

با autoregressive factorization (فاکتوراسیون خود رگرسیون) ، مدل سازی زبان را می توان به مدل کردن توزیع شرطی توکن بعدی (یعنی) x با توجه به متن c (context) تقلیل داد. اگرچه ممکن است ادعا شود که یک زبان طبیعی به دلیل ترکیب بندی خود تعداد بی شماری از متن ها را مجاز می داند ، ما با در نظر گرفتن یک مجموعه محدود از متن های ممکن ، تحلیل خود را ادامه می دهیم. محدودیت زبان طبیعی نتیجه گیری ما را تحت تأثیر قرار نمی دهد.

ما با یک زبان طبیعی را به عنوان یک مجموعه متناهی از زوج های یک متن و توزیع شرطی توکن بعدی آن در نظر

$$\mathcal{L} = \{(c_1, P^*(X|c_1)), \dots, (c_N, P^*(X|c_N))\} \quad \text{می گیریم}$$

که N تعداد متن های ممکن است. برای احتساب خطا و انعطاف پذیری زبان های طبیعی فرض می کنیم $P^* > 0$ همه جا برقرار است. $\{x_1, x_2, \dots, x_M\}$ یک مجموعه شامل M توکن ممکن در زبان L را نشان می دهد. هدف مدل زبانی این است که توزیع $p_\theta(X|c)$ که توسط θ پارامتر شده با توزیع واقعی داده تطبیق یابد. زبان طبیعی L را داریم. آیا یک پارامتر θ وجود دارد که برای تمام C های درون L داشته باشیم :

$$P_\theta(X|c) = P^*(X|c)$$

برای ساخت این مدل زبانی مدل های مبتنی بر Softmax را بررسی می کنیم.

Softmax

اکثر مدل های زبانی پارامتریک از یک تابع Softmax استفاده می کنند که بر روی یک بردار متن context vector حالت پنهان (یا hidden state) h_c و یک کلمه تعبیه شده word embedding w_x برای تعریف توزیع شرطی $p_\theta(X|c)$ کار می کند. به طور خاص ، توزیع مدل معمولاً به صورت زیر نوشته می شود:

$$P_\theta(x|c) = \frac{\exp h_c^\top w_x}{\sum_{x'} \exp h_c^\top w_{x'}}$$

که h_c تابعی از c و w_x تابعی از x است. هر دو تابع توسط θ پارامتر شده اند. هم context vector (بردار متن) h_c و هم word embedding (تعبیه کلمه) w_x دارای یک بعد یکسان d هستند. حاصل ضرب داخلی $h_c^\top w_x$ را logit می نامند.

برای کمک به بحث درباره ی expressiveness مدل softmax سه معیار زیر را تعریف می کنیم:

$$H_\theta = \begin{bmatrix} h_{c_1}^\top \\ h_{c_2}^\top \\ \vdots \\ h_{c_N}^\top \end{bmatrix}; \quad W_\theta = \begin{bmatrix} w_{x_1}^\top \\ w_{x_2}^\top \\ \vdots \\ w_{x_M}^\top \end{bmatrix}; \quad A = \begin{bmatrix} \log P^*(x_1|c_1), & \log P^*(x_2|c_1) & \dots & \log P^*(x_M|c_1) \\ \log P^*(x_1|c_2), & \log P^*(x_2|c_2) & \dots & \log P^*(x_M|c_2) \\ \vdots & \vdots & \ddots & \vdots \\ \log P^*(x_1|c_N), & \log P^*(x_2|c_N) & \dots & \log P^*(x_M|c_N) \end{bmatrix}$$

$$H_\theta \in \mathbb{R}^{N \times d}, W_\theta \in \mathbb{R}^{M \times d}, A \in \mathbb{R}^{N \times M}$$

سطرهای H_θ ، W_θ و A به ترتیب با context vector ها، word embedding ها و لگاریتم احتمال توزیع داده های واقعی مطابقت دارند. H_θ به عنوان شبکه های عصبی عمیق مانند یک شبکه recurrent پیاده سازی می شود، در حالی که W_θ به عنوان lookup embedding نمونه سازی شده است. علاوه بر این یک مجموعه از ماتریس ها را مشخص می کنیم که با شیفت سطری row-wise shift بر روی A ایجاد شده اند:

$$F(A) = \{A + \Lambda J_{N,M} | \Lambda \text{ is diagonal and } \Lambda \in \mathbb{R}^{N \times N}\},$$

و $J_{N,M}$ یک ماتریس $N \times N$ است که تمام مقادیرش یک هستند. اساساً عملیات شیفت سطری row-wise shift یک عدد حقیقی دلخواه را به هر سطر A اضافه می کند. بنابراین $F(A)$ یک مجموعه نامتناهی است.

مجموعه $F(A)$ دارای دو ویژگی مهم هستند که کلید تحلیل ما می باشند.

ویژگی 1: برای هر ماتریس $A' \in F(A)$ ، اگر و تنها اگر $\text{softmax}(A') = P^*$ ، به عبارت دیگر $F(A)$ مجموعه ای تمام logit های ممکن که با توزیع واقعی داده ها تطابق دارند را تعریف می کند.

ویژگی 2: برای هر $A_1 \neq A_2 \in F(A)$ داریم: $|\text{rank}(A_2) - \text{rank}(A_1)| \leq 1$. به عبارت دیگر تمام ماتریس ها در $F(A)$ دارای مرتبه (rank) های مشابه هستند با حداکثر اختلاف 1.

بر اساس ویژگی 1 اصل (lemma) زیر را داریم:

lemma1: با توجه به پارامتر θ مدل $H_\theta W_\theta^T \in F(A)$ اگر و تنها اگر برای تمام c ها در L داشته باشیم:

$$P_\theta(X|c) = P^*(X|c)$$

حال سوال این است که آیا پارامتر θ و $A' \in F(A)$ وجود دارد بطوریکه $H_\theta W_\theta^T = A'$ این اساساً یک مسئله matrix factorization است.

ما می خواهیم مدل ماتریس های H_θ و W_θ را طوری یاد بگیرد که قادر به factorize کردن ماتریسی مانند $A' \in F(A)$ باشد. ابتدا توجه داشته باشید که برای اینکه یک factorization معتبر وجود داشته باشد مرتبه (rank) $H_\theta W_\theta^T$ باید حداقل برابر مرتبه (rank) A' باشد. علاوه بر این چون $H_\theta \in \mathbb{R}^{N \times d}$ و $W_\theta \in \mathbb{R}^{M \times d}$ بنابراین مرتبه (rank) $H_\theta W_\theta^T$ با اندازه embedding (تعبیه) d محدود شده است. در نتیجه اگر $d \geq \text{rank}(A')$ یک تقریب زن جهانی (universal approximator) از نظر تئوری می تواند A' را باز یابی کند. با این حال اگر $d < \text{rank}(A')$ دیگر مهم نیست خانواده تابع U چقدر expressive است و هیچ (H_θ و W_θ) نمی تواند از نظر تئوری A' را باز یابی کند.

استدلال فوق را به صورت زیر خلاصه می کنیم.

قضیه 1: اگر خانواده توابع U که یک تقریب زن جهانی (universal approximator) را داشته باشیم یک

$$P_\theta(X|c) = P^*(X|c) \quad \text{پارامتر مانند } \theta \text{ وجود دارد بطوریکه برای تمام } c \text{ های در } L$$

برقرار است اگر و تنها اگر $d > \min_{A' \in F(A)} \text{rank}(A')$

با ترکیب ویژگی 2 و قضیه 1 اکنون می توانیم مشکل softmax bottleneck را به صورت رسمی بیان کنیم.

نتیجه 1 (softmax bottleneck): اگر $1 - \text{rank}(A') < d$ برای هر خانواده تابع U و پارامتر θ حداقل یک متن

$$P_{\theta}(X|c) \neq P^*(X|c). \quad \text{context}$$

نتیجه‌ی فوق نشان می‌دهد که هنگامی که بعد d خیلی کوچک باشد softmax توان بیان (express) توزیع واقعی داده‌ها را ندارد. بدیهی است که این نتیجه‌گیری به یک زبان متناهی L محدود نمی‌شود. وقتی L نامتناهی است همیشه می‌توان یک زیرمجموعه‌ی متناهی از آن را انتخاب کرد و مشکل softmax bottleneck همچنان وجود دارد. در مرحله‌ی بعد با این فرض که زبان‌های طبیعی دارای مرتبه (rank) بالایی هستند نشان می‌دهیم که چرا softmax bottleneck یک مشکل است.

فرض می‌کنیم برای زبان طبیعی L ، ماتریس A (log probability) یک ماتریس با مرتبه‌ی بالا (high rank) باشد. اثبات این فرض اگرچه ممکن است به دلیل اینکه ما به توزیع واقعی داده‌ها دسترسی نداریم سخت به نظر برسد اما با استدلال‌های شهودی و مشاهدات تجربی توصیه شده است.

زبان طبیعی به شدت وابسته به متن (context) می‌باشد به عنوان مثال توکن “north” با احتمال زیاد در یک مقاله خبری درباره‌ی سیاست قبل از “korea” یا “Korean” می‌آید که بعید است در در یک کتاب درسی درباره‌ی تاریخ ایالات متحده باشد. ما فرض می‌کنیم چنین وابستگی ظریفی به متن باید منجر به یک ماتریس با مرتبه‌ی بالا (high rank) شود.

اگر A دارای مرتبه پایین (low rank) باشد به این معنی است که انسان به تعداد محدودی در حد چند صد مبنا احتیاج دارد و تمام معانی نحوی را می‌توان با (به طور بالقوه) نفی و میانگین وزن‌دار این مبانی ایجاد کرد. با این وجود یافتن مفهومی طبیعی در زبانشناسی و علوم طبیعی که با چنین مبناهایی مطابقت داشته باشد دشوار است که وجود چنین مبناهایی را زیر سوال می‌برد.

به صورت تجربی نشان داده شده‌است که مدل‌های زبانی مرتبه‌ی بالا (high rank) از مدل‌های مرتبه پایین (low rank) (rank) مرسوم عملکرد بهتری دارند. یادگیری مدل‌های زبانی مرتبه‌ی بالا (high rank) مهم است.

با توجه به این فرض که زبان‌های طبیعی زبان‌های طبیعی مرتبه‌ی بالا (high rank) هستند بدیهی است که softmax bottleneck ، expressiveness مدل را محدود می‌کند. در عمل embedding (تعبیه) d معمولاً در مقیاس 10^2 تنظیم می‌شود در حالی که مرتبه‌ی A (rank) احتمالاً می‌تواند به اندازه‌ی M (در مقیاس 10^5) باشد که به شدت بزرگتر از d هست.

Softmax در بهترین حالت یادگیری تقریبی با مرتبه پایین (low rank) از A است و تجربه‌ها نشان می‌دهد که چنین تقریبی توانایی مدل‌سازی وابسته به متن (context) را از نظر کمی و کیفی ندارد.

راه حل آسان (easy fix):

با رخ دادن softmax bottleneck بلافاصله برخی حل‌های آسان به ذهن می‌رسد. اول اینکه می‌توان از یک مدل غیرپارامتریک مانند Ngram استفاده کرد. مدل‌های Ngram توسط هیچ فرم پارامتری محدود نمی‌شوند، بنابراین با داشتن پارامترهای کافی می‌توانند هر زبان طبیعی را به صورت عمومی (universal) تقریب بزنند. دوم این که

امکان افزایش بعد d (مثلا تا حد برابر شدن با M) وجود دارد تا مدل بتواند یک ماتریس مرتبه‌ی بالا (high rank) A را بیان (express) کند.

با استفاده از این دو روش در مقایسه با استفاده از softmax با ابعاد پایین تعداد پارامترها به صورت چشمگیری افزایش میابد. به صورت دقیق‌تر یک مدل Ngram برای بیان A (express) نیاز به $N \times M$ پارامتر دارد که N به طور بالاقوه نامتناهی است. به همین ترتیب یک softmax با ابعاد بالا به $M \times M$ پارامتر برای تعبیه کلمه (word embedding) نیاز دارد. افزایش پارامترهای مدل به سادگی منجر به overfit شدن می‌شود. از روش Back-off برای کاهش overfit استفاده شده است. علاوه بر این از آنجا که مدل‌های یادگیری عمیق با جست‌وجوی گسترده‌ی پارامترها tune (تنظیم) می‌شوند افزایش ابعاد d فراتر از چندصد واحد به صرفه نیست. بدیهی است که بین expressiveness و generalization (تعمیم) در مدل‌سازی زبان یک trade off وجود دارد. افزایش ساده‌لوحانه expressiveness منجر به آسیب زدن به generalization می‌شود. حال به معرفی یک روش می‌پردازیم که expressiveness را بدون نیاز به افزایش بیش از حد پارامترها بالا ببرد.

MIXTURE of Softmaxes (یک مدل زبانی مرتبه بالا (high rank): MoS
MoS توزیع شرطی را به صورت زیر محاسبه می‌کند.

$$P_{\theta}(x|c) = \sum_{k=1}^K \pi_{c,k} \frac{\exp \mathbf{h}_{c,k}^{\top} \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_{c,k}^{\top} \mathbf{w}_{x'}}; \text{ s.t. } \sum_{k=1}^K \pi_{c,k} = 1$$

$\pi_{c,k}$ وزن پیشین (prior) یا ترکیبی (mixture) k امین مولفه است و $\mathbf{h}_{c,k}$ ، k امین بردار متن (context vector) است که با متن c (context) مرتبط است. به عبارت دیگر، MoS توزیع k عدد softmax را محاسبه می‌کند و از میانگین وزن‌دار آن‌ها به عنوان توزیع احتمال توکن بعدی استفاده می‌کند. ابتدا یک دسته از لایه‌های recurrent را در بالای X اعمال می‌کنیم تا یک ترتیب از حالت‌های پنهان (hidden states) به دست آوریم.

$$\pi_{c_t,k} = \frac{\exp \mathbf{w}_{\pi,k}^{\top} \mathbf{g}_t}{\sum_{k'=1}^K \exp \mathbf{w}_{\pi,k'}^{\top} \mathbf{g}_t} \quad \mathbf{h}_{c_t,k} = \tanh(\mathbf{W}_{h,k} \mathbf{g}_t)$$

که $\mathbf{W}_{h,k}$ و $\mathbf{w}_{\pi,k}$ پارامترهای مدل هستند.

این روش ساده و پیاده‌سازی آن آسان است و مزایای زیر را دارد:

1. بهتر شدن expressiveness در مقایسه با softmax: با همان بعد d از نظر تئوری بیشتر از (یا حداقل به همان اندازه) expressive است. این را می‌توان با این واقعیت نشان داد که MoS با $K=1$ به softmax تقلیل میابد. مهمتر از همه MoS بهتر A را تقریب می‌زند با :

$$\hat{\mathbf{A}}_{\text{MoS}} = \log \sum_{k=1}^K \Pi_k \exp(\mathbf{H}_{\theta,k} \mathbf{W}_{\theta}^{\top})$$

که Π_k یک ماتریس قطری $N \times N$ است که اعضای آن $\pi_{c,k}$ است.

از آنجا که \hat{A}_{MoS} یک تابع غیر خطی (log-sum-exp) از بردارهای متن (context vectors) و word embeddingها است می تواند مرتبه‌ی بالایی (high rank) دلخواهی داشته باشد در نتیجه MoS مانند softmax از محدودیت مرتبه (rank) رنج نمی برد.

2. بهتر شدن generalization (در مقایسه با Ngram): مدل های Ngram و softmax با ابعاد بالا expressiveness را بهبود می بخشند اما به خوبی تعمیم نمی یابند. در مقابل MoS به دلایل زیر مشکلات تعمیم را ندارند.

MoS فرایند مولد زیر را تعریف می کند:

یک متغیر گسسته تهفته k ابتدا از $\{1, 2, \dots, K\}$ نمونه برداری می شود و سپس توکن بعدی بر اساس مولفه‌ی softmax K ام نمونه برداری می شود. با این کار ما یک bias استقرایی را معرفی می کنیم که توکن بعدی بر اساس یک تصمیم گسسته نهفته تولید می شود که اغلب در مدل سازی زبان ها ایمن است. دوم، از آنجا که \hat{A}_{MoS} توسط یک تابع غیرخطی تعریف می شود و با مرتبه (rank) bottleneck محدود نمی شود در عمل می توان d را کاهش داد تا افزایش پارامترهایی که توسط ساختار mixture آورده می شود جبران شود. در نتیجه MoS در مقایسه با softmax دارای اندازه‌ی مدل مشابهی است و بنابراین مستعد overfit شدن نیست.

MIXTURE of CONTEXTS

یک رویکرد دیگر این است که بردارهای متن (context vectors) را قبل از استفاده در softmax ترکیب کنیم نه ترکیب کردن احتمال ها بعد از softmax مانند MoS. توزیع شرطی آن به صورت زیر پارامتر شده است

$$P_{\theta}(x|c) = \frac{\exp\left(\sum_{k=1}^K \pi_{c,k} \mathbf{h}_{c,k}\right)^{\top} \mathbf{w}_x}{\sum_{x'} \exp\left(\sum_{k=1}^K \pi_{c,k} \mathbf{h}_{c,k}\right)^{\top} \mathbf{w}_{x'}} = \frac{\exp\left(\sum_{k=1}^K \pi_{c,k} \mathbf{h}_{c,k}^{\top} \mathbf{w}_x\right)}{\sum_{x'} \exp\left(\sum_{k=1}^K \pi_{c,k} \mathbf{h}_{c,k}^{\top} \mathbf{w}_{x'}\right)},$$

که در آن $\pi_{c,k}$ و $\mathbf{h}_{c,k}$ پارامترهایی مشابه MoS دارند. علارغم شباهت ظاهری آن با MoS این مدل که از آن با عنوان ترکیبی از متن ها (mixture of contexts) MoC یاد می کنیم در واقع مانند softmax از مشکل محدودیت مرتبه (rank) رنج می برد. این واقعیت با تعریف $\mathbf{h}'_c = \sum_{k=1}^K \pi_{c,k} \mathbf{h}_{c,k}$ که پارامترسازی فوق را به

$$P_{\theta}(x|c) = \frac{\exp \mathbf{h}'_c{}^{\top} \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}'_c{}^{\top} \mathbf{w}_{x'}}$$

Softmax هم ارز است. بنابراین انجام عمل ترکیب (mixture) در فضای در فضای ویژگی ها تنها expressiveness را بیشتر می کند اما این واقعیت را تغییر نمی دهد که مرتبه (rank) $\mathbf{H}\theta\mathbf{W}^{\top}$ توسط embedding (تعبیه) d محدود شده است.

نتیجه گیری: تحت چارچوب matrix factorization ، expressiveness مدل های زبانی مبتنی بر softmax با ابعاد embedding کلمات محدود شده است. که softmax bottleneck نامیده می شود. مدل MoS پیشنهادی

expressiveness را نسبت به softmax بهبود می بخشد و در عین حال در مقایسه با مدل های غیر پارامتریک و افزایش ساده لوحانه ابعاد embedding کلمه از overfit شدن جلوگیری می کند.

2- روش اول: محاسبه مشتق تک تک پارامترها به صورت مستقیم.

$$y = \omega^T \max \{ 0, W^T x + c \} + b$$

$$J(\theta) = \frac{1}{N} \sum (y^* - y)^2$$

$$y = \omega_1 \max \{ 0, W_{11}x_1 + W_{21}x_2 + c_1 \} + \omega_2 \max \{ 0, W_{12}x_1 + W_{22}x_2 + c_2 \} + b$$

$$\frac{\partial J}{\partial W_{11}} = \begin{cases} W_{11}x_1 + W_{21}x_2 + c_1 < 0 & 0 \\ 0 & -\frac{2}{N} \omega_1 x_1 \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial W_{21}} = \begin{cases} W_{11}x_1 + W_{21}x_2 + c_1 < 0 & 0 \\ 0 & -\frac{2}{N} \omega_1 x_2 \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial c_1} = \begin{cases} W_{11}x_1 + W_{21}x_2 + c_1 < 0 & 0 \\ 0 & -\frac{2}{N} \omega_1 \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial \omega_1} = \begin{cases} W_{11}x_1 + W_{21}x_2 + c_1 < 0 & 0 \\ 0 & \frac{-2}{N} \sum (y^* - y) (W_{11}x_1 + W_{21}x_2 + c_1) \end{cases}$$

$$\frac{\partial J}{\partial w_{12}} = \begin{cases} w_{12}x_1 + w_{22}x_2 + c_2 < 0 & 0 \\ 0 & -\frac{2w_{12}x_1}{N} \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial w_{22}} = \begin{cases} w_{12}x_1 + w_{22}x_2 + c_2 < 0 & 0 \\ 0 & -\frac{2w_{22}x_2}{N} \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial c_2} = \begin{cases} w_{12}x_1 + w_{22}x_2 + c_2 < 0 & 0 \\ 0 & -\frac{2w_2}{N} \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial w_2} = \begin{cases} w_{12}x_1 + w_{22}x_2 + c_2 < 0 & 0 \\ 0 & \overbrace{\left(\frac{w_{12}x_1 + w_{22}x_2 + c_2}{N} \right)} - 2 \sum (y^* - y) \end{cases}$$

$$\frac{\partial J}{\partial b} = \begin{cases} \text{[scribbled out]} & -\frac{2}{N} \sum (y^* - y) \end{cases}$$

روش دوم : استفاده از backpropagation

$$J(\theta) = \frac{1}{N} \sum (y^* - y)^2$$

$$y = w_1 h_1 + w_2 h_2 + b$$

$$h_1 = \max \{0, w_{11}x_1 + w_{21}x_2 + c_1\}$$

$$h_2 = \max \{0, w_{12}x_1 + w_{22}x_2 + c_2\}$$

$$\frac{\partial J}{\partial y} = \frac{-2}{N} \sum (y^* - y)$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial y} \times \frac{\partial y}{\partial w_1} = \frac{-2}{N} \sum (y^* - y) h_1 =$$

$$\frac{-2}{N} \sum (y^* - y) \max \{0, w_{11}x_1 + w_{21}x_2 + c_1\}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial y} \times \frac{\partial y}{\partial w_2} = \frac{-2}{N} \sum (y^* - y) h_2 =$$

$$\frac{-2}{N} \sum (y^* - y) \max \{0, w_{12}x_1 + w_{22}x_2 + c_2\}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial y} \times \frac{\partial y}{\partial b} = \frac{-2}{N} \sum (y^* - y) \times 1 = \frac{-2}{N} \sum (y^* - y)$$

$$\frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial y} \times \frac{\partial y}{\partial h_1} = \frac{-2}{N} \sum (y - y^*) w_1$$

$$\frac{\partial J}{\partial h_2} = \frac{\partial J}{\partial y} \times \frac{\partial y}{\partial h_2} = \frac{-2}{N} \sum (y^* - y) w_2$$

$$\frac{\partial J}{\partial w_{11}} = \frac{\partial J}{\partial h_1} \times \frac{\partial h_1}{\partial w_{11}} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_{21} x_1 & \text{if } w_{11} x_1 + w_{21} x_2 + c_1 < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial J}{\partial w_{21}} = \frac{\partial J}{\partial h_1} \times \frac{\partial h_1}{\partial w_{21}} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_{11} x_2 & \text{if } w_{11} x_1 + w_{21} x_2 + c_1 < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial J}{\partial c_1} = \frac{\partial J}{\partial h_1} \times \frac{\partial h_1}{\partial c_1} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_1 & \text{if } w_{11} x_1 + w_{21} x_2 + c_1 < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial J}{\partial w_{12}} = \frac{\partial J}{\partial h_2} \times \frac{\partial h_2}{\partial w_{12}} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_2 x_1 & \text{if } w_{12} x_1 + w_{22} x_2 + c_2 < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial J}{\partial w_{22}} = \frac{\partial J}{\partial h_2} \times \frac{\partial h_2}{\partial w_{22}} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_2 x_2 & \text{if } w_{12} x_1 + w_{22} x_2 + c_2 < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial J}{\partial c_2} = \frac{\partial J}{\partial h_2} \times \frac{\partial h_2}{\partial c_2} = \begin{cases} -\frac{2}{N} \sum (y^* - y) w_2 \\ 0 \end{cases}$$

$o_1 = w_{12}x_1 + w_{22}x_2 + c_2$
 $o w$

-3

Sequential API

```

▶ model = Sequential()
##### put your implementations here #####
model.add(tf.keras.layers.Input(shape=data_train[0].shape))
model.add(tf.keras.layers.Flatten())
model.add(Dense(units=30, activation='relu'))
model.add(Dense(units=10, activation='softmax'))
##### put your implementations here #####

```

```

[5] ##### put your implementation here #####
model.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])
##### put your implementation here #####

```

همانطور که در کد مشاهده می‌شود غیر از لایه خروجی فقط از یک لایه Fully Connected یا Dense استفاده شده است. optimizer را rmsprop و loss را categorical_crossentropy انتخاب کردیم و تابع فعالسازی لایه آخر softmax می‌باشد. از accuracy نیز برای نمایش دقت مدل استفاده کردیم. دقت با فقط یک لایه مخفی Fully Connected بعد از 10 epoch به 88 درصد روی داده train و 87 درصد روی داده‌های validation و test رسیده است. دقت داشته باشید که برای افزایش دقت هم تعداد لایه‌ها را اضافه کردیم و هم فقط تعداد نورون‌های همان

یک لایه را اضافه کردیم که اولی تاثیر چندانی نداشت و دومی در صورت اضافه شدن تعداد بسیاری زیادی نورون کمی تاثیر گذار بود اما منجر به overfit شدن می شد.

Functional API

```
#### put your implementations here ####
model_input = tf.keras.layers.Input(shape=data_train[0].shape)
x = tf.keras.layers.Flatten()(model_input)
x = Dense(units=30, activation='relu')(x)
model_output = Dense(units=10, activation='softmax')(x)

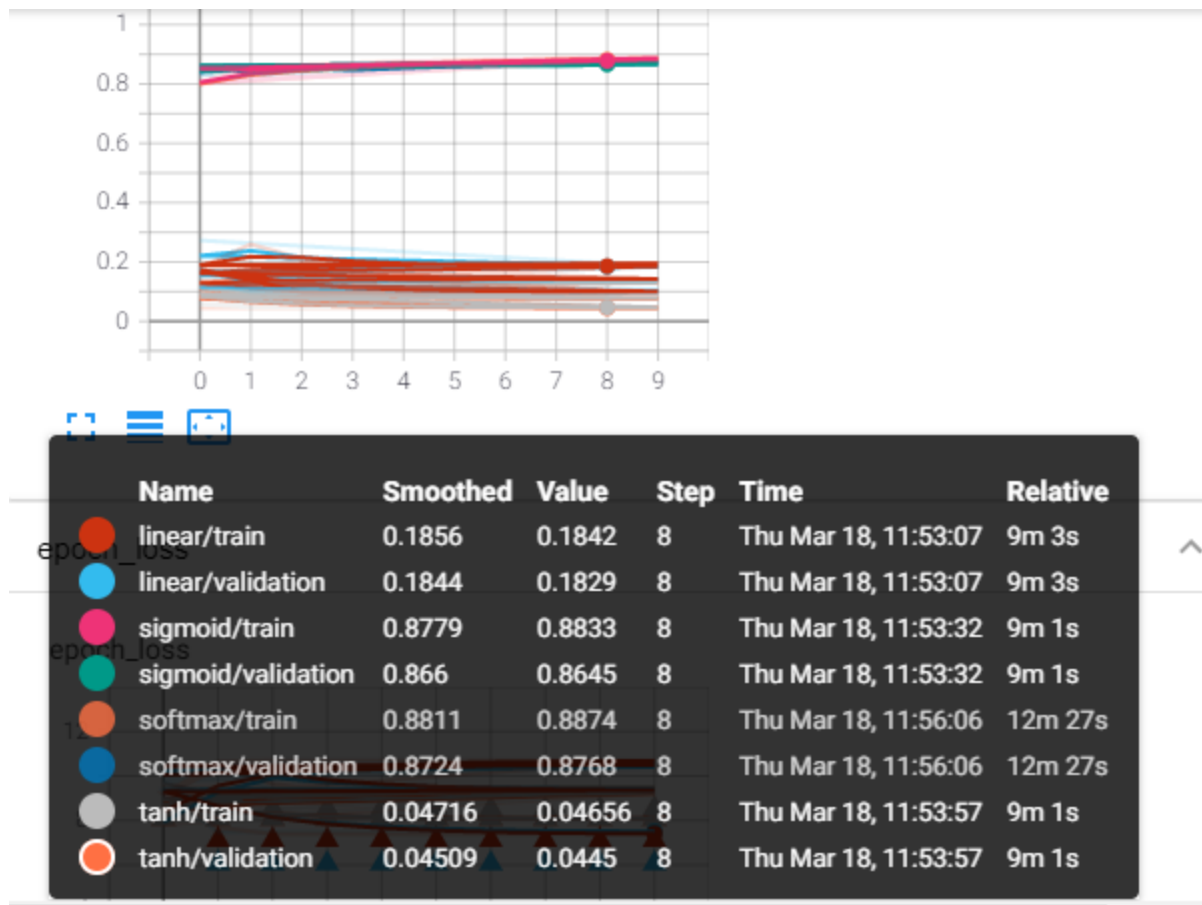
model = tf.keras.Model(model_input, model_output, name="model")
model.summary()
model.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])
history_f = model.fit(data_train, label_train,
                      epochs=10,
                      batch_size=batch_size,
                      validation_split=0.2, )
#### put your implementations here ####
```

قسمت مربوط به Functional API هیچ توضیح تازه ای درباره ی تعداد نورون ها ، تعداد لایه های مخفی ، optimizer، loss نداریم و تنها تفاوت این قسمت مربوط به چینش لایه ها می باشد که آن را در کد و در شکل فوق می بینیم.

Custom Training Loop

این قسمت هیچ توضیحی غیر از کد ندارد.

Analyzing the effect of using different activation functions



همانطور که از نمودارهای مربوط به accuracy که از tensorboard گرفته شده است مشخص است دقت برای activation function های Linear و tanh به شدت پایین می باشد مخصوصا tanh اما دقت sigmoid و softmax از 86 درصد بیشتر می باشد که بین این دو دقت softmax بهتر است.