

現場で使えるディープラーニング基礎講座

DAY7

SkillUP AI

前回の復習

- 前回は何を学びましたか？

自然言語処理における深層学習

目次

1. 自然言語処理と深層学習
2. 自然言語処理の基礎
3. word2vec
4. 系列変換モデル
5. アテンション
6. トランスフォーマー
7. 外部メモリを持つニューラルネットワーク
8. その他の話題

自然言語処理と深層学習

自然言語処理における深層学習

- 人間がコミュニケーションに用いている言葉を対象とする処理を自然言語処理と呼ぶ。
 - 自然言語以外の言語の例
 - コンピュータ言語などの形式言語
 - Pythonなどのプログラミング言語
 - エスペラントのような人工言語
- 自然言語処理の分野でも深層学習を用いた手法が多数提案されている。

伝統的な自然言語処理

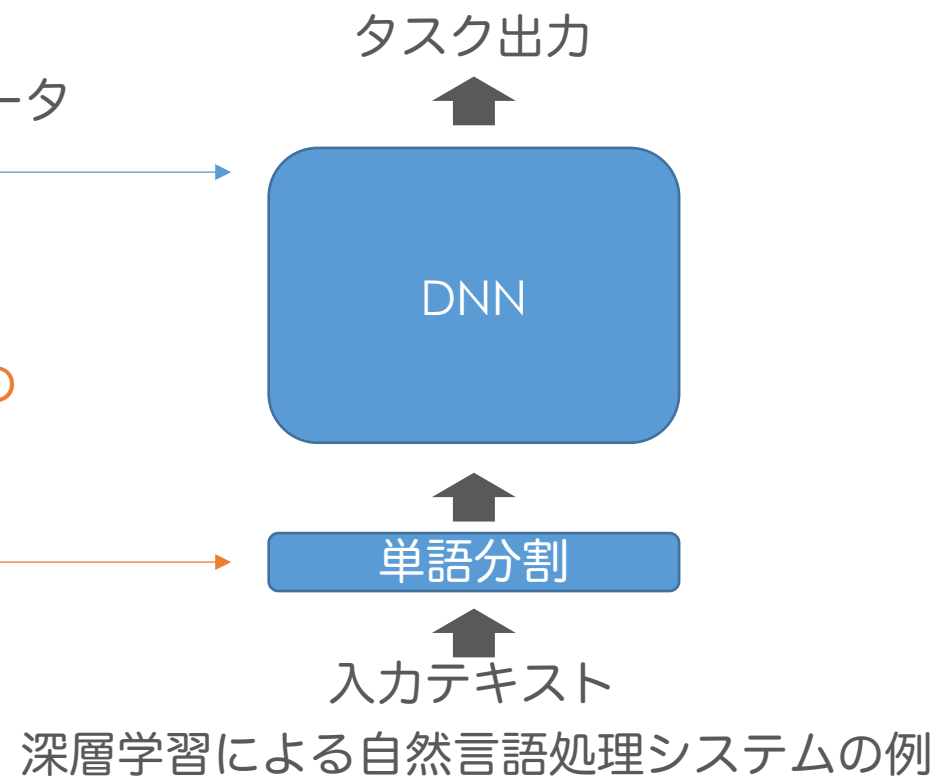
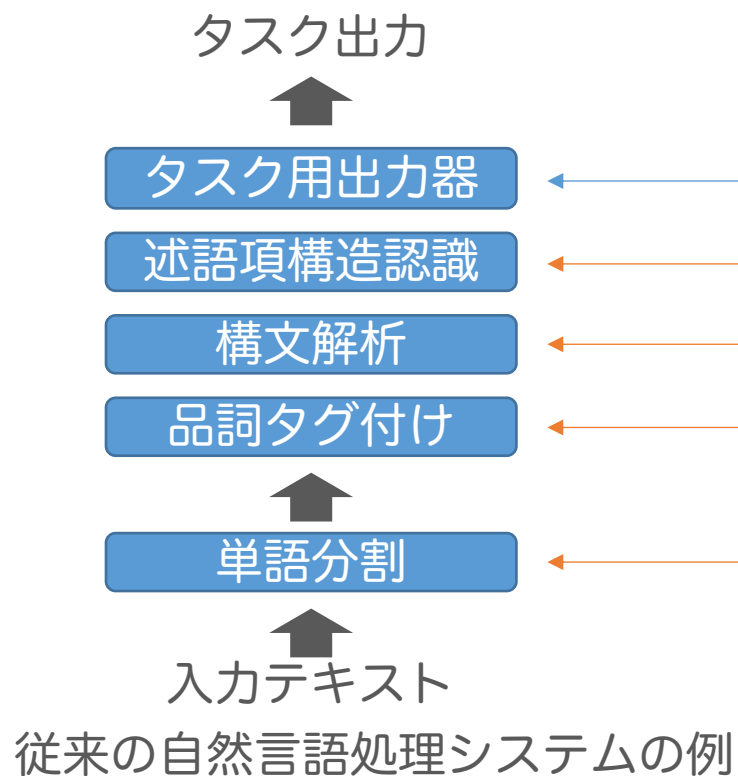
- 自然言語処理で典型的なタスクには、文章分類、機械翻訳、文章要約、質問応答、対話などがある。
- 伝統的な自然言語処理では、これらのタスクを部分問題に分解していた。
- 主な部分問題を以下に示す。

部分問題	処理内容
品詞タグ付け	単語に名詞・動詞などの文法的な役割分類(品詞)を付与する処理
単語分割	日本語などの単語に分けられていないテキストを単語列に分割する処理
語義曖昧性解消	複数の語義をもつ単語の語義を特定する処理
固有表現抽出	人名・地名・日付などを抽出する処理
構文解析	文法に基づく文の木構造を構築する処理
述語項構造認識	述語を中心とした意味構造を抽出する処理

参考:深層学習による自然言語処理, 坪井ら, 講談社

深層学習による自然言語処理システム

- 従来の自然言語処理システムでは、部分問題用学習器を組み合わせることにより、タスクを解いていた。
- 深層学習を用いた自然言語処理システムでは、1つの学習器でタスクを解くようにモデルが構築される。



深層学習による自然言語処理システム

- 深層学習を用いた自然言語処理システムでは、再帰型ニューラルネットワーク(RNN)がよく用いられる。
- 一方で、RNNを用いずに機械翻訳問題を解くというモデルも提案されている。
 - Ashish Vaswani et al. Attention is All You Need, NIPS 2017.
<https://arxiv.org/pdf/1706.03762.pdf>
- 単語の埋め込み行列を求める際には、全結合型ニューラルネットワークが用いられることが多い。
 - 例、word2vec

自然言語処理に関わる代表的なできごと

- 東西冷戦下では、**英語-ロシア語の機械翻訳**が注目される。
- 1964年~1966年、心理カウンセラーのように振る舞う**対話型プログラムELIZA**がジョセフ・ワイゼンバウムによって開発される。
- 2011年、IBMが開発した**人工知能ワトソン**が、アメリカのクイズ番組「ジョパディー」に出演し、歴代の人間チャンピオンと対戦し勝利。
- 東大入試合格を目指す**人工知能東ロボくん**。国立情報学研究所が中心となって、2011年~2016年まで開発が続けられ、偏差値57.8をマークし、私立大学に合格できるレベルまで達した。
- 2016年、**Googleのニューラル機械翻訳**がサービス開始。
- 2010年代、**スマートスピーカーが普及**(Amazon Alexa:2014年、Google Home:2016年、Apple Home Pod:2017年)

自然言語処理の基礎

形態素解析

- 自然言語を形態素まで分解する技術。
- 形態素とは、何らかの意味をもつ最小限の文字の集まり。
- 日本語での検索エンジンや自然言語処理で極めて重要な技術。
- MeCabやJUMAN、ChaSenなどが代表的な形態素解析エンジン。

すももももももものうち



すもも	も	もも	も	もも	の	うち
名詞、一般	助詞、係助詞	名詞、一般	助詞、係助詞	名詞、一般	助詞、連体化	名詞、非自立、 副詞可能

n-gram

- n-gramとは、隣り合って出現したn個の単語のこと。文字n-gramもあるが、ここでは単語n-gramのことを指している。
- n=1の場合は、1-gram(ユニグラム)。
 - 吾輩, は, 猫, で, ある
- n=2の場合は、2-gram(バイグラム)。
 - 吾輩は, は猫, 猫で, である
- n=3の場合は、3-gram(トリグラム)。
 - 吾輩は猫, は猫で, 猫である

バッグオブワーズ

- バッグオブワーズ(bag-of-words)とは、文章における単語の出現頻度を数えて、頻度ベクトルで表現する方法。
- 単語(words)をバラバラにして袋(bag)に入れた状態に例えられている。
- 頻度ベクトルには、語順の情報は残らない。
- 実際には、出現頻度が極端に低い単語や一般的すぎる単語を外してベクトル化する。
- n-gramで数えることもある。

すももももももものうち
すいかもももまるいくだもの



左から、
{すもも、もも、すいか、も、の、まるい、うち、くだもの}
の出現頻度

[1, 2, 0, 2, 1, 0, 1, 0]

[0, 1, 1, 2, 0, 1, 0, 1]

頻度ベクトル

分布仮説

- 分布仮説(distributional hypothesis)とは、「単語の意味は、その単語が出現した際の周囲の単語によって決まる」という考え方。
- 1950年代に提案された。
- 後述するword2vecは、この分布仮説を用いて分散表現を獲得していると解釈できる。

私 は 今日 会社 へ 行く

会社 という単語の意味は、 周辺の単語(文脈) から決まる

[演習] MeCabによる形態素解析

- 7_1_morphological_analysis.ipynb
 - MeCabはオープンソースの高速な形態素解析エンジンです。
 - MeCabを用いて形態素解析を行なってみましょう。

word2vec

単語を数学的に表現するには？

- 自然言語を計算機で扱うには、記号である「単語」を数学的に表現する必要がある。
- 単語を数学的に表現するにはどうすればいいか？
- 例えば、ワンホットベクトルで表現するという方法が考えられる。
 - この方法の場合、ベクトル同士の演算を行うことに意味がない。
 - しかも、語彙数が増えると非常に高次元なベクトルになってしまう。
- もっと賢い方法はないか？

分散表現

- 分散表現(単語を強調する場合は単語分散表現と呼ばれる)とは、単語を計算機上で扱うための方法の1つ。単語埋め込み(word embeddings)とも呼ばれる。
- 単語間の関連性や類似度を考慮したベクトル(埋め込みベクトル)で、ある単語を表現する。
- 埋め込みベクトルどうしは、足し算や引き算が可能。
- 単語が5つだけの場合の分散表現の例を以下に示す。
- 単語埋め込み行列を求めるには、word2vecなどを用いる。

分散表現の例

$$\begin{array}{c} \text{埋め込み行列} \\ \text{(各列は埋め込みベクトル)} \end{array} \begin{pmatrix} 0.01 & 0.02 & 0.03 & 0.04 & 0.05 \\ 0.06 & 0.07 & 0.08 & 0.09 & 0.10 \\ 0.11 & 0.12 & 0.13 & 0.14 & 0.15 \end{pmatrix} \begin{array}{c} \text{埋め込み} \\ \text{ベクトル} \end{array} \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0.02 \\ 0.07 \\ 0.12 \end{pmatrix}$$

ある単語のワンホットベクトル

単語のワンホットベクトルは、実際の問題では数万次元になったりする。一方、埋め込みベクトルは数百次元くらいに設定する。ある単語の埋め込みベクトルには、ある単語の情報が埋め込まれていると解釈できる。

分散表現

- 埋め込みベクトルは、単語間の類似度を計算する目的でも使われるが、次の何らかの分析に利用されることも多い。
- 例えば、回帰モデルやクラス分類モデルの入力に埋め込みベクトルを用いるという活用が考えられる。

word2vec

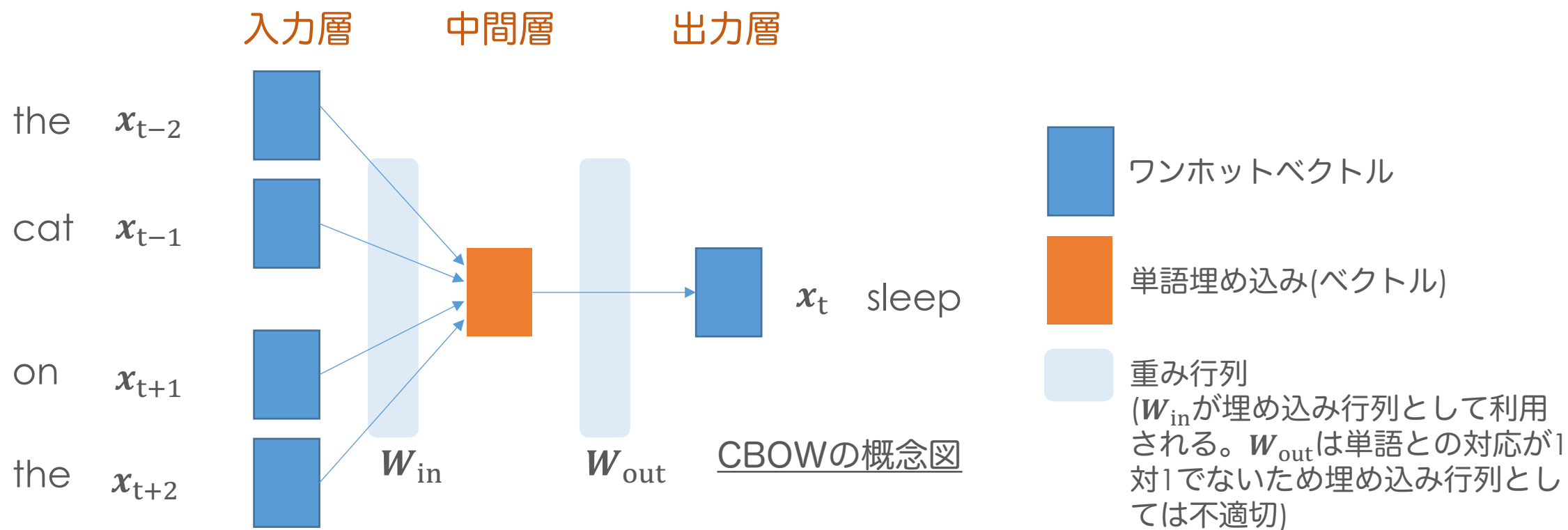
- 埋め込み行列を学習することができるツール。
- word2vecの公式ウェブサイト
 - <https://code.google.com/archive/p/word2vec/>
- word2vecをPythonから利用するには、gensimなどのライブラリを使用するのがよい。
- 「word2vecで求めた埋め込み行列を用いると、王 - 男 + 女 = 女王 のような演算が可能になる」と紹介されることが多い。
- 原著論文
 - Tomas Mikolov et al. Distributed Representations of Words and Phrases and their Compositionality. <https://arxiv.org/pdf/1310.4546.pdf>
 - Tomas Mikolov et al. Efficient Estimation of Word Representations in Vector Space. <https://arxiv.org/pdf/1301.3781.pdf>

word2vec

- word2vecには、埋め込み行列を学習する方法として、Continuous Bag-of-Words (CBOW) と skip-gramの2つのニューラルネットワークが用意されている。
- CBOWは、前後の単語から、ある単語を予測するためのニューラルネットワーク。
 - 例えば、”the cat sleep on the mat”という文章がある場合、“cat” と “on” から”sleep”を予測できるように学習を行う。
- skip-gram とは、ある単語が与えられた時に、その周辺の単語を予測するためのニューラルネットワーク。
 - 例えば、”the cat sleep on the mat”という文章がある場合、“sleep”から”cat”と”on”を予測できるように学習を行う。
- CBOWもskip-gram も、分布仮説を用いて分散表現を獲得していると解釈できる。

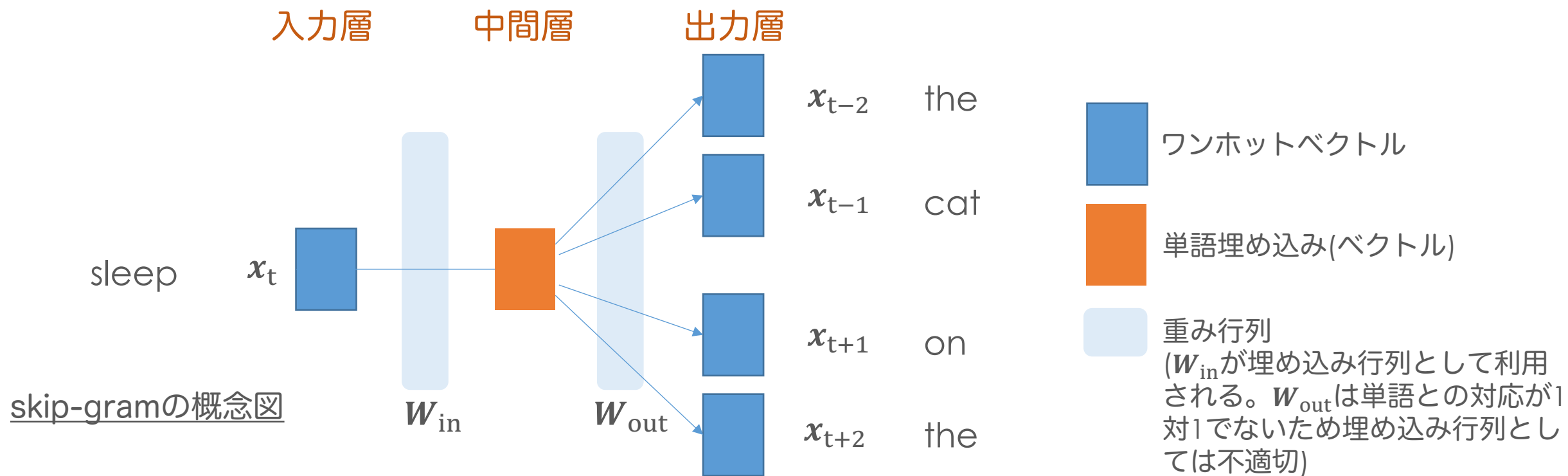
Continuous Bag-of-Words

- word2vecのアルゴリズムの1つ。
- 前後の単語から対象の単語を予測するニューラルネットワーク。
- 学習にかかる時間がskip-gramよりも短い。



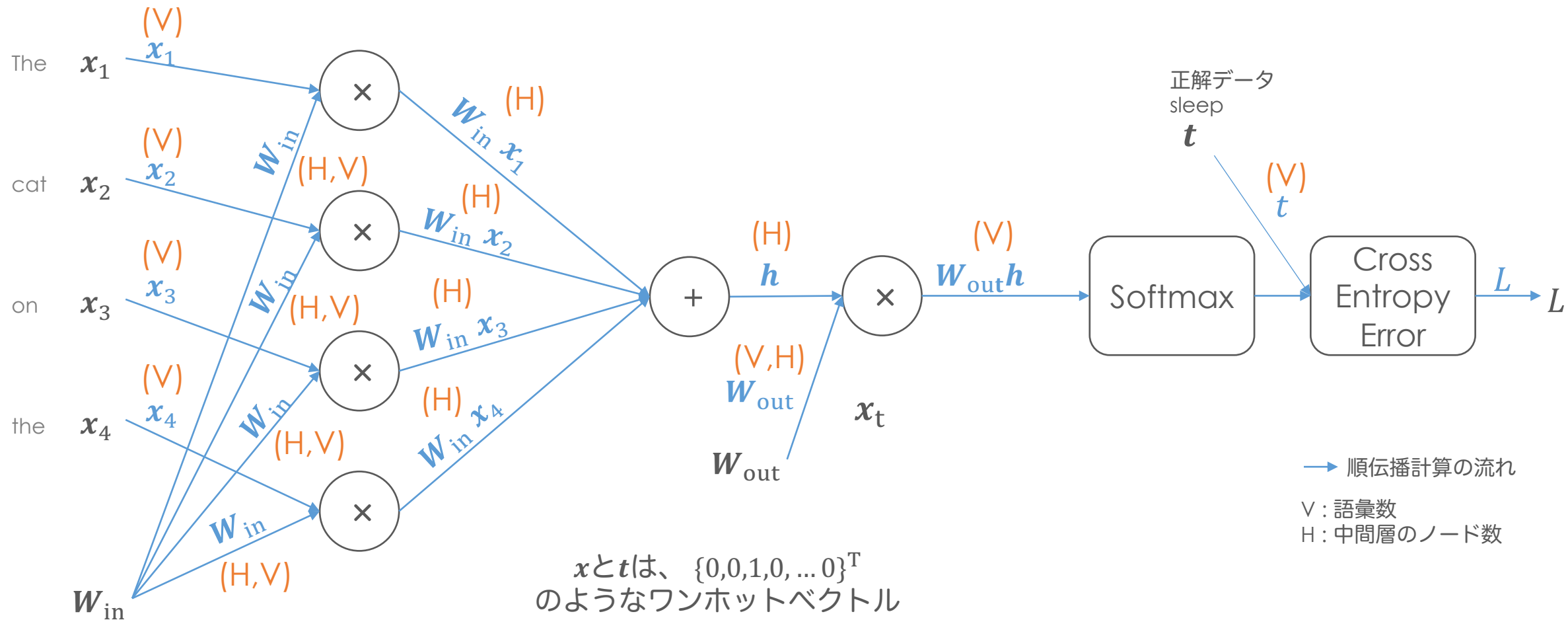
skip-gram

- word2vecのアルゴリズムの1つ。
- ある単語から、前後の単語を予測するニューラルネットワーク。
- CBOWに比べ、学習に時間がかかるが、精度が高くなる傾向がある。



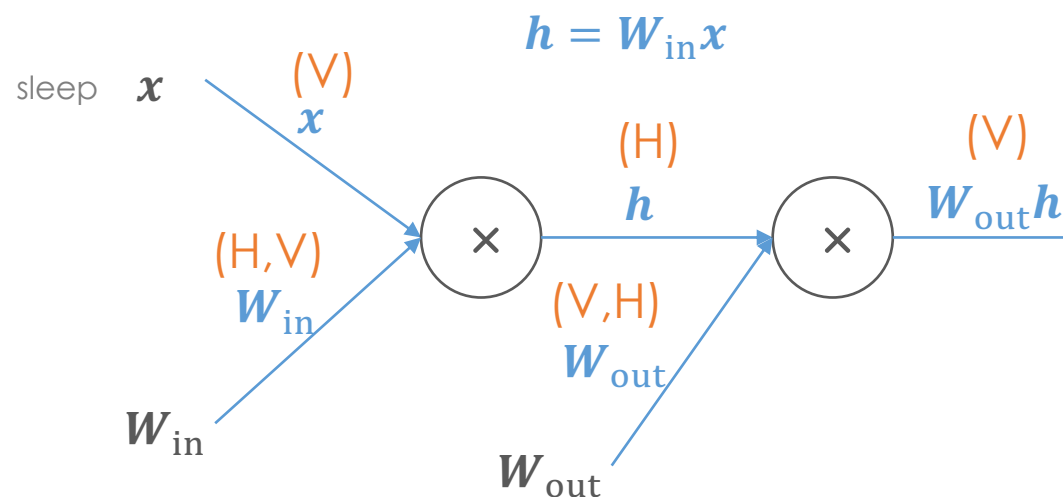
Continuous Bag-of-Wordsの計算グラフ

- CBOWの計算グラフ例を以下に示す。



skip-gramの計算グラフ

- skip-gramの計算グラフ例を以下に示す。

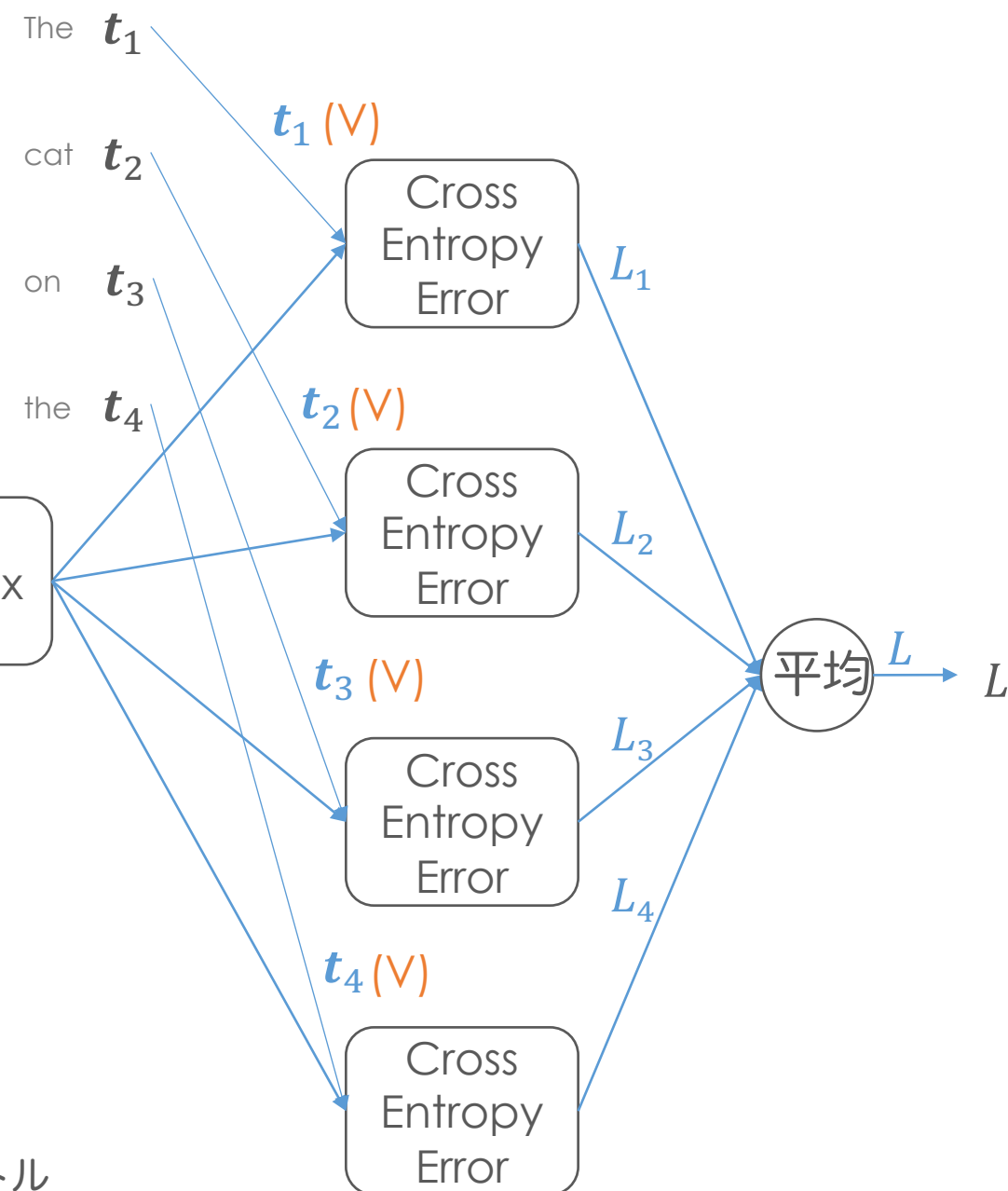


→ 順伝播計算の流れ

V : 語彙数
H : 中間層のノード数

x と t は、 $\{0,0,1,0, \dots 0\}^T$
のようなワンホットベクトル

正解
データ



コサイン類似度

- コサイン類似度とは、ベクトルどうしの類似度を測る指標。
- 埋め込みベクトルどうしの類似度を測る際にも使用される。

<コサイン類似度>

$$\text{similarity}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{x_1 y_1 + x_2 y_2 \cdots x_n y_n}{\sqrt{x_1^2 + x_2^2 \cdots x_n^2} \sqrt{y_1^2 + y_2^2 \cdots y_n^2}} = \cos \theta$$

$$\mathbf{x} = \{x_1, x_2, \cdots, x_n\}$$

$$\mathbf{y} = \{y_1, y_2, \cdots, y_n\}$$

[演習] word2vecを使って単語の分散表現を獲得する

- 7_2_word2vec.ipynb
 - word2vecを使って単語の分散表現を求めましょう。

Any Questions ?

系列変換モデル

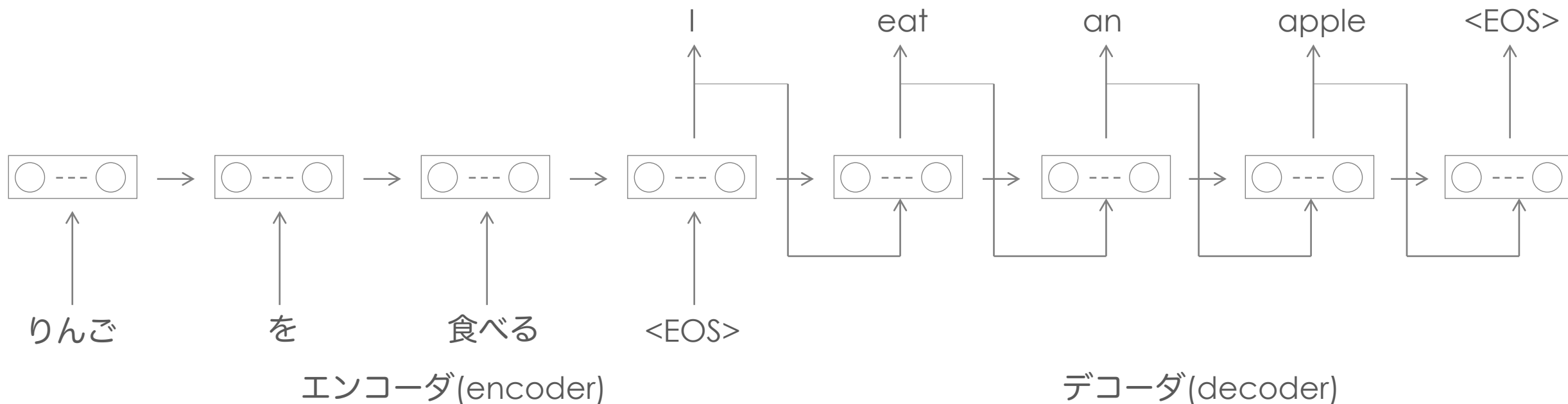
系列変換モデル

- 単語列などの系列(sequence)を受け取り、別の系列へ変換するモデルのことを系列変換モデル(sequence-to-sequence model, seq2seqモデル)という。
- 主に自然言語処理で用いられる。
- RNNエンコーダ・デコーダとも呼ばれる。
- 適応されるタスクの例
 - 機械翻訳 (翻訳元の言語から翻訳先の言語への変換)
 - 対話 (相手の発言から自分の発言への変換)
 - 質問応答 (質問文から返答文への変換)
 - 文書要約 (元文書から要約文への変換)

系列変換モデル

- 入力系列の長さと出力系列の長さは、一致していなくてもよい。
- 系列の長さは、データによって異なってもよい。

系列変換モデルの概念図



系列変換モデルの呼ばれ方

- 系列変換モデルは、seq2seqモデルと呼ばれる場合があったり、RNNエンコーダ・デコーダと呼ばれる場合があったりする。
- これは、両者の起源とする論文が異なるためである。
- 両者の概念はほぼ同じであるため、実務上は区別なく呼ばれることが多い。
- seq2seqモデル
 - Ilya Sutskever, Oriol Vinyals, Quoc V. Le. Sequence to Sequence Learning with Neural Networks. 2014. <https://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf>
- RNNエンコーダ・デコータ
 - Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. 2014. <https://www.aclweb.org/anthology/D14-1179>

[演習] seq2seqモデルの実装

- 7_3_seq2seq_trainee.ipynb
 - seq2seqモデルを実装しましょう。
 - RNN部分には、LSTMを用います。

Any Questions ?

アテンション

アテンション

- 系列変換モデルでは、系列情報をRNNで固定長のベクトルに変換するが、系列が長くなると、LSTMのような仕組みを使ってもうまく学習させるのは難しい。
 - 系列が長くなると、最初に入力された情報がデコード側まで伝播しづらくなる。
- この課題に対処する方法として、アテンションという仕組みがある。
- アテンションの仕組みを最初に提案した論文
 - Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate.
<https://arxiv.org/pdf/1409.0473.pdf>

ソフト・アテンションとハード・アテンション

- ソフト・アテンション

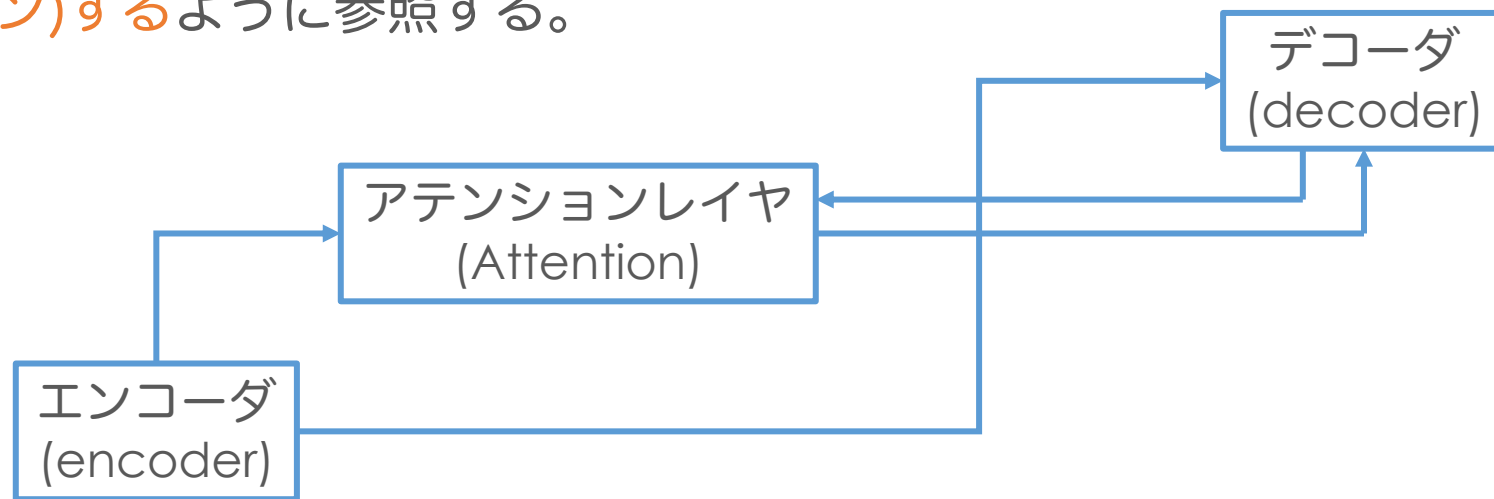
- ソフトマックス関数などで確率分布を求め、その確率分布を用いた重み付き平均をとることによって、データに注目する方法。
- seq2seqモデルで用いるアテンションは、ソフト・アテンション。

- ハード・アテンション

- ソフトマックス関数などで確率分布を求め、その確率分布に従って無作為抽出された1点だけに注目する方法。
- 無作為抽出が入っていると微分不可能なので、何らかの方法を用いて勾配を推定する必要がある。

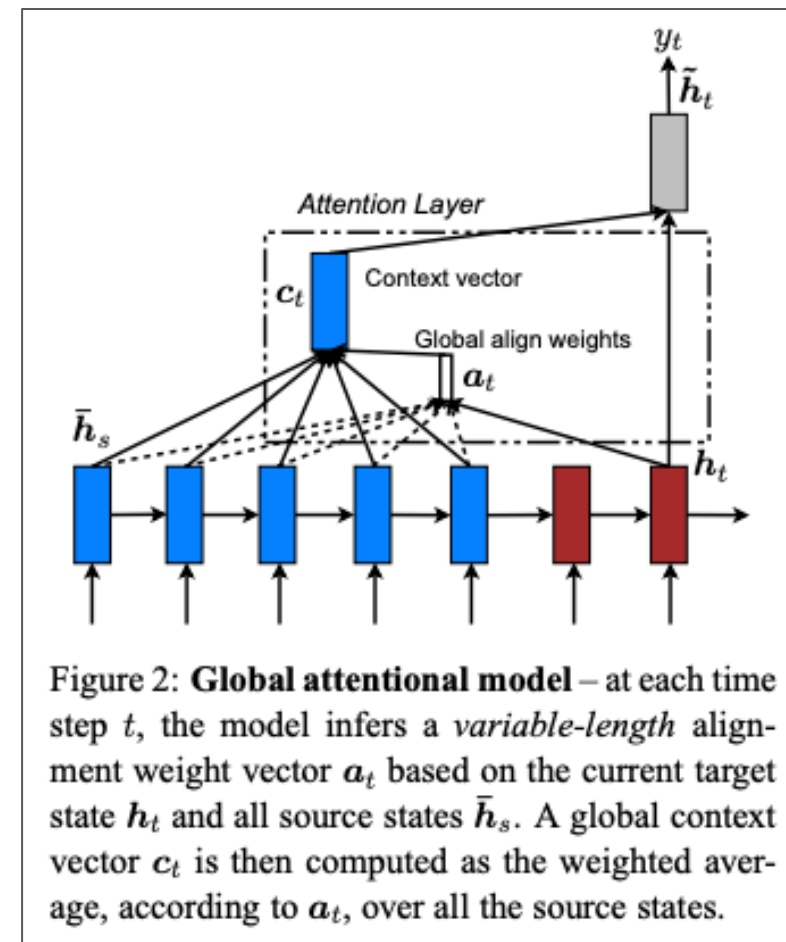
系列変換モデルにおけるアテンション

- 系列変換モデルにおいて、通常、アテンションは、一つのレイヤとして扱われる。
- アテンションレイヤは、エンコーダとデコーダの間に位置付けられる。
- 通常の系列変換モデルでは、エンコーダの最終中間状態だけがデコーダに伝わるが、アテンション付き系列変換モデルでは、全時刻の中間状態を参照する。
 - その際、全時刻の中間状態を等価に参照するのではなく、重要な中間状態をより重視(アテンション)するように参照する。



系列変換モデルにおけるアテンションの仕組み

- 系列変換モデルにおけるアテンションの仕組みを、以下の論文で提案されたグローバル・アテンションモデルを例に説明する。
- 参照論文
 - Minh-Thang Luong, Hieu Pham, Christopher D. Manning. Effective Approaches to Attention-based Neural Machine Translation.
<https://arxiv.org/pdf/1508.04025.pdf>
- 同様の解説が以下の文献にも記載されているため、そちらも参照されたい。
 - 坪井裕太, 海野裕也, 鈴木潤. 深層学習による自然言語処理. 講談社. 4章



参照論文より引用

系列変換モデルにおけるアテンションの仕組み

翻訳タスクであればあれば、例えば x が日本語、 y が英語

- 入力系列を $x = \{x_1, \dots, x_I\}$ 、目標系列を $y = \{y_1, \dots, y_J\}$ とする。
- エンコーダが計算する各時刻の中間状態ベクトルを $\{h_1^{(s)}, \dots, h_i^{(s)}, \dots, h_I^{(s)}\}$ とおく。
- デコーダが計算する各時刻の中間状態ベクトルを $\{h_1^{(t)}, \dots, h_j^{(t)}, \dots, h_J^{(t)}\}$ とおく。
- 系列変換モデルでは、最終的に、 $h_j^{(t)}$ を用いて、 y_j の確率分布を求めたい。

s は source、 t は target を意味する

$$p(y_j | y_{<j}, x) = \text{softmax}(W^{(t)} h_j^{(t)})$$

$W^{(t)}$: パラメータ

softmax: ソフトマックス関数

系列変換モデルにおけるアテンションの仕組み

- ・ エンコーダ側のRNNの遷移関数を $\psi^{(s)}$ とすると、エンコーダの中間状態ベクトルは以下のように再帰的に計算される。

$$h_i^{(s)} = \psi^{(s)}(x_i, h_{i-1}^{(s)}) \quad \psi^{(s)} \text{は、LSTMレイヤなど}$$

- ・ 通常の系列変換モデルでは、最後の中間状態ベクトル $h_I^{(s)}$ のみを使って、デコーダの計算を行う
- ・ この時、例えば、 x_1 の情報は、 $\psi^{(s)}$ が I 回適応される間、ずっと各 $h_i^{(s)}$ で保持される必要がある。
- ・ もっと直接的に入力系列の情報をデコーダ側に伝える方法はないだろうか？
-> これがアテンションの動機

系列変換モデルにおけるアテンションの仕組み

- デコーダの中間状態ベクトル $h_j^{(t)}$ もエンコーダと同様に以下で計算される。

$$h_j^{(t)} = \psi^{(t)}(y_{j-1}, h_{j-1}^{(t)}) \quad \psi^{(t)} \text{は、LSTMレイヤなど}$$

- この後、デコーダは、 $h_j^{(t)}$ を利用して、目標系列の確率分布 $p(y_j | y_{<j}, x)$ を計算するが、ここで、エンコーダが出力した中間状態ベクトル $h_i^{(s)}$ を直接利用することを考える。

系列変換モデルにおけるアテンションの仕組み

- デコーダが j 番目の単語を予測する時に、エンコーダの i 番目の中間状態ベクトル $\mathbf{h}_i^{(s)}$ の重要度を示すスカラー値 a_{ij} を計算し、 a_{ij} による $\mathbf{h}_i^{(s)}$ の重み付き平均 \mathbf{c}_j を計算する。

$$\mathbf{c}_j = \sum_i^I a_{ij} \mathbf{h}_i^{(s)} \quad \text{加重平均された中間状態}$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{i=1}^I \exp(e_{ij})}, a_{ij} \in \mathbb{R} \quad \text{中間状態の重要度}$$

$$e_{ij} = \Omega(\mathbf{h}_i^{(s)}, \mathbf{h}_j^{(t)}), e_{ij} \in \mathbb{R}$$

- 例えば、3番目の入力単語 \mathbf{x}_3 だけが重要であれば、 a_{3j} が1に近く、それ以外の a_{ij} がほとんど0になる。
- a_{ij} は、対象となる複数のベクトルの中から、重要な情報を選別するような役割を担うことになる。

系列変換モデルにおけるアテンションの仕組み

- 関数 Ω について、参照論文では以下のような複数の関数を検討している。

$$\Omega(h_i^{(s)}, h_j^{(t)}) = \begin{cases} h_i^{(s)} \cdot h_j^{(t)} \\ h_i^{(s)} \cdot W h_j^{(t)} \\ v \cdot \tanh(W[h_i^{(s)}; h_j^{(t)}]) \end{cases}$$

W, v : パラメータ

$[\cdot]$: 結合(concatenate)の記号

内積

デコーダ側中間状態ベクトルに重みをかけてから内積

エンコーダ側中間状態ベクトルとデコーダ側中間状態ベクトルを結合してから重みをかけて、 \tanh を通し、重みベクトル v と内積

系列変換モデルにおけるアテンションの仕組み

- デコーダが j 番目の単語を予測する際、 $h_j^{(t)}$ の代わりに以下で定義される $\tilde{h}_j^{(t)}$ を利用する。

$$\tilde{h}_j^{(t)} = \tanh(W_c [c_j; h_j^{(t)}])$$

W_c : パラメータ

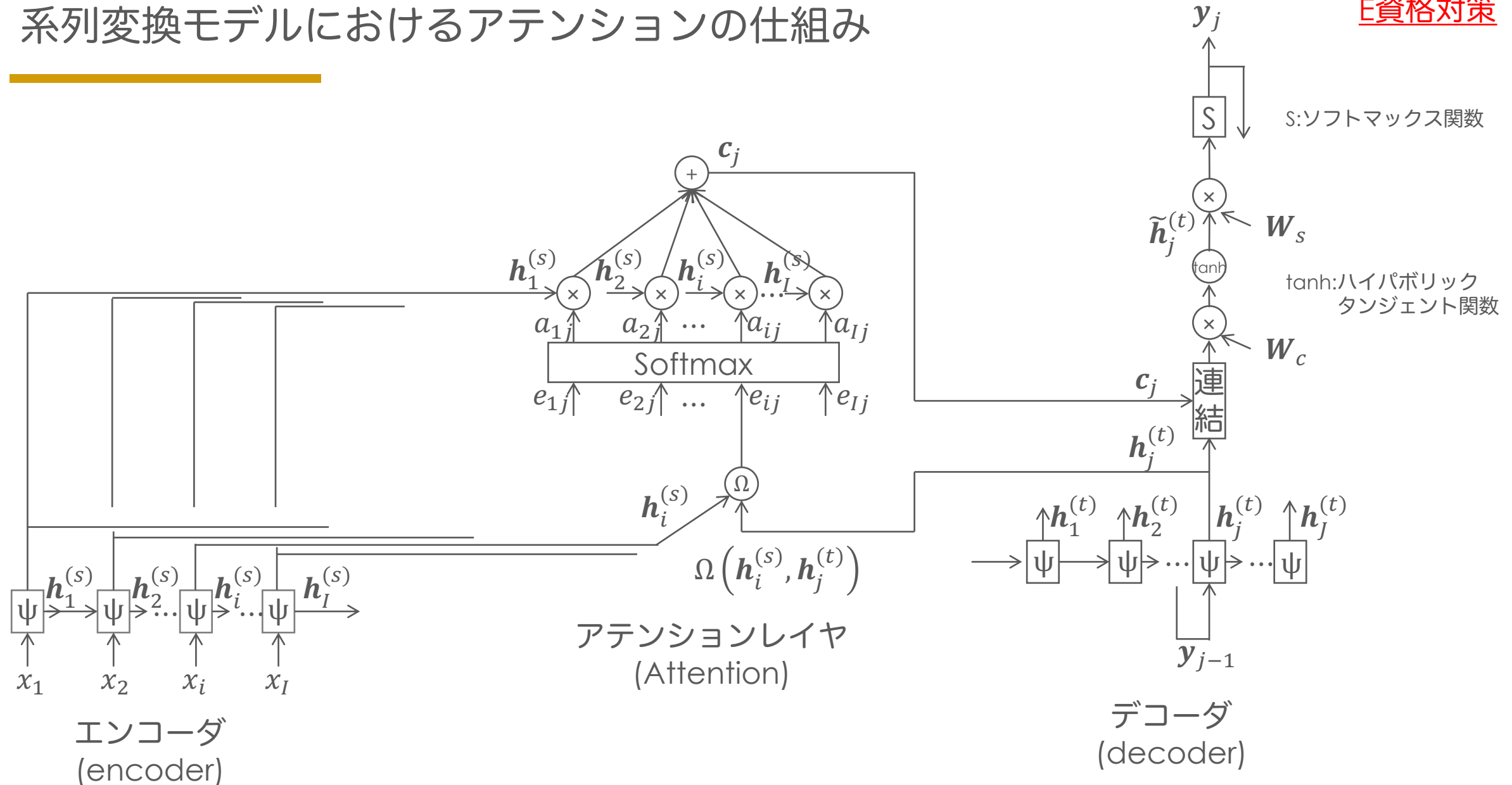
$[\cdot]$: 結合(concatenate)の記号

- デコーダが j 番目の単語を予測する際に、 $\tilde{h}_j^{(t)}$ を利用する。

$$p(y_j | y_{<j}, x) = \text{softmax}(W_s \tilde{h}_j^{(t)})$$

- これにより、エンコーダの全ての中間状態をデコーダで利用できることになる。
- このアテンションでは、全て微分可能な関数を用いているため、誤差逆伝播法で勾配を計算することが可能。

系列変換モデルにおけるアテンションの仕組み



アテンションを活用した事例

- アテンションというアイデアは、汎用的なものであり、様々なモデルで取り入れられている
 - 自然言語処理にアテンションを使った事例
 - Google's Neural Machine Translation System(GNMT)
 - トランスフォーマー
 - 外部メモリを持つニューラルネットワーク
- コンピュータビジョン分野では、ビジュアルアテンションという考え方がある
 - 例えば、SENetが参考になる
 - Jie Hu, Li Shen, Samuel Albanie, Gang Sun, Enhua Wu. Squeeze-and-Excitation Networks. CVPR 2018. <https://arxiv.org/pdf/1709.01507.pdf>

Google's Neural Machine Translation System

- Google's Neural Machine Translation System(GNMT)は、Google社が開発したニューラル機械翻訳のシステム
- 実際のサービスとして、2016年から使用されている。
- 論文
 - Yonghui Wu et al. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation.
<https://arxiv.org/pdf/1609.08144.pdf>
- Google社公式ブログ
 - <https://ai.googleblog.com/2016/09/a-neural-network-for-machine.html>

ニューラル機械翻訳とは、ニューラルネットワークモデルを用いて機械翻訳を行うこと

Google's Neural Machine Translation System

- GNMTは、エンコーダ、デコーダ、アテンションから構成されている。
- GNMTで取り入れられた工夫
 - LSTMレイヤの多層化
 - 双方向LSTM(エンコーダの1層目のみ)
 - スキップコネクション
 - 複数GPUでの分散学習
 - 予測の高速化 (量子化)

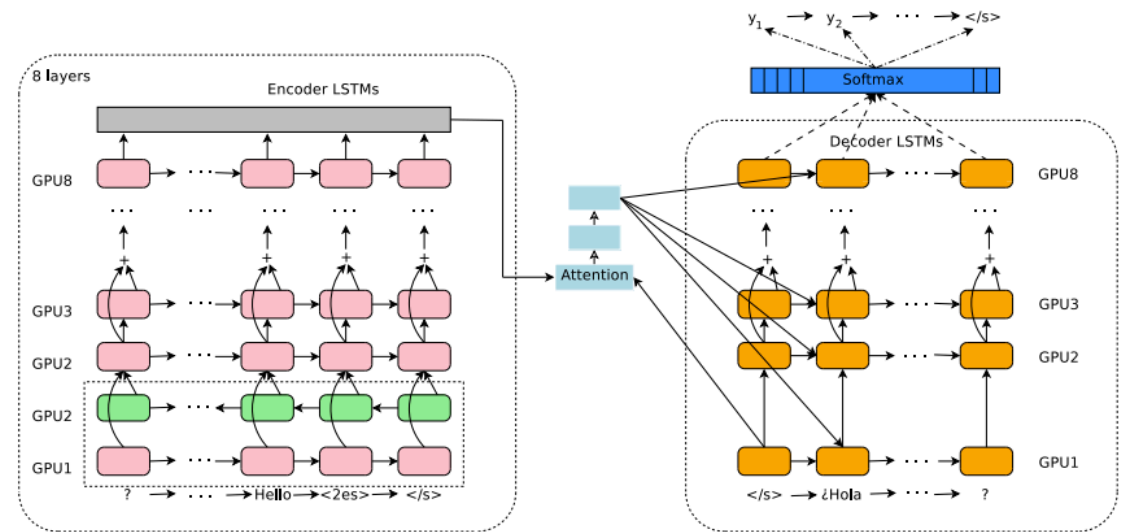
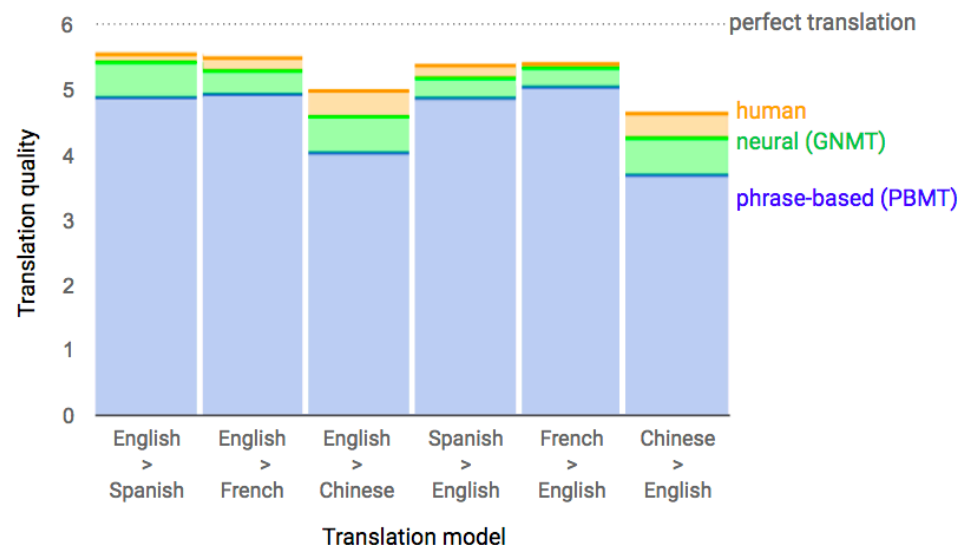


Figure 1: The model architecture of the Multilingual GNMT system. In addition to what is described in [24], our input has an artificial token to indicate the required target language. In this example, the token “<2es>” indicates that the target sentence is in Spanish, and the source sentence is reversed as a processing step. For most of our experiments we also used direct connections between the encoder and decoder although we later found out that the effect of these connections is negligible (however, once you train with those they have to be present for inference as well). The rest of the model architecture is the same as in [24].

Google's Neural Machine Translation System

- GNMTの精度評価の結果を以下に引用する。
 - 従来のフレーズベース機械翻訳に比べ、**人による翻訳に近い精度**を実現できていることがわかる。
 - PBMTは、フレーズベース機械翻訳の略。フレーズベース機械翻訳とは、文章を小さな単位に分割して翻訳し、それらを並べ替えることによって、翻訳文を生成する手法。



左図は、<https://ai.googleblog.com/2016/09/a-neural-network-for-machine.html>より引用

Data from side-by-side evaluations, where *human raters compare the quality of translations for a given source sentence*. Scores range from 0 to 6, with 0 meaning “completely nonsense translation”, and 6 meaning “perfect translation”.

[演習] アテンション付きseq2seqモデルの実装

- 7_4_TimeAttention_trainee.ipynb
 - Attentionレイヤとそれを時間方向に束ねるTimeAttentionレイヤを実装しましょう。
- 7_5_AttentionSeq2seq_trainee.ipynb
 - アテンション付きseq2seqを計算するためのクラスを実装しましょう。
 - 計算グラフは、「系列変換モデルにおけるアテンションの仕組み」の頁を参考にしましょう。

[演習] 双方向LSTMの実装

- 7_6_TimeBiLSTM_trainee.ipynb
 - 双方向LSTMを計算するためのTimeBiLSTMレイヤを実装しましょう。
- 7_7_AttentionBiSeq2seq_trainee.ipynb
 - エンコーダ側を双方向LSTMにしたアテンション付きseq2seqモデルを計算するためのクラスを実装しましょう。
 - 計算グラフは、「系列変換モデルにおけるアテンションの仕組み」の頁を参考にしましょう。

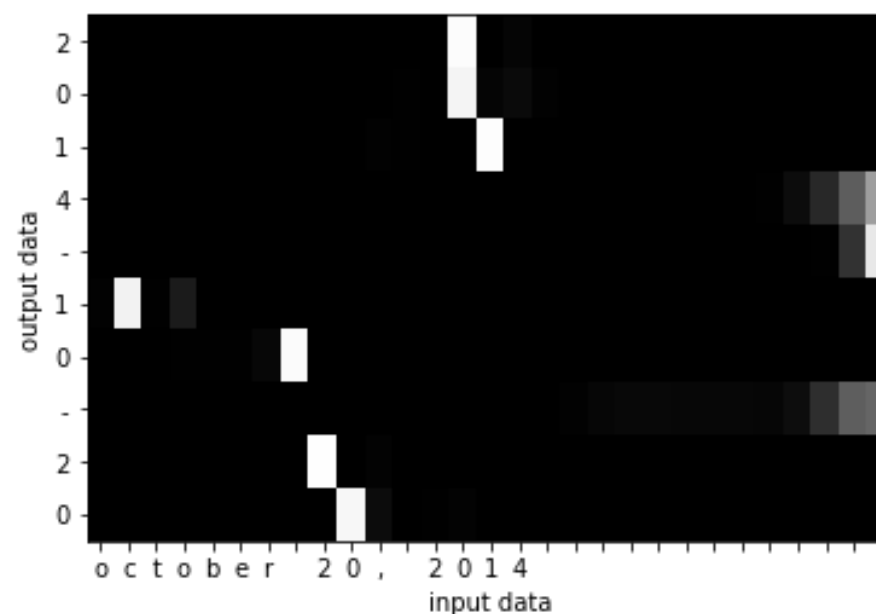
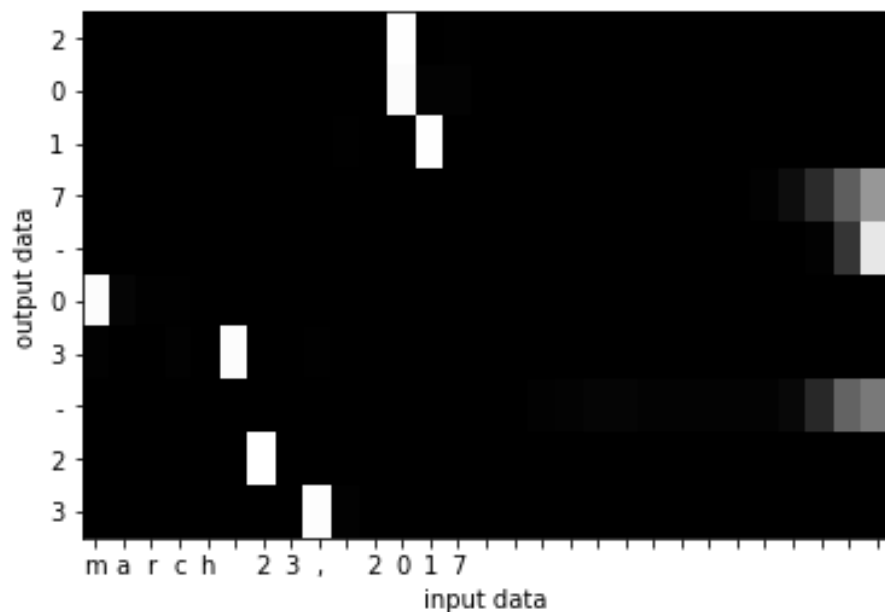
[演習] 系列変換モデルを用いた学習

- 7_8_train_with_seq-to-seq.ipynb
 - 系列変換モデルを用いた学習を行きましょう。
 - seq2seqモデル
 - アテンション付きのseq2seqモデル
 - エンコーダ側を双方向LSTMにしたアテンション付きseq2seqモデル
 - ここでのタスクは、以下のような日付フォーマット変換です。
 - september 27, 1994 -> _1994-09-27
 - August 19, 2003 -> _2003-08-19
 - 2/10/93 -> _1993-02-10
 - 10/31/90 -> _1990-10-31

[演習] アテンションの可視化

- 7_9_visual_attention.ipynb

- アテンション付きseq2seqの学習済みモデルを用いて、アテンションの結果を可視化しましょう。
- ある単語を予測する際に、その単語と関連する入力単語に注目(アテンション荷重が大きい)していたらアテンションの学習がうまくいっています。



白い部分は、アテンション荷重が大きい(1に近い)

Any Questions ?

トランスフォーマー

RNNの課題

- これまで、RNNを用いると可変長の時系列データをうまく扱えることを確認してきた。
- しかし、RNNには計算効率面で課題がある。
 - 前時刻に計算した結果を用いながら逐次的に計算を行なっていくため、計算に時間がかかる。
- こういった背景から、RNNを使わずにRNNと同等のことができないかという研究が行われている。
- その一つの成果が、トランスフォーマーである。

トランスフォーマー

- トランスフォーマーは、下記論文で提案された、RNNを使わずにアテンションのみで機械翻訳問題を解くためのモデル。
- 圧倒的な学習効率で、翻訳タスクにおいてSotA(当時最高性能)を達成。
- これ以降、機械翻訳問題はトランスフォーマーが主流になっていった。
- 原著論文
 - Ashish Vaswani et al. Attention is All You Need, NIPS 2017.
<https://arxiv.org/pdf/1706.03762.pdf>
- この論文では、アテンションに関する明瞭な定式化・解釈を与えている。
 - アテンション=辞書オブジェクト (Query, Key, Valueモデル)

アテンションの一般化

- アテンションを辞書(dictionary)オブジェクトとして解釈。
 - 質問文(Query)を入れると、参照すべき場所(Key)が決まり、その場所の値 (Value) が得られる。
 - ソフト・アテンションでは、参照すべき場所(Key)の重みが決まり、その重みに応じて値(Value)を平均化することで回答を得ていることになる。
 - KeyとValueのセットが事前知識によって構築されるものであり、メモリに相当する。

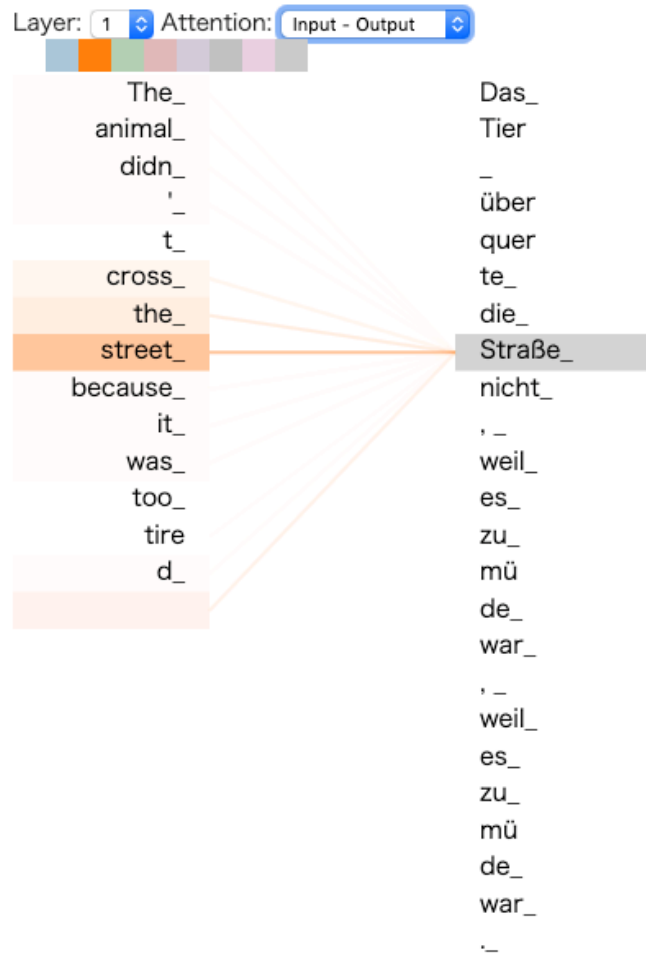
$$\text{ソフトアテンションにおける回答} = \text{Softmax}(QK^T) \cdot V$$

ソースターゲット・アテンションとセルフ・アテンション

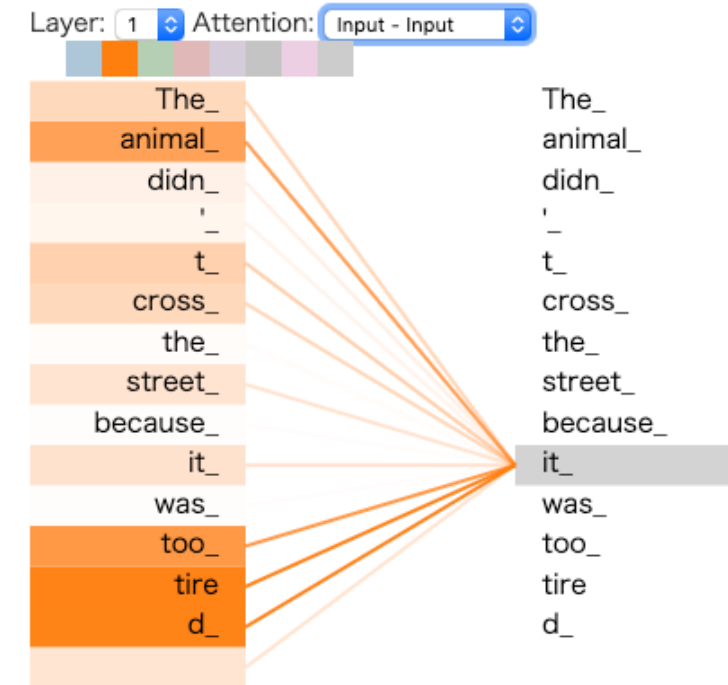
- ソースターゲット・アテンション
 - seq2seqモデルで用いていたアテンション。
 - 2つの系列間の対応関係を捉えることが目的。
 - Queryはデコーダ側、KeyとValueはエンコーダ側に位置付けられる。
- セルフ・アテンション
 - 文章内の単語間の関係を捉えることが目的。
 - Query、Key、Valueが同じ単語から生成されていく。
 - 各単語が文章中のすべての単語と比較され、その結果が文章内の全ての単語のアテンションスコア(重み)になる。
- トランスフォーマーでは、ソースターゲット・アテンションとセルフ・アテンションの両方が用いられている。

ソースターゲット・アテンションとセルフ・アテンション

ソースターゲット・アテンションの例



セルフ・アテンションの例



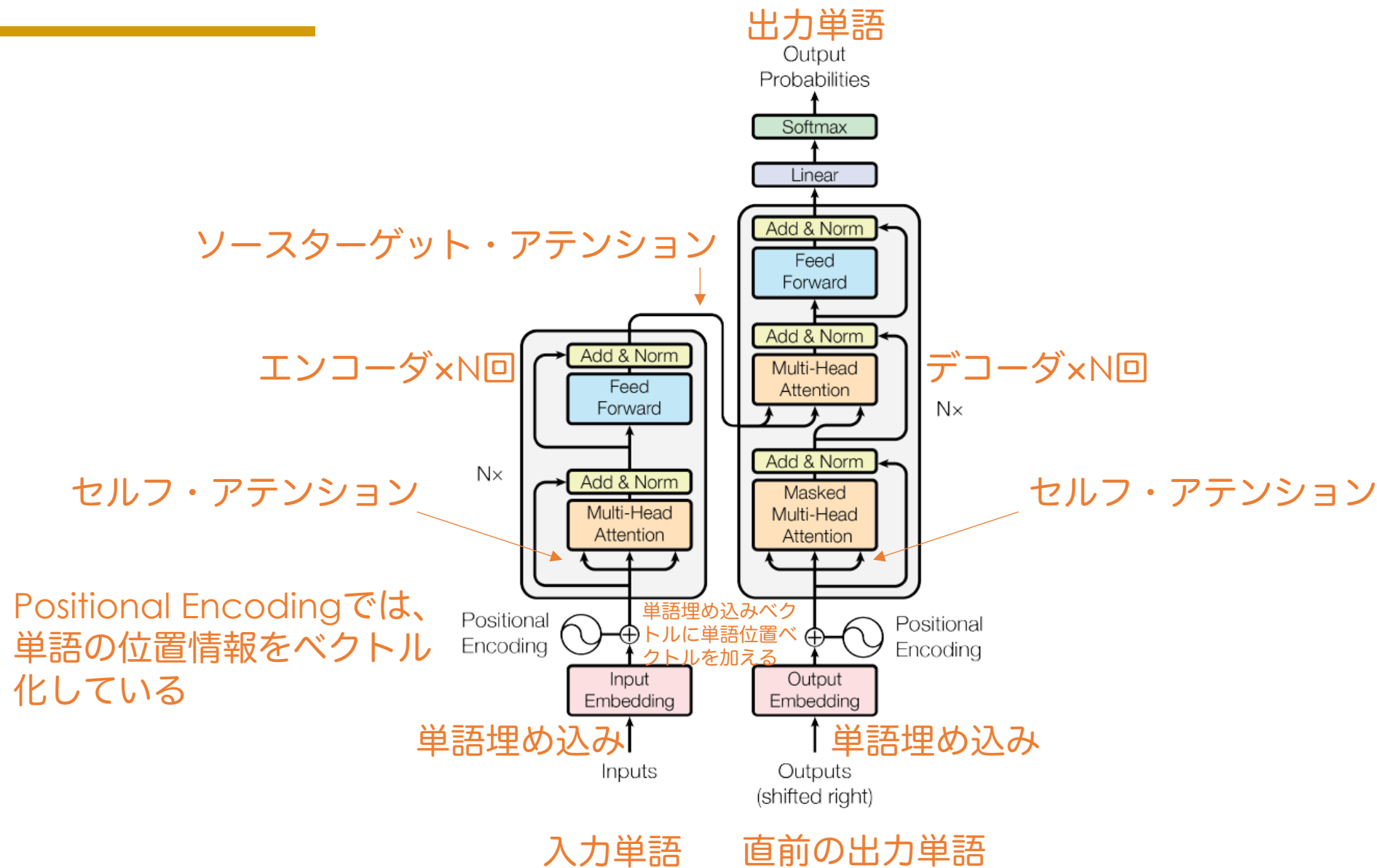
Tensor2Tensor Intro

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb#scrollTo=OJKU36QAfqOC

トランスフォーマーのネットワーク構成

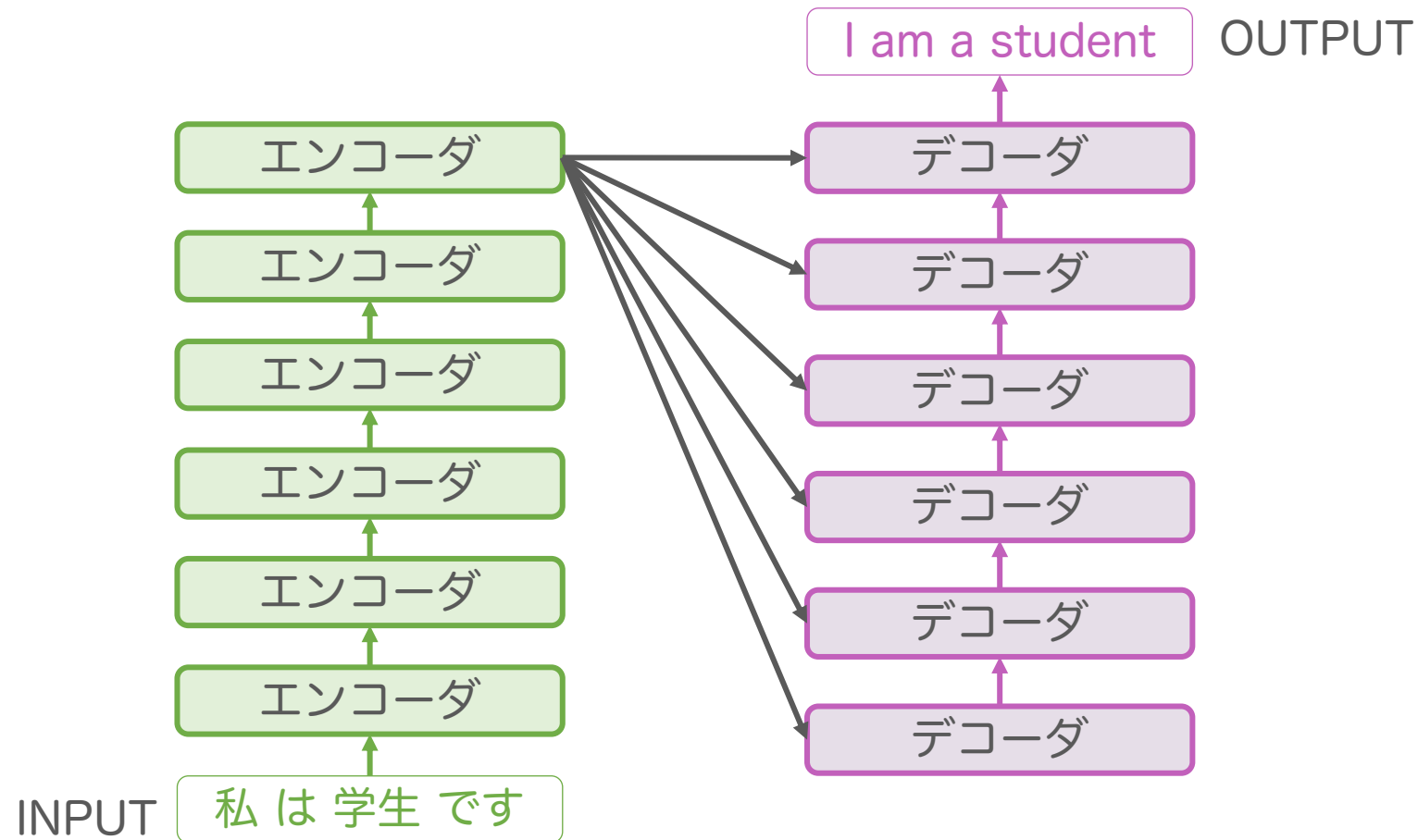
- 次頁からトランスフォーマーのネットワーク構成を見ていく。
- トランスフォーマーの仕組みは、以下のブログが非常に参考になる。
 - <http://jalammar.github.io/illustrated-transformer/>
- Pythonでの実装については、こちらが参考になる。
 - <https://qiita.com/halhorn/items/c91497522be27bde17ce>

トランスフォーマーのネットワーク構成



トランスフォーマーのネットワーク構成

- エンコーダおよびデコーダはN回(論文では6回)繰り返される。



トランスフォーマーのネットワーク構成

$$\mathbf{Z}_i = \text{softmax}\left(\frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_k}}\right) \mathbf{V}_i$$

d_k : keyベクトルの次元数

• Multi-head attentionの内部構成

1) 入力
単語

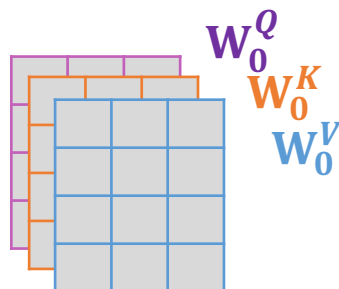
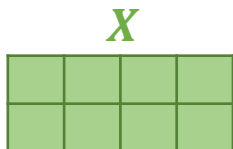
2) 単語を埋
め込む

3) ヘッド毎(論文では8ヘッド)に \mathbf{X} と \mathbf{W}_i を掛けることで、
単語埋め込みを $\mathbf{Q}_i \mathbf{K}_i \mathbf{V}_i$ に変換する

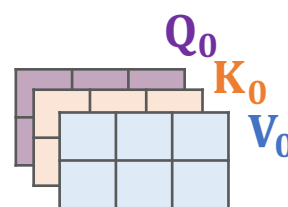
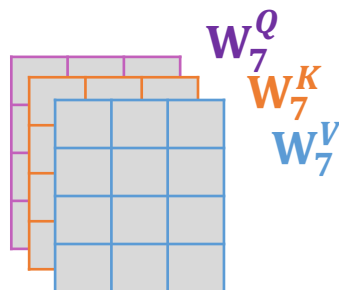
4) $\mathbf{Q}_i \mathbf{K}_i \mathbf{V}_i$ の結果を用いて、 \mathbf{Z}_i を計算する

5) 全ての \mathbf{Z}_i を列方向に結合し、 \mathbf{W}^o と掛ける

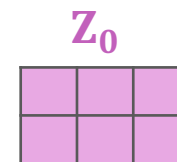
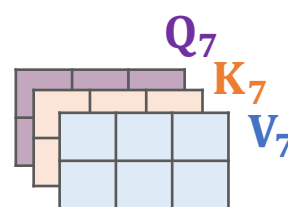
私
は



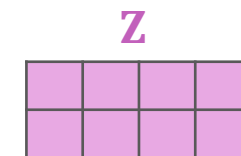
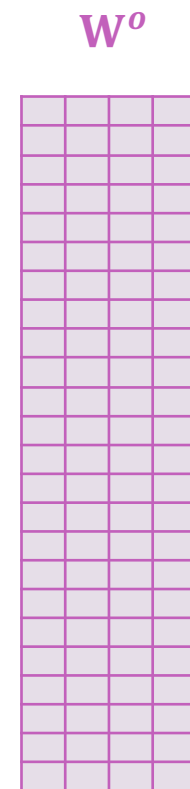
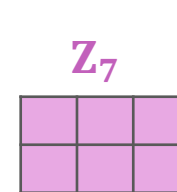
...



...



...



論文では、1つの大きなアテンションを行うよりも、小さな複数のヘッドに分けてアテンションを行うほうが性能がよかったと報告されている。

トランスフォーマーのネットワーク構成

- 単語の順番を考慮するために、ポジショナル・エンコーディング(Positional encoding)を導入する

pos は、その単語が文章の何番目かを表す数字

d_{model} は、単語埋め込みの次元数

PE は、ポジショナル・エンコーディングの計算結果

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

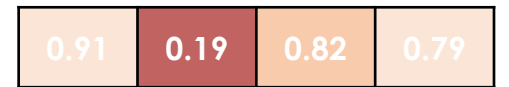
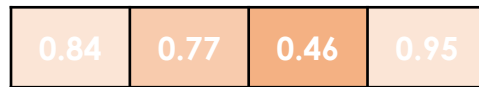
偶数番目の値はsin関数の方を用い、
奇数番目の値はcos関数の方を用いる

where pos is the position and i is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from 2π to $10000 \cdot 2\pi$. We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} .

トランスフォーマーのネットワーク構成

- ポジショナル・エンコーディングの例を以下に示す。

POSITIONAL
ENCODING



EMBEDDINGS



INPUT

私

は

学生

[演習]トランスフォーマーにおけるPositional Encoding

- 7_10_Positional_encoding.ipynb
 - トランスフォーマーにおけるPositional Encodingの計算方法を確認しましょう。

BERT

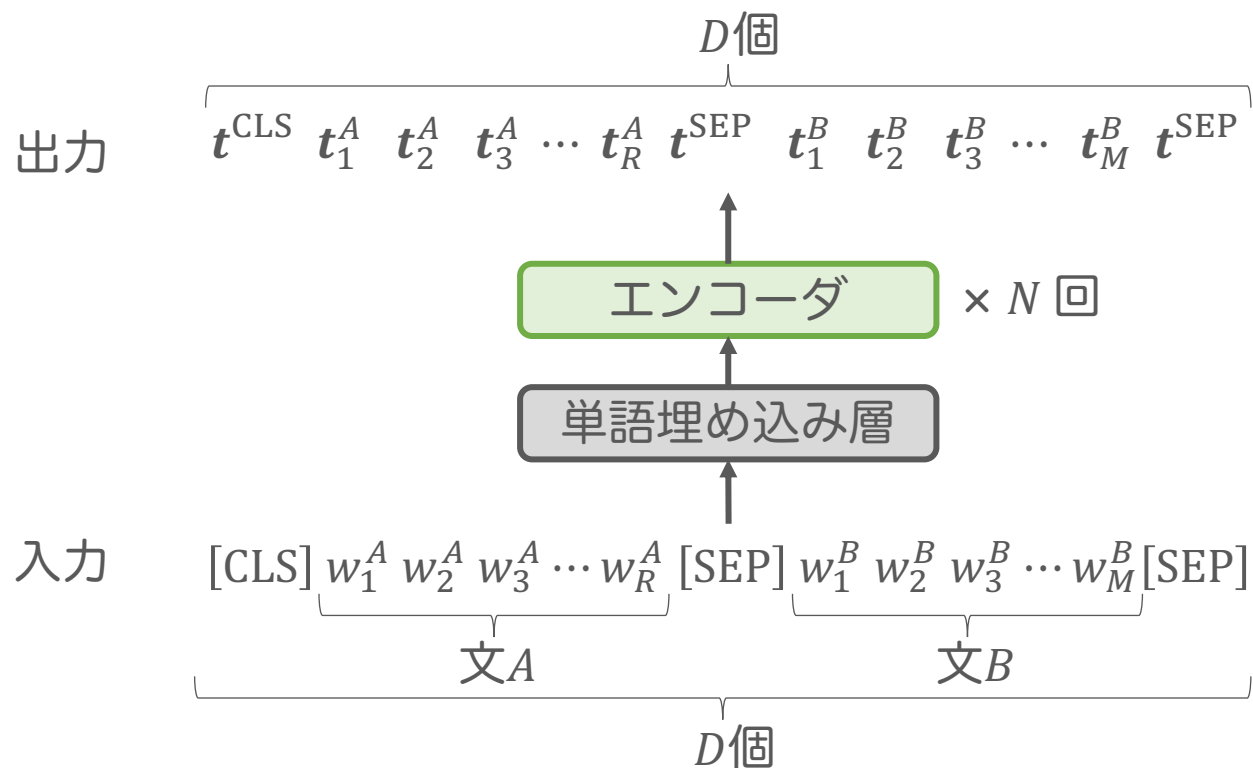
- BERT(Bidirectional Encoders' Representations from Transformer)とは、以下の論文で提案された自然言語処理用モデル。
 - Jacob Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. (Arxiv preprint, from Google, Oct. 2018).
<https://arxiv.org/pdf/1810.04805.pdf>
 - モデルには、トランスフォーマーのエンコーダ部分が用いられている。
- 感情分析・含意関係認識・質疑応答・意味等価性判別など、**11のタスクでSotA(当時最高性能)**。
 - 後に提案された以下の論文では翻訳タスクでもStoA。
 - Guillaume Lample, Alexis Conneau. Cross-lingual Language Model Pretraining.
<https://arxiv.org/pdf/1901.07291.pdf>
- **汎用的な言語表現(分散表現)を獲得できるモデルと言われている。**
 - word2vecよりも豊かな言語表現(分散表現)。
- 自然言語処理の決定的手法になるかもしれないと話題の手法。

BERTの学習手順

- BERTでは、2段階で学習を行う。
 - 事前学習により言語表現(分散表現)を獲得し、再学習によって特定のタスクへ適応させる。
- BERTの学習手順
 1. 事前学習(pre-training)
 - 単語マスク問題(Masked Language Modeling, MLM)と隣接文問題(Next Sentence Prediction, NSP)を同時に解く。
 - 教師なし学習。
 2. 再学習(fine-tuning)
 - 事前学習済みのモデルの最終層を解きたいタスク用の層に入れ替え、全ての重みを更新する。
 - これにより、特定のタスクを解けるようになる。
 - 教師あり学習。

BERTの基本構成

- BERTの基本構成を以下に示す。
- BERTは、**トランスフォーマーのエンコーダ**を利用しており、入力系列に対応する特徴量 t の系列を出力する。



t : 特徴量 $t \in \mathbb{R}^H$

w : 単語ID

H : 中間状態のサイズ(768など)

D : 1度に入力できる単語の総数 (512など)

N : エンコーダの繰り返し回数(12など)

R : 文Aの単語の個数

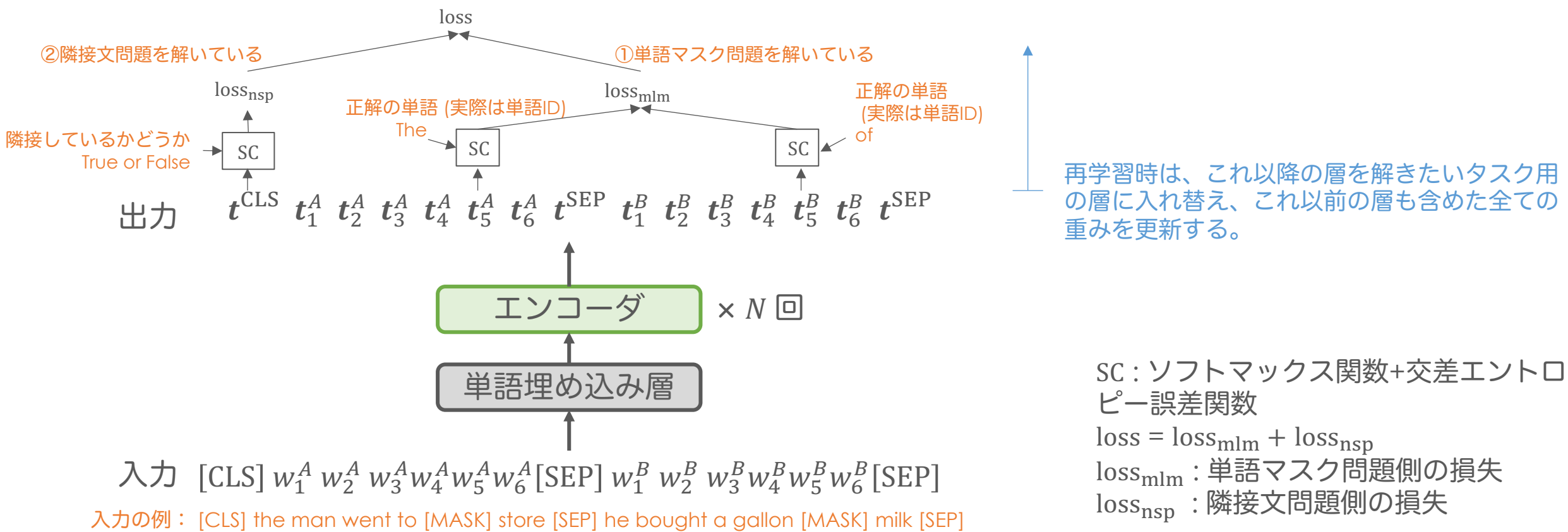
M : 文Bの単語の個数

$[CLS]$: 入力系列の先頭に必ず配置する記号

$[SEP]$: 文の区切りに配置する記号

BERTにおける事前学習

- ① 単語マスク問題: 入力系列のうち、隠された単語がなにかを予測(局所的な特徴学習)
- ② 隣接文問題: 2つの入力文が隣り合うものかどうかを判別(大域的な特徴学習)



BERTの学習済みモデル

<なぜ、双方向(bidirectional)トランスフォーマーという名称が使われたのか？>
BERT以前のモデルであるOpenAI GPTでは、トランスフォーマーのデコーダを用いて、1期先単語を予測するタスクを解くことによって言語表現を獲得していた。このモデルでは過去の単語だけを使って将来の単語を予測しているため、OpenAI GPTを単方向(unidirectional)なモデルと呼ぶことにした。これに対しBERTは、トランスフォーマーのエンコーダを用いて、過去の単語も将来の単語も使って、[mask]の単語を予測するので、双方向(bidirectional)なモデルと呼ぶことにした。

- 英語の学習済みモデル

- 公式リポジトリ

- <https://github.com/google-research/bert>

- 日本語の学習済みモデル

- 京都大学 黒橋・村脇研究室

- <http://nlp.ist.i.kyoto-u.ac.jp/index.php?BERT%E6%97%A5%E6%9C%AC%E8%AA%9EPretrained%E3%83%A2%E3%83%87%E3%83%AB>

- 情報通信研究機構 データ駆動知能システム研究センター

- <https://alaginrc.nict.go.jp/nict-bert/index.html>

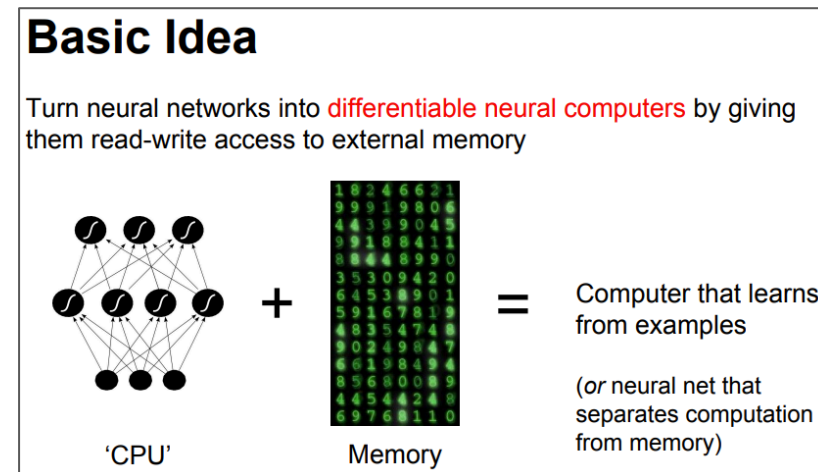
Any Questions ?

外部メモリを持つニューラルネットワーク

外部メモリを持つニューラルネットワークとは

- 近年、外部メモリを持つニューラルネットワークがいくつか提案されている。
- 外部メモリを持つNNでは、通常のLSTMなどとは異なり、メモリとモデルを切り離して考える。これは、PCにおいて、CPUとメモリが別々に存在していることに例えられる。
- メモリを切り離すことにより、通常のLSTMなどよりも、表現力を豊かにできたり(必要な箇所のみ更新または取り出しができる)、多くの情報を蓄えられるようになる。
- メモリに対して、必要な情報を読み書きするアドレスを決める機構のことをアテンション (attention) という。

引用：Differentiable Neural Computers Alex Graves, Greg Wayne & many others
<http://people.idsia.ch/~rupesh/rnnsymposium2016/slides/graves.pdf> (リンク切れ)



外部メモリを持つニューラルネットワークの例

- 外部メモリを持つニューラルネットワークの例
 - Jason Weston et. al., Memory Networks, 2014
 - <https://arxiv.org/pdf/1410.3916.pdf>
 - <http://www.thespermwhale.com/jaseweston/icml2016/icml2016-memnn-tutorial.pdf>
 - Graves et. al., Neural Turing Machines, 2014
 - <https://arxiv.org/pdf/1410.5401.pdf>
 - Ankit Kumar et. al., Ask Me Anything: Dynamic Memory Networks for Natural Language Processing, 2015
 - <https://arxiv.org/pdf/1506.07285.pdf>
 - Sainbayar Sukhbaatar ,End-To-End Memory Networks, 2015
 - <https://arxiv.org/pdf/1503.08895.pdf>
 - Alex Graves et. al., Differentiable Neural Computers, 2016
 - Neural Turing Machine (NTM) の後継。より複雑なタスクを解けるようになった。
 - <https://www.nature.com/articles/nature20101>
 - https://www.nature.com/articles/nature20101.epdf?author_access_token=lmTXBI8aWbYxYQ51Plys8NRgN0jAjWel9jnR3ZoTv0MggmpDmwIjGswxVdeocYSurJ3hxupzWuRNeGvvXnoO8o4jTJcnAyhGuZzXJ1GEaD-Z7E6X_a9R-xqJ9TfJWBqz

外部メモリを持つニューラルネットワークでできることの例

- 外部メモリを持つニューラルネットワークでできること

Joe went to the kitchen. Fred went to the kitchen. Joe picked up the milk.
Joe travelled to the office. Joe left the milk. Joe went to the bathroom.
Where is the milk now? **A: office**
Where is Joe? **A: bathroom**
Where was Joe before the office? **A: kitchen**

Sheep are afraid of wolves.
Cats are afraid of dogs.
Mice are afraid of cats.
Gertrude is a sheep.
What is Gertrude afraid of? **A:wolves**

引用：Jason Weston et. al., Memory Networks, 2014

上：<https://arxiv.org/pdf/1410.3916.pdf>

下：

<http://www.thespermwhale.com/jaseweston/icml2016/icml2016-memnn-tutorial.pdf>

外部メモリを持つニューラルネットワークの精度比較

- 各手法の精度比較を以下に示す。

NTM:Neural Turing Machines

DNC : Differentiable Neural Computers

MemN2N : End-To-End Memory Networks

DMN : Dynamic Memory Networks

bAbI Results

Task	bAbI Best Results						
	LSTM (Joint)	NTM (Joint)	DNC1 (Joint)	DNC2 (Joint)	MemN2N (Joint) ²¹	MemN2N (Single) ²¹	DMN (Single) ²⁰
1: 1 supporting fact	24.5	31.5	0.0	0.0	0.0	0.0	0.0
2: 2 supporting facts	53.2	54.5	1.3	0.4	1.0	0.3	1.8
3: 3 supporting facts	48.3	43.9	2.4	1.8	6.8	2.1	4.8
4: 2 argument rels.	0.4	0.0	0.0	0.0	0.0	0.0	0.0
5: 3 argument rels.	3.5	0.8	0.5	0.8	6.1	0.8	0.7
6: yes/no questions	11.5	17.1	0.0	0.0	0.1	0.1	0.0
7: counting	15.0	17.8	0.2	0.6	6.6	2.0	3.1
8: lists/sets	16.5	13.8	0.1	0.3	2.7	0.9	3.5
9: simple negation	10.5	16.4	0.0	0.2	0.0	0.3	0.0
10: indefinite knowl.	22.9	16.6	0.2	0.2	0.5	0.0	0.0
11: basic coreference	6.1	15.2	0.0	0.0	0.0	0.1	0.1
12: conjunction	3.8	8.9	0.1	0.0	0.1	0.0	0.0
13: compound coref.	0.5	7.4	0.0	0.1	0.0	0.0	0.2
14: time reasoning	55.3	24.2	0.3	0.4	0.0	0.1	0.0
15: basic deduction	44.7	47.0	0.0	0.0	0.2	0.0	0.0
16: basic induction	52.6	53.6	52.4	55.1	0.2	51.8	0.6
17: positional reas.	39.2	25.5	24.1	12.0	41.8	18.6	40.4
18: size reasoning	4.8	2.2	4.0	0.8	8.0	5.3	4.7
19: path finding	89.5	4.3	0.1	3.9	75.7	2.3	65.5
20: agent motiv.	1.3	1.5	0.0	0.0	0.0	0.0	0.0
Mean Err. (%)	25.2	20.1	4.3	3.8	7.5	4.2	6.4
Failed (err. > 5%)	15	16	2	2	6	3	2

引用：Differentiable Neural Computers Alex Graves, Greg Wayne & many others

<http://people.idsia.ch/~rupesh/rnnsymposium2016/slides/graves.pdf> (リンク切れ)

bAbIプロジェクトとは

- bABIプロジェクトとは、文章理解や文章意味付けの自動化を目指すプロジェクト。
- Facebook AI Research が主導している。
- このプロジェクトでは、言語理解を測るためのテキストデータセットを提供している (bAbIタスク)。
- <https://research.fb.com/downloads/babi/>

外部メモリを持つニューラルネットワーク

- 外部メモリを持つニューラルネットワークの代表的なモデルとして、Memory NetworksとEnd-to-End MemNN(MemN2N)を紹介する。

Memory Networks

- Memory Networksは、**記憶とそれ呼び起こす仕組み**の一般的なモデル。
- どのような構成で質疑応答システムに適したモデルが作れるかを提案している。
- Memory Networksの各コンポーネントにニューラルネットワークを用いたものをMemory Neural Networks(MemNN)という。
- 原著論文
 - Jason Weston, Sumit Chopra, Antoine Bordes. Memory Networks. ICLR 2015.
<https://arxiv.org/pdf/1410.3916.pdf>

Memory Networks

- 質疑応答システムをDNNで実現する方法を考える。
- 質疑応答では事前に与えられた知識を何らかの形で蓄えておかなければならない。

事前知識=学習データ (事実; fact)

机の上には美味しそうないちごとバナナが置いてあります。
太郎君はバナナを食べたそうにしています。
でも、太郎君はお母さんに止められているのでバナナを食べられません。
なぜなら太郎君は先週の健康診断で体重が増えてしまっていたからです。

質問(入力文)

Q. 机の上にはバナナの他に何がおいてありますか？

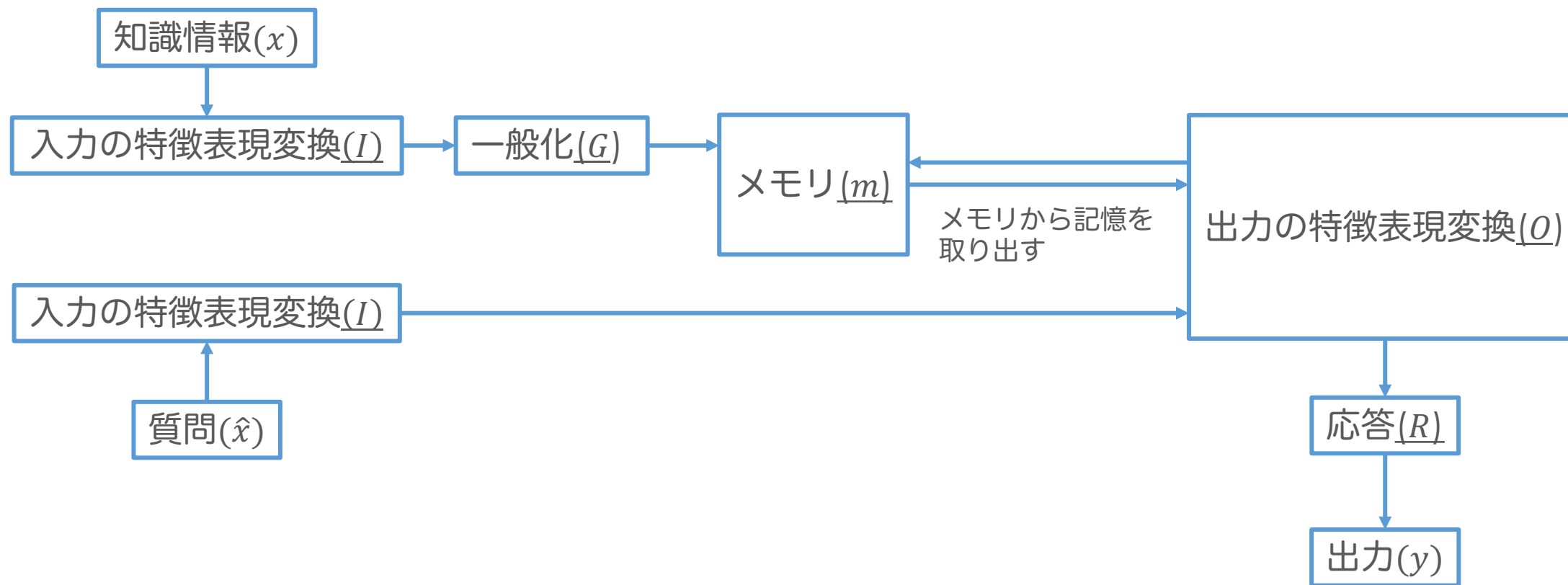
回答(出力)

A. いちご

Memory Networks概観

- Memory Networksの主な構成要素は5つ
 - メモリ(m): 知識(入力された文、あるいはその特徴)を蓄えておく配列
 - 入力の特徴表現変換(I): 入力された知識を特徴表現に変換する部分
 - 入力の一般化(G): 新たな入力を用いてメモリMを更新する
 - 出力の特徴表現変換(O): 質問から回答に対応する特徴表現を作成する部分
 - 応答(R): O を用いて実際に文としての回答を生成する部分
- 知識情報: 入力の知識文 x
 - 特徴表現に変換 $h = I(x)$, 知識としてメモリに格納: $m_i = h$
- 質問: 入力の質問文 \hat{x}
 - 質問文を特徴表現に変換 $\hat{h} = I(\hat{x})$, メモリから知識を引いて回答 $y = R(O(\hat{h}, m))$

Memory Networksの構成要素



I, G, O, Rモジュール

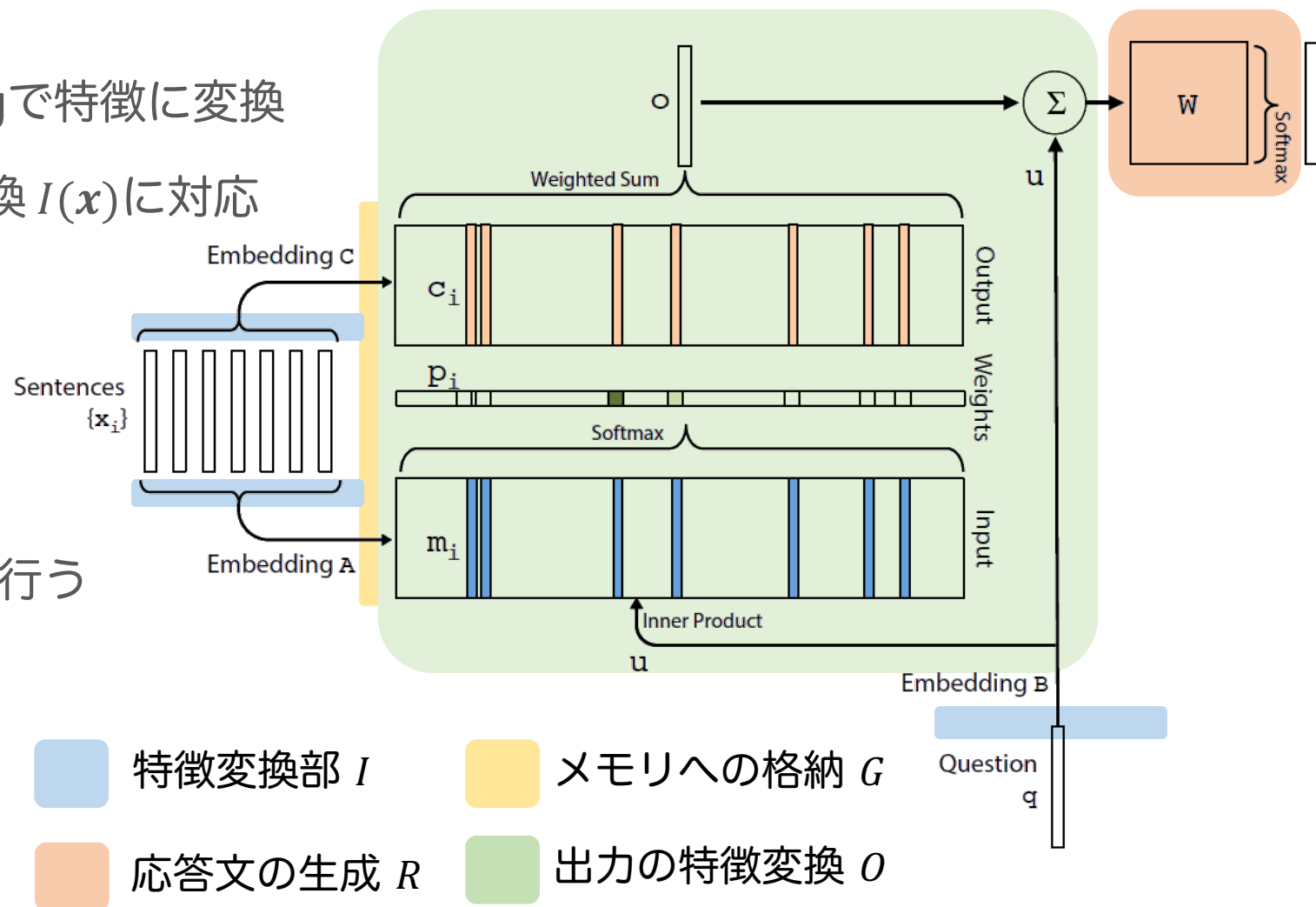
- I (Input)モジュールでは生の文字列を特徴表現に変換する。
 - Bag of Words, Word2Vecなど、なんでもよい
 - この I によってモデルが行う文章の表現が定まる
- G (Generalize)モジュールでは知識をメモリ m に格納する。
 - m は特徴表現 $I(x)$ を値に持つただの配列 $m = (m_i)_{i \in I}$
 - G モジュールは m の格納位置を決定関数 $H(x)$ で決めて代入(保存)
 - $m_{H(x)} = I(x)$
 - 最も単純には後ろに追加していくだけの格納:
$$m_{H(x)} = I(x) \text{ where } H(x) = N, N = N + 1$$
- O モジュールは格納された知識の中から、上位 k 個の回答に役立つメモリ m_i を見つけてくる。
- R モジュールは求められた $O(x)$ から回答文 r を生成する

End-To-End Memory Networks

- Memory Networksの考え方を微分可能なニューラルネットワークで実現したモデル
 - End-to-End MemNN => MemN2N と略される
- MemNNでは、ハード・アテンションでメモリから記憶を取り出すことを考えていた。
 - MemN2Nでは、ソフト・アテンションを用いる。これにより、モデル全体を1つのニューラルネットワークで構築することが可能になった。
- 原著論文
 - End-To-End Memory Networks. Sainbayar Sukhbaatar et al. NIPS 2015.
<https://arxiv.org/pdf/1503.08895.pdf>

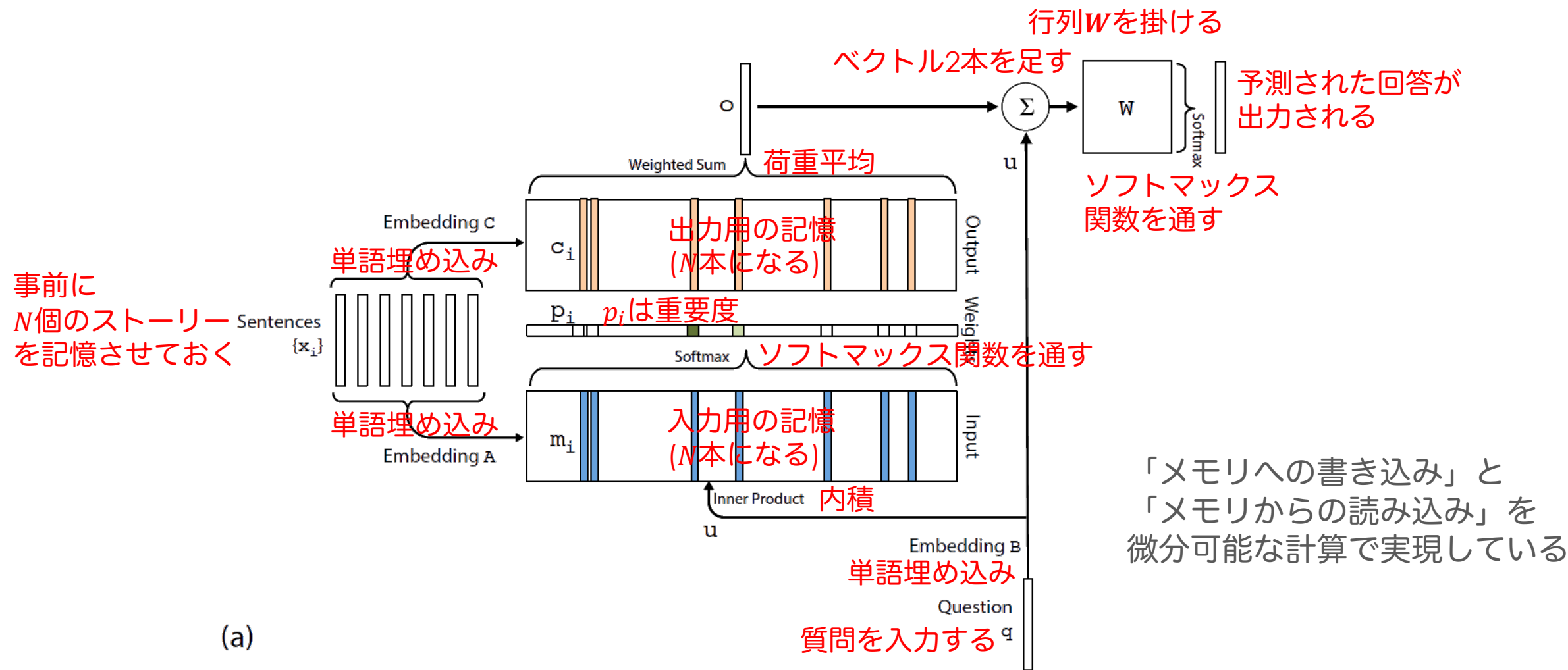
MemN2N概観

- 入力された文はWord embeddingで特徴に変換
 - Memory Networksの特徴変換 $I(x)$ に対応
- MemN2Nでは入力用と出力用で異なる表現のメモリを持つ
 - 入力用 m_i , 出力用 c_i
 - それぞれ異なるembeddingを行う
- 出力文を生成する O モジュールはソフト・アテンションを行う
- 生成された o と質問文の表現 $B(x)$ から回答を生成する



(原論文 Fig. 1から引用・一部改変)

MemN2Nの構成



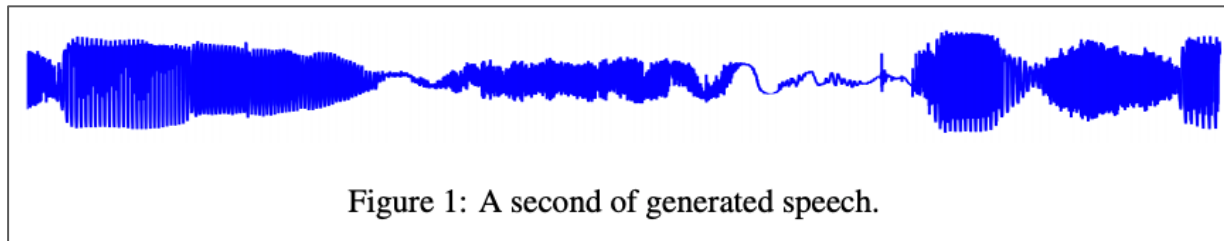


Any Questions ?

その他の話題

WaveNet

- Google DeepMind社が考案したオーディオ生成モデル。
- Google Homeなどで利用されている。
- 原著論文
 - Aaron van den Oord , WAVENET: A GENERATIVE MODEL FOR RAW AUDIO, <https://arxiv.org/pdf/1609.03499.pdf>
- 以下のサイトにアクセスし、WaveNetによって生成されたオーディオを聞いてみよう。
 - <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>

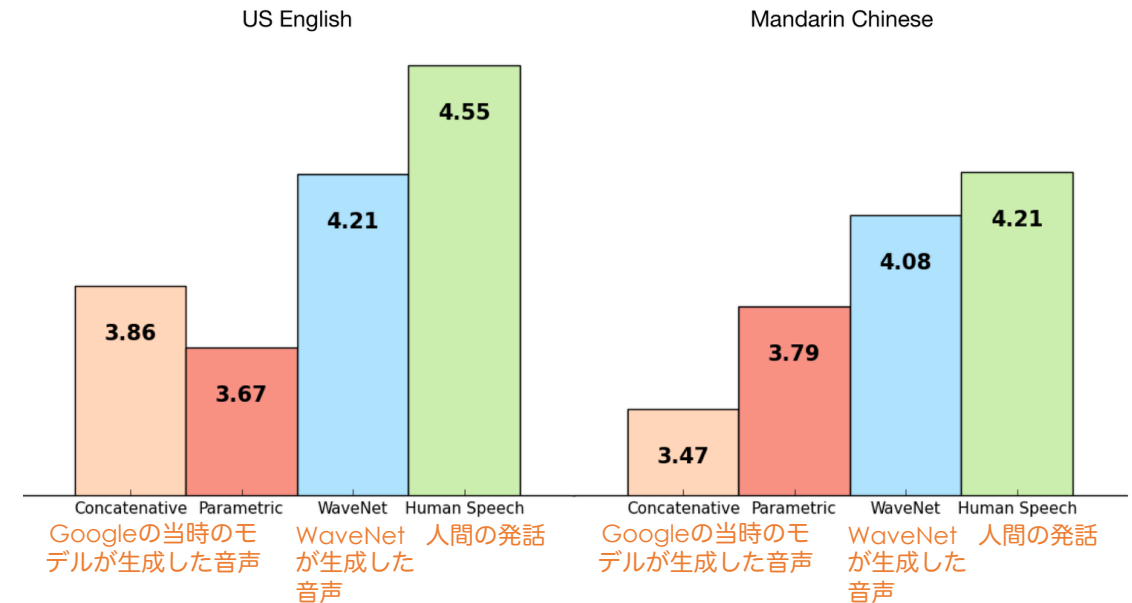


WaveNet

- WaveNetの特徴

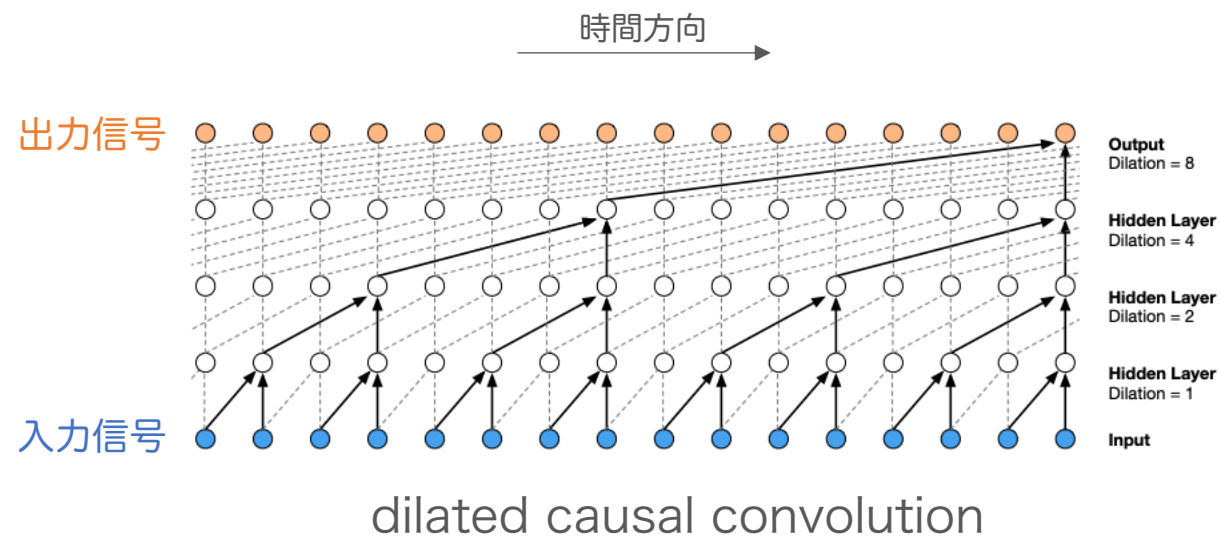
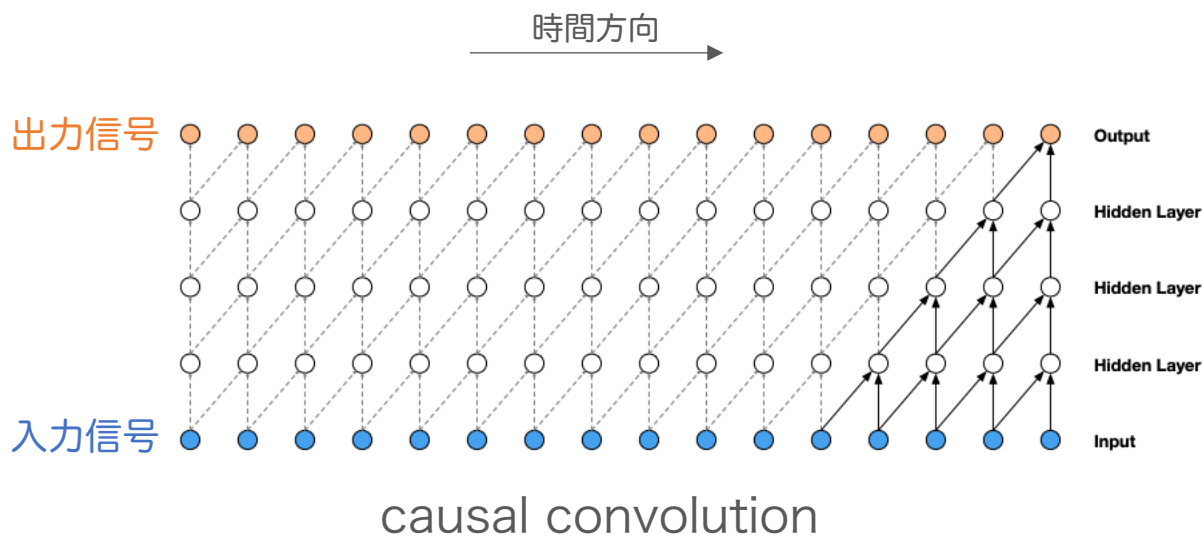
- 入力されたテキストから人間らしい自然な音声を生成できる(text-to-speech, TTS)。
- 発話者の条件を付けて学習させることで、1つのモデルで様々な声を生成できる。
- 音声認識や音楽生成も可能。

人による主観的な評価結果
(スコアは1~5の値をとる)



WaveNet

- WaveNetのアーキテクチャの一番の特徴は、dilated causal convolution である。
 - causal convolutionとは、過去の入力データだけを使って畳み込む方法。
 - dilated causal convolutionは、causal convolutionの受容野を広げた畳み込み法。
 - 少ないパラメータで長い時間の特徴を捉えることができる。
 - dilated causal convolution は、RNNのような再帰結合ではないため、RNNに比べ時間方向の計算を効率的に行える。

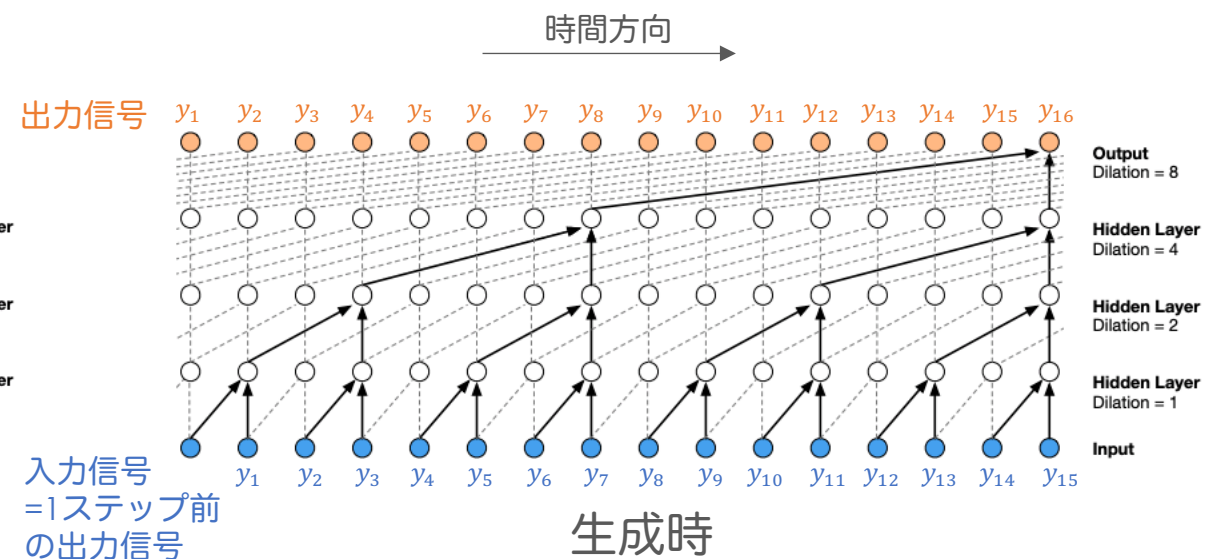
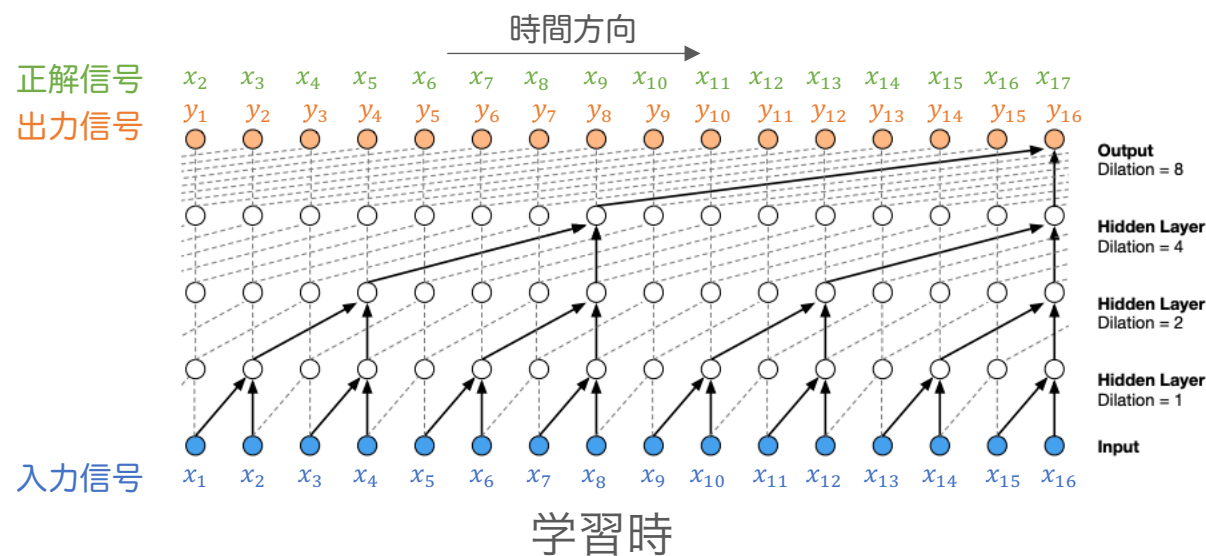


WaveNet

- WaveNetの学習と生成

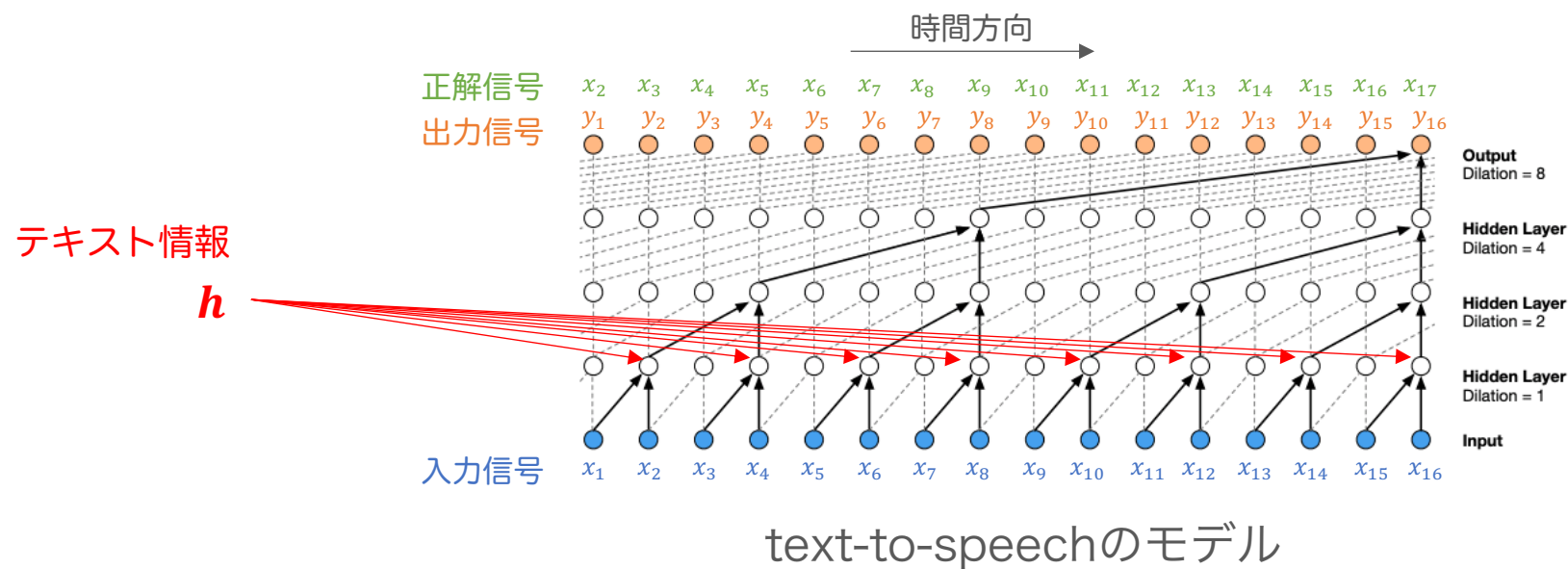
- 学習時は、1ステップ先の信号を予測できるようにフィルタなどのパラメータを学習する
- 生成時は、1ステップ目から順番に生成していく。
 - tステップ目の信号を生成したら、その生成値をt+1ステップ目の入力に持っていき、t+1ステップ目の信号を生成する。これを終わりまで繰り返す。
- 生成する際の仕組みは、こちらのアニメーションがわかりやすい。

- https://lh3.googleusercontent.com/Zy5xK_i2F8sNH5tFtRa0SibLp_CU7QwzS2iB5nf2ijlf_OYm-Q5D0SgoW9SmfbDF97tNEF7Cmxal-o6oLC8sGlrJ5HxWNk79dL1r7Rc=w1440-rw-v1



WaveNet

- テキストから音声を生成する(text-to-speech)モデルを作る場合は、**テキスト情報 h** を畳み込みノードに与えるようにする



- パープレキシティ(perplexity)は、言語モデルを評価するための指標。
- 1以上の値をとり、値が小さいほどモデルの性能が良いことになる。
- 計算式

$$L = -\frac{1}{N} \sum_n \sum_k^K t_{nk} \log y_{nk}$$
$$p = e^L$$

L : クロスエントロピー誤差

N : データ数

K : 出力層ノード数

t_{nk} : データ n のノード k の正解値

y_{nk} : データ n のノード k の出力値

p : パープレキシティ

パープレキシティ

- パープレキシティ(perplexity)は、確率の逆数という意味を持っており、正解となる単語を選ぶ難しさ、と解釈できる。
- 正解の単語に対応する出力値(確率)を y とすると、 $e^{-\log y} = (e^{\log y})^{-1} = \frac{1}{y}$ であるため、パープレキシティは、確率の逆数という意味になる。
- $y = 0.1$ のとき、その確率の逆数は10である。これは、10個の単語候補の中から正解となる単語を見つけるくらい難しい、と解釈できる。

私 は 今日 ? へ 行った

正解 : 会社

モデルAによる予測 : 会社である確率 0.1

モデルBによる予測 : 会社である確率 0.5

モデルAにとっては、10個の単語候補の中から正解となる単語を見つけるくらい難しい

モデルBにとっては、2個の単語候補の中から正解となる単語を見つけるくらい難しい

- BLEU(bilingual evaluation understudy)スコアは、機械翻訳されたテキストを評価するための指標。
 - 機械翻訳文が人手で作成された正解翻訳文にどれだけ近いかを測る
 - BLEUの意味を直訳すると、「バイリンガルによる評価の代役」
 - 機械翻訳文の評価をバイリンガルにやってもらうと、時間とコストがかかるので、その置き換えを目指すという意味が込められている
- 提案論文
 - Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a Method for Automatic Evaluation of Machine Translation. <https://www.aclweb.org/anthology/P02-1040.pdf>
- 参考文献
 - Thushan Ganegedara. Natural Language Processing with TensorFlow
 - <https://www.packtpub.com/application-development/natural-language-processing-tensorflow>
- 参考サイト
 - <https://cloud.google.com/translate/automl/docs/evaluate?hl=ja#bleu>

BLEUスコア

- BLEUスコアでは、**適合率**が用いられる。
 - 機械翻訳文中の単語が人手翻訳文中に入っていたら、**適合率**が高くなる。
- BLEUスコアは必ず**0-1の値**になる。
 - BLEUスコアが0の場合、機械翻訳文に人手翻訳文と一致する部分がない（品質が低い）ことを意味する。
 - BLEUスコアが1の場合、機械翻訳文が人手翻訳文と完全に一致している（品質が高い）ことを意味する。

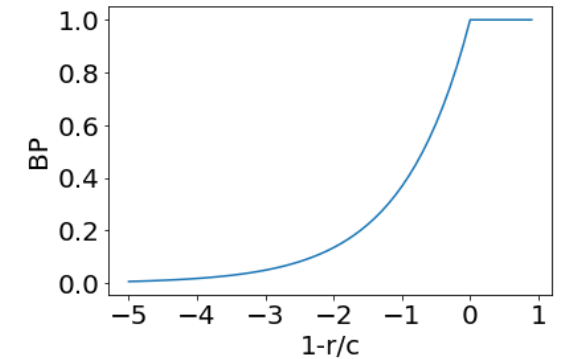
- BLEUスコアの算出方法を以下に示す。

$$\text{BLEU} = \text{BP} \times \exp \left(\sum_{n=1}^N w_n \log p_n \right)$$

$N = 4, w_n = \frac{1}{4}$
で計算されることが多い

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}$$

機械翻訳文が短いほど、BLEUが良くなる傾向があるため、短いことに対してペナルティを課す。
 p_n は適合率であるが、このペナルティによって、再現率を考慮していると解釈することもできる。



$$p_n = \frac{\sum_{C \in \{\text{Candidates}\}} \sum_{n\text{-gram} \in C} \min(\text{Count}(n\text{-gram}), \text{Max_Ref_Count}(n\text{-gram}))}{\sum_{C' \in \{\text{Candidates}\}} \sum_{n\text{-gram}' \in C'} \text{Count}(n\text{-gram}')$$

機械翻訳文が複数あり、そこから一つずつ取り出す その機械翻訳文に含まれるn-gramを1つずつ取り出す 機械翻訳文に含まれる該当n-gramの個数 人手翻訳文毎の該当n-gramの個数の最大値

機械翻訳文が複数あり、そこから一つずつ取り出す その機械翻訳文に含まれるn-gramを1つずつ取り出す 機械翻訳文に含まれる該当n-gramの個数

BP : 機械翻訳文が短いことへのペナルティ (brevity penalty)

N : N -gramまで考慮する

p_n : 修正された n -gram 適合率

w_n : 重み

c : 機械翻訳文の単語数

r : 人手翻訳文の単語数

- 修正された適合率の計算例1

標準の適合率だと、機械翻訳文の良さを正しく評価できないことがわかる

Candidate:
The the the the the the the.

Reference 1:
The cat is on the mat.

Reference 2:
There is a cat on the mat.

標準の *1-gram* 適合率 = $\frac{7}{7}$

修正された *1-gram* 適合率 = $\frac{\min(7,2)}{7} = \frac{2}{7}$

Count("the") = 7

Max_Ref_Count("the") = $\max(2,1) = 2$

- 修正された適合率の計算例2

この場合は、標準の適合率と修正された適合率が同じになる

Candidate:
The cat sat on the mat.

Reference 1:
The cat is on the mat.

Reference 2:
There is a cat on the mat.

$$\text{標準の } 1\text{-gram} \text{ 適合率} = \frac{(2+1+0+1+1)}{6} = \frac{5}{6}$$

$$\text{修正された } 1\text{-gram} \text{ 適合率} = \frac{\min(2,2)+\min(1,1)+\min(1,0)+\min(1,1)+\min(1,1)}{6} = \frac{5}{6}$$

Count("the") = 2

Count("cat") = 1

Count("sat") = 1

Count("on") = 1

Count("mat") = 1

Max_Ref_Count("the") = $\max(2,1) = 2$

Max_Ref_Count("cat") = $\max(1,1) = 1$

Max_Ref_Count("sat") = $\max(0,0) = 0$

Max_Ref_Count("on") = $\max(1,1) = 1$

Max_Ref_Count("mat") = $\max(1,1) = 1$

BLEUスコア

<内容語>

名詞・形容詞・動詞・副詞のように実質的な内容を表す言葉

<機能語>

代名詞・前置詞・接続詞・助動詞・限定詞 (a, an, the, your, their) などのように、文法的な関係や話し手の事態のとらえ方を表す言葉

- BLEUスコアの特徴

- コーパスベースの指標

- BLEUスコアは、個々の文の評価に使用してもうまく機能しないことがわかっており、通常はコーパス全体で評価する。
 - 異なるコーパス間や言語間でBLEUスコアを比較することは推奨されていない。

- 内容語と機能語が区別されない

- 「a」や「your」のような機能語が抜けることは、文意を理解する上でそれほど重要でないが、「Tokyo」と「Kyoko」を間違えることと、等価に扱われてしまう。

- 文の意味や文法の正しさまでは十分に評価できない

- 「not」のような単語が1つ抜けると、文の意味が正反対になるが、BLEUスコアではそれを考慮できない。

Any Questions ?

[グループワーク] seq2seqモデルの計算グラフ

- Notebook演習で実装した「エンコーダ側を双方向LSTMにしたアテンション付き seq2seqモデル」について、計算グラフを描きましょう。
- 「系列変換モデルにおけるアテンションの仕組み」の頁を参考にしましょう。
- 描くのは、順伝播のみで構いません。
- 2~3名のグループに分かれて、上記課題に取り組みましょう。(20分)
- 最後に、グループごとに発表していただきます。(15分)

対面講義でこのグループワークを円滑に行うため、
予習の段階で必ずNotebook演習を行なって来てください。

講座の時間が余ったら

- 今回の復習をします。
- 次回の予習をします。