# Lecture Notes for
# **Machine Learning in Python**

[ 👨‍🏫 , 👨‍💻 , 🐍 , 👨‍🔬 ]

## Professor Eric Larson
## **Sequential CNN and Transformers**

**In progress** lecture replacing detailed implementation of recurrent networks

# Lecture Agenda

- Logistics
  - Grading Update
  - Sequential Networks due **Last Day of Finals**
- Agenda
  - CNNs for Sequential Processing
  - Transformers

# Class Overview, by topic

**Table Data Visualization**

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

**Dimension Reduction and Image Processing**

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

**Linear and Logistic Regression**

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

**Neural Networks and Back Prop.**

Numpy
Detailed mathematics for NN operations
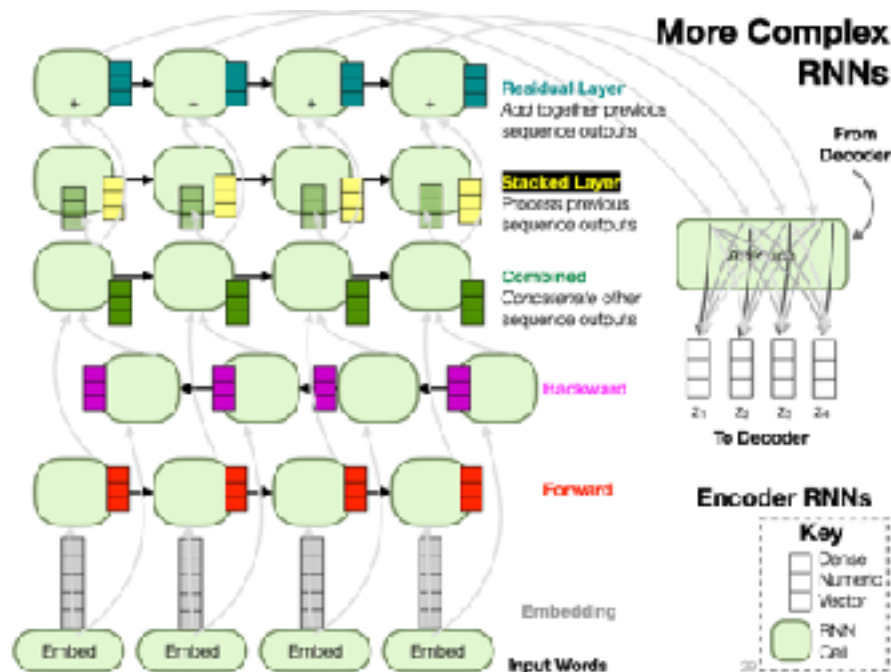
**Wide and Deep Networks**

**Convolutional Networks**

**Sequential Networks**
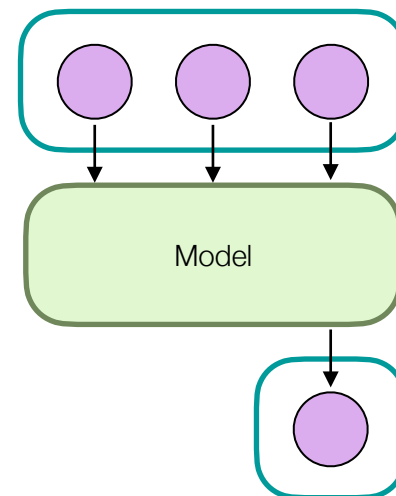
Keras, Tensorflow
Intuition, Detailed implement.

**Ethics in Language Models**

ConceptNet
Case studies

# CNNs for Sequences

# CNNs for Sequences

The    quick    brown    fox    jumped    over    the    lazy    dogs

Embed

Many Filters

Filter Size determines how many elements can be looked at simultaneously (i.e., phrases, here kernel length == 3)

Filter Output

Each channel is output from separate filter

Perform Pooling across sequence

Stride of pooling determines downsampling here stride==2

# CNNs for Sequences

- RNNs are not inherently parallelized or efficient at remembering based on state vector, but CNNs can be run in parallel groups

Input Sequence

Filtering + Activation

Output Activations

…

Pooling

Pooled Activations

…

Filtering + Activation + Pooling

Second Conv Layer

…

**Can we do something smarter here?**

Flatten

Fully Connected Layer

Predictions

- Everything we learned in 2D CNNs can be applied to 1D CNNs…
- Residuals, separable convolution, squeezing, everything

Input Sequence

Filtering + Activation

Output Activations

...

Pooling

Pooled Activations

...

Filtering + Activation + Pooling

Second Conv Layer

...

Mean across channels

Variance across channels

Min/Max across channels

**Does it matter how long the input sequence was?**

Concatenate + Flatten

Fully Connected Layer

Predictions

# CNNs for Sequences

Input Sequence

Filtering + Activation

Output Activations

…

Pooling

Pooled Activations

…

Filtering + Activation + Pooling

Second Conv Layer

…

Mean across channels

Variance across channels

Min/Max across channels

**This is also a popular method for processing voice characteristics!**

Concatenate + Flatten

Fully Connected Layer

Predictions

The Sequential CNN
IMdB sentiment analysis

```
13a. Sequence Basics [Experimental].ipynb
```

# Transformers

**Dr Simone Stumpf** @DrSimoneS... · 13h  ···
God grant me the confidence of an average
machine learning expert.

- Recurrent networks track state using an "updatable" state vector, but this takes processing iterative

- Attention mechanism (in RNNs) already takes a weighted sum of state vectors to generate new token in a decoder

- … so why not just use attention on a transformation of the embedding vectors? **Do away with the recurrent state vector all together?**

# Attention is All You Need

- **Continued Motivation**:
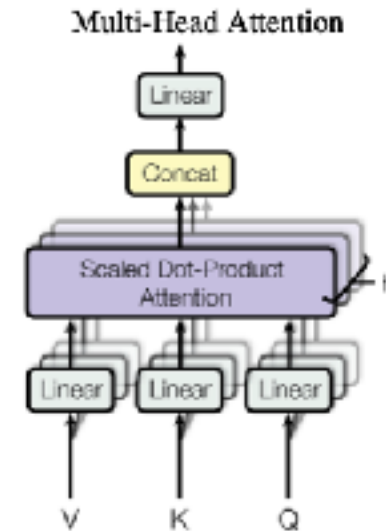  - RNNs are not inherently parallelized or efficient at remembering based on state vector
  - CNNs are not resilient to long-term word relationships, limited by filter size
- **Transformer Solution:**
  - Build attention into model from the **beginning**
  - Compare all words to each other through **self-attention**
  - Define a notion of "**position**"in the sequence
  - ***Should be resilient to long term relationships and be highly parallelized for GPU computing!!***

Multi-Head Attention

Scaled Dot-Product Attention

for each word

more than one
Q,K,V use in document

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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

# Transformer: in more detail

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $X_1$ ▭▭▭▭ | $X_2$ ▭▭▭▭ |

Multiply These
(in parallel)

**Outputs of Matrix Multiplications:**

| | | |
|---|---|---|
| Queries | $q_1$ ▭▭ | $q_2$ ▭▭ |
| Keys | $k_1$ ▭▭ | $k_2$ ▭▭ |
| Values | $v_1$ ▭▭ | $v_2$ ▭▭ |

Learned Matrices

$W^Q$

$W^K$

$W^V$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

46

# Transformer: in more detail

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ( $\sqrt{d_k}$ )
in visual, $d_k = 3$

Softmax

Softmax
X
Value

Sum

**Thinking** **Machines**

$x_1$  $x_2$

Calc. q, k, v for each word

$q_1$  $q_2$

$k_1$  $k_2$

$v_1$  $v_2$

$q_1 \bullet k_1 = 112$   $q_1 \bullet k_2 = 96$

Divide   Divide

14   12

Softmax

0.88   0.12

Multiply   Calc weights   Multiply

$v_1$  $v_2$

weighted sum for all words in document

Sum

$z_1$ attention for word 1   $z_2$ attention for word 2

Straight forward to do this operation in matrix form:

X   WQ   Q

Thinking Machines   $q_1$ $q_2$   $\leftarrow d_k \rightarrow$

X   WK   K

Thinking Machines   $k_1$ $k_2$   $\leftarrow d_k \rightarrow$

X   WV   V

Thinking Machines   $v_1$ $v_2$   $\leftarrow d_v \rightarrow$

Q   KT   V

$\text{softmax}\left( \dfrac{\phantom{xx} \times \phantom{xx}}{\sqrt{d_k}} \right)$

Z

$=$   $z_1$ $z_2$

Size of W matrices:
$W^V$: |Embed Size| x $d_v$
$W^{Q,K}$: |Embed Size| x $d_k$

Size of Q,K,V:
|Seq Len| x $d_v$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

Thinking Machines

$X$
Head 0
Head 1
Head 7

$W_0^Q$ $W_0^K$ $W_0^V$
$W_1^Q$ $W_1^K$ $W_1^V$
$W_7^Q$ $W_7^K$ $W_7^V$

$Q_0$ $K_0$ $V_0$
$Q_1$ $K_1$ $V_1$
$Q_7$ $K_7$ $V_7$

$Z_0$ $Z_1$ $Z_7$

concat $\times$ $W^O$ $=$ $Z$

one row for each word

Num columns determined by $W^O$

$Z_0$ $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$ $Z_6$ $Z_7$

$\times W^O = Z$

(words x $Z_{concat}$)

($Z_{concat}$ x Embed)   (words x Embed)

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

Prediction

⊞⊞⊞⊞⊞⊞

Fully Connected

Global Average

Layer Normalization

Fully Connected

Layer Normalization

```
tf.keras.layers.MultiHeadAttention(
    num_heads,      (Number of heads $Z_1$-$Z_7$)
    key_dim,        (size of query/key $d_k$)
    value_dim,      (size of each $d_v$)
    output_shape,   (Embed size of Z)
    …
```

$W^Q$

$\leftarrow d_k \rightarrow$

$W^K$

$\leftarrow d_k \rightarrow$

$W^V$

$\leftarrow d_v \rightarrow$

Add & Norm

Position-wise FFN

Add & Norm

Multi-Head Attention

$L\times$

Residual Connect

Positional Encodings ⊕

?

Token Embedding

Inputs

$$LN(Z_{row}) = \alpha \frac{Z_i - \mu_Z}{\sqrt{\sigma_Z^2 + \epsilon}} + \beta$$

Learn to normalize the rows of Z

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

x $W^o$ = Z

(words x $Z_{concat}$)

($Z_{concat}$ x Embed)

(words x Embed)

# Transformer: Positional Encoding

- Objective: add notion of position to embedding
- Attempt in original paper: add sin/cos to embedding
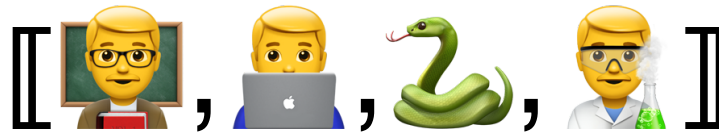- But could be anything that encodes position, like:

0, 1, 2, 3, 4, 5, … T          The quick brown fox …

Position Embed

Word Embed

Add

"Brown" embedding at Third position

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

The Transformer
        and 20 news groups with GloVe

```
13a. Sequence Basics [Experimental].ipynb
```

# Lecture Notes for
# **Machine Learning in Python**

[ 👨‍🏫 , 👨‍💻 , 🐍 , 👨‍🔬 ]

## Professor Eric Larson
## **Sequential CNN and Transformers**