

Paying attention to Attention

Piotr Mazurek

April 15, 2020

Agenda

- 1 Introduction to Attention
- 2 Introduction to Neural Networks
- 3 Basics of the Attention mechanism
- 4 Deep dive into Attention
- 5 Attention applications

Introduction to Attention

Basic intuition

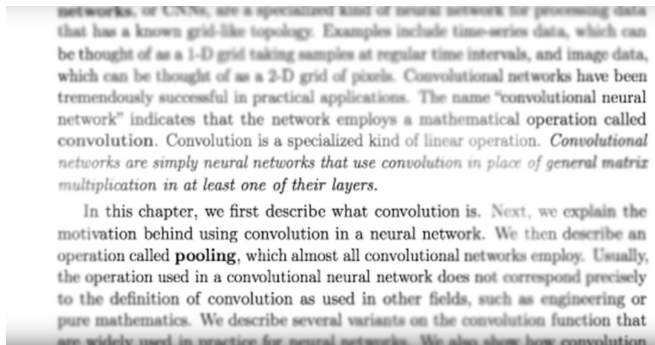


Figure: We teach a neural network to focus on some part of the data¹

¹<https://www.youtube.com/watch?v=W2rWgXJBZhU&t=180s>

Introduction to Attention

Who, where and when?

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau

Jacobs University Bremen, Germany

KyungHyun Cho **Yoshua Bengio***

Université de Montréal

Figure: Idea of the Attention first time mentioned (three times in the two consecutive lines, but still counts), Sept. 2014[1]

Who, where and when?

Attention Is All You Need

Google Brain

avaswani@google.com

Google Brain

noam@google.com

Google Research

nikip@google.com

Google Research

usz@google.com

Google Research

llion@google.com

University of Toronto

aidan@cs.toronto.edu

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

Figure: Attention is all You need paper, Dec. 2017[3]

Introduction to Attention

Boom in Attention papers

The screenshot shows the arXiv search interface. At the top, the Cornell University logo and 'arXiv' text are visible. A search bar contains the word 'attention', and a 'Search' button is to its right. Below the search bar, there are options to 'Show abstracts' and 'Hide abstracts'. A dropdown menu shows '50 results per page' and 'Sort results by Relevance'. The search results list the first entry: '1. arXiv:2004.00910 [pdf, other] [less AB] [cs.LG] [cs.SD] [eess.SP] Improving auditory attention decoding performance of linear and non-linear methods using state-space model'. The authors listed are 'Authors: Ali Aroudi, Tobias de Tellez, Simon Dudoit'. The abstract text follows: 'Abstract: Identifying the target speaker in hearing aid applications is crucial to improve speech understanding. Recent advances in electroencephalography (EEG) have shown that it is possible to identify the target speaker from single-trial EEG recordings using auditory attention decoding (AAD) methods. AAD methods reconstruct the attended speech envelope from EEG recordings, based on a linear least-squares...'. A 'More' link is provided at the end of the abstract. At the bottom of the abstract, it says 'Submitted 2 April 2020; originally announced April 2020.'

Figure: Attention - a hot research topic

Introduction to Attention

Better start paying attention now

Not All Attention Is Needed: Gated Attention Network for Sequence Data

Lanqing Xue,¹ Xiaopeng Li,^{2*} Nevin L. Zhang^{1,3}

¹The Hong Kong University of Science and Technology, Hong Kong

²Amazon Web Services, WA, USA

³HKUST-Xiao Joint Lab, Hong Kong

lxueaa@cse.ust.hk, xiaopel@amazon.com, lzhang@cse.ust.hk

Figure: Not All Attention Is Needed, Dec. 2019

[12]

Neural Networks - basic introduction

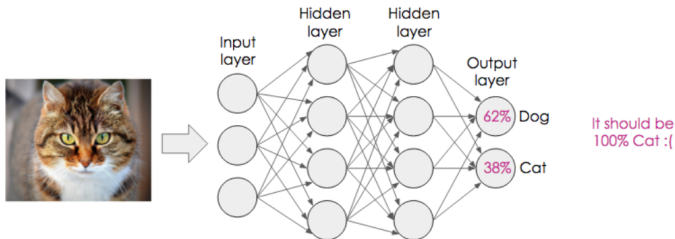


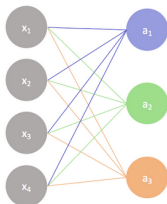
Figure: Not trained neural network²

²<https://www.fromthegenesis.com/artificial-neural-network-part-5/>

Neural Networks

Linear Algebra

Input layer Output layer



A simple neural network

$$\begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b \\ b \\ b \end{bmatrix} = \begin{bmatrix} w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \end{bmatrix} \xrightarrow{\text{activation}} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

Figure: How matrices are used to compute values in neurons ³

³<https://www.jeremyjordan.me/intro-to-neural-networks/>

Neural Networks

Algebraic intuition

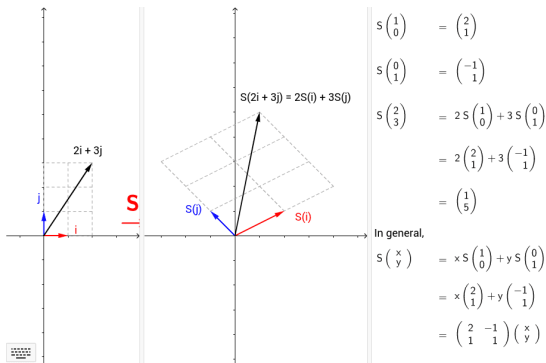


Figure: Linear Transformation ⁴

Mathematical Formalization

Oversimplified

Neural network in general

$$f(x, \theta) = \hat{y} \quad (1)$$

where:

x : input data

Θ : model parameters

\hat{y} : distribution of probability over set of classes

Mathematical Formalization

Oversimplified

Neural network in general

$$f(x, \theta) = \hat{y} \quad (1)$$

where:

x : input data

Θ : model parameters

\hat{y} : distribution of probability over set of classes

Example

$$f(x, w, b) = \sigma(x^T w + b) \quad (2)$$

How to find θ

How to find θ

$$\underset{\theta}{\text{minimize}} \ J(\theta) \quad (3)$$

For example we can use SGD optimizer

Loss Function

How to find θ

How to find θ

$$\underset{\theta}{\text{minimize}} J(\theta) \quad (3)$$

For example we can use SGD optimizer

Example loss function

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \theta)\|^2 \quad (4)$$

Latent space

Auto-encoder

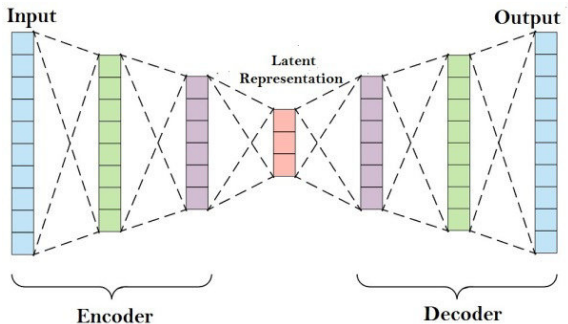


Figure: Example auto-encoder⁵

⁵<https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798>

Latent space

The big picture

- Used to find the hidden representation of data
- By training model on a particular set of data we create some kind of general knowledge

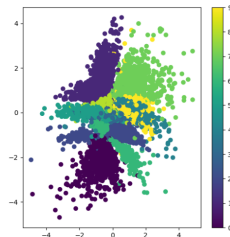


Figure: Finding information in chaos

Recurrent Neural Networks (RNNs)

Oversimplified

- Used for sequential data.
E.g for text analysis, video recognition etc.

Recurrent Neural Networks (RNNs)

Oversimplified

- Used for sequential data.
E.g for text analysis, video recognition etc.
- Have so-called "local memory" or "state memory"

Recurrent Neural Networks (RNNs)

Oversimplified

- Used for sequential data.
E.g for text analysis, video recognition etc.
- Have so-called "local memory" or "state memory"
- Computationally expensive to train

Recurrent Neural Networks (RNNs)

Oversimplified

- Used for sequential data.
E.g for text analysis, video recognition etc.
- Have so-called "local memory" or "state memory"
- Computationally expensive to train
- Problems with vanishing gradient

Recurrent Neural Networks (RNNs)

Oversimplified

- Used for sequential data.
E.g for text analysis, video recognition etc.
- Have so-called "local memory" or "state memory"
- Computationally expensive to train
- Problems with vanishing gradient

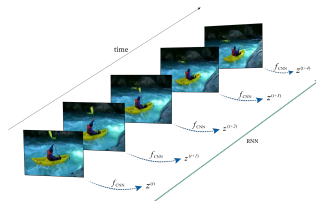
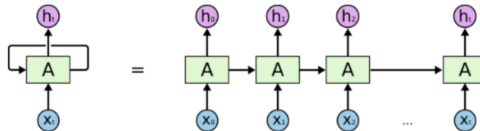


Figure: Sequence of frames^a

^a<https://awesomeopensource.com/project/HHTseng/video-classification>

Recurrent Neural Networks (RNNs)

Oversimplified



An unrolled recurrent neural network.

Figure: Folded and unfolder RNN ⁶

⁶<https://towardsdatascience.com/understanding-rnn-and-lstm-f7cdf6dfc14e>

Types of RNN

Vector to sequence model

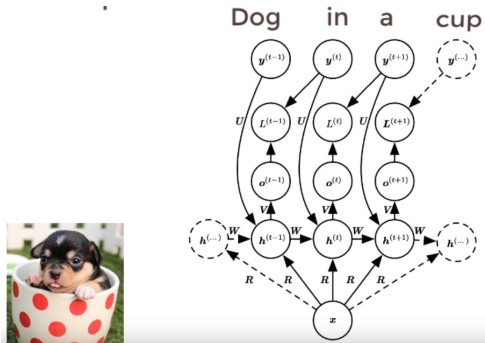


Figure: From an image (vector) generate a caption(sequence) ⁷

⁷<https://www.youtube.com/watch?v=TQQ1Zhbc5ps&t=610s>

Types of RNN

Sequence to vector model

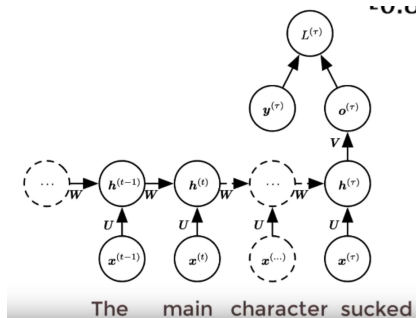


Figure: Take a review (sequence) and predict whether it is positive or negative (vector)⁷

Types of RNN

Sequence to sequence (seq2seq) model

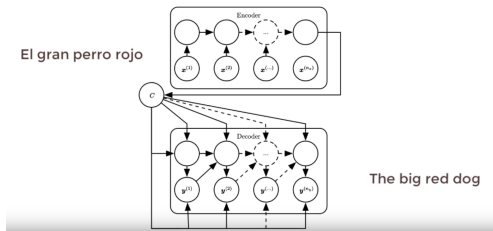


Figure: Translate a sentence in Spanish (sequence) into a sentence in English (sequence)⁷

Attention basic concept - recap

networks, or CNNs, are a specialized kind of neural network for processing data that has a known grid-like topology. Examples include time-series data, which can be thought of as a 1-D grid taking samples at regular time intervals, and image data, which can be thought of as a 2-D grid of pixels. Convolutional networks have been tremendously successful in practical applications. The name "convolutional neural network" indicates that the network employs a mathematical operation called **convolution**. Convolution is a specialized kind of linear operation. *Convolutional networks are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.*

In this chapter, we first describe what convolution is. Next, we explain the motivation behind using convolution in a neural network. We then describe an operation called **pooling**, which almost all convolutional networks employ. Usually, the operation used in a convolutional neural network does not correspond precisely to the definition of convolution as used in other fields, such as engineering or pure mathematics. We describe several variants on the convolution function that are widely used in practice for neural networks. We also show how convolution

Figure: We teach a neural network to focus on the part of the data¹

Attention basic framework

EEAP

1 Embed

Attention basic framework

EEAP

- 1 Embed
- 2 Encode

Attention basic framework

EEAP

- 1 Embed
- 2 Encode
- 3 Attend

Attention basic framework

EEAP

- 1 Embed
- 2 Encode
- 3 Attend
- 4 Predict

Pet problem to work with

Encoder

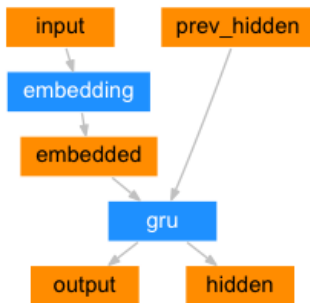


Figure: Our encoder model⁸

⁸https://pytorch.org/tutorials/intermediate/seq2seq_translation_tutorial.html

Pet problem to work with Attention

- We calculate attention weights based on **decoder input** and a hidden state
- We multiply it with **encoder outputs** to create a weighted combination
- It will help the decoder to "find out" which part of encoder output is "responsible" for which part of decoder output

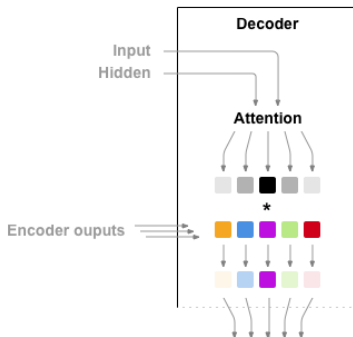


Figure: Way of computing attention⁸

Pet problem to work with

Decoder

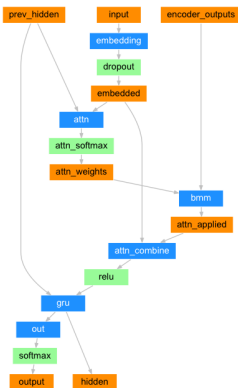


Figure: Our attention-decoder model⁸

Time to code

Now we will do a little bit of live coding

Where can You find the code?

```
https:  
//github.com/tugot17/paying-attention-to-attention
```

Attention applications

Automatic image captioning

- Nearly the same solution as previously
- Instead of encoder output, latent space from a pre-trained state-of-the-art network (e.g Inception V3)
- Attention method similar as before (e.g Bahdanau Attention)

Prediction Caption: the person is riding a surfboard in the ocean <end>

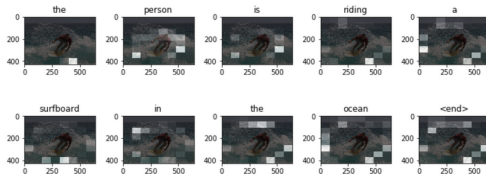


Figure: Show, Attend and Tell: Neural Image Caption Generation with Visual Attention[6]

Attention applications

Attention UNet

- State-of-the-art image segmentation solution
- Improving model sensitivity and accuracy by attaching attention gates on top of the standard U-Net

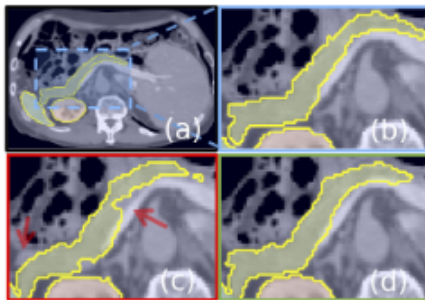


Figure: Segmentation for medical images with U-Net[7]

Attention applications

Latex code generation

- Seq2seq very similar to the Image Captioning
- Uses attention based RNN to generate Latex Code
- Real world application: <https://mathpix.com/>

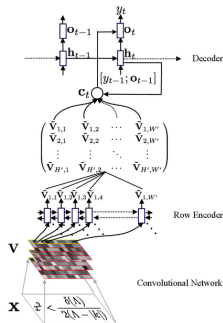


Figure: What You Get Is What You See:A Visual Markup Decompiler[8]

Attention applications

Many others, mainly some kind of seq2seq models

- Bert[5], RoBerta[9], ALBert[10], *Bert*
- GPT 2[2]
- Transformer - Attention is all You need[3]
- Residual Attention Network for Image Classification[4]
- TreeGen: A Tree-Based Transformer Architecture for Code Generation[11]
- and many others ...

Summary

Many others, mainly some kind of seq2seq models

- Attention = teach Your NN
to on what part of data
should it be focused on

Summary

Many others, mainly some kind of seq2seq models

- Attention = teach Your NN to on what part of data should it be focused on
- Attention is all You need - we can solve a wide variety of problems using attention

Summary

- Attention = teach Your NN to on what part of data should it be focused on
- Attention is all You need - we can solve a wide variety of problems using attention
- We can use nearly same solution for many problems
- Embed, Encode, Attend, Predict Framework

Summary

- Attention = teach Your NN to on what part of data should it be focused on
- Attention is all You need - we can solve a wide variety of problems using attention
- We can use nearly same solution for many problems
- Embed, Encode, Attend, Predict Framework
- Hot research topic - It's hard to stay up to date but it is relatively easy to publish

Summary

- Attention = teach Your NN to on what part of data should it be focused on
- Attention is all You need - we can solve a wide variety of problems using attention
- We can use nearly same solution for many problems
- Embed, Encode, Attend, Predict Framework
- Hot research topic - It's hard to stay up to date but it is relatively easy to publish
- There are several methods for calculating attention (e.g Bahdanau Attention, Luong attention, etc.)

References I

- [1] Y. Bengio D. Bahdanau, K. Cho.
Neural machine translation by jointly learning to align and translate.
2014.
- [2] A. Radford et al.
Language models are unsupervised multitask learners.
2019.
- [3] A. Vaswani et al.
Attention is all you need.
2017.

References II

- [4] F. Wang et al.
Residual attention network for image classification.
2019.
- [5] J. Devlin et al.
Bert: Pre-training of deep bidirectional transformers for
language understanding.
2018.
- [6] K. Xu et al.
Show, attend and tell: Neural image caption generation with
visual attention.
2016.

References III

- [7] O. Oktay et al.
Attention u-net: Learning where to look for the pancreas.
2018.
- [8] Y. Deng et al.
What you get is what you see: a visual markup decompiler.
2016.
- [9] Y. Llu et al.
Roberta: A robustly optimized bert pretraining approach.
2019.

References IV

- [10] Z. Lan et al.
Albert: A lite bert for self-supervised learning of language representations.
2019.
- [11] Z. Sun et al.
Treegen: A tree-based transformer architecture for code generation.
2019.
- [12] N. L. Zhang L. Xue, X. Li.
Not all attention is needed: Gated attention network.
2019.

Thank You for Your attention