



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

## یادگیری ژرف

نیم سال اول ۰۲-۰۱

استاد: دکتر حمید بیگی

گردآورندگان: علی سطوتی - محمدعلی صدرایی جواهری - امیرحسین عاملی - محدثه میربیگی

بررسی و بازبینی: مهران حسینزاده

تمرین چهارم

یادگیری بازنمایی و مکانیزم توجه

مهلت ارسال: ۹ دی ۱۴۰۱

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر تمرین ها بدون کسر نمره تا سقف ۱۰ روز (تا سقف ۳ روز برای هر تمرین) وجود دارد. محل بارگزاری جواب تمرین ها بعد از ۵ روز بسته خواهد شد و پس از گذشت این مدت، پاسخ های ارسال شده پذیرفته نخواهند شد. همچنین، به ازای هر روز تأخیر غیر مجاز ۱۰ درصد از نمره تمرین به صورت ساعتی کسر خواهد شد.
- هم کاری و هم فکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت هم فکری و یا استفاده از هر منابع خارج درسی، نام هم فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.
- پاسخ تمامی سوالات (تئوری و عملی) را در یک فایل فشرده به صورت `DL_Hw4_[firstName]_[lastName]_[StudentId]` نام گذاری کرده و ارسال کنید.

سوالات نظری (۶۰ نمره)

### ۱. Autoencoder - 1 (۱۰ نمره)

فرض کنید مدل های خطی با دو وزن ماتریس داریم: یک رمزگذار  $W_1 \in \mathbb{R}^{K \times m}$  و یک رمزگشا  $W_2 \in \mathbb{R}^{m \times k}$ . مدل سنتی autoencoder بازنمایی با ابعاد کم  $n$  نقطه را با داده های آموزشی  $X \in \mathbb{R}^{m \times n}$  یاد می گیرد. این یادگیری با کمینه سازی تابع زیر صورت می گیرد.

$$L(W_1, W_2; X) = \frac{1}{n} \|X - W_2 W_1 X\|_F^2$$

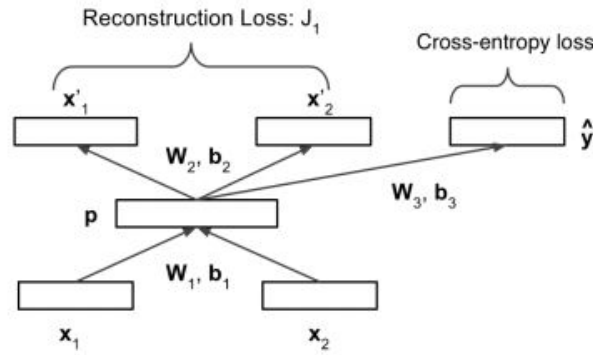
فرض کنید که برای ماتریس  $1/n X X^T$  مقادیر ویژه  $\sigma_1 > \dots > \sigma_k > \sigma_{k+1} = 0$  وجود دارند. حال L2-regularized linear autoencoder را با تابع هدف زیر در نظر بگیرید:

$$L(W_1, W_2; X) = \frac{1}{n} \|X - W_2 W_1 X\|_F^2 + \lambda \|W_1\|_F^2 + \lambda \|W_2\|_F^2$$

برای حل این رابطه گفته شده اگر لاابدا مقدار غیر صفر کوچک باشد، آیا ما یک inductive bias برای پیدا کردن ماتریس  $W_2$  داریم؟ بحث نمایید.

### ۲. Autoencoder - 2 (+Word Embedding) (۱۴ نمره)

یکی از راه های ترکیب بردارها جمعشان با هم است. در این سوال می خواهیم راه دیگری با استفاده از autoencoder امتحان کنیم. شکل زیر نشان دهنده این راه حل است.



در این autoencoder دو بردار ورودی  $x_1$  و  $x_2 \in \mathbb{R}^{D_x \times 1}$  به هم متصل می‌شوند و یک بردار  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  می‌سازند. بردار والد و ماتریس  $W$  به صورت زیر محاسبه می‌شوند.

$$p = \text{ReLU}(W_1 x + b_1) \in \mathbb{R}^{D_p \times 1}$$

$$W_1 = \begin{bmatrix} W_{11} & W_{12} \end{bmatrix}$$

در حین آموزش از بردار والد برای بازسازی بردار ورودی استفاده می‌شود.

$$x' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} = W_2 p + b_2 \in \mathbb{R}^{2D_x \times 1}$$

بردارهای  $x'_1$ ،  $x'_2$  بردارهای بازسازی شده هستند. تابع هزینه بازسازی برابر با فاصله اقلیدسی بین ورودی‌ها و بازسازی‌ها در حین آموزش است.

$$J_1 = \frac{1}{2} \|x' - x\|^2 \in \mathbb{R}$$

برای تحلیل احساسات از بردار والد برای پیش‌بینی استفاده می‌شود تا کلاس مورد نظر  $\hat{y}$  پیش‌بینی شود.

$$\hat{y} = W_3 p + b_3 \in \mathbb{R}^{D_c \times 1}$$

برای تحلیل احساسات ما سه کلاس داریم. شبکه با تابع هزینه cross entropy آموزش می‌بیند.

$$J_2 = CE(y, \hat{y} \in \mathbb{R})$$

$Y$  یک بردار برچسب one-hot است. شبکه با جمع دو تابع هزینه گفته شده آموزش می‌بیند:  $J = J_1 + J_2$

(الف) چگونه بردارهای بازسازی شده و تابع هزینه بازسازی به یادگیری بردار والد کمک می‌کنند؟

(ب) ابعاد هر یک از وزن‌ها و بایاس‌های شبکه را بیابید؟ تعداد کل پارامترهای مدل چند تا است؟

(ج) چگونه می‌توان شبکه‌ای فقط با کپی از این ماژول autoencoder داشت که بتواند برچسب برای تحلیل احساسات برای یک جمله کامل را انجام دهد؟

### ۳. مکانیزم توجه (۱۸ نمره)

مکانیزم توجه برای از بین بردن گلوگاه اطلاعات بین رمزگذار و رمزگشا معرفی شده است. به این صورت که به جای آخرین بردار نهان رمزگذار، رمزگشا به تمام بردارهای نهان رمزگذار دسترسی دارد. این مکانیزم به صورت زیر فرموله می شود و در هر گام شبکه ی تکرارشونده ی رمزگشا مورد استفاده قرار می گیرد:

$$a_t(s) = \frac{\exp \text{score} \left( h_d^{(t)}, h_e^{(s)} \right)}{\sum_{s'} \exp \text{score} \left( h_d^{(t)}, h_e^{(s')} \right)}$$

$$c_t = \sum_{s'} a_t(s') h_e^{(s')}$$

$$\hat{h} = \tanh W_c \left[ c_t; h_d^{(t)} \right]$$

$$y_t = \text{softmax} \left( W_s \hat{h} \right)$$

که در آن  $h_d^{(i)}$  بردار نهان رمزگشا،  $h_e^{(i)}$  بردار نهان رمزگذار و  $y_t$  خروجی گام  $t$  ام رمزگشا می باشد. تابع  $\text{score} \left( h_d^{(t)}, h_e^{(s)} \right)$  را می توان به سه روش زیر تعریف کرد:

$$\text{score} \left( h_d^{(t)}, h_e^{(s)} \right) = \begin{cases} h_d^{(t)T} h_e^{(s)} & \text{dot} \\ h_d^{(t)T} W_a h_e^{(s)} & \text{general} \\ v_a^T \tanh W_a \left[ h_d^{(t)}; h_e^{(s)} \right] & \text{tanh layer} \end{cases}$$

(آ) این سه تابع را از نظر توان مدل کردن، هزینه ی محاسباتی و عبور گرادیان در مرحله بازانتشار خطا مقایسه کنید. شما کدام یک را برای یک شبکه Seq2Seq انتخاب می کنید؟

(ب) در ادبیات یادگیری عمیق، دو کار پژوهشی دو مکانیزم توجه ارائه داده اند که جز رایج ترین کارهای این حوزه می باشد: ۱ - **مکانیزم ۱** - ۲ - **مکانیزم ۲**. این دو ساختار را با هم مقایسه کنید و تفاوت های آن را ذکر کنید. کدام یک توانایی مدل کردن بیشتری دارد؟

(ج) یکی از مشکلات رایج مکانیزم توجه، مخصوصا هنگامی که متن ورودی در طرف رمزگذار طولانی باشد، عدم توانایی این مکانیزم در پرداختن به تکه های مختلف متن ورودی است. به طور مثال ممکن است در تمامی گام های رمزگشا، مکانیزم توجه فقط به یک یا دو کلمه ی خاص امتیاز بسیار بالایی بدهد و فقط آن ها را در نظر بگیرد. در این صورت مدل قادر نخواهد بود که از تمامی متن ورودی استفاده کند. برای حل این مشکل چه راهکاری پیشنهاد می دهید؟ توضیح دهید.

#### ۴. Transformer (۱۸ نمره)

یکی از مشکلاتی که transformerها دارند این است که مرتبه هزینه محاسباتی و هزینه ذخیره سازی عملیات self-attention دارای عبارت  $N^2$  می باشد. این مرتبه باعث می شود که آموزش این شبکه روی داده های طولانی مانند کتاب مشکلزا باشد. دلیل این امر عملگر Softmax می باشد که برای محاسبه شباهت دو بردار استفاده می شود. در این تمرین قصد داریم به بررسی یک راهکار جایگزین برای این مورد پردازیم. یکی از این راهکارها استفاده از مکانیزم های توجهی کرنلی می باشد.

اگر ورودی را  $x \in \mathbb{R}^{N \times F}$  و ماتریس های مکانیزم توجه را  $W_Q \in \mathbb{R}^{F \times D}$ ,  $W_K \in \mathbb{R}^{F \times D}$  و  $W_V \in \mathbb{R}^{F \times M}$  در نظر بگیریم. می توان این عملیات را به صورت زیر نوشت:

$$Q = xW_Q, K = xW_K, V = xW_V$$

$$V' = \text{softmax} \left( \frac{QK^T}{\sqrt{D}} \right) V$$

حال با تعریف  $\text{sim}(Q_i, K_j) = \exp \left( \frac{Q_i^T K_j}{\sqrt{D}} \right)$  می توان این عبارت را به فرم زیر بازنویسی کرد:

$$V'_i = \frac{\sum_{j=1}^N \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^N \text{sim}(Q_i, K_j)}$$

(آ) مرتبه زمانی و حافظه مورد نیاز برای محاسبه عملگر self-attention بالا را براساس پارامترهای  $N, D, M$  محاسبه کنید.

(ب) یکی از توابعی که می توان جایگزین  $\text{sim}(Q_i, K_j)$  کرد، کرنل توجه چندجمله ای می باشد. عبارت جایگزین را برای حالت درجه دو (Quadratic) بنویسید.

(ج) برای کرنل مرتبه بخش قبل، بردار ویژگی  $\phi(\cdot)$  را بنویسید

(د) حال با توجه به رابطه  $K(q, k) = \phi(q)^T \phi(k)$ ، رابطه نهایی  $V_i^t$  (طبق توضیحات بالا) را بازنویسی کنید و مرتبه زمانی رابطه و مرتبه حافظه مورد نیاز را برای آن محاسبه کنید. با مقایسه این مرتبه ها با مرتبه های رابطه قبلی، در چه شرایطی استفاده از این رابطه بهتر از رابطه قبلی می باشد؟

سوالات عملی (۴۰+۱۰ نمره)

#### ۱. بردارهای معنایی (۲۵ نمره)

در این تمرین عملی قرار است شما با استفاده از بردارهای glove چهار کار جالب در زمینه NLP را انجام دهید. در بخش اول و دوم شما با استفاده از روابط ریاضی برداری چند ویژگی جالب در بردارهای کلمات مشاهده

می‌کنید. در بخش سوم و چهارم قرار است یک طبقه‌بندی را آموزش دهید. داده ورودی شما متون نقد فیلم کاربران است و هر نقد بر اساس محتوایش می‌تواند نقد مثبت یا منفی باشد. در این تمرین شما باید به دو روش مختلف مدلی را برای طبقه‌بندی آموزش دهید و این دو روش را با هم مقایسه کنید. در بخش سوم مدل شما باید با روش میانگیری بین بردارهای کلمات موجود در نقد، بردار معادل هر نقد را پیدا کند و با استفاده از مدل MLP آن را طبقه‌بندی کند. بخش زیادی از کد بخش سوم زده شده است و فقط یکی از تابع‌های آن نیاز به تکمیل دارد. در بخش چهارم شما باید یک مدل LSTM برای این موضوع آموزش دهید. در زمینه تعداد لایه‌ها یا اندازه بردار مدل انتخاب دست خودتان است. همچنین می‌توانید با دیدن کدهای روش اول برای زدن کدهای روش دوم ایده بگیرید. همچنین شما باید مدل بخش چهارم را با دو تنظیم مختلف آموزش دهید. در تنظیم اول لایه embedding فریز می‌شود و در تنظیم دوم خیر و در نهایت باید عملکرد مدل نهایی را در این دو حالت مقایسه کنید و در نوتبوک خود توضیح دهید. برای حل این سوال به نوتبوک ضمیمه شده `explore_word_embeddings.ipynb` مراجعه کنید و همه‌ی پاسخ‌های خواسته شده را داخل همان نوتبوک ارائه دهید.

## ۲. طبقه‌بندی اسپم (۲۵ نمره)

در این تمرین به پیاده‌سازی یک طبقه‌بند هرزنامه خواهیم پرداخت. دادگان موجود مجموعه دادگان پیامک‌های هرزنامه هستند که دارای برچسب مناسب می‌باشند. برای پیاده‌سازی طبقه‌بند از مدل‌های زبانی از پیش آموزش دیده و یک نوع خاص از خانواده ترنسفورمر به نام برت (BERT) استفاده خواهیم کرد. برای حل این سوال به نوتبوک ضمیمه شده `spam_classification.ipynb` مراجعه کنید و همه‌ی پاسخ‌های خواسته شده را داخل همان نوتبوک ارائه دهید.