

機械学習における自然言語処理技術

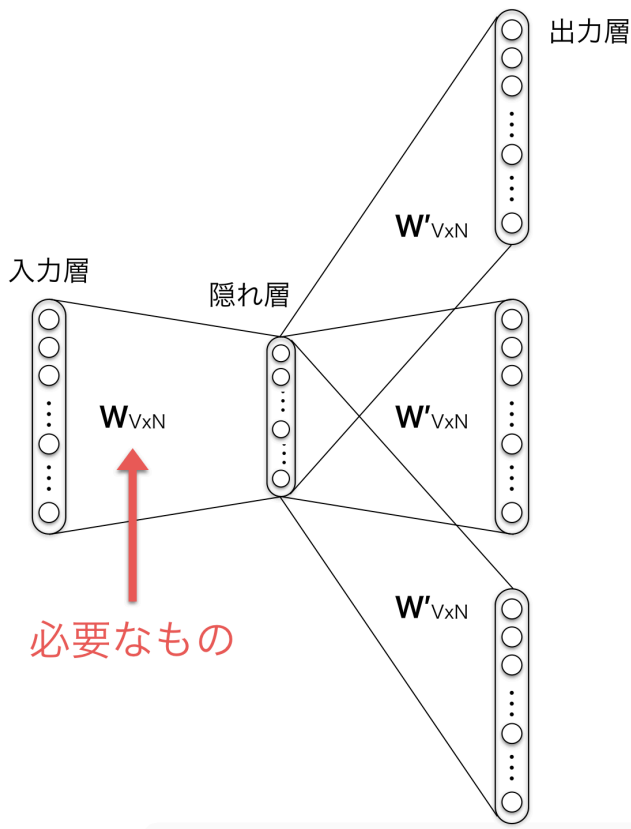
Embedding

one-hot表現

- 最もシンプルなembedding手法.
- 表現したい語彙をリストに表現して、各単語を表現する次元を準備する。表現したい文章に含まれているか否かのベクトルで表現する.

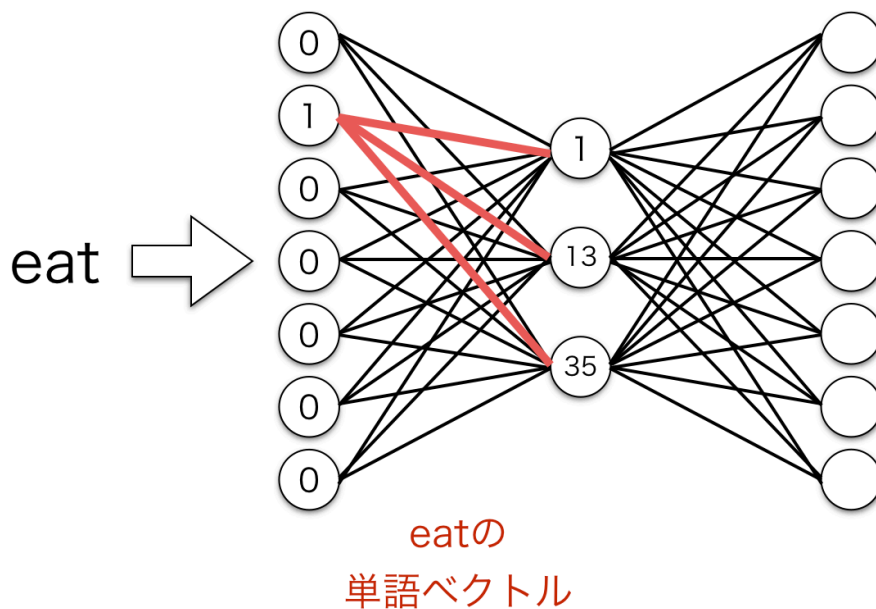
Word2Vec

- 論文: [Efficient Estimation of Word Representations in Vector Space](#)
- 大量のテキストデータを解析して、各単語の意味をベクトルで表現する方法.
- その中でもSkip-Gramモデル(ある単語の周辺に出現する単語の出現確率を計算する)が主に使われる
- Skip-Gram は2層のニューラルネットワークであり隠れ層は一つだけ、隣接する層のユニットは全結合している.

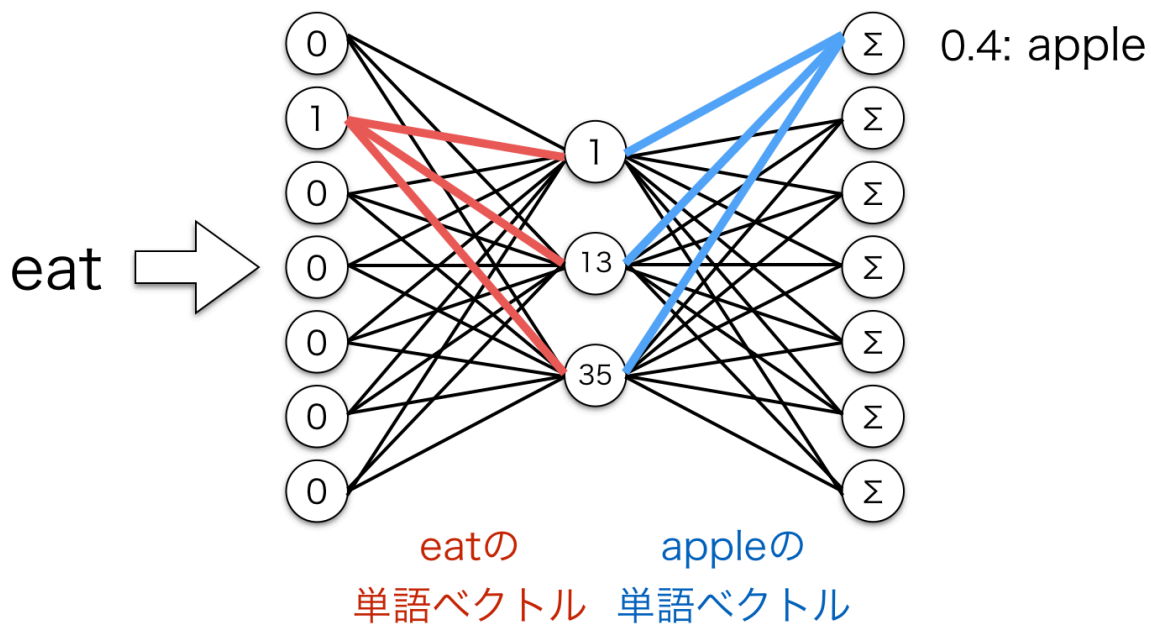


- 目的関数を設定して、2層のニューラルネットワークを構築するが、**word2vecにおいて必要なものは、モデル自体ではなく隠れ層の重み**であることに注意。
- 入力としてある単語、出力にその周辺単語を与えてニューラルネットワークを学習させることで、「意味が近い(=意味ベクトルの距離が近い)時は周辺単語の意味ベクトルもまた距離が近いはず」という仮説に基づいたembedding表現を得ることができる。

$$[0 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0] \times \begin{bmatrix} 74 & 29 & 98 \\ 1 & 13 & 35 \\ 65 & 31 & 37 \\ 96 & 88 & 84 \\ 45 & 94 & 96 \\ 21 & 88 & 9 \\ 6 & 78 & 94 \end{bmatrix} = [1 \quad 13 \quad 35]$$



- 上図のように、one-hotベクトルを入力として与えてあげれば、実際に対象の単語ベクトルを抽出する際は内積ではなくインデックスを使って抽出すれば良いだけなので単語数や隠れ層の次元を気にすることなくモデルを構築することができる



- 出力層は上図のよう。対象の単語の重みベクトルを得たあと、中間層から出力層の単語の重みベクトルとの内積をとっている。つまり単語同士の内積が出力となっていることがわかる。

LSTM

- 論文: [Sequence to Sequence Learning with Neural Networks](#)
- LSTM(Long short-term memory)は、RNN(Recurrent Neural Network)の拡張として1995年に登場した、時系列データ(sequential data)に対するモデル、あるいは構造(architecture)の1種。その名は、Long term memory(長期記憶)とShort term memory(短期記憶)という神経科学における用語から取られている。LSTMはRNNの中間層のユニットをLSTM blockと呼ばれるメモリと3つのゲートを持つブロックに置き換えることで実現されている。

Hochreiterの勾配消失問題

- 当時のRNNの学習方法は、BPTT(Back-Propagation Through Time)法とRTRL(Real-Time Recurrent Learning)法の2つが主流で、その2つとも完全な勾配(Complete Gradient)を用いたアルゴリズムだった
- しかし、このような勾配を逆方向(時間をさかのぼる方向)に伝播させるアルゴリズムは、多くの状況において「爆発」または「消滅」することがあり、結果として長期依存の系列の学習が全く正しく行われなかったという欠点が指摘されてきた
- Hochreiterは自身の修士論文(91年)において、時間をまたいだユニット間の重みの絶対値が指定の(ごくゆるい)条件を満たすとき、その勾配はタイムステップ t に指数関数的に比例して消滅または発散することを示した。
- これはRNNだけではなく、勾配が複数段に渡って伝播する深いニューラルネットにおいてもほぼ共通する問題らしい。
- 例えば、単体のユニット u から v への誤差の伝播について解析する。ステップ t における任意のユニッ

ト u で発生した誤差が q ステップ前のユニット v に伝播する状況を考えたとき、誤差は以下に示すような係数でスケールする.

$$\frac{\partial v_v(t-q)}{\partial v_u(t)} = \begin{cases} f'_v(net_v(t-1))w_{uv} & q = 1 \\ f'_v(net_v(t-q)) \sum_{l=1}^n \frac{\partial v_v(t-q+1)}{\partial v_u(t)} w_{lv} & q > 1 \end{cases}$$

- $l_q = v$ と $l_0 = u$ を使用して、

$$\frac{\partial v_v(t-q)}{\partial v_u(t)} = \sum_{l_1=1}^n \cdots \sum_{l_{q-1}=1}^n \prod_{m=1}^q f'_{l'_m}(net_{l'_m}(t-m))w_{l'_m l_{m-1}}$$

- 上式より、以下の場合にはスケール係数は発散し、その結果としてユニット v に到着する誤差の不安定性により学習が困難になる.

$$|f'_{l'_m}(net_{l'_m}(t-m))w_{l'_m l_{m-1}}| > 1.0 \quad \text{for all } m$$

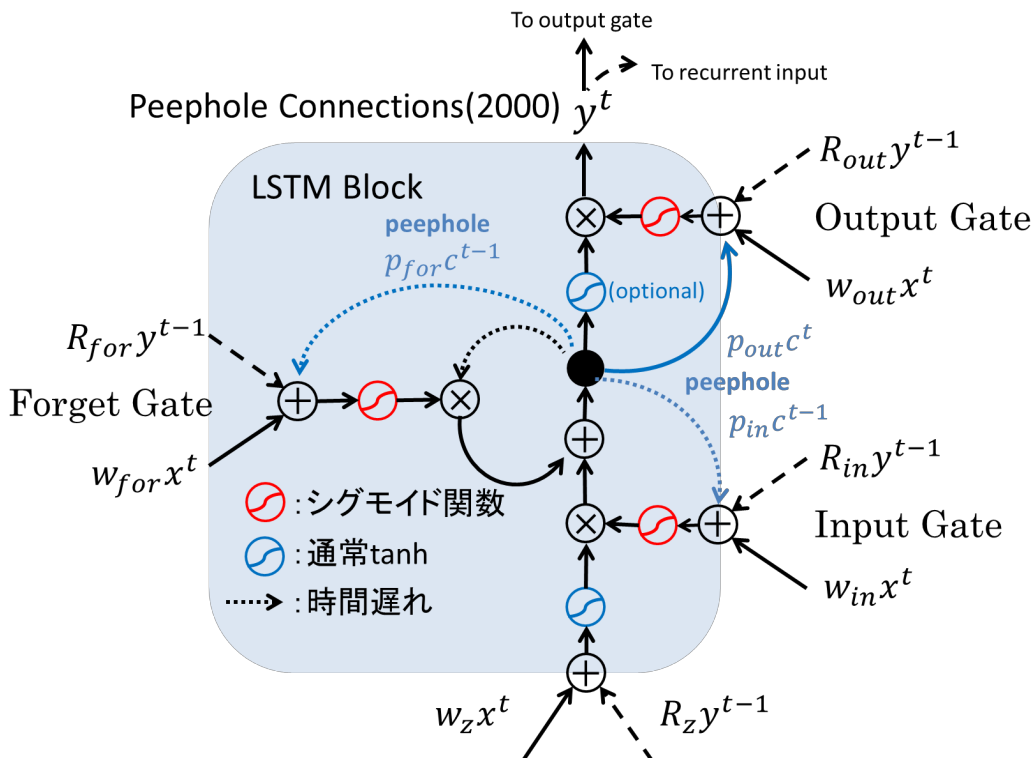
- 一方、以下の場合にはスケール係数は q に関して指数関数的に減少する.

$$|f'_{l'_m}(net_{l'_m}(t-m))w_{l'_m l_{m-1}}| < 1.0 \quad \text{for all } m$$

- これらの問題を解決するために考案されたのがLSTM

LSTMモデル

- R と W は重み行列



LSTMの順伝播計算

$$\bar{z}^t = W_z x^t + R_z y^{(t-1)} + b_z, z^t = g(\bar{z}^t)$$

$$\bar{i}^t = W_{in} x^t + R_{in} y^{(t-1)} + p_{in} \odot c^{t-1} + b_{in}, i^t = \sigma(\bar{i}^t)$$

$$\bar{f}^t = W_{for} x^t + R_{for} y^{(t-1)} + p_{for} \odot c^{t-1} + b_{for}, f^t = \sigma(\bar{f}^t)$$

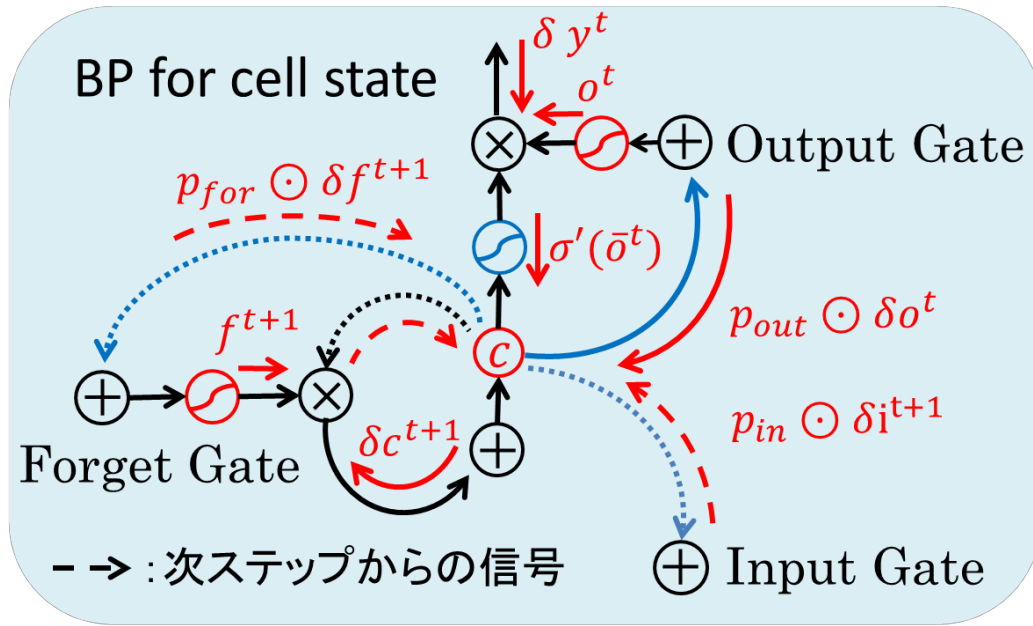
$$c^t = i^t \odot z^t + f^t \odot c^{t-1}$$

$$\bar{o}^t = \sigma(W_{out} x^t + R_{out} y^{(t-1)} + p_{out} \odot c^t + b_{out}), o^t = \sigma(\bar{o}^t)$$

$$y^t = o^t \odot h(c^t)$$

$$s.t. \sigma(x) = \text{sigmoid}(x) = \frac{1}{1+e^{-x}}, g(x) = h(x) = \tanh(x)$$

逆伝播



$$\delta y^t = \Delta^t + R_z^T \delta z^{t+1} + R_{in}^T \delta i^{t+1} + R_{for}^T \delta f^{t+1} + R_{out}^T \delta o^{t+1}$$

$$\delta o^t = \delta y^t \odot h(c^t) \odot \sigma'(\bar{o}^t)$$

$$\delta c^t = \delta y^t \odot o^t \odot h'(c^t) + p_{out} \odot \delta o^t + p_{in} \odot \delta i^{t+1} + p_{for} \odot \delta f^{t+1} + \delta c^{t+1} \odot f^{t+1}$$

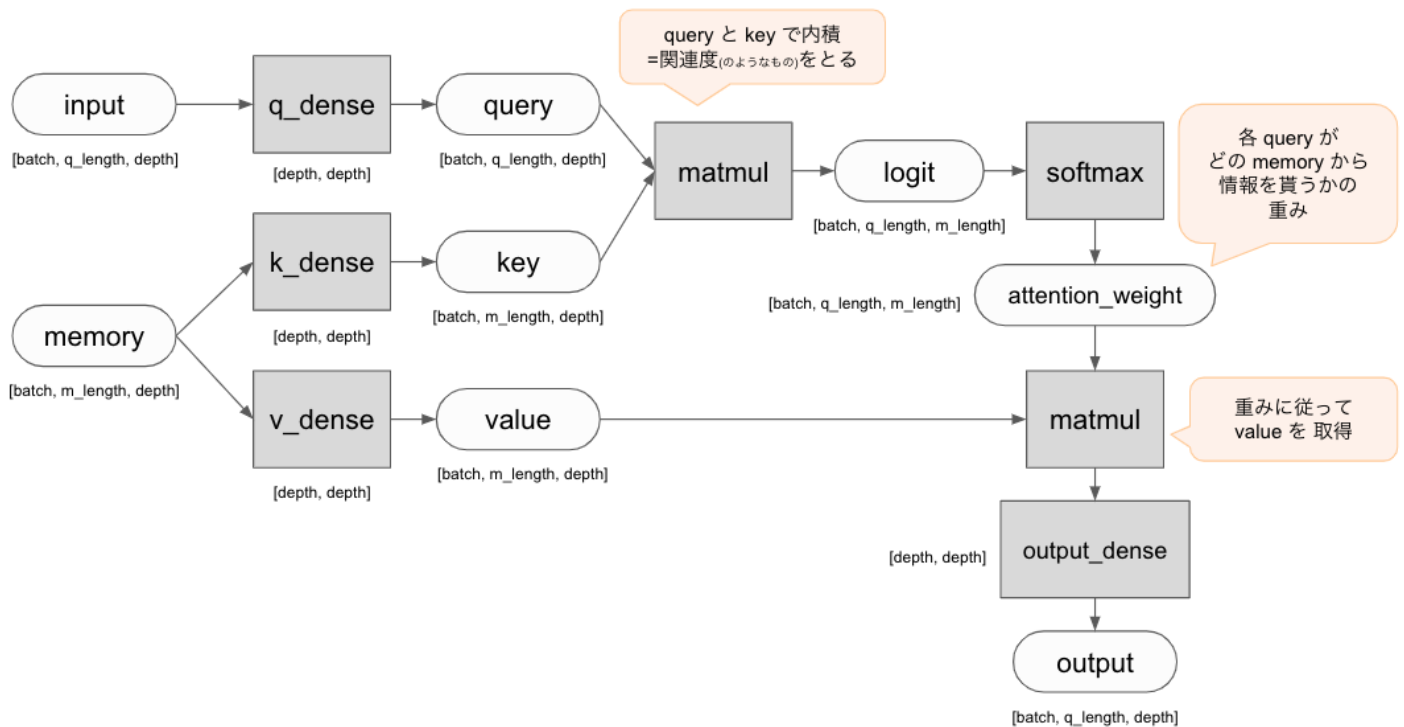
$$\delta f^t = \delta c^t \odot c^{t-1} \odot \sigma'(\bar{f}^t)$$

$$\delta i^t = \delta c^t \odot z^t \odot \sigma'(\bar{i}^t)$$

$$\delta z^t = \delta c^t \odot i^t \odot g'(\bar{z}^t)$$

Attention

- 論文: [Effective Approaches to Attention-based Neural Machine Translation](#)



- Attentionの基本は*query*と*memory*(*key*, *value*).
- Attentionとは*query*によって*memory*から必要な情報を選択的に引っ張ってくること. *memory*から情報を引っ張ってくるときには, *query*は*key*によって取得する*memory*を決定し, 対応する*value*を取得する.

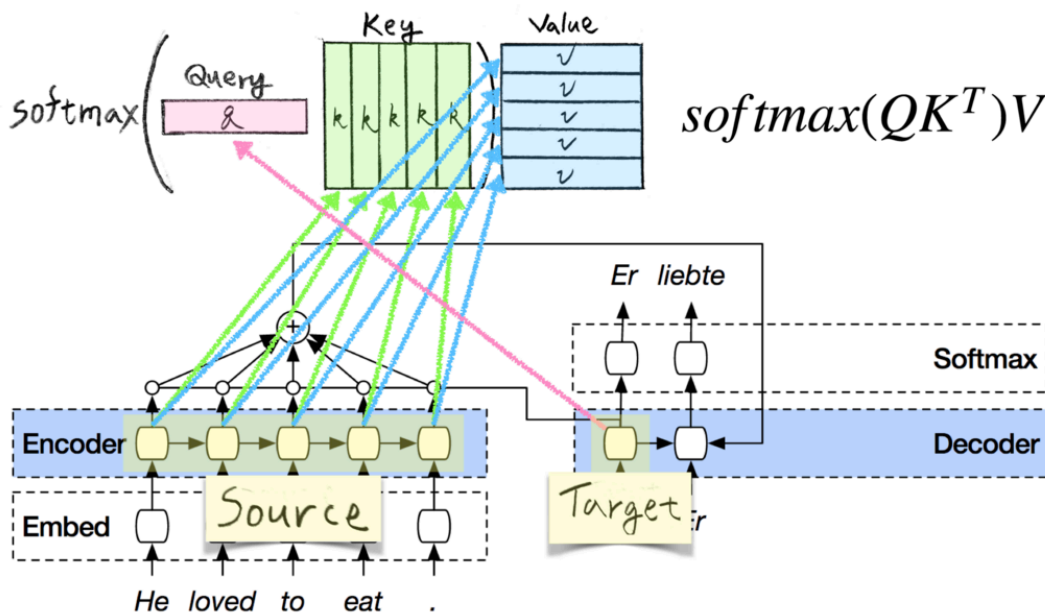
Encoder-Decoderにおけるattention

- 一般的なEncoder-Decoderの注意はエンコーダの隠れ層を*Source*, デコーダの隠れ層を*Target*として次式によって表される.

$$\text{Attention}(\text{Target}, \text{Source}) = \text{Softmax}(\text{Target} \cdot \text{Source}^T) \cdot \text{Source}$$

- より一般化すると*Target*を*query*(検索クエリ)と見做し, *Source*を*Key*と*Value*に分離する.

$$\text{Attention}(\text{query}, \text{Key}, \text{Value}) = \text{Softmax}(\text{query} \cdot \text{Key}^T) \cdot \text{Value}$$

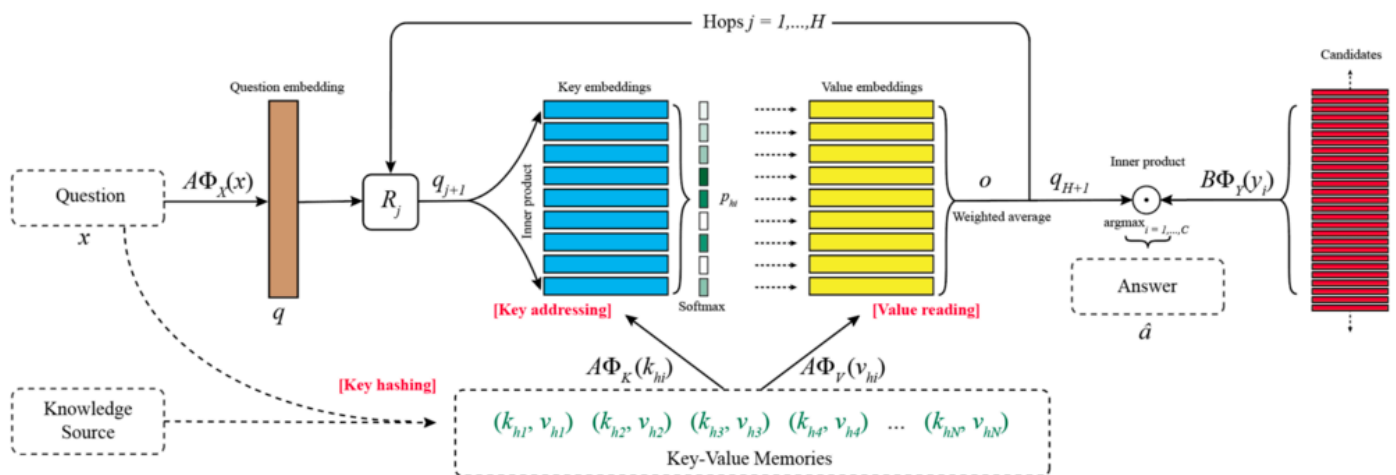


- この時 *Key* と *Value* は各 *key* と各 *value* が一対一対応する key-value ペアの配列, つまり辞書オブジェクトとして機能する.
- *query* と *Key* の内積は *query* と各 *key* の類似度を測り, *softmax* で正規化した注意の重み (Attention Weight) は *query* に一致した *key* の位置を表現する. 注意の重みと *Value* の内積は *key* の位置に対応する *value* を加重和として取り出す操作である.
- つまり注意とは *query* (検索クエリ) に一致する *key* を索引し, 対応する *value* を取り出す操作であり, これは辞書オブジェクトの機能と同じである. 例えば一般的な Encoder-Decoder の注意は, エンコーダのすべての隠れ層 (情報源) *Value* から *query* に関連する隠れ層 (情報) *value* を注意の重みの加重和として取り出すことである.

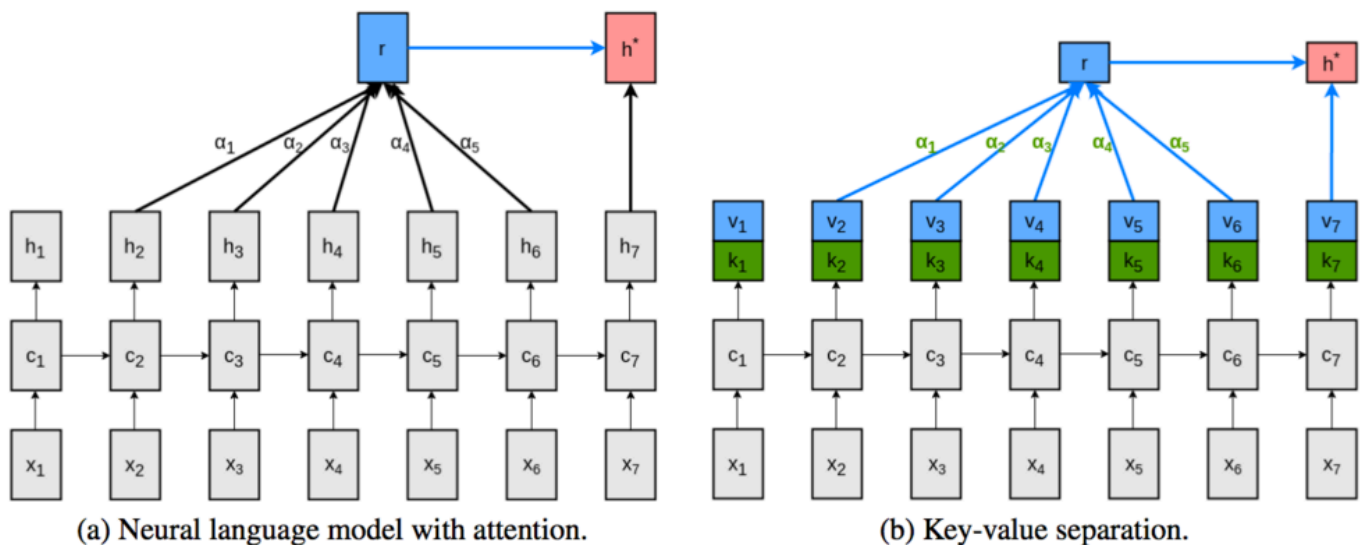
query の配列 *Query* が与えられれば, その数だけ key-value ペアの配列から *value* を取り出す.

MemoryをKeyとValueに分離する意味

- key-value ペアの配列の初出は End-To-End Memory Network [Sukhbaatar, 2015] であるが, *Key* を Input, *Value* を Output (両方を合わせて Memory) と表記しており, 辞書オブジェクトという認識はなかった.
- (初めて辞書オブジェクトと認識されたのは [Key-Value Memory Networks](#) [Miller, 2016] である.)



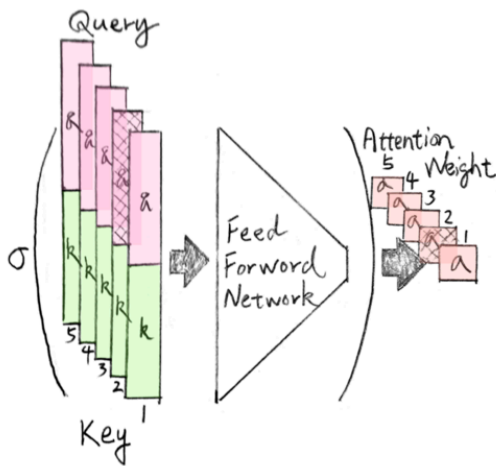
- Key-Value Memory Networks では key-value ペアを文脈 (e.g. 知識ベースや文献) を記憶として格納する一般的な手法だと説明している。 **MemoryをKeyとValueに分離することでkeyとvalue間の非自明な変換によって高い表現力が得られる** という。 ここでいう非自明な変換とは、例えば「keyを入力してvalueを予測する学習器」を容易には作れない程度に複雑な (予測不可能な) 変換という意味である。
- その後、言語モデルでも同じ認識の手法 [Daniluk, 2017] が提案されている。



attentionのweightの算出方法

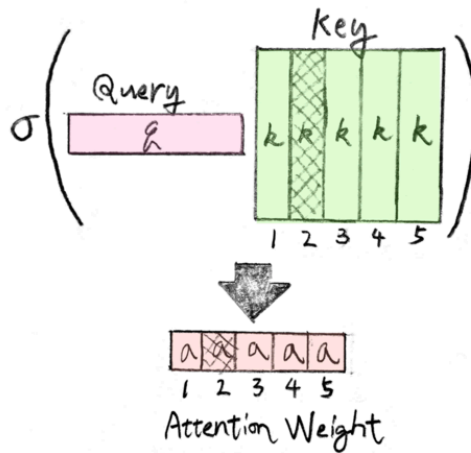
加法注意 (Additive Attention)

$$\text{softmax}(\text{FFN}([Q; K]))$$



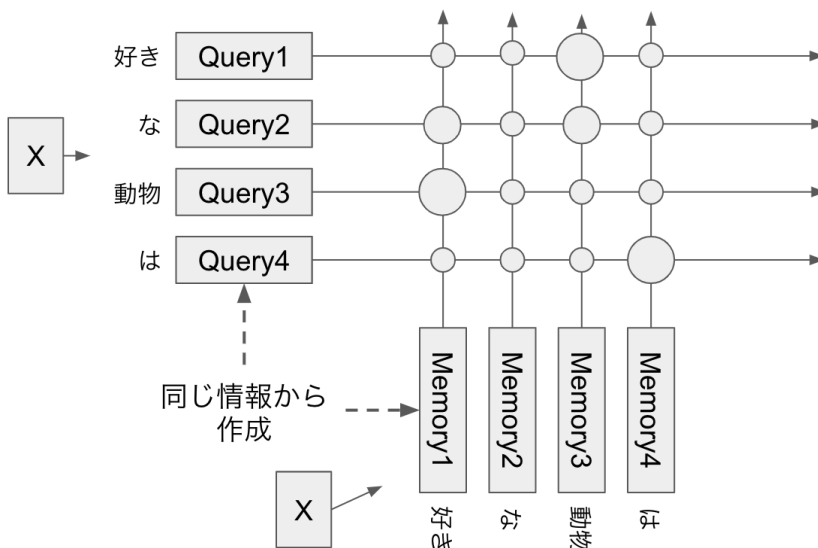
内積注意 (Dot-Product Attention)

$$\text{softmax}(QK^T)$$



- 加法注意と内積注意があり、加法注意は一層のフィードフォワードネットワークで重みを算出する一方、内積注意はattentionの重みをqueryとkeyの内積で算出する。こちらは前者に比べてパラメータが必要ないため、効率よく学習ができる。

self-attention



- $input(query)$ と $memory(key, value)$ すべてが同じTensorを使うAttention
- Self-Attentionは言語の文法構造であったり、照応関係（its が指してるのは Law だよなとか）を獲得するのに使われているなどと論文では分析されている
- 例えば「バナナが好き」という文章ベクトルを自己注意としたら、以下のような構造になる。