



- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- همکاری و همفکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت همفکری و یا استفاده از هر منابع خارج درسی، نام همفکران 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.

## سوالات نظری (۵ + ۶۰ نمره)

### مسئله ۱. (۱۰ 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 تکرارشونده رمزگشا مورد استفاده قرار می گیرد:

$$a_t(s) = \frac{\exp \text{score}(h_d^{(t)}, h_e^{(s)})}{\sum_{s'} \exp \text{score}(h_d^{(t)}, h_e^{(s')})} \quad (1)$$

$$c_t = \sum_{s'} a_t(s') h_e^{(s')} \quad (2)$$

$$\hat{h} = \tanh W_c [c_t; h_d^{(t)}] \quad (3)$$

$$y_t = \text{softmax}(W_s \hat{h}) \quad (4)$$

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  $\text{score}(h_d^{(t)}, h_e^{(s)})$  را می توان به سه روش زیر تعریف کرد:

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

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

(ج) یکی از مشکلات رایج مکانیزم توجه، مخصوصاً هنگامی که متن ورودی در طرف رمزگذار طولانی باشد، عدم توانایی این مکانیزم در پرداختن به تکه‌های مختلف متن ورودی است. به طور مثال ممکن است در تمامی گام‌های رمزگشا، مکانیزم توجه فقط به یک یا دو کلمه‌ی خاص امتیاز بسیار بالایی بدهد و فقط آن‌ها 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 می‌دهید؟ توضیح دهید.

یکی از مشکلاتی که transformer ها 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 مشکل زا باشد. دلیل این امر عملکرد 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 بپردازیم. یکی از این راهکارها استفاده از مکانیزم های توجهی کرنلی می باشد.

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

$$V' = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$
$$V_i' = \frac{\sum_{j=1}^N \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^N \text{sim}(Q_i, K_j)} \quad (\Delta)$$

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

موفق باشید :