

$$Q6 = (0.88 + \overset{V1}{112}) \times (0.12 + \overset{V2}{96}) = 10,850.0256. \quad \text{by } V1 \text{ 和 } V2 \text{ 衍生}$$

QKV的關係和Z的關係，用單頭注意力解釋，並說明Z是怎麼derived

1. input sequences
 thinking machine
 嵌入向量
 2. embed each words
3. 計算單頭注意力
 X_1 2×4 X_2 2×4
 W^K 4×3 W^V 4×3
4. 得自注意力
 K_0 V_0 2×3
5. $\text{Softmax}(\frac{Q \cdot K}{\sqrt{d_k}}) \cdot V$ 得 Z_0 2×3
6. 把多頭Z連起來 $\times W^O$ 得 Z
 Z_0 Z_1 Z_2 Z_3 Z_4
 2×3 2×3 2×3 2×3 2×3
 2×15 4×4
 Z 2×4
 意旨計算詞在所有單頭自注意力中的Value值，成多頭的Z值。

透過將 one-head attention 的結果

彙總產生 multi-self attention,

是為了得到多個 self attention 間不同詞間相關性和上下文意完整理解。透過將詞嵌入向量得 X_1 X_2

並將詞向量 X_1 / X_2 分別乘上 W^K / W^V 得到 Head 0 的 $Q_0 / K_0 / V_0$ self-attention

→ 再 $\text{Softmax}(\frac{Q \cdot K}{\sqrt{d_k}}) \times V_0$ 得 Z_0 。

→ 把所有單頭 self-attention 的 Z 合起來 $\times W^O$ 得 Z ，就可以計算所有詞在不同 self-attention 的值。多頭

$$\text{故 } Z = V_1 \cdot \text{softmax}(0.88) + V_2 \cdot \text{softmax}(1.2)$$

= 詞1 (thinking) 在 thinking 和 machine 中產生的相關性和上下文意都被考慮。

矩陣與向量的計算

1. 一維向量 × 一維向量 先憶直
再憶橫

$$[1, 0, 0] \times [1, 2, 0]$$

將上述轉行列相乘 $\rightarrow [1, 0, 0] \times \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} = 1$

遇到矩陣相同要轉置

2. 不同大小矩陣相乘

算法一

$$\begin{bmatrix} 1 & 2 \end{bmatrix}_{1 \times 2} \times \begin{bmatrix} 1 & 3 & 2 & 1 \\ 2 & 3 & 1 & 2 \end{bmatrix}_{2 \times 4} = \begin{bmatrix} 5 & 9 & 4 & 5 \end{bmatrix}_{1 \times 4}$$

or 算法二較簡單

$$1 \times \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + 2 \times \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 6 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 5 & 9 & 4 & 5 \end{bmatrix}$$

更大的矩陣相乘也要記得前面橫的一列一列處理

算法一

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix}_{3 \times 2} \times \begin{bmatrix} 1 & 3 & 2 & 1 \\ 2 & 3 & 1 & 2 \end{bmatrix}_{2 \times 4} = \begin{bmatrix} 5 & 9 & 4 & 5 \\ 4 & 9 & 5 & 4 \\ 11 & 21 & 10 & 11 \end{bmatrix}_{3 \times 4}$$

算法二

$$1 \times \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + 2 \times \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 6 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 5 & 9 & 4 & 5 \end{bmatrix}$$

算法二

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix}_{3 \times 2} \times \begin{bmatrix} 1 & 3 & 2 & 1 \\ 2 & 3 & 1 & 2 \end{bmatrix}_{2 \times 4} = \begin{bmatrix} 5 & 9 & 4 & 5 \\ 4 & 9 & 5 & 4 \\ 11 & 21 & 10 & 11 \end{bmatrix}_{3 \times 4}$$

一個配一列

$$1 \times \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + 2 \times \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 6 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 5 & 9 & 4 & 5 \end{bmatrix}$$

2.

$$2 \times \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + 1 \times \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 6 & 4 & 2 \end{bmatrix} + \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 9 & 5 & 4 \end{bmatrix}$$

$$3 \times \begin{bmatrix} 1 & 3 & 2 & 1 \end{bmatrix} + 4 \times \begin{bmatrix} 2 & 3 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 3 & 9 & 6 & 3 \end{bmatrix} + \begin{bmatrix} 8 & 12 & 4 & 8 \end{bmatrix} = \begin{bmatrix} 11 & 21 & 10 & 11 \end{bmatrix}$$

Transformer Single Embedding 嵌入.

前由 input \rightarrow input embedding \rightarrow encoding 編碼

1. Input Embedding 嵌入 = 將單詞轉換為一系列數列 (向量), 並將向量做為處理和理解語言的輸入 (稱為編碼器), Dimensional 維度是自設的.

input Embedding \rightarrow Cat is sleeping on the mat

.54	.76	.23	.89	.1	.67
.31	.93	.68	.45	.82	.17
.40	.56	.72	.64	.20	.95
.72	.89	.18	.31	.52	.39
.62	.87	.04	.78	.34	.52

} 預設5維 / 5 dim

2. Positional Encoding 位置編碼 = 讓 Model 知道每個詞的順序, Transformer 是一整句一起理解, 並會考慮詞 & 詞互相影響 (似上下文)

* PE 有 2 公式 For even position, For odd position, 不過應該題目會給 相關程度.

3. Adding "Word Embedding" & "Positional Embedding"

input embedding		positional embedding
for ex Cat is sleeping on the mat.	+	Cat is sleeping on the mat.
.54 .76 .23 .89 .1 .67		.0 .841 .90 .141 .75 .95
⋮		⋮
* Cat is sleeping on the mat *		
.54 1.601 1.13 1.031 -0.65 -0.28 *		
⋮		

4. Multi-self attention = Matrix 相乘 QKV 考慮上文、詞間互相影響.

計算 / 考慮句子中某個詞與其他部分的相關程度.

4. Multi-self attention

所有的 Word Embedding + Positional Embedding 都相同, 不同的是 QKV 的 linear weights.

疊加 weight 再相加

Cat is sleeping on the mat

54 1.601 1.13 1.031 -0.65 -0.28

1.31 1.917 1.63 1.33 1.625 0.87

40 566 0.721 0.641 0.202 0.953

1.72 1.889 1.17 1.309 1.509 1.329

62 4.85 1.77 1.97 1.932 2.51

5x6 Matrixs

linear Weights

67 14 71 03

85 31 14 95

24 57 76 93

39 95 59 38

64 94 77 51

94 73 17 02

6x4 Matrix

* Query 5x4 Matrixs

1.71 1.38 1.52 2.64

5.27 5.13 4.62 4.72

2.19 2.13 1.60 1.58

5.76 5.12 4.53 4.22

9.32 8.11 5.54 8.05

1 Query 計算

直向相加

$$\begin{aligned}
 & 54 \times (67, 14, 71, 03) \\
 & + 1.601 \times (85, 31, 14, 95) \\
 & + 1.13 \times (24, 57, 76, 93) \\
 & + 1.031 \times (39, 95, 59, 38) \\
 & + -0.65 \times (64, 94, 77, 51) \\
 & + -0.28 \times (94, 73, 17, 02)
 \end{aligned}
 =
 \begin{pmatrix}
 362 & 076 & 383 & 016 \\
 1361 & 050 & 224 & 1521 \\
 0271 & 644 & 860 & 1051 \\
 402 & 980 & 610 & 392 \\
 -0416 & -0611 & 501 & -033 \\
 -263 & -204 & -048 & -0056
 \end{pmatrix}$$

Query

1.71 1.38 1.52 2.64

2 Key 計算

只算第一行

$$\begin{aligned}
 & 54 \times (50, 76, 28, 12) \\
 & + 1.601 \times (35, 33, 63, 25) \\
 & + 1.13 \times (79, 99, 04, 95) \\
 & + 1.031 \times (03, 09, 23, 09) \\
 & + -0.65 \times (85, 04, 009, 43) \\
 & + -0.28 \times (14, 63, 30, 21)
 \end{aligned}
 =
 \begin{pmatrix}
 027 & 410 & 151 & 0065 \\
 056 & 53 & 101 & 40 \\
 -90 & 112 & 005 & 107 \\
 003 & 109 & 24 & 09 \\
 -055 & -003 & -01 & -28 \\
 004 & -18 & -09 & -006
 \end{pmatrix}$$

Key

1.17 1.94 1.35 1.285
1.16 1.29

3. Value 計算

老師給的公式錯誤

$$\begin{aligned}
 & 54 \times (50, 76, 28, 12) \\
 & + 1.601 \times (35, 33, 63, 25) \\
 & + 1.13 \times (79, 99, 04, 95) \\
 & + 1.031 \times (03, 09, 23, 09) \\
 & + -0.65 \times (85, 04, 009, 43) \\
 & + -0.28 \times (14, 63, 30, 21)
 \end{aligned}
 =
 \begin{pmatrix}
 027 & 410 & 151 & 0065 \\
 0560 & 0528 & 1009 & 0400 \\
 0892 & 1119 & 0045 & 1074 \\
 0031 & 0093 & 0237 & 0093 \\
 -553 & -0026 & -000585 & 280 \\
 0039 & 176 & -0006 & 059
 \end{pmatrix}$$

084

反正得出的 Query / Key / Value

$$5. \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

Query 5x4

1.71 1.38 1.52 2.64
5.29 5.13 4.62 4.72
2.19 2.13 1.60 1.58
5.76 5.12 4.53 4.22
9.32 8.11 5.54 8.05

Key 5x4

1.16 1.94 1.35 1.29
4.15 3.97 2.22 3.18
1.29 1.87 0.93 1.21
3.95 4.10 2.43 2.83
5.45 5.65 4.52 4.50

成 4x5
Key 矩陣要轉置 Transpose

1.16 4.15 1.29 3.95 5.45
1.94 3.97 1.87 4.10 5.65
1.35 1.87 0.93 2.43 4.52
1.29 4.10 2.43 2.83 4.50

只求第一行

$$\begin{pmatrix} 1.71 \times (1.16 & 1.94 & 1.35 & 1.29) \\ 1.38 \times (4.15 & 3.97 & 2.22 & 3.18) \\ 1.52 \times (1.29 & 1.87 & 0.93 & 1.21) \\ 2.64 \times (3.95 & 4.10 & 2.43 & 2.83) \end{pmatrix}$$

得 5x5 Matrix

10.1 24.3 9.39 23.5 35.8
28.3 67.5 26.3 66.4 99.8
10.8 26.1 10.2 25.7 38.3
28.1 67.7 26.3 66.6 99.7
44.4 108.7 42.08 106.3 159.8

★ Scaling 縮放, 將上述 QxK 轉置結果的 5x5 $\div \sqrt{d_K}$ → dimension 5

10.1 24.3 9.39 23.5 35.8
28.3 67.5 26.3 66.4 99.8
10.8 26.1 10.2 25.7 38.3
28.1 67.7 26.3 66.6 99.7
44.4 108.7 42.08 106.3 159.8

$\frac{10.1}{\sqrt{5}} = 4.52$ 10.87 4.20 10.51 16.01
12.66 30.19 11.76 29.70 44.63
4.83 11.67 4.56 11.99 17.13
12.57 30.28 11.76 29.78 44.59
19.86 48.61 18.82 47.54 70.57

可以想成一條一個頭

$\sqrt{d_K} < \sqrt{5}$

$$\star \text{softmax} = \frac{e^{4.52}}{e^{4.52} + e^{10.87} + e^{4.20} + e^{10.51} + e^{16.01}} = 0.00001$$

橫全相加

每個重複 softmax 得 5x5 Matrix

$\text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$ 5x5 Matrix

Value 5x4 = Softmax 值 x V = Z

0.00001 0.0054 0.0007 0.00404 0.9904 1.20 1.80 0.32 1.90
0 0.0005 0 0.00003 0.9999 2.95 4.52 3.45 4.99
0.00045 0.004 0.0033 0.0036 0.9922 0.84 1.58 1.14 1.83
0 0.0005 0 0.00037 0.9999 3.17 4.64 3.85 5.17
0 0.0003 0 0.00001 1 3.79 7.19 5.03 7.70

只算一行

$$\begin{pmatrix} 0.00001 \times (1.20 & 1.80 & 0.32 & 1.90) \\ 0.0054 \times (2.95 & 4.52 & 3.45 & 4.99) \\ 0.0007 \times (0.84 & 1.58 & 1.14 & 1.83) \\ 0.00404 \times (3.17 & 4.64 & 3.85 & 5.17) \\ 0.9904 \times (3.79 & 7.19 & 5.03 & 7.70) \end{pmatrix}$$

得 5x4 Matrix

3.78 7.16 5.01 7.67
3.79 7.19 5.03 7.70
3.78 7.17 5.02 7.68
3.79 7.19 5.03 7.70
3.79 7.19 5.03 7.70

一個單頭注意力

得 5x4 Matrix

三注意力

得 3x4 Matrix

⇒ 5x12 Matrix

Transformer = 序列輸入 序列輸出 模型結構 \rightarrow input \rightarrow encoder \rightarrow context \rightarrow decoder

based on RNN, 但透過「self-attention」自注意力機制考慮了句子中所有單字相關性和上下文, 可以同時理解句子的所有部分, 對整理語意更精準分析, 克服 RNN 中長期依賴和計算效率不佳的限制
long-term dependency

★ Transformer 由編碼器 encoder 和解碼器 decoder 組成

★ Encoder 編碼器 \rightarrow 將輸入序列轉為連續形式, 以獲得上下文意義
輸出通常為一系列嵌入向量, 稱「隱藏狀態」或「上下文」。

★ Decoder 解碼器 \rightarrow 將 Encoder 輸出的隱藏狀態反覆計算成 token
可能有三種 Type 的 Transformer

★ 1. Encoder-only Transformers \rightarrow 強在可深入理解和解釋文本, 透過將輸入序列轉為各種數位形式, ex. BERT / R.BERTa / DistilBERT
使用雙向注意力 \rightarrow 考慮前 / 後文的內容。

★ 2. Decoder-only Transformer 此 Model 強在強大的產出、表述能力, 透過反覆計算預測下一個 P_{max} 的單詞來完成輸出序列, GPT。
而此 Model 採 Autoregressive attention, 僅考慮已輸出的上文。

★ 3. Encoder-decoder Transformer \rightarrow 同時擁有輸入文本、理解語意、輸出回應的能力, 如 T5 / BART, 可用於翻譯和產生摘要。

★ 標記化文本

★ Tokenizer 分詞器 \rightarrow 使用 Tokenizer 將單詞轉為數位, 因機器只懂數位, 而每個詞都有自己的編號。

Apple \rightarrow 107
is \rightarrow 38
my \rightarrow 51

Tokenizer
 \rightarrow Input
Embedding
嵌入

Embedding
嵌入

107 0000
38 0000
51 0000
 $\xleftrightarrow{512}$

Position Embedding

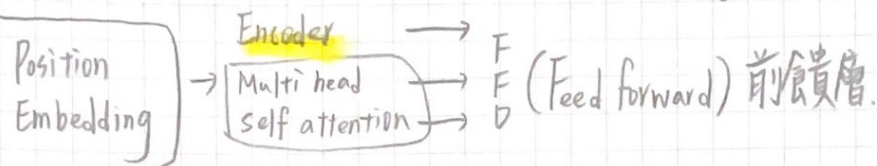
0000
0000
0000

107 0000
38 0000
51 0000
 $\xleftrightarrow{512}$

Position Embedding

0000
0000
0000

- ★ Embedding 嵌入 → 負責將被標記成數字編號的詞轉成向量，一個標記為一個向量，而這些向量被存在一個「向量嵌入空間」，同時也會將單詞的位置紀錄下來 (Position Embedding)，方便文意理解。



- ★ Encoder 編碼器 → 嵌入序列會 input 到編碼器中，這些 input 先經過 Multi head self-attention 及 層注意力層，再到 Feed Forward，並輸入向量 (logits vector)。
- ★ Multi head self-attention → 句子中各個詞會互相影響，要用不同角度去分析語意，所以是「Milti head 多頭」。
- ★ Feed-Forward 前饋 → 由兩層全連接神經網路組成，單獨處理嵌入的序列，並由 Transformer 最終輸出層生成 logit，通常 FFD 使用 Gelu 做為激活函數 (Gelu 讓數據可在線性和非線性中轉換)。
- ★ Decoder 解碼器 → 和編碼器相同，通常有多個 Decoder 組成，Decoder 同樣將嵌入序列的 input 放到 Multi head self-attention 再將 input 傳到 FFD，Decoder 和 Encoder 主要差別在於，Decoder 解碼器有 2 個 注意力子層。Masked multi-head self-attention 蒙面多頭注意力機制：避免偷看答案，確保生成的 token 僅基於「過去輸出 (已產生的內容) 和當前預測」的 token。多頭同樣代表詞間互相影響，並讓 Model 考慮上下文意。
- ★ Encoder-decoder attention → 讓 Model 專注在原文和譯文之間的轉換，此 decoder 會在輸出時考慮上文以及目前為止輸入的內容，並輸出此時最相關的 Output (Token) (P Max)。
- ★ Attention 機制：幫助機器提升翻譯質量。
- ★ Transformer = 讓模型訓練並行，提升訓練效率。