PAPER

Van der Maaten & Hinton 2008. Visualizing data using t-SNE

Word embeddings can learn feature vector for words. So instead of using one-hot feature vector for words, we ~~want~~ can use feature vectors for words from word-embedding algorithms.

## W2-L2

* FACE ENCODING is similar to word embeddings (from CNN)

## W 2-L3

PAPER: Mikolov et al 2013, Linguistic regularities in continuous space word representation

surprising results for word embeddings

## W2-L4

* When you learn word embedding, you end-up learning embedding matrix ~~vector~~

Embedding Matrix $\in R^{300 \times |V|}$ : 300 D vector for word

|V| is total number of words in dictionary

## W2-L5   Learning word embeddings

Paper Bengio et al. 2003, A neural probabilistic language model

→ Neural Language Model

For language model: You can take context (last 4-words) & build language model.

For word embedding: You may want to pick last & future words to learn meaningful word embedding.

Paper: Mikolov et al. 2013: Efficient estimation of word representation in vector space

Word2VEC SkipGram Model :

→ SkipGram Model ( problem is denominator )
→ CBOW ( picks surrounding words for learning target word )

( write notes from 06:40 )

GOAL    How to learn word embeddings . Lets take a look at SkipGram model.

Sentence: "I want a glass of orange juice to go along with cereal"

Lets pick ~~target~~ context word ~~these~~ & then its target from window of length "t" words

Context          target
orange           juice

Let say    Vocab size = 10,000

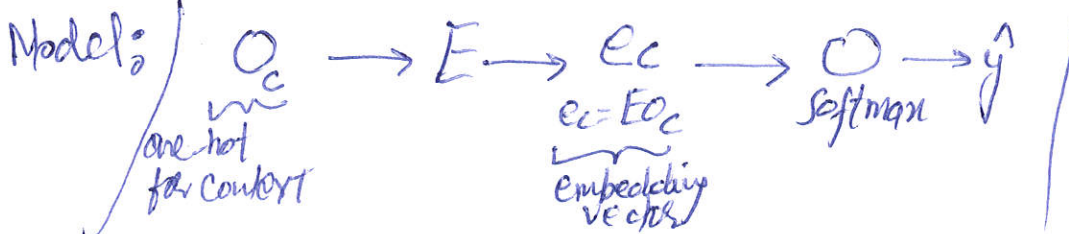TASK    Context "c" ("orange")  ⟶  Target "t" ("joice")
                      6257                      4834

Here   E = embedding Matrix
       = 300 × 10,000  ( Dim of word embeddg × Vocab size )



Model: $O_c \rightarrow E \rightarrow E_c \rightarrow \bigcirc \rightarrow \hat{y}$

one hot for context         $e_c = E O_c$         softmax
                            embedding vector

Softmax:

$$P(target \mid context) = \frac{e^{\theta_t^T e_c}}{\sum\limits_{j=1}^{19,000} e^{\theta_j^T e_c}}$$

$\theta_t$ = parameter associated with output "t"

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

$$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow 4834 \; (target) = juice"$$

Problem

Denominator of softmax $= \sum\limits_{j=1}^{19,000} e^{\theta_j^T e_c}$

It is solved by hierarchical softmax. Have a look for more info. This is one idea to speed up softmax.

Paper: Mikolov et al 2013. Distributed representations of words & phrases and their compositionality
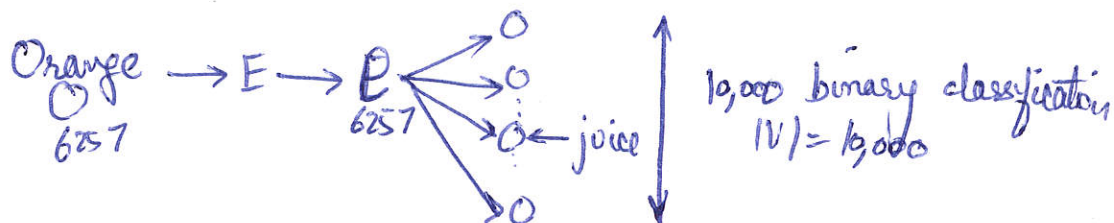
In this lecture:
- Context/target words are training input and $y=1$ if context/target are true pair

| Context | Target | $y$ |
|---------|--------|-----|
| Orange | juce | 1 |
| orange | King | 0 |

- Negative samples are sampled from corpus
- Logistic regression is applied for ~~ear~~ word pairs

$$P\left(y=1 \mid c=\text{Orange}, t=\text{juce}\right) = \sigma\left(\theta_t^T e_c\right)$$

Orange $\rightarrow E \rightarrow e$ 6257
6257

10,000 binary classification
$|V| = 10,000$

---

W2-L8

## Glove word vectors
↳ Another word embedding algorithm

Pennington et al 2014 Glove: Global vectors for word representation
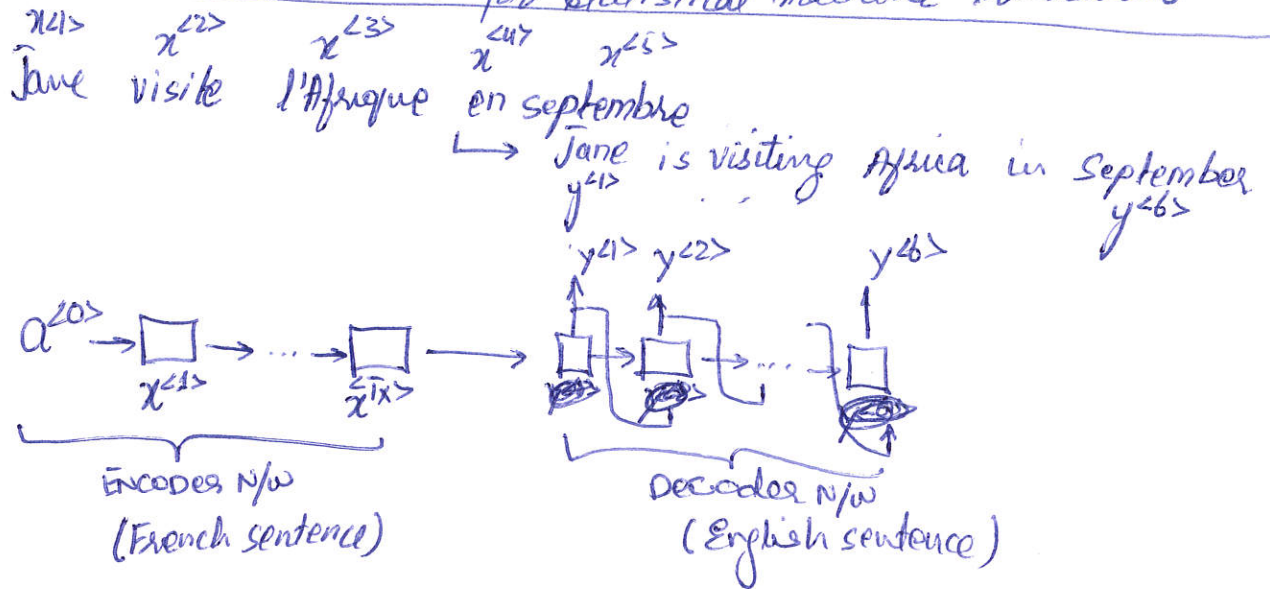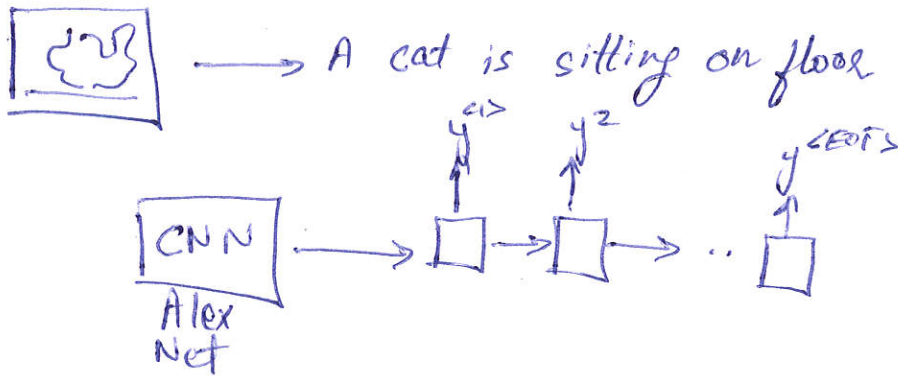
---

W2-L9

## Sentiment Classification

- Simple NN
- RNN
} for sentiment classification

---

W2-L10

Paper: Bolukbashi et al 2016: Man is to computer programmer as woman is to homemaker? Debiasing word embeddings

## SEQUENCE to sequence model

PAPER 1) SutsRever et al. 2014. Sequence to Sequence learning with neural network

2) Cho et al. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation

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

Jane visite l'Afrique en septembre

$\longmapsto$ Jane is visiting Africa in September
$y^{<1>}$                                $y^{<6>}$

$y^{<1>}$  $y^{<2>}$          $y^{<6>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$
$\quad\quad x^{<1>} \quad\quad x^{<Tx>}$

$\underbrace{\qquad\qquad}$  ENCODER N/w (French sentence)

$\underbrace{\qquad\qquad}$  DECODER N/w (English sentence)

## Image Captioning

$\boxed{\text{(image)}} \longrightarrow$ A cat is sitting on floor

$y^{<1>}$        $y^{2}$            $y^{<EOS>}$

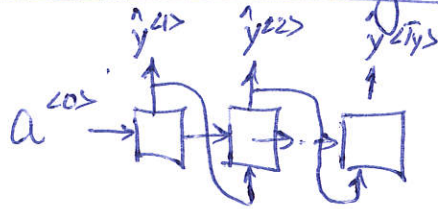$\boxed{\text{CNN}} \longrightarrow \square \rightarrow \square \rightarrow \cdots \square$
Alex Net

PAPERS 1) Mao et al 2014. Deep captioning with multimodal RNN
2) Vinyals et al 2014. Show & tell: Neural image caption generators
3) Karpathy & fei fei 2015: Deep visual-semantic alignments for generating image descriptions

# Picking the most likely sequence

In machine translation, we want to pick the most likely sequence. In language model, we sampled words randomly.
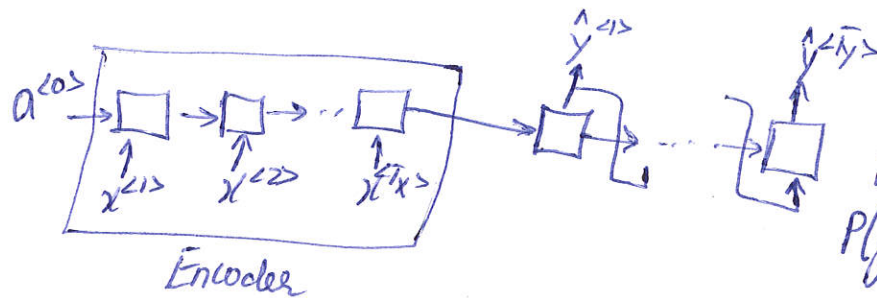
## Machine translation as building a conditional language Model

language model:



Here we were calculating
$P(y^{<1>}, \dots y^{<T_y>})$

Machine translation:



Encoder

Here we want
$P(y^{<1>} \dots y^{<T_y>} | x^{<1>}_{\dots x^{<T_x>}})$

* To find the best sequence $P(y^{<1>}, \dots, y^{<T_y>} | x)$, we will use beam search algorithm.

---

W3-L3        Beam search        (see the video)

W3-L4        Refinements to Beam search

In beam search, we would like to get the most likely sequence $(y)$ such that

$$P(\hat{y}^{<1>}, y^{<2>}, \dots y^{<T_y>} | x) = \underset{y}{\text{argmax}} \left[ P(\hat{y}^{<1>} | x) \times P(\hat{y}^{<2>} | x, \hat{y}^{<1>}) \right.$$
$$\times P(y^{<3>} | x, y^{<1>}, y^{<2>}) \dots$$
$$\left. \times P(y^{<T_y>} | x, \hat{y}^{<1>}, \dots y^{<T_y-1>}) \right]$$

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

# Refinements to beam search

## 1. Length normalization

maximize log probabilities to avoid underflow

2.)
$$\frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log P(y^{<t>}|x, y^{<1>}, \ldots, y^{<t-1>}) \quad \alpha = [0, 1]$$

Normalize log-likelihood objective

---

W3-25    Error-Analysis on beam Search

&rarr; who is at fault
1) Beam Search or 2) RNN ?

---

W3-26    BLEU score (Bilingual evaluation understudy)

Paper:    Papineni et al. 2002. A method for automatic evaluation of machine
translation (very influential)

BLEU score &rarr; for machine translation
&rarr; for captioning system
&rarr; can be used when text is generated
and multiple good text is available
as good text prediction

---

W3-27    Attention Model Intuition

Paper:    Bahdanau et al 2014. Neural machine translation by jointly
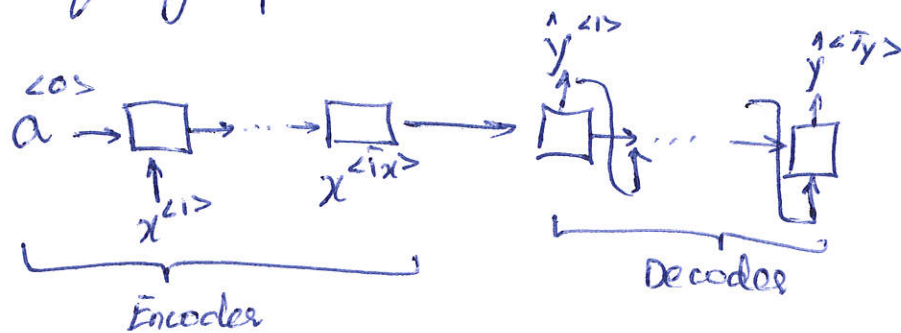learning to align & translate.

---

W3-28    Attention_Model

Paper: 1) Bahdanau et.al 2014. Neural machine translation ....
2) Xu et.al. 2015. Show attention and tell: neural image caption generation
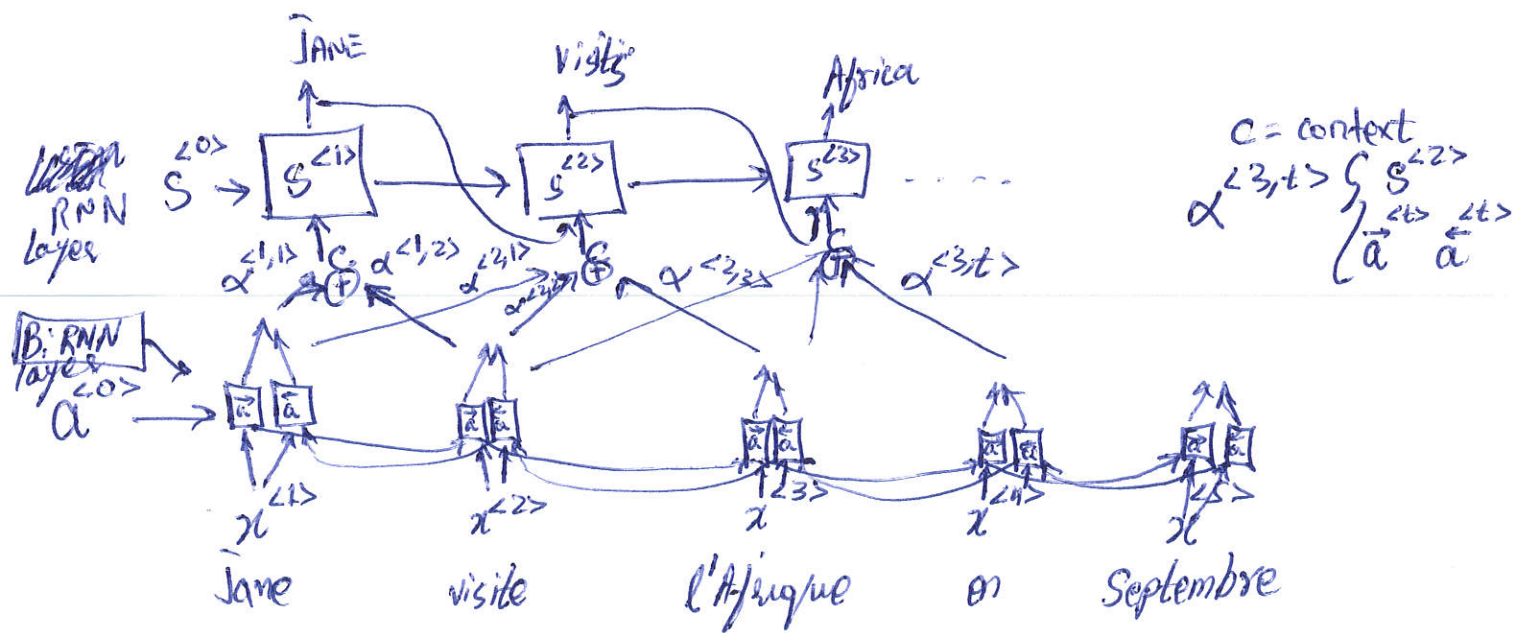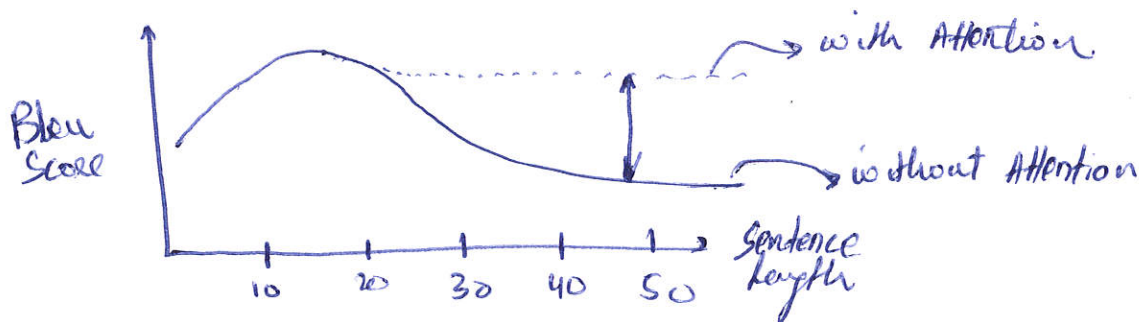with visual attention

# The Problem of long Sequences



Encoder / Decoder diagram

French   Jane s'est rendue en Afrique en ..... aussi.

↓ translate to english

Eng   Jane went to Africa last .... too.

To translate sentence, we take word & its neighbors, translate it to target word rather than reading all text & then translate it to english

## PROBLEM



Bleu Score vs Sentence Length — with Attention / without Attention



$c$ = context

$\alpha^{<3,t>} \begin{cases} s^{<2>} \\ \overrightarrow{a}^{<t>} \quad \overleftarrow{a}^{<t>} \end{cases}$

LISTEN
RNN
Layer

$s^{<0>} \rightarrow s^{<1>} \rightarrow s^{<2>} \rightarrow s^{<3>}$ ....

JANE   Visits   Africa

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

Bi RNN Layer
$a^{<0>}$

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

Jane   visite   l'Afrique   en   Septembre

here: $a^{<t'>} = (\vec{a}^{<t'>}, \overleftarrow{a}^{<t'>})$ Activation features from BiRNN at timestep $t'$

The context $"C"^{<t>}$ is the weighted sum of activation features at time step $(t)$ or:

$$\Rightarrow c^{<1>} \text{ or } = \frac{1}{t'} \alpha^{<1,t'>} a^{<t'>}$$
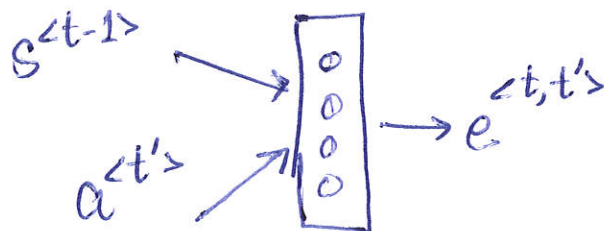
where $a^{<t'>} = (\vec{a}^{<t'>}, \overleftarrow{a}^{<t'>})$

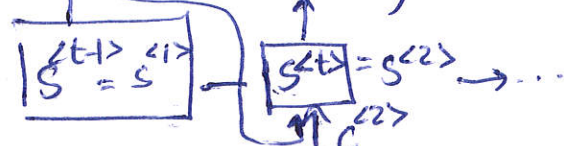And $\frac{1}{t'} \alpha^{<1,t'>} = 1$

$\alpha^{<t,t'>} =$ amount of "attention" $y^{<t>}$ should pay to $a^{<t'>}$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\frac{1}{t'=1}^{T_x} \exp(e^{<t,t'>})}$$
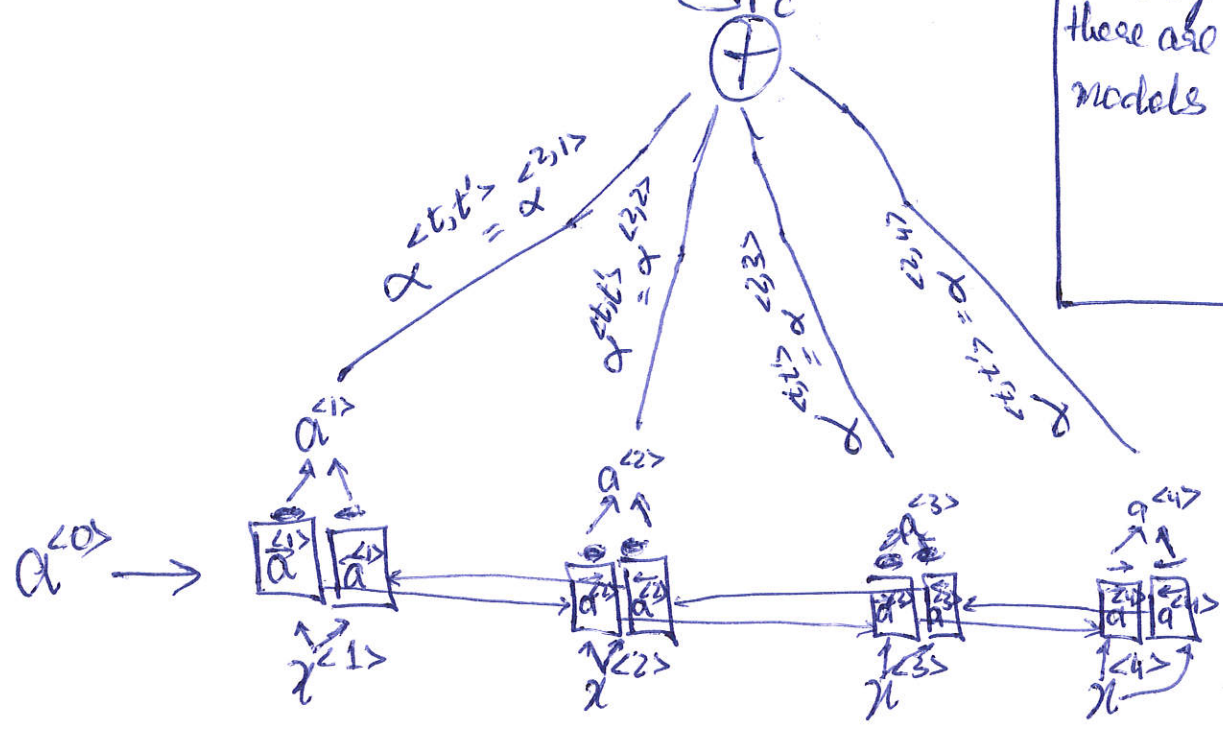
How to get $e^{<t,t'>}$? One way is to get train small NN & learn a function using FF network & believe in gradient descent
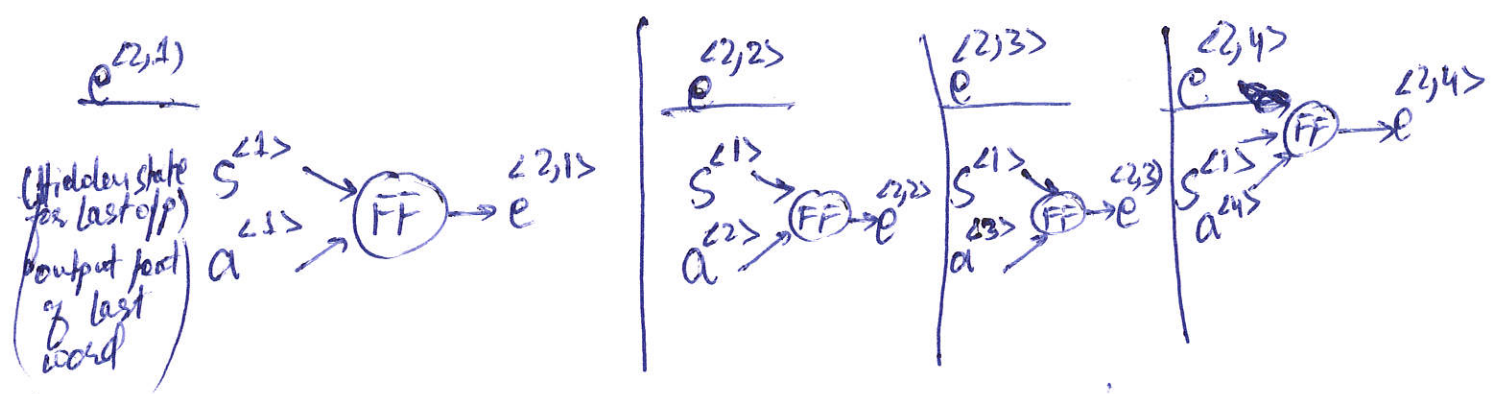
For $y^{<t>} = y^{<2>}$ i.e for 2nd word (we generate output word by word)

$y_1^{<t>} = y^{<1>}$

$y^{<t>} = y^{<2>}$

$s^{<t-1>} = s^{<1>}$

$s^{<t>} = s^{<2>} \rightarrow \dots$

$c^{<2>}$

Here we generate translation word by word although these are encoder-decoder models for translation



$\alpha^{<t,t'>} = \alpha^{<2,1>}$

$\alpha^{<t,t'>} = \alpha^{<2,2>}$

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

$\alpha^{<2,4>} = \alpha^{<2,4>}$

$a^{<1>}$

$a^{<2>}$

$a^{<3>}$

$a^{<4>}$

$\alpha^{<0>} \rightarrow$

$y^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<4>}$

$$\Rightarrow \quad \alpha^{<2,1>} = \frac{\exp\left(e^{<2,1>}\right)}{\exp\left(e^{<2,1>}\right) + \exp\left(e^{<2,2>}\right) + \exp\left(e^{<2,3>}\right) + \exp\left(e^{<2,4>}\right)}$$

$e^{<2,1>}$

(hidden state for last o/p)  $s^{<1>}$

(output word of last word)  $a^{<1>}$

$\rightarrow$ (FF) $\rightarrow e^{<2,1>}$

$e^{<2,2>}$

$s^{<1>}$

$a^{<2>}$ (FF) $\rightarrow e^{<2,2>}$

$e^{<2,3>}$

$s^{<1>}$

$a^{<3>}$ (FF) $\rightarrow e^{<2,3>}$

$e^{<2,4>}$

$s^{<1>}$
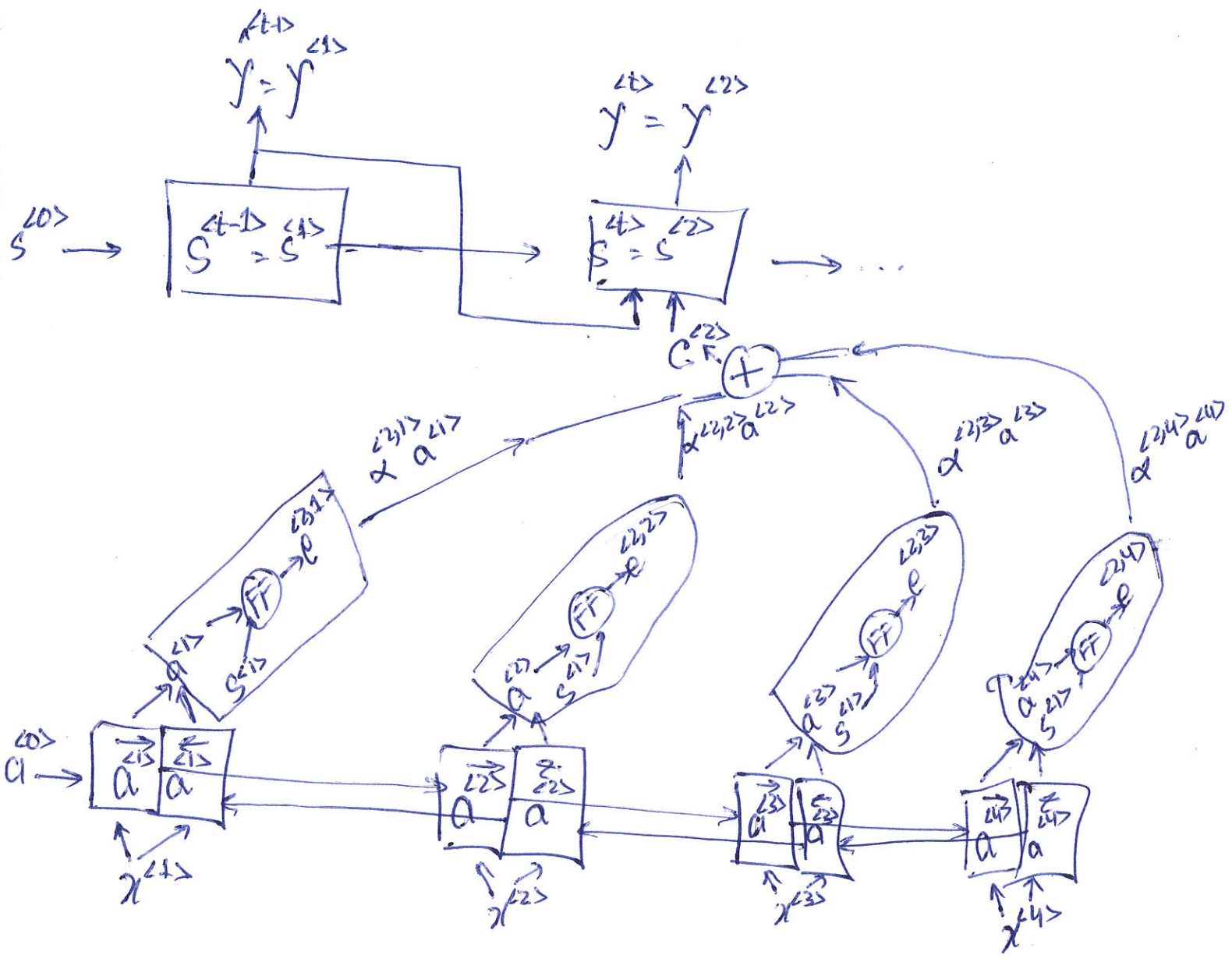
$a^{<4>}$ (FF) $\rightarrow e^{<2,4>}$

The only downside of this algorithm is it runs in quadratic time i.e if input is $T_x$ & o/p sequence is $T_y$ then number of attention terms will be $(T_x)(T_y)$ & we will run FF $(T_x)(T_y)$ times

And $c^{<2>} = \alpha^{<2,1>} a^{<1>} + \alpha^{<2,2>} a^{<2>} + \alpha^{<2,3>} a^{<3>} + \alpha^{<2,4>} a^{<4>}$

For $y^{<t>} = y^{<2>}$   45-D

# Keras Implementation

$$a \in R^{(m, Tx, 2 \times n\_a)}$$

$$s^{t-1} \in R^{(m, n\_s)}$$

Speech recognition

Commercial Products trained on $\geq$ 100,000 hours

* CTC cost for speech recognition
  (connectionist temporal classification)

Graves et al. 2006. Connectionist Temporal Classification: Labelling unsegmented sequence data with RNN

---

W3-L10   Trigger word detection