

الف) با افزایش batch size، تعداد قدم‌هایی که به سمت مینیمم

می‌گیریم لازم (از) کاهش می‌یابد. زیرا با افزایش batch size، نرخ یادگیری

ثابت، در $batch_size$ بالا قدم‌های دقیق‌تری به مینیمم

ما را می‌رساند. اما با افزایش batch size، تعداد قدم‌های لازم برای رسیدن به مینیمم

کمتر می‌شود. (از طرف دیگر، اگر batch size را افزایش دهیم، سرعت هرگز به دلیل کاهش

میانگین کاهش نمی‌یابد.)

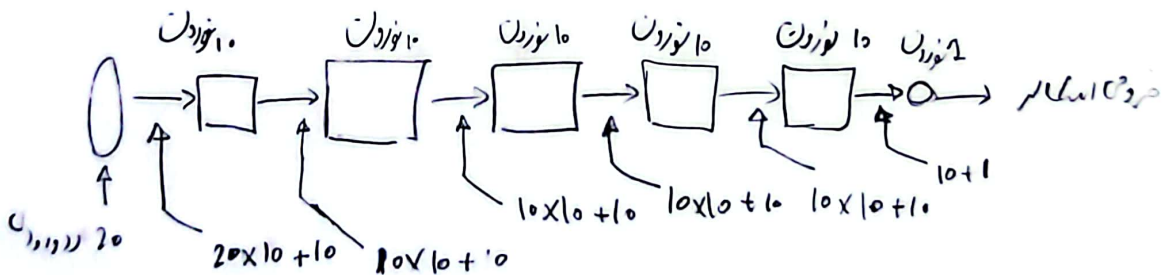
ب) با توجه به اینکه batch normalization میانگین و واریانس (ریک mini batch) انجام می‌شود، لذا تخمین دقیق از میانگین

واریانس واقعی حاصل می‌شود. همچنین این به این دلیل است که میانگین واریانس داده‌ها در هر بار (توسط دو پارامتر)

مهارسم به دست می‌آید. اگر یک نویز به داده اصلی اضافه می‌کنیم.

ج) در صورتی که ریک شبکه FL با تعویض ضرایب می‌شود وزن W (اولیم) به صورت مثبت بزرگ باشد، شبکه در ناحیه استیج

قرار می‌گیرد و وزن W به سمتی تغییر می‌کند. زیرا (رایج) می‌شود در مقادیر بسیار بزرگ مشتق آن نزدیک به صفر می‌شود.



$$2 \times 10 + 4 \times 10 \times 10 + 5 \times 10 + 10 + 1 = 661$$

این محاسبات شامل:
 - 2x10: وزن بین لایه‌های ورودی و پنهان اول
 - 4x10x10: وزن بین لایه‌های پنهان اول و دوم
 - 5x10: وزن بین لایه‌های پنهان دوم و سوم
 - 10: بایاس لایه پنهان دوم
 - 10x10+10: وزن بین لایه‌های پنهان سوم و چهارم
 - 4x10+10: وزن بین لایه‌های پنهان چهارم و پنجم
 - 5x10+10: وزن بین لایه‌های پنهان پنجم و خروجی
 - 10+1: بایاس لایه خروجی

$$\hat{y} = \text{softmax}(W_S X + b_S) \rightarrow \begin{cases} z_1 = W_{S1} X + b_{S1} \\ z_2 = W_{S2} X + b_{S2} \end{cases} \rightarrow \begin{cases} \hat{y}_1 = \text{softmax}(z_1) \\ \hat{y}_2 = \text{softmax}(z_2) \end{cases}$$

$$\hat{y}_1 \geq y_2 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \geq \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \rightarrow \frac{e^{z_1} - e^{z_2}}{e^{z_1} + e^{z_2}} \geq 0 \rightarrow \frac{1 - e^{z_2 - z_1}}{1 + e^{z_2 - z_1}}$$

$$\xrightarrow{z_1 + z_2 = t} \frac{1 - e^{-t}}{1 + e^{-t}} \geq 0 \rightarrow \frac{2 - (1 + e^{-t})}{1 + e^{-t}} \geq 0 \rightarrow \frac{2}{1 + e^{-t}} - 1 \geq 0 \rightarrow \sigma(t) \geq 0.5$$

$t = z_1 - z_2$

بنابراین اگر دو خروجی طبقه بند softmax را به صورت $z_1 - z_2$ تویین کنیم رویکرد مناسب طبقه بند اول خواهد بود پس ارزش تویین را پارامترهای تویین در توانایی طبقه بند (در حل مسائل پیچیده) توانسته است یعنی:

$$W_1 x + b_1 - W_2 x + b_2 = W_3 x + b_3$$

(2) الف) در این مقاله کمیته به سبب است یادگیری گفته می شود که هر طبقه بند را باید به ترکیب کرده و یک نتیجه با دقت بالاتر ارائه دهد. بنابراین در خروجی نتایج طبقه بندهای k کلاس یک لایه مخفی کمیته قرار می دهیم تا ترکیب خطی خروجی های k طبقه بند را محاسبه کرده و وزن های کمیته را یاد بگیرد.

$$y_{com}^k(x) = \sum_{n=1}^N W_{nk} y_{nk}(x)$$

$y_{nk}(x)$: خروجی های لایه کمیته (k کلاس)
 W_{nk} : وزن خروجی کلاس k ام طبقه بند n ام در لایه کمیته
 N : اعمال شعوبه
 $y_{com}^k(x)$: خروجی طبقه بند n ام در کلاس k

در هر طبقه بند از رویکرد متفاوتی در پیش برداشته شده است به گونه ای که خطای هر طبقه بند با بقیه نامیسته است.

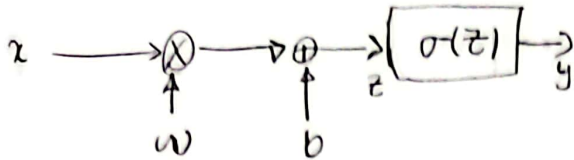
ب) در پیش برداشته داده ها (از روش های data alignment استفاده شده است:

(1) صریح distortion که وضوح تصویر را تغییر می دهند. (σ و γ)

(2) پرسش های نامساوی را دسته β

(3) داده کردن داده ها در محورهای افقی یا عمودی با پارامترهای γ و λ

$$z = wx + b, y = \sigma(z), L = \frac{1}{2} (y - t)^2, R = \frac{1}{2} w^2, L_{reg} = L + \lambda R \quad (3)$$



$$\frac{\partial L_{reg}}{\partial w} = \frac{\partial L}{\partial w} + \lambda \frac{\partial R}{\partial w} = \sigma'(wx+b)(1-\sigma'(wx+b))(y-t)x$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z} \times \frac{\partial z}{\partial w}$$

(1) (2) (3)

$$\textcircled{1}: \frac{\partial L}{\partial y} = \frac{1}{2} \times 2 \times (y-t) \quad \textcircled{2}: \frac{\partial y}{\partial z} = \frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}} \right) = \frac{e^{-z} + 1 - 1}{(1+e^{-z})^2} = \sigma(z) - \sigma^2(z)$$

$$\textcircled{3}: \frac{\partial z}{\partial w} = x$$

$$\Rightarrow \frac{\partial L}{\partial w} = (\sigma(wx+b))x(\sigma(wx+b) - \sigma^2(wx+b)) \times x = \sigma^2(wx+b)(1-\sigma(wx+b))(y-t)x$$

$$\frac{\partial L_{reg}}{\partial x} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z} \times \frac{\partial z}{\partial x} + \lambda \frac{\partial R}{\partial x} = \sigma^2(wx+b)(1-\sigma(wx+b))w(y-t)$$

(1) (2) (3) = w

$$\frac{\partial L_{reg}}{\partial b} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z} \times \frac{\partial z}{\partial b} + \lambda \frac{\partial R}{\partial b} = \sigma^2(wx+b)(1-\sigma(wx+b))(y-t)$$

(1) (2) = 1 = 0

$$\frac{\partial L_{reg}}{\partial \lambda} = R = \frac{1}{2} w^2, \quad \frac{\partial L_{reg}}{\partial t} = \frac{\partial L}{\partial t} = t - y = t - \sigma(wx+b)$$

ب. اگر چه نورون ها یکسان وزن ها شوند، همگی خروجی یکسانی تولید نکنند و شبکه نمی تواند مسائل متنوع را حل کند و ممکن است نام
 بهترین (یک local minimum یا saddle point گیر کند) سه دلیل در صورت تعدادی انجام میشود.

* در صورتیکه از وزن های با مقادیر اولیه بالا استفاده کنیم باعث افزایش حجم محاسبات میشود و ممکن است باعث
 انفجار gradient exploding شود.

$$\begin{aligned}
 & w^{(0)} = 0.1 \\
 & \mu = 0.1 \\
 & \left\{ \begin{aligned} & x^{(1)} = 2, x^{(2)} = 1, x^{(3)} = 2, x^{(4)} = 1 \\ & b^{(0)} = 0 \\ & t^{(0)} = 0, t^{(1)} = 1, t^{(2)} = 0, t^{(3)} = 1 \end{aligned} \right\} \rightarrow \left\{ \begin{aligned} & w^{(1)} = w^{(0)} - \mu \frac{\partial L_{reg}}{\partial w^{(0)}} = 0.1 - 0.1 \sigma^2(0.2) (1 - \sigma(0.2)) \times (\sigma(0.2) - 0) = 0.1372 \\ & b^{(1)} = b^{(0)} - \mu \frac{\partial L_{reg}}{\partial b^{(0)}} = 0.1 \sigma^2(0.2) (1 - \sigma(0.2)) \times (\sigma(0.2) - 0) = 0.1372 \\ & w^{(2)} = 0.1372 - 0.1 \sigma^2(0.1386) (1 - \sigma(0.1386)) \times (\sigma(0.1386) - 1) = 0.0014 \\ & b^{(2)} = 0.014 - 0.1 \sigma^2(0.1386) (1 - \sigma(0.1386)) \times (\sigma(0.1386) - 1) \times 0.1372 = -1.9 \times 10^{-4} \\ & w^{(3)} = 0.1811, b^{(3)} = 0.0018 \\ & w^{(4)} = 0.1879, b^{(4)} = -0.00024 \end{aligned} \right.
 \end{aligned}$$

(4) الف) الگوریتم ADAM به منظور یکارگیری نرخ یادگیری متغیر، ارتباطات خاص استفاده می شود:

$* g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1}) \rightarrow$ در این خطا گرایان تابع هزینه نسبت به ضرایب
 از زمان است محاسبه شده و در متغیر g_t ذخیره می شود.

$* m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t \rightarrow$ (رایت) مرحله متغیر momentum مرتبه اول
 از روی مقادیر قبل آن و گرایان فعلی با توجه به وزن β_1 آپدیت می شود.

$* v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \rightarrow$ در این مرحله متغیر مرتبه دوم که به آن velocity
 هم گفته می شود، از روی مقادیر قبلی خودش و مجذور
 متغیر گرایان مرحله فعلی آپدیت می شود (با وزن β_2)

$* \hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}, \hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$ این دو term در واقع با بایاس مربوط به مقادیر m_t و v_t
 را در مرحله اولیه آپدیت آن ها از بین می برد.

$* \theta_t \leftarrow \theta_{t-1} - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \rightarrow$ در نهایت با توجه به مقادیر بدست آمده از \hat{m}_t و \hat{v}_t وزن
 شبکه با نرخ ثابت α و ضریب متغیر $\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$ آپدیت می شوند
 (با وجود متغیر ϵ ، جلوگیری از صفر شدن مخرج کسر است).

(ب) با توجه به اینکه $\beta_1, \beta_2 = 0.9$ هستند مقادیر اولیه m_0 و v_0 برابر صفر (در نظر گرفته می شوند) لذا در مراحل اولیه آپدیت شدن m_t و v_t ، آن ها به سمت صفر با بایاس دارند و سرعت همگرای به دلیل کوچک شدن ضرایب گرایان با بایاس است. لذا عبارات \hat{m}_t و \hat{v}_t می توانند از این اتفاق جلوگیری کرده و با تقسیم $\frac{1}{1 - \beta_1^t}$ و $\frac{1}{1 - \beta_2^t}$ سرعت همگرای را در زمان ها اولیه افزایش دهند. رفته رفته با افزایش t $\hat{m}_t \approx m_t$ و $\hat{v}_t \approx v_t$ خواهد شد.

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \epsilon \nabla f(\underline{w}) \Rightarrow \underline{w}^{(t+1)} = \underline{w}^{(t)} - \epsilon \underline{H} \underline{w}^{(t)} = \underline{w}^{(t)} (\underline{I} - \epsilon \underline{H})$$

$$= \underline{H} \underline{w}^{(t)}$$

(1)

$$\underline{w}^{(1)} = \underline{w}^{(0)} (\underline{I} - \epsilon \underline{H}) \rightarrow \underline{w}^{(2)} = \underline{w}^{(1)} (\underline{I} - \epsilon \underline{H}) = \underline{w}^{(0)} (\underline{I} - \epsilon \underline{H})^2 \rightarrow$$

$$\rightarrow \underline{w}^{(t)} = \underline{w}^{(0)} (\underline{I} - \epsilon \underline{H})^t$$

ضرب از راست

(2)

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} (\underline{I} - \epsilon \underline{H}) \rightarrow \underline{w}^{(t+1)} = \underline{w}^{(t)} (\underline{I} - \epsilon \underline{Q} \underline{\Lambda} \underline{Q}^T) \xrightarrow{\uparrow} \underline{w}^{(t+1)} \underline{Q} = \underline{w}^{(t)} (\underline{Q} - \epsilon \underline{Q} \underline{Q} \underline{\Lambda} \underline{Q}^T \underline{Q})$$

$$= \underline{I}$$

ضرب از چپ

$$\underline{Q}^T \underline{w}^{(t+1)} \underline{Q} = \underline{w}^{(t)} (\underline{Q}^T \underline{Q} - \epsilon \underline{Q}^T \underline{Q} \underline{\Lambda} \underline{Q}^T \underline{Q}) \rightarrow \underline{Q}^T \underline{w}^{(t)} \underline{Q} = \underline{w}^{(0)} (\underline{I} - \epsilon \underline{\Lambda})^t$$

$$\underline{Q}^T \underline{Q} = \underline{I} \quad \underline{Q}^T \underline{Q} \underline{\Lambda} \underline{Q}^T \underline{Q} = \underline{\Lambda}$$

شرط همگرا: $|\underline{I} - \epsilon \underline{\Lambda}| < \underline{I}$

$$\rightarrow |1 - \epsilon \lambda_i| < 1 \quad \forall i=1, \dots, n \xrightarrow{\uparrow \text{تبدیل به عدد حقیقی}} |1 - \epsilon \lambda_{\max}| < 1 \rightarrow \boxed{0 < \epsilon < \frac{1}{\lambda_{\max}}}$$

$$\nabla f(\underline{w}) = \underline{H} \underline{w} \rightarrow \nabla^2 f(\underline{w}) = \underline{H} \rightarrow (\nabla^2 f(\underline{w}))^{-1} = \frac{1}{2} \underline{H}^{-1}$$

(1) رتقانیوت:

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \epsilon \left((\nabla^2 f(\underline{w}))^{-1} \nabla f(\underline{w}) \right) = \underline{w}^{(t)} - \epsilon \underline{w}^{(t)} = \underline{w}^{(t)} (1 - \epsilon)$$

$$= \frac{1}{2} \underline{H}^{-1} = 2 \underline{H} \underline{w}^{(t)}$$

$$\rightarrow \underline{w}^{(t)} = \underline{w}^{(0)} (1 - \epsilon)^t \xrightarrow{\epsilon=1, \text{ در صورتی که } \epsilon=1} \underline{w}^{(t)} = 0 \rightarrow \text{یک لحظه دیگر!}$$

$$\downarrow$$

$$\text{شرط همگرا: } 0 < \epsilon < 1$$

(3) در مسائلی که با تابع هزینه پیچیده روبرو هستیم روش آر (ن ماتریس همسین) که از مشتق مرتبه 2 درست می آید و معکوس گرفتن از آن ~ لحاظ محاسباتی بسیار پیچیده می شود. (از طرفی برای افزایش مقادیر ماتریس همسین با محدودیت ^{ماتریس} در بردار هستیم. همچنین در صورتی که در یک Saddle Point قرار بگیریم، معکوس ماتریس همسین مقادیر بزرگی به خود می گیرد و محاسبات را پیچیده تر می کند.

$$J_1 = \frac{1}{2} \left(y_d - \sum_{k=1}^n (\omega_k - \varepsilon_k) x_k \right)^2$$

(a) (b)

$$= \frac{1}{2} \left(y_d - \sum_{k=1}^n \omega_k x_k + \sum_{k=1}^n \varepsilon_k x_k \right)^2 = \underbrace{\frac{1}{2} \left(y_d - \sum_{k=1}^n \omega_k x_k \right)^2}_{= J_0(\omega)} - \left(\sum_{k=1}^n \varepsilon_k x_k \right) \left(y_d - \sum_{k=1}^n \omega_k x_k \right) + \left(\sum_{k=1}^n \varepsilon_k x_k \right)^2$$

$$E \left\{ \frac{\partial J_1}{\partial \omega_i} \right\} = E \left\{ \frac{\partial J_0(\omega)}{\partial \omega_i} \right\} - \underbrace{\frac{\partial}{\partial \omega_i} \left(\sum_{k=1}^n E \{ \varepsilon_k \} x_k \right) \left(y_d - \sum_{k=1}^n \omega_k x_k \right)}_{=0} + \frac{\partial}{\partial \omega_i} \left(\sum_{k=1}^n \sum_{k_2=1}^n E \{ \varepsilon_{k_1} \varepsilon_{k_2} \} x_{k_1} x_{k_2} \right)$$

$$= E \left\{ \frac{\partial J_0(\omega)}{\partial \omega_i} \right\} + \frac{\partial}{\partial \omega_i} \left(\sum_{k=1}^n \underbrace{E \{ \varepsilon_k^2 \} x_k^2}_{= \alpha \omega_k^2} + \sum_{k=1}^n \sum_{\substack{k_2=1 \\ k_2 \neq k_1}}^n \underbrace{E \{ \varepsilon_{k_1} \varepsilon_{k_2} \} x_{k_1} x_{k_2}}_{= E \{ \varepsilon_{k_1} \} E \{ \varepsilon_{k_2} \} = 0} \right)$$

$$\rightarrow E \left\{ \frac{\partial J_1}{\partial \omega_i} \right\} = E \left\{ \frac{\partial J_0(\omega)}{\partial \omega_i} \right\} + 2\alpha \omega_i x_i^2$$

ب) با توجه به افزودن شدن عبارت $2\alpha \omega_i x_i^2$ به مشتق متوجه می‌شویم که نشان‌دهنده اعمال شدن regularization به شبکه است. که توسط dropout انجام می‌شود.