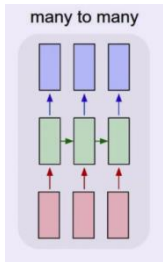


תרגיל 3

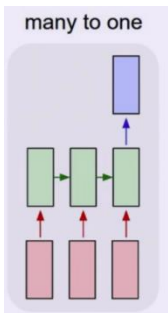
חלק תיאורטי

שאלה 1

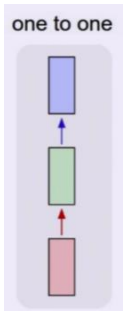


א. כאן קלט מהווה סדרת אותות ופלט זה סדרת מילים. בחירת מילה מתאימה בכל רגע זמן אמורה להתבצע ע"י התחשבות במילים הקודמות שנבחרו, כלומר מה שהוקלט לפני משפיע על הפלט הנוכחי. לכן, הרשת צריכה לעבד סדרת סיגנלים ולפלוט סדרת מילים עם התחשבות בעבר \leftarrow עדיף להשתמש ב- many to many. בנוסף, יש לציין שאורך הטקסט תלוי ישירות באורך ההקלטה \leftarrow נבחר **many to many RNN** מהסוג הבא:

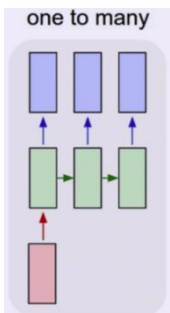
ב. אורך שאלה יכול להשתנות. ניתן לענות על השאלה תוך שימוש במשפט אחד או כמה משפטים. למרות זאת, כדי לענות על שאלה יש להבין קשר בין כל מילה בה. לכן, עיבוד ספרתי לא מתאים כאן, אלא יש לענות ע"י חישוב attention \leftarrow יש להשתמש ב- **many to many transformer**.



ג. קלט: טקסט \leftarrow ממד הקלט משתנה (טקסט יכול להכיל מספר משפטים באורכים שונים). פלט זו מילה שמתארת את רגש הטקסט \leftarrow הרשת עוברת מרב ממד אחד ממד \leftarrow יש להשתמש ב- **many to one RNN**.



ד. קלט תמיד מהווה תמונה אחת ופלט הרשת זה תיוג מחלקה בודדת \leftarrow הבעיה שקולה ל class prediction עם דחיסת תמונה לטנזור \leftarrow יש להשתמש ב- **one to one CNN**.



ה. ממד הקלט הינו קבוע, מכיוון שהרשת מקבלת מילה אחת בלבד. למרות זאת, יכול להיות שלמילה אין תרגום מדויק בשפה אחרת, אלא היא שקולה לאוסף מילים. עקב זאת, ממד הפלט אינו קבוע \leftarrow יש לבחור **one to many RNN**. אם יוצאים מנקודת הנחה שלכל מילה יש מילה שקולה אחת בכל שפה אחרת, אז ניתן להשתמש ב- **one to one FNN**.

שאלה 2

א. אורך טקסט משתנה, גודל תמונה משתנה גם אינו קבוע

← יש לבנות רשת מסוג "many to many".

שלב א': עיבוד קלט.

נשתמש ב-ELMo. ראינו בתרגול ש-ELMo מתחשב בהקשר מילה ובסדר מילים. יהיו בו שלוש שכבות: raw, shallow, deep. בסוף תהליך העיבוד, הטקסט יהווה צירוף לינארי שכבות ה"ל" גודל הקלט אחרי העיבוד יהיה שווה למספר סופי.

שלב ב': מעבר למרחב לטנטי.

ידוע כי קלט מתאר תמונות מסוג מסוים ← ניתן להוציא מהמידע הטקסטואלי נירונים הכי חשובים ורלוונטיים. עקב זאת הרשת תצטרך לדעת להוציא כמות סופית של נירונים שמתארים תמונות ← על מנת לעשות זאת נשתמש ב-mapping.

שלב ג': הכנת מידע לפני שחזור תמונה

סדרת שכבות קונבולוציה ו-upsampling אמורה לקבל מידע בצורת טנזור $\{C, H, W\}$ ← נשתמש בשכבות fully connected.

שלב ד': בניית פלט.

נעביר את המידע דרך סדרת שכבות קונבולוציה ו-upsampling כדי לקבל בסוף תמונה. המעבר הזה שקול ל-"VGG16" בסדר הפוך.

ב. נניח שארבעת הוקטורים יודעים. נקרא להם $X = \{x_1, x_2, x_3, x_4\}$. נניח שממד d.

שלב א': עיבוד קלט.

נשתמש ב-ELMo. יהיו בו שלוש שכבות: raw, shallow, deep. בסוף תהליך העיבוד, הטקסט יהווה צירוף לינארי שכבות ה"ל" ← הקלט יראה כטנזור $\{\text{num of words}, d\}$.

שלב ב': חישוב משקל הattention של כל מילה בטקסט עבור כל וקטור לטנטי. נגדיר שלוש מטריצות W_Q, W_K, W_V , כך ש- W_Q מקבלת וקטורים לטנטיים ו- W_K, W_V מקבלות קידודי טקסט.

$$attention_{i,t} = softmax\left(\frac{Q(x_i)K(w_t)}{\sqrt{d}}\right)$$

שלב ג': חיבוש תוכן לכל רבע תמונה:

$$Context_{x_i} = \sum_{t=1}^{Num\ of\ words} attention_{i,t} \cdot W_V(w_t)$$

שאלה 3

א. נשתמש בנוסחה הבאה:

$$\dim_{out} = \text{floor}\left\{\frac{\dim_{in} + 2 \cdot padding - kernel}{stride}\right\} + 1$$

נקבל:

$$\begin{aligned} & \begin{bmatrix} 128 \\ 128 \end{bmatrix} \overrightarrow{p=0, k=3, s=2} \\ & \begin{bmatrix} 63 \\ 63 \end{bmatrix} \overrightarrow{p=2, k=5, s=1} \\ & \begin{bmatrix} 63 \\ 63 \end{bmatrix} \overrightarrow{p=1, k=3, s=1} \\ & \begin{bmatrix} 63 \\ 63 \end{bmatrix} \overrightarrow{p=0, k=5, s=2} \begin{bmatrix} 30 \\ 30 \end{bmatrix} \end{aligned}$$

לקן פלט נראה כטנזור 30 על 30 על 1.

ב. receptive field – כמות נירונים שרואה כל נירון בשכבת הקונבולוציה האחרונה. נשתמש בנוסחה:

$$RF_l = RF_{l-1} + (k_l - 1) \prod_{i=1}^{l-1} s_i$$

נקבל:

$$\begin{aligned} RF_0 &= 1 \\ RF_1 &= 1 + (3 - 1) \cdot 1 = 3 \\ RF_2 &= 3 + (5 - 1) \cdot (1 \cdot 2) = 11 \\ RF_3 &= 11 + (3 - 1) \cdot (1 \cdot 1 \cdot 2) = 15 \\ RF_3 &= 15 + (5 - 1) \cdot (1 \cdot 1 \cdot 2 \cdot 1) = 23 \end{aligned}$$

כלומר: receptive field = 23 x 23.

שאלה 4

יהי אוסף וקטורים $X = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$ שהוא מטריצת פיצר'רים.

יהיו שלוש מטריצות טרנספורמר: W_Q, W_K, W_V .

נניח שהחלפנו סדר וקטורים ע"י הכפלת X במטריצה אחרת בשם P , כאשר $P \in \{0,1\}^{n \times n}$

טענה: $V^{out}(XP) = V^{out}(X)$

הוכחה:

$$Q = XW_Q, K = XW_K, V = XW_V$$

בעבור XP :

$$Q_{new} = (XP)W_Q, K_{new} = (XP)W_K, V_{new} = (XP)W_V$$

נציב בנוסלה לחישוב attention:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$$

$$\begin{aligned} Attention(XP) &= Attention(Q_{new}, K_{new}, V_{new}) = softmax\left(\frac{(XP)W_Q[(XP)W_K]^T}{\sqrt{d}}\right)XW_VP = \\ &= softmax\left(\frac{XW_QP[XW_KP]^T}{\sqrt{d}}\right)XW_VP = softmax\left(\frac{XW_QPP^TW_K^TX^T}{\sqrt{d}}\right)XW_VP \end{aligned}$$

נשים לב שהמטריצה $P \in \{0,1\}^{n \times n}$ רק משנה סדר באוסף $X \leftarrow PP^T = I$

לכן נקבל:

$$\begin{aligned} Attention(XP) &= softmax\left(\frac{XW_QPP^TW_K^TX^T}{\sqrt{d}}\right)XW_VP = softmax\left(\frac{XW_QIW_K^TX^T}{\sqrt{d}}\right)XW_VP \\ &= softmax\left(\frac{QK^T}{\sqrt{d}}\right)VP = Attention(Q, K, V)P = Attention(X)P \end{aligned}$$

$$V^{out}(X) = \sum_i Attention(x_i)v_i = Attention(X)V$$

$$V^{out}(XP) = Attention(XP)V = Attention(X)PV =$$

$$\sum_i Attention(x_i)v_i = \sum_i Attention(x_i)v_i = V^{out}(X)$$

same vectors, but another order

כעת נסתכל על הבעיה כאשר יש חשיבות לסדר. נעשה זאת ע"י הוספת וקטורי מיקום:

$$E = e_1, e_2, \dots, e_n$$

טענה: $V^{out}(XP + E) \neq V^{out}(XP + EP)$

הוכחה:

$$Q_{new} = (XP + E)W_Q, K_{new} = (XP + E)W_K, V_{new} = (XP + E)W_V$$

$$\begin{aligned}
Attention(XP + E) &= softmax\left(\frac{(XW_Q P + EW_k)(XW_V P + EW_V)^T}{\sqrt{d}}\right)(XW_V P + EW_V) \\
&= softmax\left(\frac{(XW_Q P + EW_k)(P^T W_V^T P^T + W_V^T E^T)}{\sqrt{d}}\right)(XW_V P + EW_V)
\end{aligned}$$

$$\begin{aligned}
Attention((X + E)P) &= Attention(XP + EP) \\
&= softmax\left(\frac{((XP + EP)W_K)((XP + EP)W_K)^T}{\sqrt{d}}\right)(XP + EP)W_V \\
&= softmax\left(\frac{(XW_K P + EW_K P)(XW_K P + EW_K P)^T}{\sqrt{d}}\right)(XW_V P + EW_V P) \\
&= softmax\left(\frac{(XW_K P + EW_K P)(P^T W_K^T X^T + P^T W_K^T E^T)}{\sqrt{d}}\right)(XW_V P + EW_V P)
\end{aligned}$$

$$\begin{aligned}
&Attention(XP + E) \neq Attention((X + E)P) \\
\rightarrow \sum_i Attention(x_i p_i + e_i) v_i &\neq \sum_i Attention((x_i + e_i) p_i) v_i \\
&\rightarrow V^{out}(XP + E) \neq V^{out}(XP + EP)
\end{aligned}$$