

الف) درس NNLM برداشتی کلمات بعینت nonlinear باید که جمع شوند یعنی ~~تفاوت~~ بین اطلاعاتی که باقیمانده کلمات و تفسیرهای معنایی برای دارند تفاوت در نظر گرفته می شود در حالی که در روش

Continuous bag of words صرفاً جمع افعال می هم و اصلاً به سوالی توجه نداریم.

بافت ۲) درس ~~جمله~~ Continuous bag of words پیچیده ای از کلمات را استخراج می کنیم و سعی می کنیم هم به هم وصل کنیم

به راست توجه کنیم در حالی که در NNLM به کلمات بعدی توجه نمی داریم و صرفاً به کلمات قبل توجه داریم.

مثلاً من در بسیاری از موارد غذای خودم را از یک صبح می خورم در اینجا کلمه غذا هیچ خاصی دارد context ندارد یا من بسیاری از موارد غذای خودم را از یک صبح می خورم زیرا صرفاً به قبل توجه داریم.

$$\begin{aligned} \text{merge of word vectors} &\rightarrow N \times D \\ \text{calculating } h &\rightarrow N \times D \times h \\ \text{calculating } \hat{g} \text{ from } h &\rightarrow H \times V \end{aligned} \Rightarrow \begin{aligned} O(NDH + HV) &\Rightarrow ND \\ O(HV) \end{aligned}$$

hierarchical softmax استفاده کنیم یا <sup>۲</sup>نوار Negative sampling

۱) توانیم به جای softmax از

استفاده کرد

۲) ~~استفاده کرد~~

$$\begin{aligned}
 z_1 &= xw + b \\
 z_2 &= hv + d \Rightarrow \delta_1 = \hat{y} - y \\
 \delta_2 &= \delta \frac{\partial z_2}{\partial h} = v^T \delta_1 \\
 \delta_3 &= \delta_2 \frac{\partial h}{\partial z_1} = \delta_2 (1 - \tanh^2(z_1)) \\
 \Rightarrow \frac{\partial C_E}{\partial x} &= \delta_3 \frac{\partial z_1}{\partial x} = w^T \delta_3
 \end{aligned}$$

(1) >

الف (2) جانب ادبی word2vec : سرعت کند و سنجی محاسباتی دقتی روش CCBM از لحاظ هزینه محاسباتی برنگرد  
بهینه تر عمل می کند

نمای word2vec : می تواند آلاش پیچیده را استخراج کند و به سادگی تبدیل Scalable هست.

ب خیر این مسئله دقتی به سادگی شخص است لزوماً وقتی ضریب دوبردار یکی شود معنای یک ل بودن  
بردارها نیست می توان به سادگی بردارها را اسکالر کرد  $\frac{1}{m}$  ضرب کرد یا می توان آن را به اندازه  $\alpha$  درج  
rotate کرد

ج چون هزینه محاسباتی زیادی دارد (اگر سینه  $V$  یا همان vocabulary بزرگ باشد)  $\Rightarrow$  راهکار

Negative Sampling (1)

Hierarchical Softmax (F)

Glove >> word2vec

از لحاظ حافظه Glove خیلی بهتر است و در عمل آن سبکتر است

Glove روی CCBM (concurrency count-based methods) ها آموزش داده می شود.

اولاً هیچ کدام توانایی تشخیص کلمات ستفاد ندارند البته علت آن این است که کلمات ستفاد در فضای vector

بسیار به یکدیگر نزدیک هستند چون اغلب در context های مشابهی می شوند.

ناتوانی: در یادگیری کلمات خارج زبان training بسیار ضعیف شده و قطارهای کلماتی که در زبان train بوده اند خوب عمل می کنند.

محدب نیست زیرا اولاً اگر بصورت یک ن برداریم بردارهای  $Permute$  کنیم تابع هزینه یک ن خواهد بود

مثلاً فرض کنید از embedding بردارهای این بابت می گوید تمام بردارهای تمام ابعاد یک ن باشد و هزینه اقتضای می باشد = محدب نیست (زیرا اگر محدب می بود کامل داشتیم)

ناتوانی: فرض کنید  $R$  و  $\tilde{R}$  را بیان کنیم باز هم هزینه اقتضای می باشد نه معلوم می شود محدب نیست بصورت کلی اگر  $R$  و  $\tilde{R}$  مستقیماً تعیین داده شوند هزینه کامل می باشد.

۳ الف) در این روش، برداری به سمت آموختن embedding از بردار که معنی ~~position~~ position  
آن تخصیص می‌دهیم و در نتیجه position به فضای  $\mathbb{R}^d$  منتقل می‌گردد که در position های مختلف  
embedding های مختلفی خواهد داشت.

ب) در انتقال دفعه ها (Transformer) در واقع embedding قبل از وارد شدن به encoder  
با ماتریس position ترکیب شده که باعث می‌شود مکان خط شود.

ج) در عمل Encoder تمام معادله هر جمله input (یا همان پیرلر hidden) باید در یک بردار ذخیره  
شود بنابراین هر جمله اندازه جمله در ورودی بزرگتر شود و تعدادی اطلاعات سخت‌تری شود و نهایتاً در یک threshold اطلاعات  
هر رنجی شود که به آن فراموشی گویند ولی در کانتینر توجه قسمت های مهم جمله در معنی به یکدسته شدن می‌شود  
و توجه بیشتری به آن می‌شود (البته به جای vector نیز از matrix استفاده می‌کنند)