

سوال اول) ثابت کنید فرمول تابع هزینه ی naive-softmax همانند تابع هزینه ی cross-entropy بین y و \hat{y} است.

ابتدا باید متغیر های نظیر هم در دو رابطه را بشناسیم:

y بردار خروجی ای است که باید به آن برسیم و در آن درایه ی مربوط به کلمه ی مرکزی ۱ و برای بقیه ی کلمات ۰ هستند. در واقع نمایش $one-hot$ برای کلمه ی مرکزی است.

\hat{y} بردار خروجی مدل است که هر درایه ی آن شامل احتمال $context$ بودن یکی از کلمات واژه های موجود را برای کلمه ی مرکزی است.

\hat{y}_o درایه ی مربوط به کلمه ی مرکزی در بردار \hat{y} است که برابر با احتمال $P(o|c)$ است.

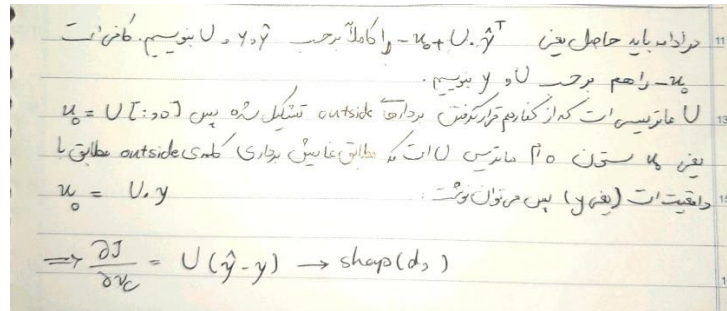
این را هم میدانیم که $\log(1) = 0$.

از فرمول $cross-entropy$ شروع می کنیم و سپس عبارت مربوط به کلمه ی عبارت مرکزی را از بقیه ی عبارت ها که صفر میشوند جدا میکنیم و با جایگذاری مقادیر معادل y_o به فرمول naive-softmax میرسیم:

$$\begin{aligned} J_{cross-entropy} &= - \sum_{w \in Vocabulary} y_w \log(\hat{y}_w) = -y_o \log(\hat{y}_o) - \sum_{\substack{w \in Vocabulary, \\ w \neq o}} y_w \log(\hat{y}_w) \\ &= -\log(\hat{y}_o) = -\log(P(O = o | C = c)) = J_{naive-softmax} \end{aligned}$$

سوال دوم) (i) محاسبه ی مشتق جزئی $J_{cross-entropy}(v_c, o, U)$ نسبت به v_c :

$$\begin{aligned} \frac{\partial J(v_c, o, U)}{\partial v_c} &= \frac{\partial}{\partial v_c} (-\log(P(o|c))) = -\frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \\ &= -\frac{\partial}{\partial v_c} (\log(\exp(u_o^T \cdot v_c)) - \log(\sum_{w \in V} \exp(u_w^T \cdot v_c))) \\ &= \underbrace{-\frac{\partial}{\partial v_c} \log(\exp(u_o^T \cdot v_c))}_{I} + \underbrace{\frac{\partial}{\partial v_c} \log(\sum_{w \in V} \exp(u_w^T \cdot v_c))}_{II} = ? \\ I &= -\frac{\partial}{\partial v_c} u_o^T \cdot v_c = -u_o^T \Rightarrow v_o \\ \text{chain rule} \\ II &= \frac{\frac{\partial}{\partial v_c} \log(\sum_{w \in V} \exp(u_w^T \cdot v_c))}{\frac{\partial}{\partial v_c} \sum_{w \in V} \exp(u_w^T \cdot v_c)} \cdot \frac{\partial \sum_{w \in V} \exp(u_w^T \cdot v_c)}{\partial v_c} \\ &= \frac{1}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \cdot \sum_{w \in V} \exp(u_w^T \cdot v_c) \cdot u_w \\ &= \frac{\sum_{w \in V} \exp(u_w^T \cdot v_c) u_w}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = \frac{\sum_{w \in V} \exp(u_w^T \cdot v_c) u_w}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = \text{Softmax}(u_o^T \cdot v_c) \\ \text{Softmax}(u_o^T \cdot v_c) &= \frac{\sum_{w \in V} \exp(u_w^T \cdot v_c) u_w}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = \frac{\sum_{w \in V} \exp(u_w^T \cdot v_c) u_w}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = \frac{\sum_{w \in V} \exp(u_w^T \cdot v_c) u_w}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \\ &= -u_o + \sum_{w \in V} \frac{u_w \exp(u_w^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = -u_o + \sum_{w \in V} u_w \cdot P(w|c) \\ &= -u_o + U \cdot \hat{y}^T \\ dx^1 &\leftarrow U \quad dx^1 \leftarrow u_o \\ dx^1 &\leftarrow \hat{y}^T \quad dx^1 \leftarrow v_c \end{aligned}$$



سوال دوم(ii) حاصل سوال قبل چه زمانی صفر میشود؟

برای پیدا کردن مقادیر VC که در آن گرادیان برابر با صفر است، ما باید معادله زیر را حل کنیم:

$$U * (\hat{y} - y) = 0$$

این معادله دارای دو حالت است: زمانی که U غیرقابل معکوس کردن باشد، که این اتفاق می تواند در صورتی رخ دهد که ستون های U با یکدیگر رابطه خطی داشته باشند. به عبارت دیگر، اگر دو یا بیشتر از دو ستون U با هم رابطه خطی داشته باشند، آنگاه ماتریس U مقدار کامل نیست و معادله $U * (\hat{y} - y) = 0$ بی نهایت جواب خواهد داشت.

در این مورد، مشتق جزئی Jnaive-softmax نسبت به VC می تواند به ازای مقادیر متعدد VC، که با نقاط بهینه ی محلی متفاوت تابع مطابقت دارد، صفر باشد. با این حال، اگر U دارای ستون های مستقل خطی باشد، معادله تنها راه حل $\hat{y} - y = 0$ را خواهد داشت، به این معنی که مشتق فقط زمانی صفر خواهد بود که $\hat{y} = y$. به این معنی VC در حال حاضر در حد مطلوب است و خروجی مدل با خروجی واقعی برابر شده و نیازی به آپدیت بیشتر نیست.

سوال دوم(iii) نحوه ی تاثیر گرادیان بدست آمده روی VC:

این مشتق تفاضل بین دو عبارت است: $U * y$ و $U * \hat{y}$

- عبارت $U * \hat{y}$ یک مجموع وزن دار از بردارهای کلمات **outside** است، جایی که وزن ها براساس اختلاف بین احتمالات پیش بینی و احتمالات برچسب واقعی داده شده اند. هر بردار **outside** بر اساس سهم آن در پیش بینی واژه مرکزی C وزن گذاری می شود. این عبارت به وکتور VC کمک می کند که به سمت بردارهای **outside** که احتمالاً بیشترین تکرار را در متن دارند حرکت کند. به طور کلی، این عبارت وکتور VC را تشویق می کند تا بهترین پیش بینی را برای واژگان **outside** داشته باشد.
- عبارت $U * y$ صرفاً یک مجموع وزنی از بردارهای بیرونی است، که در آن وزن ها با احتمالات برچسب واقعی داده می شوند. هر بردار **outside** با توجه به سهم واقعی خود در **context** کلمه مرکزی C وزن دهی می شود. این عبارت به عنوان **baseline** عمل می کند و نشان دهنده سهم هر کلمه **outside** در بافت کلمه مرکزی C است. با کم کردن این عبارت از عبارت قبلی، گرادیانی به دست می آوریم که VC را از جهت هایی که به **context** کلمه مرکزی کمک نمی کنند، دور میکند و به سمت جهت هایی که در یادگیری پیش بینی کلمات **outside** در متن موثر تر هستند، می کشد.

بنابراین، هنگامی که این گرادیان از کلمه $vector\ v$ کم می شود، بردار را به گونه ای به روز می کند که توانایی آن را برای پیش بینی کلمات بیرونی در متن کلمه مرکزی بهبود می بخشد. در حالی که حساسیت آن را نسبت به جهت هایی که به زمینه مرتبط نیستند نیز کمتر می کند.

همچنین میتوان اینگونه تفسیر کرد که گرادیان بدست آمده را می توان به عنوان یک عبارت تصحیح تفسیر کرد که کلمه مرکزی بردار vc را به منظور بهبود پیش بینی مدل تنظیم می کند. عبارت اول $(\hat{y} - y)$ نشان دهنده خطا یا عدم تطابق بین مقادیر هدف پیش بینی شده و واقعی است. عبارت دوم U نشان دهنده سهم هر بردار کلمه متنی نسبت به خروجی پیش بینی شده است. با کم کردن گرادیان از بردار کلمه فعلی، می توانیم $embedding$ را به گونه ای به روز کنیم که $context$ کلمات اطراف را بهتر نشان دهد.

سوال دوم (iv) در بسیاری از موارد word embedding، از بردار های نرمال شده ی L2 به جای بردار خام آنها استفاده میشود. مسئله ی رده بندی عبارت ها به صورت مثبت و منفی را در نظر بگیرید. چه زمانی روش نرمال سازی L2 اطلاعات مفید را از بین میبرد؟ چه زمانی نه؟

نرمال سازی L2 معمولاً در برنامه های کاربردی $downstream$ استفاده می شود زیرا می تواند به کاهش اثر طول های مختلف بردارهای ورودی کمک کند و مقایسه و تجزیه و تحلیل جاسازی های مختلف را آسان تر می کند. با این حال، مواردی وجود دارد که عادی سازی L2 به طور بالقوه می تواند اطلاعات مفیدی را برای تسک $downstream$ طبقه بندی عبارات به عنوان مثبت یا منفی از بین ببرد. یکی از این موارد زمانی است که نرمال سازی L2 برای $embedding$ هایی اعمال می شود که به طور خطی مستقل نیستند، مانند سناریویی که $ux = auv$ برای کلماتی که $x \neq v$ و مقدار اسکالری مانند α . در این حالت، نرمال سازی $embedding$ ها منجر به این می شود که هر دو بردار جهت یکسانی داشته باشند، عملاً دو بردار در یک بردار جمع شوند و تمایز بین آنها از بین برود. این می تواند منجر به از دست رفتن اطلاعات مهمی شود که برای طبقه بندی عبارات به عنوان مثبت یا منفی مرتبط است.

از سوی دیگر، زمانی که تعبیه ها به صورت خطی مستقل هستند، نرمال سازی L2 اطلاعات مفیدی را برای تسک $downstream$ از بین نمی برد. در واقع، می تواند با کاهش تأثیر فیچرهای نامربوط در بردارهای ورودی و تأکید بر ویژگی های مهم، به بهبود عملکرد طبقه بندی کمک کند. در چنین مواردی، نرمال سازی L2 می تواند به ویژه هنگام برخورد با $embedding$ های با ابعاد بالا که ممکن است مستعد $overfit$ یا $underfit$ باشند، مفید باشد. اگر فاصله نسبی بین بردارهای کلمه مهم باشد، نرمال سازی L2 می تواند اطلاعات مفیدی را برای کار $downstream$ از بین ببرد. به عنوان مثال، اگر دو کلمه x و y را در نظر بگیریم که به روش خاصی به یکدیگر مرتبط هستند (مثلاً "king" و "queen")، بردارهای نرمال شده L2 آنها ممکن است از نظر جهت شبیه تر از همتایان غیرنرمال شده ی آنها شوند. این می تواند تشخیص رابطه صحیح بین کلمات را برای تسک های $downstream$ مانند قیاس کلمات یا تشابه دشوارتر کند.

با این حال، نرمال سازی L2 می تواند در زمینه های دیگر نیز مفید باشد که بزرگی مطلق بردارهای کلمه مهم نیست. به عنوان مثال، در تحلیل احساسات، جهت بردارهای کلمه ممکن است مهمتر از بزرگی آنها باشد. در چنین مواردی،

نرمال سازی $L2$ ممکن است به بهبود عملکرد کار $downstream$ با کاهش تأثیر نقاط پرت یا فیچر های نامربوط کمک کند. علاوه بر این، نرمال سازی $L2$ همچنین می تواند به بهبود پایداری و تعمیم مدل با کاهش تأثیر نمونه های آموزشی خاص یا نویز کمک کند.

سوال سوم) مشتقات جزئی $J_{naive-softmax}(v_c, o, U)$ را با توجه به هر یک از بردارهای کلمه "outside"، uw محاسبه کنید.

$$J_{naive-softmax}(v_c, o, U) = -\log(\hat{y}_o) \rightarrow \hat{y}_o = \frac{\exp(u_o^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)}$$

V : Vocabulary

$$J = -u_o^T \cdot v_c + \log\left(\sum_{w \in V} \exp(u_w^T \cdot v_c)\right)$$

$$\frac{\partial J(v_c, o, U)}{\partial u_w} = \frac{\partial}{\partial u_w} \left(u_o^T \cdot v_c - \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right)$$

$$= \frac{\partial (u_o^T \cdot v_c)}{\partial u_w} + \frac{\partial \left(\log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right)}{\partial u_w} = -v_c (\gamma - \hat{y})$$

۱, $w \neq o$:

$$\frac{\partial J}{\partial u_w} = 0 + \frac{v_c \exp(u_w^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} = v_c \cdot \hat{y}_w$$

۲, $w = o$:

$$\frac{\partial J}{\partial u_w} = -v_c + v_c \cdot \hat{y}_w$$

سوال چهارم) مشتقات جزئی $J_{naive-softmax}(v_c, o, U)$ نسبت به U .

$$\frac{\partial J(v_c, o, U)}{\partial U} = v_c (\hat{y} - \gamma)^T = \left[\frac{\partial J}{\partial u_1} \quad \frac{\partial J}{\partial u_2} \quad \frac{\partial J}{\partial u_3} \quad \dots \quad \frac{\partial J}{\partial u_{|Vocabulary|}} \right]$$

سوال پنجم) مشتق Leaky ReLU نسبت به x :

$$0 < \alpha < 1 \quad f(x) = \max(\alpha x, x) \quad \begin{cases} \alpha x & ; x < 0 \\ x & ; x \geq 0 \end{cases}$$

۱) $x < 0$: $\frac{df(x)}{dx} = \frac{d(\alpha x)}{dx} = \alpha$

در $x=0$ مشتق نداریم چون مشتق یکتا نیست

۲) $x \geq 0$: $\frac{df(x)}{dx} = \frac{d(x)}{dx} = 1$

در $x=0$ مشتق نداریم را فرض کردیم.

$$\Rightarrow f'(x) = \begin{cases} \alpha & ; x < 0 \\ 1 & ; x \geq 0 \end{cases}$$

سوال ششم) مشتق sigmoid نسبت به x :

$$\sigma = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x + 1}$$

$$\sigma'(x) = \frac{e^x(e^x + 1) - e^x e^x}{(e^x + 1)^2} = \left(\frac{e^x}{e^x + 1} \right) \left(\frac{e^x + 1 - e^x}{e^x + 1} \right) = \sigma(x)(1 - \sigma(x))$$

سوال هفتم) (i) تکرار سوال دو و سه برای neg-sample:

$$\frac{\partial J_{\text{neg-sample}}}{\partial v_c} = -\frac{\partial}{\partial v_c} (\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)))$$

$$= (\sigma(u_0^T v_c) - 1) u_0 + \sum_{k=1}^K (1 - \sigma(-u_k^T v_c)) u_k = (\sigma(u_0^T v_c) - 1) u_0 + \sum_{k=1}^K \sigma(u_k^T v_c) u_k$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_0} = \frac{\partial}{\partial u_0} (\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)))$$

$$= \frac{\sigma(u_0^T v_c) \cdot (1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} \cdot v_c = -(1 - \sigma(u_0^T v_c)) v_c = (1 + \sigma(u_0^T v_c)) v_c$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_k} = -\frac{\partial}{\partial u_k} (\log \sigma(u_0^T v_c) - \sum_{k=1}^K \log \sigma(-u_k^T v_c))$$

$$= 0 - \frac{\partial}{\partial u_k} \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) = \frac{\sigma(-u_k^T v_c) \cdot (1 - \sigma(-u_k^T v_c)) \cdot (-v_c)}{\sigma(-u_k^T v_c)}$$

$$= + (1 - \sigma(-u_k^T v_c)) v_c$$

سوال هفتم) (ii): عبارت مشترک در هر سه مشتق جزئی است. پس میتوانیم آن را یکبار در مشتق جزئی J نسبت به vc حساب کنیم و ذخیره کنیم و در محاسبه ی دو مشتق جزئی دیگر از آن استفاده کنیم.

سوال هفتم) (iii): همانطور که می بینیم، بیشترین تکراری که باید انجام دهیم $[K]$ است، اما در سافت مکس ساده باید روی تمام واژگان تکرار کنیم. در واقع محاسبه ی مخرج سافت مکس در J ساده هزینه ی محاسباتی زیادی برای **vocabulary** بزرگ دارد. بخاطر همین کلماتی میدانیم که در **context** کلمه ی مرکزی هستند را جدا گانه حساب میکنیم و به صورت رندوم از بین کلماتی که میدانیم در **context** کلمه ی مرکزی نیستند انتخاب کرده و با فرمول جداگانه حساب میکنیم تا هزینه ی محاسباتی و زمانی و حافظه ای کمتری برایمان داشته باشد.

سوال هشتم)

$$\begin{aligned}
 S \in S, \omega_s = \omega_k & \Rightarrow u_k = u_s \\
 \frac{\partial}{\partial u_k} J_{\text{neg-sample}} &= -\frac{\partial}{\partial} \left(\log \sigma(u_o^T v_c) - \sum_{k=1}^K \log \sigma(-u_k^T v_c) \right) \\
 &= -\sum_{s=1}^{|S|} \left(\frac{1}{\sigma(-u_k^T v_c)} \left(\frac{\partial}{\partial u_k} \sigma(-u_k^T v_c) \right) \right) = -\sum_{s=1}^{|S|} \frac{\sigma(-u_k^T v_c)(1 - \sigma(u_k^T v_c)) \cdot (-v_c)}{\sigma(-u_k^T v_c)} \\
 &= \sum_{s=1}^{|S|} (1 - \sigma(u_k^T v_c)) \cdot v_c = |S| (1 - \sigma(u_k^T v_c)) v_c
 \end{aligned}$$

سوال هشتم (ii)

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} = \frac{\partial}{\partial U} \sum J(v_c, w_{t+j}, U)$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial U} J(v_c, w_{t+j}, U)$$

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_c} = \frac{\partial}{\partial v_c} \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, U)$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial v_c} J(v_c, w_{t+j}, U)$$

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_w} = \frac{\partial}{\partial v_w} \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, U)$$

$$= \frac{\partial}{\partial v_w} J(v_c, w, U)$$