

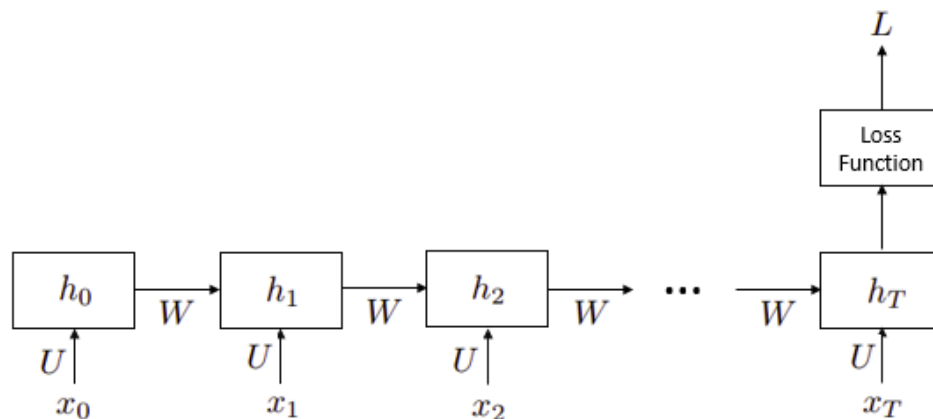


- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- فایل پاسخ را به همراه تمام کدها در یک فایل فشرده و با عنوان $STD - HW3 \#$ در سایت Quera.ir بارگذاری نمایید.
- بخش‌های پیاده‌سازی مربوط به هر سوال را در فایل مربوطه با شماره‌ی آن سوال قرار دهید
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.
- سوالات خود را از طریق [Quera](http://Quera.ir) مطرح کنید.

سوالات نظری (۷۰ نمره)

مسئله‌ی ۱. (۱۵+۵ نمره)

(بخش ۱) با توجه به شبکه عصبی بازگشتی شکل زیر به سوالات پاسخ دهید. دقت کنید که برای سادگی تمام مقادیر یعنی ورودی‌ها و وزن‌ها و خروجی مقادیر اسکالر هستند. همچنین فرض کنید تمام توابع فعالساز σ هستند.

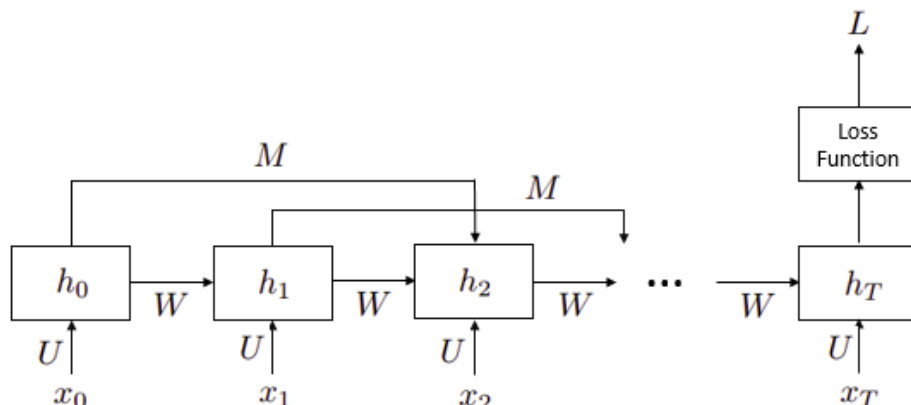


- (آ) ابتدا گرادیان h_t یعنی $\frac{\partial L}{\partial h_t}$ را بر حسب گرادیان h_{t+1} یعنی $\frac{\partial L}{\partial h_{t+1}}$ بنویسید. (۱ ≤ t ≤ T - ۱) (۳ نمره)
- (ب) حال از رابطه قسمت قبل استفاده کرده و به شکل زنجیر وار گرادیان h_0 را بر حسب گرادیان h_T بنویسید. (۲ نمره)

(بخش ۲) حال می‌خواهیم روش‌هایی برای جلوگیری از محوشدگی و انفجار گرادیان را معرفی و تحلیل کنیم.

- (آ) یکی از روش‌های مهم جلوگیری از محوشدگی و انفجار گرادیان مقداردهی اولیه صحیح وزن‌های شبکه است. توضیح دهید حداکثر مقدار اولیه W چند باشد تا فارغ از ورودی مطمئن باشیم که از همان ابتدا انفجار گرادیان رخ ندهد. (راهنمایی: یک حد بالا برای گرادیان h_0 پیدا کنید.) (۵ نمره)

(ب) یکی از راه‌های جلوگیری از محوشدگی گرادیان استفاده از skip-connection ها است. شکل زیر را در نظر بگیرید که در آن هر h_t علاوه بر h_{t+1} به h_{t+2} هم متصل است. حال دوباره گرادیان h_t را برحسب گرادیان h_{t+1} و h_{t+2} نوشته و توضیح دهید چرا اینکار تا حد خوبی باعث کاهش اثر محوشدگی گرادیان می‌شود. (۵ نمره) ($1 \leq t \leq T-2$)



(ج) یکی از راه‌حل‌های جلوگیری از انفجار گرادیان، برش گرادیان^۱ است که این خودبه‌دو زیرراه‌حل برش توسط مقدار^۲ و برش توسط اندازه^۳ تقسیم می‌شود. این دو را جداگانه توضیح دهید. برتری برش توسط اندازه را به برش توسط مقدار را توضیح دهید. (۵ نمره امتیازی)

مسئله‌ی ۲. (۲۵+۱۰ نمره)

در این مسئله می‌خواهیم با مفاهیمی در تولید دنباله در شبکه‌های Seq2Seq و مزایا و معایب آن‌ها آشنا شویم. (بخش ۱) در بخش اول می‌خواهیم مفهوم teacher forcing را بررسی کنیم. برای تولید دنباله ما می‌توانیم یک استراتژی خام اولیه در نظر بگیریم، می‌توان برای تولید نشانه^۴ $t+1$ توسط رمزگشای^۵ زمان $t+1$ ، نشانه تولید شده توسط شبکه در زمان t را به عنوان ورودی به دیکودر زمان $t+1$ بدهیم اما این حالت مشکلاتی دارد.

(آ) ابتدا توضیح دهید این مشکلات چه چیزهایی هستند و سپس روش teacher forcing را توضیح داده و بگویید که teacher forcing چگونه این مشکلات را برطرف می‌کند. (۵ نمره)

(ب) مشکل اصلی teacher forcing موضوعی به نام exposure bias است. این مشکل را توضیح دهید. (۵ نمره)

(ج) یکی از راه‌حل‌های مشکل exposure bias تکنیک scheduled sampling است، این تکنیک را توضیح داده و بگویید این تکنیک چگونه باعث کاهش اثر exposure bias می‌شود. (۵ نمره)

(بخش ۲) حال در بخش دوم مسئله می‌خواهیم بر روی الگوریتم جستجوی موجی^۶ تمرکز کنیم. این الگوریتم در تقابل با الگوریتم حریصانه برای تولید دنباله در زمان رمزگشایی مطرح می‌شود.

(آ) ابتدا تفاوت دو الگوریتم جستجوی موجی و الگوریتم حریصانه برای تولید دنباله را بیان کنید. (۵ نمره)

^۱ gradient clipping
^۲ clipping by value
^۳ clipping by norm
^۴ token
^۵ decoder
^۶ beam search

(ب) در الگوریتم جستجوی موجی ابرپارامتری بنام k وجود دارد که حداکثر تعداد شاخه‌های جستجوی ما در هر زمان را نشان می‌دهد. توضیح دهید که کاهش بیش از حد k باعث چه مشکلاتی می‌شود. همچنین توضیح دهید افزایش بیش از اندازه k چه مشکلاتی بوجود می‌آورد. (۵ نمره امتیازی)

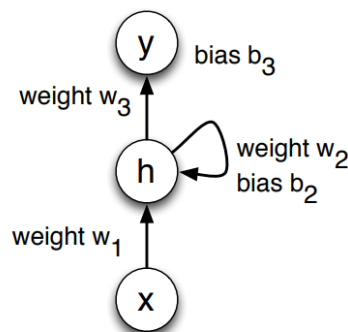
(بخش ۳) حال در بخش سوم مسئله می‌خواهیم به موضوع دیگری برای تولید دنباله پردازیم. در الگوریتم حریصانه همیشه کلمه با بیشترین احتمال در لایه‌ی softmax به عنوان کلمه خروجی انتخاب می‌شد، اما روش دیگری برای این کار وجود دارد و آن انتخاب تصادفی کلمه خروجی براساس احتمال‌های لایه softmax است.

(آ) توضیح دهید که مزایای این حالت به حالت انتخاب کلمه با بیشترین احتمال چیست. (۵ نمره)

(ب) برای این اساس دو روش sampling بنام‌های pure sampling و top-k sampling معرفی می‌شوند تفاوت این دو روش نمونه برداری را توضیح دهید. اثرات و مزایا و معایب زیاد یا کم کردن k در top-k sampling را شرح دهید. (۵ نمره امتیازی)

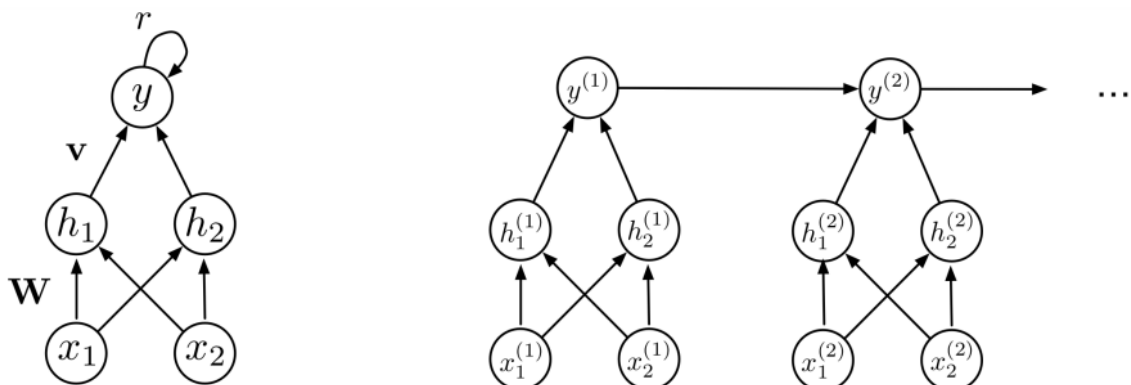
مسئله‌ی ۳. (۱۰ نمره)

یک شبکه بازگشتی به صورت مقابل را در نظر بگیرید. وزن‌ها و بایاس‌ها را به گونه‌ای تعیین کنید که در هر دنباله‌ای از اعداد تا زمانی که ورودی شبکه ۱ باشد، خروجی شبکه یک باقی بماند و به محض اینکه ورودی شبکه به صفر تغییر کند خروجی شبکه صفر شده و صفر باقی بماند. برای مثال خروجی شبکه به ازای ورودی ۱۱۱۰۱۰۱ برابر با ۱۱۱۰۰۰۰ می‌باشد.



مسئله‌ی ۴. (۵ نمره)

یک شبکه بازگشتی بصورت مقابل را در نظر بگیرید. فرض کنید این شبکه دو دنباله از اعداد صفر و یک را دریافت کرده و اگر دو دنباله برابر بودند عدد ۱ و در غیر اینصورت عدد صفر را به عنوان خروجی بر می‌گرداند.



$$\mathbf{h}^{(t)} = \phi(\mathbf{W}\mathbf{x}^{(t)} + \mathbf{b})$$

$$y^{(t)} = \begin{cases} \phi(\mathbf{v}^\top \mathbf{h}^{(t)} + ry^{(t-1)} + c) & \text{for } t > 1 \\ \phi(\mathbf{v}^\top \mathbf{h}^{(t)} + c_0) & \text{for } t = 1, \end{cases} \quad \phi(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

ماتریس \mathbf{W} یک ماتریس 2×2 و b و v بردارهای دو بعدی و c و r و c_0 مقادیر اسکالر می باشد. آن ها را به گونه ای تعیین کنید که شبکه کارکرد تعریف شده را داشته باشد. (راهنمایی: خروجی $y^{(t)}$ در هر لحظه نشان می دهد آیا دو دنباله تا آن لحظه برابر بوده اند یا خیر. لایه مخفی اول نشان میدهد آیا دو ورودی در لحظه t صفر بوده اند یا خیر و لایه مخفی دوم نشان می دهد آیا دو ورودی در لحظه t ، ۱ بوده اند یا خیر.)

سوالات عملی (۵۰ نمره)

مسئله ۵. (۲۰ نمره)

در این سوال می خواهیم با استفاده از شبکه LSTM یک دسته بندی بر روی دیتاست Yelp انجام دهیم. نوت بوک Q۵ را باز کرده و سلول های حاضر را اجرا کرده تا داده ی آموزش و اعتبارسنجی شما آماده شود. توجه داشته باشید که باید به عنوان ورودی کلمات به شبکه از بردارهای از پیش آموزش دیده Glove استفاده کنید. پس از آموزش مقدار امتیاز f_1 داده های اعتبارسنجی را برای هر epoch رسم کنید. در طراحی شبکه و ابرپارامترهای آن آزاد هستید.

مسئله ۶. (۲۵+۵ نمره)

در این تمرین هدف پیاده سازی دو شبکه LSTM و GRU و پیش بینی بازار سهام بوسیله آن ها می باشد. به نوت بوک Q۶ مراجعه شود.