

باسمه تعالی



## یادگیری ژرف

تمرین سری چهارم  
RNN, LSTM, Word2Vec

نام استاد:  
دکتر سلیمانی

نام دانشجو:  
مهدی یار شهبازی  
۹۵۱۰۶۳۹۷

تاریخ تحویل:  
جمعه ۱۳۹۸/۲/۲۰

## LSTM

### پردازش دیتا

ما ابتدا مصراع‌های شعر فردوسی را از هم جدا کردیم و بعد از تعداد کارکترهای تمام مصراع‌ها را شمردیم و بیشترین تعداد کارکتر را به عنوان طول جمله در نظر گرفتیم و اگر جمله‌ای کمتر از این طول را داشت به آن صفر اضافه کردیم. در نهایت یک دیکشنری تعریف کردیم که متاظر هر کارکتر یک عدد به ما نسبت بدهد. حال ماتریس داده‌های ما آماده است و تنها کاری که لازم است بخشی از آن را به عنوان Train و بخشی دیگر را به عنوان Test در نظر بگیریم.

### مدل

در این تمرین از کتابخانه‌ی TensorFlow استفاده کردیم. برای ساختن مدل از کدهای داخل اینترنت کمک گرفته شد. در ساخت مدل با شهود اینکه برخی از حروف امکان زیادی دارد کنار هم بیایند و برخی دیگر احتمال کمی دارند، از embedding استفاده کردیم و ۴۰ کارکتر موجود در دیکشنری را به ۲۵ کارکتر بردیم. دقت داریم که مدل دو بخش برای Train و Test دارد. ما تلاش کردیم به نوعی این دو قسمت را یکی بکنیم، اما با مشکلات زیادی روبه رو شدیم. برای بهینه کردن از الگوریتم Adam استفاده کردیم. همچنین در هنگام آموزش به مرور learning rate را کاهش دادیم.

### خروجی‌ها

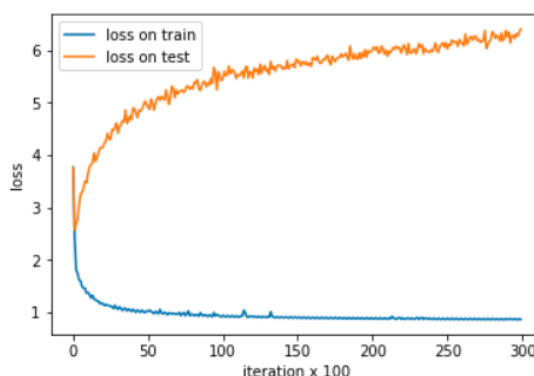
در اینجا تعدادی از جمله‌هایی را که مدل ما ساخته است مشاهده می‌کنیم.

۱. به دیدار او بر سر از درد و جفت

۲. سر تخت شاهی به ایران به دشت

باقی مصراع‌ها در نوت‌بوک مربوط به این قسمت موجود است.

این نمودار با شهود ما هم‌خوانی دارد چرا که هر چقدر مدل داده‌های آموزش را بهتر یاد بگیرد، نسبت به داده‌هایی که آن‌ها را ندیده است بیگانه تر



میشود.

## Attention

### توضیح مختصر

با استفاده از مکانیزم توجه، به جای این که ابتدا کل اطلاعات حالت های نهان در تمام گام های زمانی را فقط در یک بردار ثابت کد کنیم و سپس با استفاده از دیک کردن آن بردار، خروجی ها تولید شود، می توان برای تولید هر خروجی، به حالات نهان همه گام های زمانی نگاه کرد و با ترکیب وزن دار آن ها در هر لحظه، اطلاعاتی را که بیشتر برای تولید خروجی آن لحظه مورد نیاز است انتخاب کرد. همانطور که ذکر شد، با ترکیب وزن دار حالات نهان لحظه های مختلف، هر لحظه میتوان به ورودی هایی که برای تولید خروجی آن لحظه اهمیت بیشتری دارند، وزن بیشتری نسبت به دیگر ورودی ها داد تا تاثیر بیشتری داشته باشند بنابراین مشکل اول که کد کردن کل اطلاعات در یک بردار بود، با تولید بردارهای متفاوت برای هر خروجی با ترکیب های متفاوت حالات نهان حل میشود. همچنین مشکل دوم، با وزن دهی بیشتر به ورودی های مهمتر قابل حل است.

## پردازش دیتا

دقیقا مانند قبل.

## مدل

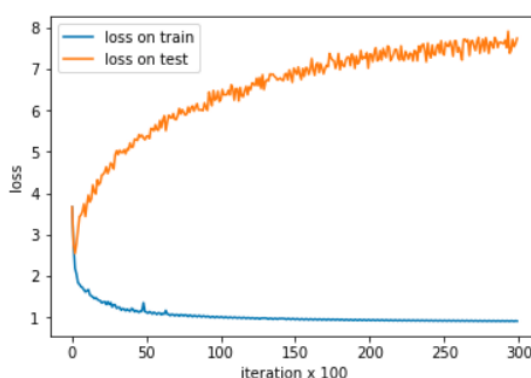
دقیقا مانند قبل با این تفاوت که یه ماژول مربوط به Attention اضافه میشود.

## خروجی ها

در اینجا تعدادی از جمله‌هایی را که مدل ما ساخته است مشاهده میکنیم.

۱. به پیش سپه را به دل بر نه اوی

۲. به ایران به دیدار بر سر بود



اگر زیادت‌تر شدن تابع هزینه را روی نمونه‌های تست را مشاهده‌ای از بیشتر یادگرفتن داده‌های آموزش بگیریم، آنگاه گویی مدل با مکانیزم توجه بهتر یاد می‌گیرد و همچنین با توجه به بیت‌هایی که به وسیله‌ی مدل مجهز به این مکانیزم تولید شده‌است بنظم معنای بیشتری دارند. مصراع‌های داخل نوبت‌بک برای دو روش به ازای ورودی‌های یکسان هستند، پس به نوعی خواسته قسمت امتیازی را ارضا میکنند.

## Word2Vec

با توجه به حجم محاسبات بالای مدل skip gram عادی، از مدل negative sampling استفاده کردیم. توضیحات پیاده‌سازی و دیتا نیز قسمت‌های بعد ارایه میشود.

## پردازش دیتا

در این قسمت با یک الگوریتم نسبتاً ناهینه دیتاها را به فرم ماتریس مجاورت در می‌آوریم. در ادامه می‌خواهیم برای هر لغت خاص، یک همسایه و تعدادی غیر همسایه انتخاب کنیم و مدل را بر این اساس آموزش دهیم.

## مدل

اگر می‌خواهیم فاز آموزش را بر این مبنا قرار دهیم که برای هر ورودی را براساس بردار One Hot به شدت حجم بالایی برای محاسبات نیاز داریم پس ایده‌ای که استفاده میکنیم، که البته از اینترنت گرفته شده است، تنها index هر ورودی را می‌گیریم و سپس تابع هزینه را حساب میکنیمو اینکار مدل ما را بسیار بهینه میکند و هر چقدر که بخواهیم میتوانیم تابع هزینه را کاهش دهیم.

## خروجی‌ها

در صورت تمرین از ما خواسته شده بود تا ۵ همسایه نزدیک ایران، رستم، خردمند، گلاب و سیستان را پیدا کنیم. خروجی های ما به ترتیب به قرار زیر است: و تابع هزینه:

```
[ 'ایران' 'دوربیند' 'نوشتست' 'سربگاشت' 'سیاهست' 'برداشت' ]  
[ 'رستم' 'وکارکرد' 'سالارنو' 'نجستم' 'کردگارمکان' 'پرستیز' ]  
[ 'خردمند' 'فروزان' 'سرسوی' 'شتابنده' 'سپردندان' 'مشکن' ]  
[ 'گلاب' 'نیکوگمان' 'خونبان' 'بخشیدی' 'اژرتر' 'بکوبند' ]  
[ 'سیستان' 'هممی' 'افسرت' 'خوانهای' 'کاسپان' 'سهد' ]
```

