# Sequence Models:

⊕ Recurrent Neural Networks (RNN)

\* Seq. Models:

| | | |
|---|---|---|
| Speech recognition | $X \longrightarrow y$ | input $X$ is a sequence |
| Music generation | $X \longrightarrow y$. | |

DNA sequence analysis    AGCCTCCAGCCTCC ...

Machine Translation

Video activity recognition

Name entity recognition...

\* Notations

EX:   $x$ :   Harry Potter & Hermoine grange invented a new spell.
$\underbrace{\qquad}_{x^{<1>} \;\; x^{<2>} \;\; x^{<3>}}$  $x^{<t>}$ ---- $x^{<9>}$ $T_x = 9$

$y$ :    1    1    0    1    1    0    0    0    0

(persons name) $y^{<1>}$  $y^{<2>}$  ----  $y^{<t>}$ ..... $y^{<9>}$  $T_y = 9$

$x^{(i)<t>} \rightarrow t^{th}$ elment of $i^{th}$ Training ex

$T_x^{(i)} \rightarrow$ length of i/p seq in $i^{th}$ Training ex.

$x^{<9>} = a$.

Vocabulary



$\begin{bmatrix} a \\ aaron \\ \vdots \\ and \\ \vdots \\ harry \\ \vdots \\ potter \\ \vdots \\ zulu \end{bmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ 367 \\ \vdots \\ 4075 \\ \vdots \\ 6830 \\ \vdots \\ 10,000 \end{matrix}$

10,000 words :

$x^{<1>}$  $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0. \end{bmatrix} \rightarrow 4075.$   . . .   $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \\ 1 \\ \\ 0 \end{bmatrix}$

one hot vectors.

use one-hot representation

<unk> : for words not in vocabulary  <unknown token>.

\* Recurrent Neural N/w (RNN) model :

• why not a standard NN?
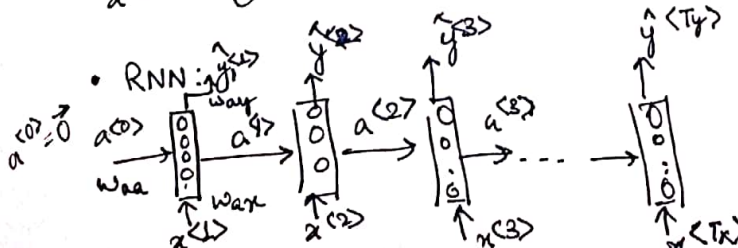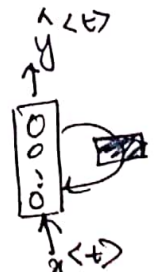
$x^{<1>}$  ○
$x^{<2>}$  ○
$\vdots$   $\vdots$
$x^{<T_x>}$  ○

problems :
- i/p and o/p can be different lengths in diff examples
- Doesn't share features learned across diff pos^n of text.
- very large i/p layer.

• RNN :



$T_x = T_y$.    other notation :

- Scans left to right

-only earlier layers are used to predict : problematic as first few letters cant be names.

Ex: can confuse b/w Teddy & Teddy Roosevelt    sol$^n$ : bidirectional RNN.

object      name .

$$a^{<1>} = g_1(w_{aa} a^{<0>} + w_{ax} x^{<1>} + b_a) \leftarrow \text{tanh/Relu}$$

$$\hat{y}^{<1>} = g_2(w_{ya} a^{<1>} + b_y) \leftarrow \text{sigmoid}$$

$$\boxed{\begin{array}{l} a^{<t>} = g(w_{aa} a^{<t-1>} + w_{ax} x^{<t>} + b_a) \\ \hat{y}^{<t>} = g(w_{ya} a^{<t>} + b_y) \end{array}} \Rightarrow$$

$$a^{<t>} = g(w_a [a^{<t-1>}, x^{<t>}] + b_a)$$

Ex: $w_{aa}$ $\underset{100,100}{}$    $a^{<t-1>}$ $\underset{100}{}$    $w_{ax}$ $\underset{(100,10000)}{}$

$$w_a = [\underset{\overleftarrow{100}}{w_{aa}} \vdots \underset{\overleftarrow{10000}}{w_{ax}}] \updownarrow 100$$

$$\overleftarrow{10100}$$

$$\hat{y}^{<t>} = g(w_y a^{<t>} + b_y)$$

Thus,

$$\boxed{\begin{array}{l} a^{<t>} = g(w_a [a^{<t-1>}, x^{<t>}] + b_a) \\ \hat{y}^{<t>} = g(w_{ya} a^{<t>} + b_y) \end{array}}$$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

100

10000 $\Big\}$ 10100

$\hat{y}^{<T_y>}$   $w_y, b_y$ given to all

**\* Backpropagation through time :**



$w_a, b_a$

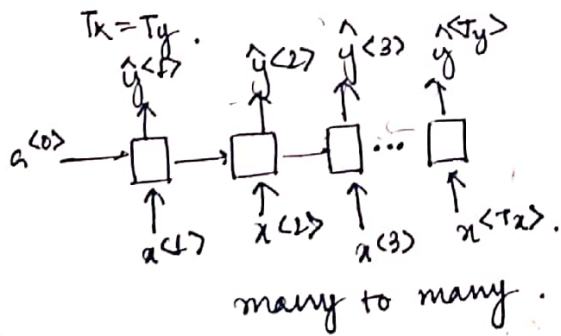$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<t>}) \log(1 - \hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>}) \leftarrow$$

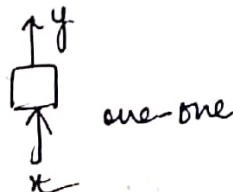Here, it is Backpropagation through time.

**\* Different types of RNN :**

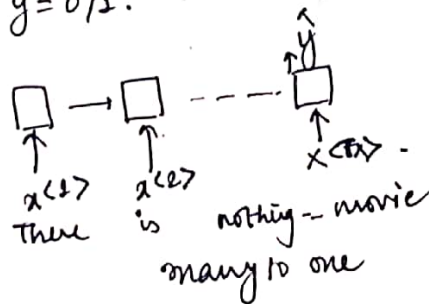$T_x$ and $T_y$ might not always be same.

• Examples of RNN architectures :

$T_x = T_y$.



$a^{<0>}$

many to many.

Ex:



one-one

EX: sentiment classification :
X = text
y = 0/1.



x$^{<1>}$  x$^{<2>}$
There is nothing - movie

many to one

∴ Music generation:

$\rightarrow y^{<1>} y^{<2>} \dots y^{<Ty>}$

$x = \phi$

$x^{<t>}$ many-one architecture



. Machine Translation.



encoder / decoder.

$\hat{y}^{<1>}$ ... $\hat{y}^{<Ty>}$

---

## * Language model and sequence generation :-

• Language modelling :

Speech recog.

- The apple & pair salad —(1)     $P(1) = 3.2 \times 10^{-13}$
- The apple & pear salad. —(11)    $P(2) = 5.7 \times 10^{-10}$

$P(\text{sentence}) = ?$     $y^{<1>}, y^{<2>}, \dots, y^{<Ty>}$

Training set : large corpus of eng text.

Cats average 15 hrs each day <EOS> Tokenise this + 1 extra token <EOS>
$y^{<1>}$  $y^{<2>}$,  ———   $y^{<7>}$ $y^{<9>}$    $x^{<t>} = y^{<t-1>}$    end of sentence.

The Egyptian Mau is beard of cat <EOS>
P(any word or unk)         <UNK>

Training :



$x^{<9>} = y^{<8>}$
day.

$L(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$

$$L = \sum_t L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

## * Sampling novel sequences :
P(sec word, given first word is the)
$P(\_ | \text{the})$

Sampling



<EOS>

much longer sequences.

$\rightarrow P(a) \, P(aaron) \, P(cat) \dots P(zulu) = np.random.choice.$

character level RNN:
Vocabulary $[a, b, \dots, z, 0, \dots, 9, ., , , '', ;, !, . — . A \dots Z]$

has own adv & disadv.

# *Vanishing gradients with RNNs :

The [cat] which already ate ......... , [was] full

The [cats] — — — — ..... [were] full.

derivative can ↑se or
↑se exponentially
due to very deep RNN.

vanishing grads are problematic. can cause numerical overflow
→ sol$^n$: gradient clipping

# * gated Recurrent Units (GRU)
— makes RNN better at capturing long range connections.

papers : cho et al, 2014 : on properties of
neural M/c translat-
Encoder - decoder
approach.



chung-et-al, 2014 : Empirical Eval.
of gated RNN on
seq. modelling

• GRU (simplified)



$c$ = memory cell

$$c^{<t>} = a^{<t>}$$

$$\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1-\Gamma_u) * c^{<t-1>}$$
(elementwise multipl$^n$)

$c^{<t>} = \frac{1}{0} \begin{matrix}(singular)\\(plural)\end{matrix}$ ----- = 1

The [cat], which --------- [was] full

$\Gamma_u \overset{u}{=} 0$   $\Gamma_u = 0$ — $\Gamma_u = 0$
(whende        memorizes
you update).    cat is
                similar.

* → multiplic$^y$/ mat convoluti

* → multiplic$^y$ mat convoluti
(as same dim$^n$)

# Full GRU:

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c).$$

$\left\{\begin{matrix}\tilde{h}\\u\\r\\h\end{matrix}\right.$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_c)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1-\Gamma_u) + c^{<t-1>}.$$

↳ alternate notations.

# *long short Term Memory : (LSTM)
— more powerful than GRU.

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c]$$

upd. gate $\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$

forget gate $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$ — variation →

o/p gate $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_0)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * c^{<t>}$$
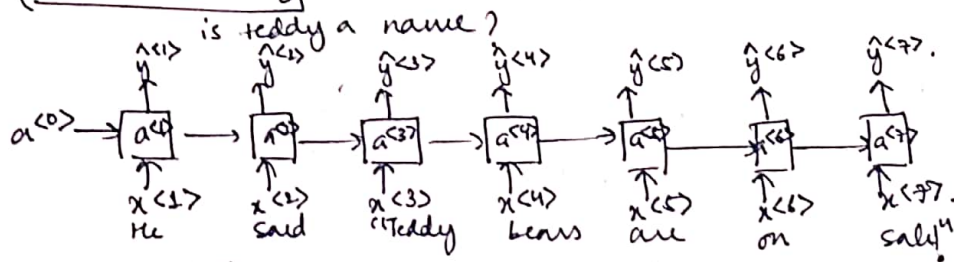
Paper : Hoschreiter &
Schmidhuber
1997.
LSTM.

$c^{<t-1>}$ replaces $x^{<t>}$.
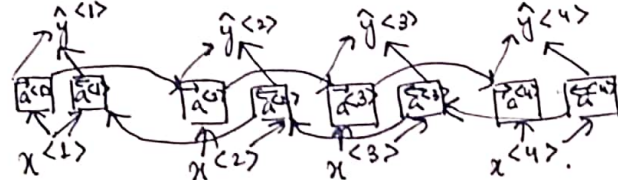- peephole connection

# * Bidirectional RNN

- Takes info from later & earlier inputs.

Ex:

He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great Prez!"

is teddy a name?



$$\hat{y}^{<t>} = g(W_y[\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t>}] + b_y)$$

• BRNN.

also 4 RU/ LSTM.



# * Deep RNNs :-

EX:

multiple RNNs stacked together.



$$a^{[2]<3>} = g(W_a^{[2]}[a^{[2]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

# ⊕ Natural language processing & word Embeddings :-

## # Word Embeddings.

## * Word Representation :

• $V = [a, aaron, ....., zulu, <UNK>]$    $|V| = 10,000$

1-hot representation

weakness :
- Treats the word as object into itself, and doesn't allow algo to easily generalise across words.

Featurised representation: word embedding

| | Man (5391) | Woman (9853) | King (4014 | Queen) (7157) | Apple (456) | Orange (6257) | |
|---|---|---|---|---|---|---|---|
| gender | −1 | 1 | −0.95 | 0.97 | 0.00 | 0.01 | paper: van der Maaten |
| Royal | 0.01 | 0.02 | 0.93 | 0.95 | −0.01 | 0.00 | and Hinton, |
| Age | 0.03 | 0.02 | 0.7 | 0.69 | 0.03 | −0.02 | 2008. |
| food. | 0.04 | 0.01 | 0.02 | 0.07 | 0.95 | 0.97. | Visualise data usin |
| : | | | | | | | t-SNE |

300

$e_{9853}$

↳ 300 D vector depresents man

$e_{5391}$

t-SNE. → objects are embeddings.

similar objects grouped together.

## * using word embeddings :

Ex: Named entity recog $\underline{\text{ex}}$ :



Sally   Johnson   is   an   orange farmer   } one hot encoding fails

Robert   Lin   is   a   apple farmer / durian cultivator   } word embedding succeeds.

steps : Transfer lear & word embeddings.   to learn with word embedding

1. learn word embeddings from a large text corpus (1-100B words)
   (or download pre-trained embedding online) (300D dense feature vector)
2. Transfer embedding to new task with smaller Tr.S.
3. optional : Continue to fine tune word embeddings with new data.

— word embeddings generalise algorithms much better.

## * Properties of word embeddings :

Analogies :

Man → woman , king → ?

$e_{man}$   $e_{woman}$
(embedding vectors).

$$ \Rightarrow \quad e_{man} - e_{woman} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} $$

$$ e_{king} - e_{queen} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} $$

Paper: Mikolov et. al. 2013;
linguistic Regularities in continuous space word representation

Analogies using word vectors :



300D

$e_{man} - e_{woman} \approx e_{king} - e_?$

Find word w : $\arg\max\limits_{w}$ similarity $(e_w, e_{king} - e_{man} + e_{woman})$

tSNE: 300D → 2D (parallelogram relationships may hold true (as above))

• Cosine similarity :

$sim (e_w, e_{king} - e_{man} + e_{woman})$

$$ sim(u,v) = \frac{u^T v}{\|u\|_2 \|v\|_2} $$



$\cos\phi = 1 \quad \phi = 0°$   similar
$\cos\phi = 0 \quad \phi = 90°$   different.

## * embedding matrix :

a   aaron   ... orange ...   zulu   <unk>

$\to$ 6257.



300D ↓   ← 10,000 →

$0_{6257}$ = $\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \to 6257 \\ \vdots \\ 0 \end{bmatrix}$ 10,000

$E \cdot O_j = e_j$ = embedding for word

$E_{(300, 10K)} \cdot O_{10K, 1} = \begin{bmatrix} 0 \\ \vdots \\ x \\ \vdots \\ 1 \end{bmatrix}$ 300,1

$= e_{6257}$

# Learning word embeddings : word2vec & Glove

## * Learning word embeddings :

I want a glass of orange

$O_{4343} \rightarrow E \rightarrow e_{4343}$
$O_{9665} \rightarrow E \rightarrow e_{9665}$
$O_1 \rightarrow E \rightarrow e_1$
$O_{3852} \rightarrow E \rightarrow e_{3852}$
$O_{6163} \rightarrow E \rightarrow e_{6163}$
$O_{6257} \rightarrow E \rightarrow e_{6257}$

$w^{[2]}, b^{[2]}$

→ softmax (10,000)

↑ $w^{[1]}, b^{[1]}$

1800 (300 × 6 words) stacked

paper : Bengio 2003
A neural probabilistic language model.

* g...

Other context/target pairs :

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

[a glass of orange] = context
[juice] = target
[to go along with my] = context

context : last 4 words   or   last 1 word [orange] ?   or   Nearby one word [skip-gram] → glass ?

## * Word2Vec

• Skip grams (allows to create supervise learning task)

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

| Context | Target |
|---------|--------|
| orange | Juice . |
| orange | glass |
| orange | my . |

paper: Mikolov 2013
Efficient estimation
of word representation
in vector space

# ...

### Model

Vocab = 10,000 words
context   c ("orange") → Target t   (juice)
$\quad$ 6257 $\qquad$ 4834

$O_c \rightarrow E \rightarrow e_c \rightarrow O \xrightarrow{softmax} \hat{y}$

$\qquad$ 10000

$L(\hat{y}, y) = -\sum_{i=1}^{10,000} y_i \log \hat{y}_i$

softmax : $p(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$

$\theta \rightarrow$ param associated with o/p t.

$y \rightarrow$ one hot vector

### Problems :

→ computational speed.

sol^n : Hierarchical softmax

5k tree may be developed as non symmetrical, common words on top, whereas non common on lower level.

### How to sample context c ?

the, of, a, and, to, ...
orange, apple, durian ...

paper : Mikolov 2013,
Distributed representation
of words & phrases &
their composition

## * Negative sampling :

| X |  | Y |
|---|---|---|
| Context | word | target ? |
| orange | juice | 1 ← +ve ex |
| orange | king | 0 |
| orange | book | 0 ← zero (-ve) example |
| orange | the | 0 |
| orange | of | 0 |

↑ c   ↑ t   ↑ y

k = 5-20 for smaller DSets
k = 2-5 for larger Data sets.

$$p(y=1|c,t) = \cancel{\sigma(\theta_t^T e_c)} \; \sigma(\theta_t^T e_c) \leftarrow$$

instead of 10000 D softmax, treat as 10000 binary classification problem

* **glove word vectors**: (global vectors for word representation)

  `I want a glass of orange juice to go along with my cereal` ———— paper: pennington, 2014 global vectors for word representation

  $c, t$.

  $X_{ij} = $ # times $i$ appears in context of $j$
  
  $c \nearrow \quad \nwarrow j \qquad \uparrow t \qquad\qquad \uparrow c.$

  maybe: $X_{ij} = X_{ji}$ (symmetric case)

  model:
  minimize $\displaystyle\sum_{i=1}^{10000}\sum_{j=1}^{10000} f(X_{ij})(\theta_i^T e_j + b_i - b_j' - \log X_{ij})^2$  $\quad 0 \log 0^u = 0$

  $\qquad\qquad\qquad\qquad \theta_t^{\,t} \quad e_c^{\,c}$  
  $\qquad\qquad\qquad\qquad \theta_t^{\,t} e_c^{\,u}$

  $\theta_i, e_j$ are symmetric.

  weighting term
  $f(X_{ij}) = 0$ if $X_{ij} = 0$.  $\qquad e_w^{(final)} = \dfrac{e_w + \theta_w}{2}$

. A note on on the featurisation view of word embeddings.

# Applications using word embeddings:

\*<u>Sentiment</u> classification :-



$x \longrightarrow y$

| $x$ | | $y$ |
|---|---|---|
| The dessert is excellent | $\longrightarrow$ | 4 stars . |
| service was slow | $\longrightarrow$ | 2 stars |
| Good for a quick meal, nothing special | $\longrightarrow$ | 3 stars |
| Completely lacking in good state, good ... | $\longrightarrow$ | 1 star. |

· Simple sentiment class. model :

The    desert    is    excellent $\longrightarrow$ 4 stars.
8928    2468    4694    3180
 ↓        ↓       ↓       ↓
 E        E       E       E
 ↓        ↓       ↓       ↓
$e_{8928}$  $e_{2468}$  $e_{4694}$  $e_{3180}$.

$\hookrightarrow$ [ Avg. | 300 D ]

↓ softmax 1-5
ŷ

· RNN for sentiment classification (many to one)  $a^{<0>}$

Completely $\rightarrow$ E $\rightarrow$ $e_{1852}$ $\rightarrow$ $a^{<1>}$
lacking $\rightarrow$ E $\rightarrow$ $e_{4966}$ $\rightarrow$ $a^{<2>}$
in $\rightarrow$ E $\rightarrow$ $e_{4427}$ $\rightarrow$ $a^{<3>}$
good $\rightarrow$ E $\rightarrow$ $e_{3852}$ $\rightarrow$
... 
ambiance. $\rightarrow$ E $\rightarrow$ $e_{830}$. $\cdots\rightarrow$ $a^{<10>}$ $\rightarrow$ [ softmax ] $\rightarrow$ ŷ

$X \longrightarrow Y \rightarrow$ can be used to monitor comments.

problem : if "good" appears many times in bad reviews, then classifier may think this is a good review

* Debiasing word embeddings:
   gender bias, etc...

Man : Woman as King : Queen

Man : Computer Programmer as Woman : Homemaker?  ✗
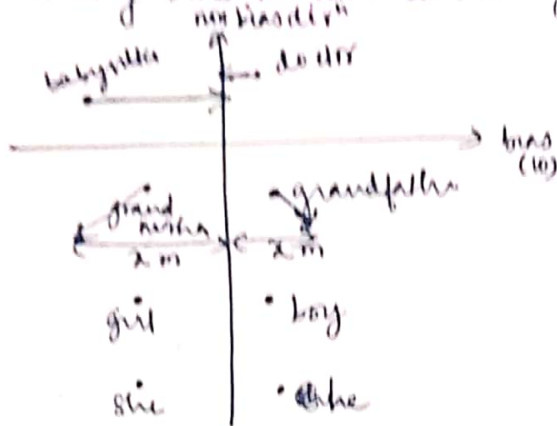
Father : Doctor as Mother : Nurse  ✗
                 ↳ gender bias          ↳ gender bias

— Word embeddings can ~~be used~~ reflect gender, ethnicity, age, sexual orientation, ..

• Addressing bias in word embeddings:
   non bias dir"

1. Identify bias direction
   { he - she
   { male - female.

         average

2. Neutralize : For every word that not
            definitional, project to get rid
                                  ↗ bias
            Ex: grandfather
                 grandmother
                 father
                 mother.

3. Equalization of pairs
      grandmother ── grandfather . ( Ex: equidistant from
            girl    ──    boy                          babysitter)

⊕ Sequence Models & Attention Mechanisms:

# Various sequence to sequence architecture:

* Basic Models:
- Translation Jane visite l'afrique en septembre
         → Jane is visiting Africa in september



        Encoding N/w                    Decoding N/w.

$y^{(1)}$  $y^{(1)}$  $y^{(2)}$  $y^{(3)}$  $y^{(5)}$  $y^{(6)}$
A cat sitting on a chair.

• Img captioning:



img    CNN    → RNN
       VGG     (Decoder N/w)
       FC vectors
       at end

* Picking the most likely sentence:
   • M/c translation as building a conditional language model.
      language model

      $a^{(0)}$ →

   M/c translation:

      $a^{(0)}$ →



                encoder-decoder

      a conditional
      lang. model

$$P(y^{<1>}, \ldots, y^{<T_y>} \mid x^{<1>}, \ldots, x^{<T_x>})$$

Ex: $P(y^{<1>}, y^{<2>}, \ldots y^{<T_y>} \mid x)$

↳ french

$\underbrace{\phantom{P(y^{<1>}, y^{<2>},}}$ English possibility.

∴ $\underset{y^{<1>}, \ldots, y^{<T_y>}}{\arg\max} P(y^{<1>}, \ldots, y^{<T_y>} \mid x)$

↳ maximised by Beam search.

• why not greedy search?

Pick first words, then pick subsequently most likely words.

$P(\hat{y}^{<1>} \mid x)$

Ex: Jane is ___ going to visit ___
            going to be visiting.

* Beam search Algorithm :  ...

step 1.



10,000 ↕
a
⋮
(in)
(jane)
(september)
⋮
zulu

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square$
        $\uparrow$           $\uparrow$
       $x^{<1>}$          $x^{<T_x>}$

$\hat{y}^{<1>}$

$B = 3$ (beam width)

→ P(①) $P(y^{<1>} \mid x)$

step 2 → evaluate 2nd word Aft first word.

[encode] →

$(in) \; \hat{y}^{<2>}$
$\hat{y}^{<1>}$

$P(y^{<2>} \mid x = "in")$

∴ $P(y^{<1>}, y^{<2>} \mid x)$
$= P(y^{<1>} \mid x)$
$P(y^{<2>} \mid x, y^{<1>})$

similarly with
jane & september.

∴ 30,000 possibilities. (B × 10,000).

cut down possibilities

in sept ✓
Jane is ✓  } → sept rejected with jane.
jane visits ✓

– 3 copies of NW for different choices of words.

(in september) ⟨ a, aaron, jane, zulu.

(jane is) ⟨ a, aaron, visiting, zulu.

In sept $\hat{y}^{<3>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \square \rightarrow \square \square \square \rightarrow \square$
              $\uparrow$      $\uparrow$
            $x^{<1>}$    $x^{<T_x>}$     $\hat{y}^{<3>}$

jane is $\hat{y}^{<3>}$

encode → $\square \square \square \square$

* Refinements to beam search !

• length normalisation

$\underset{y}{\arg\max} \prod_{t=1}^{T_y} P(y^{<t>} \mid x, y^{<1>}, \ldots, y^{<t-1>}) \Rightarrow x < 1$

$\underset{y}{\arg\max} \sum_{y=1}^{T_y} \log P(y^{<t>} \mid a, y^{<1>}, \ldots, y^{<t-1>})$

log $\log P(y \mid x)$ using
$P(y \mid x)$ using

$$\boxed{\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} \mid x, y^{<1>}, \ldots, y^{<t-1>})}$$  → Normalisation:

$\alpha = 0.7$ ← try different $\alpha$.

$T_y = 1, 2, 3, \ldots, 30$ ↓ Take sentences & compute this, highest value is

kept.

• Beam width B?

     large B : better result          small B : worse result
           slower                     faster.

     on common system    B = 10, 100
                     research sometimes uses B = 1000, ...

→ unlike BFS or DFS, Beam search runs faster but not guaranteed to find exact maximum (or arg $\max_y P(y \mid x)$)

**\* Error analysis in beam search:**

     EX:   Jane visite l'Afrique en septembre
           Human : Jane visits Africa in september $(y^*)$        → RNN
           Algorithm: Jane visited Africa last september $(\hat{y})$     → Beam search. B↑
                                               ↓ for more data.

RNN computes $P(y^*\mid x) \gtrless P(\hat{y}\mid x)$

Case 1:
   $P(y^*\mid x) > P(\hat{y}\mid x)$
     Beam search chose $\hat{y}$, but $y^*$ higher val.
            ∴ Beam search is at fault.

Case 2:
   $P(y^*\mid x) < P(\hat{y}\mid x)$
     $y^*$ better trans than $\hat{y}$, but RNN predicted $P(y^*\mid x) < P(\hat{y}\mid x)$
         RNN model at fault.

• Process

| Human | Algo | $P(y^*\mid x)$ | $P(\hat{y}\mid x)$ | At fault? |
|---|---|---|---|---|
| | | $2\times10^{-10}$ | $1\times10^{-10}$ | B |
| | | | | R |
| | | | | R |
| | | | | B |
| | | | | R |
| | | | | ⋮ |

fraction of error due to beam search vs due to RNN

if Beam search resp, B res

else if RNN is at fault, train deeper N/w.

**\* Attention model intuition :** → translates part of a long sentence at a time.

• Problem of long sentences :

Bleu score



Sentence length.

Paper: Bahdanau, 2014
     Neural M/c translation by jointly learning to align & translate

**BRNN**



$s^{<0>}$ →

$\hat{a}^{<a>}$

$\alpha^{<1,1>}$ $\alpha^{<1,2>}$ $\alpha^{<1,3>}$

$a^{<0>}$ →

$a^{<1>}$   $x^{<2>}$   $x^{<3>}$   $x^{<4>}$   $x^{<5>}$

jane   visite   l'Afrique   en   September

— Refer diagram to screenshot.

---

**\* Attention Model :**

$$\left( \underset{\rightarrow}{a}^{<t>} , \underset{\leftarrow}{a}^{<t>} \right) = a^{<t>}$$

context 'c' fed to upper RNN depends on attention params $\alpha$

$$\sum_{t'} \alpha^{<1,t'>} = 1.$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$$

$\alpha^{<t,t'>} \rightarrow$ amt of attention that $y^{<t>}$ must be paying to $a^{<t'>}$

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

paper : Bahdanau, 2014
Neural m/c translation by jointly learning to align & translate.

$s^{<t-1>}$
$a^{<t'>}$ → $\longrightarrow \dfrac{e^{<t,t'>}}{\alpha^{<t,t'>}}$

paper : Xu , 2015
show attend & tell :
Neural img caption generation with visual attention.

---

Ex)

July 20th, 1969 $\longrightarrow$ 1969-07-20

23 April , 1564 $\longrightarrow$ 1564-04-23.

---

**# Speech recognition — Audio data —**

**\* Speech recognition :**

$x$ (audio clip) $\xrightarrow{\phantom{xxx}}$ $y$ transcript.

"the quick brown fox"

pressure.



$\downarrow$ convert to spectograms ...

— once upon a time, speech recog. used to be built by "phonemes". (de mink brau)

$\rightarrow$ Hand engineered.

**# CTC cost for speech recog.**
(connectionist temporal classification)

$\rightarrow$ "the quick brown fox"



$a^{<0>} \rightarrow$   $x^{<1>}$  $x^{<2>}$   ...   $x^{<1000>}$

$\hat{y}^{<1>}$ $\hat{y}^{<2>}$ ... $\hat{y}^{<1000>}$

the_q
ttt_h_eee----U----999-⏘blank
          space

paper : Graves, 2006
connectionist temporal classification : labeling unsegmented sequence data with RNN

**\*. Trigger word detection:**

Ex: Alexa, Hey siri, Okay google, ...



set labels ①

→ multiple 1's slightly evens out ratio of 1's to 0's.

↳ trigger word finished