

Training set - 80%.

Development set - 10% (Frequently used test data).

Test set - 10%.

$$PP(w) = \sqrt[N]{\frac{1}{P(w, w_2, \dots, w_N)}}$$

$$\text{Entropy } H(x) = -\sum_{x \in X} P(x) \log_2 P(x)$$

$$= \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

$$P(w_1, w_2, \dots, w_N) = P(w_1) P(w_2 | w_1) \dots P(w_N | w_{N-1})$$

PP of Unigram  $\rightarrow$  PP of Bigram  $\rightarrow$  PP of Trigram

Lower PP, more the information

OOV rate

Out of vocabulary

PP's should be compared across language models with the same vocabularies.

Smoothing / Discounting: Done to avoid getting 0 probes for unseen events.

$$P(w_i) = \frac{c_i + 1}{N + V}$$

$$c_i^* = \frac{(c_i + 1)N}{N + V}$$

$$d_c = \frac{c^*}{N + V}$$

$$P_{\text{Laplace}}(w_n | w_{n-1}) = \frac{c(w_n, w_{n-1}) + 1}{\sum_w c(w_n, w) + 1}$$

Add-k smoothing:

$$P(w_n | w_{n-1}) = \frac{c(w_n, w_{n-1}) + k}{c(w_{n-1}) + kV}$$

Held-out estimation:

$$P(w_1, \dots, w_N) = \frac{T_r}{N_r N}$$

## Gross-estimation:

$$P_{\text{G}}(w_1, \dots, w_n) = \frac{T_Y^{0+} + T_Y^{10}}{(N_0^0 + N_0^1) N}$$

Back-off vs. Interpolation

## Katz back-off:

$$P_{\text{Bo}}(w_n | w_{n-1}, \dots, w_1) = \begin{cases} P^*(w_n | w_{n-1}, \dots, w_1) & \text{if } c(w_n, w_{n-1}, \dots, w_1) > 0 \\ \alpha(w_n, w_{n-1}) \cdot P_{\text{Bo}}(w_{n-1} | w_{n-2}, \dots, w_1) & \text{otherwise} \end{cases}$$

Good steering:  $c(w_n, w_{n-1}, \dots, w_1) = 0$  if  $w_n = w_{n-1}$

$$\gamma^* = \frac{(r+1) E(N_{r+1})}{E(N_r)}$$

$$P(w_1, \dots, w_n) = \frac{\gamma^*}{N}$$

Kneser-Ney smoothing: (Absolute discounting)

$$P_{\text{Absolute discounting}}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) - d}{\sum_v c(w_{i-1}, v)} + \eta(w_{i-1}) P(w_i)$$

$$P_{\text{CONTINUATION}}(w) \propto |\{v : c(vw) > 0\}|$$

$$\eta(w_{i-1}) = \frac{|\{v : c(vw) > 0\}|}{|\{u^l, w^l : c(u^l w^l) > 0\}|}$$

$$P_{\text{KN}}(w_i | w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - d, 0)}{c(w_{i-1})} + \eta(w_{i-1}) P_{\text{CONT}}(w_i)$$

$$\eta(w_{i-1}) = \frac{d}{\sum_v c(w_{i-1}, v)} |\{w : c(w_{i-1} w) > 0\}|$$

Witten-Bell smoothing:  $N \rightarrow \text{no. of tokens}$   
 $T \rightarrow \text{no. of types (not v)}$

$$z = \sum_{i:c_i=0} 1$$

Good Turing count for most frequent word (like "the") is undefined

$$c^* = \frac{(C+1)N_{c+1}}{N_c}$$

$$P_c = \frac{c^*}{N}$$

jii  $\Rightarrow$  Joint if. independent

$$\text{jii} = N_e (p(x) e(p(y))) \text{ for } xy$$

$N$  bins

Typically 3 bins per decade.

Simple Good-Turing: (Monte Carlo example)

$$\log(N_i) = a + b \log i$$

a, b are fit by linear regression

Closed class  $\Rightarrow r^* > 1$

$$r^* = \frac{(r+1)N_{r+1}}{Nr} = \frac{(r+1)A(r+1)^b}{Ar^b} \quad \begin{array}{l} \textcircled{1} \quad r^* < r \\ \textcircled{2} \quad \frac{r^*}{r} \uparrow \text{as } r \uparrow \end{array}$$

$$= r \left(1 + \frac{1}{r}\right)^{b+1}$$

$$r^* < r \text{ if } b+1 < 0 \Rightarrow b < -1$$

Low frequencies  $\xrightarrow{\text{Good}} \text{Turing estimates}$

Higher "  $\Rightarrow$  SGT

$$\begin{aligned} \text{Prob. of unseen unigrams} &= \frac{T}{Z(T+N)} \quad \text{B. a.} \\ \text{" seen "} &= \frac{c_i}{T+N} \end{aligned}$$

$$\begin{aligned} \text{Prob. of unseen bigrams } p^*(w_i|w_x) &= \frac{T(w_x)}{Z(w_x)(N(w_x) + T(w_x))} \\ &= \frac{c(w_x, w_i)}{N(w_x) + T(w_x)} \end{aligned}$$

$$Z(\omega) = V - T(\omega)$$

$\hookrightarrow$  seen bigram types starting with  $\omega$ .

$N(\omega)$  = No. of bigram tokens starting with  $\omega$

## POS Tagging

Assign POS to each word in a sentence

## Discourse & Coherence

## Semantics extraction

## Parsing

## Chunking

↑ Increased complexity

## POS tagging Morphology

## Named Entity Recognition

Company : Boeing Co.

Location : Wall Street

Person : Alan Mulally

NA : No entity

Local vs contextual constraints

↓  
More likely in  
general

$x(i) \rightarrow \text{Input}$        $y(i) \rightarrow \text{Output}$

Eg:  $x(i) \rightarrow \text{The dog laughs}$

$y(i) \rightarrow \text{DT NN VB}$

Conditional models :  $P(y|x)$

Generative models :  $P(x, y) = P(y) \cdot P(x|y)$

learn

Poster

$$P(y|x) = \frac{P(y) \cdot P(x|y)}{P(x)}$$

Discriminative model

$$f(x) = \operatorname{argmax}_y P(y|x)$$

$$= \operatorname{argmax}_y \frac{P(y) \cdot P(x|y)}{P(x)}$$

$$= \operatorname{argmax}_y P(y) \cdot P(x|y)$$

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n)$$

↓  
Sentence

↓  
Tag Sequence

$$\underset{y_1, \dots, y_n}{\operatorname{argmax}} P(x_1, x_2, \dots, x_n, y_1, \dots, y_n)$$

Trigram HMM's:

$V \rightarrow$  sentence  $(x_1, x_2, \dots, x_n)$

$S \rightarrow$  Tag Sequence  $(y_1, \dots, y_n)$

$y_{n+1} = \text{STOP}$

$$P(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = \prod_{i=1}^{n+1} q(y_i | y_{i-2}, y_{i-1}) \prod_{i=1}^n e(x_i | y_i)$$

↑ Trigrams      ↑ Emission  
↑ Parameters

$$e(x_i | y_i) = \frac{P(x_1, \dots, x_n | y_1, \dots, y_{n+1})}{\#(x_i, y_i)}$$

$$q(v_t | DT, JJ) = \frac{\gamma_1 \#(DT, JJ, VT)}{\#(DT, JJ)} + \frac{\gamma_2 \#(JJ, VT)}{\#(JJ)} + \frac{\gamma_3 \#(VT)}{\#(V)}$$

$\gamma_1 + \gamma_2 + \gamma_3 = 1$

$$e(\text{base} | VT) = \frac{\#(\text{base}, VT)}{\#(VT)}$$

$$e(x, y) = 0 \quad \forall y \quad \text{if } x \text{ is never seen in data}$$

Frequent words :  $\#(\text{words}) > 5$

Low frequency words : Put into corresponding classes

Tg: Profits : First word

, : Punctuation

$$e(\text{Profits} | NN) = e(\text{First word} | NN)$$

No. of possible sequences =  $|S|^n$

### Viterbi Algorithm:

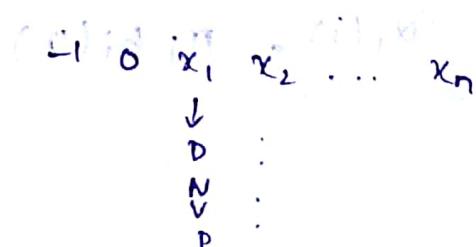
$$\delta(y_1, y_0, y_1, \dots, y_k) = \prod_{i=1}^k q(y_i | y_{i-2}, y_{i-1}) \cdot \prod_{i=1}^k e(x_i | y_i)$$

$\pi(k, u, v) = \text{Max. probability of a tag sequence ending in tags } u, v \text{ at position } k.$

Define  $S_k$  for  $k = 1 \dots n$  as the set of possible tags at position  $k$ .

$$S_{-1} = S_0 = \{\#\}$$

$$S_k = S \quad \forall k \in \{1 \dots n\}$$



$$\pi(k, u, v) = \max_{(y_F, y_0, \dots, y_k) : y_{k-1}=u, y_k=v} r(y_1, y_0, \dots, y_k)$$

$$\pi(\top, P, D) \quad S = \{D, N, V, P\}$$

w u v

The man saw the dog with the telescope

$$\pi(0, *, *) = 1$$

For any  $k \in \{1 \dots n\}$ , for any  $u \in S_{k-1}$  and  $v \in S_k$ ,

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$

Algo:

Set  $\pi(0, *, *) = 1$

$$S_{-1} = S_0 = \{\#\} ; S_k = S \text{ for } k = \{1 \dots n\}$$

for  $k = 1 \text{ to } n$

for  $u \in S_{k-1}, v \in S_k$

$$\pi(k, u, v) = \max_{w \in S_{k-2}} (\pi(k-1, w, u) \times q(v|w, u) \times e(x_k|v))$$

$$\text{return } \max_{u \in S_{n-1}, v \in S_n} (\pi(n, u, v) \times q(\text{STOP}|u, v))$$

$$(y_{n-1}, y_n) = \underset{(u, v)}{\operatorname{argmax}} (\pi(n, u, v) \times q(\text{STOP}/u, v))$$

$$\text{For } k = n-2, \dots, 1 \quad y_k = b_p(k+2, y_{k+1}, y_{k+2})$$

$$P(o|\gamma)$$

$$P(o|\gamma) = \sum_n P(o|q, \gamma) P(q|\gamma)$$

$$\alpha_t(i) = P(o_1, \dots, o_t, q_t = s_i | \gamma)$$

$$\alpha_t(i) = \pi_i b_i(o_t)$$

Words with frequency < 100 serve as unknown words.

F → ~~adjective~~ function words

R → Adverb

Chaotic systems

POS tagging ⇒ Sequence labeling task

$W = \{w_1, w_2, \dots, w_n\}$  ⇒ Observation sequence

$T = \{t_1, t_2, \dots, t_n\}$

$$P(T|W) = \frac{P(T) P(W|T)}{P(W)}$$

Markov assumption

Limited Horizon

Time invariance

$$P(x_t=i|x_{t-1}=j) = P(x_t=i|x_0=j)$$

Hidden Markov model (HMM) :

Set of states =  $S$

Alpha Output =  $V$

Transition prob.,  $A = \{a_{ij}\}_{i,j \in S}$ ,  $\pi = (A, B; \pi_0)$

Emission prob.,  $B = \{b_j(o_k)\}_{j \in S, k \in V}$

Initial state prob.

①  $P(O|\pi) \rightarrow$  Likelihood of a sequence (F-B)

② Given  $O$ , get best  $Q = \{q_1, \dots, q_n\}$  (state sequence) (Viterbi)

③ Adjust  $\pi$  to best maximize  $P(O|\pi)$  given  $O$ .

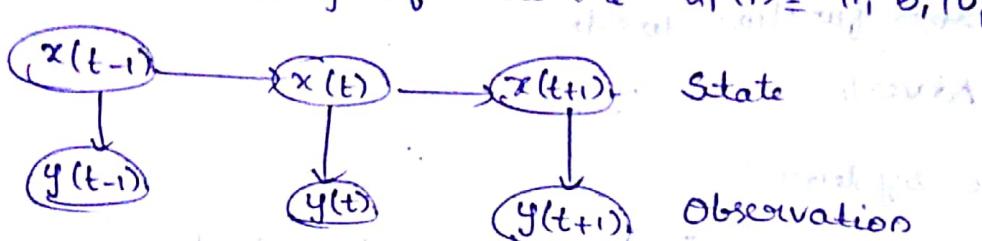
( $\pi, Q$ ) (Baum-Welch) (Re-estimation)

$$P(O|\pi) = \sum_Q P(O|Q, \pi) P(Q|\pi)$$

in terms of transition  
in terms of emission

Forward procedure: Prob. of being at state  $s_t$  given history of evidence  $\alpha_t(i) = \prod_j b_i(o_j)$

Belief state: Prob. of state at a certain time given history of evidence  $\alpha_t(i) = \prod_j b_i(o_j)$



$$\begin{aligned}\alpha_t(x_t) &= P(x_t | y_{1:t}) \\ &= \sum_{x_{t-1}} P(x_t, x_{t-1}, y_{1:t}) \\ &= \sum_{x_{t-1}} P(y_t | x_t, x_{t-1}, y_{1:t-1}) P(x_t | x_{t-1}, y_{1:t-1}) P(x_{t-1}) \\ &= P(y_t | x_t) \sum_{x_{t-1}} P(x_t | x_{t-1}) \cdot \alpha_{t-1}(x_{t-1}) \\ &\quad \downarrow \quad \downarrow \quad \downarrow \\ &\text{Emission} \quad \text{Transition}\end{aligned}$$

$$P(O|\lambda) = \sum_{i=1}^N \alpha_t(i) \stackrel{\alpha_t(j)}{\rightarrow} = \sum_{i=1}^N \alpha_t(i) \times a_{ij} \times b_j(o_{t+1})$$

Backward procedure:  $\beta_t(i) = P(O_{t+1} \dots O_T | q_{t+1} = s_i, \lambda)$

$$\begin{aligned}\beta_t(i) &= P(q_{t+1} = s_i | y_{1:t} | x_t) = \frac{P(x_t, y_{1:t})}{P(x_t)} \\ \beta_t(x_t) &= \sum_{x_{t+1}} \beta_{t+1}(x_{t+1}) \cdot P(x_{t+1} | x_t) = P(y_{t+1} | x_{t+1}) \\ \beta_t(i) &= \sum_{j=1}^N a_{ij} \cdot \beta_{t+1}(j) b_j(o_{t+1}) = \beta_t(i)\end{aligned}$$

Viterbi algorithm: Find most likely sequence of states  $s_i$  given  $O$  and  $\lambda$

Baum-Welch algo: Maximum likelihood estimate of  $\lambda$

$$\begin{aligned}\xi_t(i, j) &= P(q_t = s_i, q_{t+1} = s_j | O, \lambda) \\ &= \alpha_t(i) \cdot a_{ij} b_j(o_{t+1}) \cdot \beta_{t+1}(j) \cdot P(O|\lambda)\end{aligned}$$

$$P(O|\lambda) = \sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_t \gamma_t(i)}$$

$$\bar{b}_j(k) = \frac{\sum_{t=1}^T \gamma_t(j)}{\delta t \cdot O_t = v_k}$$

$$\bar{\pi}_i = \gamma_t(i) / \sum_t \gamma_t(j)$$

logistic regression estimate parameters of MNIST

Naive Bayes

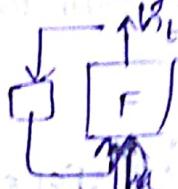
Discriminative

Normal neural networks don't scale well to images.  
So, we CNN's

3D input 3D output

Recurrent NN  $\rightarrow$  used for sequence of text (Sequential data)

In RNN, output of a hidden layer is fed back into the same layer.

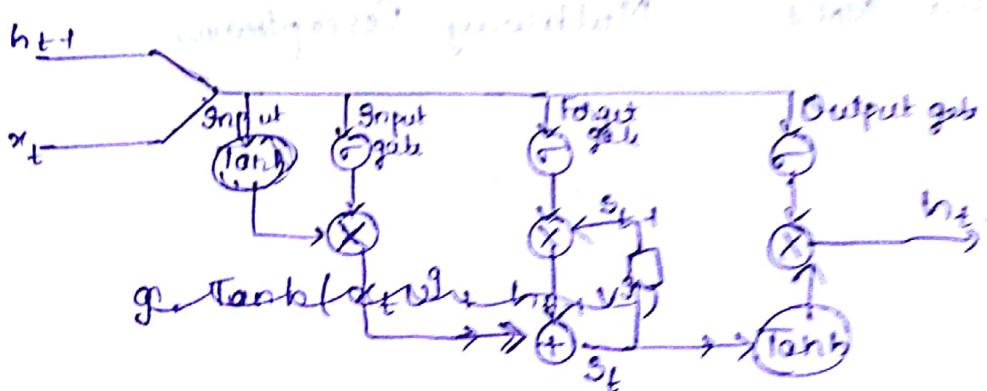


$$h_t = \sigma(Ux_t + Vh_{t-1})$$

$$= \sigma(Ux_t + V(\sigma(Ux_{t-1} + V(\sigma(Ux_{t-2}))))))$$

$$\frac{\partial L}{\partial U} = \frac{\partial L}{\partial \text{out}_3} \cdot \frac{\partial \text{out}_3}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial U}$$

Vanishing Gradient problem. So, we use LSTM.



$$s_t = s_{t-1} \circ f + g \circ i$$

$$g = \text{Tanh}(x_t v^g + h_{t-1} v^g + b^g) \quad \text{Squeezed input} \quad \left. \begin{array}{l} \text{Element wise} \\ \text{product} \end{array} \right\}$$

$$i = \sigma(b^i + x_t v^i + h_{t-1} v^i) \quad \text{Input gate} \quad \left. \begin{array}{l} \text{acts as} \\ \text{gate} \end{array} \right\}$$

$$f = \sigma(b^f + x_t v^f + h_{t-1} v^f) \quad s_{t-1} \circ f \quad \left. \begin{array}{l} \text{weight for input} \\ \text{acts as weight for} \end{array} \right\}$$

$s_{t+1}$

$$o^t = \sigma(b^o + x_t U^o + h_{t-1} V^o)$$

$$h_t = \tanh(s_t) \circ o$$

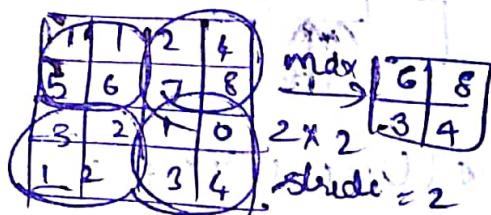
This way, LSTM controls what is "input, remember output" using gates.

$$\frac{\partial s_t}{\partial s_{t-1}} = f.$$

CNN's: Used for images

Convolution layer: Input matrix  $\times$  Convolved with filter matrix. Use ReLU on the obtained matrix.

Pooling layer: Downsampling.  
Max pooling  
Avg "



Fully connected layer: Flattens matrix into vector & feed into neural network

Feed forward NN's  $\Rightarrow$  Multiway Perceptrons

NP → Proper Noun

NP → Det Nominal

Nominal → Noun / Nominal Noun

S → NP VP

NP → Pronoun / Proper noun / Det Nominal

Nominal → Nominal Noun / Noun

NP → Verb / Verb NP / Verb NP PP / Verb PP

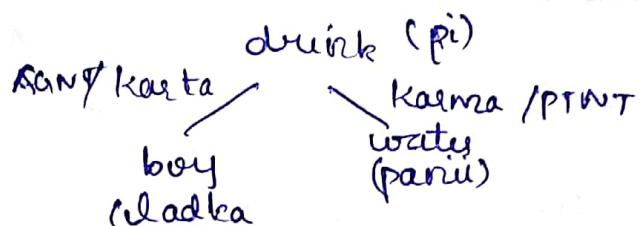
PP → Preposition NP.

S → Aux NP VP

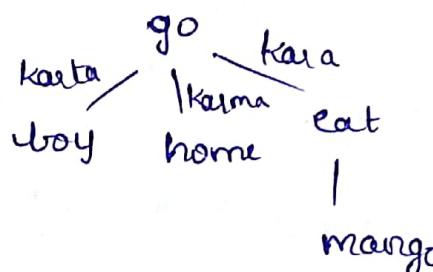
S → Wh - NP VP

S → Wh - NP Aux NP VP

Verb is head & noun is modifier



Temporal precedence (Verb - Verb modification)



NP → Det (adj)\* n  
      |  
      VP adj

VP → V [NP] [NP] PP\* [AdvP]

S → AdvP NP VP

AdvP → adv | VP

## Viterbi algorithm:

$$\textcircled{1} \quad \delta_t(i) = \pi_i b_i(o_t)$$

$$\psi_t(i) = 0$$

$$\textcircled{2} \quad \delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t)$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] \quad \leftarrow \text{backtrace}$$

$$\textcircled{3} \quad p^* = \max_{1 \leq i \leq N} [\delta_T(i)]$$

$$q^* = \operatorname{argmax}_{1 \leq i \leq N} [\delta_T(i)]$$

$$\textcircled{4} \quad q_t^* = \psi_{t+1}(q_{t+1}^*)$$

Chunks are non-recessive & non-exhaustive  
Constituents are recursive

Algorithms that try to learn  $P(y|x)$  directly  
(like logistic regression) are discriminative  
GDA, perceptron, etc.

$P(x|y) \& P(y) \Rightarrow$  Generative

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{d/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

$$\text{Cov}(x) = \Sigma$$

$\Sigma \uparrow$  compressed.

$P(y) = \phi^y (1-\phi)^{1-y} \Rightarrow$  Bernoulli,  $P(x|y) \Rightarrow$  Gaussian

$$\begin{aligned} \text{ll}(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m P(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) P(y^{(i)}; \phi) \end{aligned}$$

maximize ll

$$\phi = \frac{1}{m} \sum_{i=1}^m I\{y^{(i)} = 1\}$$

$$\mu_0 = \frac{\sum_{i=1}^m I\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m I\{y^{(i)} = 0\}}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

$$P(y=1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{e^{\theta^T x}}{1 + e^{\theta^T x}}$$

function of  $\phi, \Sigma, \mu_0, \mu_1$

Modeling assumptions correct  $\rightarrow$  GDA is correct

else LR is better

but if  $y \in \mathbb{R}$ , GDA

NB:

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) P(x_2 | y, x_1) P(x_3 | y, x_1, x_2) \dots P(x_n | y)$$

$$\phi_{i|y=1} = P(x_{i+1} = 1 | y = 1)$$

$$\phi_{i|y=0} = P(x_i = 1 | y = 0)$$

$$\phi_y = P(y = 1)$$

$$L(\phi_y, \phi_j|y=0, \phi_j|y=1) = \prod_{i=1}^m P(x^{(i)}, y^{(i)})$$

$$\phi_{j|y=1} = \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m I\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m I\{y^{(i)} = 0\}}$$

$$\phi_j = \frac{\sum_{i=1}^m I\{y^{(i)} = 0\}}{m}$$

$$\phi_y = \frac{\sum_{i=1}^m I\{y^{(i)} = 1\}}{m} \quad \text{Laplace smoothing}$$

$2^n$  possible values of  $x$

$2^n$  " "  $y$

$2^{n+1}$  parameters

$2^{n+1} - 1$  independent parameters

$2(2^n - 1) \rightarrow 2n$  bits

Naive Bayes

Bayes

$$\phi_{k|y=1} = \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} I\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{\sum_{i=1}^m I\{y^{(i)} = 1\} n_i + 1}$$

$$J(\omega) = \sum_{\ell=1}^L Y^\ell \ln P(Y^\ell = 1 | X^\ell, \omega) + (1 - Y^\ell) \ln P(Y^\ell = 0 | X^\ell, \omega)$$

$$\begin{aligned} \textcircled{1} \quad P(Y=1|X) &= \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)} \\ &= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \\ &= \frac{1}{1 + \exp\left(\ln\left(\frac{1-\pi}{\pi}\right) + \sum_i \ln \frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)} \\ &\leq \sum_i \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}} \\ &\leq \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}} \\ &\leq \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \end{aligned}$$

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

$$\textcircled{2} \quad \omega \leftarrow \operatorname{argmax}_{\omega} \prod_{\ell=1}^L P(Y^\ell = 1 | X^\ell, \omega)$$

$$\omega \leftarrow \operatorname{argmax}_{\omega} \sum_{\ell} \ln P(Y^\ell = 1 | X^\ell, \omega)$$

$$\begin{aligned} J(\omega) &= \sum_{\ell=1}^L Y^\ell \ln \frac{P(Y^\ell = 1 | X^\ell, \omega)}{P(Y^\ell = 0 | X^\ell, \omega)} + \ln P(Y^\ell = 0 | X^\ell, \omega) \\ &= \sum_{\ell=1}^L Y^\ell (w_0 + \sum_i w_i x_i^\ell) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^\ell)) \end{aligned}$$

$$\frac{\partial J(\omega)}{\partial w_i} = \sum_{\ell=1}^L x_i^\ell (Y^\ell - \hat{P}(Y^\ell = 1 | X^\ell, \omega))$$

$$\textcircled{3} \quad \omega \leftarrow \operatorname{argmax}_{\omega} \sum_{\ell} \ln P(Y^\ell = 1 | X^\ell, \omega) - \frac{\lambda}{2} \|\omega\|^2$$

$$\sim \sum_{\ell} \ln P(Y^\ell = 1 | X^\ell, \omega) + \ln P(\omega)$$

↓  
mean GD.

$$\textcircled{1} \quad w_{ji} \leftarrow w_{ji} + \gamma \sum_i^L (s_i^h (\delta(y^d - y_i) - \hat{P}(y^d = y_i | x^o, w)) - \eta \lambda w_{ji})$$

$$a_i \quad s_i \quad w_{ji} \leftarrow s_j^h \quad , \quad w_{ij} \leftarrow s_i^h$$

$$E(m) = \frac{1}{2} \sum_{k=1}^K [b_k(m) - f_k^o(x_k^o)]^2$$

$$x_k^o = \sum_{j=1}^J w_{kj} \cdot s_j^h$$

$$s_j^h = f_j^h(x_j^h)$$

$$x_j^h = \sum_{i=1}^I w_{ji} s_i$$

$$s_i = x_i = a_i(m)$$

$$\begin{aligned} \frac{\partial E(m)}{\partial w_{kj}} &= \frac{1}{2} \cdot \frac{\partial}{\partial w_{kj}} [b_k - f_k^o(\sum_j w_{kj} s_j^h)]^2 \\ &= -(b_k - f_k^o) \times f_k^o \cdot s_j^h \end{aligned}$$

$$\frac{\partial E(m)}{\partial w_{ji}} = - \sum_{k=1}^K (b_k - f_k^o) \cdot \frac{\partial}{\partial w_{ji}} f_k^o(\sum_j w_{kj} s_j^h)$$

$$= - \sum_{k=1}^K (b_k - f_k^o) \cdot f_k^o \cdot w_{kj} \cdot \frac{\partial s_j^h}{\partial w_{ji}}$$

$$\frac{\partial s_j^h}{\partial w_{ji}} = f_j^h \cdot s_i$$

Short term memory  $\Leftarrow$  Activation dynamics  
 Long term " "  $\Leftarrow$  Synaptic weights

$S \rightarrow NP VP$

$NP \rightarrow \text{det (adj)* noun}$

$VP \rightarrow V\bar{V}[NP][NP] PP^*$

$PP \rightarrow \text{Pprep NP}$

$\boxed{PP \rightarrow NP \text{ ago}}$

verb following preposition

$VG \rightarrow (\text{Verb})^* V$

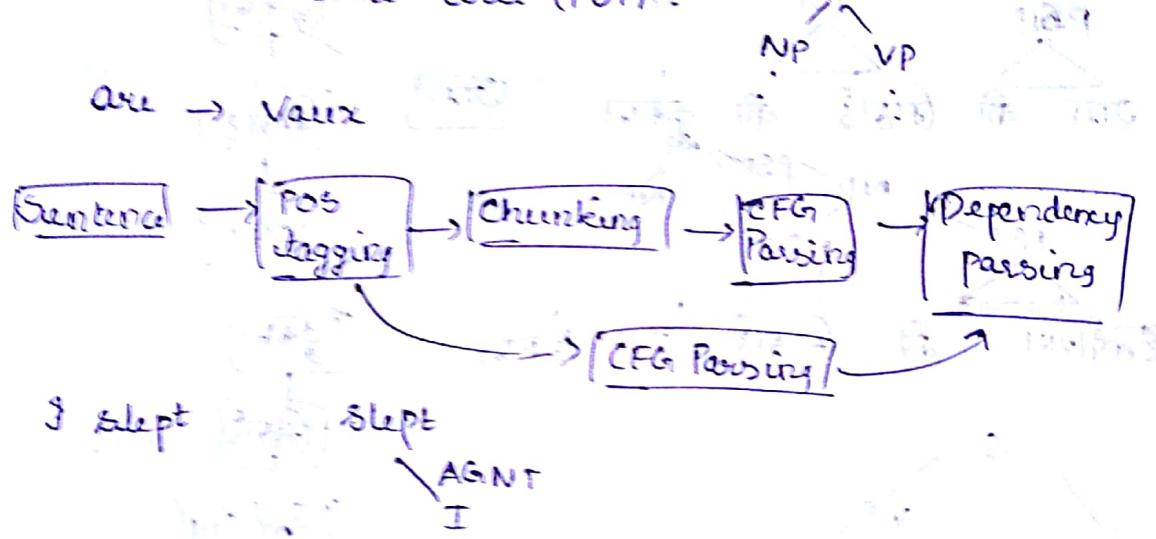
$NP \rightarrow NP PP$

Dependency tree: AGNT, PATNT, RECIP,

head

modifiers

Phrase structure tree (PST):

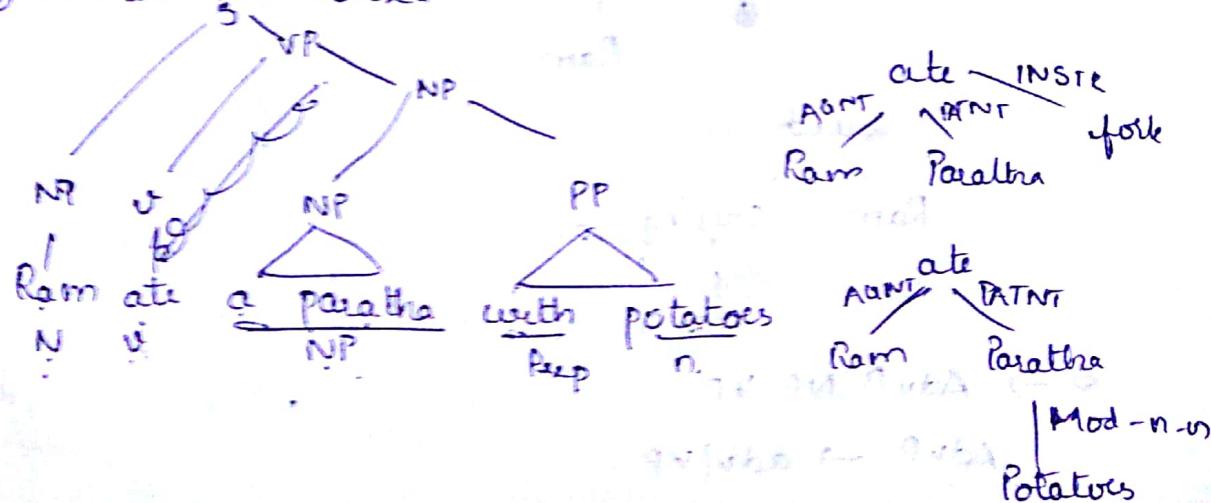


Tense can't be learnt by animals

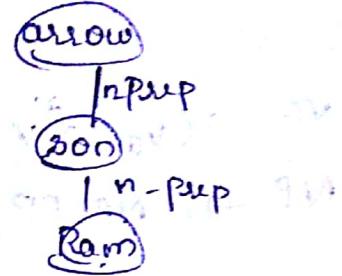
① Statistical based - Rules, grammar, semantics

② Statistical learning features

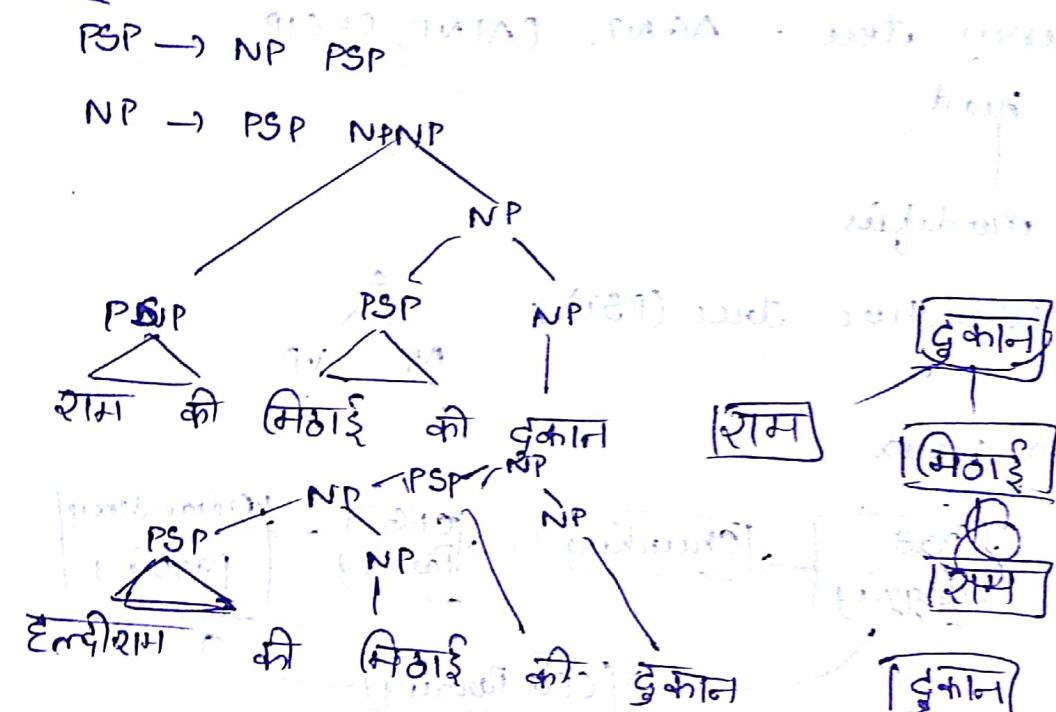
③ Neural net based



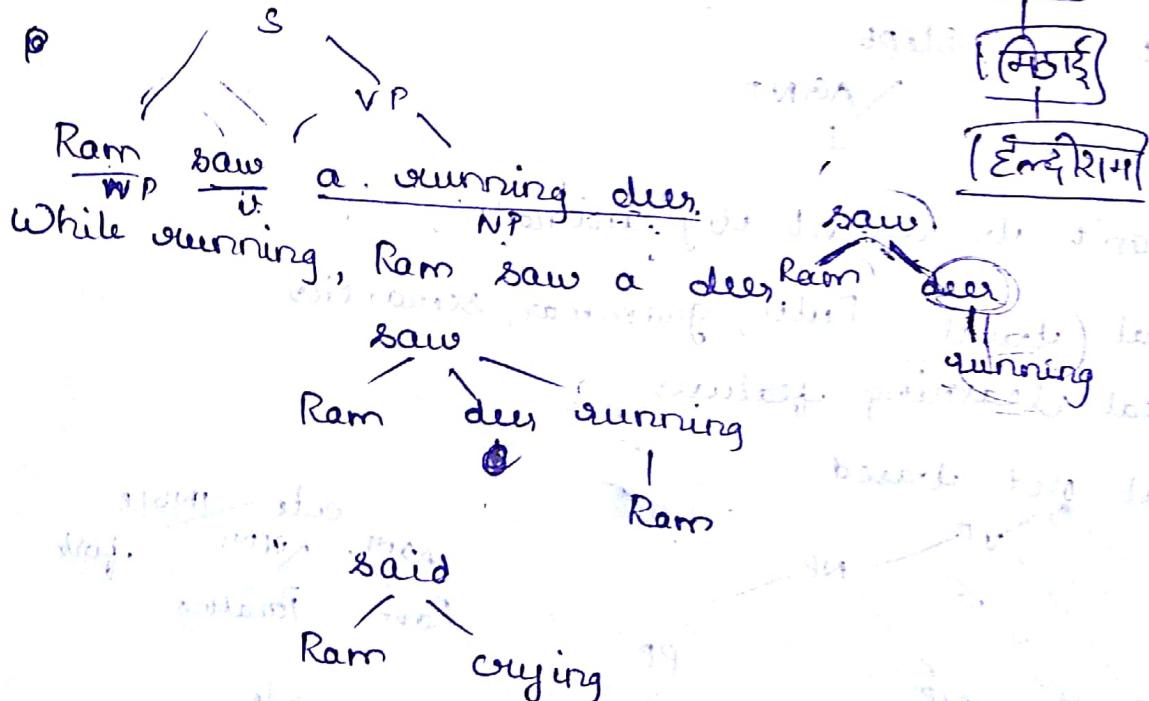
## arrow of son of Ram



Hindi:

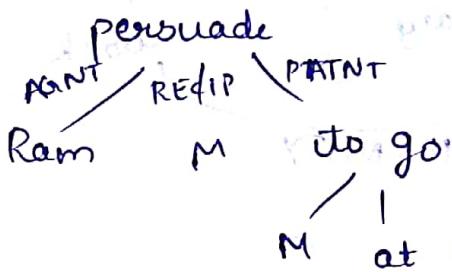
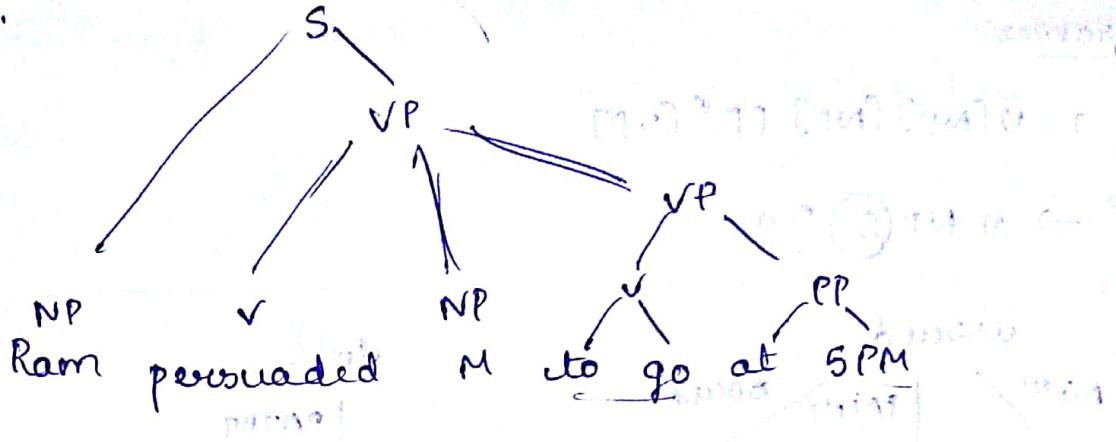


②



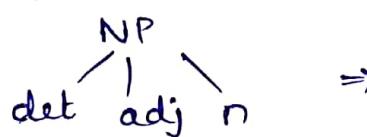
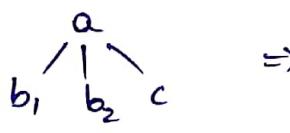
S → AdvP NP VP

AdvP → adv / vp



**persuade** - **RECIP** control

**promise** - **AANT**



$\overline{\text{NP}} \Rightarrow \overline{\text{PATNT}}$

If  $\overline{\text{VP}}$  is present,  $\overline{\text{NP}} \Rightarrow \overline{\text{PATNT}}$

without  $\overline{\text{VP}}$  many  $\overline{\text{NP}} \Rightarrow \overline{\text{RECIP}}$

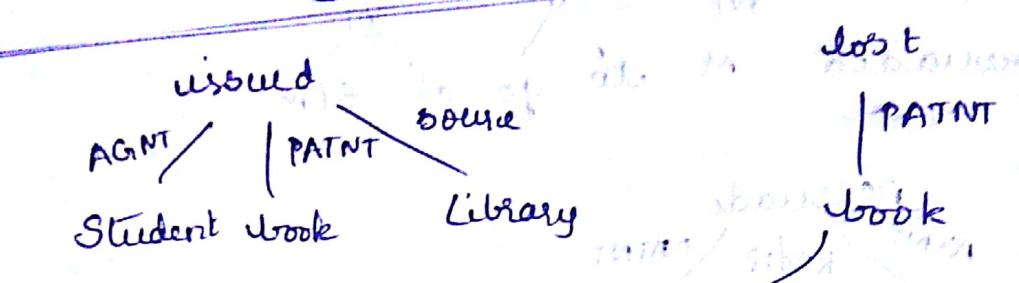
Active - Passive  $\Rightarrow$  MM tree same

PST different

Verb frames:

$VP \rightarrow V[NP][NP] PP^*[VP]$

$VP \rightarrow V NP S ??$



which  $\in$  IDENTITY

~~AGNT~~  $\swarrow$  ~~S~~  $\searrow$   
~~PATNT-JNV~~

issued  
Student book library

$\textcircled{1}$   $S \rightarrow NP VP$   
 $\textcircled{2}$   $NP \rightarrow NP S\text{-relative clause}$   
 (rel-cl)  
 $S\text{-rel-cl} \rightarrow Wh S'$

$S' \rightarrow NP VP$

$VP \rightarrow V[NP'][NP'] PPX$

$NP' \rightarrow \emptyset$

The book which the student issued from the library is lost

$NP$   $\text{R-cl}$   $NP$

$S$

$NP$   $\text{PATNT}$   $NP$   $PP$   $NP$

The student issued the book from the library

$NP$   $NP$   $NP$

That book is lost

## Object relativization

Subject

Shift rule A

(S → S' relativized with)

verb moving without a fo-circumlocution after moving verb

$S \leftarrow [NP] VP$

agreement fails if

$S \leftarrow Wh \& Aux S'$

verb moving, and for maintaining adverbial wh-dependency

at the same position with the same agreement node

24/12/2018

After Mid2

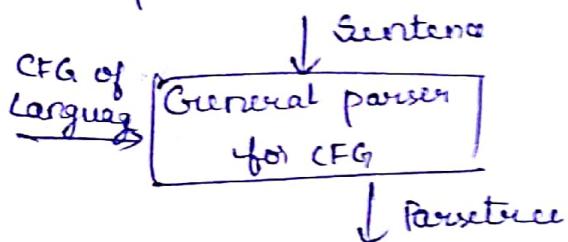
(RS)

Context-free grammar (CFG)

Parser gives the analysis of a sentence given CFG of that language.

After the ~~introduction~~ invention of CFG's, parsers for each language needn't be written separately.

General parser is built.



The child plucked a flower  
det n v det n.

R1 :  $S \rightarrow NP VP$

R2 :  $NP \rightarrow \text{det } n$

R3 :  $PP \rightarrow \text{prep } NP$

R4a :  $VP \rightarrow v$

R4b :  $VP \rightarrow v \text{ NP}$

R4c :  $VP \rightarrow v \text{ NP } NP$

R4d :  $VP \rightarrow v \text{ NP } NP \text{ PP}^*$

:  
Bottom-up Parsing algorithm:

The child plucked a flower  
det n v det n

OPEN	Remaining Sentence	Remaining Rules	
-	1	ALL	
det	2	ALL	Shift
det n	3	ALL	Shift
NP	3	R3	Reduce
NP v	4	<del>R4b</del>	Shift
NP VP	4	<del>R4b</del>	Reduce
NP VP	4	<del>R2</del>	Reduce
S	4	<del>R2</del>	POP X

NP V det	5	ALL - Shift
NP V det n	6	ALL - Shift
NP, V, NP	6	R3 - Reduce
NP VP	6	R4C - Reduce
S	6	R2, Success ✓

Probabilistic CFG:

Rules, Probabilities

$$A \rightarrow \alpha | P_1$$

$$B \rightarrow \beta | P_2$$

$$P(A \rightarrow \alpha | A) = P_1$$

Introduce probability with lexical background.

LHS + Head of that phrase.

Dependency parser  $\Rightarrow$  NLP

Valency - demand charts

sleep

AGNT

verb word (noun)

DET

DET

pluck

AGNT

verb

AV

AV

## C-parsers      CFG Parsing

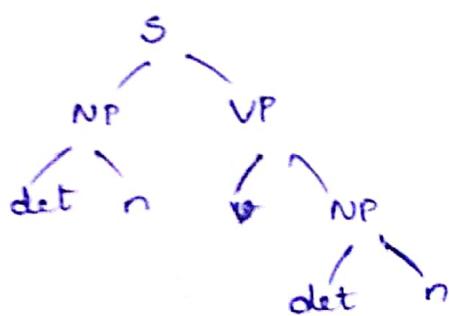
Bottom-up parsers: Generally fast, Less recursive.

Top-down parsers

The boy plucked the flower  
det n      v      det n

Stack	Open	Remaining sentence	Remaining rules
1	nil	1	ALL
1'	det	2	ALL
1''	det n	3	R3
2	NP	3	ALL
2'	NP V	4	R4b
3	NP VP	4	R2
4 x POP	S	4	ALL
3 x POP	NP VP	5	ALL
2'''	NP V det	5	ALL
2'''	NP V det n	6	ALL R3
3	NP V NP	6	R3 R4c
4	NP VP	6	R4c R2
5	S	6	<del>ALL</del> ← Record success

Record:



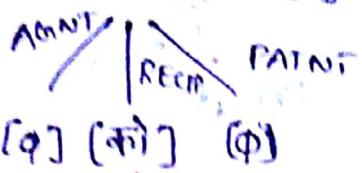
PCFG

## Dependency parsing:

मरी गली को लिखा द्यति तो Root  
न न postprep न वाक्यावृत्ति

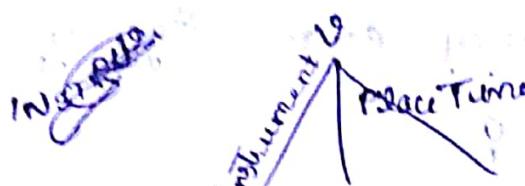
constraint graph

Mandatory  $\rightarrow [S-X]$



Di-branchable  
[A] [B] [C]

Optional :



[Q1] [Q2] [Q3]

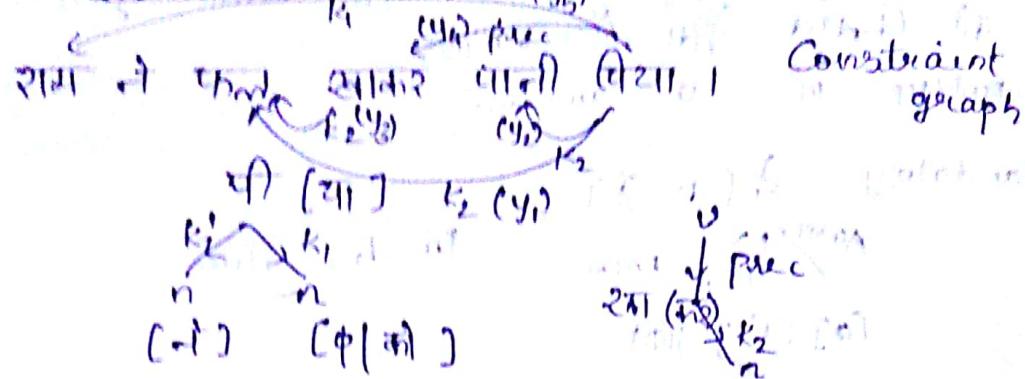
Solution graph (parse tree):

1. Not more than one incoming edge.
2. For each outgoing edges out of a node, there should be exactly one labelled edge corresponding to mandatory edges.
3. At most one optional edge.

A solution graph is a subgraph of constraint graph satisfying ①, ② & ③ at the same time

most branched-most (constraint graph) approach

## Parsing using Integer Programming:



A variable for every edge. E.g.,  $y_1$ ,

$$y_i = 1 \text{ or } 0$$

Not included in solution graph

Constraints:  $y_5 = 1$  {Mandatory edges initial}

$$y_1 + y_2 = 1$$

{ $y_1$  is mandatory and  $y_2$  is optional}

$$y_3 \leq 1$$
 {Optional}

$y_5 = 1$  {Incoming edges}

$$y_1 + y_3 = 1$$

$\Rightarrow \text{का}(y_1)$ :  $k_1$  mandatory in  $n[\phi]$

$$k_2 \quad " \quad n[\phi | k_1]$$

$\text{का}(y_2)$ :  $k_1$  mandatory in  $\phi$

$k_2$  optional in  $n[\phi | k_1]$

पौरा mandatory in  $\phi$

Treebank (Penn treebank) Phrase-structured trees

10,000 sentences  $\rightarrow$  tree

MALT Parser - Naive

{ Fundamentally }

MST Parser (Max. Spanning tree)

{ different }

$$\sum x_{ik} = 1$$

Global maximization

Local optimization (Uses Shift-Reduce algorithm)

Internally uses SVM's

## Feature bundle / Feature structure:

NP  
 a NP sheep  
 Attribute Value  
 children - n - [num plural]  
 root, child

a : [num singular]  
 specificity indefinite

sheep : [root sheep]

NP  $\rightarrow$  det n

NP.num  $\Leftarrow$  n.num

NP.num = det.num

NP.head  $\Leftarrow$  n.root

NP.spec  $\Leftarrow$  det.spec

Feature constraints

S  $\rightarrow$  NP VP

S.num = VP.num } Subject VP agreement

NP.num = VP.num

VP  $\rightarrow$  v NP

VP.num = v.num

VP.gender = v.gender

$\Rightarrow FS(VP) = FS(v)$

Feature structure

VP.head = v.root.

## Dependency structure:

S  $\rightarrow$  NP VP

S.AGNT = FS(NP)

S.PATNT = VP.PATNT

Feature structure

VP  $\rightarrow$  v NP

S.RECIP = VP.RECIP

v.Trans = Transitive verb

S.head = VP.head.

v.Trans = Yes/No | Distrans.

VP.PATNT = FS(NP)

VP.num = v.num

Lexical Function Grammar (LFG)  $\rightarrow$  Head-driven phrase structure grammar

Machine translation

HPSG  $\Rightarrow$  Head-driven Phrase Structured Grammar

Head-driven phrase structure grammar  
is a type of generative grammar that uses heads to determine the syntactic structure of a phrase.

Head-driven phrase structure grammar  
is a type of generative grammar that uses heads to determine the syntactic structure of a phrase.

3/8/2018

3/8/20  
The task of separating tokens in a given input sequence  
1. if any punctuation marks is counted as a token

Morphological	Talib-aan
↑	v. v
Lexical level	Student Plural

## Type, Token

Vocabulary word

High type ratio  $\Rightarrow$  Higher  
tokens  $\Rightarrow$  inflection

## Morphological change?

Inflection: Change in suffixes (can be added one after another).

## Tokenization challenges:

① Hyphenation : state-of-the-art, hard-working

Q Number : 9.8 Are generally we separate 9 and 8 but here no.

### ③ Dates

## ④ Abbreviations

## ⑤ Punctuations.

© URI's

## ⑦ Sentencifications

## Regular expressions

## N-grams and Zipf's Law:

Markov assumption : Word is affected only by its

"Within local context" (last few words).

## Uruguay

$$\frac{c(w_i)}{N} p(w_i | w_{i-1}, w_{i-2}) =$$

## Bigrams

$$\frac{c(w_i, w_i)}{c(w_{j-1})} \quad \frac{c(w_i, w_{i-1}, w_{i-2})}{c(w_{i-1}, w_{i-2})}$$

## Reliability vs Discrimination

Class estimation  $\langle N, v \rangle$   $\langle v, N \rangle$

1/1/2018

## lecture-3 (MM)

$$P(o_1, o_2, \dots, o_n) = P(o_1 | o_{k+1}, \dots, o_n) P(o_2 | o_{k+1}, \dots, o_n) \dots P(o_n | o_{k+1}, \dots, o_n)$$

Markov assumption:

$$\approx P(o_i) \prod P(o_i | o_{i-1})$$

$$\approx P(o_i | o_0) \prod P(o_i | o_{i-1})$$

$$\approx \pi_i P(o_i | o_{i-1})$$

Perplexity & Entropy

Smoothing & Backoff

Uniform Distribution : Highest entropy

$$\text{entropy } H(X) = -\sum_{x=1}^n P(x) \log_2 P(x)$$

$$\text{PP}(X) = 2^{H(X)}$$

Perplexity

Zipf's law: The product of frequency of words & their rank is approximately constant.



Grammatical content

$$X = UDV^T$$

$n \times k$  matrix  $k \times k$   $\xrightarrow{\text{Diagonal}}$   
 $\downarrow$  Orthogonal

Singular value decomposition

Softmax function

$$\frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\frac{d \sigma_x}{dx} = \frac{-1}{(1 - \sigma_x)^2}$$

A context predicts a word.

## Lecture - 6 (NLM)

Spelling correction (18<sup>th</sup>)

Project → Grammar error detection

N-grams and Smoothing:

- ① Laplace smoothing: Initialize the matrix with 1's instead of 0's

$$P_{\text{lap}}(w_1, w_2, \dots, w_n | w_0) = \frac{c(w_1, \dots, w_n, w_0) + 1}{N + B}$$

$N$  = No. of seen n-grams

$B$  = No. of possible n-grams

- ② Lidstone's law:

$$P_{\text{lid}}(w_1, w_2, \dots, w_n | w_0) = \frac{c(w_1, \dots, w_n, w_0) + \gamma}{N + \gamma B}$$

- ③ Held-out estimation:

$$\hat{P}(w_1, \dots, w_n) = \frac{T_i}{N_i N}$$

$$P(w_i, w_0) = \frac{T_i}{N} = \frac{T_i + N_i}{N_i + N} = \frac{T_i}{N_i N}$$

- ④ Cross-over estimation

$$P(w_1, \dots, w_n) = \frac{(T_i^0 + T_i^{10})}{(N_i^0 + N_i^{10}) N}$$

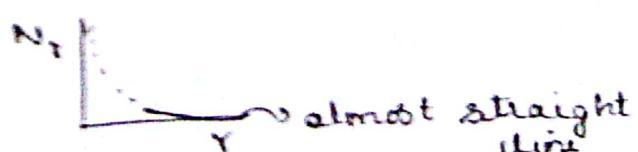
- ⑤ Good-Turing smoothing

$$\gamma^* = \frac{(\gamma+1) E(N_{i+1})}{E(N_i)}$$

$$P(w_1, \dots, w_n) = \frac{\gamma^*}{N}$$

$\gamma-1, \gamma, \gamma+1$  are all similar

$\frac{N_{i+1}}{N_i}$  tells us a lot (Sudden ↑ or ↓)



Back-off is "using lower order N-grams for estimating higher orders"

Katz back-off

Linear interpolation

language Modeling : Predicts the probability with which a sentence belongs to a language.

24/12/2018

## Lecture - 5 (Ganesh)

### Noise Bayes

Functional words: Highest frequent words (top 100)

Token: No. of words

Type: No. of unique words

$$P(\text{word} | \text{class}) = \frac{\# \text{ Word in class}}{\# \text{ Tokens in class}}$$

$$\log P(\text{sentence} | \text{class}) = \sum_{i=1}^N \log P(w_i | \text{class})$$

$$P(\text{class} | \text{sentence}) = \frac{P(\text{sentence} | \text{class}) P(\text{class})}{P(\text{sentence})}$$

$$P(\text{senti} | \text{senti}) = P(\text{senti} | \text{senti}) P(\text{senti})$$

$$\text{Example: } \sum \log(P(w_i | \text{senti})) + \log P(\text{senti})$$

\* Tokenization as classification problem

Word boundary, Sentence boundary, char

(, ; ! ?)

28/8/2018

## Lecture - 6 (MN)

Good turing :  $c^* = \frac{(c+1) N_{c+1}}{N_c}$

Church & Gale:  $j_{ii} \Rightarrow$  Join if independent

$N_c$  can be zero. Then?

So, put  $c-5$  to  $c+5$  in a bucket.

So,  $N_c$  contains  $\{c-5, \dots, c, \dots, c+5\}$

(Read this paper) Great one.

Simple good turing:  $N_d = \alpha g^b$

Pabsolute discounting ( $w_i | w_{i-1}$ ) =  $\frac{c(w_{i-1}, w_i) - d}{c(w_{i-1})} + \delta(w_{i-1}) P_{\text{M}}$

Kneser-Ney smoothing:

Pcontinuation( $w_i$ )  $\propto |\{w_{i-1} : c(w_{i-1}, w_i) > 0\}|$

$$P_{\text{cont}}(w_i) = \frac{|\{w_{i-1} : c(w_{i-1}, w_i) > 0\}|}{\sum |\{w_{i-1}' : c(w_{i-1}', w_i) > 0\}|}$$

$$(it's not quite p_{\text{cont}}(w_i), it's p_{\text{cont}}(w_i | w_{i-1}))$$

$P_{\text{KN}}(w_i | w_{i-1}) = \max(c(w_{i-1}, w_i) - d, 0) + \delta(w_{i-1}) P_{\text{cont}}(w_i)$

$$\delta(w_{i-1}) = \frac{d}{c(w_{i-1})} |\{w_i : c(w_{i-1}, w_i) > 0\}|$$

Witten-Bell smoothing:

$T(w_x)$  = No. of seen  $n$ -gram types starting with  $w_x$ .

$Z(w_x)$  = No. of unseen  $n$ -gram types starting with  $w_x$ .

$N(w_x)$  = No. of  $n$ -gram tokens starting with  $w_x$

$$P^*(w_i | w_x) = \frac{T(w_x) * 1}{Z(w_x) (N(w_x) + T(w_x))} \Rightarrow \text{Prob. of unseen unigrams}$$

$$P^*(w_i | w_x) = \frac{c(w_x, w_i)}{N(w_x) + T(w_x)} \Rightarrow \text{Prob. of seen unigrams}$$

## Lectures - 7

Learn hierarchy of languages (Regular vs Context free)  
 → Context sensitive ..... → Natural language

### POS Tagging & chunking:

Statistical model

Deep learning model

- ① Distributions      (e.g., n-gram, Representation)
- ② Assumptions      (e.g., local architecture)

Use statistical & DL together. Not provided together often get better results

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

$$= \prod_i P(t_i | w_1, w_2, \dots, w_n) \quad (\text{Assumption:})$$

Tags are independent of each other

$$= \prod_i P(t_i | w_i) \quad (\text{Assumption: Tag of a word is independent of all other words})$$

Parameters  $A, B, \pi$  are fixed for probabilities  
with respect to time.

Forward procedure

Backward procedure.

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t)$$

$$B_t(i) = \sum_{j=1}^N \beta_{t+1}(j) a_{ij} b_j(o_{t+1})$$

Expectation  $\rightarrow$  First moment

$$E_t(i, j) = \frac{\alpha_t(i) \cdot \beta_{t+1}(j) \cdot a_{ij} b_j(o_{t+1})}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}$$

$\sum_t E_t(i, j) \Rightarrow$  How many times  $j$  is seen.  
when  $i$  is seen.

$\sum_t \sum_j E_t(i, j) \Rightarrow$  How many times  $i$  is seen.

$$(i.e.) \gamma_t(i) = \sum_j E_t(i, j)$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} E_t(i, j)}{\sum_t \gamma_t(i)}$$

$$\bar{b}_i(k) = \frac{\sum_{t=1}^T \gamma_t(i)}{\sum_{t=1}^T \gamma_t(j)}$$

$$\bar{\pi}_i = \gamma_i(i)$$

## Sentiment Analysis

N-hot encoding

NB model for SA

Log Reg for SA

Hypothesis

$\pi_{sentiment} = \text{softmax}(\theta^T x)$

$\theta^T x = \sum_{i=1}^n \theta_i x_i + b$

Generative vs Discriminative

Sentiment analysis  $\begin{cases} +, -, \text{neutral}, \text{surprise} \\ +, -, \text{neutral}, \text{negative}, \text{positive} \end{cases}$

N-hot encoding : Represent sentences as 1-hot vector

(and a fixed number of words)

$$P(\text{sentiment} | \text{sentence}) = \frac{P(\text{Sentence} | \text{Sentiment}) P(\text{Sentiment})}{P(\text{Sentence})}$$

$$P(+|x) = h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\theta = (\theta_0, \theta_1, \dots, \theta_n)$$

21/10/2018

## Word2vec

Continuous bag of words (CBOW):

Skip gram (SG): Not possible in practical. Only at max 1 word can be predicted

Neural networks don't work with discrete data.

Use Real-valued vectors:

Feed-forward NN's

Recurrent NN's

Why

Softmax? Small changes in val will bring huge difference in  $e^{\text{val}}$ .

Link it back to viterbi

$$h_i = \sigma(w_{L,i} + u_{x,i})$$

$$\text{softmax}(x) = e^{x_i} / \sum e^{x_j}$$

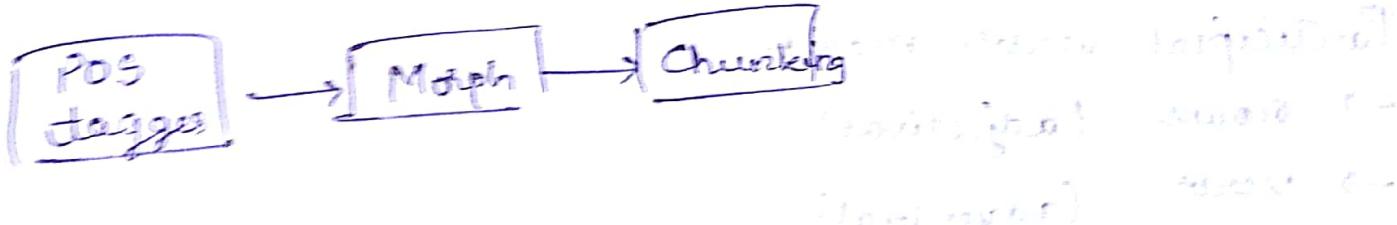
Brevity  $\leftrightarrow$  Ambiguity

Natural language: Easy to convey thoughts  
Just say half, listeners will understand whole  
meaning. Whereas in programming language,  
you need to define say things at the  
finest level.

Rule-based

Statistical methods

Artificial neural networks



26/10/2018

AGNT  $\Rightarrow$  Agent

R said that Mohan is owing  
 AGNT AGENT Sentence  
 Main verb

R  $\xrightarrow{\text{Agent}}$  said PATNT  $\rightarrow$  Patient  
 owing | AGNT M

Why only one tree for a sentence? Why not 2/3?  
 Sentence is a one connected structure

Participial verbs modify  
 → nouns (adjectival)  
 → verbs (adverbial)

NP  $\Rightarrow$  det (adj)\* n.

NP  $\Rightarrow$  det (Adj Phrase)\*  
 Adj Phrase  $\rightarrow$  Adj  
 $\rightarrow$  VP

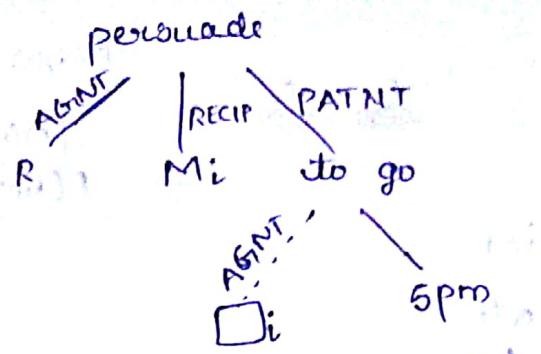
VP  $\Rightarrow$  v [NP] [NP] PP\* [adv]  
 [AdvP]

PL  $\Rightarrow$  Possessor  
 AdvP  $\rightarrow$  adv | VP

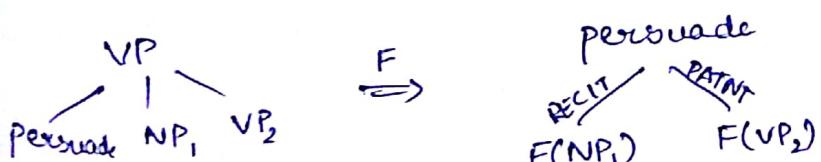
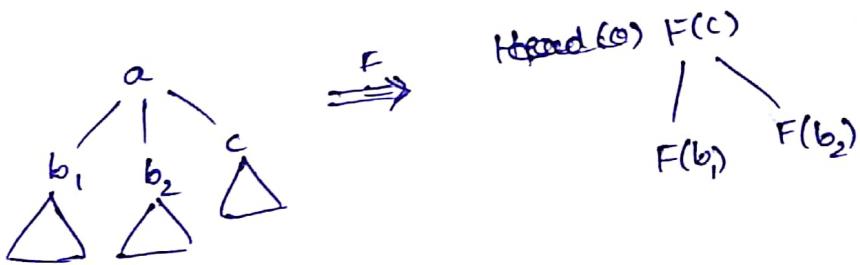
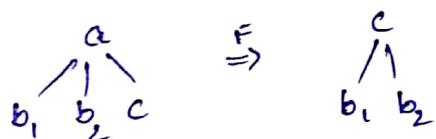
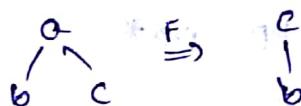
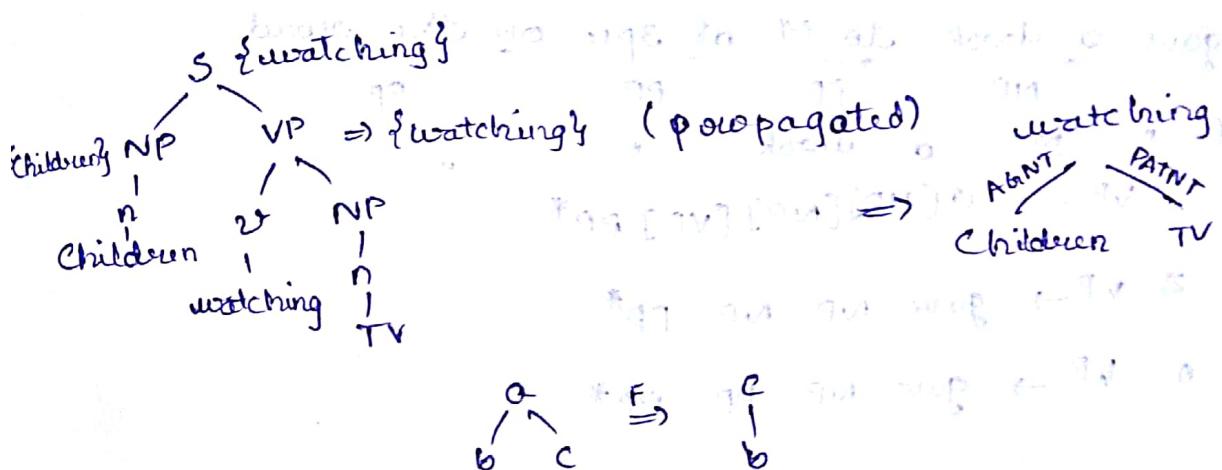
R persuaded M  $\xrightarrow{\text{to go}}$  at 5 PM  
 R persuaded M to go

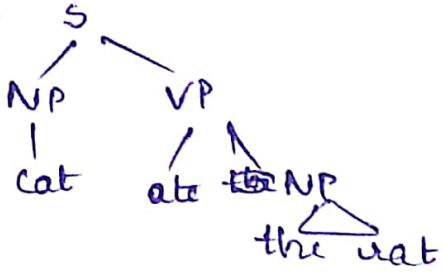
VP  $\rightarrow$  v [NP] [NP] [VP] PP\*

persuade  
 R / M / go



persuade = obj/RECIP structure  
 promise = AGNT control





$$F(VP) = \text{ate}$$

PATNT  
cat

5

1) R gave a book to M at 3pm on the road

2) R " M a book "

VP → [NP][NP][VP] PP\*

VP → gave NP NP PP\*

VP → gave NP PP PP\*

The child plucked the growing flower

The child plucked the flower growing in the garden

The child plucked the flower walking in the garden

Relative clauses

Verb w