

① الف) در رابطه  $W = X^{-1}y$  ، ممکن است  $X$  یک ماتریس مربعی نباشد و در نتیجه معکوس پذیری نیست.

اگر از تابع هزینه  $J = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$  نسبت به  $W$  مشتق بگیریم و برای صفر قرار دهیم، پاسخ فرم بسته به صورت زیر در دست خواهد آمد:

$$W = (X^T X)^{-1} X^T y$$

ب) عدم وجود  $auto\ correlation$  که یکی از فرضیات سنجش رگرسیون خطی است به معنای زنی عدم وجود ارتباط (Correlation) بین داده‌ها و خطاهاست.

برای بررسی وجود یا عدم وجود این ارتباطی یک تست آماری به نام  $durbin\ watson$  وجود دارد.

← خطی که با کمک رگرسیون خطی به داده‌ها فیت می‌شود، دارای است  $y = 0.972x + 12.702$

← در این تست آماری مطابق جدول موجود برای  $\alpha = 0.05$  مقادیر بحرانی زیر را داریم:  $d_L = 0.1879$ ،  $d_U = 1.33$

← در این تست آماری  $d$  از رابطه زیر به دست می‌آید:

$$d = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}$$

جدول زیر را تشکیل می‌دهیم:

n	x	y	$\hat{y}$	error	error <sup>2</sup>
1	48.9	40.8	40.628	0.172	0.11383
2	50.4	42.5	42.108	0.392	0.15364
3	52.9	44.9	42.33	0.27	0.0729
4	55	46.1	44.38	-0.28	0.0784
5	56.8	47.4	48.13	-0.73	0.5329
6	58.8	49.1	50.09	-0.99	0.9801
7	61.2	51.5	52.43	-0.93	0.8649
8	62.5	53.5	53.7	-0.2	0.04
9	64.7	54.2	55.84	-0.64	0.4096
10	68	57.3	56.4	0.9	0.81

مجموع جذورات خطا دارای است با 47.6، 3

$$\varepsilon_2 - \varepsilon_1 = -0.1028 \quad \text{میان} \rightarrow 0.1002$$

$$\varepsilon_3 - \varepsilon_2 = -0.118 \rightarrow 0.1022$$

$$\varepsilon_4 - \varepsilon_3 = -0.100 \rightarrow 0.1202$$

$$\varepsilon_5 - \varepsilon_4 = 0.110 \rightarrow 0.1022$$

$$\varepsilon_6 - \varepsilon_5 = -0.106 \rightarrow 0.1212$$

$$\varepsilon_7 - \varepsilon_6 = 0.124 \rightarrow 0.1042$$

$$\varepsilon_8 - \varepsilon_7 = 0.102 \rightarrow 0.128$$

$$\varepsilon_9 - \varepsilon_8 = 0.106 \rightarrow 0.1212$$

$$\varepsilon_{10} - \varepsilon_9 = 0.18 \rightarrow 0.142$$

مجموع جنوریات استقلال خطاها برای آبی: 1.96

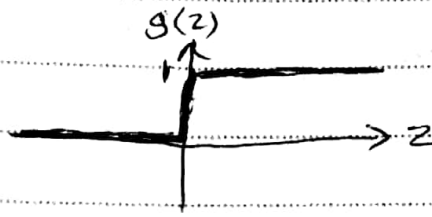
$$d = \frac{\sum_{t=2}^n (\varepsilon_t - \varepsilon_{t-1})^2}{\sum_{t=1}^n \varepsilon_t^2} = \frac{1.96}{2.472} = 0.08$$

$$d = 0.08 < d_L = 0.159 \Rightarrow \text{دارد positively autocorrelated error}$$

$$F-d = 2.25 \nless d_u = 1.32 \Rightarrow \text{نیست Negatively autocorrelated error}$$

## ۵) توابع فعال سازی (الف)

## Binary step •



← برای  $z \neq 0$  مشکل خود را برایشان را دارد زیرا در این نقطه

گرادیان دارد صفر است

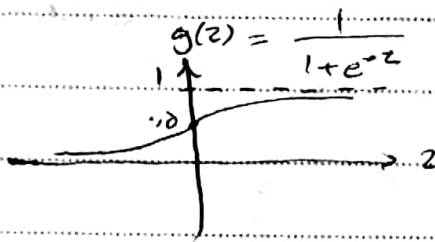
$$g(z) = U(z) = \begin{cases} 1 & ; z > 0 \\ 0 & ; z < 0 \end{cases} \quad g'(z) = U'(z) = \begin{cases} 0 & ; z > 0 \\ \delta(0) & ; z = 0 \\ 0 & ; z < 0 \end{cases}$$

← برای  $z \neq 0$  مشتق بی نهایت وی در  $z = 0$  مشتق ناپیوسته است گوییم مشتق نرم افزایشی گرفت و در این

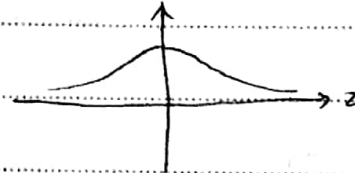
نقطه مشتق را برای سبب این مشتقات را در جیب قرار داد

← از نظر هزینه محاسباتی عملاً "کمترین هزینه ممکن" را دارد و مشکل از این لحاظ وجود ندارد

← Zero-centered نیست



$g'(z)$



## Sigmoid •

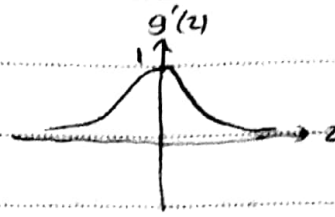
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}} \quad g'(z) = \sigma'(z) = \frac{1}{1 + e^{-z}} \left( \frac{e^{-z}}{1 + e^{-z}} \right) = \sigma(z) (1 - \sigma(z))$$

← برای  $z \rightarrow +\infty$  مشکل خود را برایشان وجود دارد زیرا در این به صفر میل می کند

← هم از این برای نقاط مشتق بی نهایت است

← از نظر هزینه محاسباتی دارای پیچیدگی است و استفاده از آن ممکن است باعث کند شدن فرایند یادگیری شود

← Zero-centered نیست



tanh

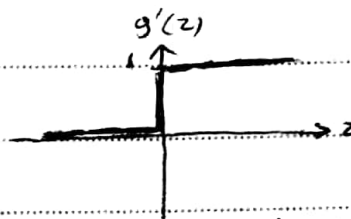
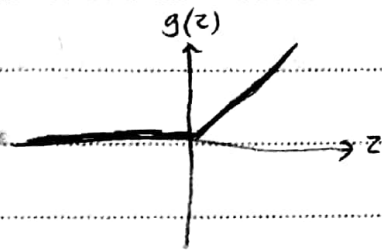
$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{2}{1 + e^{-2z}} - 1 = \text{tanh}(z) - 1 \quad g'(z) = 1 - (\tanh(z))^2$$

← برای  $z \rightarrow +\infty$  شکل مشتق گرایی وجود دارد زیرا گرایی به صفر میل می‌کند

← از آنجایی که مشتق می‌پذیرد است

← از نظر حسابی دارای پیچیدگی است و ممکن است باعث کند شدن روند یادگیری شود

← zero-centered است



Relu

$$g(z) = \text{Relu}(z) = \max(0, z)$$

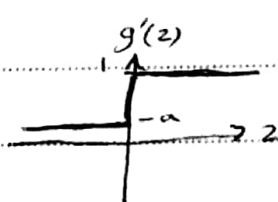
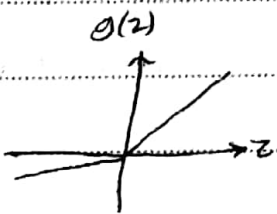
$$g'(z) = \begin{cases} 1 & z > 0 \\ 0 & z = 0 \\ 0 & z < 0 \end{cases}$$

← برای  $z < 0$  ممکن است گرایی وجود دارد زیرا گرایی به صفر میل می‌کند

← برای  $z \neq 0$  مشتق می‌پذیرد است. در  $z = 0$  مشتق ناپذیر است که می‌تواند منجر به مشکل در یادگیری شود

← از نظر حسابی بسیار ساده است

← zero-centered نیست



Leaky Relu

$$g(z) = \text{LRelu}(z) = \begin{cases} z & z > 0 \\ \alpha z & z < 0 \end{cases}$$

$$g'(z) = \begin{cases} 1 & z > 0 \\ \alpha & z = 0 \\ -\alpha & z < 0 \end{cases}$$

← شکل مشتق گرایی وجود ندارد

← از نظر حسابی بسیار ساده است

← zero-centered نیست



sigmoid

یا تابع فعال ساز Swish به صورت زیر تعریف شده است:  $g(x) = x \sigma(\beta x)$

$$\beta = 0 \Rightarrow g(x) = x/2, \quad \beta = 1 \Rightarrow g(x) = x \sigma(x), \quad \beta \rightarrow \infty \Rightarrow g(x) \rightarrow \text{Relu}(x)$$

مشتق این تابع به صورت زیر است:  $g'(x) = \sigma(\beta x) + \beta x \sigma(\beta x) (1 - \sigma(\beta x))$

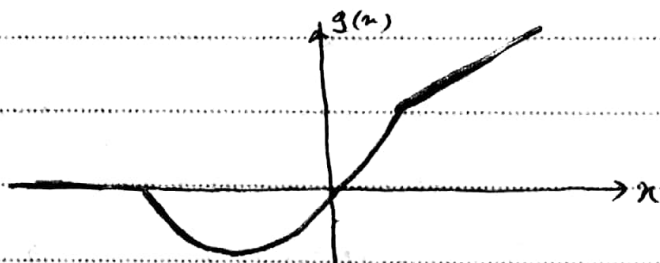
\* از آنجایی که وقتی  $\beta = 0$  آنگاه  $g(-x) = -g(x)$  و  $g(x)$  به ازای  $\beta = 0$  zero-centered است.

\* از روی مقدارهای  $x$  و  $\beta$  در یافت که وقتی  $x \rightarrow \infty$  مشکل همیستف کردن دیدن وجود دارد زیرا اگر  $x$  به منفی میل کند لازم به ذکر است که با افزایش  $\beta$  مشتق زودتری به منفی میل کند.

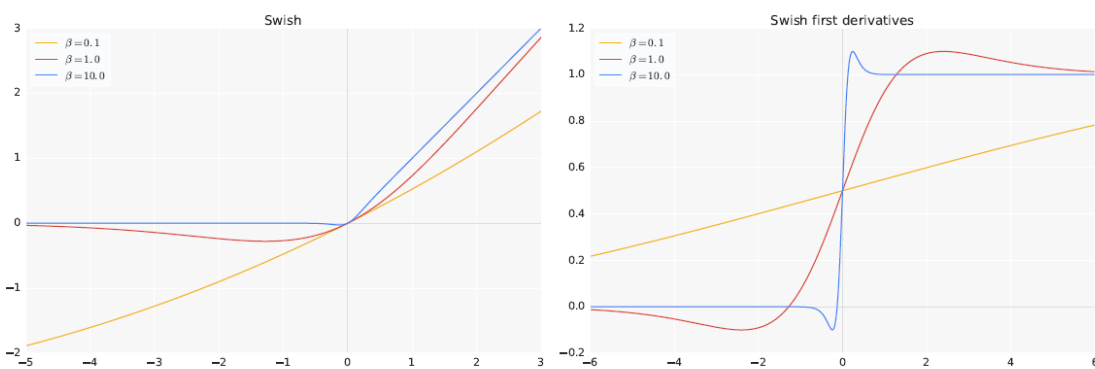
\* از نظر عملکرد Swish بر روی دیتاست های CIFAR-10, CIFAR-100 و IMAGENET در تمامی موارد عملکردی بهتر یا مشابه Relu داشته است اما در کارهای ترکیبیاتی بر روی دیتاست WMT2014 در برخی موارد Relu عملکرد بهتری داشته است.

تابع فعال ساز hard swish به صورت زیر تعریف شده است:  $g(x) = \frac{x}{4} \text{Relu}\sigma(x+3)$

$$\text{Relu}\sigma(x) = \begin{cases} 0 & ; x < 0 \\ x & ; 0 \leq x \leq 6 \\ 6 & ; x > 6 \end{cases}$$



\* از آنجایی که در رابطه تابع فعال ساز Swish از تابع sigmoid استفاده می شود، هزینه محاسباتی بیشتری نسبت به استفاده از تابع فعال ساز hard swish دارد. بنابراین جودیلان گفت که در دستگاه های که منبع محاسباتی کمتری دارند می توانند از hard swish استفاده کرد اما می بیند از hard swish استفاده کرد و با این حال که هزینه محاسباتی کمتری دارد، وقتی مشابه Swish دارد.



③ الف)  $RMSprop \rightarrow$  منکی

(۱) زیاد در سطح افتا زودتر همگرا می شود. نه اسیانات آن در یک جهت است. بای این می توان با برداشتن گام ها بلندتر زودتر به جواب رسید.

$Momentum \rightarrow$  آبی

از آنجای که در این روش گرایان های یکی نیز لحاظ می شوند. می توان گفت که از جهت میان گرایان ها استفاده می کند. بای این جهت به  $gd$  نوسانات کمتری دارد و زودتر همگرا می شود.

$Nesterov Momentum \rightarrow$  قرمز

زیاد نقطه در حدی که گرایان در آن جا به هم می رسد را پیش بینی می کند و برای سرعت افزایات کمتری نسبت به  $Momentum$  دارد. بای این زودتر از  $Momentum$  همگرا می شود.

$Gradient descent \rightarrow$  سبز

از آنجای که  $gd$  فقط از اس گرایان در نقطه فعلی جهت می گیرد و در جهت آن حرکت می کند. ممکن است دچار نوسانات زیادی شود.

(۲) برای این به روشی در بالا ذکر شده از سبب روش های ذکر شده می توان به نیاز داشت به چنین جابجیا را می که نیاز به تنظیم شدن دارند ذکر کرده در واقع تعداد جابجیا را بیشتر می شود.

(۳) روشی  $Adam$  از آنجای که ترکیب روش های  $Momentum$  و  $RMSprop$  است به روشی

$Momentum$  از برای روش  $RMSprop$  نیز به می آید. بای این می توان گفت که در روش  $Momentum$

مکن دچار مینیم های محلی شویم. از با روشی  $Adam$  تلاش بیشتری برای فرار از این مینیم های محلی

وجود دارد و به سمت  $Flat Minimum$  حرکت می کنیم.

در مورد روشی  $bias correction$  از آنجای که مقدار دین اولیه برای  $m$  و  $v$  داریم، به کمک تقسیم کردن و

$1 - \beta^t$  از این مقدار دین اولیه را کم می کنیم و حاصل آن را به  $\beta$  و  $\beta^2$  با این تقسیم می یابیم.

۳ ب) Multi-task learning زبان (اطلاعی شود که یک مدل شبکه عصبی می‌کند که بتواند چندین شبکه را به صورت همزمان یاد بگیرد. این تکنیک زمانی استفاده می‌شود که شبکه‌های که می‌خواهیم یاد گرفته شوند، بتوانند از low-level feature های مشترک بهره ببرند. در این صورت است که کارایی یک شبکه عصبی multi-task بهتر از چندین شبکه عصبی single task می‌باشد. از آنجایی که تمامی شبکه‌ها از این ویژگی‌های سطح پایین مشترک استفاده می‌کنند، می‌گوئیم که این ویژگی Shared parameter هستند.

دلیل اینکه این روش می‌تواند باعث افزایش generalization مدل شود این است که می‌توان از داده‌های بیشتری برای آموزش شبکه استفاده کرد. زیرا می‌توانیم از داده‌های استفاده کنیم که شاید برای شباهت یک شبکه باشد اما در نهایت از آنجایی که از ویژگی Shared parameter استفاده می‌کنیم، باعث می‌شود که عملکرد مدل روی دیگر شبکه‌ها نیز افزایش یابد. در داتس‌های توان از ویژگی‌های سطح پایین شباهت شبکه‌ها که مشترک هستند بهره‌مند شویم.

۴ ج)

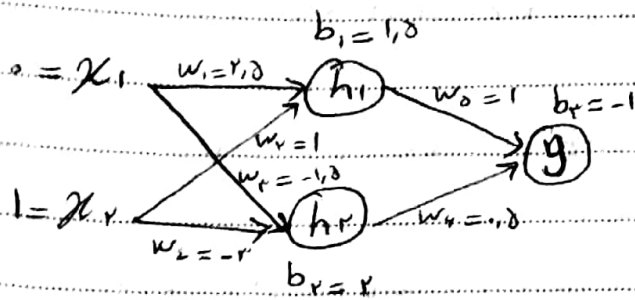
(۱) از آنجایی که در روشی drop out به صورت تصادفی در هر iteration داده‌های را حذف می‌کنیم می‌توان گفت که Complexity شبکه کاهش می‌یابد و در نتیجه مدل نیز کم‌کم ضعیف می‌شود. این یعنی که با افزایش عمق، پیچیدگی و قدرت تقسیم‌بندی مدل.

(۲) در ensemble learning چندین مدل با قدرت کمی کمتر آموزش می‌بینند و نتیجه نهایی براساس نظر اکثریت مدل‌ها است. در روشی drop out نیز به طور مشابه در هر iteration یک مدل کم‌قدرت را آموزش می‌دهیم. با حذف کردن تصادفی داده‌های (داده‌های پنهان) در نهایت می‌توان گفت میانگینی از این مدل‌های کم‌قدرت را استفاده می‌کنیم و در نتیجه نهایی نظر اکثریت مدل‌ها تأثیرگذار خواهد بود.

(۳) استفاده از drop out حتی در train به این صورت است که برای هر لایه یک احتمال نگهداری داده‌های آن لایه را در نظر می‌گیریم و به این صورت به یک شبکه خلاصه می‌رسیم. اما در زمان test به چه حاله که از drop out استفاده نکنیم و تمامی داده‌های هر لایه را نگه داریم و در زمان test می‌خواهیم از تمام قابلیت مدل و بارهای آن‌ها که شبکه یاد گرفته است استفاده کنیم.



## Back propagation (٢)



$$\alpha = 0.1, \quad y = 1 \quad (\text{target})$$

Forward:  $\text{net}_{h_1} = w_1 x_1 + w_3 x_2 + b_1 = 1.0 \times 0.5 + 1 \times 1 + 1.0 = 2.5$

$$\text{out}_{h_1} = \sigma(\text{net}_{h_1}) = \frac{1}{1 + e^{-2.5}} = 0.924$$

$$\text{net}_{h_2} = w_2 x_1 + w_4 x_2 + b_2 = -1.0 \times 0.5 - 1 \times 1 + 2 = -1$$

$$\text{out}_{h_2} = \sigma(\text{net}_{h_2}) = \frac{1}{1 + e^{+1}} = 0.269$$

$$\text{net}_y = w_5 \text{out}_{h_1} + w_6 \text{out}_{h_2} + b_3 = 0.8 \times 0.924 + (-1) \times 0.269 + (-1) = -0.09$$

$$\text{out}_y = \sigma(\text{net}_y) = \frac{1}{1 + e^{-0.09}} = 0.52$$

$$J = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log \text{out}_y^{(i)} + (1 - y^{(i)}) \log (1 - \text{out}_y^{(i)}) \right]$$

$$\xrightarrow{n=1} J = - (y \log \text{out}_y + (1 - y) \log (1 - \text{out}_y)) = 0.52$$

Backward:  $\frac{\partial J}{\partial w_5} = \frac{\partial J}{\partial \text{out}_y} \frac{\partial \text{out}_y}{\partial \text{net}_y} \frac{\partial \text{net}_y}{\partial w_5} = -1.961 \times 0.269 \times 0.924 = -0.489$

$$\frac{\partial J}{\partial \text{out}_y} = - \left( \frac{y}{\text{out}_y} - \frac{1-y}{1-\text{out}_y} \right) = - \left( \frac{1}{0.52} - \frac{1-1}{1-0.52} \right) = -1.961$$

$$\frac{\partial \text{out}_y}{\partial \text{net}_y} = \sigma'(\text{net}_y) = \text{out}_y (1 - \text{out}_y) = 0.52 (1 - 0.52) = 0.249$$

$$\frac{\partial \text{net}_y}{\partial w_5} = \text{out}_{h_1} = 0.924$$

$$\frac{\partial J}{\partial w_4} = \frac{\partial J}{\partial \text{out}_y} \frac{\partial \text{out}_y}{\partial \text{net}_y} \frac{\partial \text{net}_y}{\partial w_4} = -1.961 \times 0.249 \times 0.269 = -0.132$$

$$\frac{\partial \text{net}_y}{\partial w_4} = \text{out}_{h_2} = 0.269$$

$$\frac{\partial J}{\partial b_3} = \frac{\partial J}{\partial \text{out}_y} \frac{\partial \text{out}_y}{\partial \text{net}_y} \frac{\partial \text{net}_y}{\partial b_3} = -1.961 \times 0.249 \times 1 = -0.49$$

$$\frac{\partial \text{net}_y}{\partial b_3} = 1$$



$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial out_h} \frac{\partial out_h}{\partial net_h} \frac{\partial net_h}{\partial w_1} = -0.485 \times 0.404 \times 0 = 0$$

$$\frac{\partial J}{\partial out_h} = \frac{\partial J}{\partial out_y} \frac{\partial out_y}{\partial net_y} \frac{\partial net_y}{\partial out_h} = -0.991 \times 0.129 \times 1 = -0.128$$

$$\frac{\partial net_y}{\partial out_h} = w_3 = 1 \quad \frac{\partial out_h}{\partial net_h} = \delta_{net_h} (1 - \delta_{net_h}) = 0.485 (1 - 0.485) = 0.254$$

$$\frac{\partial net_h}{\partial w_1} = x_1 = 0$$

$$\frac{\partial J}{\partial w_r} = \frac{\partial J}{\partial out_h} \frac{\partial out_h}{\partial net_h} \frac{\partial net_h}{\partial w_r} = -0.128 \times 0.254 \times 1 = -0.032$$

$$\frac{\partial net_h}{\partial w_r} = x_r = 1$$

$$\frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial out_h} \frac{\partial out_h}{\partial net_h} \frac{\partial net_h}{\partial b_1} = -0.128 \times 0.254 \times 1 = -0.032$$

$$\frac{\partial net_h}{\partial b_1} = 1$$

$$\frac{\partial J}{\partial w_r} = \frac{\partial J}{\partial out_{hr}} \frac{\partial out_{hr}}{\partial net_{hr}} \frac{\partial net_{hr}}{\partial w_r} = -0.485 \times 0.194 \times 0 = 0$$

$$\frac{\partial J}{\partial out_{hr}} = \frac{\partial J}{\partial out_y} \frac{\partial out_y}{\partial net_y} \frac{\partial net_y}{\partial out_{hr}} = -0.991 \times 0.129 \times 0 = 0$$

$$\frac{\partial net_y}{\partial out_{hr}} = w_4 = 0 \quad \frac{\partial out_{hr}}{\partial net_{hr}} = \delta_{net_{hr}} (1 - \delta_{net_{hr}}) = 0.485 (1 - 0.485) = 0.254$$

$$\frac{\partial net_{hr}}{\partial w_r} = x_r = 0$$

$$\frac{\partial J}{\partial w_z} = \frac{\partial J}{\partial out_{hr}} \frac{\partial out_{hr}}{\partial net_{hr}} \frac{\partial net_{hr}}{\partial w_z} = -0.485 \times 0.254 \times 1 = -0.123$$

$$\frac{\partial net_{hr}}{\partial w_z} = x_z = 1$$

$$\frac{\partial J}{\partial b_r} = \frac{\partial J}{\partial out_{hr}} \frac{\partial out_{hr}}{\partial net_{hr}} \frac{\partial net_{hr}}{\partial b_r} = -0.485 \times 0.254 \times 1 = -0.123$$

در این وزن ها را باقی مانده از  $\text{gradient descent}$  به روز رسانی کرد:

$$\frac{\partial J}{\partial w_1} = 0 \quad \frac{\partial J}{\partial w_r} = -0.10259 \quad \frac{\partial J}{\partial w_r} = 0 \quad \frac{\partial J}{\partial w_z} = -0.10259$$

$$\frac{\partial J}{\partial w_0} = -0.444 \quad \frac{\partial J}{\partial w_4} = -0.112 \quad \frac{\partial J}{\partial b_1} = -0.10259 \quad \frac{\partial J}{\partial b_r} = -0.10259 \quad \frac{\partial J}{\partial b_r} = -0.10259$$

$$w_1^+ = w_1 - \alpha \frac{\partial J}{\partial w_1} = 2.0 - 0.1 \times 0 = 2.0$$

$$w_r^+ = w_r - \alpha \frac{\partial J}{\partial w_r} = 1 - 0.1 \times (-0.10259) = 1.010259$$

$$w_r^+ = w_r - \alpha \frac{\partial J}{\partial w_r} = -1.0 - 0.1 \times 0 = -1.0$$

$$w_z^+ = w_z - \alpha \frac{\partial J}{\partial w_z} = -2 - 0.1 \times (-0.10259) = -1.99897$$

$$w_0^+ = w_0 - \alpha \frac{\partial J}{\partial w_0} = 1 - 0.1 \times (-0.444) = 1.0444$$

$$w_4^+ = w_4 - \alpha \frac{\partial J}{\partial w_4} = -0.8 - 0.1 \times (-0.112) = -0.788$$

$$b_1^+ = b_1 - \alpha \frac{\partial J}{\partial b_1} = 1.0 - 0.1 \times (-0.10259) = 1.010259$$

$$b_r^+ = b_r - \alpha \frac{\partial J}{\partial b_r} = 2 - 0.1 \times (-0.10259) = 2.010259$$

$$b_r^+ = b_r - \alpha \frac{\partial J}{\partial b_r} = -1 - 0.1 \times (-0.10259) = -0.989741$$

$$\frac{\partial J}{\partial w_r} = \frac{\partial J}{\partial h_r} \frac{\partial h_r}{\partial z_r} \frac{\partial z_r}{\partial w_r} \quad (ب)$$

$$\frac{\partial J}{\partial h_{rj}} = \sum_{i=1}^{D_y} \frac{\partial J}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial h_{rj}}$$

$$J = - \sum_{i=1}^{D_y} y_i \lg \hat{y}_i \rightarrow \frac{\partial J}{\partial \hat{y}_i} = - \frac{y_i}{\hat{y}_i}$$

$$\left\{ \begin{aligned} \hat{y}_1 &= \frac{e^{h_{r1}}}{\sum_{i=1}^{D_y} e^{h_{ri}}} \rightarrow \frac{\partial \hat{y}_1}{\partial h_{r1}} = \frac{e^{h_{r1}} \sum e^{h_{ri}} - (e^{h_{r1}})^2}{(\sum e^{h_{ri}})^2} = \hat{y}_1 (1 - \hat{y}_1) \\ \hat{y}_r &= \frac{e^{h_{rr}}}{\sum e^{h_{ri}}} \rightarrow \frac{\partial \hat{y}_r}{\partial h_{r1}} = \frac{-e^{h_{r1}} e^{h_{rr}}}{(\sum e^{h_{ri}})^2} = -\hat{y}_1 \hat{y}_r \\ \hat{y}_r &= \frac{e^{h_{rr}}}{\sum e^{h_{ri}}} \rightarrow \frac{\partial \hat{y}_r}{\partial h_{r1}} = \frac{-e^{h_{r1}} e^{h_{rr}}}{(\sum e^{h_{ri}})^2} = -\hat{y}_1 \hat{y}_r \\ &\vdots \\ \hat{y}_{D_y} &= \frac{e^{h_{rD_y}}}{\sum e^{h_{ri}}} \rightarrow \frac{\partial \hat{y}_{D_y}}{\partial h_{r1}} = \frac{-e^{h_{r1}} e^{h_{rD_y}}}{(\sum e^{h_{ri}})^2} = -\hat{y}_1 \hat{y}_{D_y} \end{aligned} \right.$$

$$\Rightarrow \frac{\partial J}{\partial h_{r1}} = - \frac{y_1}{\hat{y}_1} \hat{y}_1 (1 - \hat{y}_1) + \frac{y_r}{\hat{y}_r} \hat{y}_1 \hat{y}_r + \frac{y_r}{\hat{y}_r} \hat{y}_1 \hat{y}_r + \dots + \frac{y_{D_y}}{\hat{y}_{D_y}} \hat{y}_1 \hat{y}_{D_y}$$

$$= \hat{y}_1 y_1 + \hat{y}_1 y_r + \hat{y}_1 y_r + \dots + \hat{y}_1 y_{D_y} - y_1 = \left( \sum_{i=1}^{D_y} y_i \right) \hat{y}_1 - y_1$$

$$\frac{\partial J}{\partial h_{rj}} = \left( \sum_{i=1}^{D_y} y_i \right) \hat{y}_j - y_j$$

به طریق مشابه می توان دریافت که:

$$\Rightarrow \frac{\partial J}{\partial h_r} = \begin{bmatrix} (\sum y_i) \hat{y}_1 - y_1 \\ (\sum y_i) \hat{y}_r - y_r \\ \vdots \\ (\sum y_i) \hat{y}_{D_y} - y_{D_y} \end{bmatrix} = \hat{y} \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}_{1 \times D_y} y - y$$



$$h_r = \text{Relu}(z_r) \Rightarrow \frac{\partial h_r}{\partial z_r} = \begin{cases} 1 & ; z_r > 0 \\ 0 & ; z_r < 0 \end{cases}$$

$$z_r = w_r h_i + b_r \Rightarrow \frac{\partial z_r}{\partial w_r} = h_i, \quad \frac{\partial z_r}{\partial b_r} = 1$$

$$\Rightarrow \frac{\partial J}{\partial w_r} = \begin{cases} (\hat{y} [1 \ 1 \dots 1]_{1 \times D_y} y - y) h_i^T & ; z_r > 0 \\ 0 & ; z_r < 0 \end{cases}$$

$$\Rightarrow \frac{\partial J}{\partial b_r} = \begin{cases} (\hat{y} [1 \ 1 \dots 1]_{1 \times D_y} y - y) & ; z_r > 0 \\ 0 & ; z_r < 0 \end{cases}$$

$$\frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial h_r} \frac{\partial h_r}{\partial z_r} \frac{\partial z_r}{\partial h_i} \frac{\partial h_i}{\partial z_i} \frac{\partial z_i}{\partial w_i}$$

$$\Rightarrow \frac{\partial J}{\partial w_i} = \begin{cases} w_r^T (\hat{y} [1 \ 1 \dots 1]_{1 \times D_y} y - y) x_i^T & ; z_i, z_r > 0 \\ 0 & ; \text{else} \end{cases}$$

$$\frac{\partial J}{\partial b_i} = \frac{\partial J}{\partial h_r} \frac{\partial h_r}{\partial z_r} \frac{\partial z_r}{\partial h_i} \frac{\partial h_i}{\partial z_i} \frac{\partial z_i}{\partial b_i}$$

$$\Rightarrow \frac{\partial J}{\partial b_i} = \begin{cases} w_r^T (\hat{y} [1 \ 1 \dots 1]_{1 \times D_y} y - y) & ; z_i, z_r > 0 \\ 0 & ; \text{else} \end{cases}$$