Hung-yi Lee



1

## Sophisticated Input

• Input is a vector

• Input is a set of vectors





this

is a

cat

To learn more: <a href="https://youtu.be/X7PH3NuYW0Q">https://youtu.be/X7PH3NuYW0Q</a> (in Mandarin)

## Vector Set as Input



https://medium.com/analytics-vidhya/socialnetwork-analytics-f082f4e21b16

## Vector Set as Input

 Graph is also a set of vectors (consider each node as a vector)



http://www.twword.com/wiki/%E5%8 8%86%E5%AD%90

## Vector Set as Input

 Graph is also a set of vectors (consider each node as a vector)

 $H = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & \dots \end{bmatrix}$  $C = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & \dots \end{bmatrix}$  $O = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & \dots \end{bmatrix}$ 



## What is the output?

• Each vector has a label.



### **Example Applications**



## What is the output?

• Each vector has a label.



• The whole sequence has a label.





## What is the output?

• Each vector has a label.

focus of this lecture

seq2seq



• The whole sequence has a label.



Model decides the number of labels itself.



## Sequence Labeling

Is it possible to consider the context?

FC Fullyconnected How to consider the whole sequence? a window covers the whole sequence?











Find the relevant vectors in a sequence







# $\begin{array}{l} \underline{\textit{Self-attention}}\\ \text{on attention scores} \end{array} \quad \begin{array}{l} b^1 = \sum \alpha'_{1,i} v^i \end{array}$



18















### Multi-head Self-attention Different types of relevance



### **Multi-head Self-attention** Different types of relevance



### Multi-head Self-attention Different types of relevance





## Positional Encoding

- No position information in self-attention.
- Each position has a unique positional vector e<sup>i</sup>
- hand-crafted
- learned from data



Each column represents a positional vector  $e^i$ 



29

#### Methods Inductive Data-Driven Parameter Efficient https://arxiv.org/abs/ Sinusoidal (Vaswani et al., 2017) X 1 2003.09229 Embedding (Devlin et al., 2018) х Relative (Shaw et al., 2018) х This paper ✓ (a) Sinusoidal (b) Position embedding i = 1i = 1Position Position i = 252i = 252Feature dimension Feature dimension i = 1i = 1Position Position i = 252i = 252Feature dimension Feature dimension 30

#### Table 1. Comparing position representation methods

(c) FLOATER

(d) RNN

## Many applications ...





### Transformer

https://arxiv.org/abs/1706.03762

BERT

https://arxiv.org/abs/1810.04805

Widely used in Natural Langue Processing (NLP)! 31

https://arxiv.org/abs/1910.12977

## Self-attention for Speech



If input sequence is length L

Attention in a range



## Self-attention for Image



Source of image: https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix\_fig15\_282798184



### **DEtection Transformer (DETR)**



#### https://arxiv.org/abs/2005.12872

## Self-attention v.s. CNN



CNN: self-attention that can only attends in a receptive field

> CNN is simplified self-attention.

Self-attention: CNN with learnable receptive field

Self-attention is the complex version of CNN.

## Self-attention v.s. CNN





On the Relationship between Self-Attention and Convolutional Layers https://arxiv.org/abs/1911.03584

## Self-attention v.s. CNN

### Good for more data

Self-attention



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale https://arxiv.org/pdf/2010.11929,pdf



Attention

https://arxiv.org/abs/2006.16236

## To learn more about RNN .....



https://youtu.be/xCGidAeyS4M

(in Mandarin)



https://youtu.be/Jjy6ER0bHv8 (in English)

## Self-attention for Graph



Consider **edge**: only attention to connected nodes





## Self-attention for Graph

• To learn more about GNN ...



https://youtu.be/eybCCtNKwzA (in Mandarin)



https://youtu.be/M9ht8vsVEw8 (in Mandarin)

### To Learn More ...

### Long Range Arena: A Benchmark for Efficient Transformers https://arxiv.org/abs/2011.04006



### Efficient Transformers: A Survey https://arxiv.org/abs/2009.06732

