

Numpy - Attention score

Hoàng-Nguyên Vũ

1 Giới thiệu về Transformer?

Trong lĩnh vực trí tuệ nhân tạo, **Transformer** là một kiến trúc mạnh mẽ và hiệu quả trong việc xử lý ngôn ngữ tự nhiên (NLP). Một trong những yếu tố quan trọng tạo nên sức mạnh của Transformer là cơ chế Attention. Bài viết này sẽ giúp các bạn mới học về NumPy hiểu rõ cơ chế Attention trong Transformer thông qua các phép toán cơ bản.

2 Attention là gì?

Attention cho phép mô hình tập trung vào các phần khác nhau của chuỗi đầu vào khi tạo ra mỗi phần tử của chuỗi đầu ra. Điều này giúp mô hình đánh giá tầm quan trọng của mỗi từ trong câu và cải thiện khả năng hiểu ngữ cảnh.

3 Công thức Attention

Công thức cho cơ chế Attention được biểu diễn như sau:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (1)$$

Trong đó:

- Q là ma trận query
- K là ma trận key
- V là ma trận value
- d_k là kích thước chiều của vector key
- softmax là hàm softmax, được tính toán như sau:

$$\text{Hàm kích hoạt Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)} \quad (2)$$

4 Chi tiết từng bước

4.1 Tính toán ma trận điểm số

Ma trận điểm số được tính bằng cách nhân ma trận query Q với ma trận key K chuyển vị, sau đó chia cho căn bậc hai của kích thước chiều của vector key d_k :

$$\text{scores} = \frac{QK^T}{\sqrt{d_k}} \quad (3)$$

4.2 Áp dụng hàm Softmax

Sau khi tính toán ma trận điểm số, chúng ta áp dụng hàm softmax để chuyển đổi các điểm số thành các trọng số xác suất:

$$\text{attention_weights} = \text{softmax}(\text{scores}) \quad (4)$$

4.3 Tính toán tổng có trọng số

Cuối cùng, chúng ta nhân các trọng số attention với ma trận value V để có được đầu ra của cơ chế Attention:

$$\text{output} = \text{attention_weights} \cdot V \quad (5)$$

5 Kết luận

Công thức Attention trong Transformer giúp mô hình tập trung vào các phần quan trọng của chuỗi đầu vào khi tạo ra đầu ra, giúp cải thiện khả năng hiểu ngữ cảnh của mô hình. Việc hiểu và áp dụng công thức này là cơ sở để làm việc với các mô hình ngôn ngữ hiện đại.

6 Bài tập

Hãy viết hàm tính attention score của câu đầu vào: "Tôi thích học AI" với tập Vocab được cho sẵn như sau:

```
1 import numpy as np
2
3 # Giúp kết quả giống nhau trên mọi máy khi tạo random
4 np.random.seed(42)
5
6 # Bước 1: Tạo từ điển và mã hóa từ
7 vocab = {
8     "Tôi": 0,
9     "thích": 1,
10    "học": 2,
11    "AI": 3
12 }
13
14 # Số lượng từ vựng
15 vocab_size = len(vocab)
16
17 # Kích thước vector embedding
18 embedding_dim = 4
```

```

19
20 # Khởi tạo ma trận embedding ngẫu nhiên
21 embedding_matrix = np.random.rand(vocab_size, embedding_dim)
22
23 # Chuỗi đầu vào được mã hóa thành các vector embedding
24 input_seq = np.array([embedding_matrix[vocab[word]] for word in ["Tôi", "thích", "học", "AI"]])
25 print("Chuỗi đầu vào (đã mã hóa):\n", input_seq)
26
27 # Bước 2: Khởi tạo các ma trận trọng số cho Q, K, V
28 W_q = np.random.rand(embedding_dim, embedding_dim)
29 W_k = np.random.rand(embedding_dim, embedding_dim)
30 W_v = np.random.rand(embedding_dim, embedding_dim)
31
32 # Tính toán Q, K, V
33 Q = np.dot(input_seq, W_q)
34 K = np.dot(input_seq, W_k)
35 V = np.dot(input_seq, W_v)
36
37 print("Ma trận Query Q:\n", Q)
38 print("Ma trận Key K:\n", K)
39 print("Ma trận Value V:\n", V)
40
41 # Bước 3: Tính toán Attention score
42 scores = # Your code here #
43
44 # Chia cho căn bậc hai của kích thước chiều của vector key
45 d_k = # Your code here #
46 scores = scores / np.sqrt(d_k)
47
48 print("Điểm số:\n", scores)
49
50 # Bước 4: Áp dụng hàm softmax
51 def softmax(x):
52     # Your code here #
53
54 attention_weights = softmax(scores)
55
56 print("Trọng số Attention:\n", attention_weights)
57
58 # Bước 5: Tính toán tổng có trọng số của các value
59 output = # Your code here #
60
61 print("Đầu ra:\n", output)
62

```

Kết quả:

```

1 Đầu ra:
2 [[1.1796727  0.9217873  1.45815851 1.33770808]
3  [1.09435999 0.83851284 1.39301108 1.243537  ]
4  [1.1515616  0.89314997 1.43845963 1.30643711]
5  [1.11065168 0.85384713 1.40560424 1.26208253]]

```