

מבואות

מידול שפה - בהינתן רצף מילים מה היא המילה הבאה: $\mathbb{P}(w_1|\text{START}) \cdot \mathbb{P}(w_2|w_1) \cdot \dots \cdot \mathbb{P}(w_n|w_1, \dots, w_{n-1})$
קורפרוס - "החתול נפל אבל הוא חתל על הרגליים"
המילים "החתול" ו-"הוא" מתייחסות לאותה הישות.
קאנפורמה - קודם אומרים את הביטוי ואז את האזכור.
קאנפורמה - קודם אומרים את האזכור ואז את הביטוי.

מרחק עריכה \ ליוניטטי -

מטריקה להשוואה בין מחרוזות:

- הוספת תו	↑
- מחיקת תו	→
- יישור מחרוזות והשוואה	↗

לבחור את המינימאלי מבין השלושה.

- כל תו מלבד שורה חדשה	cat file
- מילה, ספרה, מרחב לבן	lw ld ls
- לא מילה, ספרה, מרחב לבן	lW lD lS
- אחד a, b או c	[abc]
- לא a, b או c	[^abc]
- תו בין a לg	[a-g]
- מחרחיל \ נגמר במחרוזת	^abc\$
- תחילה בגבולות \ לא מילה בגבולות	^b\b
- 0 או יותר, 1 או יותר, 0 או 1	a* a+ a?
- בדיוק חמש, שתיים או יותר	a{5} a{2,}
- בין אחד לשלוש	a{1,3}
- תואם abt ולא cdt	ab cd
- הופך תו מיוחד לתו רגיל	\. \. \. \.

- כותב את כל הקובץ "file" למסוף	cat file
- מציג את הקובץ "file"	more file
- עמוד בכל פעם	less file
- גרסה מוגזנת יותר של "more", אך פחות נפוצה	head -30 file
- מציג את 30 השורות הראשונות	tail -30 file
- מציג את 20 השורות האחרונות	wc file
- סופר שורות, מילים ותווים בקובץ	grep '[A-M]' file
- מדפיס את השורות המכילות אותיות גדולות בטווח A עד M	sort file
- ממין את הקובץ בסדר אלפביתי	sort -n file
- ממיין את הקובץ באופן מספרי. 12 או אחרי 2.	sort -r file
- ממין את הקובץ בסדר הפוך.	sort -u file
- מסיר שורות כפולות, ומבטיח שכל שורת פלט היא ייחודית.	uniq file
- מחק שורות סמוכות אחת או יותר, מציא רק אחת מהם.	sort file.txt
- חושם לפי כל שורה את גודל הבלוק.	tr 'A-Z' 'a-z'
- מחליף uppercase ל-lowercase-	tr -d ' '
- מחק רווחים	tr -dc ' '
- משאיר רק רווחים. c עבור complement	tr -s ' '
- מחקב רצף תויו רווח להית רווח אחד	diff file1 file2
- משווה בין שני קבצים	tr -sc 'A-Za-z' '\n' < file1 tr 'A-Z' 'a-z'
- מביא תמניות tokens לללא תלות uppercase	[upper box] sort uniq

תמניות Token - כל המילים בטקסט כולל כפיליות.

תמנית Type - אוצר המילים בטקסט, מופע ייחודי.

צורת יסוד lemma - "להלכת" - הלך, "are, is" -

בגעול stem - ניהוש צורת יסוד עבור שפה לא מוכרת: "stories" - stori

Type Token Ration (TTR)	תמניות	תמניות
התפלגות מילים בטקסט היא ידפיסאית -		
עבור המילה w_n הניתת בשכיחות: $\alpha_n \propto \frac{1}{w_n}$		
סיווג מסמכים		
מטריצת בלבול -		
	חזיו	אמיתי
	+	+
	FN	TP
	TN	FP

Accuracy - כמה המודל צודק $\frac{TP+TN}{TP+TN+FP+FN}$
Precision - כמה המודל פוגע $\frac{TP}{TP+FP}$
Recall - כמה המודל תופס $\frac{TP}{TP+FN}$
מטריקת F - ממוצע הרמוני בין P ו- R : $F_\beta = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$
כדי לחזק את R מגדילים את β .
סיווג להרבה תגים (מאקרו) - חישוב P, R, F עבור כל אחד מהתגים ואז נמצא את F . מטיב עם התגים הקטנים.
סיווג להרבה תגים (מיקרו) - נחשב P, R עבור הכל ביחד ונחשב F אחד, התגים הגדולים מקבלים יותר משקל.
פיצ'רים אפשרים במערכת סיווג: ספירת מילים, רשימות מילים המגדרות בעולות תכונה-ניתן לספור הופעות, מספר מופעות, סימני פיסוק, כבנים תחביריים.

n-gram (מידול שפה) - מחשיבים רק את $n - 1$ המילים שמקדם למילה ה-ית בשביל לחזות אותה: $\mathbb{P}(w_k|w_{k-n+1}, \dots, w_{k-1}) = \frac{\text{count}([w_{k-n+1}, \dots, w_{k-1}, w_k])}{\sum_{w'} \text{count}([w_{k-n+1}, \dots, w_{k-1}, w'])}$

חוק בייס - $\mathbb{P}(c|d) = \frac{\mathbb{P}(d|c)\mathbb{P}(c)}{\mathbb{P}(d)}$
אלגוריתם Naive base -
קלט: מסמך d , אוסף תגים $c = \{c_1, \dots, c_k\}$
קלט: אימון $\{(d_1, c_1), \dots, (d_n, c_n)\}$
פלט: פונקציית מסוג $c \rightarrow d$
עבור $\mathbb{P}(c_i) = \frac{\#c_i \ln \text{train}}{\text{size of train}(n)}$
מניחים כי המסמך הוא bag of words $d \equiv (w_1, \dots, w_m)$
ועבור $\mathbb{P}(d_j|c_i) = \mathbb{P}((w_1, \dots, w_m)|c_i) = \prod_{j=1}^m \mathbb{P}(w_j|c_i)$

כאשר $\mathbb{P}(w_j|c_i) = \frac{\#w_j \ln c_i}{\text{\#all words in } c_i}$
נקבל פונקציית המסוג: $\hat{c} = \text{argmax}_{i \in \{1, \dots, k\}} (\mathbb{P}(c_i) \cdot \prod_{j=1}^m \mathbb{P}(w_j|c_i))$
ובגוף: $\hat{c} = \text{argmax}_{i \in \{1, \dots, k\}} (\log \mathbb{P}(c_i) + \sum_{j=1}^m \log \mathbb{P}(w_j|c_i))$

אם מילה לא מופיעה באחד התגים אז בעיה, כי האקסמג יאיפוס וכן נבצע החלקה עם 1: $\mathbb{P}(w_j|c_i) = \frac{1 + \#w_j \ln c_i}{1 + \text{\#all words in } c_i}$

מערכת סיווג

מרכיבים -
1. ייצוג לקלט (פיצ'רים) f
2. פונקציית הסיווג מחזירה \hat{y}
3. פונקציית הפסד
4. אלגוריתם למידה (SGD)

סיגמיויד - $\sigma(z) = \frac{1}{1+\exp(-z)}$
תכונות סיגמיויד - $\sigma(-z) = 1 - \sigma(z)$
גרסיה לוגיסטית- פונקציית סיווג $\hat{y}^{(i)} = \sigma(f \cdot \theta^{(i-1)})$
אלגוריתם SGD -
נחשב $\hat{y}^{(i)} = \sigma(f \cdot \theta^{(i-1)})$ ונעדכן אם צריך.
פונקציית ההפסד:

$L_{CE} = -\log(\mathbb{P}(y|x)) = -[y \cdot \log(\sigma(\theta \cdot f + b)) + (1 - y) \cdot \log(1 - \sigma(\theta \cdot f + b))]$
נגזר פונקציית הפסד: $\frac{\partial L_{CE}}{\partial \theta^{(i-1)}} = (y^{(i)} - y) \cdot f$
ונכל לעדכן: $\theta^{t+1} = \theta^t - \eta \cdot \nabla L_{CE}(\theta^{(t-1)})$
גרורליזציה - לבלום במפורש את ערכי הפרמטרים וגם להתמודד עם overfitting: $L = L_{CE} + L_{reg}$
מרהשיא. $L_{reg} = ||\theta||_2^2 - L_2$
מרהשיא. $L_{reg} = ||\theta||_1 - L_1$
מענישה מודלים שיותר מדי פרמטרים שלהם 0.

מודל overfitting - התאים את עצמו יותר מידי למדגם אימון (או שפשוט אין בו מספיק פרמטרים)

גרסיה מורבת קלאסים - נחשב את הציון לכל קלאס ובסוף נגרמל בסכום כל הציונים. נחזיק גם וקטור משקולות עבור כל קלאס. פונקציית ההפסד:

$$L_{CE} = -\log \left(\sum_{k=1}^K \exp(w_k \cdot x + b_k) \right)$$

נגזר פונקציית הפסד: $\frac{\partial L_{CE}}{\partial w_{k,i}} = - \left(\frac{\exp(w_k \cdot x + b_k)}{\sum_{k=1}^K \exp(w_k \cdot x + b_k)} \right) \cdot x_i$
כל הפרמטרים מתעדכנים שכן softmax לא פגע לא מחזיר 0 או 1.

פרסטנטון - ננסה לייצר רק מבחין בין דוגמאות חיוביות ושליליות, אם צדקנו אין עדכון אחרת נעדכן לפי הקלאס הנכון.

תיגו רצפים

בעיות - תיגו חלקי דבר, תיגו תכונת מורפחתביריות, זיהוי יישויות.

תהנית מורפחתביריות - מין, גיל, זמן, בנין, גוף, ויחידורים, voice (active/passive).

חלקי דבר -
open class:
שם תואר adj, תואר הפועל adv, שם עצם noun, פועל verb, שם פרטי pron, קריאות itn (היי, וי),
closed class:
מילות יחס aux, פעלי עזר aux, מילים שמחברות בין חלקי משפט שווים cconj, מילות חיבור משעבודות (כלומר מחברות בין משפט פחות חשוב לאחד יותר חשוב) sconj, כינוי גוף pron, מה שיש רק באנגלית particles (wake up), מספרים num, מילת יחס prep.
other class:
ממלים sym, כל השאר x.

PTB - הקורפוס המסיבי הראשון הגדול באנגלית.

תיגו רצפים:
דיוק ברמת המילה (accuracy)
 F_1 - רצף התגים (מאקרומיקרו)

מטריקות זיהוי יישויות:
טעינו איפשהו ביישות – טעות ככולה
חישוב F_2 לכל היישויות או לאחת
סכימת BIO - תחילת ישות, אומצע O לא חלק.
שרשרת מרקוב - מניחים כי כל מצב תלוי רק במצב שקדם לו: $\mathbb{P}(q_i|q_{i-1}, \dots, q_{i-1})$
נעטרך גם את ההסתברות להיות במצב הראשון π .
מודל מרקוב נסתר HMM - מודל גרנטרי, אחיה רק את המילה (הפלט) ונאחזתו רצף להבין איזה מצב (חלק דבר) יצר אותה. המצבים הם ה"נסתרים". מניחים כי פלט תלוי רק במצב שלו.

ההסתברות של רצף התגים בהינתן רצף המילים:

$$\hat{t}_{1:n} = \text{argmax} \mathbb{P}(t_1, \dots, t_n|w_1, \dots, w_n) \approx \text{argmax} \prod_{i=1}^n \mathbb{P}(w_i|t_i) \mathbb{P}(t_i|t_{i-1})$$

אלגוריתם ויטרבי - תכנון דינמי, נרצה למצוא את השקדם הכי סביר בין המצבים שייצח את התצפיות (המילים במסמך).

קבוצת N המצבים (מיקום המילה במשפט)	$Q = q_1, \dots, q_N$
מטריצת מעברים כך a_{ij} שווה להסתברות לעבור ממצב i למצב j	$A = a_{11}, \dots, a_{NN}$
קבוצת T התצפיות (מילים) מתוך אוצר המילים	$O = o_1, \dots, o_T$
סיכוי קבלת תצפית o_t במצב q_t	$B = b_t(o_t)$
התפלגות התחלתית מעל המצבים	$\pi = \pi_1, \dots, \pi_N$

ונכל למצוא:

$$v_t(j) = \max_{q_1, \dots, q_N} \mathbb{P}(q_1, \dots, q_{t-1}, o_1, \dots, o_t, q_t - j | \lambda) \\ = \max_i v_{t-1}(i) \cdot a_{ij} \cdot b_j(o_t)$$

ונעבדו בגוף:

$$v_t(j) = \max_i \log v_{t-1}(i) + \log a_{ij} + \log b_j(o_t)$$

דוגמה -

$$O = \{o_1 = \text{שרה}, o_2 = \text{שיר}, o_3 = \text{שמח}\}$$

$$Q = \{q_1 = V, q_2 = N, q_3 = M, q_4 = P, q_5 = J\}$$

$$B: \text{Emissions}(\log)$$

	שרה	שיר	שמח
V	-3.912	-6.215	-6.908
N	-5.298	-4.605	0
M	-3.507	-3.912	-6.908
P	0	-11.513	0
J	0	-3.219	0

→	V	N	M	P	J
V	-2.996	-1.050	-1.386	-1.609	-1.897
N	-0.511	-2.996	-2.996	-2.303	-1.609
M	-0.511	-2.303	-2.303	-1.897	-2.996
P	-2.996	-0.693	-1.204	-4.605	-1.966
J	-0.693	-2.996	-2.303	-1.605	-1.897

התפלגות התחלתית π :

	V	N	M	P	J
V	-1.897	-1.386	-1.386	-1.204	-2.996

חישוב:

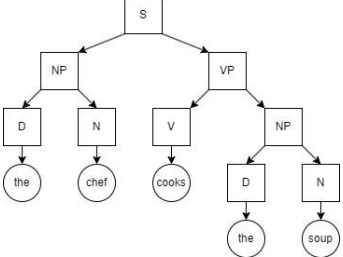
←	שרה	שיר	שמח
V	-1.897	-1.897	-3.912
N	-5.809	-6.685	-6.685
M	-5.288	-5.288	-6.685
P	-1.386	-1.386	-4.893
J	-3.507	-3.507	-4.893

בשורה עבור "שמח" נפעל בדומה, ונקבל כי הערך המקסימלי $-19.799 = V_4(5) = 5$ אשר התקבל ממצב N בתור תא \leftarrow לכן, תיגו האלגוריתם יהיה שמח=שיר=N, שרה=V, שיר=M.

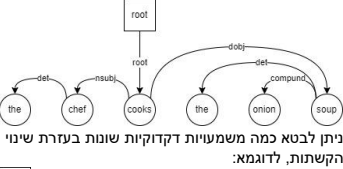
CRF - מודל דיסקרמיניטיבי. כל תג תלוי בתג הקודם במיקום עצמו ובכל הקלט (לצורך פיצ'רים)

מידול מבנים תלויים (תחביר)

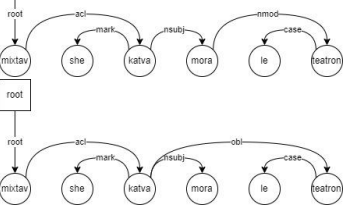
דקדוק חרי הקשר CFG - מוגדר על ידי:
קבוצת הסמלים הלא סופיים (צ'רופים/תגים) N
קבוצת הסמלים הסופיים (מילים) Σ
כללי גזירה מהצורה $A \rightarrow B$ כאשר A לא סופי
סמל התחלתי וחלק $N \ni S$
מוסקין CNF - כל כלל כולל מימין שני סמלים לא סופיים או מילה אחת סופית.



עצי תלויים -



ניתן לבטא כמה משמעות דקדוקיות שונות בעזרת שינוי הקשקות, לדוגמה:



הטיליות של עץ תלויים - תכונה לפיה קשקות העץ לא חוצות זו את זו.

סוגי קשקות בעץ תלויים -

קשקות בין שושא (פרדיקט) לארגומנט ליבה:

- נושא שמני nsbdj

- נושא פסוקי csbdj

← "ללמוד הרבה" ("ללמוד הרבה")

- מושא ישיר dobj

- מושא עקיף iobj

← "he gave her the book" ("her")

- פסוקים משלימה ccomp\lxcomp

← "לסדר את החדר" ("לסדר את החדר")

- קשקות בין שושא למרכיבי עזר:

- נספחים שמניים obl

← "הוא הלך הביתה אתמו" ← המילה אתמוול היא לא

מרכיב ליבה הכרחי של המילה הלך, כלומר גם בלעדיה המפשט תקין)

- פסוקית אופן advcl (במהירות, ביעילות)

- פועל עזר aux cop

- קשקות בתוך צירופים שמניים:

- תוניתו הדיוע determiner כמו the

- סמיכות nmod כמו "כיתת לימוד"

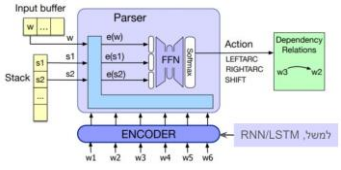
- case ony

ניתוח תלויים במעברים - החלטות מקומיות 3 מבנים:

- buffer: מחזיק את כל המילים שטרם טיפלו

- stack: מועמדים לקבלת קשת

- list of edges: הקשקות שחושפנו לפעך



מעברים -

- shift: העברת המילה מהזחלת buffer stack

- left-arc: הוספת קשת מהמילה העליונה במחסנית לזו

התחתית ונזכיר את השנייה

- right-arc: הוספת קשת מהמילה השנייה לעליונה

והוצאת העליונה (אי אפשר לקבלת קשקות חצות, רק הטליות)

- parser: לומד איזו החלטת לבצע בהיתן:

- מצב של המערכת (חץ, מחסנית, קשקות קיימות)
- החלטת מועמדת (הזח, קשת-מין, קשת-שמאל)
- eager** - יוצרים קשקות בין ראש stack התחילת buffer. כאשר מייצרים קשקות מינה ברגע שאפשר בלי להיפטר מהמילה העליונה (עם הגדרת פעולה חדשה להוציא מראש המחסנית reduce).
- oracle** - קלט: עץ מוכן, אוסף קשקות לא סדור. פלט: סדרת החלטות שהובילה ליצירת העץ מהקלט. תיגנת הציונים היא מעברים ויא בעיית סיווג שניתן לבצע ע"י נסכתל על עץ שלם ניתן ציון כולל. מתבסס על ההנחה שציון קשת לא תלוי בקשקות האחרות. נגדיר ציון לכל קשת, לכל זוג מצתים, לכל הכיוונים, וכלל נוג קשת אפשרית וגם לכל העולה מעדכנים את $\theta(n^3)$
- מציאת עץ פורש מקסימלי** - ניתן ציון לכל הקשקות, נמצא לכל צומת את הקשת המקסימלית הנכנסת אליו. אם יש מעגל:

- נפחית מכל קשת הנכנסת לצומת את המקסימלי שנסכס לצומת
- נכוון את המעגל לצומת בודד ונקרא
- נרחיב מחדש את העץ הנוצר ונבחר את הקשת לפי העץ המוקטן

ציון התאמה ללא תויות (UAS) - אחוז המילים שקיבלו

את הראש הנכון.

ציון התאמה עם תויות (LAS) - אחוז המילים שהוקצו לזן הראש הנכון ותג התלות הנכון (או "תויות"). לכן ממחיר יותר UAS ומתקיים UAS>LAS.

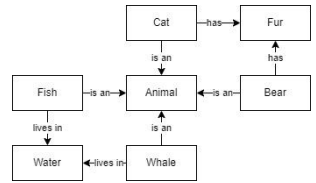
שיכוני מילים (word embeddings)

שיכוני מילים - וקטור השיכונוים של מילה מסוימת במרחב הקוטרי.
פוליטמיה - מילה אחת עם כמה משמעויות (עכבר).
מילה אחת, הקשורים שונים. וקטור ממוצע של וקטורים רחוקים.

מילים נרדפות - כמה מילים עם אותה משמעות (ורחלובנה). כמה מילים, אותו הקשר. וקטורים קרובים.
מילים דומות - מתארות חפיצא או רעיונות דומים בעולם (לבלחחחל). כמה מילים, הקשורים דומים. וקטורים קרובים עם קשר אלגברי.

מילים קשורות - מתארות חפיצא או רעיונות הקשורים לאותו שדה מסמני (לבלחחלונה). מילים בהקשר זו או זו.

וקטורים עם קואורדינטות מסוימות דומות.
רשתות סמנטיות - מילים (lemmas) ממופות למובנים (senses) והמובנים מקושרים ביניהם בצורת גרף.



האם אוסף השינויים באמת מייצג יחס דמיון לשוני אמין?

- ניתן לשאול אנשים עד כמה המילים דומות. זה סובייקטיבי.
- נשתמש בדמיון קוסינוס: $\frac{v \cdot w}{\|v\|_2 \|w\|_2} = \cosine(v, w)$
- שילי עובר skip-gram ו-PMI.

מטריצת מילים/מסמכים - ניצג כל מילה לפי אוסף המסמכים שהיא מופיע בו וכמה בכל מסמך. חסרונות:

- וקטרים ארוכים לכל מילה.
- מילה כמו the תהיה דומה לשאר המילים functionable ויהיה קשה ללמוד עליהן מנה.

מטריצת שכנויות - נספור הופעות של זוגות מילים באותו מסמך. w היא המילה שמעניינת אותנו, c הוא ההקשר.

PMI - זוהי נוסחה שעוזרת לנו להתגבר על כך שחלק מהמילים שכיחות יותר וחלק מהמילים פחות בכך שאנחנו מנרמלים את ההסתברויות:

$$\mathbb{P}(w, c) = \log_2 \frac{\mathbb{P}(w, c)}{\mathbb{P}(w) \cdot \mathbb{P}(c)}$$

כאשר: $\mathbb{P}(w, c)$ - ההסתברות ש- w תופיע במסמך עם הקשר c ; $\mathbb{P}(w)$ - ההסתברות ש- w תופיע במסמך; $\mathbb{P}(c)$ - ההסתברות ש- c תופיע במסמך. כשהתוצאה תהיה שווה לאפס זה אומר שיש אי-תלות בין המילים.

PPMI - ביצוע RELU על PMI (מבטיח רק חיוביות)

TF-IDF - גישה מקובלת באחזון מידע, מהפשים

שאלתה במנוע חיפוש ורוצים לדעת איזה מסמך להחזיר:

$$TF_IDF_i(t, d) = \log(1 + f_{t,d}) - \log\left(\frac{N}{n_t}\right)$$

כאשר:

$f_{t,d}$ - כמות הפעמים שהמילה t מילה מתוך השאלתה הופיעה במסמך d

n_t - כמות המסמכים שהמילה t הופיעה בהם

N - כמות המסמכים

החיסור נועד לתת משמעות למילים נדירות

שינון ע"י חיזוי הקשרים קומי (skip-gram) - נפעיל אלגוריתם חיזוי, אבל נשמור רק את המשקולות הנלמדים.

האלוריתם:

בהינתן מילה w , ננסה לחזות את השכנות שלה c

בהסתברות גבוהה לכל הניתן, ונעבור כך על כל המילים בקורפוס. בעצם נחשב מה ההסתברות ש- w מופיעות ביחד $|w, c|$.

לכל מילה יהיה וקטור מטרה בתפקידה כחזזה, ווקטור הקשר בתפקידה כנחזית.

היפר פרמטרים - רחב החלון (בכל פעם נדגום רק "חלון" בסביבת המילה וכן נגדיר מה יהיו מילות ההקשר)

נמנת הדוגמאות השליליות לכל דוגמה חיובית מימד וקטור השינון

α - עבור דגימת המילים השליליות מתוך התפלגות שמוסתחת לטובת מילים נדירות:

$$\mathbb{P}_\alpha(c) = \frac{(\#c)^\alpha}{\sum_c (\#c)^\alpha}$$

$\alpha = 0.75$ זה ערך טוב.

יעד חיזוי - נשתמש ברגרסיה גליסטית על מנת לחשב את ההסתברות של $\mathbb{P}(w|c)$ כאשר על כל דוגמה חיובית ניקח כמה שליליות:

$$\mathbb{P}(c + |w, c_{pos}) \cdot \prod \mathbb{P}(c - |w, c_{neg})$$

פונקציית ההפסד:

$$L_{CE} = - \left[\log(\sigma(c_{pos} \cdot w)) + \sum_{i=1}^k \log(\sigma(-c_{neg} \cdot w)) \right]$$

נגזור פונקציית הפסד:

$$\frac{\partial L_{CE}}{\partial c_{pos}} = [\sigma(c_{pos} \cdot w) - 1] \cdot w$$

$$\frac{\partial L_{CE}}{\partial c_{neg}} = [\sigma(c_{neg} \cdot w)] \cdot w$$

$\nabla L_{CE} = [\sigma(c_{pos} \cdot w) - 1] \cdot c_{pos} + \sum_{i=1}^k [\sigma(-c_{neg_i} \cdot w)] \cdot c_{neg_i}$

ונכלי לעדכן: $w^{(t+1)} = w^{(t)} - \eta \cdot \nabla L_{CE}(w^{(t)})$

מטריקות לשערוך מודלי דמיון - קלט: שתי רשימות מדורגות של חזונו הנשפטים.

פלט: מספר בטוח $[1, -1]$, שמתאר את הקשר שבין הרשימות

pearson correlation - $\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma(X) \cdot \sigma(Y)}$

אפשר לכוון ולהרחיב את התחום. רגיש לפיזור הערכים בשונים בתוך X .

spearman correlation - $r_s = \rho(R(X), R(Y)) = \frac{\text{cov}(R(X), R(Y))}{\sigma(R(X)) \cdot \sigma(R(Y))}$

עבור R פונקציית דירוג שממיינת את הצינון של כל אחת מהרשימות.

רשת ניורונים

רשת בהיזן קדמי FFN - המון שכבות של ניורונים פשוטים בתוספת פונקציית אקטיביצה ביניהם.

פרספטורן רב שכבות MLP - אותו דבר רק שהניורון הפשוט הוא פרספטורן.

אקטיביות אפסיות -

- RELU - בחירה טובה, גזירה קלה
- סיגמויד - ערכים בין 0 ל-1
- tanh - ערכים בין -1 ל-1

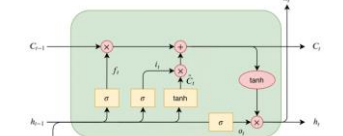
אימון הרשתות - מציאת הנגזרות החלקיות לפי loss ומפעפעים אחורה (back propagation) או כלל

(השרשרת).

מגבלות הFFN - התייחסות לכל הקלט בבת אחת, אין חשיבות לסדר קלט, מוגבל לחיזוי תג אחד.

רשת ניורונים נשנית RNN - בFFN לא עובר מידע בין שכבות עקבות, ואילו בRNN כן.

LSTM - גרסה משופרת לRNN. מכיל מרכיבים שנועדו לשמור על מצבים רחוקים (long short-term memory).



GRU - פשרה בין RNN לLSTM. מהיר בהרבה מ-LSTM.

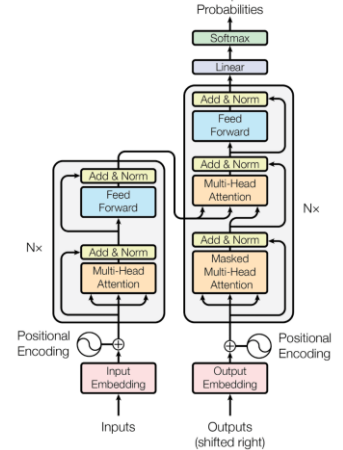
Bidirectional RNN - יודע להבדיל בין כיוונים. לכל כיוון יש את הפרמטרים שלו וה"חיבור" נעשה ברמה הבאה.

רולריציט dropout - בכל שלב אימון נאפס חלקים אקראיים בשכבה ונתעלם מהם.

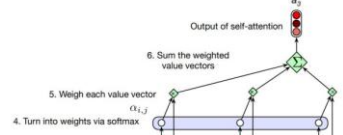
רשת אוטורגרסיבית - רשת שנעה קדימה, כלומר לומדת רק מקלט העבר ואינה יכולה להסתכל לעתיד.

רובוטריקים

מבנה הרובוטריק -



self-attention -



softmax -

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)}, \quad i \in [1, \dots, k]$$

השלים בself-attention -

- השוואה בין הווקטורים בקלט ומתן ציון לדמיון שבהים.
- נרמול הצינונים שקיבלנו ע"י שימוש בsoftmax.
- חישוב הפלט הנוחכי

כל מילה שואלת כמה היא צריכה לקחת מכל מילה ברמת ההיצוג הבאה.

Query - המילה המתשאלת

Key - המילה כגורם השואה

Value - המילה כמרכיב של שכבה הבאה

$$\text{SelfAttention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) \cdot V$$

כאשר:

$$Q = X \cdot W_q, \quad \dim(W_q) = d \times d_k$$

$$Q = X \cdot W_k, \quad \dim(W_k) = d \times d_k$$

$$Q = X \cdot W_v, \quad \dim(W_v) = d \times d$$

לדוגמה: בהינתן 3 מילים x_1, x_2, x_3 ואם מעוניינים ב- y_1, y_2, y_3 בהתאמה, נקבל לדוגמה עבור y_3 :

$$y_3 = \sum_{i=1}^3 (x_i \cdot k(x_i)^T \cdot v(x_i))$$

הקשרים השיריים (residual connection) - מתנים לשכבות העליונות גישה ישירה למידע המגיע מהשכבות הנמוכות.

$$z = \text{LayerNorm}(x + \text{MultiAttention}(x))$$

$$y = \text{LayerNorm}(z + \text{FFN}(z))$$

שכבת הנורמליזציה - נחשב שונות ותוחלת:

$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i, \quad \sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$

ואז ננרמל: $\hat{x} = \frac{(x-\mu)}{\sigma}$

צומי מרובה ראשים (multi-head attention) - כמה שכבות של self-attention שנקראות ראשים, לכל ראש יש סט מטריצות משלו. כל ראש יכול למצוא קשתות שונות בין המילים במשפט.

(מוסיפים לשכבת הצומי קלט מקימי המילים בעזרת positional embeddings).

למידה בהעברה (transfer learning) - למידת מערכת למטרה מסוימת ושימוש באלמנטים הנלמדים שלה למשימה אחרת (כמו skip-grams).

fine tuning - עדכון הפרמטרים שנלמדו למשימה המקורית על גבי המשימה החדשה.

מקודד (encoder) - לוקח טקסט (או אובייקט אחר) ומחזיר וקטור

מעננה (decoder) - לוקח וקטור (או כמה) ומחזיר טקסט

מקודד רובוטריק - כמו קודם אבל יכול לראות את כל הקלט

מעננה רובוטריק - יכול להסתכל רק אחורה בפלט, אבל יכול להסתכל על כל הקלט בעזרת דגירת cross-attention.

הלימדה כאן באמצעות קרוס אנטרופי רגיל פחות log

הסתברות של המילה. נכריח את המודל בכל שלב לקחת את המילה הבאה האמינית / חיפוש אלומה.

שני סוגים של מודלים ברובוטריק -

- מודל שפה: מסתכל רק אחורה
- מודל חיזוי או מקודד (למידה מעוברת autoregressive)

מודל שפה דו כיווני **ELMO** - חיבור של מודל קדימה ומודל אחורה.

BERT - רובוטריק כמודל שפה דו כיווני. נסתיי לחלק מהמילים את עצמן. אחוז ההסתברות המקובל הוא 15%.

יכול לשמש כמקודד רגורם.



הנדסת פרומפטים

למידה בתוך הפרומפט (ICL) - הפרומפט עצמו מכיל דוגמאות, המודל משיט על-סמך מה שראה בהן.

לדוגמה: "בהינתן טקסט ביקורת של לקוח, אני רוצה שתגיד לי האם הוא נהנה. לדוגמה: ביקורת: [ביקורת1], הלקוח [קולא] נהנה."

ההקשר לבדו אמור לתת למודל את היכולת לייצר הכללה. ICL צריך להיות חסכוני ולכן יעיל, כדי שכל הדוגמאות יכנסו בתוך חלון ההקשר. עשוי להיות איטי, כיוון שצריך לתת מחדש את כל הדוגמאות בתוך ההקשר לכל טקסט שאנחנו רוצים לתייג, בפרסד.

שרשרת מחשבה (chain of thought) -

בקונם לבקש תהנית ישר לשאלה, ננחה את המודל לענות "צעד-צעד", עובד מצין בעייתות חשבון (ומשם התפתח). עובד בעיקר במודלים בקנה מידה גדול (100 מיליארד פרמטרים+).

אתיקה והוגנות

השפעה (influence) - X - התחיל להשתמש הרבה במילה, הרשת הברתית הקרובה ל- X מאמצת אותה.

הומופיליה (homophily) - הרבה משתמשים מכירים מילה מסוימת ומתחברים לקהילה קרובה דרך השימוש המשותף.

דמיון מבני (structural equivalence) - צמתים שה"תפקידים" שלהם בתוך הרשת דומים זה לזה (אבל לא בהכרח קשורים).

הוגנות אלגוריתמית - מניחים שניתן לחלק את האוכלוסייה לפי תכונה A , שאינה רלוונטית למשימה הנלמדת (מונח מקובל הוא קבוצה הוגנת).

שונות בתוצאה - ההתפלגות החיזוי \hat{y} בהינתן A שונה מההתפלגות של המידע האמיתי בהינתן A .

שונות בשגיאה - התפלגות שגיאת החיזוי ϵ (לתוצאה מספרית, לתג סיווג מסוים, וכו') שונה עבור A מאשר עבור לא- A .

הוגנות קבוצתית - שגיאות צריכות להתפלג באופן דומה עבור קבוצות שונות.

הוגנות פרטית - פרטים בעלי תכונות (רלוונטיות) דומות צריכים לקבל יחס דומה ללא תלות בהשתייכותם לקבוצה מסוימת.

שוויון בהדדמיות - עבור קבוצות שונות, יחס "הקבלה" דומה.

שוויון ביוחסים - עבור קבוצות שונות, יחס השיגיאות החיוביות דומה (כלומר, ההסתברויות של מועמד ראוי להתקבל ושל מועמד לא ראוי להתקבל דומים עני-פני

ההתפלגות). **קיום התנאים במקביל** - ניתן להוכיח שלא ניתן לקיים את כל התנאים האלה ביחד. מצד שני, יש גם פתח לפתרון (חלקי) - הוספת הוגנות הפרטית ל-loss נוסף.

חילול טקסט

קלט - רצף מילים או משהו שכולל רצף מילים כמו תמונות או אודיו.

פלט - רצף מילים כלשהו.

קשים במעבר בין שפות -

- סדר המילים שונה
- אי תאימות במשמעות
- ביטוי בהטיה מורפולוגית (זמן וגוף)
- לעומת מילות עזר

אלמנטים חסרים (כינוי גוף)

משולש ווקואו - $\text{text} \rightarrow \text{syntax} \rightarrow \text{semantics} \rightarrow \text{interlingua}$

פלט - השכבה האחרונה במודל. ברובוטריק המימוש יהיה להכפיל-פנימית את ה- h שהגיע מהשכבה האחרונה בכל אחד מהשינויים של המילים באוצר המילים: $h_h^T \cdot u$.

שיטות חילול -

- חילול חסדי: ביצוע argmax בעייתי מכיוון שהוא צפוי, חוזר על עצמו ואפילו דטרמיניסטי.

- דגימה טהורה: דגימה בהתפלגות המוגדרת לפי וקטור האומצות שמקבל על הולג'טים. בעייתי כי קיים זנב ארוך של אוצר-הטקסטים. ההסתברות המצטברת של הזנב גבוהה ואנחנו צפויים להיתקל הרבה מאוד בטקסטים דירים, גם אם כל אחת בפרסד קטנה מאוד.

- דגימת ראש: בהינתן k ידוע מראש, ניקח את k המילים שבראש רשימת ההסתברויות הממוינת ומזרזק את היתר.

ננרמל את ההתפלגות הנונת, כלומר נחלק את ההסתברויות המקוריות של k המילים בסכום ההסתברויות שלהן.

דגימת גרעין: בבחר סף q , ונגריל מאוסף המילים אשר סכום ההסתברויות שלהם עבר את סף q . יש צורך לנרמל לאחר בחירת המילים.

טמפרטורה: מחלקים את ערכי הולג'טים בסקלאר τ (היפר-פרמטר) לפני הכניסה לsoftmax. הנוסחה המעודכנת לשלב בניית ההתפלגות היא:

$$y = \text{softmax}\left(\frac{U}{\tau}\right)$$

ברוב המקרים $\tau \in (0, 1]$ כי תוספת חום מקרבת אותנו להתפלגות אחידה וקויר מקרב להתפלגות דטרמיניסטית.

ציון perplexity - כמה המודל מופתע מכל מילה

$$\text{PPL}(x) = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log \mathbb{P}(x_i | x_{<i})\right)$$

כאשר ערך PPL נע בתחום $[1, \infty)$

