

Attentionの種類

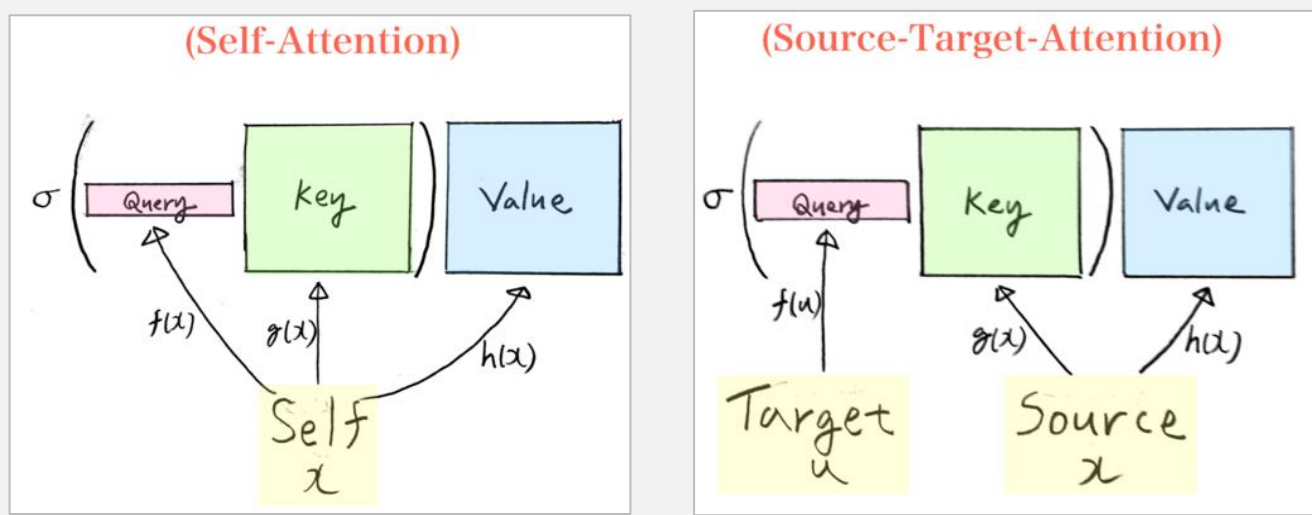
①Self-Attention

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 だよねとか) を獲得するのに使われている
Self-Attention は汎用に使える仕組みで、Transformer の Encoder 部分、Decoder 部分のどちらでも使われている。
文章分類は SourceTarget-Attention は必要なく Self-Attention のみで作ることができる。

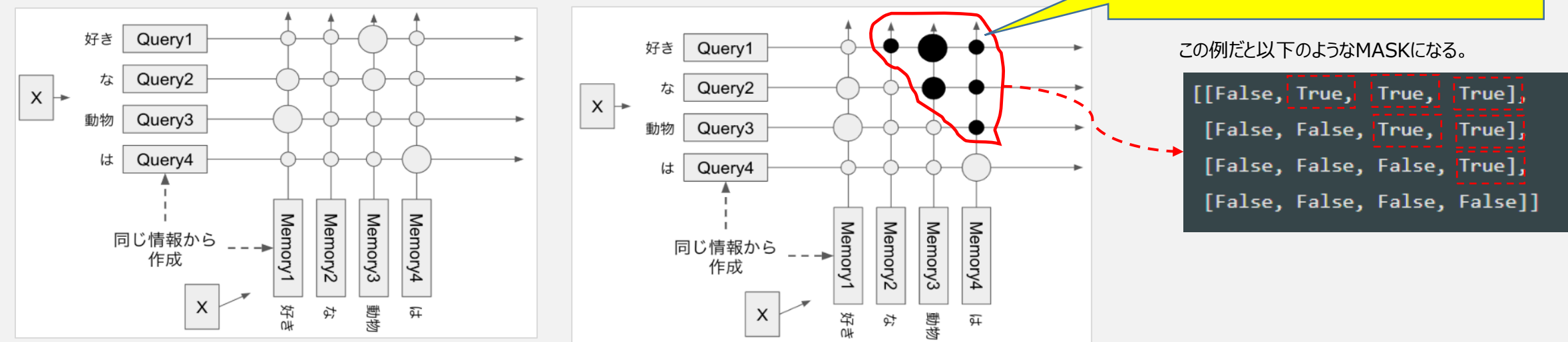
②SourceTarget-Attention

input (query) と memory (key, value) の2つが別の Tensor を使う Attention



Decoder の Self-Attention で未来の情報を参照できないようにする

Decoder は推論時はある時刻tの入力からt+1を予測し、その予測を次の入力に回してさらにそのさを予測する。(自己回帰)
一方学習のときにはすべての時刻で同時に次の時刻のトークンを予測する学習を行う。
→この時 Self-Attention で予測対象の未来の情報が得られてしまうと困る = 未来の情報をマスクする。



上の図で、query が「好き / な」までしかまだ生成されていないのに、memory の未来の情報「動物」を参照できたら困る (推論時には未来の情報は当然与えられないので)

Multi-head Attention

Multi-head Attention は、これまでのような Simple Attention をパラレルに並べたもの。
それぞれの attention を head と呼ぶ。

※論文では一つの大きな Attention を行うよりも、小さな複数の head に分けて Attention を行うほうが性能が上がったとされている。
※query, key, value をそれぞれ head_num 個に split してからそれぞれ attention を計算し、最後に concat するという仕組み

Positional Encoding

Transformer は文章内のトークンの順序を学習に使うことができないので位置情報を付加する必要がある。
各トークンがどの位置にあるのかを表すための値、Positional Encoding を各トークンを Embedding したものに足す。
Positional Encoding は次のような数式で計算される。

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

pos: 時刻
2i, 2i+1: Embedding の何番目の次元か

このような数式を選んだ理由としては、この式が相対的な位置（ PE_{pos+k} ）を PE_{pos} の線形関数で表現可能（=ニューラルネットが学習しやすいはず）だからと
のことです。