

$$z = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

$$\{ w_{ij} = 1 \mid i=1,2,3, j=1,2 \}$$

نرم: تنها یک ورودی مقدار 1 باشد  $\Leftrightarrow Z < 0$  نرم

$$1 \times 1 + 0(1) + 0(1) + b_1 = 0 \Rightarrow b_1 < -1$$

خط در جواب  $b_1 = -1.5$  نظر نرم

برای  $n_2$  نرم می بینیم تنها دو ورودی 1 باشد بنابراین نرم

نرم  $Z < 0$

$$(1)(1) + (1)(1) + (0)(1) + b_2 < 0$$

$$2 + b_2 < 0 \Rightarrow b_2 = -2.5$$

$m_1$	$m_2$	$m_3$	$n_1$	$n_2$
0	0	0	off	off
0	0	1	off	off
0	1	0	off	off
✓ 0	1	1	0.5	off
1	0	0	off	off
✓ 1	0	1	0.5	off
✓ 1	1	0	0.5	off
✓ 1	1	1	1.5	0.5

$m_1$	$m_2$	$m_3$	$Z(n_1)$	$Z(n_2)$	$Y$
0	1	1	1	0	1
1	0	1	1	0	1
1	1	0	1	0	1
1	1	1	1	1	0

$$w_{10} + b_0 \geq 0 \rightarrow w_{10} = 1, b_0 \geq -1$$

$$w_{10} + w_{20} + b_0 < 0$$

$$1 + w_{20} - 1 < 0 \rightarrow w_{20} < 0$$

$$\left\{ \begin{aligned} &w_{ij} = 1 \text{ where } i=1,2,3, j=1,2 \\ &w_{10} = 1, w_{20} = -1 \\ &b_1 = -1.5 \\ &b_2 = -2.5 \\ &b_0 = -1 \end{aligned} \right.$$

$w_{20} = -1$
$b_0 = -1$
$w_{10} = 1$

Approach 1

Simple Logistic Regression (one neuron)

$$\hat{y} = \sigma(w_n + b_e)$$

Approach 2

Simple softmax regression (two neurons)

$$\hat{y} = \text{Softmax}(w_s x + b_s) = [\hat{y}_1, \hat{y}_2]^T$$

حقوقی و حقوقی نیست یک شبکه عصبی دلیلی بر Performance نیست می باشد و شبکه های Softmax اصولاً جهت مدارهای multi-classification طراحی گردیده اند در این مثال با نری کلاسیک داریم و یک نرون خروجی به تنهایی کارآمدی لازم دارد علاوه بر دو تابع Approach 1 و 2 عادی می

باشد [اوه] می باشد پس Approach 2 ظاهراً هیچ ارجحیتی در این معنی Approach 1 ندارد

در این حالت بتوان ابطال بین در این صورت کارایی مردودی است در حالتی حقوقی Approach 2 نیست است

در این صورت کارایی مردودی است در حالتی حقوقی Approach 2 نیست است

در Approach 2 نرون دو خروجی را تعیین می کند

مقدار خروجی در این نرون

$$S(w_s x + b_s) = \frac{e^{w_{s2}x + b_{s2}}}{e^{w_{s2}x + b_{s2}} + e^{w_{s1}x + b_{s1}}}$$

صورت و مخرج را بر تقسیم می کنیم

$$\frac{1}{1 + e^{\frac{w_{s1}x + b_{s1} - w_{s2}x - b_{s2}}{1}}}$$

$$-(w_{s2}x + b_{s2}) = -[(w_{s2} - w_{s1})x + (b_{s2} - b_{s1})]$$

معادله خبری

Approach 1

$$\frac{1}{1 + e^{-(w_{s2}x + b_{s2})}}$$

نمی توانیم

$$(w_e, b_e) = (w_{s2} - w_{s1}, b_{s2} - b_{s1})$$

۱- در مورد multi-class classification به دایره زیر توجه کنید Cross entropy مناسب می باشد

$$L(y - \hat{y}) = - \frac{1}{10} \sum_{i=1}^{10} (y_i - \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i))$$

در این تابع  $\hat{y}$  داریم

$$L(y - \hat{y}_i) = \begin{cases} \log \hat{y}_i & \text{if } y_i = 1 \\ \log (1 - \hat{y}_i) & \text{if } y_i = 0 \end{cases}$$

نکته: به ازای هر نمونه ورودی حتما در هر ۱۰ تره خروجی محاسبه می گردد هم خروجی ای که می بایست صفر باشد و هم بقیه خروجی ها که می بایست صفر نباشند.

نکته دیگر: می خواهیم از یک شبکه binary-class classification به یک شبکه multi-class classification

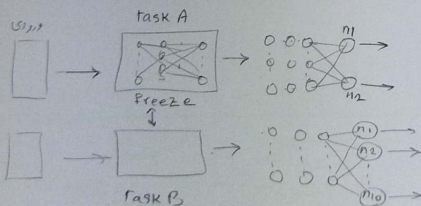
Transfer learning انجام دهیم چه پارامترهای جدیدی را باید tune کنیم؟

در هر epoch تعدادی از لایه ها را فریز می کنیم و بقیه را آپدیت می کنیم

۱- به ازای تعداد دیتاهای از لایه های میانی را Freeze کردیم می بایست تغییری در output انجام دهیم

و ۲ خروجی مدل قبل را به ۱۰ خروجی تبدیل نمائیم که Fully connected می باشد

بنابراین هر لایه آخر از  $n \times 2 + n$  پارامتر تبدیل به  $n \times 10 + n$  پارامتر می گردد

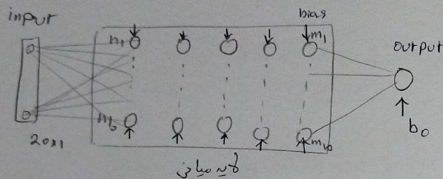


۲- به ازای تعداد epoch های کمتر از مرحله قبل می توانیم شبکه را آپدیت می دهیم اما در نهایت در ابتدا همان درختی شبکه را داریم می بایست

۳- حال قسمت Freeze شده را جدا کرده و روی چند لایه آخر علاوه بر output ها خلاصه درخت را tune می کنیم

# قسمت سوم از سوین سوال

weight sharing یا همان transfer learning معکس است باعث افزایش bias یا variance گردد  
 برای مثال اگر تعداد لایه های Freeze زیاد باشد و دیتاست جدید هم داده های زیادی داشته باشد ممکن  
 است شبکه نتواند ~~یاد~~ لایه های آزاد به خوبی train سکود و high-bias اتفاق بیافتد  
 همچنین اگر دیتاست داده های کمی داشته باشد و تعداد لایه های Freeze کم باشند شبکه جدید  
 با داده های دیتاست جدید adapt می سکود و high-variance اتفاق می افتد بنابراین انتخاب  
 تعداد لایه های Freeze بسیار مهم است.

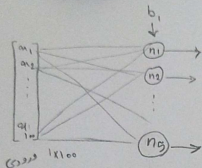


# قسمت چهارم از سوین سوال

$$\begin{aligned}
 & 20 \times 10 + 10 && \text{تعداد پارامترهای ورودی به لایه اول میانی} \\
 & + && \\
 & 4 (10 \times 10 + 10) && \text{لایه های میانی} \\
 & + && \\
 & 10 \times 1 + 1 && \text{لایه خروجی}
 \end{aligned}$$

$$\boxed{661}$$

Method 1:



Fully connected

$$g(m_1 w_{11} + m_2 w_{21} + \dots + m_{100} w_{100} + b_1) \rightarrow \text{output of } n_1$$

$$100 + 1 = 101$$

$$101 \times 5 = 505$$

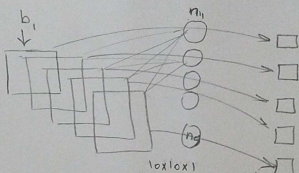
تعداد پارامترها به ازای هر خروجی

تعداد کل پارامترها

Method 2:



ورودی (10x10x1)



5 فیلتر

output

$$S(m_1 w'_{11} + m_2 w'_{21} + \dots + m_{100} w'_{100} + b_1) \rightarrow \text{output of } n_1$$

$$100 + 1 = 101$$

$$101 \times 5 = 505$$

تعداد پارامترها به ازای هر خروجی

تعداد کل پارامترها

باید این را در نظر بگیریم که تعداد فیلترها می باشد. آموزش دیده شود در هر دو مدل می باشد. لذا از این جهت هر دو مدل می باشد. activation fun. می باشد و همچنین

شکل ورودی

باز به این که Softmax عملکرد بهتری نسبت به Logistic دارد. threshold نیز دارد. به نظری رسید علی رقم پارامترهای

فیلترها و معادلات فیلترها در آموزش هر دو مدل. هنوز مدل 2 CNN بهتر عمل نماید (در روی تصویر + Softmax Func)



$$L_{CE}(\hat{y}, y) = - \sum_{i=1}^n y_i \log \hat{y}_i$$

$y = (\text{dog}, \text{cat}, \text{iguana}, \text{mouse})$

a)  $\hat{y} = (0.25, 0.25, 0.3, 0.2)^T \Rightarrow L_{CE} = - \log 0.3$   
 $y = (1, 0, 1, 0)^T$

نقشه اول (a)

$\hat{y} = (0, 0, 0.4, 0.6)^T \Rightarrow L_{CE} = - (0 \log 0 + 0 \log 0 + 1 \log 0.4 + 0 \log 0)$   
 $= - \log 0.4$

نقشه دوم (b)

$$- \log 0.4 < - \log 0.3$$

نقشه سوم (c)

خطای نقشه دوم را کمتر محاسبه کرده در صورتیکه در نقشه دوم misprediction داریم و این دلیل است که میزان خطای نقطه دوم در خروجی مورد انتظار محاسبه می‌کند. زیرا مقدار 1 در جبهه خروجی‌ها صفری نباشد.

$y \log \hat{y}$   
 $(0, 0, 1, 0)$

مقدار 1 برای خروجی معرف dog, cat, mouse صفری باشد بنابراین در نقشه دوم سبک قادر به تشخیص خطا در تشخیص نادرست mouse نخواهد بود.

$(1-y) \log (1-\hat{y})$  عبارت برای هر همه خروجی‌ها کدام است به فرمول مذکور عبارت

اضافه نکرد این عبارت خطا در خروجی‌های دیگر وقتی  $y=0$  است محاسبه می‌کند.

$$= \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

نقشه چهارم (d)

استفاده از دقت یا accuracy به جای تابع cross-entropy loss ایده خوبی نیست زیرا در مرحله backpropagation نیاز به تابعی داریم که قابل مشتق‌گیری باشد تا بتوانیم آن را به همراه update نمود علاوه بر تابع accuracy طبق تقریب شده در سال 1991 هر روزی با Example به نتایج قابل محاسبه نمی‌باشد و به جای try کردن کل دقت حاصل محاسبه می‌باشد.

(e) تابع Softmax یک تابع احتمال به صورت زیر می باشد

$$S(y_i) = \frac{e^{y_i}}{e^{y_1} + e^{y_2} + \dots + e^{y_i} + \dots + e^{y_K}}$$

وقتی که تعداد خروجی ها K باشد

این تابع هیچ وقت 0 یا 1 نمی گردد زیرا طبق بهمان خلف

اما تابع نمایی هیچ وقت صفر نمی گردد  $\times$

if  $\left\{ \begin{array}{l} (S(y_i) = 1 \Rightarrow e^{y_j} = 0 \quad \forall j \neq i \\ (S(y_i) = 0 \Rightarrow e^{y_i} = 0 \end{array} \right. \quad \times$

بنابراین  $\hat{y}_i$  که وارد تابع loss function می گردد هیچ وقت 1 نمی گردد پس  $y_i \log \hat{y}_i$  وقتی  $y_i = 1$  محاسبه می شود  $\log \hat{y}_i$  مخالف صفر است (یعنی  $\hat{y}_i \neq 1$ ) و صحت خطای ندارند به صفر نزدیک می شود اما هیچگاه 0 صفر نمی گردد.

(a) تعداد پارامترهای شبکه A برای Tune شدن کمتر از پارامترهای شبکه B می باشد

(b) شبکه B هیچ ارجحیتی نسبت به شبکه A ندارد با توجه به اینکه activation Function در هر دو شبکه خطی می باشد عمق بیشتر شبکه در آموزش ترجیح غیر خطی تأثیری ندارد و بی فایده است



$$(i) \delta_1 = \frac{\partial F(m)}{\partial z_1^2} = \frac{\partial F(m)}{\partial a_1^2} \times \frac{\partial a_1^2}{\partial z_1^2} = \omega_1^3 \times a_1^2 (1 - a_1^2)$$

$$(ii) \delta_2 = \frac{\partial F(m)}{\partial z^2} \begin{cases} \rightarrow \frac{\partial F(m)}{\partial z_1^2} = \omega_1^3 \times a_1^2 (1 - a_1^2) \\ \rightarrow \frac{\partial F(m)}{\partial z_2^2} = \omega_2^3 \times a_2^2 (1 - a_2^2) \end{cases}$$

$$(iii) \delta_3 = \frac{\partial F(m)}{\partial z^1} = \begin{cases} \rightarrow \frac{\partial F(m)}{\partial z_1^1} = \frac{\partial F(m)}{\partial a_1^2} \times \frac{\partial a_1^2}{\partial z_1^2} \times \frac{\partial z_1^2}{\partial a_1^1} \times \frac{\partial a_1^1}{\partial z_1^1} + \frac{\partial F(m)}{\partial a_2^2} \times \frac{\partial a_2^2}{\partial z_2^2} \times \frac{\partial z_2^2}{\partial a_1^1} \times \frac{\partial a_1^1}{\partial z_1^1} \\ = \omega_1^3 \times a_1^2 (1 - a_1^2) \times \omega_{11}^2 \times a_1^1 (1 - a_1^1) + \omega_2^3 \times a_2^2 (1 - a_2^2) \times \omega_{21}^2 \times a_1^1 (1 - a_1^1) \\ \rightarrow \frac{\partial F(m)}{\partial z_2^1} = \frac{\partial F(m)}{\partial a_2^2} \times \frac{\partial a_2^2}{\partial z_2^2} \times \frac{\partial z_2^2}{\partial a_1^1} \times \frac{\partial a_1^1}{\partial z_2^1} + \frac{\partial F(m)}{\partial a_1^2} \times \frac{\partial a_1^2}{\partial z_1^2} \times \frac{\partial z_1^2}{\partial a_2^1} \times \frac{\partial a_2^1}{\partial z_2^1} \\ = \omega_2^3 \times a_2^2 (1 - a_2^2) \times \omega_{22}^2 \times a_2^1 (1 - a_2^1) + \omega_1^3 \times a_1^2 (1 - a_1^2) \times \omega_{12}^2 \times a_2^1 (1 - a_2^1) \end{cases}$$

$$(iv) \delta_4 = \frac{\partial F(m)}{\partial w_{11}} = \frac{\partial F(m)}{\partial a_1^2} \times \frac{\partial a_1^2}{\partial z_1^2} \times \frac{\partial z_1^2}{\partial a_1^1} \times \frac{\partial a_1^1}{\partial z_1^1} \times \frac{\partial z_1^1}{\partial w_{11}} +$$

$$\frac{\partial F(m)}{\partial a_2^2} \times \frac{\partial a_2^2}{\partial z_2^2} \times \frac{\partial z_2^2}{\partial a_1^1} \times \frac{\partial a_1^1}{\partial z_1^1} \times \frac{\partial z_1^1}{\partial w_{11}} =$$

$$\omega_1^3 \times a_1^2 (1 - a_1^2) \times \omega_{11}^2 + a_1^1 (1 - a_1^1) \times \omega_{11} +$$

$$\omega_2^3 \times a_2^2 (1 - a_2^2) \times \omega_{21}^2 \times a_1^1 (1 - a_1^1) \times \omega_{11}$$