

# 学生・若手研究者のための BERTワークショップ

RIKEN AIP 中山功太

# 本ワークショップの目的

---

## 目的

BERTについての知識を得る。

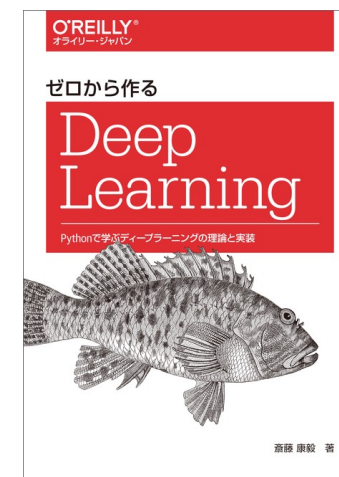
BERTを用いて文書分類タスクを試す(Pytorch実装)。

## より発展的な理解のために

本ワークショップでは深層学習の知識がなくともデモの実行は可能ですが、より深く理解したい場合は、「[ゼロから作るDeep Learning](#)」を読んだり、以下の深層学習や誤差逆伝播を数学的に解説した動画の視聴をお勧めします。

<https://youtu.be/xzzTYL90M8s>

<https://youtu.be/0itH0iDO8BE>



# 深層学習についての説明

# ザクっとわかる深層学習

深層学習は一般的に予測と正答ラベルの間の誤差を最小化することでモデルを学習する。

入力 吾輩は猫である。 名前はまだ無い。

深層学習モデル e.g. BERT  
(計算グラフ)

予測

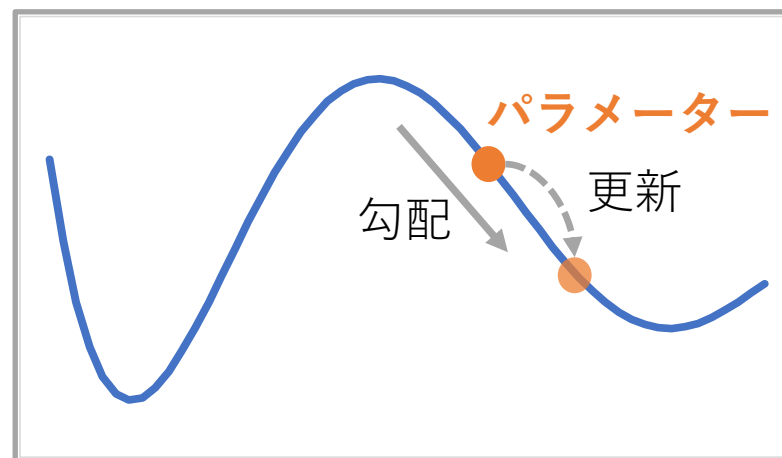
夏目漱石	太宰治
0.3	0.7

誤差

正答ラベル

夏目漱石	太宰治
1	0

誤差逆伝播+パラメーター更新  
深層学習モデルの誤差を微分し  
誤差が小さくなる方向にパラメーターを更新する。  
⇒オプティマイザーが勝手にやってくれる。



# バッチ入力

入力

吾輩は猫である。名前はまだ無い。
私は、その男の写真を三葉、見たことがある。
親譲の無鉄砲で小供の時から損ばかりしている。
...

} バッチサイズ

深層学習モデル e.g. BERT  
(計算グラフ)

バッチ入力とは  
複数の入力をもとに誤差を最小化することで学習  
を安定化させるテクニック。

予測

夏目漱石	太宰治
0.3	0.7
0.8	0.2
0.4	0.6
...	...

誤差  
(基本的には  
誤差の平均)

夏目漱石	太宰治
1	0
0	1
1	0
...	...

正答  
ラベル

# BERTについての説明

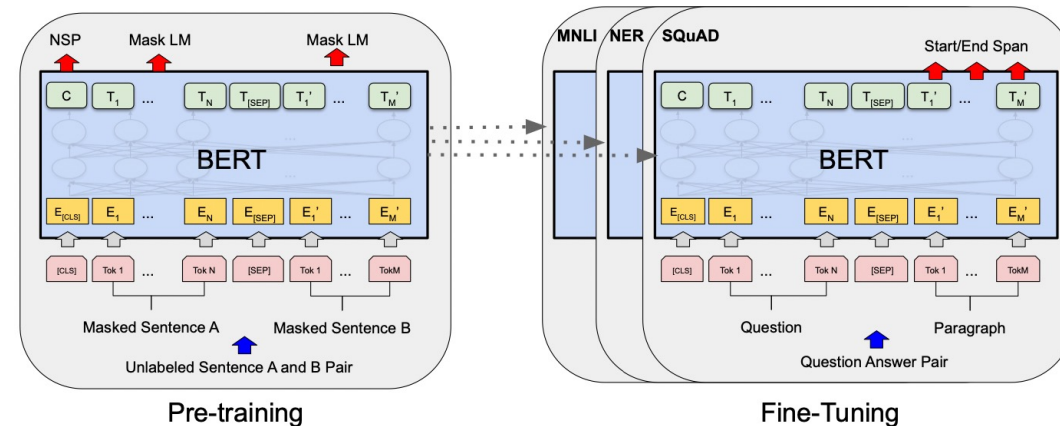
# BERTとは

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (2018)  
で提案された深層学習モデル。2022/7/22現在 44471引用!!!

簡単に、**さまざまな自然言語処理タスクに適用可能**であり、実際に多くの研究で用いられている。

以下の順序でBERTについて説明していく。

1. 入力文章の前処理
2. BERTの構造
  - A) 埋め込み
  - B) 自己注意機構(Self-Attention Mechanism)
  - C) 残差接続(Residual Connection)
  - D) 層正規化(Layer Normalization)
3. 事前学習
  - A) マスク言語モデリング(MLM : Masked Language Modeling)
4. 分類タスクへの適用法



# 1. 入力文章の前処理

BERTでは自然文を分割した後、ID配列に変換し入力とする。

入力文章      本日は晴天なり。



**単語分割** ※正確には単語より細かいサブワードに分割する ※##は1つ前のサブワードと結合して1つの単語を意味することを示す。

トークン      [本, ##日, は, 晴, ##天, なり, 。]



**特殊トークン結合**

特殊トークン

[CLS] : 文章の先頭を示す

[SEP] : 文章の終わり(もしくは区切り)

[PAD] : 入力長の調節用

[UNK] : 語彙リストに含まれない  
トークン

[MASK] : のちに説明

トークン      [[CLS], 本, ##日, は, 晴, ##天, なり, 。, [SEP], [PAD], ..., [PAD]]



**辞書を用いてID化**

辞書は事前に大規模データから単語頻度を用いて作成する

辞書のサイズは32000程度が用いられることが多い。

辞書に含まれなかったトークンは事前に[UNK]トークンに置き換える。

トークンID    [ 2, 108, 28486, 9, 4798, 28849, 297, 8, 3, 0, ..., 0 ]



# 1. 入力文章の前処理

BERTは各トークンIDに対応する位置ID、文章ID、注意マスクも入力に取る。

トークンID	2	108	28486	9	4798	28849	297	8	3	0	...	0
位置ID	0	1	2	3	4	5	6	7	8	9	...	N-1
文章ID	0	0	0	0	0	0	0	0	0	0	...	0
注意マスク	1	1	1	1	1	1	1	1	1	0	...	0

各入力トークンの位置をモデルが把握するために使用する。0~入力長N-1までの連番をとる。

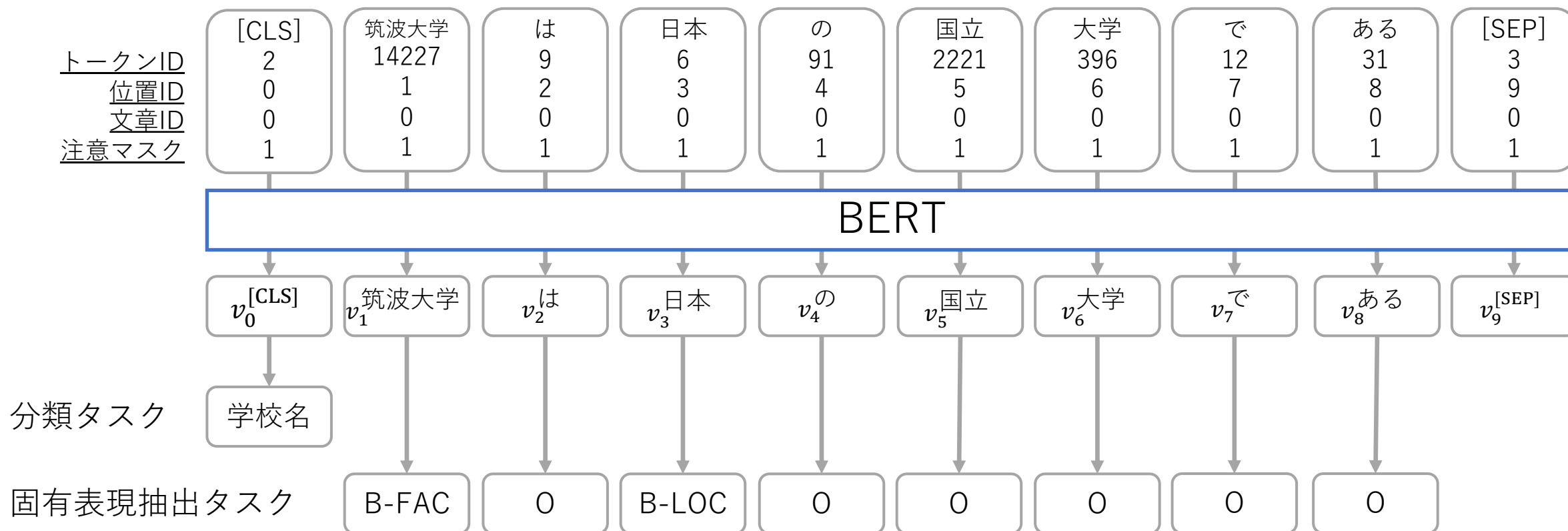
複数入力文章がある場合にモデルが文章のスパンを把握するために使用する。分類等単一文章が入力の場合は0のみをとる。

意味を持たない[PAD]トークンの位置をモデルが把握するために使用する。[PAD]の場合0、それ以外の場合1をとる。

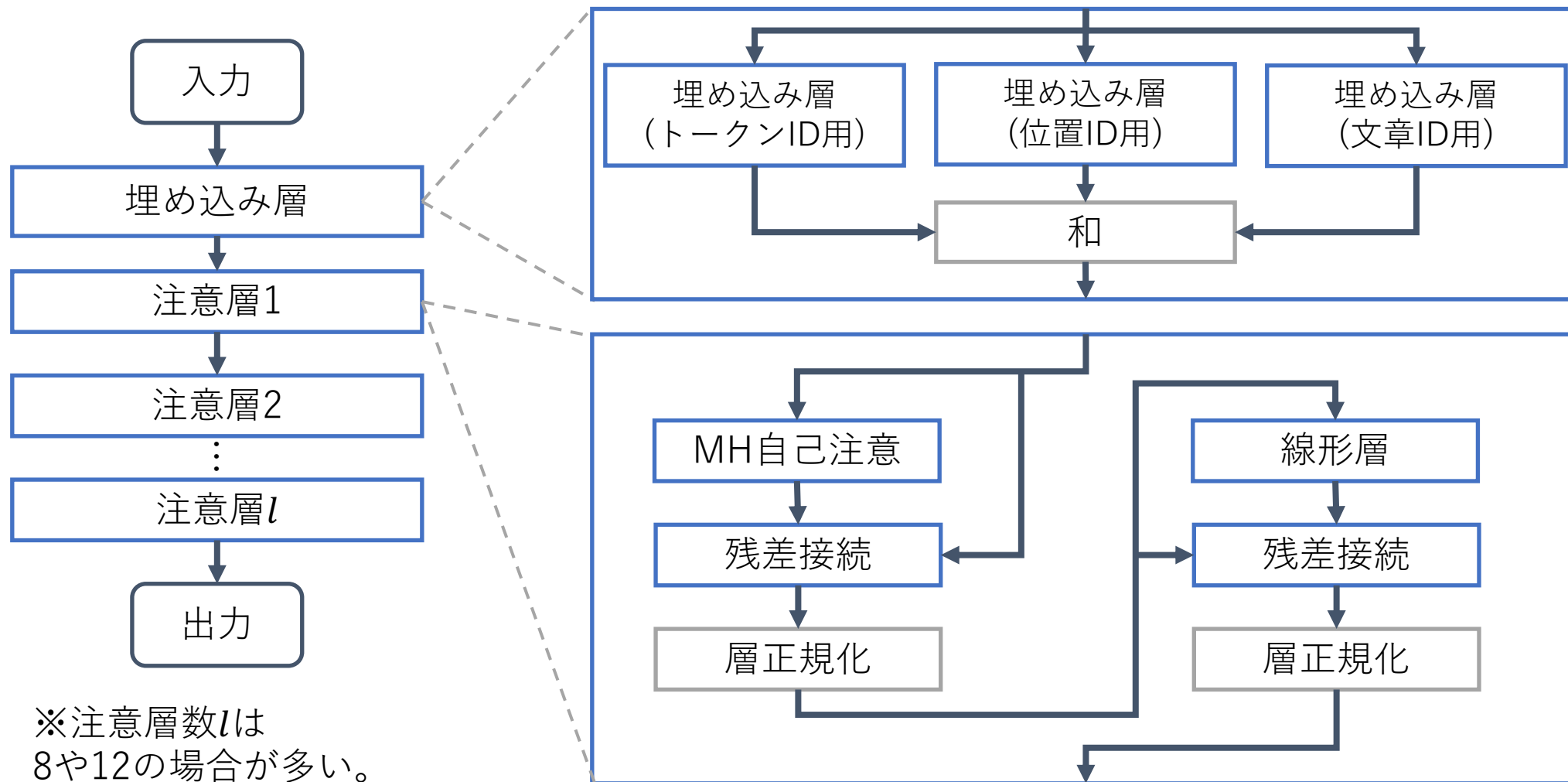
入力長N

## 2. BERTの構造

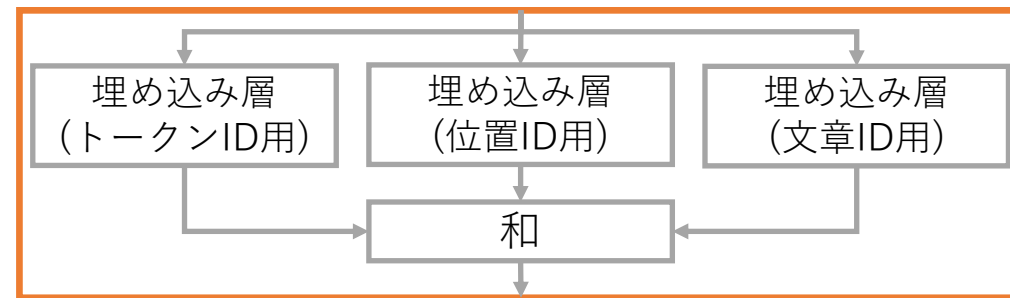
BERTは各トークンごとにベクトルを出力する。  
多くの場合この出力ベクトルをもとにさまざまなタスクを解くことになる。



## 2. BERTの構造



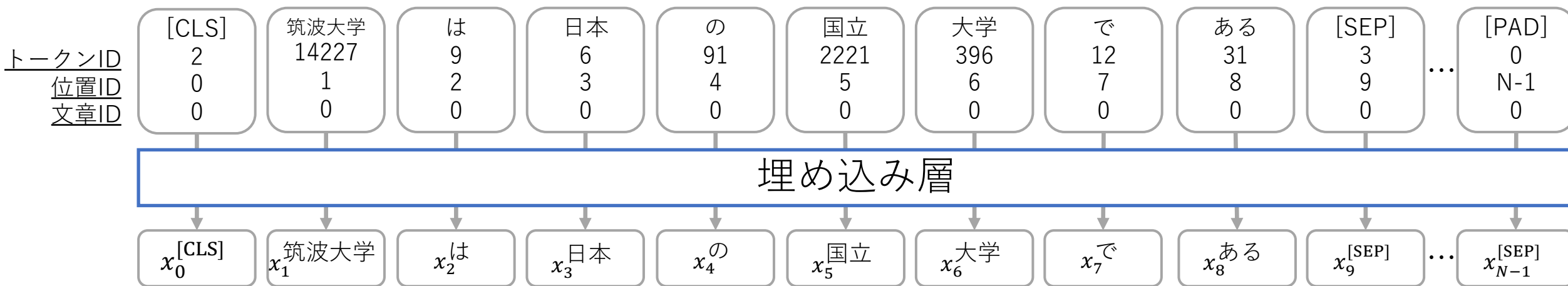
# 2.A 埋め込み



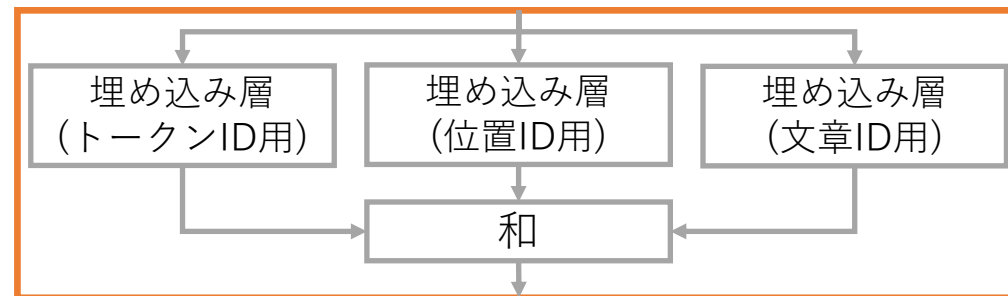
## 埋め込みとは

トークンIDなどの整数値はニューラルネットでは扱いにくいいため、連続値ベクトルに変換することが多い。この連続値ベクトルのことを埋め込みと呼ぶ。

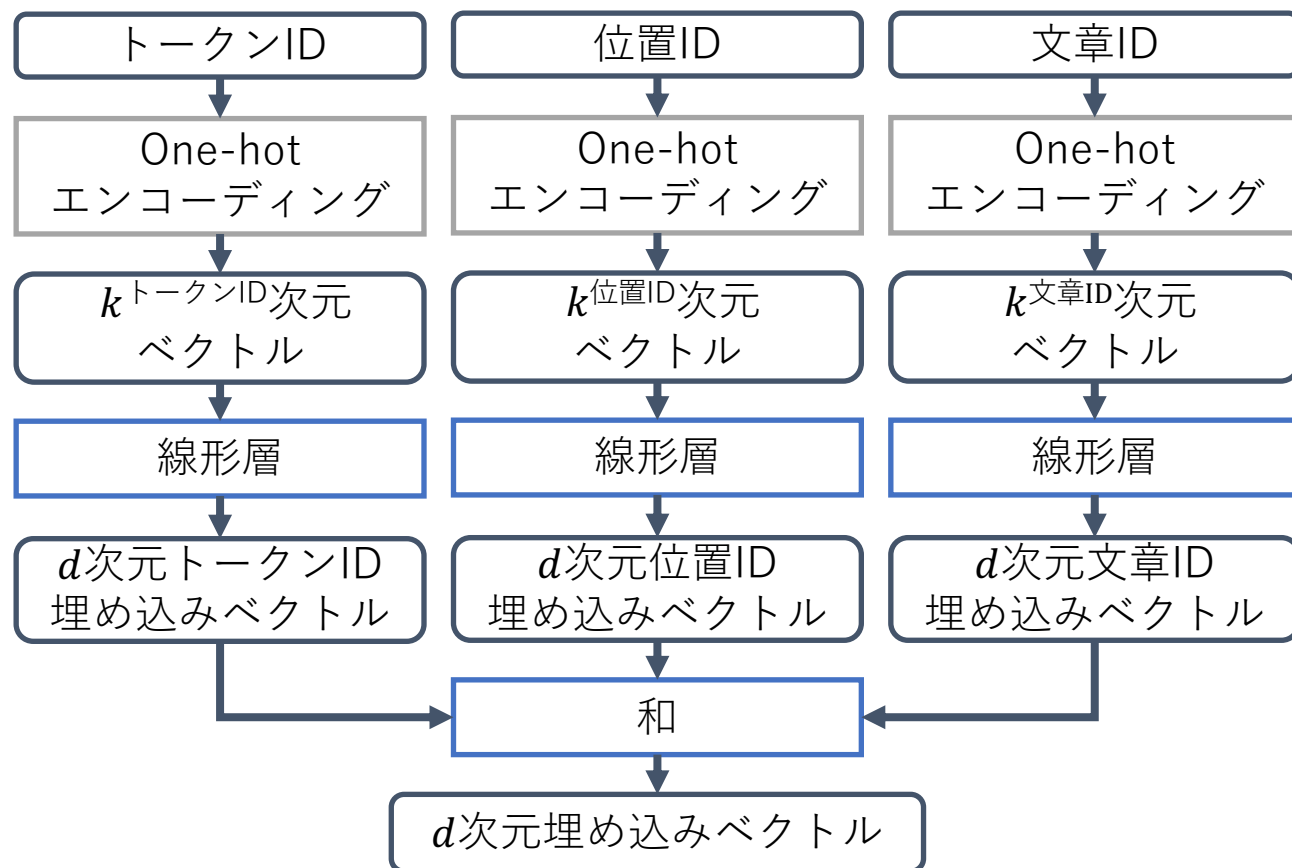
各トークンごとに $d$ 次元の埋め込みを得る。



## 2.A 埋め込み



各埋め込み層は  
One-hotエンコーディングと線形層からなる。



One-hotエンコーディングとは  
対応するインデックスのみ1、それ以外0のベクトルに変換する。  
次元数 $k$ は、入力整数値の範囲に対応する。  
 $k$ トークンID $\Rightarrow$ 全体の語彙数  
 $k$ 位置ID $\Rightarrow$ 入力長  
 $k$ 文章ID $\Rightarrow$ 1入力に含まれる文章数

# 2.B 自己注意機構

## Self-Attention Mechanism

自己注意(Self-Attention)とは  
 入力ベクトルを、**文脈を加味した**ベクトルに変換する機構

以前から、自己注意はLSTMの補助等で活用されていたが、Attention Is All You Need(2017)で提案されたTransformerでモデルの核として利用されて以来画像処理等でも広く使用されるようになった。

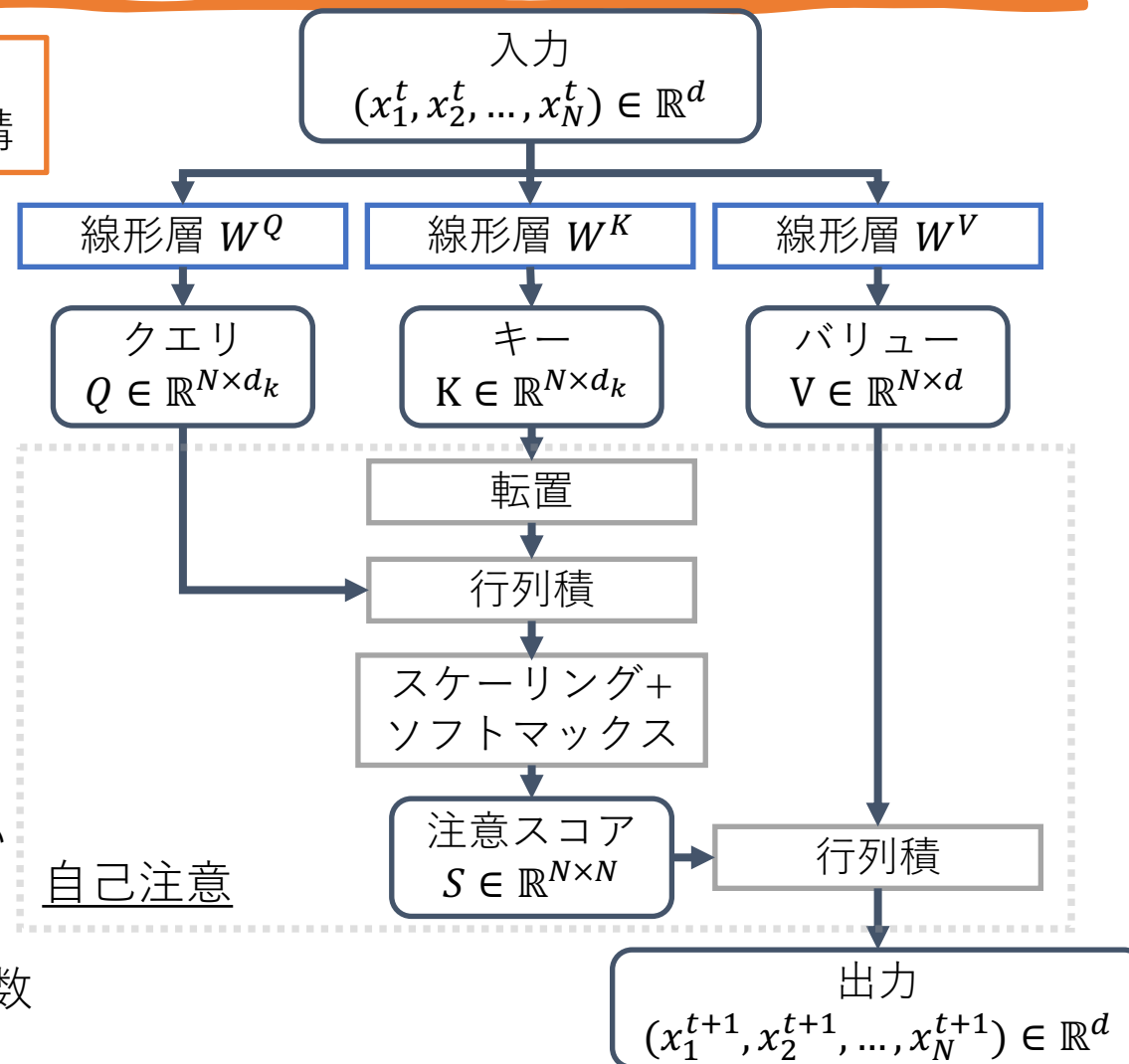
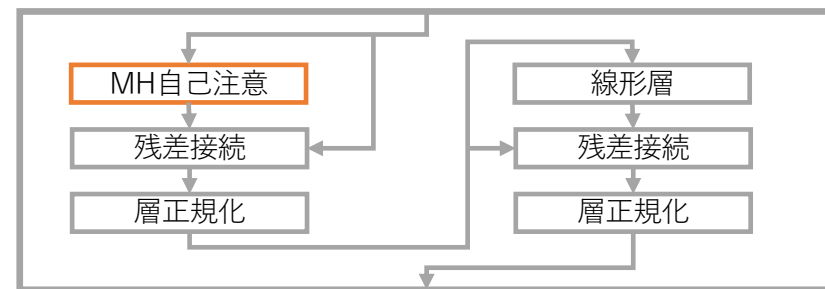
自己注意を数式で表すと

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$d_k$  : Query, Keyのベクトル次元数

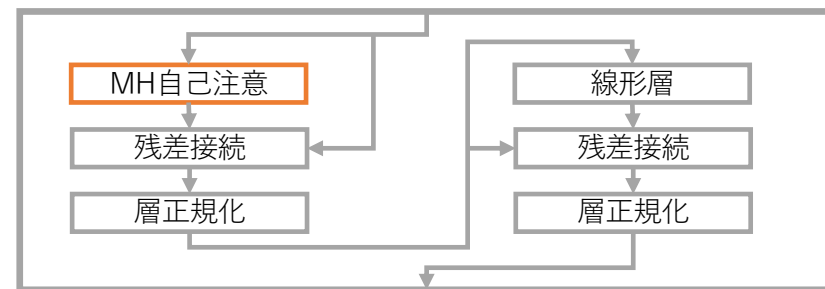
$\sqrt{d_k}$  : 次元数 $d_k$ が大きいと内積値 $QK^T$ も大きくなりやすい次元数の影響を調節するスケール因子

ソフトマックスは値を0~1の範囲、合計1に正規化する関数



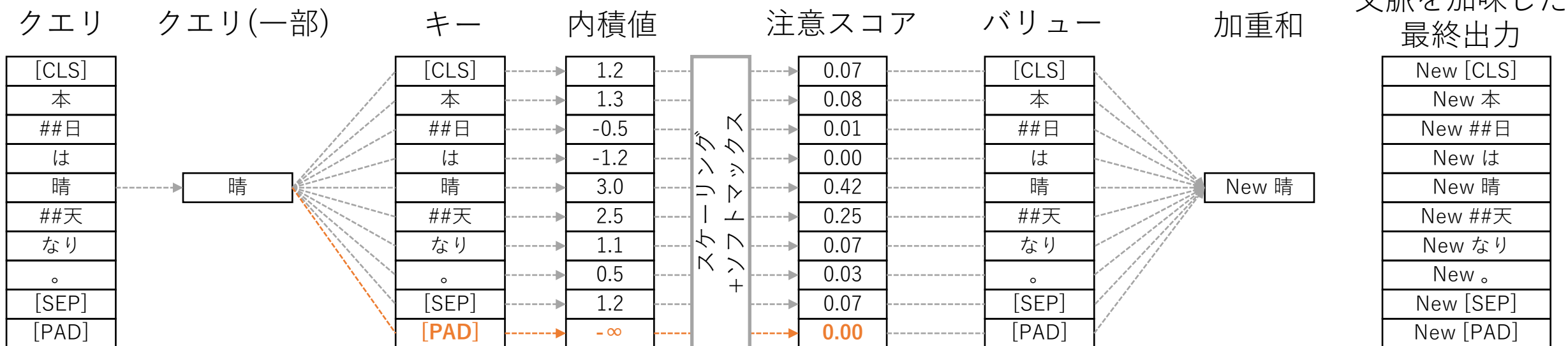
# 2.B 自己注意機構

## Self-Attention Mechanism



どのように文脈を加味したベクトルを得ているのか図解的に説明していく。

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



わかりやすさのため、一旦クエリのうち1ベクトルのみで考える。

各キーと内積を取る

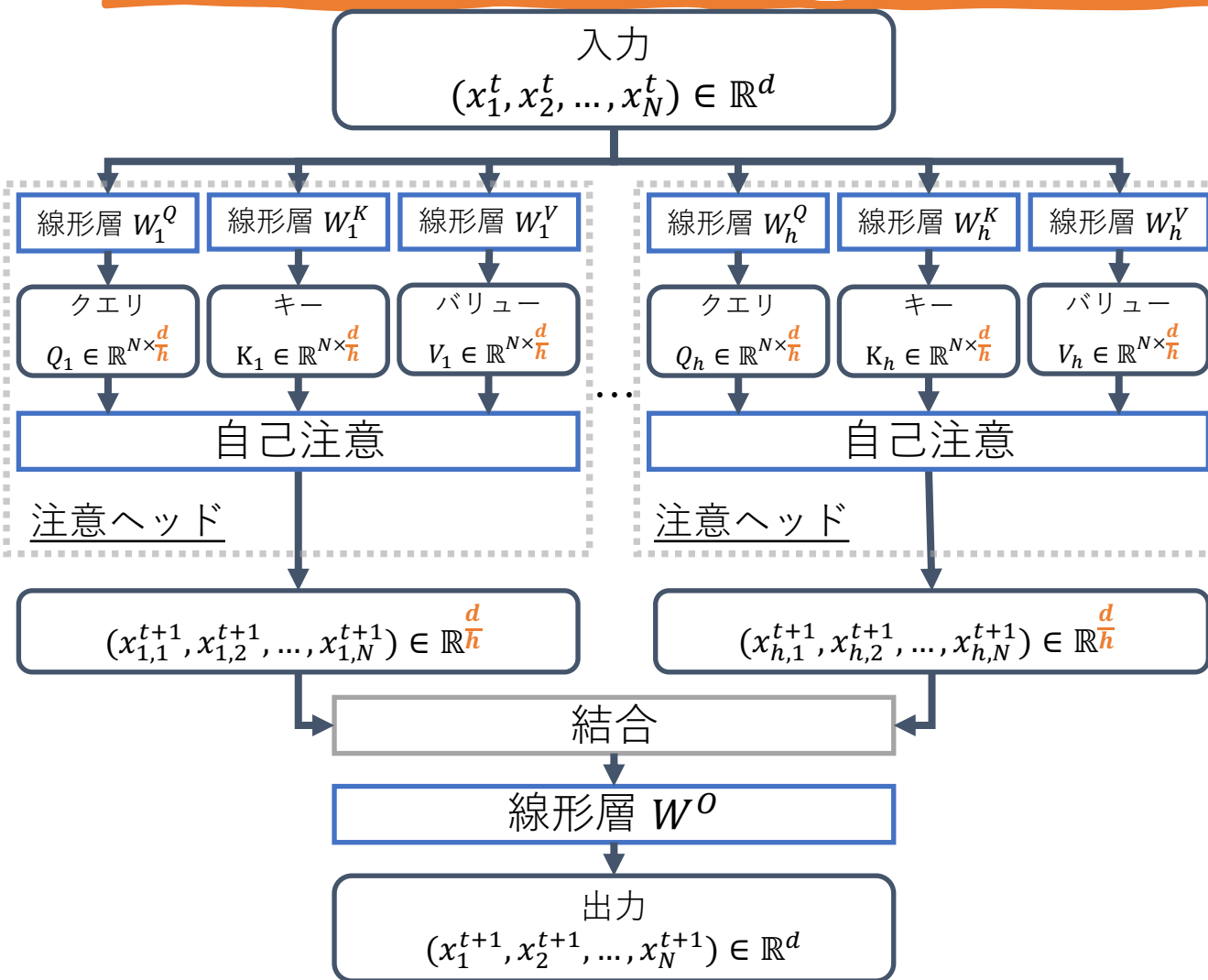
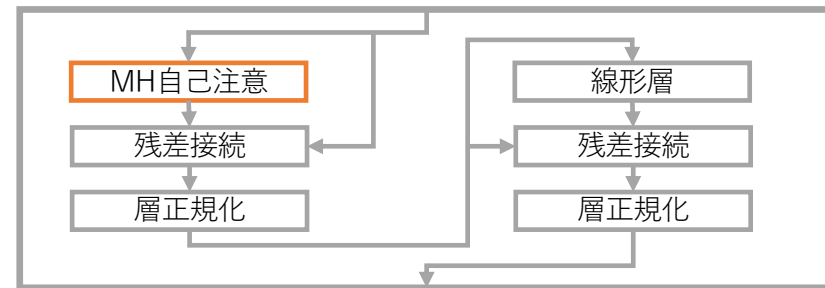
[PAD]とのスコアは $-\infty$ で置き換える。  
(この時に注意マスクが使用される)  
⇒ソフトマックスを通した際に値が限りなく0に近くなる。

注意スコアをもとにバリューの加重和を取る。  
⇒**文脈を加味した埋め込み**を得ることができる。

全てのクエリで同じ処理をすることで、最終出力を得る。

# 2.B マルチヘッド自己注意

Multi-Head Self-Attention



マルチヘッド自己注意とは  
 複数の自己注意を並列で行うことで、単一の場合より**多くの文脈関係を扱うことができる**ようになる手法。

線形層を用いて、入力 $N \times d$ を $N \times \frac{d}{h}$ 次元に圧縮。  
 ここで $h$ は $d$ の約数。  
 $h$ 個の注意ヘッドを用いることで出力 $h \times N \times \frac{d}{h}$ を得た後、結合し $N \times d$ とする。  
 結合位置を隠すため線形層に通し最終出力を得る。  
 式で表すと以下のようなになる。

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

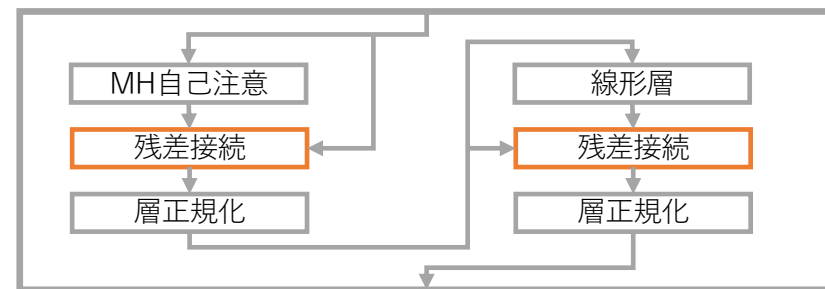
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$W_i^Q, W_i^K, W_i^V$  :  $i$ 番目の頭の次元圧縮用の線形層の重み  
 $W^O$  : 結合の境目を隠すための線形層の重み



# 2.C 残差接続

## Residual Connection



### 残差接続とは

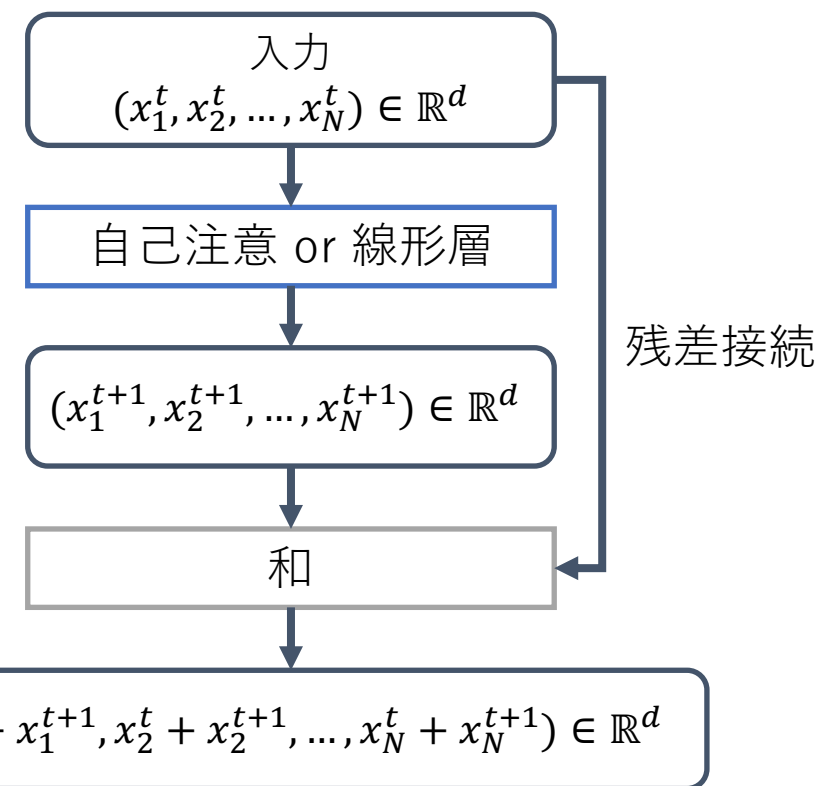
深い(つまり多層の)ネットワークを**安定して学習**するためのテクニック

Deep Residual Learning for Image Recognition(2015)で提案された。

何かしらのネットワーク(今回の場合は自己注意か線形層)の入力  $x^t$  と出力  $x^{t+1}$  の和  $x^t + x^{t+1}$  を最終出力とする。

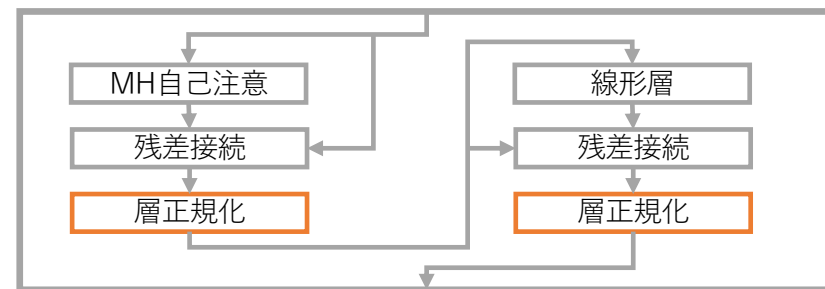
この場合、モデルは入力との差分  $x^{t+1}$  を学習すれば良いため全体のネットワークが深くなっても**学習が安定**する。

やや詳しく説明するとニューラルネットは層が深くなればなるほど、誤差逆伝播の際に勾配が消失しやすくなることが知られている。残差接続はこの勾配消失を防ぐ効果を持っており、結果として深いネットワークの学習の安定に繋がる。



# 2.D 層正規化

## Layer Normalization



### 層正規化とは

ネットワーク出力の分布変動を抑制することで、学習を**安定化&高速化**できる手法。

Layer Normalization(2016)で提案された。

1次元に結合した残差接続後の出力： $a^t = \text{concat}(x_1^t, x_2^t, \dots, x_N^t)$ ,  $a^t \in \mathbb{R}^{Nd}$

平均と分散を用いて、出力を平均0、分散1に正規化する。

$$\text{平均: } \mu^t = \frac{1}{H} \sum_{i=1}^H a_i^t$$

$$\text{分散: } \sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^t - \mu^t)^2}$$

$$h^t = \frac{\mathbf{g}}{\sigma^t} \odot (a^t - \mu^t) + \mathbf{b}$$

$\odot$ : 要素積

$(\mathbf{g}, \mathbf{b}) \in \mathbb{R}^{Nd}$ : 学習可能なベクトル

$$(x_1^{t'}, x_2^{t'}, \dots, x_N^{t'}) = \text{decompose}(h^t)$$

# 3 事前学習

## Pre-Training

### 事前学習とは

モデルを大規模データセットで学習した後、解きたいタスクに転用することで**収束を早めたり精度を向上させる**ことができる学習手法。

膨大な計算コストを必要とするため、一般的には配布されている事前学習済みのBERTモデルを使用することが多い。

本ワークショップでも東北大が公開している日本語事前学習済みBERTを使用する。

BERTの事前学習は

A. マスク言語モデリング (Masked Language Modeling)

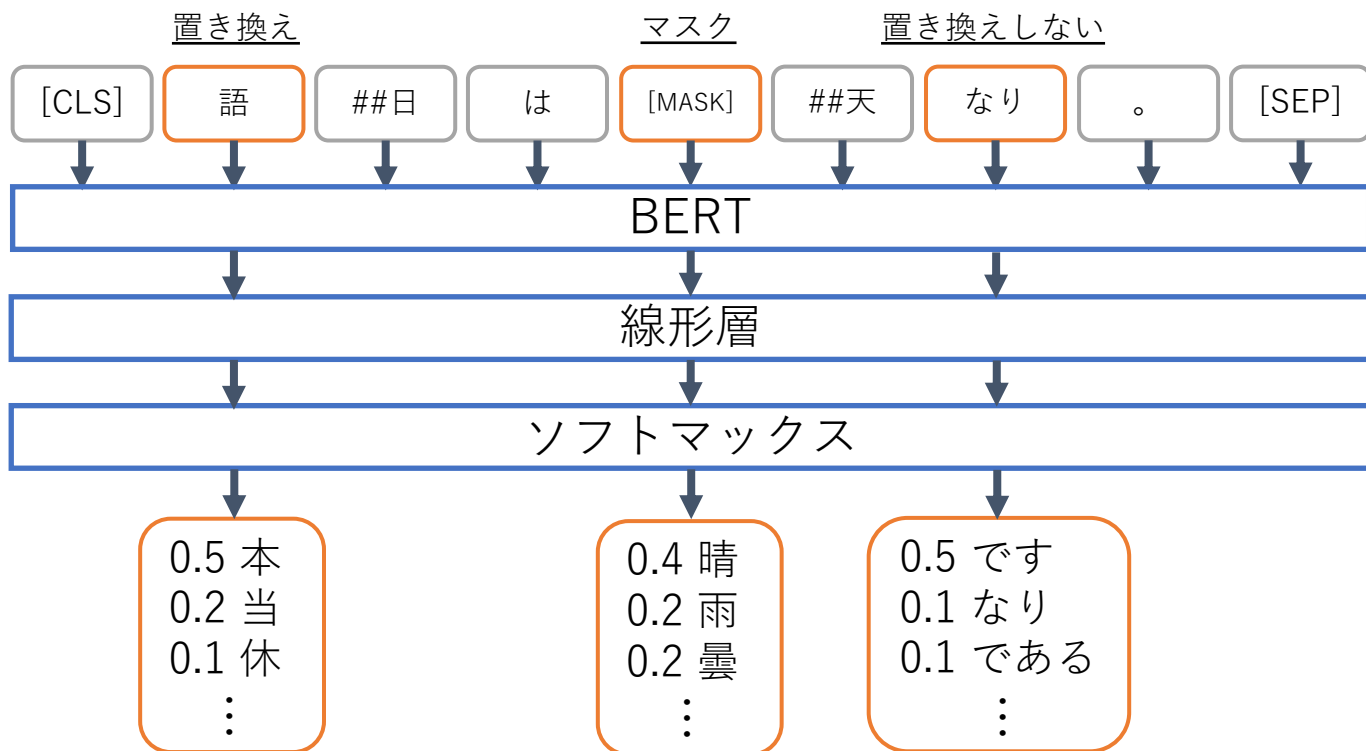
B. 次文予測 (Next Sentence Prediction) ←あまり効果がないことがわかっているため説明を省く※  
からなる。

※RoBERTa: A Robustly Optimized BERT Pretraining Approach(2019)

# 3.A マスク言語モデリング

## Masked Language Modeling

一定の確率で単語を[MASK]に置き換えて単語を当てるタスク。  
Wikipediaなどの大規模コーパスで事前学習を行う。



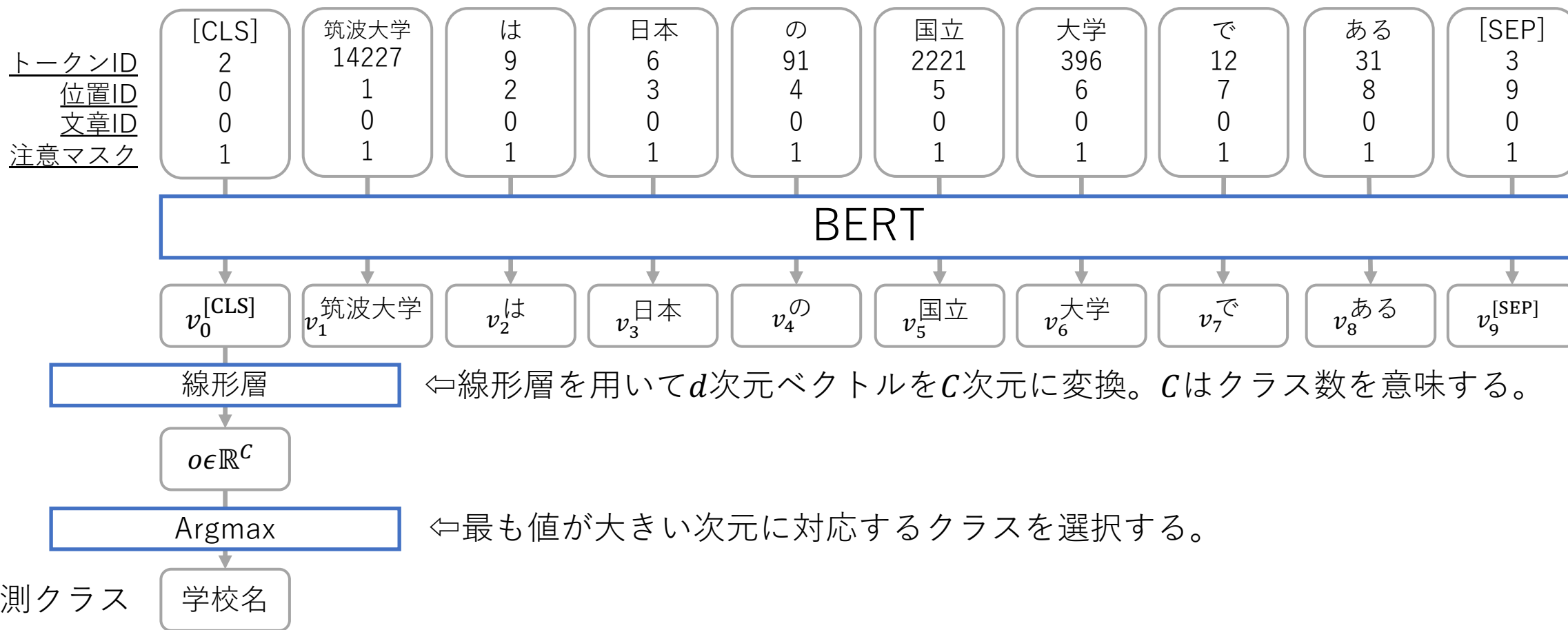
全体の15%の単語が対象となる。

そのうち、  
80%を[MASK]トークン  
10%を別のトークン  
10%を同じトークン  
に置き換える。

BERTは対象の単語が何であったかを当てる必要がある。

# 4 分類タスクへの適用

分類タスクの場合、[CLS]トークンに対応する出力ベクトルを用いてラベル予測を行う。



# 森羅分類データセットの説明

# 森羅分類データセットとは

今年度開催している森羅共有タスクのデータセットの一つ。  
日本語Wikipediaを拡張固有表現に分類する。

日本語Wikipedia



ラベルは複数付与される場合がある。  
例えば、映画化された書籍の場合、  
書物名、映画名が付与される。  
全体の約3%に複数ラベルが付与されている。

## 拡張固有表現階層(約200カテゴリー)

<b>名前</b> 名前_その他 人名 神名 <b>生物呼称名</b> 動物呼称名 動物呼称名_その他 競走馬名 植物呼称名 <b>組織名</b> 組織名_その他 国際組織名 公演組織名 家系名 民族名 民族名_その他 国籍名 競技組織名 競技組織名_その他 競技連盟名 競技リーグ名 競技団体名 法人名 法人名_その他 非営利団体名 / 企業名 企業グループ名 政治的組織名 政治的組織名_その他 政府組織名 / 政党名 / 内閣名 軍隊名 <b>地名</b> 地名_その他 GPE GPE_その他 市区町村名 / 都道府県州都部名 国名 地域名 地域名_その他 大陸地域名 / 国内地域名 地形名 地形名_その他 温泉名 / 山地名 / 島名 / 河川名 湖沼名 / 海洋名 / 湾名 天体名 天体名_その他 天体部分名 / 銀河名 / 恒星名 惑星_衛星名 / 星座名 アドレス アドレス_その他 郵便住所	<b>施設名</b> 施設名_その他 施設部分名 ダム名 遺跡名 遺跡名_その他 墳墓名 FOE FOE_その他 軍事基地名 / 城名 / 宮殿名 公共機関名 / 宿泊施設名 医療機関名 / 学校名 研究機関名 / 取引所名 発電所名 / 公園名 商業施設名 / 競技施設名 美術博物館名 / 動物園名 遊園施設名 / 劇場名 宗教施設名 <b>交通施設名</b> 交通施設名_その他 停車場名 / 鉄道駅名 空港名 / 港名 / 道路施設名 鉄道施設名 路線名 路線名_その他 道路名 / 鉄道路線名 / 運河名 航路名 / トンネル名 / 橋名 <b>イベント名</b> イベント名_その他 <b>自然現象名</b> 自然災害名 自然災害名_その他 地震名 / 水害名 催し物名 催し物名_その他 祭礼名 / 選挙名 / 競技会名 展示会名 / 会議名 事故事件名 事故事件名_その他 交通事故名 / 社会事件名 戦争名 <b>病気名</b> 病気名_その他 動物病気名	<b>プロダクト名</b> プロダクト_その他 株名 / 便名 / 識別番号 / サービス名 / ブランド名 ソフトウェア名 / 情報機器名 玩具名 / 楽器名 / 衣類名 医薬品名 / キャラクター名 <b>作品名</b> 作品名_その他 絵画名 / 番組名 / 映画名 公演名 / 音楽名 / 映像作品名 <b>出版物名</b> 出版物名_その他 新聞名 / 雑誌名 / 書物名 ゲーム名 ゲーム名_その他 電子ゲーム名 <b>食べ物名</b> 食べ物名_その他 料理名 <b>武器名</b> 武器名_その他 火器名 乗り物名 乗り物名_その他 車名 / 列車名 / 飛行機名 船名 / 宇宙船名 / 軍用車両名 軍用機名 / 艦艇名 <b>主義方式名</b> 主義方式名_その他 罪名 / 等級名 / 賞名 / 勳章名 貨幣名 / 技術名 / 規格名 / 制度名 試験名 / 主義思想名 / 文化名 宗教名 / 学問名 / 理論名 / 流派名 競技名 / 政策計画名 <b>規則名</b> 規則名_その他 条約名 / 法令名 称号名 称号名_その他 地位職業名 <b>言語名</b> 言語名_その他 国語名 <b>単位名</b> 単位名_その他 通貨単位名 <b>バーチャルアドレス名</b> バーチャルアドレス名_その他 チャンネル名 / 電話番号 電子メール / URL	<b>自然物名</b> 自然物名_その他 元素名 / 化合物名 / 鉱物名 生物名 生物名_その他 空想生物名 / 細菌_ウイルス名 真菌類名 / 軟体動物名 節足動物名 / 昆虫類名 / 魚類名 両生類名 / 恐竜名 / 爬虫類名 鳥類名 / 哺乳類名 / 植物名 <b>生物部位名</b> 生物部位名_その他 動物部位名 / 植物部位名 <b>色名</b>	<b>時間表現</b> 時間表現_その他 時間 時間_その他 時刻表現 / 日付表現 / 曜日表現 時代表現 <b>期間</b> 期間_その他 時刻期間 / 日数期間 / 週数期間 月数時間 / 年数期間	<b>数値表現</b> 数値表現_その他 <b>金額表現</b> <b>株指標</b> <b>ポイント</b> <b>割合表現</b> <b>倍数表現</b> <b>頻度表現</b> <b>年齢</b> <b>学齢</b> <b>序数</b> <b>順位表現</b> <b>経度緯度</b>	<b>個数</b> 個数_その他 人数 組織数 場所数 場所数_その他 国数 施設数 プロダクト数 イベント数 自然物数 自然物数_その他 動物数 / 植物数	<b>CONCEPT</b> <b>IGNORED</b>
---	--	--	--	--	---	--	----------------------------------

# 森羅分類データセットとは

学習データ：分類済み日本語Wikipedia(2019年時点)

予測データ：未分類日本語Wikipedia(2021年時点)

学習データは100万件近く存在する。

計算コストの都合上、ワークショップでは5000件をランダムにサンプリングし学習データとする。

また、簡易化のため約3%の複数ラベルを除去し、単一ラベルのデータセットとして扱う。

本ワークショップでは学習、予測、リーダーボードによる評価まで行う。

Rank	Team Name	Submitted on	Description	Micro-F1 (Public)
1	運営-中山	2022/06/01	RoBERTaベースライン <a href="https://github.com/k141303/Shinra2022">https://github.com/k141303/Shinra2022</a>	95.2224
2	森羅2022実行委員会	2022/06/24	ベースラインシステム	94.2519
3	運営-中山-森羅ワークショップ	2022/07/31	森羅ワークショップで使用したコードを用いて学習したBERTの結果です。	87.1518
4	M ando	2022/08/03	デフォルト	86.3202



# いざ実践

# まずは動作確認

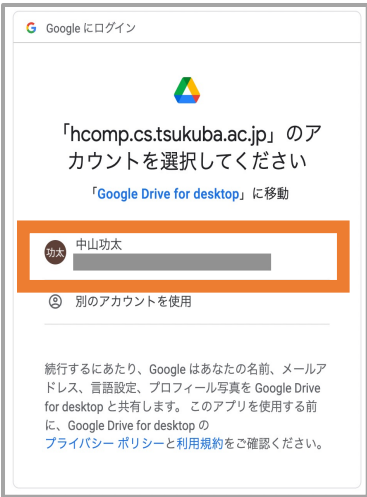
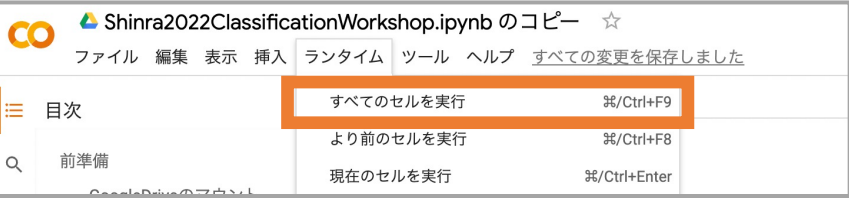
①Colabを開きます。(Googleアカウントが必要です。)

<https://colab.research.google.com/drive/1M2xbF4gqCH5O7LnmnAvVOkbBO22SiX8K?usp=sharing>

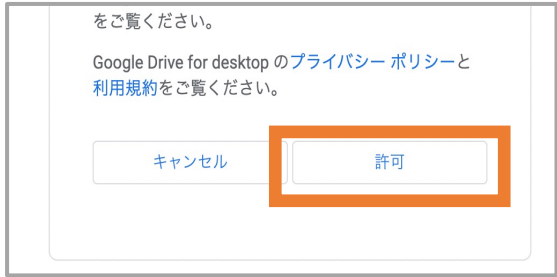
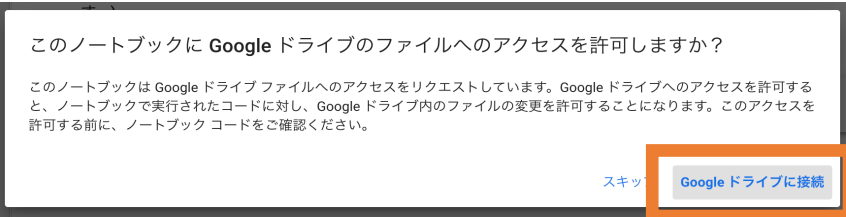
②自分のGoogleDriveにコピーします。(コードの編集が可能になります。)



③ランタイム>全てのセルを実行を押す



④ポップアップが出るのでGoogleドライブに接続



④5分ほど待って最後行のコードの横に緑のチェックがつけば成功

```
0 秒
▶ def save_json(output_path, data):
  data = map(json.dumps, data)
  with open(output_path, "w") as f:
    f.write("\n".join(data))
```

もし動かない場合はSlack #bert\_workshop2022チャンネルでご相談ください🙏