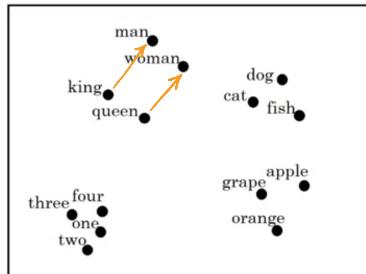


$$\begin{aligned} \ell_{\text{man}} - \ell_{\text{woman}} &\approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\ \ell_{\text{king}} - \ell_{\text{queen}} &\approx \begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

THE MAIN DIFFERENCE BETWEEN MAN/WOMAN AND KING/QUEEN IS THEIR GENDER FEATURE.
 $\ell_{\text{man}} - \ell_{\text{woman}} \approx \ell_{\text{king}} - \ell_{\text{queen}}$



* You want to find the word (w) that maximizes the similarity between (ℓ_w) and $\ell_{\text{king}} - \ell_{\text{man}} + \ell_{\text{woman}}$.

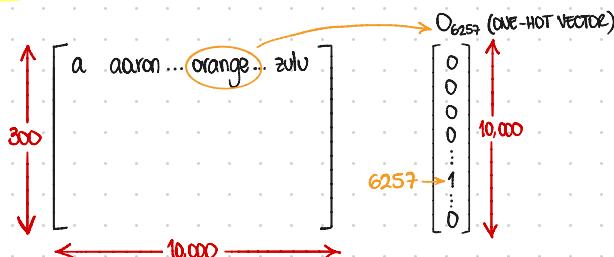
$$\arg \max_w \text{Sim}(\ell_w, \ell_{\text{king}} - \ell_{\text{man}} + \ell_{\text{woman}})$$

* You can also find the similarity between two vectors with the Cosine Similarity:

IF u AND v ARE VERY SIMILAR,
THEN THE INNER PRODUCT
(NUMERATOR) WILL BE VERY LARGE

$$\text{Sim}(u, v) = \frac{u \cdot v}{\|u\| \|v\|}$$

EMBEDDING MATRIX



$$E \cdot O_{6257} = e_{6257}$$

$(300, 10k)$ $(10k, 1)$ $(300, 1)$

$$E \cdot e_j = e_j$$

EMBEDDING FOR WORD j

LEARNING WORD EMBEDDINGS

I want a glass of orange _____.

I e_{4343} \longrightarrow E \longrightarrow e_{4343}

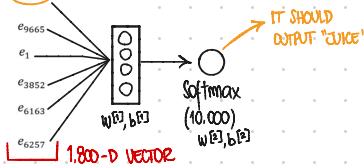
want e_{9665} \longrightarrow E \longrightarrow e_{9665}

a e_1 \longrightarrow E \longrightarrow e_1

glass e_{3852} \longrightarrow E \longrightarrow e_{3852}

of e_{6163} \longrightarrow E \longrightarrow e_{6163}

orange e_{6257} \longrightarrow E \longrightarrow e_{6257}



I want a glass of orange juice to go along with my cereal.

CONTEXT

TARGET

CONTEXT

* you can select ± contexts to feed your NN to predict the target

- last 4 words
- 4 words on the left and right
- word before the target
- nearby one word (skip gram) 

WORD2VEC MODEL

* rather than having the context be always the previous 4 words or the following 4 words, you randomly pick a word to be the context word. Then you pick randomly another word within some window (± 5 or ± 10 words from the context word) and that becomes the target.

I want a glass of orange juice to go along with my cereal.

T CONTEXT T

T T

* you want to learn the mapping from some word c to some target t

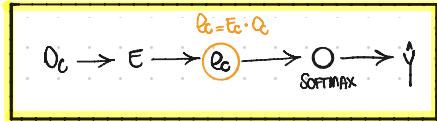
context c ("orange") \longrightarrow target t ("juice")

6257

4834

X \longrightarrow Y

Skip-Gram Model



$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{1000} e^{\theta_j^T e_c}}$$

θ_t = parameter associated with output t .

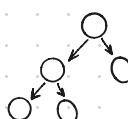
LOSS FUNCTION

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{1000} y_i \log \hat{y}_i$$

SOFTMAX

* disadvantages of skip-gram model:

- computational speed: for the softmax function, you need to carry out a sum over all 10,000 words in the vocabulary. SOLUTION?



→ HIERARCHICAL SOFTMAX CLASSIFIER:

- you run the softmax on subsets of the vocabulary until you narrow it down to the target

* how to choose/sample context c ? \rightarrow Sample uniformly at random from your training corpus.

NEGATIVE SAMPLING

* you sample a context and target word and associate the pair with a label of 1 (they go well together) or a label of 0 (they don't go well together).

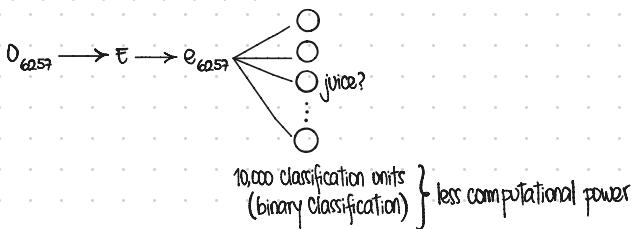
<u>CONTEXT</u>	<u>WORD</u>	<u>TARGET</u>
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

K= 5-20
SMALL SETS
K= 2-5
LARGE SETS

* then you create a supervised learning problem where you input the training pairs along with their labels.

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

we use ① for logistic regression



* how do you sample the negative examples?

- Sample according to the empirical frequency of words in your corpus.
DISADVANTAGE: you end up with an overrepresentation of words like "the", "of", "a", etc.
- Sample negative examples uniformly at random. **DISADVANTAGE:** not very representative of the distributions of English words.
- THIS WORKS BEST!

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

Something in between the previous two examples

GLOVE WORD VECTORS

* global vectors for word representation

$$x_{ij} = \# \text{times } \textcircled{j} \text{ appears in the context } \textcircled{i}$$

\downarrow
 t

\downarrow
 c

$$\text{Min} \left[\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2 \right]$$

WEIGHTING TERM
 $f(x_{ij}) = 0 \text{ IF } x_{ij} = 0$

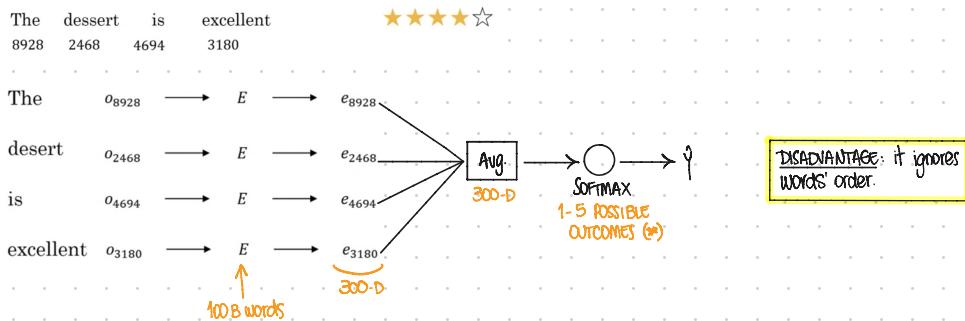
SHOULD BE RANDOMLY INITIALIZED
AT THE BEGINNING OF TRAINING

SENTIMENT CLASSIFICATION

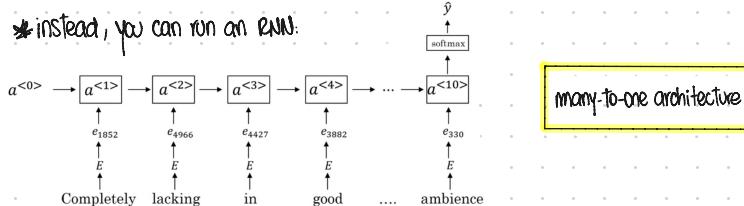
x	y
The dessert is excellent.	★★★★★☆
Service was quite slow.	★★☆☆☆☆
Good for a quick meal, but nothing special.	★★★★☆☆
Completely lacking in good taste, good service, and good ambience.	★☆☆☆☆☆

* you may not have a big training set (10,000 - 100,000 words)

* Simple sentiment classification problem:



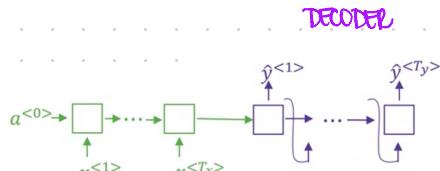
* instead, you can run an RNN:



SEQUENCE TO SEQUENCE MODELS

- * mostly used for machine translation

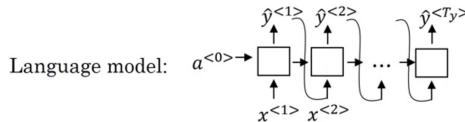
$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
 Jane visite l'Afrique en septembre
 → Jane is visiting Africa in September.
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$



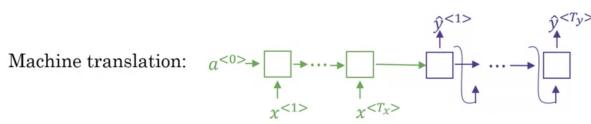
- * it can also be used for image captioning

PICKING THE MOST LIKELY SENTENCE

- * machine translation → building a conditional language model



* you estimate the probability of a sentence:
 $P(y^{<1>} \dots, y^{<Ty>})$



* "conditional language model": you estimate the probability of an English translation given the French sentence provided as input.

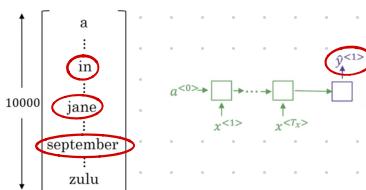
$$P(y^{<1>} \dots, y^{<Ty>} | x^{<1>} \dots, x^{<Tx>})$$

- * for machine translation you don't want to sample at random. You want to find the sentence that maximizes the conditional probability.

BEAM SEARCH

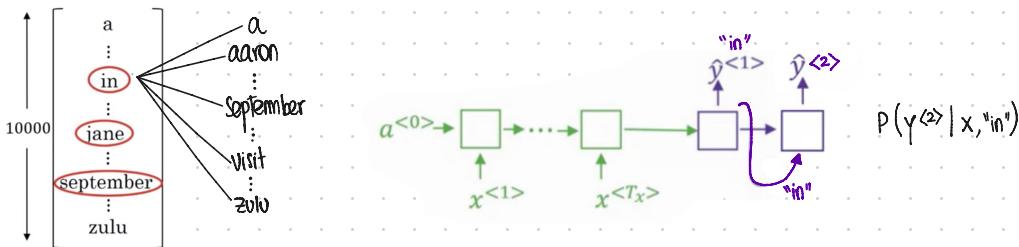
- * steps for the beam search algorithm

STEP ①



* you run the input French sentence through the encoder network and the first step will decode the network. Then, it will output overall 10,000 possibilities. It will output "in", "Jane", and "September" as the three most likely choice of the first word.

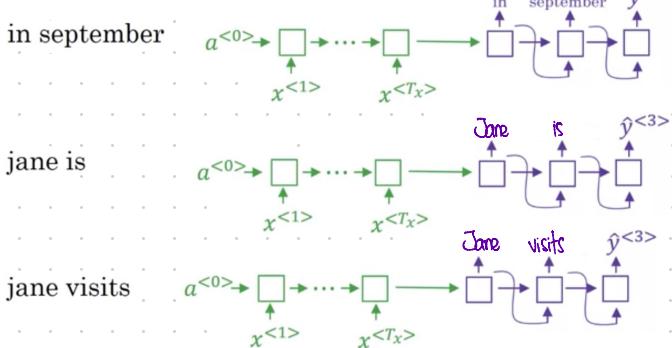
Step II



* in general, we want to find $P(y^{(2)} | x, y^{(1)}) = P(y^{(2)} | x) * P(y^{(2)} | x, y^{(1)})$

Step III

* you found that the most likely pairs are:



BEAM SEARCH REFINEMENT

#1. length normalization:

- $\arg \max \prod_{t=1}^T P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$ → BECAUSE SOME PROBABILITIES MAY BE ALMOST 0, WE IMPLEMENT \log FOR THE ALGORITHM TO WORK BETTER.
- $\arg \max \sum_{t=1}^T \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$ ← ONE OTHER THING THAT MAY IMPROVE THE ALGORITHM IS TO NORMALIZE IT BY THE NUMBER OF WORDS IN THE TRANSLATION. IT ALSO REDUCES THE PENALTY FOR OUTPUTTING LONG TRANSLATIONS.
- $\frac{1}{T} \sum_{t=1}^T \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$ ←

#2 how to set B?

- large B: better results, but slower
 - small B: worse results, but faster
- } commonly 10

ERROR ANALYSIS ON BEAM SEARCH

Jane visite l'Afrique en septembre.

Human: Jane visits Africa in September.

Algorithm: Jane visited Africa last September.

} we try to see why this \oplus in translation took place. Was it the RNN or the Beam search?

* in this case, increasing ② might not influence at all.

#1. CASE 1 $P(y^*|x) > P(\hat{y}|x)$

Beam search chose \hat{y} , but y^* attains higher $P(y|x)$.
Beam search is at fault!!

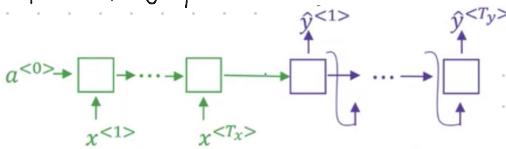
#2. CASE 2 $P(y^*|x) \leq P(\hat{y}|x)$

y^* is a better translation than \hat{y} , but the RNN predicted $P(y^*|x) < P(\hat{y}|x)$.
RNN model is at fault !!

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}	BEAM!
		1×10^{-10}	2×10^{-10}	RNN

ATTENTION MODEL INTUITION

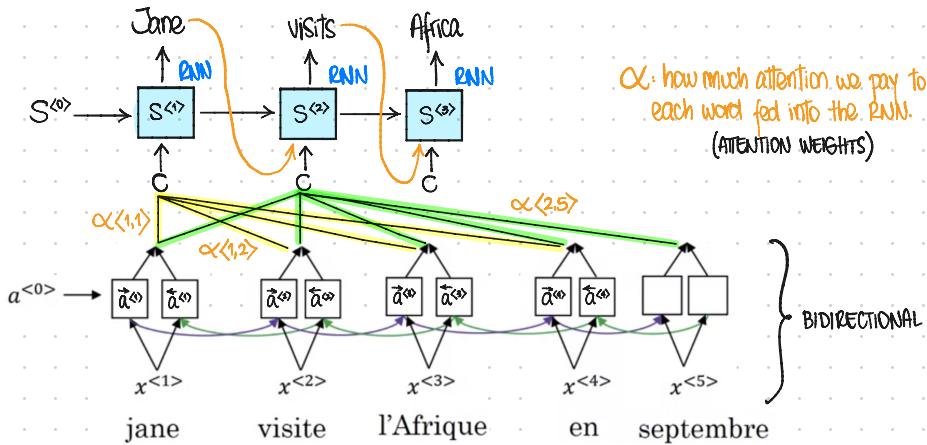
* it helps solve the problem of long sequences



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.

* it translates sentences more like humans do: a part of the sentence at a time.

* it pays attention to part of the sentence while generating a translation



Q: how much attention we pay to each word fed into the RNN.
(ATTENTION WEIGHTS)

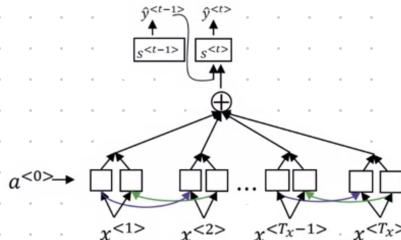
* the attention weights satisfy the condition $\sum_t \alpha^{(t,t)} = 1$

* the context element C satisfies the condition $C = \sum_t \alpha^{(t,t)} a^{(t)}$ $\rightarrow a^{(t)} = (\vec{a}^{(t)}, \bar{a}^{(t)})$

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

amount of "attention" $y^{(t)}$
should pay to $a^{(t)}$

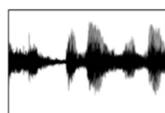
$$s^{(t-1)} \rightarrow e^{(t,t')}$$



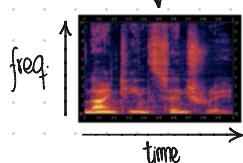
* the model runs in quadratic cost

* it also works for image captioning

SPEECH RECOGNITION

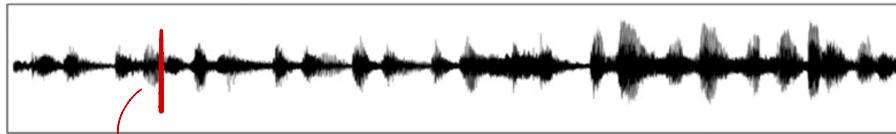
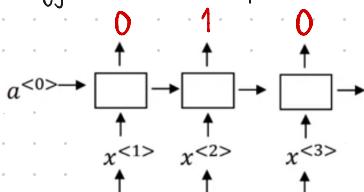


* reasonable size of data: 3,000 - 30,000 hours



SPECTROGRAM \rightarrow colors: amount of energy

* trigger word detection system:



If your trigger word is identified there,
your labels for this unit will be 1 and the rest 0.

* downsize → it creates an imbalanced dataset with more 0 than 1.