

sent  
→ Answer  
wayfinding

— we in between threads

tion of syntax

Not a finite automata  
(no memory)

Bottom up parsing / left

Parse tree

J

positive parse tree

A. individual parse tree probability  
on them to get the probability of sentence.

Ancestors  
Ancestral  
→  
Caching

$a + ab$ , matched one in between the two

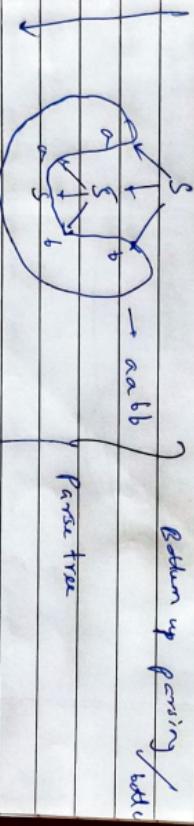
#  
Statistical Parsing:  
↳ verification of syntax

$S \rightarrow a S b$

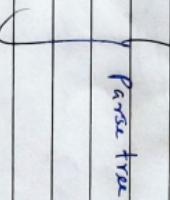
$S \rightarrow ab$  ↗ not a finite automata  
(no memory)

$L = \{a^n b^n : n \geq 1\}$

Top down



Bottom up parsing / etc.



Parse tree

Given a sentence  
generate all possible parse tree  
calculated individual parse tree probabilities  
sum them to get the probability of sentence.

Computers  $\xleftarrow[\text{in language}]{\text{communicate}}$  People

Programming language  
→ deterministic / control flow  
→ processing in own

Date: \_\_\_\_\_  
Page: \_\_\_\_\_

## #

### Introduction:

Temp slide page: 8

Application:

- Translation
- Text to speech / speech to text
- OCR (Character recognition)
- Sentiment Analysis (Opinion opinion)
- Chatbots
- Summarization
- Information retrieval

Human  $\rightarrow$  Machine

Computer vision

Search engines  
logic to knowledge representation  
NLP

Expert systems

- Phonetics and phonology  $\rightarrow$  (back to native language)
- Morphology  $\rightarrow$  (animal = early + noun)
- Lexical analysis  $\rightarrow$  (dog  $\rightarrow$  noun) (adjective)
- Syntactic analysis  $\rightarrow$  (Part Tagging) (Parsing) (Noun)
- Semantic analysis  $\rightarrow$  (word sense disambiguation) (Lexical substitution) NLP
- Pragmatics  $\rightarrow$  model user intention
- Discourse  $\rightarrow$  processing of sequence of sentences (use reasoning by world knowledge)

\* Ambiguity in NLP:

John put the carrot on the plate and ate it.



### Turing test : Text under-

→ Extract relations between

\* Anaphora resolution / coreference:

John put the carrot on the plate and ate it.

ML Ayo:

|                       | Morphology             | POS tagging | HMM, MaxEntropy, CRF, SVM, Logistic Reg., RNN                     |
|-----------------------|------------------------|-------------|---|
| * Question Answering  | lexical                | lexical     | of memory   |
| * Text summarization  |                        |             |   |
| * Sentiment Analysis  | sentiment              |             | Statistical / corpus based :                                      |
| * Machine translation | chunking (verb center) | center      | ML methods need to acquire knowledge from annotated short corpora |
|                       | paraphrasing           |             |   |
|                       | summarizing            |             |   |
|                       | discourse              |             |   |
|                       | coreference            |             |   |

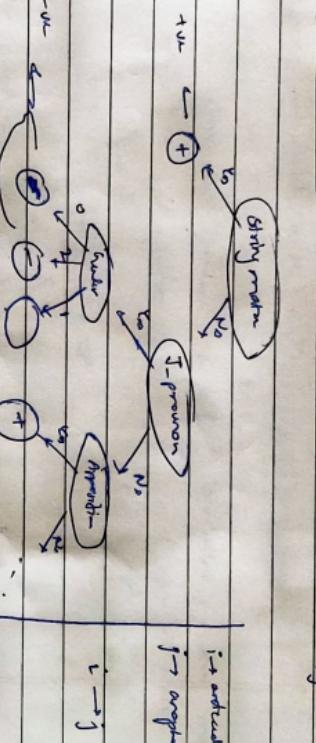
opposition?



### → Features:

- string match, alias, number agreement, gender agreement
- pronoun, defn. clitic (NO type)  $\xrightarrow{\text{defn.}} \text{gender}$
- distance (6 bins against binomial) ( $0, 1, 2, \dots$ )  $\xrightarrow{\text{distance}}$  number
- Semantic class agreement

→ Model: Decision Tree:



→ Note: Decoding: (right to left, consider each constituent C until we find a definite reading for C)

Now the model is already trained, we need to test it on a text.

→ pair of noun phrases  $i, j$ , check if  $\text{Prob}(i | c_{ij}) > \text{Prob}(i | c_{ik})$

if true  $\rightarrow$  they agree or can other else may don't

High context: right to left, consider each constituent & select the higher probability

Evaluation: Precision, Recall, F-score

- Finally returns confidence score (disjointed sets)
- complete confidence matrix



True bank (multiple pass draw)  $P(\text{O} \rightarrow \text{B}|\text{A}) = \frac{\text{Count}(\text{O} \rightarrow \text{B})}{\text{Count}(\text{A})}$



# PCFn : Max. likelihood Training

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$F_1 = \frac{2 \cdot \text{TP}}{2 \cdot \text{TP} + \text{FP} + \text{FN}}$$

# Performance Resolution  
(anglophile vs. ")

John put the carrot ~~the~~ on the plate and ate the ~~the~~ <sup>the</sup> antecedent

(act of carrying back)

Cataphora → anaphora precede antecedent

Because she -- -- -- Many words ahead -- --

→ Constraints:

\* Anaphora-resolution : i) rule based  
ii) statistical based / ML based

→ number agreement  
→ gender "  
→ person / can  
→ grammatical agreement  
→ selectional agreement

→ Preprocessing :

→ pos Tagger

→ Noun phrase identification  
→ Named Entity recognition

→ N-best noun phrase extraction

Encoding

↓ (current - training instance) : noun phrase (anaphora) linked  
current pair → the label turned  $\left[ \begin{array}{c} \text{U5} \\ \text{Tm} \end{array} \right]$   
↳ (all other pairs of statements)  
↳ the answer → new

#### # Sentence (domain independent) :

- local context : preceding & succeeding words
- pos tags (current & surrounding tokens)
- word length (shorter words → many n-grams), frequent → most n-grams
- prefix / suffix (pos length - strings)
- uppercase / lowercase / Digit
- Gazetteer → NE list
- word orientation
- previous occurrence (the class label)

#### # Statistical Parsing :

PFFN : Probabilistic Context Free Grammar

e.g. sentence : Book in flight through' Human

$P(\text{Sentace}) = \sum_{\text{Prob}} (\text{all possible parse-tree})$

$$S \rightarrow NP VP \quad 0.8 \\ S \rightarrow AdjNP VP \quad 0.1 \\ S \rightarrow VR \quad 0.1$$

$$P(\text{parse tree}) = \text{product of probability of production}$$

Perf. tasks :

i) Observation likelihood : Probabilities (evidence)

ii) Most likely derivation : for all parse tree derivations - select the one with max probability : use CKY parser

iii) Maximum likelihood training



\* Inside Algorithm: CKY parser with slight modification -  
(according to forward step of Viterbi) combine probability of multiple derivations  
of any configuration using addition instead of multiplication

(Used for observation likelihood)

## NAVELLS

- Max Entropy
- big, - new, - unique
- ~~est.~~ ~~intra-class~~

information about begin  
Viven Richards plays  
Data Page

I/O → begin  
Viven Richards plays  
Data Page

### # NER C: Name Entity Recognition & Classification:

Extract paper from

class into categories

like person, organization, location, etc.  
(Name - date, time, money, percent)

→ take time

→ statistical based

→ global features

→ local features

NNM Basic NERL :  $\hat{y} = \operatorname{argmax}_{\{P(w_i|t)\}} P(t)$

Bigram model :

prob. of a tag depends only on previous tag

$P(t) = P(t_1) \cdot P(t_2|t_1) \cdot \dots \cdot P(t_n|t_{n-1})$

unigram

$P(w|t) = P(w_1, t_1), P(w_2, t_2), \dots, P(w_n, t_n)$

2nd order NNM:

$P(t) = \mathbb{E} P(t_i | T_{i-2}, T_{i-1})$

Correct tag depends on

precision 2 tags

(Trigram)

# Ensemble Learning with Ensemble Algo

RNN vs → Initialization, fitness computation, selection,

CRF2 uses over, mutation, tournament

Ensemble tags:

- max - predict - avg,

- average, - softmax

Ensemble - data - diversity

Global tags

## ① L continued

# CNN:

Example: image  $\rightarrow$   $224 \times 224 \times 3 \rightarrow$  RGB channels  
- filter  $\rightarrow 11 \times 11 \times 3 \xrightarrow{\text{left}} 3$  (left of filter & image  
- padding = 0       $\xrightarrow{\text{f}} \xrightarrow{\text{f}} \xrightarrow{\text{f}}$  spatial extent  
- stride = 4  
# filters = 36

$$\text{Output} \rightarrow W_2 \times H_2 \times D$$

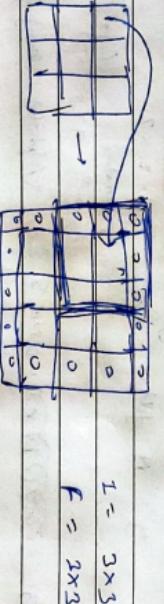
Each filter will produce a 2D convolution / feature map  
 $\rightarrow D = \# \text{ filters}$       ( $D = 3^2 - 2^2 \cdot \text{since}$   
 $\text{Apk} = \frac{1}{2} \cdot \text{filter size}$ )

$$W_2 = W_1 - f + 2P + 1$$

$\xrightarrow{S}$

$$\therefore \text{output} = 224 - 11 + 0 + 1$$

$$= \xrightarrow{56} \xrightarrow{4} \xrightarrow{56} \xrightarrow{36} \xrightarrow{36} \xrightarrow{36} \xrightarrow{36}$$



$$\text{padding} = 1$$

$$\text{for } W_2 = 0, \quad W_2 = W_1 - f + 2P + 1$$

$$0 \times \cancel{10} \xrightarrow{2} P = \frac{f-1}{2} = \frac{2-1}{2} = 1$$

$$\text{eg: edge detector kernel}$$
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

- # Hyperparameter:
  - batch normalization after each conv. layer
  - newton / 2 update
- ↓
  - SGD + Momentum (0.9)
  - learning rate : 0.1
- ↓ mini-batch size  $\rightarrow$  256
- ↓ weight decay  $\rightarrow$   $10^{-5}$
- ↓ no dropout used

# Lecture 12 :





### # Graph Net:

→ at each layer, instead of using multiple filters of same dimension, we filter of different dimensions & then conduct the feature merge (with adequate padding)

$$\rightarrow \text{eg: } (H \times W \times D) \xrightarrow{\text{layer}} (F \times F \times D) \xrightarrow{\text{pool}}$$

$$\text{Output} \rightarrow (H - F + 1) \times (W - F + 1) \times D$$

$$\# \text{ computations per each pixel} = (F \times F \times D)$$

→ To reduce # computations: → Use  $(1 \times 1)$  convolutions

$$F \times F \times D \xrightarrow{\text{use } (1 \times 1) \text{ convolution with } D \text{ filters}} \text{Output} \rightarrow (F \times F \times D)$$

→ At the end before FCN → as:  $1024 \times 7 \times 7 \Rightarrow 1024 \times 1$

(to prevent too much parameters)



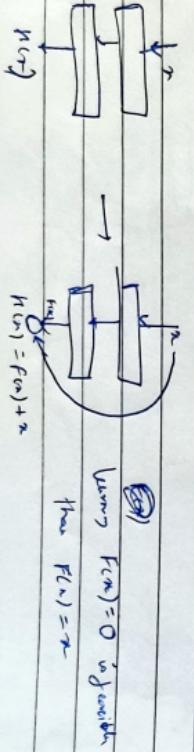
$\rightarrow 1 \times 1$  layer performs then AlexNet  $1024 \xrightarrow{g=36}$

$\rightarrow$  2x more computation (multiple layers: conv, maxpool, inception)

### # ResNet:

→ deeper CNN should perform better but performed worse as it couldn't learn identity function well

→ problem lies in initialization values



$$f(x) = f(x) + r$$

Maxpol doesn't change the depth or add new parameters  
It's just a max operation  
(max pooling operation is perfectum)  
(not across the depth)

Date \_\_\_\_\_  
Page \_\_\_\_\_

#### # Alex Note:

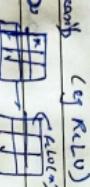
Input → Conv → Maxpool → Linear → Maxpool → Conv → Loss

→ maxpool → FC1 → FC2 → Output  
 $(2 \times 2 \times 246) \rightarrow (1094) \rightarrow (4386) \rightarrow (1000)$

# layers = 5 + 3 = 8  
(max pool not counted as  
layer)

\* 27. 55 million params  
(params in FCL :  $4M + 16M + 4M \approx 24M$ )

Note: \* For every convolution operation, non-linearities (e.g. ReLU)  
\* also applies to the volume.



\* Now the pooling operation can be performed.

# ZFNet: Some # layer, more parameters — better

# Vgg Net: Each filter is 3x3  
16 layers

Input → C C M C C M C C M C C M F F F  
↓  
maxpool  
↓  
filter

# Parameters in FC layer  $\approx 12.2M$

Non-Ps → 16M

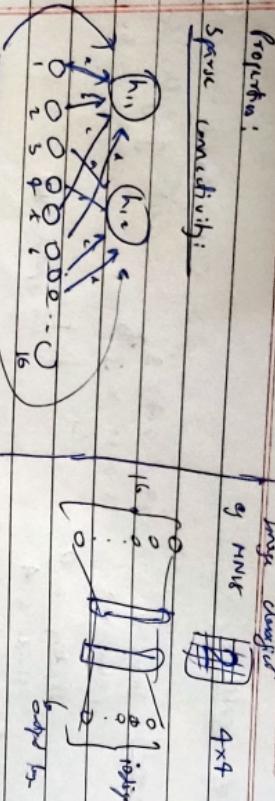


(CNN → multiple convolution filter → multiple layer)



#

Projection:  
Sparse connectivity:



filter  $\rightarrow$  2x2

$g_{11} = 2$

not all the pixels

$l_{11} \rightarrow 3, 4, 7, 8$

To increase receptive field → increase depth

iii) weight sharing :

Using only 1 filter, no. of weight (parameter) =  $2 \times 2 = 4$  for first layer.

→ underlying

sel' → multiple filters (kernel)

$y : N$  feature  $\rightarrow N \times L \times 2$  (  $N \times f_N$  )

padding (for 1<sup>st</sup> hidden layer)

(addition)  $N \times f_N \times D$   $\rightarrow$   $N \times f_N \times D$  (for 1<sup>st</sup> hidden layer)

g: ~~CONV~~ DP. kernels : (Alphab. 'convolution')

$$I : \left( 4 \times 14 \times 6 \right) \xrightarrow[s=1, p=5]{k=16, r=0} (10 \times 10 \times 16)$$

$$\text{feature} = (5 \times 5 \times 6) * 16 = 2400$$

filter size  $\#$  filter fully connected

$$(10 \times 10 \times 16) \xrightarrow{\text{maxpool}} 5 \times 2 \xrightarrow{k=16, p=0} (5 \times 5 \times 16)$$

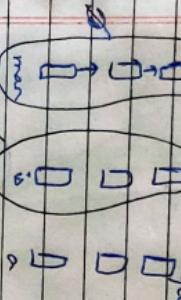
padding = 0

$$f_{\text{output}} = \frac{1}{16} \sum_{i=1}^{16} (x_i - \bar{x})^2$$

$$= 4000 \times 128 + 128 \left( \frac{128 \times 128}{16} \right) + 84$$

### # RNN:

- input size is not fixed
- input are dependent on each other
- task network performs the same task



↳ Many time networks perform same task (input → word)

$$y_t = f(s_t) \quad (\text{current output at every time step})$$

(input → tag)

$$s_t = \sigma(Ux_t + Ws_{t-1} + b)$$

$$y_t = f(s_t) \quad (\text{output activating})$$

$$s_t = \sigma(Ux_t + Ws_{t-1} + b)$$

(n-dimensional input)

$$x_t \in \mathbb{R}^n \quad (\text{n-dimensional input})$$

$$y_t \in \mathbb{R}^k \quad (\text{tag K class})$$

$$U \in \mathbb{R}^{n \times n}$$

$$V \in \mathbb{R}^{n \times k}$$

$$W \in \mathbb{R}^{k \times n}$$

$$\int_{\text{out}} \int_{\text{out}}$$

↳ We have many neural networks  $\theta$

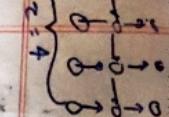
$$\text{Total Loss} \rightarrow L(\theta) = \sum_{t=1}^T l_t(\theta)$$

$$L(\theta) = -\log(y_{te})$$

(probability of true channel and then step t)

HAWELLS

$$E_t = \sum_{k=1}^n (\hat{y}_k - y_k)^2$$



$$E_t = \frac{1}{2} \sum_{t=1}^T \sum_{k=1}^n (\hat{y}_{kt} - y_{kt})^2$$

↓  
total loss  
computation  
↓  
loss  
from all  
outputs)

$$y_{kt} = \mathbf{U}_{kt} + \mathbf{W}\phi(h_{k-1})$$

$$\mathbf{y}_t = \mathbf{V}\phi(\mathbf{h}_t)$$

Now to find the neurons  $U, V, W$

$$\text{compute } \frac{\partial E}{\partial w} = \frac{\partial E}{\partial v} = \frac{\partial E}{\partial u}$$

$$\frac{\partial E}{\partial w} = \sum_{t=1}^T \frac{\partial E_t}{\partial w} = \sum_{t=1}^T \sum_{k=1}^n \frac{\partial E_t}{\partial h_{kt}} \cdot \frac{\partial h_{kt}}{\partial w}$$

$$\frac{\partial h_{kt}}{\partial w} = \frac{\partial \phi_r}{\partial h_{k-1}} \cdot \frac{\partial h_{k-1}}{\partial w} = \frac{\partial \phi_r}{\partial h_{k-1}} \cdot \frac{\partial h_{k-1}}{\partial h_{k-1}}$$

$$\phi_r : h_i = \mathbf{U}_{ri} + \mathbf{w}^T \phi(h_{i-1})$$

$$\frac{\partial h_i}{\partial h_{i-1}} = \frac{\partial \mathbf{U}_{ri}}{\partial h_{i-1}} + \mathbf{w}^T \frac{\partial \phi'(h_{i-1})}{\partial h_{i-1}} = \frac{\partial \mathbf{U}_{ri}}{\partial h_{i-1}} + \mathbf{w}^T \frac{\partial \phi'(h_{i-1})}{\partial h_{i-1}}$$

$$\frac{\partial E}{\partial h_{ri}} < 1 \Rightarrow \frac{\partial h_{ri}}{\partial h_{i-1}} \rightarrow 0 \quad (\text{explosion})$$

then  $\frac{\partial E}{\partial v} \rightarrow 0$  or the weight matrix will stop  
upgrading its weights or

$$(w' = w - \eta \Delta w \Rightarrow w' = w) \quad (\text{stagnation})$$

If  $\frac{\partial h_i}{\partial h_{i-1}} > 1$  then  $\frac{\partial h_i}{\partial h_{i-1}} \rightarrow \infty \rightarrow$  exploding gradient

for  $\frac{\partial h_i}{\partial h_{i-1}} < 1$  then gradient is  $\frac{\partial h_i}{\partial h_{i-1}} \rightarrow$  KILLED weights update directly

$\frac{\partial h}{\partial v}$  where ( $v$ ) = tensor ( $d \times d \times d$ )  
(and)

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\text{Op } h_t = U_{M_t} + M_t \Phi(h_{t-1})$$

$$\Phi(h_t) \neq \Phi(h_{t-1})$$

$$\frac{\partial h_t}{\partial v} \neq \Phi(h_{t-1})$$

$h_t(h_{t-1})$  also depends on  $v$

$\frac{\partial h_t}{\partial v} = \frac{\partial U_{M_t}}{\partial v} + \frac{\partial M_t}{\partial v} \Phi(h_{t-1})$   
Now again note  
approximate  
invariance  
(trust all other  
inputs as constant)

→ solution next, objective cost, selection forget

WITH → solves the problem of vanishing gradient  
•  $s_t = f_t \odot s_{t-1} + b \cdot \tilde{s}_t$  → comes all info from the start  
info for morphed boy and morphed

solution forget selection next forget  
 $s_t = f_t \odot s_{t-1} + b \cdot \tilde{s}_t$

$s_t$  com

$$s_t = \sigma(w_i^T h_{t-1} + v_i^T r_{t-1} + b)$$

$$f_t = \sigma(w_f^T h_{t-1} + v_f^T r_{t-1} + b)$$

$$r_t = \sigma(w_r^T h_{t-1} + v_r^T r_{t-1} + b)$$

$$h_t = \sigma(w_h^T h_{t-1} + v_h^T r_{t-1} + b)$$

or forget

In GAN → ~~GAN~~ → forget get & loss given are tied  
→ ~~GAN~~ domain in ~~GAN~~  $s_{t-1}$  & not  $h_{t-1}$   
for  $s_t = \sigma(w_s^T h_{t-1} + v_s^T r_{t-1} + b)$   $s_t = (1 - \eta) \odot s_{t-1} + \eta \odot s_t$   
 $r_t = \sigma(w_r^T h_{t-1} + v_r^T r_{t-1} + b)$   $+ i_{40} \tilde{s}_t$

I'm going to town

#

Attention:

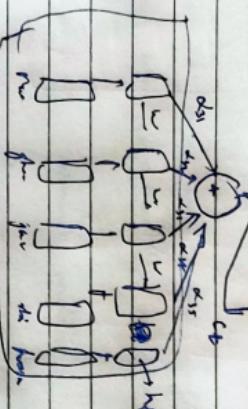
$y_{t-1}$

$y_t$

$s_{t-1}$

$s_t$

$y_t$



The output "going" depends on the input words "jar" and "plus" so at  $t=3$ ,  $d_{ht}$  and  $d_{yt}$  should have more weight

$$\text{Encoder: } h_t = \text{RNN}(h_{t-1}, x_t)$$

$$s_0 = h_T$$

$$\text{Decoder: } \hat{c}_{tj} = F(s_{t-1}, h_j)$$

(combination of  $m_j$  from input word at timestamp  $t$ )

$$\hat{c}_{tj} = c_{tj} = V^T \tanh(W_h h_{t-1} + W_c h_j)$$

$c_{tj}$  is scalar  
for learnable parameters, plus  $\mu_{tj}$

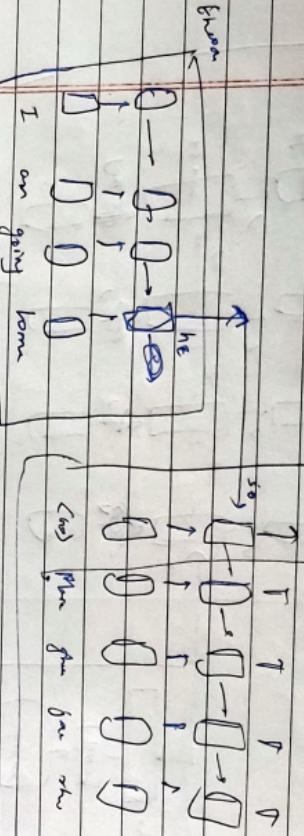
$$\text{decoder softmax } (\hat{c}_{tj}) = \frac{\exp(c_{tj})}{\sum_{k=1}^K \exp(c_{tk})}$$

\* Attention is itself a FFNN with the params  $W_h, W_c, U_a$

$$c_{tj} = \sum_{j=1}^J d_{tj} \cdot h_j$$

$$d_{tj} = F_{\text{Att}}(s_{t-1}, [c_{tj-1}, c_t])$$

A Medium translation:



Lesson

ENGLISH

ht = RNN (ht-1, xt)

so = ht

st = RNN (st-1, e(y^t))

or

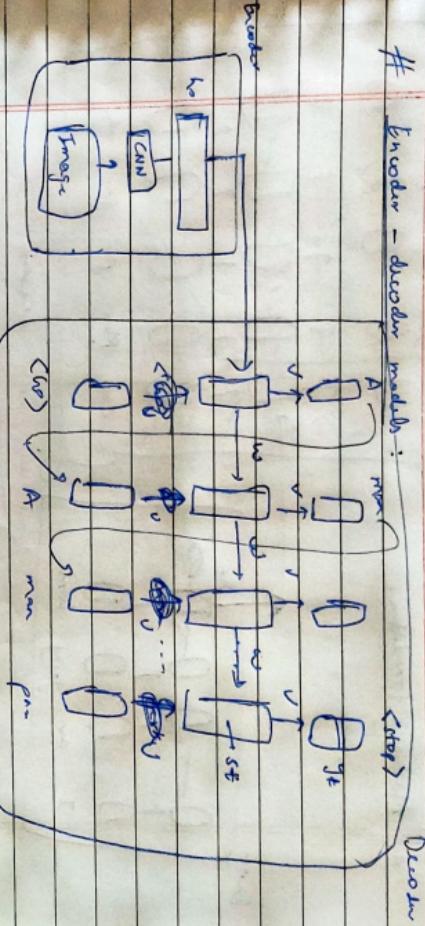
st = RNN (st-1, [ht, e(y^t)])

conduction of output from (t-1) to encoding ht

i.e. ht is passed in every neural network



# Encoder - decoder module:



→

At every time stamp, we have to predict probability distribution of the  $y_t = j | y_{t-1}, I$  → probability distribution of the  $j^{th}$  word at timestamp  $t$

$$= P(y_t | s_{t-1}, f_{T(t)}) \quad (j \in \text{English vocabulary})$$

( $s_t$  encodes all info from 0 to  $t-1$ )

→

② Input at timestamp  $t =$  output predicted at time stamp  $t-1$  (eg. can workforce to convert input to vector)

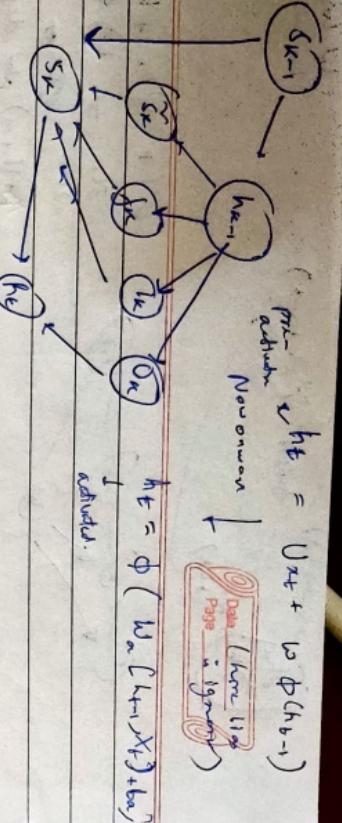
→ Rather than feeding the encoded image at  $t-1$ , we can feed it to all neural networks

$$\text{Output } \delta^t = \text{softmax} \left( w s_{t-1} + b_m + b \right)$$

→

$$(y_t, \delta^t \rightarrow \sum \text{cross entropy}) \quad P(y_t = j | y_{t-1}) = \delta_j^t \text{softmax}(V \delta^t + C_j)$$

(across all time steps)



$$\text{antiderivative} \sqrt{ht} = \frac{1}{3}ht^{\frac{3}{2}} + C \quad (\text{check!})$$

Starting from SK-1 Akkei we reached SK and HK

Different Notation :  
 We have  $M_{t-1} = h_{t-1}$  — memory cell

Updation rule :  $\hat{Y}_t = \sigma(\text{W}_0 [m_{t-1}, X_t] + b_0)$   
 $\propto$  original value  
 $\times$  current value

~~simplyfied~~  
G.R.D  
~~then~~ even if  $F_1 \rightarrow 0$  due to uniting gradient,  $m_1$  want to 0

Fall GRD:  
 Relevanzgrade:  $r_t = \sigma(\text{W}_t [m_{t-1}, x_t] + b_t)$   
 (mit  $m_t$ )  $k_t = r_t \odot m_t$

LSTM:  $\sigma$  ( $b_t^h [h_{t+1}, u_t] + b_t^c$ )

$$\begin{cases} \sigma = \sigma_0 (\mu_0 [k_{\text{eff}}, n_0] + \Delta) \\ \sigma = \sigma_0 (\mu_0 [k_{\text{eff}}, n_0] + \Delta_0) \end{cases}$$

$$m_t = \text{tanh} \left( \beta m [h_{t+1}, x_t] + b_m \right)$$

$$m_t = f_0 \odot m_t^+ + f_1 \odot m_t^-$$

$$h_t = f_0 \odot m_t$$

Falk GRO:

1

Curri

231

187

1

1

1

1

1

1

10

1

1

# NLP Continued

Date: 5/1/2018  
Page: 1

## # Sediment Analysis:

→ Aspect based (fine grained)      Surface / sediment based (coarse grained)  
(bottom, camera, display))      (sound phone))

→ Polarity → +ve, -ve, neutral, conflict

→ Aspect from extraction, not from identifying the polarity

→ CNN based sediment analysis: gram, 37mm → etc

→ Hybrid: CNN + optimized feature set (very GA based optimization)  
+ Train SVM using optimized feature & trained CNN

Optimized features → sentence, punctuation, etc

## # RNN / LSTM based.

→ LSTM

→ Target dependent LSTM

→ Target correction LSTM

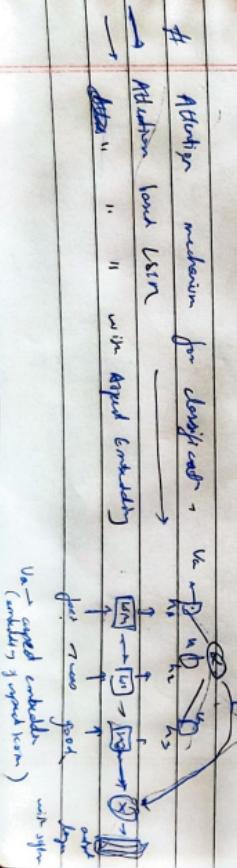
src → target

Out of Vocabulary words → translate into target language

RNN → to perform history in target bilingual

embedding

way skip from unknown (source bilingual word embedding (B2S2B))  
method share the same vector space w/ target RNN

→ Attention mechanism for bilingual → 

HAVELLS

Maximum likelihood training : softmax  $\rightarrow$  PCFG  
— unpruned

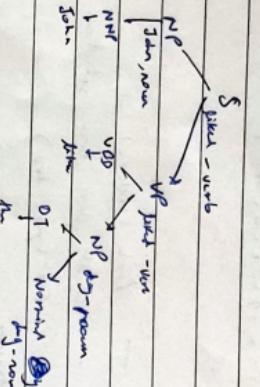


### \* lexicalization of PCFG:

- include head word in LHS non-terminal  
most central to the phrase

+ POS

e.g.



Notes: only for non-terminals

$A \leftrightarrow B \leftrightarrow C \leftrightarrow \dots$

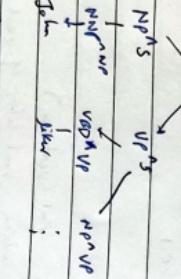
Head for  
 $NP \rightarrow \text{verb}$   
 $NP \rightarrow \text{noun}$   
 $VP \rightarrow \text{verb}$   
 $VP \rightarrow \text{noun}$   
 $NP \rightarrow \text{adj-noun}$   
 $NP \rightarrow \text{adv-noun}$   
 $NP \rightarrow \text{prep-noun}$   
 $NP \rightarrow \text{nominal}$

f. Collins' parser — rules

+ splitting non-terminals :

- more conditions like
- e.g. parent annotation

e.g.  $VRS \rightarrow VRS$



+ Evaluating: correct tree T

Computed tree P

$$\text{Pr}(T) = \frac{1}{Z} \prod_{i=1}^n \text{Pr}(T_i | \text{conditions}_i = z_i)$$

$$\text{Pr}(P) = \frac{1}{Z} \prod_{i=1}^n \text{Pr}(P_i | \text{conditions}_i = y_i)$$

# correct conditions =  $Z$

# correct conditions =  $Z$

## # Statistical Parsing :

5

→ Probabilistic Context free grammar

$$\begin{array}{l} \text{Prob (sentence)} = \sum \text{Prob (all parse trees)} \\ \rightarrow \text{Prob (parse tree)} = \prod (\text{all transition prob}) \\ = 0.5 \times 0.9 \times 0.1 \times 0.2 \times 0.3 \end{array}$$

\* PCFG parser,

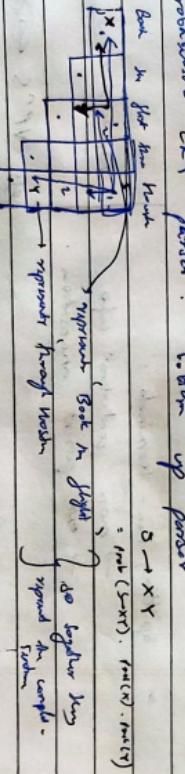
i) Observation likelihood  $\rightarrow$  Prob (sentence)

ii) Most likely derivation  $\rightarrow$  of all parse tree, select one with max Prob (parse tree)

iii) Maximum likelihood training.

$$\begin{array}{ll} \# \text{ CNF : } & A \rightarrow BC \quad | \quad A \rightarrow BX \quad | \quad \text{no non-term. domain} \\ B \rightarrow \epsilon & \longrightarrow \quad X \rightarrow \epsilon \\ C \rightarrow t & \quad B \rightarrow \epsilon \quad | \quad \text{prob} = 1 \\ D \rightarrow j & \quad C \rightarrow t \quad | \quad \text{terminal} \\ & \quad D \rightarrow j \end{array}$$

# Probabilistic CKY parser : bottom up parser



•  $\rightarrow$  terminals: Num : 0.1

Num - Num = 1

# Infill Algorithm:

for observation token  $\rightarrow$  Prob (sentence)

$\rightarrow$  In case parser : combine probability of each cell (instead of max)  $\rightarrow$  in case of most likely derivation.

q: cell S: represents the complete sentence so no observation needed

$\rightarrow$  Prob (of cell S)

## # Word Embedding:

text  $\rightarrow$  vector (vector space mode)

Non-sentential representation  $\rightarrow$  each word  $\rightarrow$  one hot vector

Problem  $\rightarrow$  no semantics captured, high dimension of vector, sparse

Sentential word representation:

i) Co-occurrence matrix

I like big banan

I like NLP

like NLP

like big

big NLP

big banan

banan NLP

banan big

big big

big NLP

NLP NLP

NLP big

big NLP

$$\text{Window size} = 1 \quad \begin{matrix} \text{like word} \\ \text{big word} \\ \text{banan word} \\ \text{NLP word} \end{matrix}$$

|       |   |   |   |
|-------|---|---|---|
| I     | 0 | 2 | 0 |
| like  | 2 | 0 | 1 |
| big   | 0 | 1 | 0 |
| banan | 0 | 0 | 1 |

symmetric

very high dimension, sparse

ii) context vector

skip gram: predict context words given current word

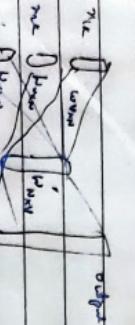
minimize  $J = -\log p(\text{like word}, \text{big word}, \text{banan word}, \dots, \text{NLP word})$

$$J = -\log p(\text{like word}, \text{big word}, \text{banan word}, \dots, \text{NLP word})$$

Window size =  $m$

→ continuous bag of words: predict current word based on other context words

minimize  $J = -\log p(\text{like word}, \text{big word}, \text{banan word}, \dots, \text{NLP word})$



like

big

banan

NLP

output

big

banan

NLP

## # Word Embedding:

text  $\rightarrow$  vector (vector space model)

$\rightarrow$  Non-syntactic representation  $\rightarrow$  each word  $\rightarrow$  one lat vector  
Problem  $\rightarrow$  no semantics captured, high dimension of vector, sparse

$\rightarrow$  Semantic word representation:

i) Co-occurrence matrix

I like big lemons

I like NP

|      | 1 | like | NP | after |
|------|---|------|----|-------|
| 1    | 0 | 2    | 0  |       |
| like | 2 | 0    | 1  |       |
| NP   | 0 | 1    | 0  |       |

$\rightarrow$  symmetric

$\rightarrow$  very high dimension - sparse

ii) word2vec:

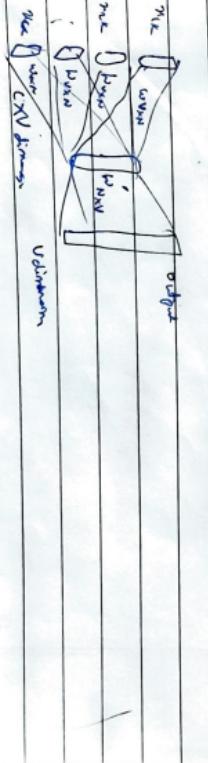
$\rightarrow$  skip gram: predict context words given current word

minimize  $\sum_{i=1}^n \sum_{j=1}^{m_i} \log p(w_j | w_i)$

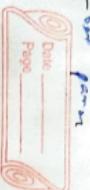
$$w_i \rightarrow \boxed{w_i} \rightarrow J = -\log p(w_1, w_2, \dots, w_m | w_i)$$

$\rightarrow$  Output: predicted context word based on the context word

minimize  $J = -\log p(w_1 | w_i, \dots, w_{i-1}, w_{i+1}, \dots, w_m)$



→ Parse reading : build selected parse from N-best parse



→ Human Parsing:

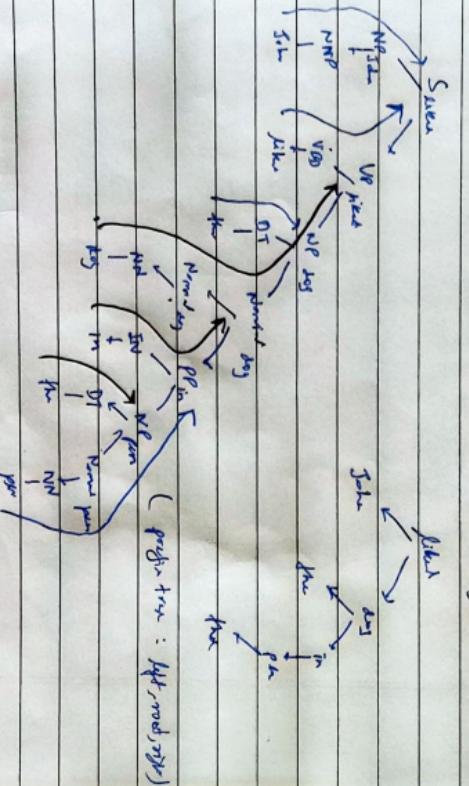
→ Gender path sentence : The old man the star

→ Center embeddings → n-ary expressions

→ Dependency grammar :

(from parse tree)

Dependency graph



lexical knowledge NN:

lexical knowledge NN:

Program : some sort / holding last 4 hr. morning (or 7-8)

→ heterogeneity: one body used → to  
start for another  
e.g.: common foraging

stand for another.

single word, multiple

whale meaning . y: b

question | instrument : relational situation

*Hypernym / Hyperonym* (superclass - subclass) (parent - child)

卷之三

logically under  $V_2$  :

Temporal Inclusion : VI temporally include verb or Ep:

**→ homonym:** Verb in specific form & its more general form

$V_1 = \frac{\text{stone}}{\text{work}} / \text{teller}$

can → what / why: need, for

→ Radiation: series of gradual changes