{mlp1}} \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$, $\mathbf{b}_l^{\text{mlp1}} \in \mathbb{R}^{d_{\text{mlp}}}$
- $\mathbf{W}_l^{\text{mlp2}} \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$, $\mathbf{b}_l^{\text{mlp2}} \in \mathbb{R}^{d_e}$

$\mathbf{W}_f \in \mathbb{R}^{d_e \times d_f}$, $\mathbf{b}_f \in \mathbb{R}^{d_f}$, $\gamma, \beta \in \mathbb{R}^{d_f}$, the final linear projection and layer-norm parameters.

$\mathbf{W}_u \in \mathbb{R}^{d_e \times N_V}$, the unembedding matrix

Input and pre-training

BERT — Encoder architecture in detail

Input and pre-training

Pre-training

Downstream tasks and fine-tuning

BERT: Tokenization

Tokenizing into a multilingual WordPiece inventory

- Recall that WordPiece units are sub-word units
- 30,000 WordPiece units (newer models 110k units, 100 languages)

Implications: BERT can "consume" any language

BERT: Input representation

- Each WordPiece token from the input is represented by a **WordPiece embedding** (randomly initialized)
- Each position from the input is associated with a **positional embedding** (also randomly initialized)
- Input length limited to **512** WordPiece tokens, using `<PAD>`ding
- Special tokens
 - The first token is always a special token **[CLS]**
 - If the task involves two sentences (e.g., NLI), these two sentences are separated by a special token **[SEP]**; also special two **segment position embeddings**

BERT: Input representation summary

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\#ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

Pre-training

BERT — Encoder architecture in detail

Input and pre-training

Pre-training

Downstream tasks and fine-tuning

BERT: Self-supervised multi-task pre-training

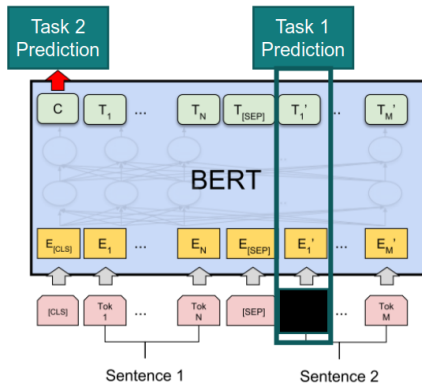
Prepare two auxiliary tasks that need no labeled data

Task 1: Cloze-test task

- Predict the masked WordPiece unit (multi-class, 30k classes)

Task 2: Consecutive segment prediction

- Did the second text segment appeared after the first segment? (binary)



BERT: Pre-training data generation

Take the entire Wikipedia (in 100 languages; 2,5 billion words)

To generate a single training instance, sample two segments (max combined length 512 WordPiece tokens)

- For Task 2, replace the second segment randomly in 50% (negative samples)
- For Task 1, choose random 15% of the tokens, and in 80% replace with a [MASK]

BERT: Pre-training data – Simplified example

Input = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

- <PAD>ding is missing
- The actual segments are longer and not necessarily sentences (just spans)
- The WordPiece tokens match full words here

BERT: pre-training by masked language modeling

```
1: function ETRAINING( $\{\mathbf{x}_n\}_{n=1}^{N_{\text{data}}}$  seqs,  $\theta$  init. params;  $p_{\text{mask}} \in (0, 1)$ ,  $N_{\text{epochs}}$ ,  $\eta$ )
2:   for  $i \in [N_{\text{epochs}}]$  do
3:     for  $n \in [N_{\text{data}}]$  do
4:        $\ell \leftarrow \text{length}(\mathbf{x}_n)$ 
5:       for  $t \in [\ell]$  do
6:          $\tilde{\mathbf{x}}_n[t] \leftarrow \text{<mask\_token> with prob. } p_{\text{mask}}, \text{ otherwise } \mathbf{x}_n[t]$ 
7:        $\tilde{T} \leftarrow \{t \in [\ell] : \tilde{\mathbf{x}}_n[t] = \text{<mask\_token>}\} \triangleright \text{Indices of masked tokens}$ 
8:        $\mathbf{P}_{\theta} \leftarrow \text{ETRANSFORMER}(\tilde{\mathbf{x}}_n | \theta)$ 
9:        $\text{loss}_{\theta} \leftarrow - \sum_{t \in \tilde{T}} \log \mathbf{P}_{\theta}[t, \mathbf{x}_n[t]]$ 
10:       $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}_{\theta}$ 
11:   return  $\theta$ 
```

Simple example explaining lines 6–7 (masking)

$(\text{The} \quad \text{cat} \quad \text{sat}) \rightarrow \mathbf{x}_n = (21 \quad 11987 \quad 5438)$ (Indices in V)

Random masking (index of `<mask_token>` = 50001):

1. For $t = 1$, the random outcome is "mask"
2. For $t = 2$, the random outcome is "keep"
3. For $t = 3$, the random outcome is "mask"

$$\tilde{\mathbf{x}}_n = (50001 \quad 11987 \quad 50001), \tilde{T} = \{1, 3\}$$

Explaining line 9 (negative log likelihood)

$$\left(\text{The} \quad \text{cat} \quad \text{sat}\right) \rightarrow \mathbf{x}_n = \begin{pmatrix} 21 & 11987 & 5438 \end{pmatrix}, \tilde{\mathbf{x}}_n = \begin{pmatrix} 50001 & 11987 & 50001 \end{pmatrix}, \tilde{T} = \{1, 3\}$$

$$\mathbf{P}_\theta \leftarrow \text{ETRANSFORMER}(\tilde{\mathbf{x}}_n | \theta)$$

$$\mathbf{P}_\theta = \begin{pmatrix} 0.001 & 0.0007 & \dots & 0.0003 \\ 0.0013 & 0.0065 & \dots & 0.0001 \\ 0.079 & 0.015 & \dots & 0.0001 \end{pmatrix}$$

$\mathbf{P}_\theta \in (0, 1)^{\ell_x \times N_v}$, where each row of \mathbf{P} is a distribution over the vocabulary

Explaining line 9 (negative log likelihood), $t = 1$

$$\mathbf{x}_n = (21, 11987, 5438), \tilde{\mathbf{x}}_n = (50001, 11987, 50001), \tilde{T} = \{1, 3\}$$

$$\mathbf{P}_\theta = \begin{pmatrix} 0.001 & \dots & 0.0041_{21} & \dots & 0.0003 \\ \vdots & & & & \end{pmatrix}$$

For $t = 1$, the model should learn to predict "The" (index 21)

$$\text{Gold: } \mathbf{y} = (0, 0, \dots, 1_{21}, \dots, 0) \in \mathbb{R}^{N_v}$$

$$\text{Pred: } \hat{\mathbf{y}} = \mathbf{P}_\theta[1, :] = (0.001, \dots, 0.0041_{21}, \dots, 0.0003) \in \mathbb{R}^{N_v}$$

Categorical cross entropy loss (Lec. 4)

$$\begin{aligned} L(\hat{\mathbf{y}}, \mathbf{y}) &:= - \sum_{k=1}^K \mathbf{y}_{[k]} \log(\hat{\mathbf{y}}_{[k]}) \\ &= -1 \cdot \log(\hat{\mathbf{y}}[21]) = -\log(\mathbf{P}_\theta[1, 21]) \\ &= -\log(\mathbf{P}_\theta[1, \mathbf{x}_n[1]]) = -\log(\mathbf{P}_\theta[t, \mathbf{x}_n[t]]) \end{aligned}$$

Explaining line 9 (negative log likelihood), $t = 3$

$$\mathbf{x}_n = (21, 11987, 5438), \tilde{\mathbf{x}}_n = (50001, 11987, 50001), \tilde{T} = \{1, 3\}$$

$$\mathbf{P}_\theta = \begin{pmatrix} \vdots & \dots & \dots \end{pmatrix}$$

For $t = 3$, the model should learn to predict "sat" (id 5438)

Categorical cross entropy loss

$$\begin{aligned} L(\hat{\mathbf{y}}, \mathbf{y}) &:= - \sum_{k=1}^K \mathbf{y}_{[k]} \log(\hat{\mathbf{y}}_{[k]}) \\ &= -1 \cdot \log(\hat{\mathbf{y}}[5438]) = -\log(\mathbf{P}_\theta[3, 5438]) = -\log(\mathbf{P}_\theta[t, \mathbf{x}_n[t]]) \end{aligned}$$

Sum over all masked token positions in \tilde{T} gives us line 9:

$$\text{loss}_\theta \leftarrow - \sum_{t \in \tilde{T}} \log \mathbf{P}_\theta[t, \mathbf{x}_n[t]]$$

Downstream tasks and fine-tuning

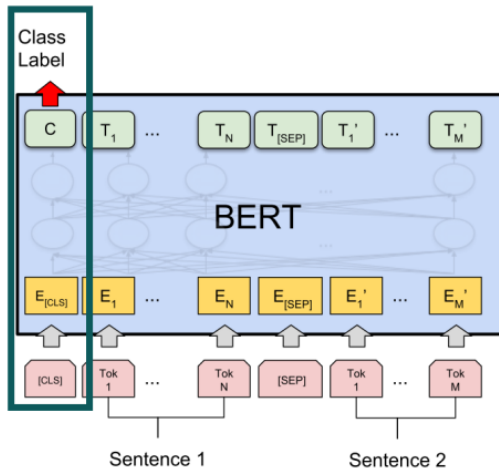
BERT — Encoder architecture in detail

Input and pre-training

Pre-training

Downstream tasks and fine-tuning

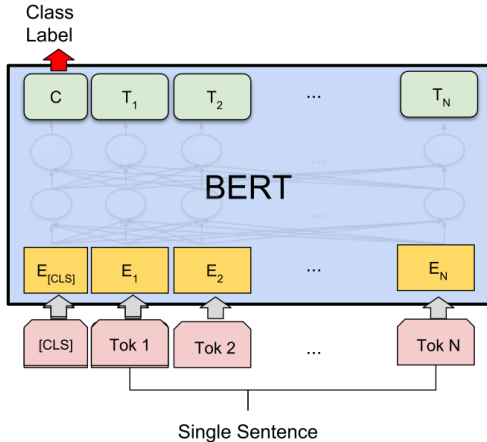
BERT: Representing various NLP tasks



That explains the special [CLS] token at sequence start

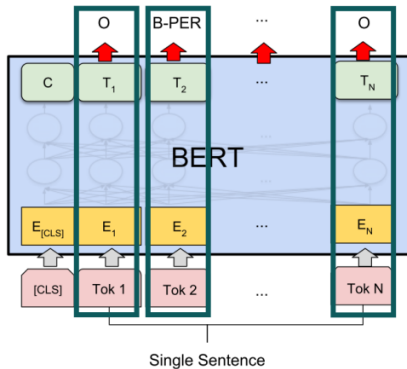
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

BERT: Representing various NLP tasks



(b) Single Sentence Classification Tasks:
SST-2, CoLA

BERT: Representing various NLP tasks



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Not conditioned on surrounding predictions

BERT pre-training time

Pretraining BERT took originally 4 days on 64 TPUs¹

Once pre-trained, transfer and “fine-tune” on your small-data task and get competitive results

P. Izsak, M. Berchansky, and O. Levy (2021). “**How to Train BERT with an Academic Budget**”. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, pp. 10644–10652

¹Can be done more efficiently, see, e.g., Izsak, Berchansky, and Levy (2021)

Recap

BERT stays on the shoulders of many clever concepts and techniques, mastered into a single model

What do we know about how BERT works?

“BERTology has clearly come a long way, but it is fair to say we still have more questions than answers about how BERT works.” — Rogers, Kovaleva, and Rumshisky (2020)²

A. Rogers, O. Kovaleva, and A. Rumshisky (2020). **“A Primer in BERTology: What We Know About How BERT Works”**.

In: *Transactions of the Association for Computational Linguistics* 8, pp. 842–866

²Highly recommended reading!

License and credits

Licensed under Creative Commons
Attribution-ShareAlike 4.0 International
(CC BY-SA 4.0)



Credits

Ivan Habernal

Content from ACL Anthology papers licensed under CC-BY
<https://www.aclweb.org/anthology>