



Finite state Transducers (FSTs)

- FSA

Accept/Reject

| spell checking

- Mapping b/w upper language and lower language
- If successful
→ returns the upper side symbols on the path as result.
- If unsuccessful
→ returns nothing (reject)
- FST have some input labels.
- FST have input : output pairs on labels.

Q → finite set of states

Σ

Δ → Output alphabet

$q_0 \in Q$

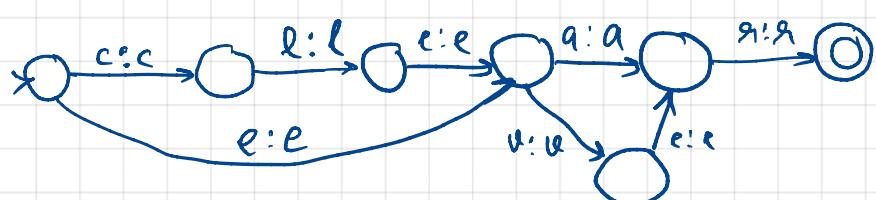
$F \subseteq Q$

$\delta(q, w) : Q \times \Sigma^* \rightarrow 2^Q$

$\delta(q, w) : Q \times \Sigma^* \rightarrow 2^{\Delta^*}$

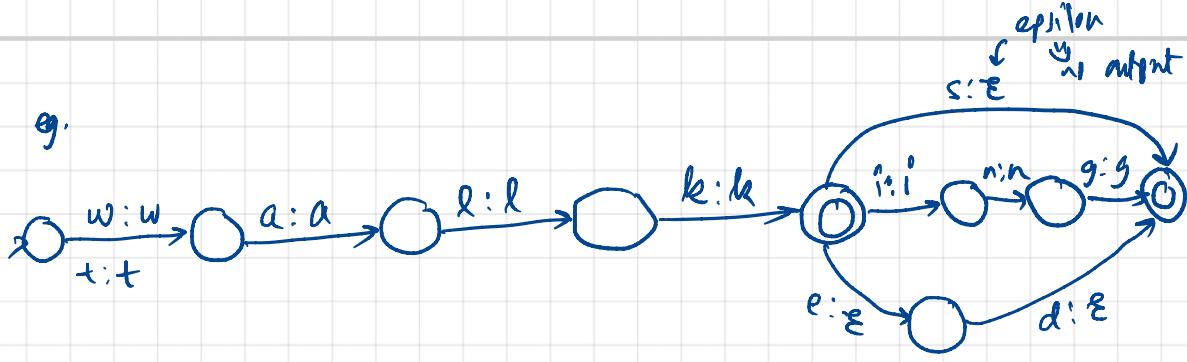
surface form ← input → cats
 lexical form ← output → cat + N + P

e.g.



input → clear , output → clear

input → ever , output → ever



can map

talk, talks, talked, talking to talk

and vice versa.

Operations on FST:

1. Inversion

$$T \rightarrow T^{-1}$$

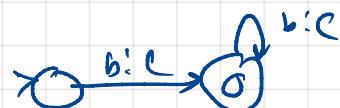
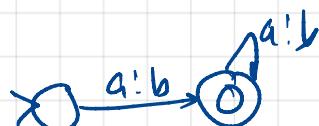
$$T = \{(a, a_1), (a, a_2), (b, b_1), (c, a)\}$$

$$T^{-1} = \{(a_1, a), (a_2, a), (b_1, b), (a, c)\}$$

The inversion of transducer $T \rightarrow T^{-1}$

simply switches the input and output labels

2. Composition



$$T_1 \circ T_2(x) = T_1(T_2(x))$$

Probabilistic language modeling

Please turn your homework ?
 in the ✓ ✗

N-gram model

↳ predicts next word from previous N-1 words

An N-gram is a N-token sequence of words.

Bigram is a sequence of 2 words.

("please turn", "turn your", "your homework")

- Smoothing
- Add One
 - Witten Bell
 - Good Turing
 - Interpolation
 - Back off

Witten Bell shared probability mass

$$\hookrightarrow \text{unseen} \rightarrow \frac{T(c)}{N(c) + T(c)}$$

$$\text{seen} \rightarrow \frac{N(c)}{N(c) + T(c)}$$

} Bigrams starting with c, for example in slide

Unsmoothing ??

→ To get counts from probabilities :

↳ known from Witten Bell

$$P(w_2|w_1) = \frac{c(w_2|w_1)}{N(w_1)} \leftarrow \text{needed}$$

$$c(w_2|w_1) = P(w_2|w_1) * N(w_1)$$

$$\uparrow \quad = \frac{1}{z(w_1)} \cdot \frac{T(w_1)}{N(w_1) + T(w_1)} \cdot N(w_1)$$

counts of
unseen types

$$= \frac{T(w_1)}{z(w_1)} \cdot \frac{N(w_1)}{N(w_1) + T(w_1)}$$

Good Turing Discounting

N_c grams which occur once in the corpus

use the frequency of singletons to reestimate the frequency of zero-count bigrams, trigrams, ...

N_c = Number of N -grams that occur c times
("frequency of frequency")

N_0 \rightarrow Number of bigrams with count 0.

N_1 \rightarrow Number of bigrams with count 1.

N_c can be considered as bin \rightarrow stores N -grams occurring with frequency c .

$$N_c = \sum_{x: \text{count}(x)=c} 1$$

$$c^* = (c+1) \frac{N_{c+1}}{N_c} \quad \} \text{ formula to reestimate counts}$$

$$P_{GT}^* (\text{things with zero freq in training}) = \frac{N_1}{N} \quad \} \text{ Good Turing probability}$$

N_1 \rightarrow number of items in bin 1

N \rightarrow total no. of items in training set

Eg. Fishing in a lake with 8 species

bass - 0	unseen species - bass, catfish probability that next fish caught is bass or catfish $= P_{GT}^*(\text{bass, catfish}) = \frac{N_1}{N} = \frac{3}{18}$
carp - 16	
catfish - 0	
eel - 1	
perch - 3	
salmon - 1	
trout - 1	
whitefish - 2	
<u>total - 28</u>	

$$\begin{aligned} N_0 &= 2 \\ N_1 &= 3 \\ N_2 &= 1 \\ N_3 &= 2 \\ N_4 &= 0 \\ N_5 &= 0 \end{aligned}$$

$$\begin{aligned} N_6 &= 0 \\ N_7 &= 0 \\ N_8 &= 0 \\ N_9 &= 0 \\ N_{10} &= 1 \end{aligned}$$

unseen (Bam or Catfish)		trout
c	0	1
MLE	$P = \frac{0}{18}$	$\frac{1}{18}$
c^*		$(1+2) * \frac{N_2}{N_1} = \frac{2}{3}$
p_{GT}^*	$\frac{3}{18}$ ↳ divided among all unseen events	$\frac{(2/3)}{18}$

Revised count c^* for trout was calculated
from $c=1$ to $c^* = 0.67$

Interpolation (Linear)

$$P(w_n | w_{n-2} w_{n-1})$$

If there is no example of the trigram $w_n w_{n-1} w_{n-2}$ in the corpus, we use bigram to estimate its probability.

↓

$$P(w_n | w_{n-1})$$

If bigram is not present, then do unigram.

↙

$$P(w_n)$$

$$P(w_n | w_{n-2} w_{n-1}) = \lambda_1 P(w_n | w_{n-2} w_{n-1}) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n)$$

$$\boxed{\sum_i \lambda_i = 1}$$

The λ values need to be learned from a defined corpus.

Back off N-gram Modelling

→ Katz Backoff

If N-gram we need has 0 count, we approximate it by
backing off to N-1 gram, recursively.

$$P_{\text{Katz}}(w_n | w_{n-N+1}^{n-1}) = \begin{cases} p^*(w_n | w_{n-N+1}^{n-1}) & \text{if } c(w_{n-N+1}^n) > 0 \\ \alpha(w_{n-N+1}^{n-1}) \cdot P_{\text{Katz}}(w_n | w_{n-N+2}^{n-1}) & \text{otherwise} \end{cases}$$

for trigram x, y, z

$$P_{\text{Katz}}(z|x,y) = \begin{cases} p^*(z|x,y) & \text{if } c(x,y,z) > 0 \\ \alpha(x,y) P_{\text{Katz}}(z|y) & \text{else if } c(x,y) > 0 \\ p^*(z) & \text{otherwise} \end{cases}$$

for bigram y, z

$$P_{\text{Katz}}(z|y) = \begin{cases} p^*(z|y) & \text{if } c(y,z) > 0 \\ \alpha(y) p^*(z) & \text{otherwise} \end{cases}$$

Discounted probability

$$p^*(w_n | w_{n-N+1}^{n-1}) = \frac{c^*(w_{n-N+1}^n)}{c(w_{n-N+1}^{n-1})}$$

$$\text{original MLE} = \frac{c(w_{n-N+1}^n)}{c(w_{n-N+1}^{n-1})}$$

p^* is slightly smaller than the original MLE.
This leftover probability may be distributed to lower order
N-grams using the parameter α .

Left over probability mass sum of all patterns for which count is non-zero.

$$p(w_{n-N+1}^{n-1}) = 1 - \underbrace{\sum_{w_n} p^*(w_n | w_{n-N+1}^{n-1})}_{\substack{\vdots \\ c(w_{n-N+1}^n) > 0}}$$

Distribute β among patterns with zero counts

$$\begin{aligned} d(w_{n-N+1}^{n-1}) &= \frac{p(w_{n-N+1}^{n-1})}{\sum_{w_n: c(w_{n-N+1}^n) = 0} p_{\text{ideal}}(w_n | w_{n-N+2}^{n-1})} \\ &= \frac{1 - \sum_{w_n: c(w_{n-N+1}^n) > 0} p^*(w_n | w_{n-N+1}^{n-1})}{1 - \sum_{w_n: c(w_{n-N+1}^n) > 0} p^*(w_n | w_{n-N+2}^{n-1})} \end{aligned}$$

Evaluating N-Gram Models

(a) Extrinsic evaluation

- embed it in an application and measure performance of the application.

(b) Intrinsic evaluation

- measure quality of model irrespective of application

e.g. perplexity (PP)

for a test set $w = w_1 w_2 \dots w_N$

$$PP(w) = \boxed{p(w_1 w_2 \dots w_N)^{-1/N}}$$

\uparrow probability (calculated using any language model)

Pos tagging

Problem: Determine best tag sequence for a sentence.

$$T = t_1, t_2, t_3, \dots, t_n \text{ (tags)}$$

$$W = w_1, w_2, w_3, \dots, w_n \text{ (sentence)}$$

Input : The lead paint is unsafe

Output : The/DT lead/NP Paint/VB is/V unsafe/ADJ

Approach: 1) Rule based

or 2) Hidden Markov Model based tagger

HMM

↳ Components of the Model:

$$Q = \{q_1, q_2, q_3, \dots, q_N\} \rightarrow \text{set of } N \text{ states}$$

$$A = \{a_{11}, a_{12}, a_{13}, \dots, a_{ij}, \dots, a_{nn}\} \rightarrow \text{transition probability matrix}$$

↳ represents probability of moving from state i to state j .

$$\sum_{j=1}^N a_{ij} = 1 \quad \forall i$$

$$O = \{o_1, o_2, \dots, o_T\} \rightarrow \text{sequence of } T \text{ observations each drawn from a vocabulary } V = \{v_1, v_2, \dots, v_k\}$$

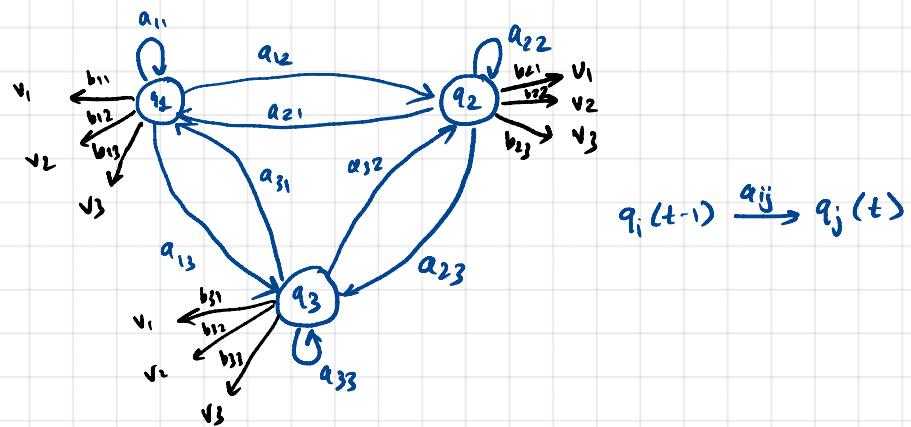
$$B = b_i(o_t) \rightarrow \text{Emission probabilities}$$

↳ probability of observation O_t being generated from state i .

$$\pi = \{\pi_1, \pi_2, \dots, \pi_N\} \rightarrow \text{initial probability distribution}$$

$\pi_i \rightarrow$ probability that Markov chain starts from state Q_i

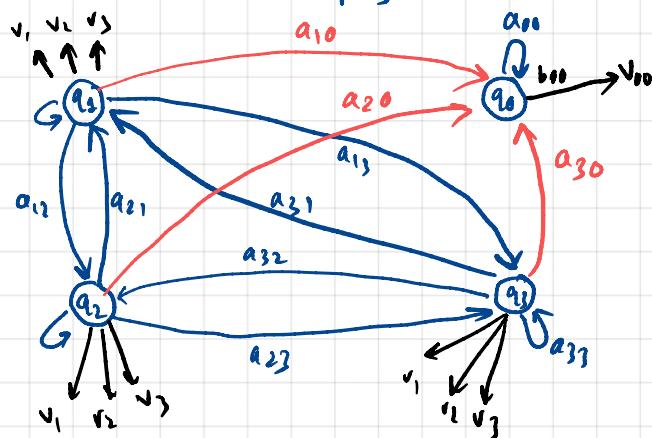
$$\sum_{i=1}^N \pi_i = 1$$



$\sum_j a_{ij} = 1 \quad \forall i \quad \rightarrow$ sum of all transition probabilities of a state is 1.

$\sum_k b_{jik} = 1 \quad \forall j \quad \rightarrow$ sum of all emission probabilities for all hidden states is 1.

Once HMM reaches accepting state, it emits v_0 .



→ Central Issues / Basic Problems

- (a) Evaluation Problem
- (b) Decoding Problem
- (c) Learning Problem

Evaluation Problem

→ forward Algorithm

Given a HMM θ and observation sequence v^T ,
determine the likelihood

$$P(v^T | \theta)$$

↪ probability of emitting v^T from machine θ .

Known quantities: $\theta, O, a_{ij}, b_{ik}$

Decoding Problem

→ Viterbi Algorithm

Given an observation sequence and HMM θ ,
find the best hidden state sequences α .

once we
know α ,
we can find
the tag.

Learning Problem

Estimate a_{ij}, b_{ik} from a large annotated training set

Evaluation Problem

Given: $\delta, \alpha, \alpha_{ij}, b_{ik}$

$$P(v^T | \alpha) = \sum_{n=1}^{n_{\max}} P(v^T | \alpha_n^T) * P(\alpha_n^T) \rightarrow ①$$

$\alpha_n^T \rightarrow$ denotes a sequence from δ_T

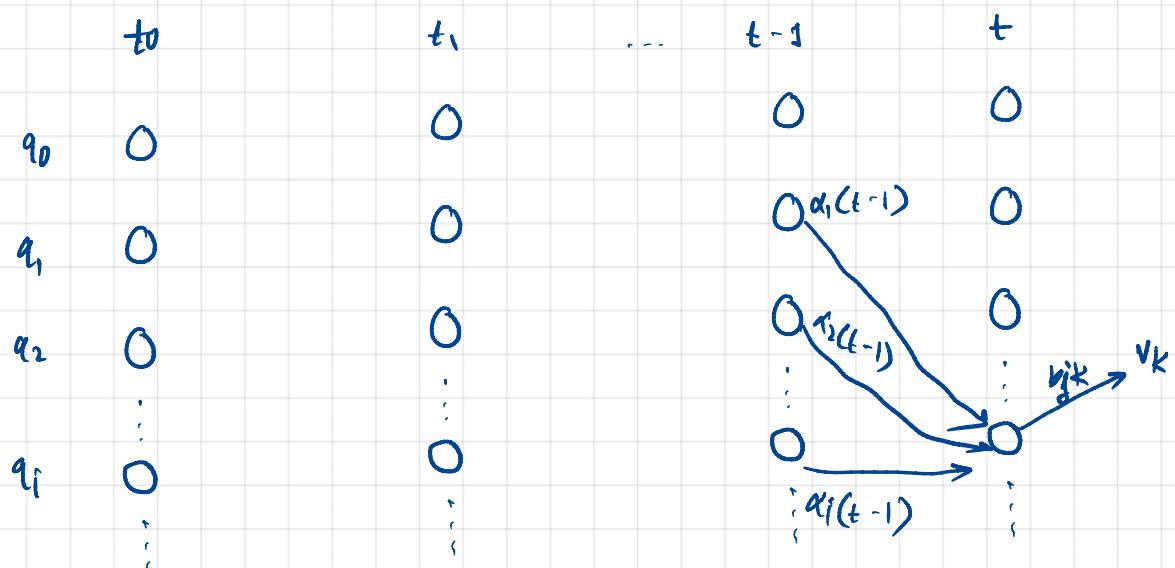
$n_{\max} \rightarrow$ no. of all possible sequences

$$\alpha_n^T = \{ \delta(1), \delta(2), \dots, \delta(T) \} \quad T \text{ can be thought of as time-stamps.}$$

$$P(\alpha_n^T) = \prod_{t=1}^T P(\delta(t) | \delta(t-1)) \rightarrow ②$$

$$P(v^T | \alpha_n^T) = \prod_{t=1}^T P(v(t) | \delta(t)) \rightarrow ③$$

$$\therefore P(v^T | \alpha) = \sum_{n=1}^{n_{\max}} \prod_{t=1}^T P(v(t) | \delta(t)) . P(\alpha(t) | \alpha(t-1))$$



At time t , the probability of the previous state being q_i is $\alpha_i(t-1)$

Note that $\alpha_i(t-1)$ is not the same as a_{ii} which is the transition probability.

$$\xrightarrow{\text{probability that machine will be in state } q_j \text{ after generating symbol } v_t} \alpha_j(t) = \begin{cases} 0 & t=0, q_j \neq \text{initial state} \\ 1 & t=0, q_j = \text{initial state} \\ \sum_{i=1}^N \alpha_i(t-1) a_{ij} b_{jk} v(t) & \text{otherwise} \end{cases}$$

forward algorithm:

1. Initialization

$$t \leftarrow 0, a_{ij}, b_{jk}, v^T, \alpha_j(0)$$

2. for $t \leftarrow t + 1$

$$\alpha_j(t) = b_{jk} v(t) \sum_{i=1}^N \alpha_i(t-1) \cdot a_{ij}$$

3. until $t = T$

4. Return $P(v^T | \alpha) \leftarrow \alpha_0(T)$



Probability that it reaches final state q_0 and emits v_0 at time T .

Consider an example.

$V^5 \rightarrow$ sequence of five symbols = $\{V_3, V_1, V_2, V_1, V_3\}$

$\alpha_3(3) \rightarrow$ what is the probability that HMM will be in state Q_3 after generating first 3 symbols $\langle V_3, V_1, V_2 \rangle$

$\alpha_2(4) \rightarrow$ in state 2 after generating 4 symbols $\langle V_3, V_1, V_2, V_1 \rangle$

e.g. Hidden state : Q_1, Q_2, Q_3 and final state Q_0

Visible state : V_0, V_1, V_2, V_3, V_4

$$a_{ij} = \begin{bmatrix} Q_0 & Q_1 & Q_2 & Q_3 \\ 1 & 0 & 0 & 0 \\ 0.2 & 0.3 & 0.1 & 0.4 \\ 0.2 & 0.5 & 0.2 & 0.1 \\ 0.7 & 0.1 & 0.1 & 0.1 \end{bmatrix} \quad \sum_j a_{ij} = 1 \quad \forall i$$

$$b_{jk} = \begin{bmatrix} V_0 & V_1 & V_2 & V_3 & V_4 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0.3 & 0.4 & 0.1 & 0.2 \\ 0 & 0.1 & 0.1 & 0.7 & 0.1 \\ 0 & 0.5 & 0.2 & 0.1 & 0.2 \end{bmatrix} \quad \sum_k b_{jk} = 1 \quad \forall j$$

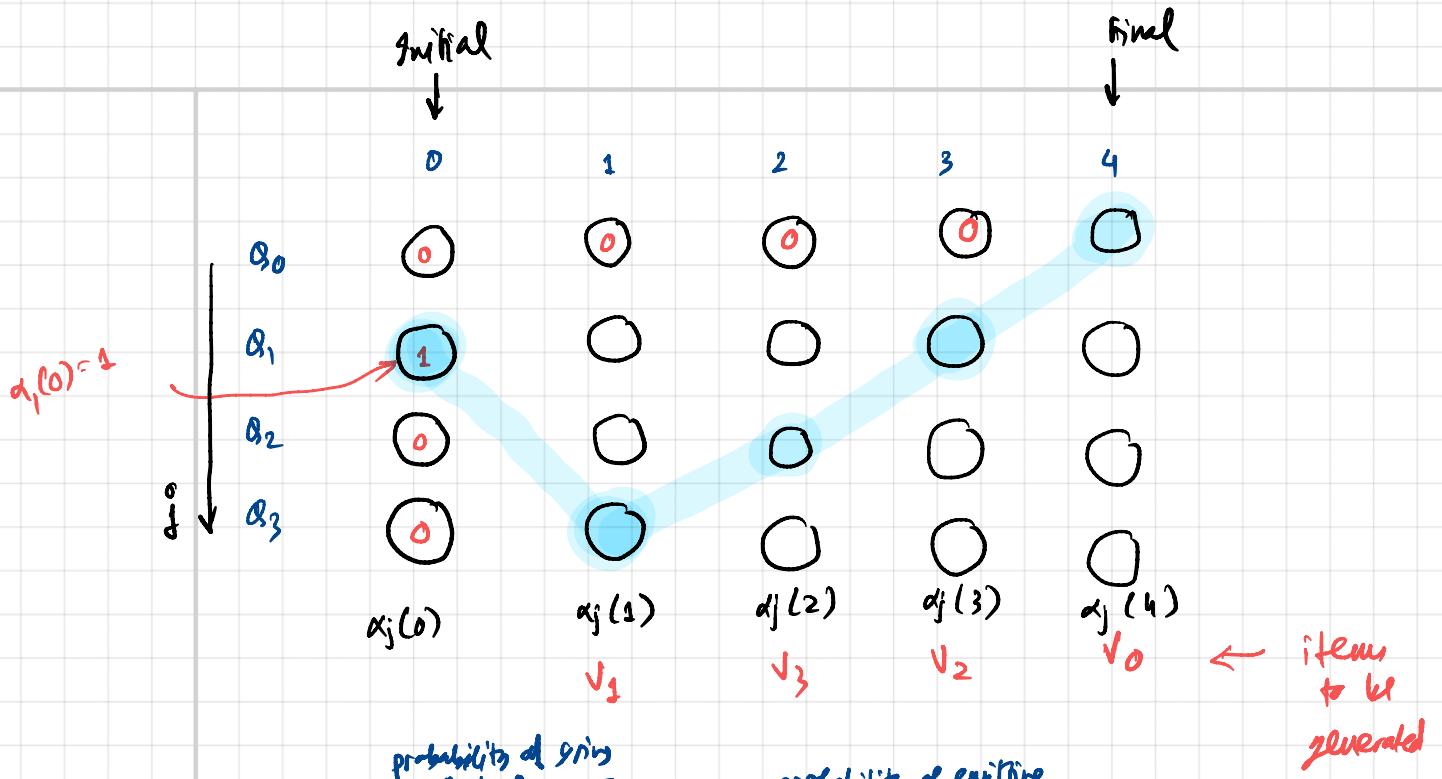
$$V^4 = \{V_1, V_3, V_2, V_0\}$$

what is $P(V^4 | \sigma)$?

Assume, at t=0, initial state Q_1



$\kappa_1(0) = 1$



$$\begin{aligned} \alpha_1(1) &= [\alpha_1(0) * a_{11} * b_{11}] \\ &\quad + [\alpha_2(0) * a_{21} * b_{21}] = 0.09 + 0 + 0 = 0.09 \\ &\quad + [\alpha_3(0) * a_{31} * b_{31}] \end{aligned}$$

$$\alpha_2(1) = [1 * 0.1 * 0.1] + 0 + 0 = 0.01$$

$$\alpha_3(1) = 0.2$$

$$\alpha_1(2) = 0.0052$$

$$\alpha_2(2) = 0.0343$$

$$\alpha_3(2) = 0.0057$$

$$\alpha_1(3) = 0.0019$$

$$\alpha_2(3) = 0.0008$$

$$\alpha_3(3) = 0.0012$$

$$\alpha_1(4) = 0$$

$$\alpha_2(4) = 0$$

$$\alpha_3(4) = 0$$

$$\alpha_0(4) = 0.00138$$

Path $\rightarrow Q_1 \rightarrow Q_3 \rightarrow Q_2 \rightarrow Q_1 \rightarrow Q_0$

[by finding max probability in each time step]

$$P(V_4 | \theta) = 0.00138$$

evaluation problem

↑
decoding problem

At time $t=0$,

most probable state θ_1

At time $t=1$

most probable state θ_3

$t=2 \rightarrow \theta_2$

$t=3 \rightarrow \theta_1$

$t=4 \rightarrow \theta_0$

This can be
considered as
sentence

↳ Output sequence: v_1, v_3, v_2, v_0

↳ Hidden state sequence: $\theta_1, \theta_3, \theta_2, \theta_1, \theta_0$

This can
be considered
as pos tags

→ Algorithm for Decoding:

Given v^T , find most probable hidden state sequences.

1. Initialization: $\text{Path} = \{\emptyset\}$, $t \leftarrow 0, j \leftarrow 0, a_{ij}, b_{jk}$

2. for $t \leftarrow t+1$

 for $j \leftarrow j+1$

$$\alpha_j(t) \leftarrow b_{jk} v(t) \sum_{i=1}^N \alpha_i(t-1) \cdot a_{ij}$$

 until $j=N$

$$j' \leftarrow \arg \max_j \alpha_j(t)$$

 Append $\theta_{j'}$ to $\boxed{\text{Path}}$

 until $t=T$

3. Return $\boxed{\text{Path}}$

Assumption:

All states are
interconnected and
reachable except θ_0
(with each other)

(because otherwise when
we do \max , that
path may not exist.)

HMM Learning Problem

Forward-Backward algorithm
(Baum-Welch algorithm)

→ Evaluation Problem

$$P(V^T | \theta) = \sum_{\alpha_0}^{n_{max}} P(V^T | \alpha_0^T) * P(\alpha_0^T)$$

Estimate a_{ij} and b_{jk}

→ Forward Algorithm

$$\alpha_j = \begin{cases} 0 & t=0 \text{ } \& q_j \neq \text{initial state} \\ 1 & t=0 \text{ } \& q_j = \text{initial state} \\ [\sum \alpha_i(t-1) a_{ij}] b_{jk} v(t) & \text{otherwise} \end{cases}$$

→ Backward Algorithm

$$\beta_i(t) = \begin{cases} 0 & q_i(t) \neq q_0 \text{ at } t=T \\ 1 & q_i(t) = q_0 \text{ at } t=T \\ [\sum \beta_j(t+1) a_{ij}] b_{jk} v(t+1) & \text{otherwise} \end{cases}$$

Define $\gamma_{ij}(t)$

↪ Probability of transition from $q_i(t-1)$ to $q_j(t)$

$$\gamma_{ij}(t) = \frac{\alpha_i(t-1) a_{ij} b_{jk} \beta_j(t)}{P(V^T | \theta)}$$

$$\hat{a}_{ij} = \frac{\text{Expected no. of transitions from state } i \text{ to } j}{\text{expected no. of transitions from state } i}$$

$$\text{Expected no. of transitions from } q_i(t-1) \text{ to } q_j(t) = \sum_{t=1}^T \gamma_{ij}(t)$$

$$\text{Expected no. of transitions from } q_i \text{ to any state} = \sum_{t=1}^T \sum_k \gamma_{ik}(t)$$

∴ Refined transition probability

$$\hat{a}_{ij} = \frac{\sum_{t=1}^T \gamma_{ij}(t)}{\sum_{t=1}^T \sum_k \gamma_{ik}(t)}$$

$\hat{b}_{jk} = \frac{\text{Expected no. of times in state } j \text{ and observing the symbol } v_k}{\text{Expected no. of times in state } j}$

$$\hat{b}_{jk} = \frac{\sum_{t=1, k \in V(t) \cap v(k)}^T \sum_l \gamma_{jl}(t)}{\sum_{t=1}^T \sum_l \gamma_{jl}(t)}$$

↓ Refined emission probability

Forward-Backward algorithm

1. Initialize A & B .
2. Iterate until convergence
 - Define $\delta_{ij}(t)$
 - compute \hat{a}_{ij}
 - compute \hat{b}_{jk}
3. Return A & B

Syntactic Analysis

- (a) Top-down Parser
- (b) Bottom-up Parser
- (c) Dynamic Programming Parsing algorithm (Early, CYK)
- (d) Probabilistic Parsing

Context-Free Grammar $\rightarrow G = (V, \Sigma, P, S)$

For CYK, CFG has to be converted to its Chomsky Normal Form.

$$\text{e.g. } A \rightarrow BCD \quad | \quad \begin{array}{l} A \rightarrow BX \\ X \rightarrow CD \end{array}$$

Some rules for English:

- I prefer a morning flight.
 Subject Noun Verb
- show the latest fore.
 Verb
- Yes/No questions
- "uh.." questions

$$\begin{array}{ll} S \rightarrow NP VP \\ S \rightarrow X_1 VP \\ X_1 \rightarrow Aux NP \end{array}$$

$$\begin{array}{ll} S \rightarrow book | include | prefer \\ S \rightarrow Verb NP \\ S \rightarrow X_2 PP \\ S \rightarrow Verb PP \end{array}$$

$$\begin{array}{ll} NP \rightarrow I | she | me \\ NP \rightarrow Houston \\ NP \rightarrow Det Nominal \\ Nominal \rightarrow book | flight | meal | money \\ Nominal \rightarrow Nominal Noun \end{array}$$

$$\begin{array}{ll} VP \rightarrow book | include | prefer \\ VP \rightarrow X_2 PP \\ X_2 \rightarrow Verb NP \\ VP \rightarrow Verb PP \\ VP \rightarrow VP PP \\ PP \rightarrow preposition NP \end{array}$$

$$Det \rightarrow the | a$$

$$\begin{array}{ll} Noun \rightarrow book | flight | meal | money \\ Verb \rightarrow book | include | prefer \\ Proper noun \rightarrow Houston \\ Aux \rightarrow Does \\ Preposition \rightarrow from | to | on | through \end{array}$$

Using the rules above, something like
 "Book the flight through Houston" \rightarrow sentence length = $n = 5$
 is derivable with multiple parse trees.

We develop the CKY algorithm below where we fill the upper triangular portion of a $(n+1) \times (n+1)$ matrix using dynamic programming, where $n = \text{sentence length}$

CKY - Cocke Kasami Younger Algorithm

function CKY-Parse (words, grammar)

```

for j ← from 1 to length(words) do           // diagonal
    for all rules { A | A → words[j] ∈ Grammar }
        table[j-1, j] ← table[j-1, j] ∪ A

    for j ← length(words) down to 0 do
        for i ← from j-2 down to 0 do
            for k ← i+1 to j-1 do
                for all { A | A → BC ∈ grammar and
                            B ∈ table[i, k] and
                            C ∈ table[k, j] }

                    table[i, j] ← table[i, j] ∪ A
    
```

from this table,
 multiple parse trees \longrightarrow
 are possible if grammar
 is ambiguous.

return table

e.g. sentence : Book the flight through Houston $\rightarrow \text{len}(\text{words}) = 5$
 grammar : defined in previous page.

Book	the	flight	through	Houston
S ← VP Verb Nominal noun	[0,1]	S ← VP Verb Nominal noun	S ← VP Verb Nominal noun	S ← VP Verb Nominal noun
Det ← NP	[0,2]	[0,3]	[0,4]	[0,5]
[1,2]	[1,3]	[1,4]	[1,5]	[1,6]
new neutral	[2,3]	[2,4]	[2,5]	[2,6]
	[3,4]	[3,5]	[3,6]	[3,7]
		[4,5]	[4,6]	[4,7]
			[4,7]	[4,8]
				[4,9]

Fill-up order of cells is written as (1,1), (1,2), (1,3), (1,4), (1,5), (2,2), (2,3), (2,4), (2,5), (2,6), (3,3), (3,4), (3,5), (3,6), (3,7), (4,4), (4,5), (4,6), (4,7), (4,8), (5,5), (5,6), (5,7), (5,8), (5,9).

For rule $A \rightarrow BC$, check if B belongs to cells on left and C belongs to cells below, then add A to cell.

Probabilistic CFG Parsing (PCFG)

$$G = (\Sigma, \Sigma, P, S)$$

Σ = set of non-terminals

Σ = set of terminals

P = set of production rules

$$A \rightarrow \beta[P]$$

$$\beta \in \{\Sigma \cup \Sigma\}^*$$

S = start symbol

$$P(A \rightarrow \beta) \quad \text{or} \quad P(\beta | A)$$

$$\sum_{\beta} P(A \rightarrow \beta) = 1$$

: P is probability of β being derivable from A .

Example PCFG:

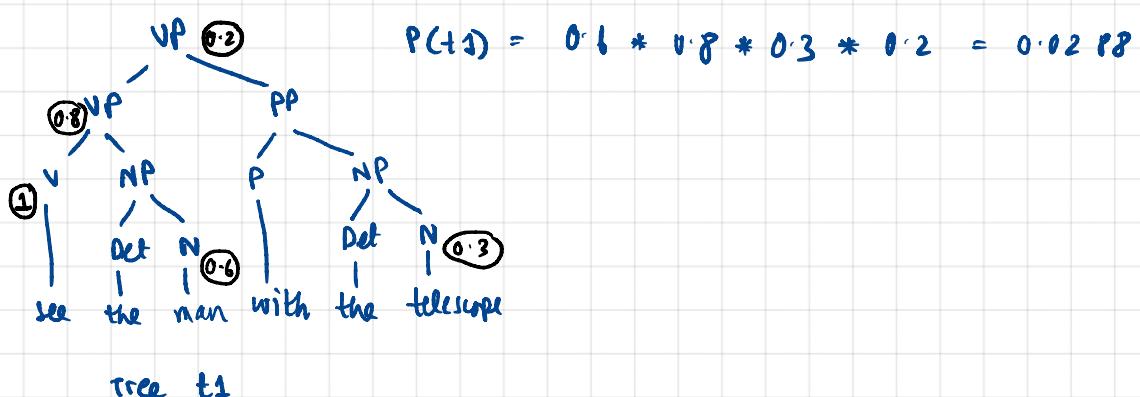
$VP \rightarrow V$	NP	$[0.8]$
$VP \rightarrow VP$	PP	$[0.2]$
$NP \rightarrow Det$	N	$[1]$
$Det \rightarrow$	the	$[1]$
$P \rightarrow$	with	$[1]$
$N \rightarrow$	man	$[0.6]$
$N \rightarrow$	telescope	$[0.3]$
$PP \rightarrow P$	NP	$[1]$
$N \rightarrow N$	PP	$[0.1]$
$V \rightarrow$	see	$[1]$

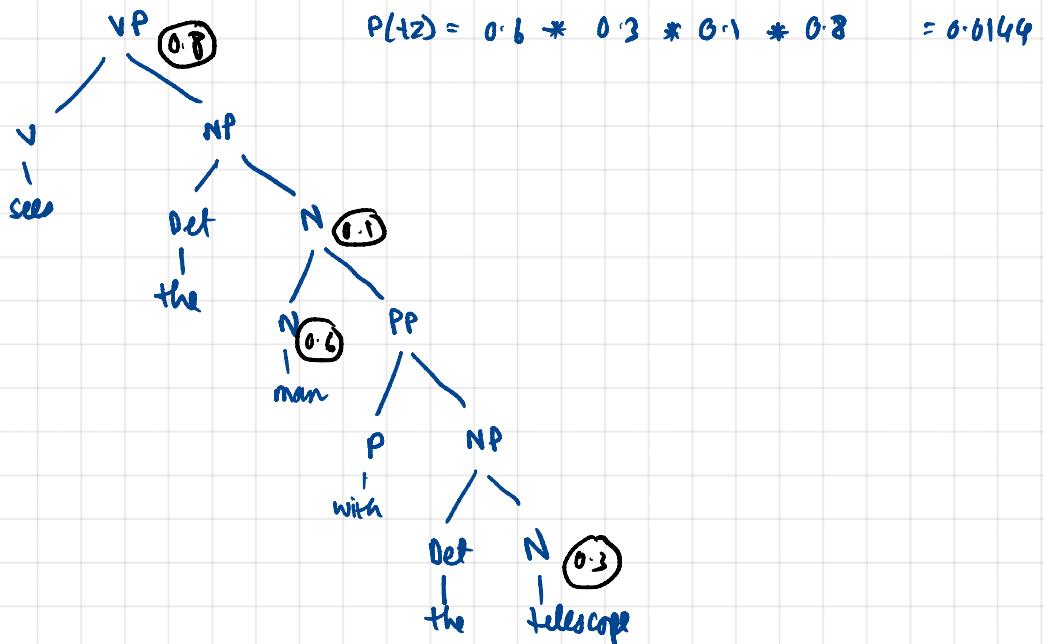
VP is the start symbol

Probability of a parse tree

= product of the probabilities of the rules used to generate the parse tree.

e.g. sentence \rightarrow see the man with the telescope





Tree t_2

consistent \rightarrow sum of all probabilities of all sentences in grammar the language is equal to 1.

$$P(VP, \text{sentence}) = P(t_1) + P(t_2) + \dots = 1$$

Example PCFG:

$S \rightarrow A [0.4]$ $S \rightarrow B [0.6]$ $A \rightarrow a [1]$ $B \rightarrow B [1]$	$\} \rightarrow$ when infinite derivations are possible, grammar is inconsistent
--	--

Probabilistic CKY parsing of PCFG:

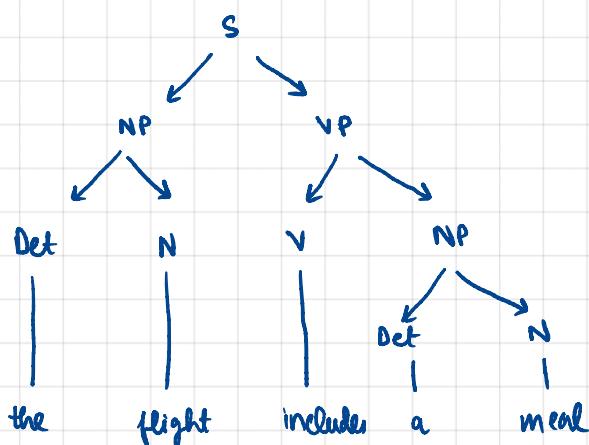
$S \rightarrow NP VP [0.8]$
$NP \rightarrow Det N [0.3]$
$VP \rightarrow V NP [0.2]$
$V \rightarrow \text{includes} [0.05]$
$Det \rightarrow \text{the} [0.4]$
$Det \rightarrow a [0.4]$
$N \rightarrow \text{meal} [0.01]$
$N \rightarrow \text{flight} [0.02]$

sentence: the flight includes a meal

the	flight	includes	a	meal
Det : 0.4	NP : $0.3 \times 0.4 \times 0.02 = 0.0024$			S : $0.8 \times 0.0024 \times 0.0000012 = 2.3 \times 10^{-9}$
[0,1]	[0,2]	[1,3]	[0,4]	[1,5]
	N : 0.002			
			V : 0.05	VP : $0.2 \times 0.05 \times 0.0012 = 0.0000612$
			[2,3]	[2,5]
			[2,4]	
				Det : 0.4
				NP : $0.3 \times 0.4 \times 0.01 = 0.0012$
				[3,4]
				N : 0.01
				[4,5]

No rules exist to fill up the remaining cells.

Homework → Assignment Write the algorithm for backtrace to find the most probable parse tree.

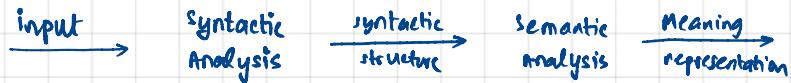


Learn PCFG rule probabilities

$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{count}(\alpha \rightarrow \beta)}{\sum_{\gamma} \text{count}(\alpha \rightarrow \gamma)} = \frac{\text{count}(\alpha \rightarrow \beta)}{\text{count}(\alpha)}$$

A corpus has to be used where there are annotated parse trees for each sentence.

3/5/22



Properties of Knowledge Base:

Meaning representation can be in terms of FOL.

can be represented using Lambda notation.

- referentiality

e.g. Does Spice Island serve veg food?
↳ serve (Spice Island, veg food)

- ambiguity

e.g. Give me the book?
↳ which book?

- canonical form

- inference and variables

e.g. serves (x , vegfood).

Examples:

Every one likes chocolate,

$$\forall x \text{ person}(x) \Rightarrow \text{Likes}(x, \text{chocolate})$$

Someone likes chocolate,

$$\exists x \text{ person}(x) \wedge \text{Likes}(x, \text{chocolate})$$

Everyone likes chocolate unless they are allergic to it.

$$\forall x (\text{Person}(x) \wedge \neg \text{allergic}(x, \text{chocolate})) \Rightarrow \text{Likes}(x, \text{chocolate})$$

A restaurant near ICSI serves meat food.

$$\exists x (\text{Restaurant}(x) \wedge \text{serves}(x, \text{MeatFood}) \wedge \text{near}(\text{westm}(x), \text{Location}(ICSI)))$$

Lambda notation:

$$\lambda x. P(x)$$

Lambda reduction:

$$\begin{cases} \lambda z. P(z) (A) \\ P(A) \end{cases}$$

e.g. $\lambda x. \lambda y. \text{Near}(x, y) (\text{Kolkata})$

$$\hookrightarrow \lambda y. \text{Near}(\text{Kolkata}, y)$$

Howrah is near
Kolkata.

e.g. $\lambda y. \text{Near}(\text{Kolkata}, y) (\text{Howrah})$

$$\hookrightarrow \text{Near}(\text{Kolkata}, \text{Howrah})$$

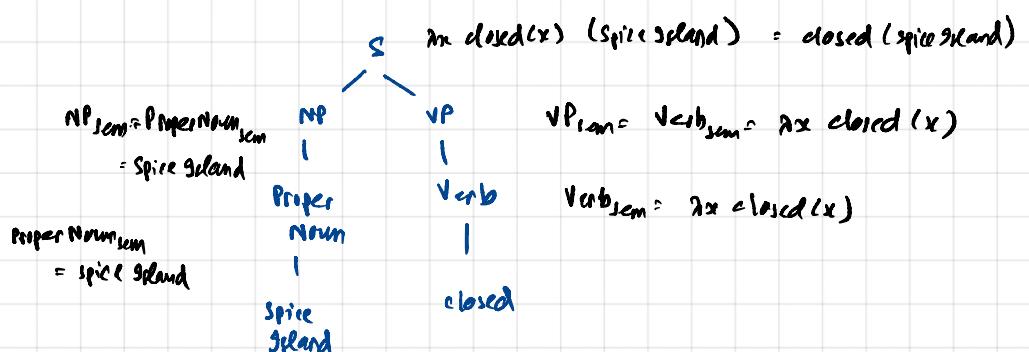
Syntax driven Semantic analysis:

$$A \rightarrow \alpha_1, \alpha_2, \dots, \alpha_n \quad \{ f[\text{Ls.sem}, \dots, \text{Rk.sem}] \}$$

Example:

$$\begin{array}{ll} SP \rightarrow NP \ VP & \{ VP_{sem}(NP_{sem}) \} \\ NP \rightarrow Proper\ Noun & \{ Proper\ Noun_{sem} \} \\ Proper\ Noun \rightarrow Spice\ Island & \{ Spice\ Island \} \\ VP \rightarrow Verb & \{ Verb_{sem} \} \\ Verb \rightarrow closed & \{ \lambda x. \text{closed}(x) \} \end{array}$$

e.g. Spice Island closed.



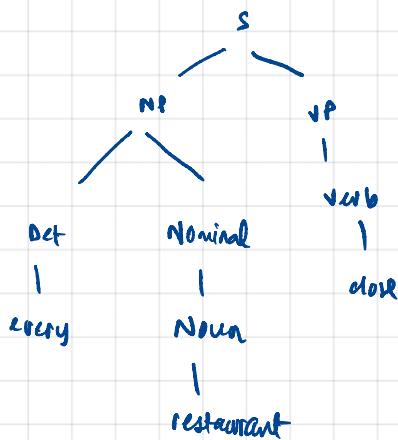
Example:

$SP \rightarrow NP \ VP$	$\{ VP_{sem} (NP_{sem}) \}$
$NP \rightarrow Det \ Nominal$	$\{ Det_{sem} (Nominal_{sem}) \}$
$Det \rightarrow every$	$\{ AP \ \forall x \ P(x) \Rightarrow \alpha(x) \}$
$Nominal \rightarrow Noun$	$\{ Noun_{sem} \}$
$Noun \rightarrow restaurant$	$\{ \forall x \ Restaurant(x) \}$
$verb \rightarrow close$	$\{ \exists e \ closed(e) \wedge closed_thing(e, x) \}$

q. Every restaurant is closed.

Ans: $\forall x \ Restaurant(x) \Rightarrow \exists e \ closed(e) \wedge closed_thing(e, x)$

H.W.



$$Det.sem = \forall P \ \forall B \ \forall x \ P(x) \Rightarrow \alpha(x)$$

$$Noun.sem = \forall x \ Restaurant(x)$$

$$Nominal.sem = \forall x \ Restaurant(x)$$

$$NP.sem = \underline{\forall P \ \forall B \ \forall x \ P(x) \Rightarrow \alpha(x)} \quad (\underline{\forall x \ Restaurant(x)})$$

$$= \forall B \ \forall x \ \underline{\forall x \ Restaurant(x)} \Rightarrow \alpha(x)$$

$$Verb.sem = \forall x \ \exists e \ closed(e) \wedge closed_thing(e, x)$$

5/5/21

df_t = Number of documents in the collection that contain the term t .

N = Total number of documents in the collection.

$$idf_t = \log \frac{N}{df_t}$$