<div dir="rtl">

| Structured Prediction |
| :---: |
| שיעורי מולב 3 |

עוסקים בבעיות חיזוי מובנות. ה-labels שונות יש להן מבנה כלשהו.

- חיזוי על מבנה של פלט.

- חיזוי מתוך מרחב פלטים (לפעמים אינסופי) של מבנים אפשריים.

נתמקד בחיזוי supervised.

קורא: Natural Language Parsing

דף: מבנה

דף: על פי מבנה את הקלטים במשפט, היות מילים ציוני, ותלויות

את מה יהא הקלטים אנחנו מזהה דרך כל?

1. נוסח חיזוי הקלטים ספציפית.

2. התיך את מבנים על דף של ימינה: חלק אחד מבנה מעין

ולא מסוד עובר את התץ ובו לא נתין את ה מבנים מראיין.

כוון: נחלק את הבעיה לאא, צריך הבעיה יותר פשוט-

אל אני את רק אומרת! ♢ = כל רק אפשר על בחלק חוק

בלחלק הבעיה כין את הבעיה.

<u>Sequence Prediction Problems</u> : (נתחיל עם הס 2 קלים-)

I <u>Named Entity Recognition</u> :

הבעה: בהנתן משפט, נרצה לסוון את חלק החלמים לפי שאני

ולהיך מתוך סוג את המילה כלל (... ,Location, Company)

יתא התבנית? בעיה מפורית. לנו אותו מתא לעת את הקלטים על לי ש

אלא מין מבנה ומלי.

התגית מפורית.

1. לסך התון בלונד לא נתפרש כי: המון את החיזוי כי לי אפשר ולנתן

מה שמנויו על משפט בתוך 20 הקלטים ע קלאין 21.

2. עם החיזוי הפרוד לא נדל prediction ויא גבוה לא פוה אלא.

</div>

II   <u>Part of Speech Tagging</u>: נרצה לתת לכל מילה (שם עצם) ( noun, verb, ...)

א. דוגמה: נרצה לחלק את (כל) המילים שייכני (כזן).

דוגמה: נרצה להבין לאיזה סוג מילה שייך כל פעם בכל מקום.

⟵ נרצה להבין את הקונטקסט של ההקשר.

אבל איך יש לנו? ?

1. (הסתברויות של מילים) סביר שיכירו ההסתברות המילה    saw
מופיע יותר כפועל ובויכין (הסתברות גבוהה).

2. מילים שכנות: יש רצפים של תגי נישה שהם יותר או פחות סבירים,
את (כזן) רוצה ת.ה.ע את להבחין את הקונטקסט שלו.

⟵ 3. מאפיינים מורפולוגיים: 'ly-' ⟶ adverb     סוף פתיח (סיור)

ואפשר לבנות מודל של המילה מהמבנה שלה.

הערה: ואפשר להגיע לדיוק גבוה סביר מאוד (∼90%) אם נמצא רצף מין
מילים בקורפוס. את שלשול תלקים מה ב-100 הנפרם ...

אבל ... נרצה להגדיר את המודל (הסתברותי) שלנו

הגדרה: סדרה (sequence) של משתנים אקראיים $Y_1, ..., Y_n$
היא שרשרת מרקובית ( Markov Chain), (הומוגנית, מסדרה ראשונה)

אם    $\forall i \quad P(Y_i | Y_1, ..., Y_{i-1}) = P(Y_i | Y_{i-1})$

$(\Longleftrightarrow) \quad Y_i \perp \{Y_1, ..., Y_{i-2}\} | Y_{i-1}$

$(\Longleftrightarrow) \quad P(Y_{1:n}) = \prod_{i=1}^{n} P(Y_i | Y_{i-1})$

(התפלגות משותפת של $Y_1, Y_2, ..., Y_n$)

הנחה נוספת: נסמן לצולי ההסתברות התחלתית $Y_0 = start$  p
שמשתמש בהם בת נא לא להציג בסים דבר - רק/נוחות.

הגדרה 2: סדרה של משתנים $X_1, ..., X_n$, $Y_1, ..., Y_n$ היא HMM
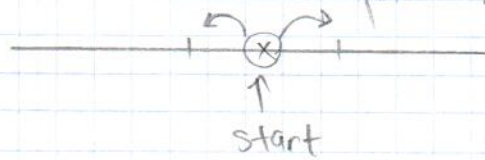( Hidden Markov Model) אם ההסתברות המשותפת שלהם נתונה:

$$P(Y_{1:n}, X_{1:n}) = \underbrace{\prod_{i=1}^{n} P(Y_i | Y_{i-1})}_{} \cdot \underbrace{\prod_{i=1}^{n} P(x_i | Y_i)}_{P(X_{1:n} | Y_{1:n})}$$

↑ מניחים ש-Y הם שרשרת מרקובית —
וכן שגם ההתפלגות של —

$P(Y_{1:n})$

start

דוגמא: הילוך אקראי על הישר:

$$P(Y_0 = 0) = 1 \qquad P(Y_i \mid Y_{i-1}) = \begin{cases} Y_{i-1}+1 & w.p.\ \frac{1}{2} \\ Y_{i-1}-1 & w.p.\ \frac{1}{2} \end{cases}$$

כלומר $\{Y_i\}_{i=1}^{n}$ הם סדרת מיקומים.

נניח שיש לנו $n$ סוכנים $X_1, \dots, X_n$ שכל סוכן מדווח לנו את המיקום, אך רק בקירוב

נאמר:

$$P(X_i \mid Y_i) = \begin{cases} Y_i & w.p.\ \frac{1}{2} \\ Y_i+1 & w.p.\ \frac{1}{4} \\ Y_i-1 & w.p.\ \frac{1}{4} \end{cases}$$

הנחות: 1. $Y_{1:n}$ הם שרשרת מרקוב.

2. (הנחה) $X_i$ פלוס רעש בכל דבר אחר חוץ מ-$Y_i$.

כלומר אם נדע את $Y_i$ אז אין עוד אינפורמציה ונסמן

$\oplus$ את $X_i$ המיקום שלו.

הנחה זו הגיונית תחת ניסיון אי התלות של $X$

כלומר נדע מיקומו שלו, (הרעש נפרד בין רעש

מעצם על הקרי לניסיון עם אחרים סביבו והיה בזה רעש קולי



זוהי הנחה 1 של
HMM

מעתה ונקרא סיכום ב: POS

$X_{1:n} = $ (הינ"ע מילים במשפט

$Y_{1:n} = $ (הינ"ע של POS של המילים במשפט $\longrightarrow$ (חלקי א' דבר

נניח שלכל $Y_{1:n}$, $X_{1:n}$ הם HMM. ונרצה לתייג

האם ההנחות מתקיימות? הם?

1. $Y_{1:n}$ הם שרשרת מרקובית: משמעות ההנחה שם היא

שכל את "תיוג" שלפנינו מילים, מספיק את "תיוג" של שתי מילים

אחרונות כדי להגיע למשפט. זה על פי ההגיון כי התיוג מאוד מקומי

אבל את כל אינו מספיק $\Longrightarrow$ (כאשר שרשר מרקובית

נאמר על א' צדדי יותר ($X$, $P$-$S$ אחרון שתאמה ...) אבל במובן

סטטיסטי, לא.

POS: נקודת התחלה של HMM (הנושא:

2. $X_i$ כאו ציר עמום בר מין ל-n. יהיה בי תמיד $Y_i$-n.

לא נאמיך! כי אם נולין על "לאן" מה המילה, אז אם ילך את הדר POS שלה, וזה תמיד עשיה יהיה מושר עלוי העינוי.

המרכים של HMM :

1. transition probabilities : על פלומים $(y,y')$, מה ההסתברות סיפר
   $t(y,y') = P(Y_i = y \mid Y_{i-1} = y')$   :כמה ? נורמל
   $\forall y' \quad \sum_y t(y,y') = 1$   :ונים להתקיים

2. emission probabilities : על ערך $Y_i$ נאשר $y$ ומהי $w$ של $X_i$,
   $e(w,y) = P(X_i = w \mid Y_i = y)$   ?מה ההסתברות
   $\forall y \quad \sum_w e(w,y) = 1$   :וגם כן לדרוש

| Bill | saw | that | man | yesterday | דוגמא: |
|------|-----|------|-----|-----------|--------|
| " | " | " | " | " | |
| $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | |
| Name | Verb | Conj | Noun | Adverb | |
| " | " | " | " | " | |
| $Y_1$ | $Y_2$ | $Y_3$ | $Y_4$ | $Y_5$ | |

$P(X_{1:n}, Y_{1:n}) = t(\text{Name}, \text{start}) \cdot t(\text{verb}, \text{Name}) \cdot \ldots$

$\cdot e(\text{Bill}, \text{Name}) \cdot e(\text{Saw}, \text{Verb}) \cdot \ldots$

לכן מה שקורה הוא זה לחשב את e או t לפניו יהיה לי ושם
מתוכנה דאטה (ונרחב לכך).

POS tags    vocabulary : Inference : HMM ב POS

$\{t(i,j)\}_{i,j\in L}$ , $\{e(w,i)\}_{\substack{w\in V \\ i\in L}}$ : HMM מודל של ים נציג לא כי יכן

של הבא ולגן סדרה $X_{1:n}$ (מילים) , (תגים):

$$y^* = \underset{y_{1:n}}{argmax}\ P\left(y_{1:n}, X_{1:n}\right)$$

$$= \underset{y_{1:n}}{argmax}\ P\left(y_{1:n}\mid X_{1:n}\right)\cdot P(X_{1:n})$$

נוכל לחשב זאת על כל נתיב של המילים הבאות: $L=\{l_1, l_2, ..., l_k\}$



$W_{y_{i-1}\to y_i} = t(y_i, y_{i-1})\cdot e(x_i, y_i)$ : כל נתיב אל מגדיר נרצה לחשב

וספונ המס חשבון של נתיב עם משקל מקסימלי אופטימלי

צריך לחשב : Viterbi Algorithm

$$\pi(t,j) = \underset{\substack{y_1,...,y_t \\ y_t = j}}{max}\ P\left(y_{1:t}, X_{1:t}\right)$$

$$\pi(t,j) = \underset{j'}{max}\left\{\pi(t-1, j')\cdot t(j, j')\cdot e(x_t, j)\right\}$$

כי נזהה את המשקל האופטימלי עד צומת t
היה משקל נתיב אופטימלי אל
כולל t-1 ואז כפל המשקלים.

$$\pi(1, j) = t(j, START)\cdot e(x_1, j)$$

כאשר המשקל את כולל צומת 1 ומגיע
מ-START

זאת נחשב עבור כל צומת נוכל למצוא נתיב אל זה.

 נגדיר MLE קריטריון : נחשב את

$$\ell\left(\{t(i,j)\}, \{e(w,j)\}\right) = \sum_{r=1}^{M}\sum_{k=1}^{n_r}\left[log\ t\left(y_k^{(r)}, y_{k-1}^{(r)}\right) + log\ e\left(x_k^{(r)}, y_k^{(r)}\right)\right]$$

עבור אוסף הדגימות $\{X_{1:n}^{(r)}, y_{1:n}^{(r)}\}_{r=1}^{M}$ , נמצא פתרון בצורה נקיה.

$t(i,j) = n_{ij}/\sum n_{kj}$ , $e(w,i) = n_{wi}/\sum n_{w'i}$ : ומה נספר ובזה

Maximum Entropy Markov Model : <u>MEMM</u> ⟵ POS

בא מאפשר יותר גמישות מ- HMM.

**שלב 1:** ב-MMM אנו נאלצים לקבוע מראש גיליון של ה-features, ואז

אנחנו מאמנים על מנת לגלות. אל על מנת פה נרצה אז רוצים או

חוסר אם נרצה - כי אם נאזה - כי אז נרצה.

כ-ם נרצה יכולים לקבוע מראש איזה מאפיינים ישפיעו.

בא נרצה מודל discriminative. אנחנו - רק נאפשר לעשות מה-X-ים שלנו

מן נלמד את X-ים למפות את ה-Y-ים.

$\Longrightarrow$ אם נרצה תלוי של X.

**נוסחה:**

$$P(Y_{1:n} \mid X_{1:n}) = \prod_{i=1}^{n} P(Y_i \mid Y_{i-1}, X_{1:n}) =$$

conditional
independance

נרמול
↓

$$= \prod_{i=1}^{n} \frac{e^{\langle \Phi(y_i, y_{i-1}, X_{1:n}, i), w \rangle}}{z(X_{1:n}, Y_{i-1})}$$

כאשר $\Phi$ זו פונקציית ה-features שלנו ו-$w$ יהי הפרמטרים של המודל

APML
פרק 4

ההסתברות הסופית ל- $\langle \Phi(y_i, y_{i-1}, X_{1:n}, i), w \rangle$ זו (מסמלת) נתן score.
$\mathbb{R}^d$ $\mathbb{R}^d$

אז מי הוא $\Phi$? כיצד אנו מגדירים את ה-features שלנו?

$\underline{\Phi}$ הוא וקטור ורכיב שונים ב-$\mathbb{R}^d$

$$\Phi(y_i, y_{i-1}, X_{1:n}, i) = \begin{pmatrix} \vdots \\ \mathbb{1}_{\{y_i, y_{i-1}\}} \\ \vdots \\ \mathbb{1}_{\{x_i, y_i\}} \\ \vdots \\ \mathbb{1}_{\{pref, y_i\}} \\ \vdots \end{pmatrix}$$

transition features
$\{y_i, y_{i-1}\}$ מאיזה תיוג לאיזה תיוג $\}$ (מספר POS)$^2$

emission features
מאיזה מילה לאיזה POS $\}$ מספר המילים × מספר POS

prefixes $p_1, \dots, p_m$
$\{pref(x_i) = p_j, y_i = y\}$ $\}$ מספר POS × m

של אותו נרצה w מספר שנאמר ל- transition features יהיו של

כל זוג תיוג אל תיוג $y_{i-1}, y_i$ נותן מ.ס. מילולית למיקום

ומן של של (אלא לזה אלא תלוי מאנא) אל המון למשל למיקום נותן.

משל שנאמר ל- emission, נלמד משל, אם התכונה למשל (man, verb)

אם משהו, זה יהיה ערך הזה נותן נותן את ותו...

Gradient Ascent :שיטה אחת לפתור זאת בעזרת

$$\nabla LL = \sum_{r=1}^{N} \sum_{i=1}^{n} \left[ \Phi(y_i^{(r)}, y_{i-1}^{(r)}, x_{1:n}^{(r)}, i) - \sum_{y'} P(y'|y_{i-1}, x_{1:n}) \Phi(y', y_{i-1}, x_{1:n}, i) \right]$$

instances ← $r$     transitions ← $i$

$\underbrace{\qquad\qquad}_{\text{Expected value}}$

## Structured Perceptron:

מקומות רבים אנו רוצים את כל ה- lables ביחד, ולכן ההבדל ל-

כולה ולא רק לרצף יחיד את המון הערך (פועל היה לא רק אך בחשבון)

מה ההבדל?

מקומות רבים אנו כבר יודעים, ורוצה את הפתרון כולו

כולן של הפתרון ובחשבון וכדומה.

על בחשבון לפתרון לדעת אומר אנחנו צריך גם קחם לזה,

ושרל בחשבון.

וגיות לפי זו מודלים score :כי נתן לפי על סוכם, פעם נקרא על ל פעם הולכים הרומים

$$Score(y_{1:n}, x_{1:n}) = \sum_{i=1}^{n} Score(y_i, y_{i-1}, x_{1:n}, i) =$$

$$= \sum_{i=1}^{n} \Phi(y_i, y_{i-1}, x_{1:n}, i)^T \cdot w$$

... MEMM-ב שעשינו אותו נגדיר פעולה זו של את פסיק)

חזרה - נגדיר את $w$ הנלמד <= פרמטרים, ואיך ניתן ואז לפתור את זה

בעזרת inference בלבד.

## Perceptron for Sequence Labeling:

Input: $\left\{ (x_{1:n}^{(r)}, y_{1:n}^{(r)}) \right\}_{r=1}^{N}$

1. $w \leftarrow \vec{0}$ /random

2. for $r = 1 \ldots N$ :

2.1    $\hat{y}_{1:n} \leftarrow$ find the seq of ys that maximizes the score = do inference

2.2    $w \leftarrow w \cdot \eta \left( \underbrace{\sum_{r=1}^{N} \Phi(y_i^{(r)}, y_{i-1}^{(r)}, x_{1:n}^{(r)}, i)}_{\text{the correct answer}} - \sum_{r=1}^{N} \Phi(\hat{y}_i^{(r)}, \hat{y}_{i-1}^{(r)}, x_{1:n}^{(r)}, i) \right)$

המחשבות אך $\nearrow$ הפעם עושים $P(y_{1:n}, x_{1:n})$

3. return $w$

הפעמים לרצף לחדש את המודל על זו

ה-$w$ כולל פעמים וכולל המודל

<u>אופטימיזציה של את הפרספטרון עובד:</u>

נתון מניחים של w, ואז בכל פעם עשה inference

ונשווה בין הפלט, לבין מה שצריך לין שהתקבל:

אם זה אותו דבר, אל תעשה כלום

אם יש הבדל בין הפלט, לבין הפלט, (עדכן את w.

זה דומה לזה שעושים ב-GD, אבל במקום לקחת את התוחלת בפועל,

לוקחים את הדגימות בזמן נתון מסויים