

سبحان رنجبر

من در کل موارد تنها از مقالات که در زمینه word2vec می باشد استفاده کردم.

سوال 1

1.1

fox, dog => \b(dog|fox)\b

\b برای این که کلمات مستقل پیدا کند و | معنی همان <<یا>> را می دهد.

1.2

\b[d|f]\w*[g|x]\b

d یا f به معنی

\w* هر تعداد کلمه می تواند این وسط باشد.

1.3

\b\w*[aeiou]{2,}\w*\b

\w*\b, \b\w*

فقط اگر کلمه باشد بررسی می کند.

[aeiou]{2,}

حداقل دو حرف صدا دار در یک کلمه می گردد

1.4

\b(\w)\w*\1\b

(\w):

اولین حرف از کلمه را در گروهی ذخیره می کند

\w*:

شامل حروف داخلی کلمه می شود

\1:

تطبیق مجدد با همان حرف اول

1.5

`\b[IL]\w*\b(?:[!]=?)`

`\b[IL]:`

کلماتی که با "I" یا "L" شروع می‌شوند.

`\w*:`

ادامه‌ی کلمه را می‌گیرد

`\b(?:[!]=?)`

تضمین می‌کند که این کلمه در انتهای جمله باشد

سوال (2)

در دنیای واقعی ما چون اشکال مختلف فعل و اسم و قید و ... داریم برای محاسبه هرکدام و ساخت بردارهای آن نیاز به فضای زیادی داریم و از لحاظ معنایی یک فعل و شکل‌های مختلف آن فرقی ندارد پس می‌ایم از دو روش Stemming و Lemmatization استفاده می‌کنیم.

Lemmatization:

فرایند دقیق‌تر است که کلمات را به ریشه لغوی (lemma) بر می‌گردانند. با استفاده از روش فرهنگ لغت و قواعد دستور زبان استفاده می‌کند تا شکل اصلی کلمه را پیدا کنند.

Stemming:

پسوندها و پیشوندها یک کلمه را حذف می‌کنیم تا به stemming برسیم. از قواعد ساده استفاده می‌کنند و تبدیل به ریشه خام می‌کند که شاید کلمه معنی دار نباشد.

Ex)

Lemmatization => خوردم -> خوردن

Stemming => خوردم -> خورد

سوال 3)

king $\rightarrow [0.1, -0.9, 0.2]$

queen $\rightarrow [0.7, 0.5, 0.3]$

apple $\rightarrow [0.8, 0.6, -0.4]$

apple $\rightarrow [0.8, 0.6, -0.4]$ این درست می باشد چون کلمه King , Queen هردو مربوط به سلطنت می باشند و انسان می باشند و به هم نزدیک تر می باشند و apple از این دو دورتر می باشد.

سوال 4

روش Bigram یعنی اتصال آتون کلمات بیت سرهم. (به صورت دوتایی)
 حد باریک تعداد کلمات به صورت تنهایی بدست می آوریم. $\leftarrow V$
 اتصال bigram با Add-one smoothing

$$P(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i) + 1}{C(w_{i-1}) + V}$$

$C \rightarrow$ تعداد حرف و کلمه

$V \rightarrow$ شمارش کلمات به صورت
 منحصر به فرد

$$P(<S> \text{the lazy fox jumps } <S>) =$$

$$= P(\text{the} | <S>) \times P(\text{lazy} | \text{the}) \times P(\text{fox} | \text{lazy}) \times$$

$$\times P(\text{fox} | \text{lazy}) \times P(\text{jumps} | \text{fox}) \times P(<S> | \text{jumps})$$

★ در این جا گی مورد معاینه کنیم و بقیه به همین گونه می باشد.

$$P(\text{the} | <S>) = \frac{C(<S> \text{the}) + 1}{C(<S>) + V} = \frac{1 + 3}{10 + 5} = \frac{4}{15}$$

$\leftarrow V$ $\downarrow C(<S>)$

$$\text{answer} = \frac{1}{13} \times \frac{4}{13} \times \frac{1}{13} \times \frac{1}{4} \times \frac{4}{15} = \frac{16}{131250}$$

سوال 5)

در این سوال قرار است از صفر word2vec را پیاده سازی کنیم در این جا چند مورد را باید بلد باشیم مثل این که چگونه از صفر مدل یادگیری عمیق را بنویسیم و قوانین پیاده سازی ان مثل forwardpass چگونه به جلو ببریم در مدل و یا backprop که چگونه مدل را اپدیت کنیم و کلیت کار ما به این شکل است داده ها را بگیریم چون کلمات ما از نوع رشته ای است باید تبدیل به بردار عددی کنیم که راحت ترین ان تبدیل به وان هات می باشد و بعد از تبدیل به بردار عدد ان را به مدل یادگیری عمیق بدهیم و روی ان آموزش ببیند (train)

در اخر هم می توانیم از activtion function مختلف استفاده کنیم که به صورت پیشفرض در word2vec از softmax استفاده می شود که داخل پرانتز در مقالات بعدی سعی شد از این تابع استفاده نشود چون درمخرج کسر باید همه موارد هر دفعه بررسی می کرد که هزینه محاسبات داشت.

و در اخر برای پیدا کردن شباهت کلمات می توانیم از cosinesimilarity و یا دوستاش استفاده کنیم و ان هایی که بردار نزدیک تری دارند ان را انتخاب کنیم.

سوال 6)

در این سوال با کتابخونه gensim قرار بود هر دو مدل word2vec را بررسی کنیم و با هم مقایسه کنیم , skip-gram , Cbow پالشی که من تو این سوال داشتم پیدا کردن درست استفاده کردن از لایبرری بود داک خیلی قوی نداشت و دیگری من سعی کردم کلمه به کلمه بخونم که رم 16 توانایش نداشت و از تی ای هم پرسیدم و فهمیدم باید جمله به جمله بخونم. براساس مقاله خود word2vec و نتایج که رسیدیم skip-gram راه حل بهتری می باشد.