

# EDAN95

## Applied Machine Learning

### Lecture 10: Sequence Prediction

Pierre Nugues

Pierre.Nugues@cs.lth.se  
[http://cs.lth.se/pierre\\_nugues/](http://cs.lth.se/pierre_nugues/)

December 1st, 2021

# Outline

In the previous lecture, we used networks to produce one output  $y$  per input vector  $\mathbf{x}$ , for instance one category per sentence.

Given an input sequence  $\mathbf{x}$ , we will now produce an output sequence:  $\mathbf{y}$ .

We will experiment four kinds of neural networks:

- ① Feed forward
- ② Recurrent
- ③ LSTM
- ④ Transformers

In the laboratory assignment, you will experiment with items 2. and 3.

# Motivation

The analysis of sentences often involves the analysis of words.

We can divide it in three main tasks:

- 1 Identify the type of word, for instance noun or verb using the classical grammar;
- 2 Identify a group or segment, for instance are these three words, *Kjell Olof Andersson*, the name of a person;
- 3 Identify the relations between two words: for instance is this group the subject of a verb? This corresponds to parsing, semantic analysis, or information extraction.

We will consider the two first tasks.

This lecture will show you how to solve the first one, part-of-speech tagging, and you will write a program for the second one, named entity recognition (NER), in the next laboratory assignment.

# Word Categorization: The Parts of Speech

Sentence:

*That round table might collapse*

Annotation:

Words	Parts of speech	POS tags
<b>that</b>	Determiner	DET
<b>round</b>	Adjective	ADJ
<b>table</b>	Noun	NOUN
<b>might</b>	Modal verb	AUX
<b>collapse</b>	Verb	VERB

The automatic annotation uses predefined POS tagsets such as the Penn Treebank tagset for English

# Ambiguity

Words	Possible tags	Example of use
<b>that</b>	Subordinating conjunction Determiner Adverb Pronoun Relative pronoun	<i>That he can swim is good</i> <i>That white table</i> <i>It is not that easy</i> <i>That is the table</i> <i>The table that collapsed</i>
<b>round</b>	Verb Preposition Noun Adjective Adverb	<i>Round up the usual suspects</i> <i>Turn round the corner</i> <i>A big round</i> <i>A round box</i> <i>He went round</i>
<b>table</b>	Noun Verb	<i>That white table</i> <i>I table that</i>
<b>might</b>	Noun Modal verb	<i>The might of the wind</i> <i>She might come</i>
<b>collapse</b>	Noun Verb	<i>The collapse of the empire</i> <i>The empire can collapse</i>

# Training Sets: The CoNLL Format

The CoNLL format is a tabular format to distribute annotated texts. This format was created for evaluations carried out by the Conference in natural language learning

The CoNLL annotation has varied much across the years. We use CoNLL-U, the latest iteration.

Annotation of the Spanish sentence:

*La reestructuración de los otros bancos checos se está acompañando por la reducción del personal*

*'The restructuring of Czech banks is accompanied by the reduction of personnel'*

# Example of Annotation (CoNLL-U)

*La reestructuración de los otros bancos checos se está acompañando por la reducción del personal*

ID	FORM	LEMMA	UPOS	FEATS
1	La	el	DET	Definite=Def Gender=Fem Number=Sing PronType=Art
2	reestructuración	reestructuración	NOUN	Gender=Fem Number=Sing
3	de	de	ADP	AdpType=Prep
4	los	el	DET	Definite=Def Gender=Masc Number=Plur PronType=Art
5	otros	otro	DET	Gender=Masc Number=Plur PronType=Ind
6	bancos	banco	NOUN	Gender=Masc Number=Plur
7	checos	checo	ADJ	Gender=Masc Number=Plur
8	se	se	PRON	Case=Acc Person=3 PrepCase=Npr PronType=Prs Reflex=Yes
9	está	estar	AUX	Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin
10	acompañando	acompañar	VERB	VerbForm=Ger
11	por	por	ADP	AdpType=Prep
12	la	el	DET	Definite=Def Gender=Fem Number=Sing PronType=Art
13	reducción	reducción	NOUN	Gender=Fem Number=Sing
14	del	del	ADP	AdpType=Preppron
15	personal	personal	NOUN	Gender=Masc Number=Sing
16	.	.	PUNCT	PunctType=Peri

# Another Example

ID	FORM	LEMMA	PLEMMA	POS	PPOS	FEAT	PFEAT
1	Battle	battle	battle	NN	NN	—	—
2	-	-	-	HYPH	HYPH	—	—
3	tested	tested	tested	NN	NN	—	—
4	Japanese	japanese	japanese	JJ	JJ	—	—
5	industrial	industrial	industrial	JJ	JJ	—	—
6	managers	manager	manager	NNS	NNS	—	—
7	here	here	here	RB	RB	—	—
8	always	always	always	RB	RB	—	—
9	buck	buck	buck	VBP	VB	—	—
10	up	up	up	RP	RP	—	—
11	nervous	nervous	nervous	JJ	JJ	—	—
12	newcomers	newcomer	newcomer	NNS	NNS	—	—
13	with	with	with	IN	IN	—	—
14	the	the	the	DT	DT	—	—
15	tale	tale	tale	NN	NN	—	—
16	of	of	of	IN	IN	—	—
17	the	the	the	DT	DT	—	—
18	first	first	first	JJ	JJ	—	—
19	of	of	of	IN	IN	—	—
20	their	their	their	PRP\$	PRP\$	—	—
21	countrymen	countryman	countryman	NNS	NNS	—	—
22	to	to	to	TO	TO	—	—
23	visit	visit	visit	VB	VB	—	—
24	Mexico	mexico	mexico	NNP	NNP	—	—
25	,	,	,	,	,	—	—
26	a	a	a	DT	DT	—	—
27	boatload	boatload	boatload	NN	NN	—	—
28	of	of	of	IN	IN	—	—
29	samurai	samurai	samurai	NN	NN	—	—
30	warriors	warrior	warrior	NNS	NNS	—	—
31	blown	blow	blow	VRN	VRN	—	—



# Designing a Part-of-Speech Tagger

We will now create part-of-speech taggers, where we will examine three architectures:

- 1 A feed-forward pipeline with a one-hot encoding of the words;
- 2 A feed-forward pipeline with word embeddings: We will replace the one-hot vectors with GloVe embeddings;
- 3 A recurrent neural network, either a simple RNN or a LSTM, with word embeddings.

# Features for Part-of-Speech Tagging

The word *visit* is ambiguous in English:

*I paid a **visit** to a friend* → *noun*

*I went to **visit** a friend* → *verb*

The context of the word enables us to tell, here an article or the infinitive marker

To train and apply the model, the tagger extracts a set of features from the surrounding words, for example, a sliding window spanning five words and centered on the current word.

We then associate the feature vector  $(w_{i-2}, w_{i-1}, w_i, w_{i+1}, w_{i+2})$  with the part-of-speech tag  $t_i$  at index  $i$ .

# Part-of-Speech Tagging

ID	FORM	PPOS	
	BOS	BOS	Padding
	BOS	BOS	
1	Battle	NN	
2	-	HYPH	
3	tested	NN	
...	...	...	
17	the	DT	
18	first	JJ	
19	of	IN	
20	their	PRP\$	
21	countrymen	NNS	Input features
22	to	TO	
23	visit	VB	Predicted tag
24	Mexico		↓
25	,		
26	a		
27	boatload		
...	...	...	
34	years		
35	ago		
36	.		
	EOS		Padding
	EOS		

# Feature Vectors

ID	Feature vectors							PPOS
	$w_{i-2}$	$w_{i-1}$	$w_i$	$w_{i+1}$	$w_{i+2}$	$t_{i-2}$	$t_{i-1}$	
1	BOS	BOS	Battle	-	tested	BOS	BOS	NN
2	BOS	Battle	-	tested	Japanese	BOS	NN	HYPH
3	Battle	-	tested	Japanese	industrial	NN	HYPH	JJ
...	...	...	...	...	...	...	...	...
19	the	first	of	their	countrymen	DT	JJ	IN
20	first	of	their	countrymen	to	JJ	IN	PRP\$
21	of	their	countrymen	to	visit	IN	PRP\$	NNS
22	their	countrymen	to	visit	Mexico	PRP\$	NNS	TO
23	countrymen	to	visit	Mexico	,	NNS	TO	VB
24	to	visit	Mexico	,	a	TO	VB	NNP
25	visit	Mexico	,	a	boatload	VB	NNP	,
...	...	...	...	...	...	...	...	...
34	ashore	375	years	ago	.	RB	CD	NNS
35	375	years	ago	.	EOS	CD	NNS	RB
36	years	ago	.	EOS	EOS	NNS	RB	.

# Architecture 1: A Feed-Forward Neural Network

We first use a feed-forward architecture corresponding to a logistic regression:

```
np.random.seed(0)

model = models.Sequential([Dense(NB_CLASSES,
                                  input_dim=X.shape[1],
                                  activation='softmax')])

model.compile(loss='sparse_categorical_crossentropy',
              optimizer=OPTIMIZER,
              metrics=['accuracy'])

model.summary()

model.fit(X, y, epochs=EPOCHS, batch_size=BATCH_SIZE)

model.save('out.model')
```

# Encoding the $\mathbf{y}$ Vector

In the previous examples, we used `categorical_crossentropy`. This requires that all the targets are encoded with one-hot vectors. For instance:

- determiner: [1, 0, 0, 0]
- noun: [0, 1, 0, 0]
- verb: [0, 0, 1, 0]
- adjective: [0, 0, 0, 1]

With `sparse_categorical_crossentropy`, we can use numerical indices:

- determiner: 1
- noun: 2
- verb: 3
- adjective: 4

We do not need to use the `to_categorical` function.

# Preprocessing

Preprocessing is more complex though: Four steps:

- 1 Read the corpus

```
train_sentences, dev_sentences, test_sentences, \
    column_names = load_ud_en_ewt()
```

- 2 Store the rows of the CoNLL corpus in dictionaries

```
conll_dict = CoNLDDictorizer(column_names, col_sep='\t')
train_dict = conll_dict.transform(train_sentences)
test_dict = conll_dict.transform(test_sentences)
```

- 3 Extract the features and store them in dictionaries

```
context_dictorizer = ContextDictorizer()
context_dictorizer.fit(train_dict)
X_dict, y_cat = context_dictorizer.transform(train_dict)
```

- 4 Vectorize the symbols

```
# We transform the X symbols into numbers
dict_vectorizer = DictVectorizer()
X_num = dict_vectorizer.fit_transform(X_dict)
```

# Code Example

Jupyter Notebook: [4.1-nn-pos-tagger.ipynb](#)



# Architecture 2: Using Embeddings

We replace the one-hot vectors with embeddings, the rest being the same  
Word embeddings are dense vectors obtained by a principal component analysis or another method.

They can be trained by the neural network or pretrained  
In this implementation:

- ① We use pretrained embeddings from the GloVe project;
- ② Our version of GloVe is lowercased, so we set all the characters in lowercase;
- ③ We add the embeddings as an `Embedding` layer at the start of the network;
- ④ We initialize the embedding layer with GloVe and make it trainable or not.

It would be possible to use a randomly initialized matrix as embeddings instead

# The Embedding Layer

```
model = models.Sequential([
    Embedding(cnt_uniq, EMBEDDING_DIM,
              input_length=2 * W_SIZE + 1),
    Flatten(),
    Dense(NB_CLASSES, activation='softmax')
])

if embedding_matrix is not None:
    model.layers[0].set_weights([embedding_matrix])
model.layers[0].trainable = True

model.compile(loss='sparse_categorical_crossentropy',
              optimizer=OPTIMIZER,
              metrics=['accuracy'])

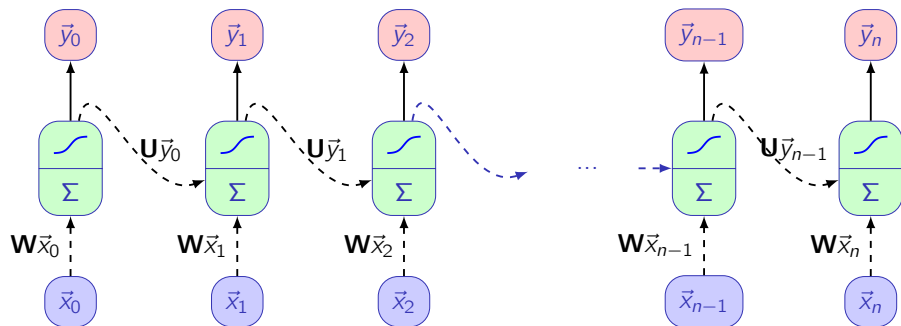
model.summary()

model.fit(X, y, epochs=EPOCHS, batch_size=BATCH_SIZE)
```

# Code Example

Jupyter Notebook: `4.2-nn-pos-tagger-embeddings.ipynb`

# The RNN Architecture



# Input Format for RNNs

The input format is different from feed forward networks.

We need to build two lists: one for the input and the other for the output

<b>y</b>	DET	NOUN	VERB	DET	NOUN
<b>x</b>	The	waiter	brought	the	meal

All the vectors in a same batch must have the same length. We pad them:

<b>y</b>	PAD	PAD	PAD	DET	NOUN	VERB	DET	NOUN
<b>x</b>	PAD	PAD	PAD	The	waiter	brought	the	meal

We could apply the padding after

# Building the Sequences

```
def build_sequences(corpus_dict, key_x='form', key_y='pos',
                    tolower=True):
    X, Y = [], []
    for sentence in corpus_dict:
        x, y = [], []
        for word in sentence:
            x += [word[key_x]]
            y += [word[key_y]]
        if tolower:
            x = list(map(str.lower, x))
        X += [x]
        Y += [y]
    return X, Y
```

At this point, we have **x** and **y** vectors of symbols

# Building Index Sequences

0 is for the padding symbol and 1 for the unknown words

```
idx_word = dict(enumerate(vocabulary_words, start=2))  
idx_pos = dict(enumerate(pos, start=2))  
word_idx = {v: k for k, v in idx_word.items()}  
pos_idx = {v: k for k, v in idx_pos.items()}
```

At this point, we have **x** and **y** vectors of numbers

# Padding the Index Sequences

We build the complete **X\_idx** and **Y\_idx** matrices for the whole corpus  
And we pad the matrices:

```
X = pad_sequences(X_idx)
```

```
Y = pad_sequences(Y_idx)
```

```
# The number of POS classes and 0 (padding symbol)
```

```
Y_train = to_categorical(Y, num_classes=len(pos) + 2)
```

`pad_sequences` can have an argument that specifies the maximal length  
`maxlen (MAX_SEQUENCE_LENGTH)`.

The padded sentences must have the same length in a batch. This is  
automatically computed by Keras



# Recurrent Neural Networks (RNN)

```
model = models.Sequential([
    Embedding(len(vocabulary_words) + 2,
              EMBEDDING_DIM,
              mask_zero=True,
              input_length=None),
    SimpleRNN(100, return_sequences=True),
    # Bidirectional(SimpleRNN(100, return_sequences=True)),
    Dense(NB_CLASSES + 2, activation='softmax')])

model.layers[0].set_weights([embedding_matrix])
# The default is True
model.layers[0].trainable = True
```

# Parameters

Keras functions have many parameters.

In case of doubt, read the documentation

A few useful parameters:

- 1 `mask_zero=True` is to tell whether or not the input value 0 is a special “padding” value;
- 2 `return_sequences=True` tells whether to return the last output in the output sequence, or the full sequence. In sequences, it is essential;
- 3 `recurrent_dropout=0.3` tells how much to drop for the linear transformation of the recurrent state.

# Code Example

Jupyter Notebook: `4.3-rnn-pos-tagger.ipynb`

# Long Short-Term Memory (LSTM)

```
model = models.Sequential([
    Embedding(len(vocabulary_words) + 2,
              EMBEDDING_DIM,
              mask_zero=True,
              input_length=None),
    Bidirectional(LSTM(100, return_sequences=True)),
    Dense(NB_CLASSES + 2, activation='softmax')])

model.layers[0].set_weights([embedding_matrix])
# The default is True
model.layers[0].trainable = True
```

# Segment Recognition

## Group detection – chunking –:

Brackets: [<sub>NG</sub> The government <sub>NG</sub>] has [<sub>NG</sub> other agencies and instruments <sub>NG</sub>] for pursuing [<sub>NG</sub> these other objectives <sub>NG</sub>] .

Tags: *The/I government/I has/O other/I agencies/I and/I instruments/I for/O pursuing/O these/I other/I objectives/I ./O*

Brackets: Even [<sub>NG</sub> Mao Tse-tung <sub>NG</sub>] [<sub>NG</sub> 's China <sub>NG</sub>] began in [<sub>NG</sub> 1949 <sub>NG</sub>] with [<sub>NG</sub> a partnership <sub>NG</sub>] between [<sub>NG</sub> the communists <sub>NG</sub>] and [<sub>NG</sub> a number <sub>NG</sub>] of [<sub>NG</sub> smaller, non-communists parties <sub>NG</sub>] .

Tags: *Even/O Mao/I Tse-tung/I 's/B China/I began/O in/O 1949/I with/O a/I partnership/I between/O the/I communists/I and/O a/I number/I of/O smaller/I ,/I non-communists/I parties/I ./O*

# Segment Categorization

Tages extendible to any type of chunks: nominal, verbal, etc.

For the IOB scheme, this means tags such as I.Type, O.Type, and B.Type, Types being NG, VG, PG, etc.

In CoNLL 2000, ten types of chunks

Word	POS	Group	Word	POS	Group
<i>He</i>	PRP	B-NP	<i>to</i>	TO	B-PP
<i>reckons</i>	VBZ	B-VP	<i>only</i>	RB	B-NP
<i>the</i>	DT	B-NP	<i>£</i>	#	I-NP
<i>current</i>	JJ	I-NP	<i>1.8</i>	CD	I-NP
<i>account</i>	NN	I-NP	<i>billion</i>	CD	I-NP
<i>deficit</i>	NN	I-NP	<i>in</i>	IN	B-PP
<i>will</i>	MD	B-VP	<i>September</i>	NNP	B-NP
<i>narrow</i>	VB	I-VP	<i>.</i>	.	O

Noun groups (NP) are in red and verb groups (VP) are in blue.

# IOB Annotation for Named Entities

CoNLL 2002		CoNLL 2003			
Words	Named entities	Words	POS	Groups	Named entities
Wolff	B-PER	U.N.	NNP	I-NP	I-ORG
,	O	official	NN	I-NP	O
currently	O	Ekeus	NNP	I-NP	I-PER
a	O	heads	VBZ	I-VP	O
journalist	O	for	IN	I-PP	O
in	O	Baghdad	NNP	I-NP	I-LOC
Argentina	B-LOC	.	.	O	O
,	O				
played	O				
with	O				
Del	B-PER				
Bosque	I-PER				
in	O				
the	O				
final	O				
years	O				
of	O				
the	O				
seventies	O				
in	O				
Real	B-ORG				
Madrid	I-ORG				
.	O				

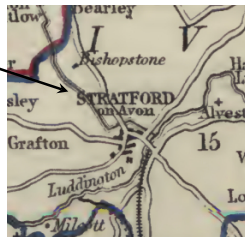
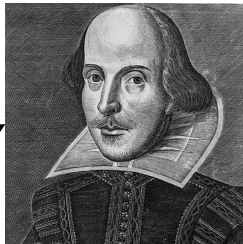
# Named Entities: Proper Nouns

William Shakespeare

was born and brought

up in

Stratford-upon-Avon





# Others Entities: Common Nouns

Meeting with  
lunchtime

our guest

on

the landing

at



# Evaluation

There are different kinds of measures to evaluate the performance of machine learning techniques, for instance:

- Precision and recall in information retrieval and natural language processing;
- The *receiver operating characteristic* (ROC) in medicine.

	Positive examples: $P$	Negative examples: $N$
Classified as $P$	True positives: $A$	False positives: $B$
Classified as $N$	False negatives: $C$	True negatives: $D$

More on the receiver operating characteristic here: [http://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](http://en.wikipedia.org/wiki/Receiver_operating_characteristic)

# Recall, Precision, and the F-Measure

The **accuracy** is  $\frac{|A \cup D|}{|P \cup N|}$ .

**Recall** measures how much relevant examples the system has classified correctly, for  $P$ :

$$\text{Recall} = \frac{|A|}{|A \cup C|}.$$

**Precision** is the accuracy of what has been returned, for  $P$ :

$$\text{Precision} = \frac{|A|}{|A \cup B|}.$$

















Recall and precision are combined into the **F-measure**, which is defined as the harmonic mean of both numbers:

$$F = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}.$$

# Evaluation

Accuracy, precision, and recall.

For noun groups with the predicted output:

Word	POS	Group	Word	POS	Group
 <i>He</i>	PRP	B-NP	 <i>to</i>	TO	B-PP
 <i>reckons</i>	VBZ	B-VP	 <i>only</i>	RB	B-NP
 <i>the</i>	DT	B-NP	 <i>£</i>	#	I-NP
 <i>current</i>	JJ	<b>B-NP</b>	 <i>1.8</i>	CD	<b>B-NP</b>
 <i>account</i>	NN	I-NP	 <i>billion</i>	CD	I-NP
 <i>deficit</i>	NN	I-NP	 <i>in</i>	IN	B-PP
 <i>will</i>	MD	B-VP	 <i>September</i>	NNP	B-NP
 <i>narrow</i>	VB	I-VP	 <i>.</i>	.	O

Accuracy =  $\frac{14}{16}$ , recall =  $\frac{2}{4} = 0.5$ , precision =  $\frac{2}{6} = 0.33$

harmonic mean =  $2 \times \frac{0.33 \times 0.5}{0.33 + 0.5} = 0.4$

# Transformers

An architecture proposed in 2018 based on the concept of **attention**  
Sometimes marketed as the ImageNet moment (See <https://ruder.io/nlp-imagenet/>)  
The similarity lies in the possibility to train semantic relations on very large corpora and memorize them in matrices  
We obtain these matrices through a **masked language model**

# Contextual Embeddings

Embeddings we have seen so far do not take the context into account  
Attention is a way to make them aware of the context.

Consider the sentence:

*I must go back to my ship and to my crew*  
*Odyssey, book I*

The word *ship* can be a verb or a noun with different meanings, but has only one GloVe embedding vector

Self-attention will enable us to compute contextual word embeddings.

# Self-Attention

In the paper *Attention is all you need*, Vaswani et al. (2017) use three kinds of vectors, queries, keys, and values. Here we will use one type corresponding to GloVe embeddings.

We compute the attention this way:

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

where  $d_k$  is the dimension of the input. The softmax function is defined as:

$$\text{softmax}(x_1, x_2, \dots, x_j, \dots, x_n) = \left( \frac{e^{x_1}}{\sum_{i=1}^n e^{x_i}}, \frac{e^{x_2}}{\sum_{i=1}^n e^{x_i}}, \dots, \frac{e^{x_j}}{\sum_{i=1}^n e^{x_i}}, \dots, \frac{e^{x_n}}{\sum_{i=1}^n e^{x_i}} \right)$$

also defined as

$$\text{softmax}(x_1, x_2, \dots, x_j, \dots, x_n) = \left( \frac{e^{-x_1}}{\sum_{i=1}^n e^{-x_i}}, \frac{e^{-x_2}}{\sum_{i=1}^n e^{-x_i}}, \dots, \frac{e^{-x_j}}{\sum_{i=1}^n e^{-x_i}}, \dots, \frac{e^{-x_n}}{\sum_{i=1}^n e^{-x_i}} \right)$$

in physics.

# The meaning of $QK$

**QK** is the dot product of the GloVe vectors. It will tell us the similarity between the words

This is analogous to cosine similarity:

	i	must	go	back	to	my	ship	and	to	my	crew
i	1.00	0.75	0.86	0.76	0.73	0.90	0.35	0.65	0.73	0.90	0.42
must	0.75	1.00	0.85	0.68	0.87	0.69	0.42	0.69	0.87	0.69	0.45
go	0.86	0.85	1.00	0.84	0.84	0.81	0.41	0.68	0.84	0.81	0.49
back	0.76	0.68	0.84	1.00	0.83	0.76	0.49	0.77	0.83	0.76	0.51
to	0.73	0.87	0.84	0.83	1.00	0.68	0.54	0.86	1.00	0.68	0.51
my	0.90	0.69	0.81	0.76	0.68	1.00	0.38	0.63	0.68	1.00	0.44
ship	0.35	0.42	0.41	0.49	0.54	0.38	1.00	0.46	0.54	0.38	0.78
and	0.65	0.69	0.68	0.77	0.86	0.63	0.46	1.00	0.86	0.63	0.49
to	0.73	0.87	0.84	0.83	1.00	0.68	0.54	0.86	1.00	0.68	0.51
my	0.90	0.69	0.81	0.76	0.68	1.00	0.38	0.63	0.68	1.00	0.44
crew	0.42	0.45	0.49	0.51	0.51	0.44	0.78	0.49	0.51	0.44	1.00



# Vaswani's attention score

The attention scores are scaled and normalized by the softmax function.

$$\text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right),$$

	i	must	go	back	to	my	ship	and	to	my	crew
i	0.36	0.05	0.07	0.05	0.04	0.19	0.01	0.02	0.04	0.19	0.01
must	0.14	0.20	0.10	0.06	0.11	0.10	0.03	0.05	0.11	0.10	0.02
go	0.18	0.09	0.14	0.09	0.08	0.13	0.02	0.04	0.08	0.13	0.02
back	0.14	0.05	0.09	0.19	0.08	0.12	0.03	0.06	0.08	0.12	0.03
to	0.11	0.11	0.09	0.09	0.15	0.08	0.04	0.07	0.15	0.08	0.03
my	0.19	0.03	0.05	0.04	0.03	0.29	0.01	0.02	0.03	0.29	0.01
ship	0.03	0.03	0.03	0.04	0.05	0.03	0.55	0.03	0.05	0.03	0.13
and	0.10	0.08	0.07	0.10	0.12	0.09	0.04	0.15	0.12	0.09	0.04
to	0.11	0.11	0.09	0.09	0.15	0.08	0.04	0.07	0.15	0.08	0.03
my	0.19	0.03	0.05	0.04	0.03	0.29	0.01	0.02	0.03	0.29	0.01
crew	0.06	0.05	0.05	0.06	0.05	0.06	0.21	0.04	0.05	0.06	0.31

# Attention

We use these scores to compute the attention.

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

For *ship*:

```
attention_ship = (0.03 * embeddings_dict['i'] +  
                  0.03 * embeddings_dict['must'] +  
                  0.03 * embeddings_dict['go'] +  
                  0.03 * embeddings_dict['back'] +  
                  0.04 * embeddings_dict['to'] +  
                  0.05 * embeddings_dict['my'] +  
                  0.55 * embeddings_dict['ship'] +  
                  0.03 * embeddings_dict['and'] +  
                  0.05 * embeddings_dict['to'] +  
                  0.03 * embeddings_dict['my'] +  
                  0.13 * embeddings_dict['crew'])
```

where the *ship* vector received 13% of its value from *crew*

# Code Example

**Experiment:** Jupyter Notebook: `4.4-attention.ipynb`  
(First part of the notebook)

# Multihead Attention

This attention is preceded by dense layers:

If  $\mathbf{X}$  represents complete input sequence (all the tokens), we have:

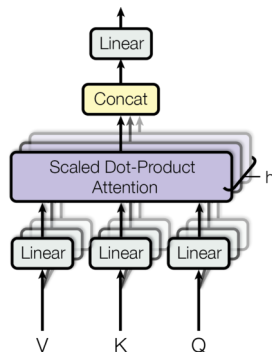
$$\mathbf{Q} = \mathbf{XW}_Q,$$

$$\mathbf{K} = \mathbf{XW}_K,$$

$$\mathbf{V} = \mathbf{XW}_V.$$

And followed by another dense layer.

In addition, most architectures have parallel attentions, where the outputs (called heads) are concatenated (multihead)



From *Attention is all you need*,  
Vaswani et al. (2017)

# Code Example

Keras has an implementation of this architecture with the `MultiHeadAttention()` layer.

**Experiment:** Jupyter Notebook: `4.4-attention.ipynb`  
(Second part of the notebook)

# Transformers

Transformers are architectures, where:

- 1 The first part of the layer is a multihead attention;
- 2 We reinject the input to the attention output in the form of an addition:

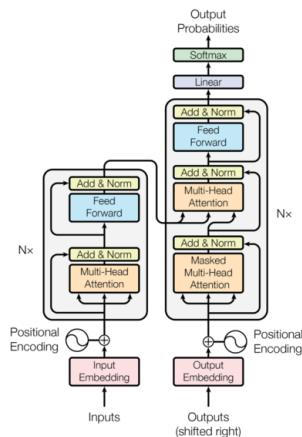
$$\mathbf{X} + \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{Q}).$$

This operation is called a skip or residual connection, which improves stability.

- 3 The result is then normalized per instance, i.e. a unique sequence, defined as:

$$x_{i,j_{norm}} = \frac{x_{i,j} - \bar{x}_{i,.}}{\sigma_{x_{i,.}}}.$$

- 4 It is followed by dense layers.



Left part, from *Attention is all you need*, Vaswani et al. (2017)

# Code Example

**Experiment:** Jupyter Notebook:

`chapter11_part03_transformer.ipynb` from Chollet's book, first part up to positional embeddings

# Positional Embeddings

BERT (the first transformer, Devlin et al. (2019)) maps each token to three embedding vectors:

- the token embedding,
- the position of the token in the sentence (positional embeddings), and
- the segment embeddings (we will skip this part).

The three kinds of embeddings are learnable vectors.

In the BERT base version, each embedding vector has 768 dimensions.

Let us consider two sentences simplified from the *Odyssey*:

*Tell me of that hero. Many cities did he visit.*

Input:	[CLS]	Tell	me	of	that	hero	[SEP]	Many	cities	did	he	visit	[SEP]
Token	$E_{[CLS]}$	$E_{\text{tell}}$	$E_{\text{me}}$	$E_{\text{of}}$	$E_{\text{that}}$	$E_{\text{hero}}$	$E_{[SEP]}$	$E_{\text{many}}$	$E_{\text{cities}}$	$E_{\text{did}}$	$E_{\text{he}}$	$E_{\text{visit}}$	$E_{[SEP]}$
Segment	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$	$E_{11}$	$E_{12}$



# Code Example

**Experiment:** Jupyter Notebook:

`chapter11_part03_transformer.ipynb` from Chollet's book, second part from positional embeddings

As a personal work and to gain a deeper understanding, you can read a tutorial in PyTorch:

<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

# Training Transformers

Transformers, such as BERT, are often trained on masked language models with two tasks:

- 1 For a sentence, predict masked words: We replace 15% of the tokens with a specific mask token and we train the model to predict them. This is just a cloze test;
- 2 For a pair of sentences, predict if the second one is the successor of the first one;

Taking the two first sentences from the *Odyssey*:

*Tell me, O Muse, of that ingenious hero who travelled far and wide after he had sacked the famous town of Troy.*

*Many cities did he visit, and many were the nations with whose manners and customs he was acquainted;*

# Masked language models

We add two special tokens: [CLS] at the start of the first sentence and [SEP] at the end of both sentences, and the token [MASK] to denote the words to predict.

We would have for the first task:

*[CLS] Tell me, O Muse, of that [MASK] hero who travelled far  
and wide [MASK] he had sacked the [MASK ] town of Troy.  
[SEP]*

For the second task, we would have as input:

*[CLS] Tell me, O Muse, of that [MASK] hero who travelled far  
and wide [MASK] he had sacked the [MASK ] town of Troy.  
[SEP] Many cities did he [MASK visit], and many were the  
[MASK nations] with whose manners [MASK and] customs he  
was acquainted; [SEP]*

where the prediction would return that the second sentence is the next one (as opposed to random sequences)

# Model size

Transformers are trained on large corpora like the colossal clean crawled corpus (<https://arxiv.org/abs/2104.08758>) and encapsulate semantics found in text in the form of numerical matrices.

This results in large models (Devlin et al., 2019):

*In this work, we denote the number of layers (i.e., Transformer blocks) as  $L$ , the hidden size as  $H$ , and the number of self-attention heads as  $A$ . We primarily report results on two model sizes:  $BERT_{BASE}$  ( $L=12$ ,  $H=768$ ,  $A=12$ , Total Parameters=110M) and  $BERT_{LARGE}$  ( $L=24$ ,  $H=1024$ ,  $A=16$ , Total Parameters=340M).*

Transformers can then act as pre-trained models for a variety of tasks.

See the list from Huggingface

Finally, an interesting reading: <https://sayakpaul.medium.com/an-interview-with-colin-raffel-research-scientist-at-google-5>