And the second s

PROPERTIES

PRESS PRINT

THE REAL

THE CLEAR DR. D. D.

even of the loss

AT SCHOOL PROVIDENT OF THE SCHOOL PROVIDENT

Attention and Transformer

Prof. Kuan-Ting Lai 2022/5/11

Attention Test: Where is the 4-leaf clover?



https://www.mentalup.co/blog/concentration-test

Attention

• Visual attention and textual attention





https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

Attention = Vector of Importance Weights

- Use weights to assign importance of input data
- Types of attention
 - -Implicit vs. Explicit
 - Max pooling is implicit attention
 - -Hard vs. Soft
 - Soft attention is described by continuous variables
 - Hard attention is described by discrete variables
 - -Additive vs. Multiplicative



Seq2seq model (Language Translation)



https://towardsdatascience.com/introduction-to-rnns-sequence-to-sequence-language-translation-and-attention-fc43ef2cc3fd 5

Training and Inference of Seq-to-Seq



6

Homonyms: Multiple-meaning Words

• Words are "Context-aware".

• date

Her favorite fruit to eat is a **date**. John took Mary out on a **date**. What is your **date** of birth?

• fall

I love cool, crisp **fall** weather. Don't **fall** on your way to the gym.

https://grammar.yourdictionary.com/for-students-and-parents/words-with-multiple-meanings.html

Additive Attention in RNN



Attention is All You Need!

Ashish Vaswani* Google Brain avaswani@google.com

Noam Shazeer* Google Brain noam@google.com ni

Niki Parmar*Jakob Uszkoreit*Google ResearchGoogle Researchnikip@google.comusz@google.com

9

Llion Jones* Google Research llion@google.com
 Aidan N. Gomez* †
 Łukasz Kaiser*

 University of Toronto
 Google Brain

 aidan@cs.toronto.edu
 lukaszkaiser@google.com

Illia Polosukhin* [‡] illia.polosukhin@gmail.com

NIPS, 2017

Google Brain & University of Toronto

The Transformer Model

Multi-head attention

- Self-attention in encoders
- Masked Self-attention in decoders
- Encoder-decoder attention
- Positional encoding



Attention Module in Transformer

• Query (Q), Key (K), Value (V) attention

Attention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



Visualizing Attention

• Tensor2Tensor Notebook

https://colab.research.google.com/github/tenso rflow/tensor2tensor/blob/master/tensor2tensor /notebooks/hello_t2t.ipynb

Inputs: The animal didn't cross the street because it was too tired

Outputs: Das Tier überquerte die Straße nicht, weil es zu müde war, weil es zu müde war.



Self-attention Example

 "animal" and "it" have a high relevancy score



Self Attention



Self-attention Pseudocode

```
Compute the dot
                                             Iterate over each token
                                                                      product (attention
                                             in the input sequence.
                                                                      score) between the
def self attention(input sequence):
                                                                      token and every
    output = np.zeros(shape=input sequence.shape)
                                                                      other token.
    for i, pivot vector in enumerate(input sequence):
        scores = np.zeros(shape=(len(input sequence),))
                                                                         Scale by a
        for j, vector in enumerate(input sequence):
                                                                         normalization
             scores[j] = np.dot(pivot vector, vector.T)
                                                                         factor, and apply
        scores /= np.sqrt(input sequence.shape[1])
        scores = softmax(scores)
        new pivot representation = np.zeros(shape=pivot vector.shape)
        for j, vector in enumerate(input sequence):
             new pivot representation += vector * scores[j]
                                                                         Take the sum
        output[i] = new pivot_representation
                                                                         of all tokens
    return output
                                                                         weighted by the
                                                That sum is
                                                                         attention scores.
                                               our output.
```

Query, Keys, Values





Multi-head Attention

- Use multiple attention modules and concatenate outputs
- Like depthwise separable convolution



Transformer Encoder

- Multi-head attention
- Dense network
- Residual connection
- Layer normalization



import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers

class TransformerEncoder(layers.Layer): **def** init (self, embed dim, dense dim, num heads, **kwargs): super(). init (**kwargs) self.embed dim = embed dim Size of the input self.dense dim = dense dim token vectors self.num heads = num heads self.attention = layers.MultiHeadAttention(Size of the inner num heads=num heads, key dim=embed dim) dense layer self.dense proj = keras.Sequential([layers.Dense(dense dim, activation="relu"), Number of layers.Dense(embed dim),] attention heads self.layernorm 2 = layers.LayerNormalization() Computation goes in call(). def call(self, inputs, mask=None): The mask that will be generated by if mask is not None: the Embedding layer will be 2D, but mask = mask[:, tf.newaxis, :] the attention layer expects to be 3D attention output = self.attention(or 4D, so we expand its rank. inputs, inputs, attention mask=mask) proj input = self.layernorm 1(inputs + attention output) proj output = self.dense proj(proj input) return self.layernorm 2(proj input + proj output) **def** get config(self): Implement config = super().get_config() serialization so config.update({ we can save the "embed dim": self.embed dim, model. "num heads": self.num heads, "dense dim": self.dense dim, 21 return config

Transformer Encoder

Layer Normalization

```
Input shape: (batch_size,
                                                            sequence_length, embedding_dim)
def layer normalization(batch of sequences):
    mean = np.mean(batch of sequences, keepdims=True, axis=-1)
    variance = np.var(batch of sequences, keepdims=True, axis=-1)
    return (batch of sequences - mean) / variance
                                                                 To compute mean and
                                                            variance, we only pool data
                                                             over the last axis (axis -1).
Compare to BatchNormalization (during training):
                                                            Input shape: (batch_size,
height, width, channels)
def batch normalization(batch of images):
    mean = np.mean(batch_of_images, keepdims=True, axis=(0, 1, 2))
    variance = np.var(batch of images, keepdims=True, axis=(0, 1, 2))
    return (batch of images - mean) / variance
                                                              Pool data over the batch axis
                                                         (axis 0), which creates interactions
                                                              between samples in a batch.
```

Transformer Encoder for Text Classification

```
vocab_size = 20000
embed_dim = 256
num_heads = 2
dense_dim = 32
```

Since TransformerEncoder returns full sequences, we need to reduce each sequence to a single vector for classification via a global pooling layer.

Different Types of NLP Models

	Word order	Context awareness
	awareness	(cross-words interactions)
Bag-of-unigrams	No	No
Bag-of-bigrams	Very limited	No
RNN	Yes	No
Self-attention	No	Yes
Transformer	Yes	Yes

```
A downside of position embeddings is that the sequence
                                                                                         length needs to be known in advance.
                                    class PositionalEmbedding(layers.Layer):
                                        def init (self, sequence length, input dim, output dim, **kwargs): <----</pre>
                                             super(). init (**kwargs)
                                                                                                                And another
                                             self.token embeddings = layers.Embedding(
                            Prepare an
                                                                                                                one for
                                                 input dim=input dim, output dim=output dim)
                            Embedding
                                                                                                                the token
                                             self.position embeddings = layers.Embedding(
                           layer for the
                                                                                                                positions
                                                 input dim=sequence length, output dim=output dim)
                         token indices.
                                             self.sequence length = sequence length
                                             self.input dim = input dim
                                             self.output dim = output dim
Positional
                                        def call(self, inputs):
                                             length = tf.shape(inputs)[-1]
                            Add both
Encoding
                                             positions = tf.range(start=0, limit=length, delta=1)
                           embedding
                                             embedded tokens = self.token embeddings(inputs)
                              vectors
                                             embedded positions = self.position embeddings(positions)
                            together.
                                             return embedded tokens + embedded positions
                                                                                                Like the Embedding layer, this
                                                                                                layer should be able to generate a
                                        def compute mask(self, inputs, mask=None):
                                                                                                mask so we can ignore padding 0s
                                             return tf.math.not equal(inputs, 0)
                         Implement
                                                                                                in the inputs. The compute mask
                         erialization
                                                                                                method will called automatically
                                        def get config(self):
                          so we can
                                                                                                by the framework, and the
                                             config = super().get config()
                           save the
                                                                                                mask will get propagated
                                             config.update({
                            model.
                                                                                                to the next layer.
                                                 "output dim": self.output dim,
                                                 "sequence length": self.sequence length,
                                                 "input dim": self.input dim,
                                             })
                                                                                                                  25
                                             return config
```

```
vocab size = 20000
sequence length = 600
embed dim = 256
num heads = 2
dense dim = 32
                                                                      Look here!
inputs = keras.Input(shape=(None,), dtype="int64")
x = PositionalEmbedding(sequence length, vocab size, embed dim)(inputs)
x = TransformerEncoder(embed dim, dense dim, num heads)(x)
x = layers.GlobalMaxPooling1D()(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(optimizer="rmsprop",
              loss="binary crossentropy",
              metrics=["accuracy"])
model.summary()
callbacks = [
    keras.callbacks.ModelCheckpoint("full transformer encoder.keras",
                                     save best only=True)
model.fit(int train ds, validation data=int val ds, epochs=20,
     callbacks=callbacks)
model = keras.models.load model(
    "full transformer encoder.keras",
    custom objects={"TransformerEncoder": TransformerEncoder,
                    "PositionalEmbedding": PositionalEmbedding})
print(f"Test acc: {model.evaluate(int test ds)[1]:.3f}")
                                                                       26
```

Transformer Encoder with Positional Embedding

IMDB Movie Review Classification



Bag-of-Words or Sequence Models?

• Chollet, et al. (http://mng.bz/AOzK)



Transformer Encoder + Decoder



TransformerDecoder



- 1. François Chollet, "Deep Learning with Python, 2nd edition", Chapter 11, 2021
- 2. <u>https://towardsdatascience.com/introduction-to-rnns-sequence-to-sequence-language-translation-and-attention-fc43ef2cc3fd</u>
- 3. Vaswani et al., "Attention is All You Need," NIPS, 2017