

 translate


deeplearning.ai

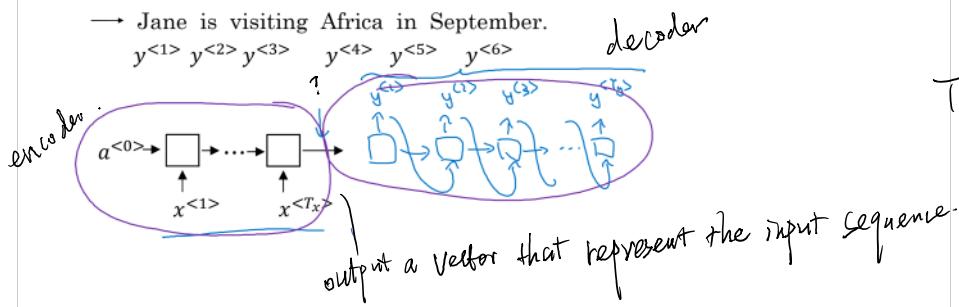
Sequence to sequence models

Basic models

Sequence to sequence model

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
 Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.



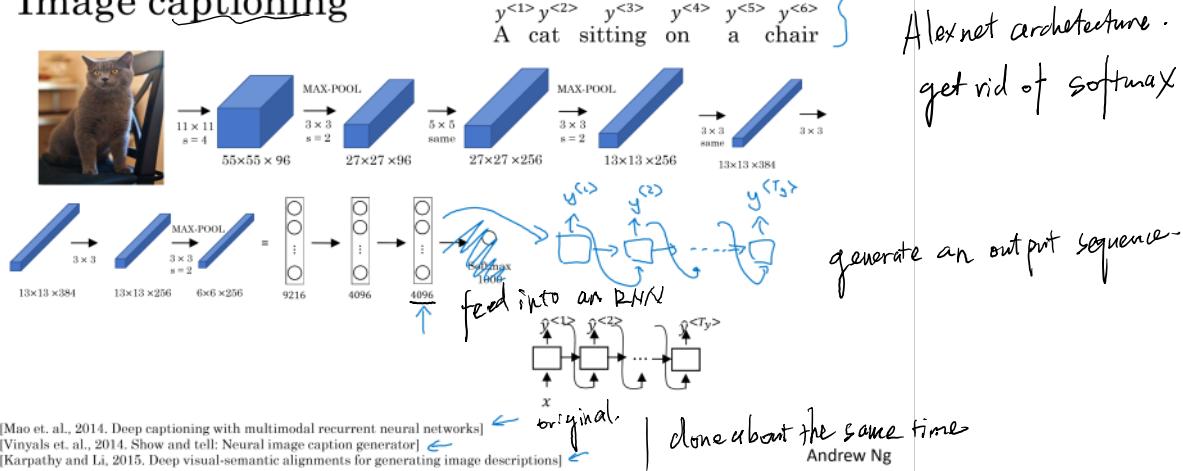
[Sutskever et al., 2014. Sequence to sequence learning with neural networks] ↩

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation] ↩

Andrew Ng

Another task this method can do.

Image captioning



There are some differences. Don't want a randomly chosen. Want Most likely. Andrew Ng
Done about the same time. (2012) · (2014).

picked

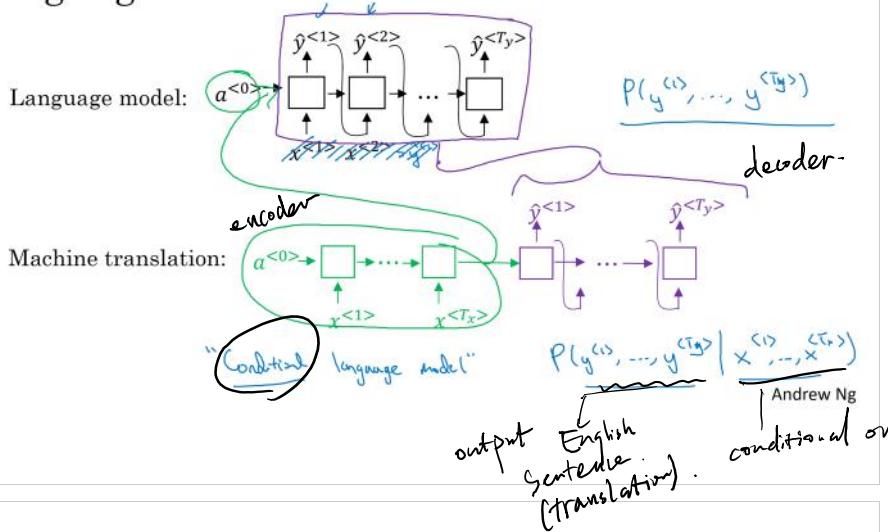


deeplearning.ai

Sequence to sequence models

Picking the most likely sentence

Machine translation as building a conditional language model



the decoder network look identical to the language model.

Differences start from the output of the encoder network.

Finding the most likely translation

Jane visite l'Afrique en septembre.

$$P(y^{(1)}, \dots, y^{(T_y)} | x)$$

Do not want random sampling.

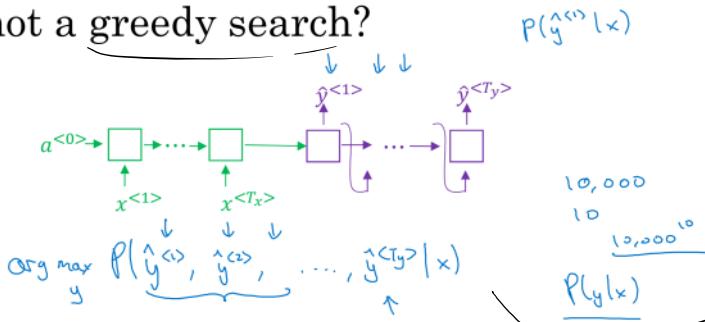
- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\arg \max_{y^{(1)}, \dots, y^{(T_y)}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

Andrew Ng

~~greedy search may end up with suboptimal~~ ~~greedy~~ top

Why not a greedy search?



greedy search - search for the most likely word.

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.

$$P(\text{Jane is going } | x) > P(\text{Jane is visit } | x)$$

goal: pick the words to maximize the joint probability.

→ Jane is going to be visiting Africa in September.
 $P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$

Andrew Ng

~~locally optimal. Not globally max.~~

The common thing to do
use approximate algorithm

The total # of combination is
exponentially large \Rightarrow can't calculate all
1000 words in dict. (or word long
sentence)
 1000^{10} sentences

beam search

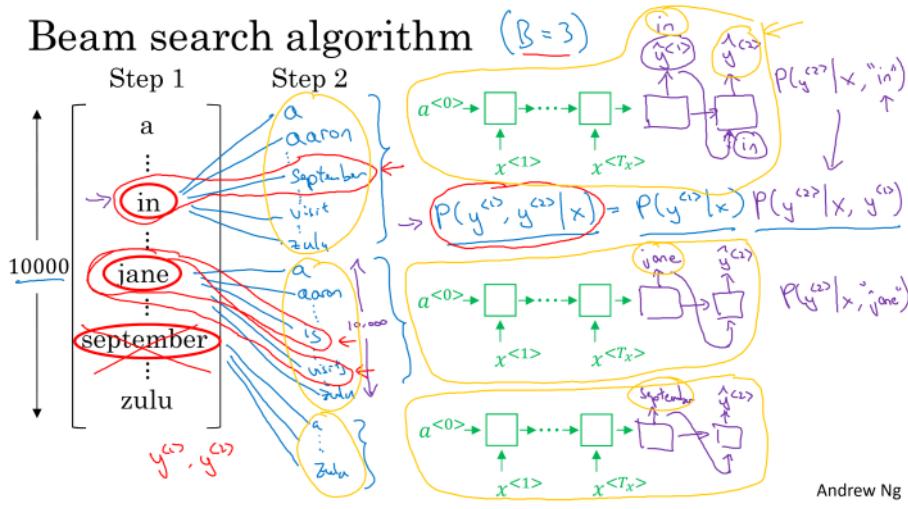
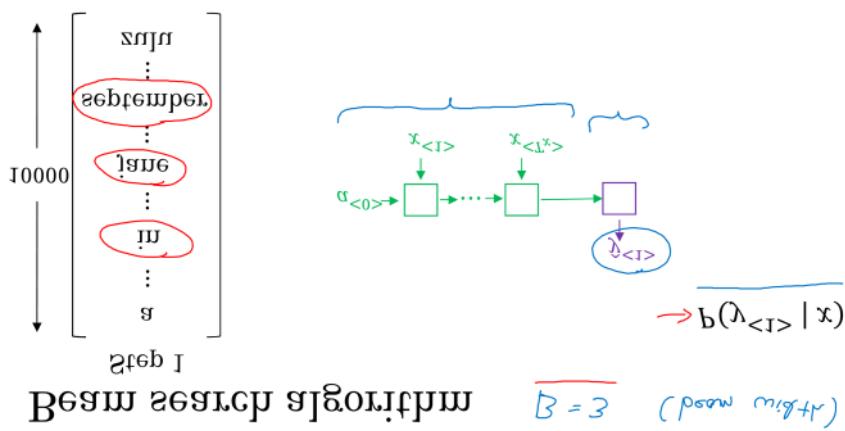


deeplearning.ai

Sequence to sequence models

Beam search

Beam width consider the 3 most likely possibilities



Keep a memory of the top 3.

$B=3$
 $B \times 10^{1000}$ possibilities
 \Rightarrow pick top B .
 Why pick top 3... path 2 is very good.
 (Is it big in the beginning?)

→ Beam width = 1. greedy search

 refine beam



deeplearning.ai

Sequence to
sequence models

Refinements to
beam search

Length normalization

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Andrew Ng

$$\log \left[\arg \max_y \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}) \right]$$

$\alpha = 0.7$

$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$

normalize by length

middle

full normalization

no normalization

$T_y = 1, 2, \dots, 30.$

$\alpha = 1$

$\alpha = 0$

$\log p_{\text{ly}}(x) \leftarrow$

$p_{\text{ly}}(x) \leftarrow$

~~\log , $\beta \rightarrow$ numerical underflow~~

If not normalize by length \Rightarrow reward
 short translation \Rightarrow not good.
 \rightarrow normalize.

(larger B , larger computational power required.)

\int *for research, sometimes used.*
very large for production system.

Beam search: faster, not guaranteed to find exact maximum.



deeplearning.ai

Sequence to sequence models

Error analysis on beam search

Example

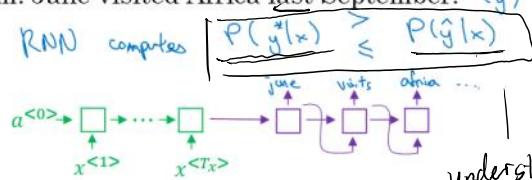
Jane visite l'Afrique en septembre.

→ RNN
→ Beam Search
B↑

Decide which part is to blame
— Always tempting to ↑B. But it may not help.

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y}) ←



bad translation, changed meaning.

Use error analysis to decide which part to focus on fixing.

understand this and know which part is problematic

Andrew Ng

Error analysis on beam search

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y})

Case 1: $P(y^*|x) > P(\hat{y}|x)$ ←

Beam search chose \hat{y} . But y^* attains higher $P(y|x)$.

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \leq P(\hat{y}|x)$ ←

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) < P(\hat{y}|x)$.

Conclusion: RNN model is at fault.

$$P(y^*|x)$$

$$P(\hat{y}|x)$$

$$\arg \max_y P(y|x)$$

highest prob $\frac{1}{2} \times 10^{-10}$, \hat{y} .

$(\hat{y}, \hat{y}, \hat{y})$ highest probability. \Rightarrow Beam search is at fault.

Andrew Ng

highest probability $\frac{1}{2} \times 10^{-10}$.
 $\Rightarrow P(y^*, y^*, y^*)$

Error analysis process

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	2×10^{-10}	1×10^{-10}	B
...	...	—	—	R
...	...	—	—	Q
				R
				R
				...

Figures out what fraction of errors are "due to" beam search vs. RNN model

Andrew Ng

Take a few examples to analyze.



blue score

Challenge: equally good translation may exist

Try to combine N-best reference
~~by up and score~~ \rightarrow $\frac{1}{N}$ $\sum_{i=1}^N \text{score}_i$ \rightarrow $\frac{1}{N}$ $\sum_{i=1}^N \log \text{score}_i$?

Sequence to sequence models



Blue score

* Evaluate how much your predicted



deeplearning.ai

Bleu score (optional)

- * Evaluate how much your predicted output is ~~the same~~ similar to one or more references
- Multiple equally good reference
- The Bleu score serve as human evaluation
- Ask human to evaluate.

Evaluating machine translation

French: Le chat est sur le tapis.

~~2/3~~ Bleu
bilingual evaluation understudy
understudy

Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: the the the the the the the.

Precision:

~~7/7~~

Modified precision:

~~a unique word~~ get credit at most twice.

one way, look at each word, see it appears in the references.
⇒ precision

A reader-friendly paper Ng recommended

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

Bleu score on bigrams

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

	Count	Count clip
the cat	2 ←	1 ←
cat the	1 ←	0
cat on	1 ←	1 ←
on the	1 ←	1 ←
the mat	1 ←	1 ←

$$\frac{4}{6}$$

~~7/17~~ ↑
~~17~~ ↑
4 bigram
↑
13 13 2 ↑

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

→ MT output: The cat the cat on the mat. (g)

$$p_1 = \frac{\sum_{\text{unigram} \in \text{g}} \text{count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram} \in \text{g}} \text{count}(\text{unigram})}$$

↑
unigram
Unigram g

$$p_n = \frac{\sum_{n\text{-gram} \in \text{g}} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram} \in \text{g}} \text{count}(n\text{-gram})}$$

↑
n-gram
n-gram g

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

$p_1, p_2 = 1.0$ → Sometimes possible even if the g not exact the same with any of the reference

total # of ngram

Bleu details

p_n = Bleu score on n-grams only

Combined Bleu score: $\text{BP} \exp\left(\frac{1}{4} \sum_{n=1}^4 p_n\right)$

$\text{BP} = \text{brevity penalty}$

$$\text{BP} = \begin{cases} 1 & \text{if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & \text{otherwise} \end{cases}$$

p_1, p_2, p_3, p_4

brevity penalty.

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

short words → high score
because higher chance all terms are ~~the~~ in the reference

⇒ So use BP to adjust for that
(penalize shorter translation)

find details here in the paper

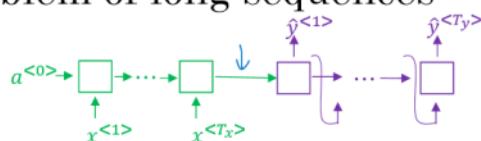


deeplearning.ai

Sequence to sequence models

Attention model intuition

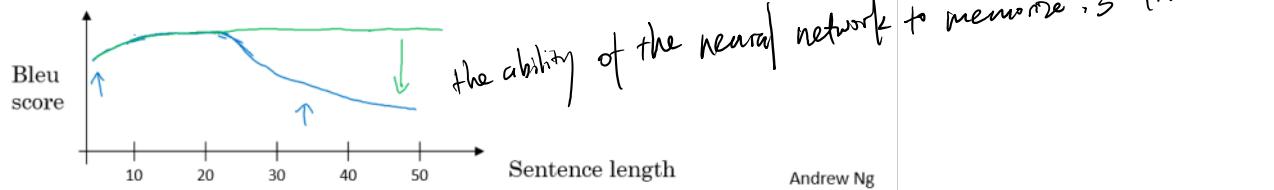
The problem of long sequences



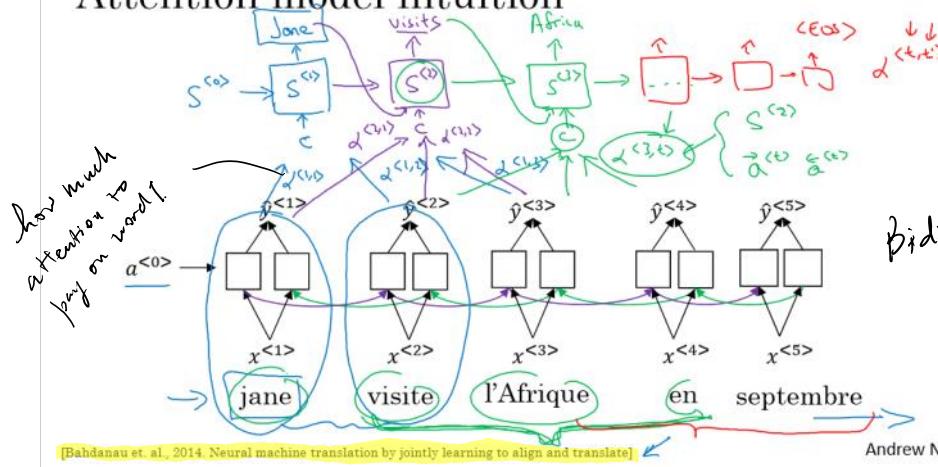
Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.

human translator will read the sentence part by part
to try to translate



Attention model intuition



[Bahdanau et al., 2014. Neural machine translation by jointly learning to align and translate.]



context. Depends on
attention. next few
words.

$a^{(t,t')}$ attention weight.

Bidirectional RNN

for each step of generation
only focus on a few words

instead of attention weights we
pay attention to $a^{(1,1)}$.

attention model 2



deeplearning.ai

Sequence to sequence models

Attention model

Attention model

How much attention should the 1st word in the English sentence

pay on all $1, 2, \dots, t'$ words in the French sentence

[Bahdanau et al., 2014. Neural machine translation by jointly learning to align and translate]

$\alpha^{(t,t')} = \text{amount of "attention" } y^{(t)} \text{ should pay to } a^{(t')}$

$$c^{(t)} = \sum_{t'=1}^T \alpha^{(t,t')} a^{(t')}$$

$$\bar{a}^{(t)} = (\vec{a}^{(t)}, \overleftarrow{a}^{(t)})$$

$$\sum_{t'=1}^T \alpha^{(t,t')} = 1$$

$$c^{(t)} = \sum_{t'=1}^T \alpha^{(t,t')} a^{(t)}$$

$$c^{(t)} = \sum_{t'=1}^T \alpha^{(t,t')} a^{(t)}$$

$$\bar{a}^{(t)} = (\vec{a}^{(t)}, \overleftarrow{a}^{(t)})$$

$\vec{a}^{(t)}$ a feature vector for time step t' of the French sentence (Input)

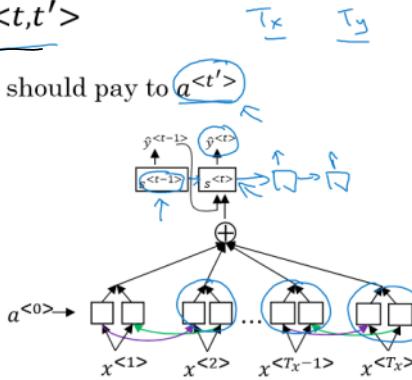
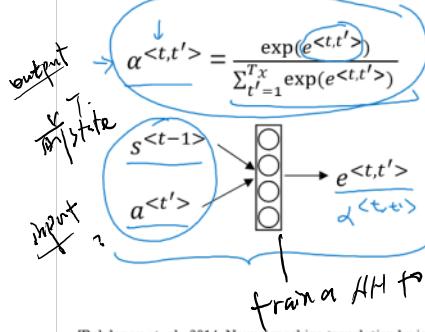
$\overleftarrow{a}^{(t)}$ weight \times attention weights $\alpha^{(t,t')}$

\vec{a} forward occurrence, \overleftarrow{a} backward occurrence

[Andrew Ng]

Computing attention $\alpha^{(t,t')}$

$\alpha^{(t,t')} = \text{amount of attention } y^{(t)} \text{ should pay to } a^{(t')}$



[Bahdanau et al., 2014. Neural machine translation by jointly learning to align and translate]
[Xu et al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng

Downside; take quadratic time

to run algorithm
words words.

$T_x \times T_y$. cost

But input & output aren't that long
so maybe acceptable

[Applied to other problem:
image captioning]

pay attention to part of a picture
when writing a caption

Normalize data

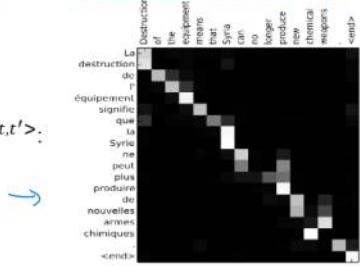
July 20th 1969 → 1969 - 07 - 20

Attention examples

July 20th 1969 → 1969 - 07 - 20

23 April, 1564 → 1564 - 04 - 23

Visualization of $\alpha^{(t,t')}$:



→ Visualization of
the attention weights

high — white in the figure



plus
produire
de
nouvelles
armes
chimiques
<end>



high — white in the figure

Andrew Ng