

NLP - EX 2

- 1) Prove that any generative probabilistic classification model that assumes that

 - the sequence of labels forms a Markov Chain
 - the generation of x_i is conditionally independent given y_i of all other x_j and y_j

is an HMM model.

- * Generative model: works with a joint probabilistic model of data. Assumes functional form for $P(X|Y)$ and $P(Y)$. Estimates parameters from data.

$$\text{Prediction } y^* = \underset{y}{\operatorname{argmax}} \text{IP}(Y|X) = \underset{y}{\operatorname{argmax}} \text{IP}(X|Y) \cdot \text{IP}(Y)$$

↑
bayes rule

- * Markov Chain : a finite state automaton with probabilistic state transitions.

Markov Assumption: next state only depends on the current state and independent of previous history given the current state.

A maximum likelihood estimator for a K^{th} order language model is : Markov

$$\text{g}_{ML} (y_m = +m_j | y_{m-1} = t_{i1}, \dots, y_{m-k} = t_{ik}) = \frac{\text{Count}(t_j, t_{i1}, \dots, t_{ik})}{\sum \text{Count}(t_j, t_{i1}, \dots, t_{ik})}$$

transition probability (label)

Input sentence $x = x_1, x_2, \dots, x_m$ x_i תואר בפונקציית
 Output sequence $y = y_1, y_2, \dots, y_{m+1}$ y_i תואר בפונקציית

The most likely tag sequence for X is:

$$\arg \max_{y_1, \dots, y_{m+1}} \text{IP}(x_1, \dots, x_m, y_1, \dots, y_{m+1})$$

פונקציית פירסינג (Piercing Function) מוגדרת כפונקציה $f(x)$ אשר מקבלת כערך כניסה נקודה (x, y) ומחזירה 1 אם נקודה (x, y) נמצאת בתחום הגרף של פונקציית $y = f(x)$, ו- 0 אחרת.

$$\begin{aligned} P(y_1, \dots, y_m) &= \prod_{i=1}^m P(y_i | y_{i-1}) \quad \leftarrow \text{Markov Chain} \\ &= \prod_{i=1}^m P(y_i | y_{i-1}, \dots, y_0) \quad \text{transition probability} \end{aligned}$$

בנוסף למכרזים נוראים נשים, מטבחיים, גזעיים, גזעים נאכליים, גזעים אקולוגיים ו遐

χ_i נסמן כטווית ה- i -הוותה של המונומטר $P(X|Y)$

לעומת ה-POS נספרו מילים ושורשים, POS כנ"ג מכוון ל-POS-NPVR (או POS->POS).

$$P(x_1, \dots, x_m, y_1, \dots, y_m) = \left(\prod_{i=1}^m q(\text{STOP}|y_m) \right) P(y_1, \dots, y_m) \cdot \prod_{i=1}^{m-1} q(y_i|y_{i-1}) e(x_i|y_i)$$

התקנים (הנעלים)
תוקן נובע סופר כפער
ה נתונים בוחנים נא- training data

ב-טנזורים כ-טנזורים
ר-טנסורס כ-טנסורס
וירוטס כ-טנתקס (*)

$$P(X_1, \dots, X_m, Y_1, \dots, Y_{m+k}) = \left(\prod_{i=1}^{m+k} P(Y_i | Y_{i-1}, \dots, Y_{i-k}) \right) \quad \text{because } Y_0 = Y_{-1} = \dots = Y_{-(k-2)} = \text{STOP}$$

$y_0 = y_{-1} = \dots = y_{-(k-2)} = \dots$ where * is START symbol $y_{m+1} = \text{STOP}$ $x(k)$

$$\arg\max_{y_1, \dots, y_m} P(x_1, \dots, x_m, y_1, \dots, y_m)$$

תורם גנרטיבי (ג'י.ג'). גן יוניברסיטאי הוא גוף ציבורי שמייסדים
הוועדה המרכזית לניהולו וניהולו בפועל.

$$States = \{ H, L \}$$

$$P(L|H) = 0.5, \quad P(H|L) = 0.4 \\ P(A|H) = 0.2, \quad P(C|H) = 0.3, \quad P$$

$$P(T|H) = 0.2, \quad P(A|L) = 0.3, \quad P(C|L) = 0.2, \quad P(G|L) = 0.2, \quad P(T|L) = 0$$

$$P(A \cap E) = 0.3, P(E \cap B) = 0.2, P(G \cap E) = 0.2, P(C \cap E) = 0.1$$

$S = ACCGTTGCA$
previous state before S was H

$$P(H|X) = P(H|H)P(S|H) + P(H|H')P(G|H')$$

$$P(H \mid A) = P(A \mid H) \cdot P(H) / [P(A \mid H) \cdot P(H) + P(A \mid \neg H) \cdot P(\neg H)]$$

$\text{IP}(\text{H}_1\text{H})\text{IP}(\text{C}_1\text{H})$ $\text{IP}(\text{H}_1\text{L})\text{IP}(\text{C}_1\text{H})$

~~$$P(H|L) \quad P(L|H)P(C|L) \quad P(L|H)P(C|L) \quad P(H|L)P(C|L)$$~~

$\text{O}(\text{AII})$ L $\text{P}(\text{CH})$

$$(T1H) \quad IP(L \rightarrow L)IP(C \sqsubset L) \quad IP(L \sqsubset L)IP(C \sqsubset L) \quad IP(L \sqsubset L)IP(G \sqsubset L)$$

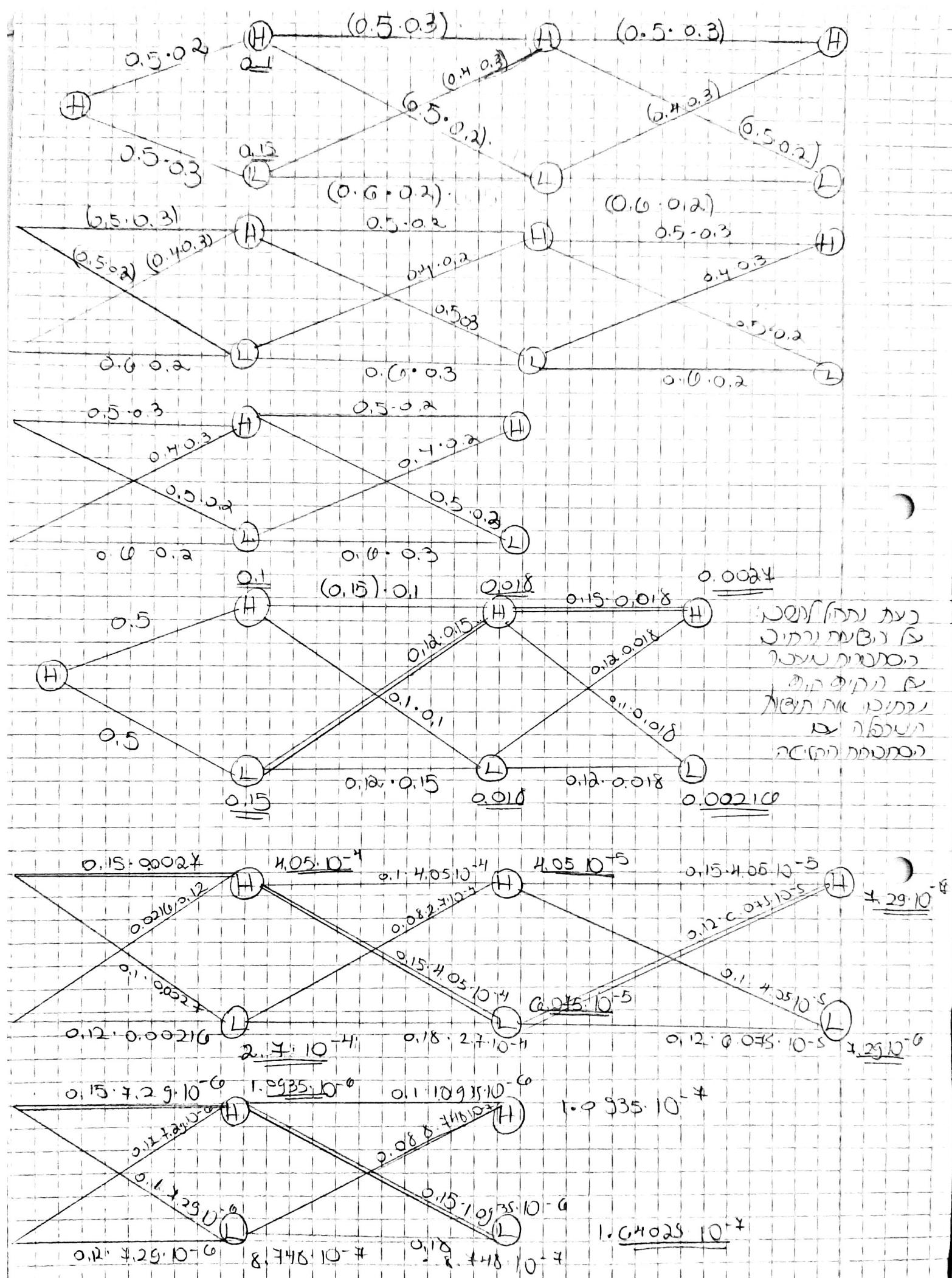
~~$P(H|L) P(H|H) P(G|H) P(H|H) P(C|H) P(H|H) P(A|H)$~~

IP(TIL) IP(HU) IP(HL) IP

~~IP(HIE)~~ ~~IP(GIH)~~ ~~IP(CIH)~~ ~~IP(~~

$$P(L|H) = P(L|G) \cdot P(G|H) + P(L|C) \cdot P(C|H)$$

$\text{IP(TIL)} \quad \text{IP(LIL)} \quad \text{IP(GLL)} \quad (\text{P(LIL}) \text{IP(CIL)}) \quad \text{IP(CLIL}) \text{IP(ALL)}$



probability: $1.64025 \cdot 10^{-7}$, LHLHLHLH : back tracking

3) Four-gram tagger:

$$IP(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = \prod_{i=1}^n q(y_i | y_{i-3}, y_{i-2}, y_{i-1}) \prod_{i=1}^n e(x_i | y_i)$$

$y_0 = y_{-1} = y_{-2} = *$ where * is the START symbol

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

$y_i \in \mathcal{K}$ for $i=1, \dots, n$ where \mathcal{K} is the set of all possible tags.

Viterbi algorithm that takes as input an integer n and finds

$$\max_{y_1, \dots, y_{n+1}, x_1, \dots, x_n} IP(x_1, \dots, x_n, y_1, \dots, y_{n+1})$$

Define a dynamic programming table

$\Pi(k, t, u, v)$ = maximum probability of a tag sequence ending in tags t, u, v at position k

Input: An integer n (not a sentence x_1, \dots, x_n),

$q(w|t, u, v)$ - transition probabilities

$e(x|s)$ - emission probabilities

Definition: Define \mathcal{K} to be the set of all possible tags.

Define $\mathcal{K}_{-2} = \mathcal{K}_{-1} = \mathcal{K}_0 = \{*\}$, $\mathcal{K}_k = \mathcal{K}$ for $k=1, \dots, n$

Define \mathcal{V} to be the set of possible words

In generative model we generate a sentence, and predict the most probable sequence of tags. In this case the output is the probability.

Initialization: Set $\Pi(0, *, *, *) = 1$

Algorithm: For $k=1, \dots, n$

For $t \in \mathcal{K}_{k-2}$, $u \in \mathcal{K}_{k-1}$, $v \in \mathcal{K}_k$

$$\Pi(k, t, u, v) = \max_{w \in \mathcal{V}, w \in \mathcal{K}_{k-3}} \{ \Pi(k-1, w, t, u) \cdot q(v|w, u) \cdot e(x_k|v) \}$$

$$\text{Return } \max_{t \in \mathcal{K}_{n-2}, u \in \mathcal{K}_{n-1}, v \in \mathcal{K}_n} \{ \Pi(n, t, u, v) \cdot q(\text{STOP} | t, u, v) \}$$

We are not given a sentence x_1, \dots, x_n , so we should find x_k from the vocabulary

Running Time: $n \cdot |\mathcal{K}|^3$ entries in Π to be filled in for $t, u, v \in \mathcal{K}$

$O(|\mathcal{K}| |\mathcal{V}|)$ time to fill one entry

for $w \in \mathcal{K}_{k-3}$ for $x_k \in \mathcal{V}$

\Rightarrow time to calculate $q(w|t, u, v) \cdot e(x_k|w)$ for all K, w, t, u, v is $O(n \cdot |\mathcal{K}|^4 |\mathcal{V}|)$

	Most likely tag baseline	HMM	Add-one smoothing	Pseudo Words	Pseudo Words & Add-one smoothing
Known Words Error Rate	6.10%	23.00%	17.60%	21.90%	19.40%
Unknown Words Error Rate	73.20%	78.30%	72.70%	52.60%	55.50%
Total Error Rate:	11.90%	27.90%	22.40%	24.60%	22.50%