

تمرین چهارم - بخش تئوری

سوال اول)

(الف)

ناپدید شدن گرادیان‌ها در بیشتر شبکه‌های عمیق و در هنگام backpropagation رخ می‌دهد. می‌دانیم در این شبکه‌ها تعداد لایه‌ها زیاد است بنابراین در تعداد مشتق‌هایی که نسبت به وزن‌های اولیه شبکه گرفته می‌شود بیشتر است. (طبق قاعده زنجیره‌ای) همچنین به خصوص اگر در این شبکه‌ها از فعال‌کننده‌هایی نظیر \tanh و sigmoid استفاده شده باشد و خروجی لایه‌ها به بازه‌ای محدود شود این اعداد کوچک به مراتب در هم ضرب می‌شوند در نتیجه گرادیان‌های لایه‌های اولیه بسیار کوچک (ناپدید) می‌شوند و در نتیجه در زمان آموزش وزن‌های این لایه‌ها نمی‌توانند به خوبی بروزرسانی شوند.

از طرفی CNN ها به طور کلی شبکه‌های عمیق‌تر نسبت به MLP ها حساب می‌شوند زیرا در آن‌ها لایه‌های feature extraction نیز داریم که در آن‌ها تعدادی زیادی لایه‌های پیچشی روی هم استک شده‌اند. لایه‌های اولیه پیچشی معمولاً مسئول extract کردن ویژگی‌های سطح پایین‌تر و ساده‌تر مانند گوشه‌ها و edge ها است. در نتیجه شبکه نمی‌تواند این ویژگی‌ها پایه‌ای را به خوبی شناسایی و یاد بگیرد و این ناتوانی در دیگر قسمت‌های شبکه نیز به صورت دومینویی ادامه پیدا می‌کند چون ویژگی‌های پایه‌ای را نتوانسته یاد بگیرد. همچنین به طور کلی اگر شبکه بسیار عمیق باشد ممکن است داده در روند پیمایش لایه‌ها معنی خودش را از دست بدهد و از هر لایه به لایه بعدی ممکن است اطلاعات بیشتری از دست برود که این خود باعث performance degradation است.

(ب)

اگر فرض کنیم در یک لایه پیچشی مقداردهی اولیه فیلترها (کرنل) همگی ۰ و یا به طور کلی یک مقدار ثابت c باشد با توجه به اینکه ورودی یکسان است و وزن‌های یکسان (همان درایه‌های ماتریس فیلتر) خروجی یکسان می‌دهند و در نتیجه گرادیان‌های آن نیز در فرایند آموزش مثل هم می‌شوند و بروزرسانی وزن‌ها متقارن (مثل هم) می‌شوند در نتیجه شبکه representation

ها و ویژگی‌های جدیدی یاد نمی‌گیرد و واحدها کانولوشنی در این لایه (به جز یکی) اضافی می‌شوند و این تقارن در طول آموزش شکسته نمی‌شود.

مقدار دهی Xavier برای توابع فعال‌ساز sigmoid و tanh طراحی شده و مقادیر اولیه با توجه به یک توزیع (معمولاً گاوسی) و با واریانس مشخص انتخاب می‌شوند. هدف این است که واریانس گرادیان‌ها در لایه‌های مختلف حفظ شود. واریانس نیز با فرمول پایین تعیین می‌شود.

$$var = 1/(n_{in} + n_{out})$$

تعداد نورون‌های ورودی و خروجی آن لایه به ترتیب n_{in} و n_{out} می‌باشد.

در مقداردهی He که به صورت ویژه برای Relu و LeakyRelu و شاخه‌های آن طراحی شده مانند Xavier از یک توزیع مانند یونیفرم یا گاوسی استفاده می‌شود و تفاوت آن در واریانس توزیع است که صرفاً در آن n_{in} در نظر گرفته می‌شود:

$$var = 2/n_{in}$$

برای تثبیت فرایند آموزش و دوری از performance degregation نیاز به این داریم باید واریانس activation ها را در لایه‌های مختلف کنترل شده نگه داریم تا گرادیان ها دچار vanish و یا حتی explode نشوند. از آنجایی که در ReLU تقریباً ۵۰ درصد واریانس از بین می‌رود زیرا مقادیر منفی ۰ می‌شوند. ما سعی می‌کنیم این واریانس از دست رفته را با ۲ در صورت جبران کنیم (واریانس بیشتر) و همینطور نبود n_{out} نیز برای این است که بیشتر واریانس در fwd pass رخ می‌دهد و رفتار variance propagation در لایه‌ها بیشتر توسط n_{in} تعیین می‌شود.

در نتیجه هر لایه تقریباً واریانسی یکسان را حفظ می‌کند و این امر توزیع‌های پایدار فعال‌سازی را در سراسر شبکه تضمین می‌کند. واریانس مناسب در فعال‌سازی به این معناست که گرادیان‌ها به خوبی توزیع می‌شوند، که این امر باعث می‌شود آموزش سریع‌تر و همگرایی روان‌تر انجام شود.

(ج)

همانطور که گفته شد دو مشکل اصلی در شبکه‌های عمیق، ناپدید شدن گرادیان‌ها و از دست رفتن اطلاعات مهم و تضعیف آن در رفتن از لایه به لایه دیگر است. بلوک‌های

residual با دادن ورودی یک لایه به خروجی‌اش به صورت مستقیم یک یا چند لایه را رد (بایپس) می‌کنند.

$$y = F(x) + x$$

اتصالات میان‌بر به گرادیان‌ها اجازه می‌دهند تا لایه‌های میانی را دور بزنند و مستقیماً از طریق مسیر باقی‌مانده به عقب جریان پیدا کنند. این کار مشکل ناپدید شدن گرادیان را کاهش می‌دهد، زیرا گرادیان‌ها می‌توانند آزادانه‌تر به لایه‌های اولیه منتقل شوند.

با افزودن ورودی مستقیماً به خروجی شبکه می‌تواند در صورت نیاز یک تابع همانی را یاد بگیرد. این امر بهینه‌سازی را آسان‌تر می‌کند، زیرا شبکه نیازی به یادگیری تبدیلات غیرضروری ندارد. همانطور که ذکر شد مشکل از بین رفتن اطلاعات هم با این تکنیک address می‌شود. سیگنال ورودی اصلی حفظ شده و به خروجی تبدیل شده اضافه می‌شود.

به طور کلی این اتصالات آموزش شبکه‌های بسیار عمیق را پایدار می‌کنند و به شبکه‌هایی با صدها لایه (مانند ResNet-152) اجازه می‌دهند به‌طور مؤثری همگرا شوند. بهبود جریان در backpropagation :

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \left(\frac{\partial F(x)}{\partial x} + 1 \right)$$

این +1 باعث می‌شود در با مشکل ناپدید شدن گرادیان مقابله کنیم. (اثر دادن ورودی مستقیم به خروجی)

Subject:

Date:

convolution Layer { Filter (kernel) size : 3×3 (سوال 2)
 num Filters : 1
 padding : valid \rightarrow no padding
 stride : 1

$$\text{output dimension} = \left\lfloor \frac{n_{in} - K + 2P}{S} \right\rfloor + 1$$

\uparrow kernel size
 \uparrow padding (0)
 \rightarrow stride

$$\Rightarrow n_{out} = \left\lfloor \frac{4 - 3 + 0}{1} \right\rfloor + 1 = 2$$

$$X = \begin{bmatrix} 1 & 2 & 1 & 0 \\ 2 & 1 & 0 & 1 \\ 1 & 0 & 2 & 1 \\ 0 & 1 & 2 & 1 \end{bmatrix}$$

$$W = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ -1 & 0 & 1 \end{bmatrix}$$

$$\frac{\partial \mathcal{L}}{\partial Y} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}$$

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

(الف)

$$Y = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix}$$

$$y_{ij} = \sum_{m=1}^{H_w} \sum_{n=1}^{W_w} X[i+m, j+n] w[m, n]$$

For stride: 1
 and no padding

$$\begin{cases} y_{11} = \sum_{m=1}^3 \sum_{n=1}^3 X[1+m, 1+n] \cdot w[m, n] \\ y_{12} = \sum_{m=1}^3 \sum_{n=1}^3 X[1+m, 2+n] \cdot w[m, n] \\ y_{21} = \sum_{m=1}^3 \sum_{n=1}^3 X[2+m, 1+n] \cdot w[m, n] \\ y_{22} = \sum_{m=1}^3 \sum_{n=1}^3 X[2+m, 2+n] \cdot w[m, n] \end{cases}$$

$$\begin{cases} y_{11} = 1 + (-2) + 0 + 0 + 1 + 0 + (-1) + 0 + 2 = 1 \\ y_{12} = 2 + (-1) + 0 + 0 + 0 + 0 + (-1) + 0 + 0 + 1 = 1 \\ y_{21} = 2 + 0 + (-1) + 0 + 0 + 0 + (-2) + 0 + 0 + 2 = 1 \\ y_{22} = 1 + 0 + 0 + 0 + 2 + (-1) + (-1) + 0 + 1 = 2 \end{cases}$$

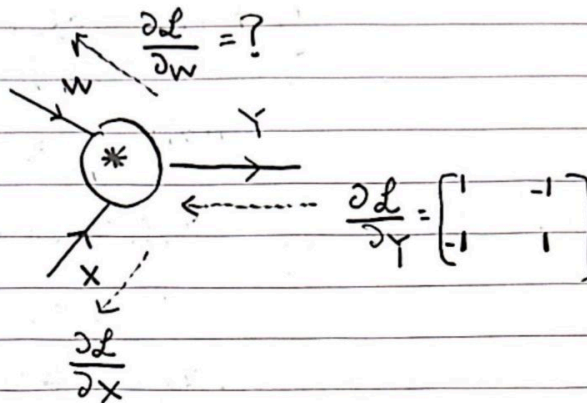
CS

Scanned with CamScanner

a.m

$$\Rightarrow Y = \begin{bmatrix} 1 & 1 \\ 1 & Y \end{bmatrix} \quad \text{end of forward pass}$$

backward pass:



we have this

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial w} \quad \xrightarrow{\text{we need this}} \quad \frac{\partial Y[i,j]}{\partial w[m,n]} = X[i+m, j+n]$$

$$\Rightarrow \frac{\partial L}{\partial w[m,n]} = \sum_i \sum_j \frac{\partial L}{\partial Y[i,j]} \frac{\partial Y[i,j]}{\partial w[m,n]} \quad \xrightarrow{\quad} \quad \frac{\partial Y}{\partial w[m,n]} = \sum_i \sum_j X[i+m, j+n]$$

$$= \sum_i \sum_j \frac{\partial L}{\partial Y[i,j]} X[i+m, j+n]$$

$$\Rightarrow \frac{\partial L}{\partial w[1,1]} = \sum_{i=1}^Y \sum_{j=1}^Y \frac{\partial L}{\partial Y[i,j]} \cdot X[i+1, j+1] = 1 + (-Y) + (-Y) + 1 = -Y$$

$$\frac{\partial L}{\partial w[1,Y]} = \sim \sim X[i+1, j+Y-1] = Y + (-1) + (-1) + 0 = 0$$

$$\frac{\partial L}{\partial w[1,0]} = \sim \sim X[i+1, j+0-1] = 1 + 0 + 0 + 1 = Y$$

$$\frac{\partial L}{\partial w[Y,1]} = \sim \sim X[i+Y-1, j+1-1] = Y + (-1) + (-1) + 0 = 0$$

$$\frac{\partial L}{\partial w[Y,Y]} = \sim \sim X[i+Y-1, j+Y-1] = 1 + 0 + 0 + Y = Y$$

$$\frac{\partial \mathcal{L}}{\partial w_{[2,2]}} = \dots \cdot X[i+2-1, j+2-1] = 0 + (-1) + (-2) + 1 = -2$$

$$\frac{\partial \mathcal{L}}{\partial w_{[2,1]}} = \dots \cdot X[i+2-1, j+1-1] = 1 + 0 + 0 + 1 = 2$$

$$\frac{\partial \mathcal{L}}{\partial w_{[1,2]}} = \dots \cdot X[i+1-1, j+2-1] = 0 + (-2) + (-1) + 2 = -1$$

$$\frac{\partial \mathcal{L}}{\partial w_{[1,1]}} = \dots \cdot X[i+1-1, j+1-1] = 2 + (-1) + (-2) + 1 = 0$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w} = \begin{bmatrix} -2 & 0 & 2 \\ 0 & 2 & -1 \\ 2 & -1 & 0 \end{bmatrix}$$

we have this

$$\frac{\partial \mathcal{L}}{\partial X} = \frac{\partial \mathcal{L}}{\partial Y} \cdot \frac{\partial Y}{\partial X} \rightarrow \text{we need this}$$

$$\frac{\partial \mathcal{L}}{\partial x_i} = \sum_k \frac{\partial \mathcal{L}}{\partial y_k} \frac{\partial y_k}{\partial x_i}$$

$$Y[i, j] = \sum_m \sum_n X[i+m-1, j+n-1] \cdot w[m, n]$$

بعنوان مثال $X[1,1]$ فقط در محاسبه $y[1,1]$ ظاهر می شود پس

$$\frac{\partial \mathcal{L}}{\partial X[1,1]} = \frac{\partial \mathcal{L}}{\partial Y[1,1]} \cdot \frac{\partial Y[1,1]}{\partial X[1,1]}$$

$$= \frac{\partial \mathcal{L}}{\partial Y[1,1]} \cdot w[1,1] = 1 \cdot 1 = 1$$

* یعنی باید ببینیم که $X[1,1]$ در کدام $Y[i,j]$ ها contribution داشته باشد *

مگر می توانیم contribution را به صورت زیر بنویسیم

contributes
→

$$X[1,1] \rightarrow Y[1,1]$$

$$X[1,2] \rightarrow Y[1,1], Y[1,2]$$

$$X[1,3] \rightarrow Y[1,1], Y[1,2]$$

$$X[1,4] \rightarrow Y[1,2]$$

$$X[2,1] \rightarrow Y[1,1], Y[2,1]$$

$$X[2,2] \rightarrow Y[1,1], Y[1,2], Y[2,1], Y[2,2]$$

$$X[2,3] \rightarrow Y[1,1], Y[1,2], Y[2,1], Y[2,2]$$

$$X[2,4] \rightarrow Y[1,1], Y[2,2]$$

$$X[3,1] \rightarrow Y[1,1], Y[2,1]$$

$$X[3,1] \rightarrow Y[2,1]$$

$$X[3,2] \rightarrow Y[1,1], X[1,2], X[2,1], X[2,2]$$

$$X[3,2] \rightarrow Y[2,1], Y[2,2]$$

$$X[3,3] \rightarrow Y[1,1], Y[1,2], X[2,1], X[2,2]$$

$$X[3,3] \rightarrow Y[2,1], Y[2,2]$$

$$X[3,4] \rightarrow Y[1,2], Y[2,2]$$

$$X[3,4] \rightarrow Y[2,2]$$

* در هر ترمینال $X[i,j]$ برای i و j وابستان ظاهر نشده در Y را در نظر نمی گیریم.

$$\frac{\partial L}{\partial X[i,j]} = \sum_k \frac{\partial L}{\partial y_k} \frac{\partial y_k}{\partial X[i,j]} \quad \text{for } k \text{ in above } (a,b)$$

$$\frac{\partial L}{\partial X[1,2]} = \frac{\partial L}{\partial y[1,1]} \cdot w[1,1] + \frac{\partial L}{\partial y[1,2]} \cdot w[1,2] = (-1) + (-1) = -2$$

$$\frac{\partial L}{\partial X[1,3]} = \frac{\partial L}{\partial y[1,1]} \cdot w[1,1] + \frac{\partial L}{\partial y[1,2]} \cdot w[1,2] = 0 + 1 = 1$$

$$\frac{\partial L}{\partial X[1,4]} = \frac{\partial L}{\partial y[1,2]} \cdot w[1,2] = 0$$

$$\frac{\partial L}{\partial X[2,1]} = \frac{\partial L}{\partial y[1,1]} \cdot w[2,1] + \frac{\partial L}{\partial y[2,1]} \cdot w[1,1] = 0 + (-1) = -1$$

$$\frac{\partial L}{\partial X[2,2]} = \frac{\partial L}{\partial y[1,1]} \cdot w[2,1] + \frac{\partial L}{\partial y[1,2]} \cdot w[2,2] + \frac{\partial L}{\partial y[2,1]} \cdot w[1,1] + \frac{\partial L}{\partial y[2,2]} \cdot w[2,2] = 1 + 1 + 1 + 1 = 4$$

$$\frac{\partial L}{\partial X[2,3]} = \frac{\partial L}{\partial y[1,1]} \cdot w[2,1] + \frac{\partial L}{\partial y[1,2]} \cdot w[2,2] + \frac{\partial L}{\partial y[2,1]} \cdot w[1,1] + \frac{\partial L}{\partial y[2,2]} \cdot w[2,2] = -3$$

sam

$$\frac{\partial \mathcal{L}}{\partial Y[c, c]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] = 1$$

$$\frac{\partial \mathcal{L}}{\partial Y[c, 1]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, 1] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, 1] = -1$$

$$\frac{\partial \mathcal{L}}{\partial Y[c, 2]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, 2] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, 2] + \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, 2] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, 2] = 0$$

$$\frac{\partial \mathcal{L}}{\partial Y[c, c]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] = 2$$

$$\frac{\partial \mathcal{L}}{\partial Y[c, c]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] = -2$$

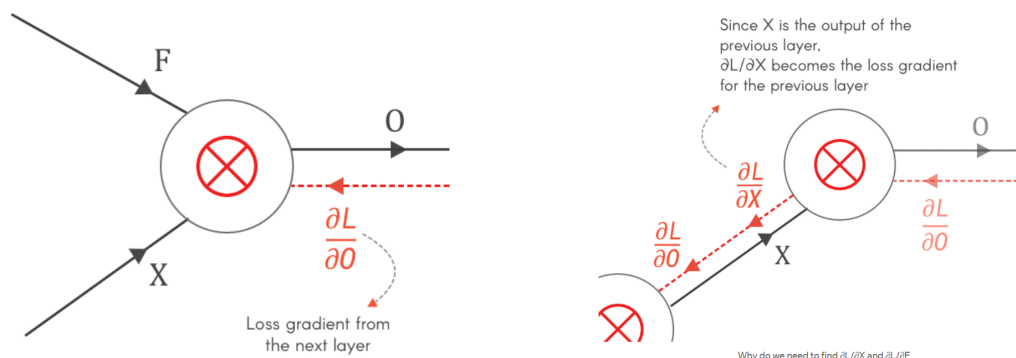
$$\frac{\partial \mathcal{L}}{\partial Y[c, 1]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, 1] = -1$$

$$\frac{\partial \mathcal{L}}{\partial X[c, 2]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, 2] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, 2] = -1$$

$$\frac{\partial \mathcal{L}}{\partial X[c, c]} = \frac{\partial \mathcal{L}}{\partial Y[c, 1]} \cdot w[c, c] + \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] = -1$$

$$\frac{\partial \mathcal{L}}{\partial X[c, c]} = \frac{\partial \mathcal{L}}{\partial Y[c, 2]} \cdot w[c, c] = 1$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial X} = \begin{bmatrix} 1 & -2 & 1 & 0 \\ -1 & 2 & -2 & 1 \\ -1 & 0 & 2 & -2 \\ 1 & -1 & -1 & 1 \end{bmatrix}$$



همینطور که در شکل‌های بالا می‌بینیم، می‌دانیم در انتشار رو به جلو خروجی هر لایه کانولوشنی ورودی لایه بعدی‌اش می‌باشد. اگر هر کدام از این لایه‌ها را به عنوان یک بلوک با ۳ عنصر ورودی، فیلتر و خروجی نگاه کنیم. به ما کمک می‌کند تا ببینیم که در انتشار رو به عقب ورودی هر بلوک، گرادیان خطای لایه بعدی نسبت به خروجی‌اش است که این خروجی با توجه به ورودی‌اش حساب شده ($conv(X, W)$) و از طرفی ورودی لایه بعد درواقع همان خروجی لایه کنونی است. پس درواقع از مقدار گرادیان خطای لایه بعد نسبت به ورودی‌اش (خروجی لایه کنونی) می‌توان به عنوان گرادیان خطا نسبت به خروجی این لایه نگاه کرد.

در شبکه‌های عصبی کانولوشنی، گرادیان تابع خطا نسبت به ورودی هر لایه نشان می‌دهد که هر پیکسل یا ویژگی در ورودی چقدر در خطای نهایی تاثیر داشته است. این گرادیان $\frac{\partial L}{\partial X}$ به عنوان سیگنال خطا به لایه‌های قبلی منتقل می‌شود. به عبارت دیگر، خروجی هر لایه به عنوان ورودی لایه قبلی در نظر گرفته شده و این روند به صورت بازگشتی ادامه می‌یابد تا وزن‌های تمام لایه‌ها برای کاهش خطای نهایی به روزرسانی شوند.

سوال سوم)

برای استفاده از RNN ابتدا مانند هر وظیفه NLP دیگر باید ورودی و درواقع اینجا کلمات را تبدیل به بردارهای عددی کنیم. مثلا می‌توانیم از الگوریتم word2vec استفاده کنیم. این دنباله از امبدینگ‌ها کلمه به کلمه وارد شبکه می‌شود. در هر گام شبکه امبدینگ کلمه کنونی و خروجی لایه hidden قبلی را دریافت می‌کند و خروجی لایه مخفی کنونی را با توجه به فرمول پایین محاسبه می‌کند:

$$h_t = f(W_h \cdot w_t + U_h \cdot h_{t-1} + b_h)$$

بعد از پردازش تمام n کلمه ورودی خروجی لایه مخفی آخر به عنوان نمایشی از کل جمله عمل می‌کند. این استیت آخر contextual meaning و احساسات بیان شده در جمله را در بر دارد و در ادامه این نمایش معنادار از جمله وارد یک لایه fully connected می‌شود و در نهایت یک لایه softmax برای پیشبینی مقادیر بین ۰ تا ۴ همینطور از تابع خطا CE استفاده می‌کنیم. به طور خاص در زمان inference (آزمایش) pipeline پایین اجرا می‌شود.

- کلمات جمله ورودی به از طریق یک الگوریتم امبدینگ به بردارهای عددی تبدیل می‌شوند.
- شبکه هر کلمه را به صورت sequential و با نگه داشتن حالت مخفی در هر گام زمانی پردازش می‌کند.
- خروجی لایه آخر مخفی final hidden state است که نمایشی معنادار از جمله ورودی است.
- این استیت مخفی پایانی وارد یک لایه fully connected می‌شود با فعالساز softmax می‌شود می‌شود.
- کلاس با بیشترین امتیاز به عنوان لیبل آن جمله خروجی داده می‌شود.

البته حالت‌های دیگری نیز از embedding وجود دارد مانند BPE یا WordPiece. این‌ها الگوریتم‌های subword tokenization هستند و هر کلمه می‌تواند تبدیل به چند بردار شود. اما دیگر مراحل ثابت است.

(الف)

- با توجه به تکنیک‌هایی که وجود دارد دو حالت را می‌توان برای تعداد استفاده شدن از softmax متصور شد. یک حالت این است که صرفاً در آخر مانند روندی که توصیف شد ۱ بار softmax را صدا بزنیم و کلاس با بیشترین امتیاز را اعلام کنیم. راه دیگر این است که در هر a گام زمانی خروجی مخفی آن لایه را به softmax (پس از گذشتن از لایه fully connected) بدهیم یا هر بار که کلمه خاص یا علامت نگارشی خاص (conditional). در این روش پس از پردازش تمام کلمات ما تعدادی خروجی softmax داریم که هر کدام sentiment آن جمله را از اول تا جای مشخصی در بر دارد. ما می‌توانیم از روش‌های مختلف مانند میانگین‌گیری برای aggregate کردن این خروجی‌ها استفاده کنیم و سپس کلاس با بیشترین امتیاز را اعلام کنیم.
- در روش اول ۱ بار softmax صدا می‌شود.

- در روش دوم بین ۱ تا n بار softmax می‌تواند صدا شود. (این روش برای جملات طولانی بهتر است زیرا کلماتی که معنای جمله را دربر دارند به صورت یونیفرم در جمله پخش نشدند و همینطور با ناپدید شدن گرادیان نیز می‌توان مقابله کرد.)
در نظر داریم اگر از الگوریتم‌های subword tokenization استفاده کنیم به ازای جمله n کلمه‌ای می‌توانیم بیش از n بردار embedding داشته باشیم و در نتیجه این امکان وجود دارد در این استراتژی بیش از n بار softmax صدا شود.

(ب)

هر خروجی softmax یا همان \hat{y}_n در شبکه موردنظر برای تحلیل احساسات، یک توزیع احتمالاتی روی کلاس‌های احساسات (مثلاً منفی، خنثی، مثبت) را نشان می‌دهد. این توزیع نشان می‌دهد که جمله (یا کلمه، بسته به استراتژی خروجی) با چه احتمالی به هر کلاس تعلق دارد. مقادیر این توزیع از طریق یک لایه کاملاً متصل و سپس یک تابع Softmax به دست می‌آید. در نهایت، کلاسی که بالاترین احتمال را دارد، به عنوان برچسب نهایی احساسات پیش‌بینی می‌شود. بسته به استراتژی استفاده‌شده (خروجی در زمان نهایی، در هر گام، یا به صورت شرطی)، این توزیع می‌تواند در یک یا چند مرحله تولید شود.
پس درواقع هر خروجی softmax نشان‌دهنده یک توزیع احتمالی روی کلاس‌های احساسات می‌باشد که می‌توان به صورت زیر نشان داد:

$$\hat{y}_n = [P(c_0|x), P(c_1|x), \dots, P(c_4|x)]$$

(ج)

اگر منظور از خروجی مدل در هر گام خروجی، خروجی لایه‌های مخفی باشد، همانطور که توضیح داده شد ورودی‌ها در هر گام زمانی t شامل امبدینگ کلمه کنونی یعنی x_t و خروجی لایه مخفی قبل یعنی h_{t-1} می‌باشد که شامل معنای جمله تا کلمه t ام است.
اگر منظور از خروجی در هر گام زمانی خروجی softmax است، خروجی آن لایه مخفی وارد یک لایه تماماً متصل می‌شود و سپس به عنوان ورودی به softmax داده می‌شود.

سوال چهارم)

(الف)

با توجه به مقاله لایه پیچشی آخر در نظر گرفته می‌شود. در این فرمول y_c امتیاز کلاس پیش‌بینی شده برای نمونه مورد نظر است. (قبل از اعمال softmax) پس درواقع امتیاز هر کانال لایه پیچشی مورد نظر حاصل global-average pooling روی گرادیان y_c نسبت به هر درایه k feature map بعد از activation است. H و W نیز ابعاد این feature map را نشان می‌دهند.

برای تفسیر این فرمول در نظر می‌گیریم که گرادیان حساب شده در این فرمول نشان‌دهنده این است که چقدر تغییر در A_{ij}^k تغییر در امتیاز کلاس انتخاب شده برای آن نمونه می‌شود. (y_c) اگر این گرادیان زیاد بود پس تاثیرپذیری زیادی دارد و برعکس. همینطور میانگین‌گیری روی کل feature map باعث می‌شود تا میزان تاثیر و تغییرات محاسبه شده محدود به برخی نقاط خاص آن feature map نباشد و نویزها را کمتر کنیم و امتیاز نهایی محاسبه شده برای وزن آن کانال را با توجه به کل آن محاسبه کنیم. همچنین با تمرکز صرفا بر روی logit کلاس پیش‌بینی شده مطمئن می‌شویم تنها نواحی مربوط به کلاس c شناسایی می‌شوند. پس به طور کلی این خاصیت گرادیان که میزان حساسیت خروجی نهایی به مقادیر میانی را اندازه‌گیری می‌کند پایه و اساس این فرمول و روش است.

(ب)

می‌دانیم فعالساز ReLU تنها مقادیر مثبت را نگه می‌دارد و تمامی مقادیر منفی feature map را با ۰ جایگزین می‌کند از طرفی به طول کلی در میانگین وزندار feature map ها که وزن‌های آن با توجه به کلاس پیش‌بینی شده c محاسبه شده مقادیر مثبت نشان‌دهنده positive contribution به کلاس c می‌باشد زیرا این پیکسل‌ها بیشترین تاثیر را روی کلاس c داشته‌اند و برعکس یعنی مقادیر منفی contribution ای به این کلاس نداشته‌اند و احتمالا contribution آن‌ها برای کلاس‌های دیگر بوده است و یا مربوط به ویژگی‌هایی است که برای تشخیص این کلاس مهم نبوده‌اند. (ذکر شده در قسمت 3 مقاله). اگر این ReLU را اعمال نکنیم برخی اوقات هیت مپ بدست آمده نواحی مربوط به کلاس‌های دیگر را هایلایت می‌کند. همینطور این حذف این مقادیر باعث می‌شود هیت مپ نهایی واضح‌تر و discriminative تر باشد و نویزها را در نظر نگیرد.