

Transformers

CS 4277 Deep Learning

Kennesaw State University

Transformers

- ▶ Convolutional networks take advantage of the relatedness of nearby pixels in images to deal with the high dimensionality of image data.
- ▶ Text data is also high-dimensional with much redundancy that can be exploited, e.g., most instances of `dog` have the same meaning.

Problems: text sequences vary in length and cannot easily be resized.

Consider:

- ▶ Please do not cut all my hair off.

shortened to:

Transformers

- ▶ Convolutional networks take advantage of the relatedness of nearby pixels in images to deal with the high dimensionality of image data.
- ▶ Text data is also high-dimensional with much redundancy that can be exploited, e.g., most instances of `dog` have the same meaning.

Problems: text sequences vary in length and cannot easily be resized.

Consider:

- ▶ Please do not cut all my hair off.

shortened to:

- ▶ Cut all my hair off.

Processing Text Data

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and service.

Dot-Product Self-Attention

Recall that NN layer $f(x)$ takes $D \times 1$ input x , applies linear transformation to input then passes result to activation function:

$$f(x) = \text{ReLU}(\beta + \Omega x)$$

A self-attention block $sa(\cdot)$ takes N inputs.