

Language Modelling

Cristina España-Bonet

UdS & DFKI, Saarbrücken, Germany

Artificial Intelligence with Deep Learning

25th April 2019

Intuition

Language Modelling is...

...anticipating coherent words.

Intuition

Language Modelling is...

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic _____

Intuition

Language Modelling is...

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic goal

Intuition

Language Modelling is...

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic goal

It's raining cats and -----

Intuition

Language Modelling is...

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic goal

It's raining cats and dogs

Intuition

Language Modelling is...

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic goal

It's raining cats and dogs

I saw falling from that tree a red -----

...anticipating coherent words.

How works your language modelling?

Messi scored a fantastic goal

It's raining cats and dogs

I saw falling from that tree a red apple

Intuition

A Language Model...

...measures fluency in sequences.

Intuition

A Language Model...

...measures fluency in sequences.

Which sequence is preferred by your language model?

I saw falling from that tree a red apple

I saw falling from that tree a clever apple

I saw falling from that tree a elephant apple

Outline

- 1 Intuition
- 2 Statistical Approach
- 3 Neural Approach
- 4 Software & References

Statistical Approach

Language Model Definition

Estimation of how probable a sequence of tokens is.

Statistical Approach

Language Model Definition

Estimation of how **probable** a sequence of tokens is.

Naïve estimation on a corpus with N sentences:

Frequentist probability

of a sentence e :

$$P(e) = \frac{N_e}{N_{\text{sentences}}}$$

Estimation of how **probable** a sequence of tokens is.

Naïve estimation on a corpus with N sentences:

Frequentist probability

of a sentence e :

$$P(e) = \frac{N_e}{N_{\text{sentences}}}$$

Problem:

- Long chains are difficult to observe in corpora.
⇒ Long sentences may have zero probability!

Statistical Approach

The n-gram approach

The language model assigns a probability $P(e)$ to a sequence of words $e \Rightarrow \{w_1, \dots, w_m\}$.

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

- The probability of a sentence is the product of the conditional probabilities of each word w_i given the previous ones
- Independence assumption: the probability of w_i is only conditioned by the n previous words

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

Statistical Approach

Example, a 4-gram model

e : All work and no play makes Jack a dull boy

$$P(e) = P(\text{All}|\phi, \phi, \phi)$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$P(e) = P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All})$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$P(e) = P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work})$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$P(e) = P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ P(\text{no}|\text{All}, \text{work}, \text{and})$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$P(e) = P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no})$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) \end{aligned}$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \end{aligned}$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ & P(\text{a}|\text{play}, \text{makes}, \text{Jack}) \end{aligned}$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ & P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \end{aligned}$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ & P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \\ & P(\text{boy}|\text{Jack}, \text{a}, \text{dull}) \end{aligned}$$

Statistical Approach

Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{aligned} P(e) = & P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work}) \\ & P(\text{no}|\text{All}, \text{work}, \text{and}) P(\text{play}|\text{work}, \text{and}, \text{no}) \\ & P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ & P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \\ & P(\text{boy}|\text{Jack}, \text{a}, \text{dull}) \end{aligned}$$

where, for each factor,

$$P(\text{and}|\phi, \text{All}, \text{work}) = \frac{N_{(\text{All work and})}}{N_{(\text{All work})}}$$

Statistical Approach

n-gram Model, Comments

- 4-gram conditions on a lot of context, so given sufficient training data these counts will be high and it will converge to the *true value*
- 4-gram many counts will be equal to zero, so we need many samples to get a good estimate

Statistical Approach

n-gram Model, Comments

- 4-gram** conditions on a lot of context, so given sufficient training data these counts will be high and it will converge to the *true value*
- 4-gram** many counts will be equal to zero, so we need many samples to get a good estimate
- 1-gram** will converge relatively quickly to their expected value, and so don't need many samples
- 1-gram** completely ignores context, and so it will converge to a less-good estimator as the number of training samples increase

Statistical Approach

Known Problems

Main problems and criticisms:

- Long-range dependencies are lost
- Still, some n -grams can be not observed in the corpus

Statistical Approach

Known Problems

Main problems and criticisms:

- Long-range dependencies are lost
- Still, some n -grams can be not observed in the corpus

Solution

Smoothing techniques:

- Linear interpolation
- Back-off models
- Discounting

Statistical Approach

Smoothing techniques

- Discounting
 - Keep part of the probability mass for unseen words

Statistical Approach

Smoothing techniques

- Discounting
 - Keep part of the probability mass for unseen words
- Back-off models
 - If n -gram is not present, go to $(n - 1)$ -gram

Statistical Approach

Smoothing techniques

- Discounting
 - Keep part of the probability mass for unseen words
- Back-off models
 - If n -gram is not present, go to $(n - 1)$ -gram
- Linear interpolation

$$P(\text{and}|\text{All, work}) = \frac{N_{(\text{All, work, and})}}{N_{(\text{All, work})}}$$

Statistical Approach

Smoothing techniques

- Discounting
 - Keep part of the probability mass for unseen words
- Back-off models
 - If n -gram is not present, go to $(n - 1)$ -gram
- Linear interpolation

$$P(\text{and}|\text{All, work}) = \lambda_3 \frac{N_{(\text{All, work, and})}}{N_{(\text{All, work})}} + \lambda_2 \frac{N_{(\text{work, and})}}{N_{(\text{work})}} + \lambda_1 \frac{N_{(\text{and})}}{N_{\text{words}}} + \lambda_0$$

Outline

- 1 Intuition
- 2 Statistical Approach
- 3 Neural Approach
- 4 Software & References

Neural Approach

First Neural Language Model [Bengio et al. 2003]

Goal

Learn with a NN the probability of sequences of words:

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1}) \quad n\text{-gram}$$

Neural Approach

First Neural Language Model [Bengio et al. 2003]

Goal

Learn with a NN the probability of sequences of words:

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1}) \quad n\text{-gram}$$

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_0, \dots, w_{i-1}) \quad \text{neural}$$

Neural Approach

First Neural Language Model [Bengio et al. 2003]

Goal

Learn with a NN the probability of sequences of words:

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1}) \quad n\text{-gram}$$

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_0, \dots, w_{i-1}) \quad \text{neural}$$

Network tasks:

- 1 learn a representation for each word (embedding)
- 2 learn the probability function from these representations

Neural Approach

First Neural Language Model [Bengio et al. 2003]

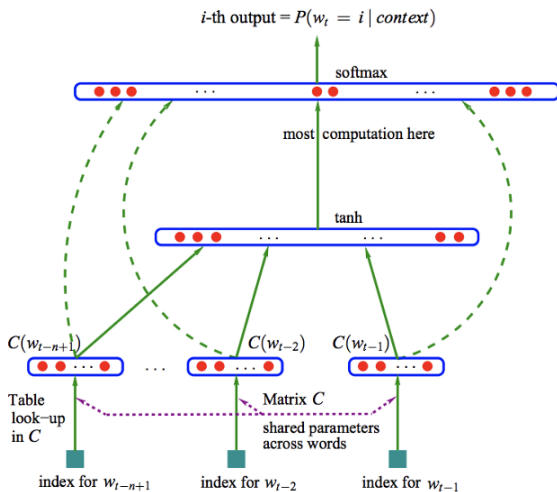


Figure 1: Neural architecture: $f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$ where g is the neural network and $C(i)$ is the i -th word feature vector.

Neural Approach

First Neural Language Model [Bengio et al. 2003]

Why it works?

- First, it does work!
Perplexity improvement wrt. n -gram
(Brown and Hansard corpora in the paper)
- Word embeddings learn semantic similarities
The dog was sleeping vs. *The cat was sleeping*

Neural Approach

First Neural Language Model [Bengio et al. 2003]

More comments

- The model scales linearly with the vocabulary
- At that time (2003) the softmax was a bottleneck

$$\text{softmax}(a_k) = \frac{e^{a_k(\mathbf{x})}}{\sum_{j=1}^T e^{a_j(\mathbf{x})}}$$

- For next slides... if we share parameters (LSTMs!) it could be sublinear in the number of words

Neural Approach

Feed Forward vs. Recurrent Language Models

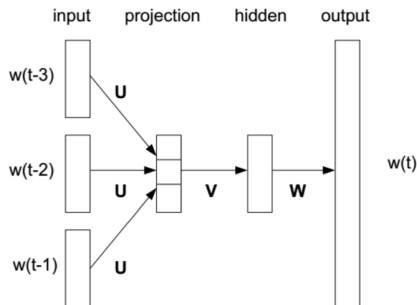
- Bengio's model is a FFN (that works great!)
- FFNs only indirectly deal with sequential data
- FF LMs are basically *soft* n -gram LMs (their history is still fixed)
- The model needs to remember a longer history, with loops
- A neural net with loops is a recurrent neural network (RNN)

Neural Approach

Recurrent Neural Language Model [Mikolov et al. 2011]

Feedforward neural network LM

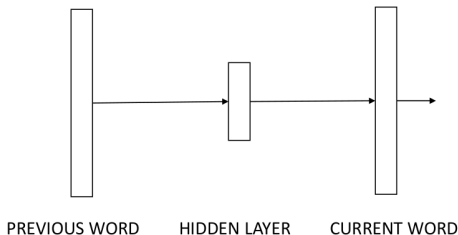
- Proposed by (Bengio et al, 2003)
- The projection layer is linear
- The hidden layer is non-linear
- Softmax at the output computes probability distribution over the whole vocabulary
- The basic model is computationally very expensive



Neural Approach

Recurrent Neural Language Model [Mikolov et al. 2011]

Towards recurrent neural network LM



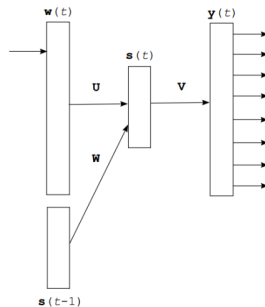
- Back to the bigram NNLM: is there a better way how to represent time than using N-1 previous words as separate inputs?

Neural Approach

Recurrent Neural Language Model [Mikolov et al. 2011]

Recurrent neural network LM

- Recurrent NNLM is about the same as the bigram NNLM
- Additional weights from the hidden layer in the previous time step
- In theory, the hidden layer can learn to represent unlimited memory



Neural Approach

Comparison of Architectures [Shi 2017]

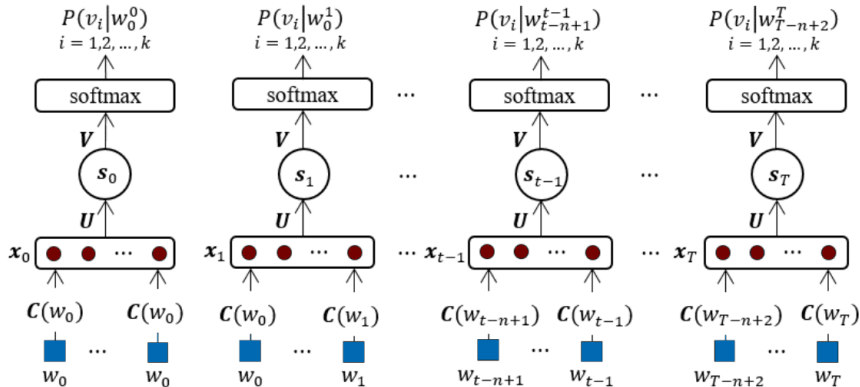


Figure 1: Feed-forward neural network language model

Neural Approach

Comparison of Architectures [Shi 2017]

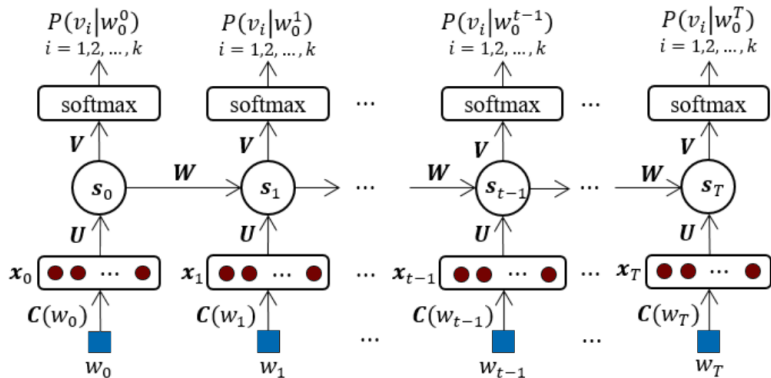


Figure 2: Recurrent neural network language model

Training of NNLMs

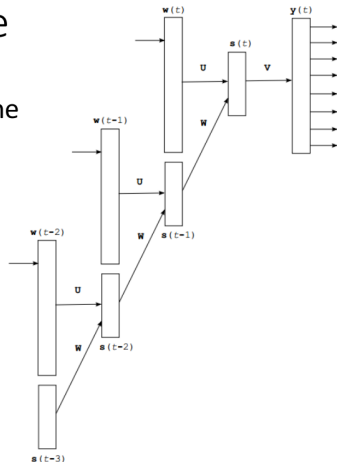
- Feedforward NNLM: the classic SGD + backpropagation
- Recurrent NNLM: the same, but the backpropagation part is more difficult to implement correctly
- The algorithm for computing gradients in RNN is called “Backpropagation through time”

Neural Approach

Recurrent Neural Language Model [Mikolov et al. 2011]

Backpropagation through time

- The intuition is that we unfold the RNN in time
- We obtain deep neural network with shared weights **U** and **W**

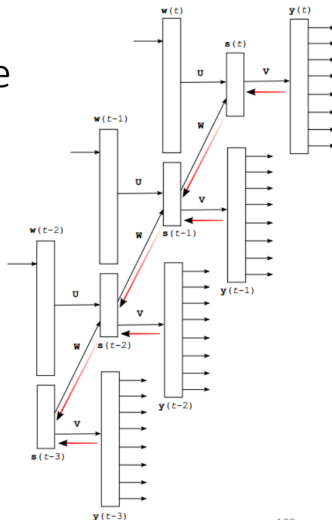


Neural Approach

Recurrent Neural Language Model [Mikolov et al. 2011]

Backpropagation through time

- We train the unfolded RNN using normal backpropagation + SGD
- In practice, we limit the number of unfolding steps to 5 – 10
- It is computationally more efficient to propagate gradients after few training examples (batch mode)



Comparison of performance on small data: Penn Treebank

- Small, standard dataset, ~1M words
- RNN outperforms FNN by about 10%

Model	Perplexity
Kneser-Ney 5-gram	141
Maxent 5-gram	142
Random forest	132
Feedforward NNLM	140
Recurrent NNLM	125

More results in: *Empirical Evaluation and Combination of Advanced Language Modeling Techniques* (Mikolov et al, 2011)

Neural Approach

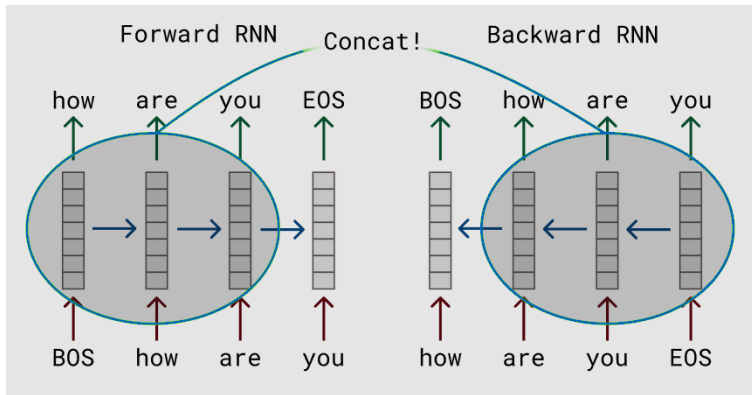
Large Scale Neural Language Model

Language models as **auxiliary task** for context aware word embeddings

Embeddings	LM Architecture	Reference
ELMo	biLSTM	[Peters et al., 2018]
OpenAI GTP	transformer decoder	[Radford et al., 2018]
Bert	biTransformer encoder	[Devlin et al., 2018]

Neural Approach

Remember from the Word Embeddings Lecture... biRNN



[https://medium.com/@plusepsilon/
the-bidirectional-language-model-1f3961d1fb27](https://medium.com/@plusepsilon/the-bidirectional-language-model-1f3961d1fb27)

Neural Approach

Stop! Acknowledgments

In the following slides **figures** are borrowed from

Lil'Log

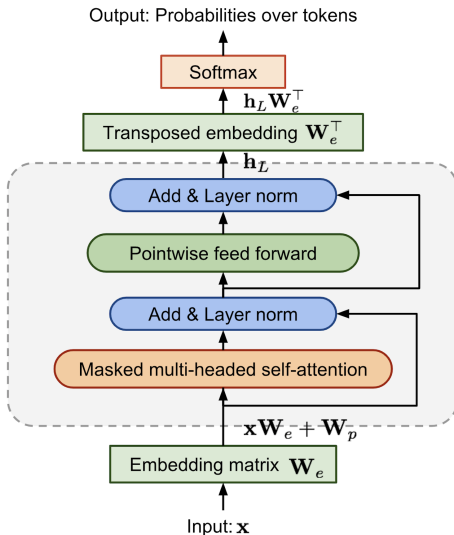
<https://lilianweng.github.io/lil-log/2019/01/31/generalized-language-models.html>

Jay Alammar

<http://jalammar.github.io/illustrated-transformer/>

Neural Approach

Transformer Decoder Language Model [Radford et al., 2018]



Transformer Block

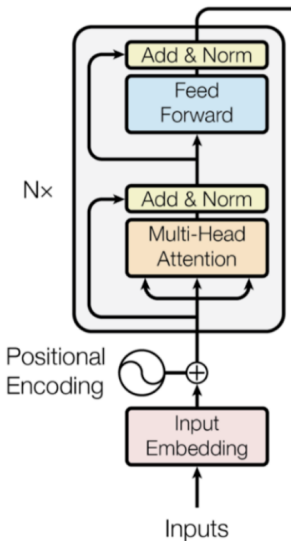
Repeat x L=12

$$\mathbf{h}_\ell = \text{transformer_block}(\mathbf{h}_{\ell-1})$$

$$\ell = 1, \dots, L$$

Neural Approach

Transformer Encoder Language Model [Devlin et al., 2018]



Neural Approach

Stop!

Have you heard about them before?

Transformer

Encoder

Decoder

Self-attention

Multi-head attention

Neural Approach

Stop!

Have you heard about them before?

Transformer

Encoder

Decoder

Self-attention

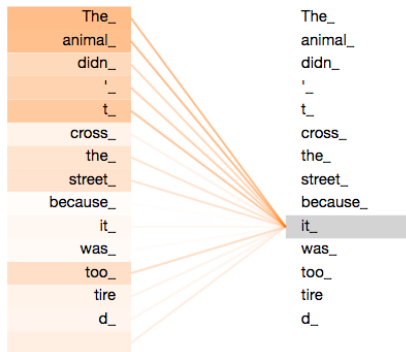
Multi-head attention

State-of-the art in **most** NLP tasks, **want details!!**

Digression: Transformer Elements

Attention vs. Self-Attention

- **Intra** endoder/decoder attention vs. **inter** encoder-decoder
- Always **alignment weights**



Digression: Transformer Elements

Self-Attention

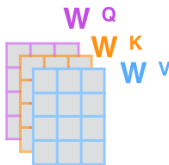
1) This is our input sentence

Thinking
Machines

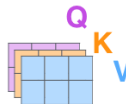
2) We embed each word



3) We multiply **X** with weight matrices



4) Calculate attention using the resulting **Q/K/V** matrices



Digression: Transformer Elements

Self-Attention

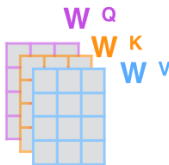
1) This is our input sentence

Thinking
Machines

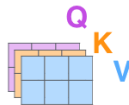
2) We embed each word



3) We multiply **X** with weight matrices



4) Calculate attention using the resulting **Q/K/V** matrices



$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \text{2x4 grid} \end{matrix} \times \begin{matrix} \text{K}^T \\ \text{4x2 grid} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \text{2x4 grid} \end{matrix} = \begin{matrix} \text{Z} \\ \text{2x4 grid} \end{matrix}$$

Digression: Transformer Elements

Self-Attention

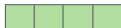
Input

Thinking

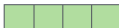
Machines

Embedding

x_1



x_2



Queries

q_1



q_2



W^Q

Keys

k_1



k_2



W^K

Values

v_1



v_2



W^V

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Thinking

x_1 

q_1 

k_1 

v_1 

Machines

x_2 

q_2 

k_2 

v_2 

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

Thinking

x_1 

q_1 

k_1 

v_1 

$$q_1 \cdot k_1 = 112$$

Machines

x_2 

q_2 

k_2 

v_2 

$$q_1 \cdot k_2 = 96$$

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

Divide by $8 (\sqrt{d_k})$

Thinking

x_1 

q_1 

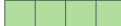
k_1 

v_1 

$q_1 \cdot k_1 = 112$

14

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

12

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Thinking

x_1 

q_1 

k_1 

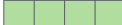
v_1 

$$q_1 \cdot k_1 = 112$$

14

0.88

Machines

x_2 

q_2 

k_2 

v_2 

$$q_1 \cdot k_2 = 96$$

12

0.12

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

Divide by $8 (\sqrt{d_k})$

Softmax

Softmax

X

Value

Thinking

x_1 

q_1 

k_1 

v_1 

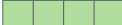
$q_1 \cdot k_1 = 112$

14

0.88

v_1 

Machines

x_2 

q_2 

k_2 

v_2 

$q_2 \cdot k_2 = 96$

12

0.12

v_2 

Digression: Transformer Elements

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax
X
Value

Sum

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

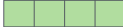
14

0.88

v_1 

z_1 

Machines

x_2 

q_2 

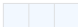
k_2 

v_2 

$q_1 \cdot k_2 = 96$

12

0.12

v_2 

z_2 

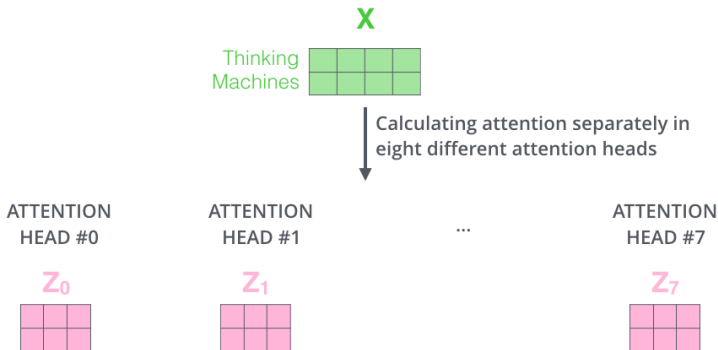
Digression: Transformer Elements

Multi-head Attention

- Self-attention with multiple heads
- Each head (weights W^Q , W^K and W^V) is initialised independently
- Project inputs into different representation subspaces
- Transformer models [Vaswani et al., 2017] use 8 heads

Digression: Transformer Elements

Multi-head Attention



Digression: Transformer Elements

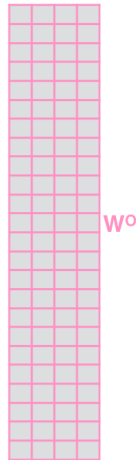
Multi-head Attention

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Digression: Transformer Elements

Multi-head Attention

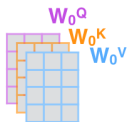
1) This is our input sentence

Thinking
Machines

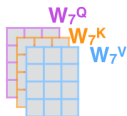
2) We embed each word



3) Split into 8 heads.
We multiply X with weight matrices



...



4) Calculate attention using the resulting $Q/K/V$ matrices



...



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



...



W^O

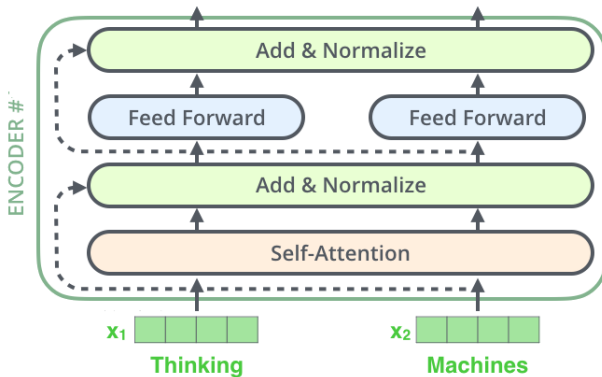


Z



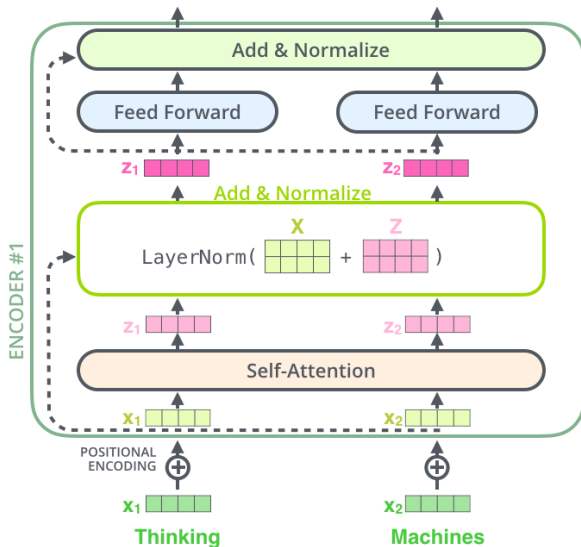
Digression: Transformer Elements

An Encoder Layer



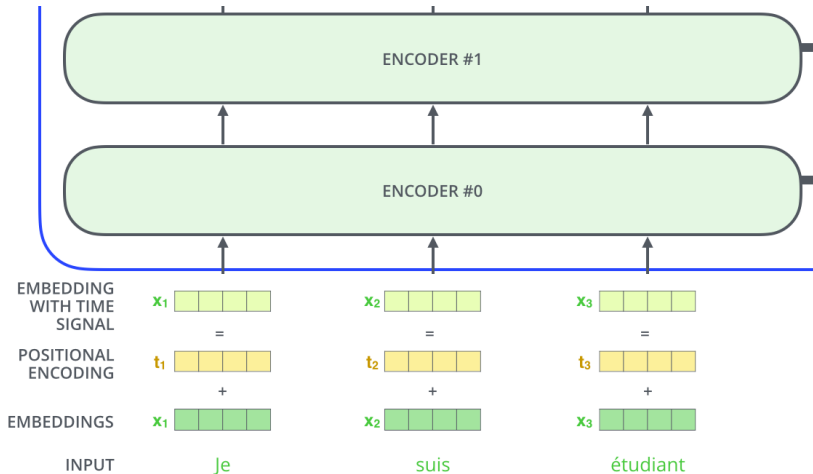
Digression: Transformer Elements

An Encoder Layer: Residuals



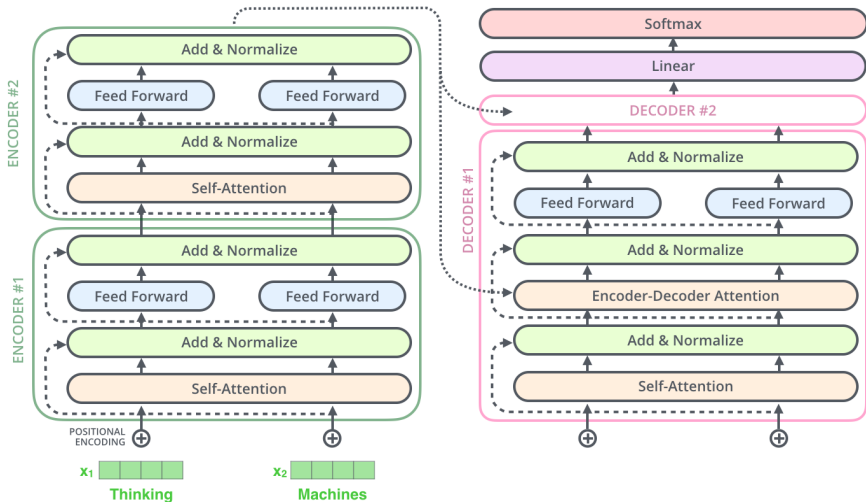
Digression: Transformer Elements

Positional Embeddings



Digression: Transformer Elements

Toy Transformer (more on the NMT lecture!)



Neural Approaches

Stop!

Transformer

Encoder

Decoder

Self-attention

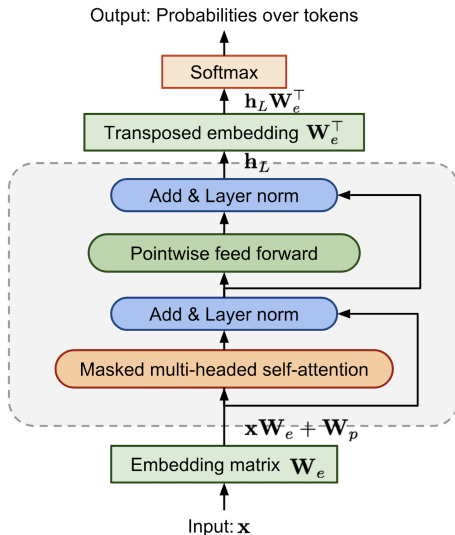
Multi-head attention

State-of-the art in **most** NLP tasks, **want details!!**

Got them!!

Neural Approaches

Transformer Decoder Language Model [Radford et al., 2018]



Transformer Block
Repeat x L=12

$$\mathbf{h}_\ell = \text{transformer_block}(\mathbf{h}_{\ell-1})$$
$$\ell = 1, \dots, L$$

OpenAI GTP(-2)

- Large transformer-based language model with 1.5 billion parameters
- Trained on a dataset of 8 million web pages (40 GB)
- Objective: predict the next word, given all of the previous words within some text

<https://openai.com/blog/better-language-models/>

Neural Approaches

Transformer Decoder Language Model [Radford et al., 2018]

SYSTEM PROMPT
(HUMAN-WRITTEN)

Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

MODEL
COMPLETION
(MACHINE-
WRITTEN, SECOND
TRY)

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back

The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.

She was carrying a pair of black and white striped gloves and a small black bag.

Neural Approaches

After Basic FFNs, RNNs and Transformers LMs...

Things we should cover (with more time!)

- LSTM cells
- Character-level LMs
- Bidirectional RNNs (ELMo)
- Combination of models
- ...

Software & References

- 1 Intuition
- 2 Statistical Approach
- 3 Neural Approach
- 4 Software & References

Software & References

Libraries & Packages

- **KenLM**

<https://kheafield.com/code/kenlm/>

- **SRILM**

<http://www.speech.sri.com/projects/srilm/download.html>

- **RNNLM**

<https://github.com/mspandit/rnnlm>

- **CSLM**

<https://github.com/hschwenk/cslm-toolkit/blob/master/README>

Software & References

References Used

- Reinhard Kneser and Hermann Ney. **Improved Backing-off for N-gram Language Modeling**. Proceedings ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing. Pages 181–184, 1995.
- Daniel Jurafsky and James H. Martin. 2018. **Speech and Language Processing**. Chapter 3: N-gram Language Models.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent and Christian Jauvin. **A Neural Probabilistic Language Model**. Journal of Machine Learning Research, volume 3, pages 1137–1155, 2003.

Software & References

References Used

- Tomas Mikolov, Anoop Deoras, Stefan Kombrink, Lukas Burget, Jan "Honza" Cernocky. **Empirical Evaluation and Combination of Advanced Language Modeling Techniques**. 12th Annual Conference of the International Speech Communication Association, pages 605–608, 2011.
- Tomas Mikolov. **Statistical Language Models Based on Neural Networks**. PhD. Thesis. 2012.
- Dengliang Shi. **A Study on Neural Network Language Modeling**. *eprint arXiv:1708.07252*

Everything fine?

Questions?

Language Modelling

Cristina España-Bonet

UdS & DFKI, Saarbrücken, Germany

Artificial Intelligence with Deep Learning

25th April 2019

Extra Slides

The language model $P(e)$ examples

Examples

`http://files.asimihsan.com/courses/nlp-coursera-2013/
notes/nlp.html`

Example 1: 3-gram language model with linear interpolation

- Consider the weights

$$\lambda_1 = \lambda_2 = \lambda_3 = 1/3; \lambda_0 = 0$$

- Consider the corpus
 - the green book STOP
 - my blue book STOP
 - his green house STOP
 - book STOP

Extra Slides

The language model $P(e)$ examples

Example 1: 3-gram language model with linear interpolation

$P(\text{book}|\text{the, green})?$

$P(\text{book}|\text{the, green}) = 0.571$

Extra Slides

The language model $P(e)$ examples

Example 1: 3-gram language model with linear interpolation

$P(\text{book}|\text{the, green})?$

$$P(\text{book}|\text{the, green}) = 0.571$$

$$\begin{aligned} & \frac{1}{3}P(\text{book}|\text{the, green}) + \frac{1}{3}P(\text{book}|\text{green}) + \frac{1}{3}P(\text{book}) = \\ & \frac{1}{3} \frac{\text{Count}(\text{the,green,book})}{\text{Count}(\text{the,green})} + \frac{1}{3} \frac{\text{Count}(\text{green,book})}{\text{Count}(\text{green})} + \frac{1}{3} \frac{\text{Count}(\text{book})}{\text{Count}()} = \\ & \frac{1}{3} \times \frac{1}{1} + \frac{1}{3} \times \frac{1}{2} + \frac{1}{3} \times \frac{3}{14} \end{aligned}$$

Extra Slides

The language model $P(e)$ examples

Example 2: Discounting

x	Count(x)	$P(w_i w_{i-1})$	Count*(x)	$P^*(w_i w_{i-1})$
the	48			
the, dog	15	15/48	14.5	14.5/48
the, woman	11	11/48	10.5	10.5/48
the, man	10	10/48	9.5	9.5/48
the, park	5	5/48	4.5	4.5/48
the, job	2	2/48	1.5	1.5/48
the, telescope	1	1/48	0.5	0.5/48
the, manual	1	1/48	0.5	0.5/48
the, afternoon	1	1/48	0.5	0.5/48
the, country	1	1/48	0.5	0.5/48
the, street	1	1/48	0.5	0.5/48

Extra Slides

The language model $P(e)$ examples

Example 2: Discounting

Discounted counts:

$$\text{Count}^*(x) = \text{Count}(x) - \text{constant}$$

Left-over probability mass:

$$\alpha(w_i - 1) = 1 - \sum_w \frac{\text{Count}^*(w_{i-1}, w)}{\text{Count}(w_{i-1})}$$

$$\text{In the example: } \alpha(\text{the}) = 10 \times \frac{0.5}{48} = \frac{5}{48}$$