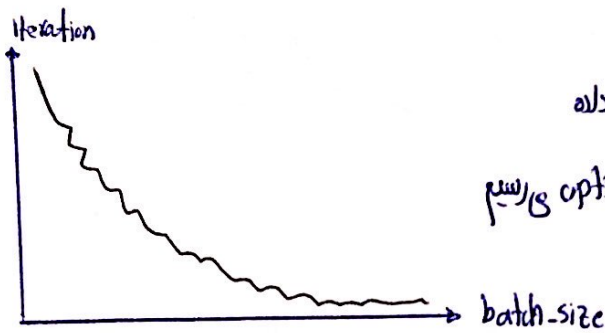


Problem 1:

a)



با افزایش تعداد داده هر بار ما تخمین بهتری نسبت به تویان کلی داده داریم و در نتیجه، گویان ما دقیق تر و سریع تر به نقطه optimal می رسیم (هنگامی سریع تر است)

b)

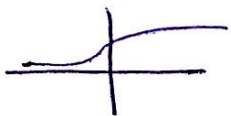
بله - در هنگام آموزش، در باب Batch normalizing، مقادیر تویان ترینی شده تا مدل ما عملکرد بهتری داشته باشد

$$\underline{x} = \{x_1, \dots, x_m\}^T \text{ input of mini batch} \rightarrow \text{normalize} \quad \hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\text{then} \rightarrow y_i = \gamma \hat{x}_i + \beta \quad \gamma, \beta \text{ hyperparameters}$$

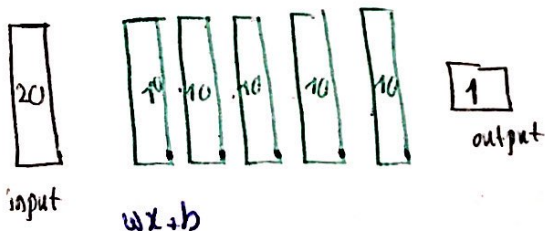
به دلیل این که تخمین میانگین (μ) و واریانس (σ^2) برای هر mini batch، متفاوت بوده و بعضی و در هر batch یکسان نیست مانند این است که تویان به داده ها اعمال کرده باشند

c) $z = \sigma(w^T x + b)$



فیر - زیرا تابع سیگموئید، انتشار افتد و قابلیت یادگیری خود را از دست می دهد و می تواند غایب خوبی برای مدلی مختلف قابل شود (Vanishing Gradient)

d)



$$10[20+1] + 10[10+1] \times 4 + 10[10+1] = 661$$

e)

خیر - برای Softmax، 2 کلاس، دقیقاً مشابه Logistic Regression ساده عملی کند

In Logistic Regression

$$\text{class 0: } \hat{y} = \frac{1}{1 + e^{-(\omega_0 x + b_0)}} \leq 0.5 \rightarrow e^{-(\omega_0 x + b_0)} \geq 1 \rightarrow -(\omega_0 x + b_0) \geq 0 \rightarrow \omega_0 x + b_0 \leq 0$$

$$\text{class 1: } \hat{y} = \frac{1}{1 + e^{-(\omega_1 x + b_1)}} \geq 0.5 \rightarrow -(\omega_1 x + b_1) \leq 0 \rightarrow \omega_1 x + b_1 \geq 0$$

①

Softmax

$$y_1 = \frac{e^{\omega_{s1}x + b_{s1}}}{e^{\omega_{s1}x + b_{s1}} + e^{\omega_{s2}x + b_{s2}}}$$

$$y_2 = \frac{e^{\omega_{s2}x + b_{s2}}}{e^{\omega_{s1}x + b_{s1}} + e^{\omega_{s2}x + b_{s2}}}$$

$$\text{class } 0 \text{ if } y_1 > y_2 \rightarrow \frac{e^{\omega_{s1}x + b_{s1}}}{e^{\omega_{s1}x + b_{s1}} + e^{\omega_{s2}x + b_{s2}}} > \frac{e^{\omega_{s2}x + b_{s2}}}{e^{\omega_{s1}x + b_{s1}} + e^{\omega_{s2}x + b_{s2}}} \rightarrow e^{x(\omega_{s1} - \omega_{s2}) + b_{s1} - b_{s2}} > 1$$

$$\rightarrow x(\omega_{s1} - \omega_{s2}) + b_{s1} - b_{s2} > 0 \rightarrow x(\omega_{s2} - \omega_{s1}) + b_{s2} - b_{s1} < 0$$

$$\begin{cases} \omega_l = \omega_{s2} - \omega_{s1} \\ b_l = b_{s2} - b_{s1} \end{cases}$$

Problem 2: a)

در روش کمینه، ما چندین مدل را روی دیتاستهای مختلف آموزش می دهیم و در نهایت، طبق روشی مانند رای اکثریت، یا میانگین یا median، خروجی مدل های مختلف، پیش بینی برای مشخصی کنیم. برای این که کمینه بهتر عمل کند خطای هر کدام از مدل ها نسبت به هم ناهمبسته باشد و در این مقاله برای رسیدن به این منظور (در pre process) عملی انجام شد. موازاً: در پردازش کمینه شده، ماتریس آموزش دهنده و ماتریس تست دهنده از ماتریس دیگری از overfitting می شود به دلیل این که خطای مدل ها مستقل اند

b)

1. Scaled data set to $[0,1]$

2. Elastic Deformations

3. Affine Transformations

4. Normalized the data set with that the width and height of the bounding box equals 20 pixels. Due to a large variation bounding box is normalized to 8 to 20

5. Deslanted Training Set. Generation of a deslanted training set by horizontally shearing the digit. The shearing magnitude is determined by $\tan(\alpha) \times d$, α is angle of the first principal component of pixel intensities with respect to the vertical axis, and d is the vertical distance from the image center

These preprocessing steps may contribute to preventing model errors.

1. Robustness of Deformation

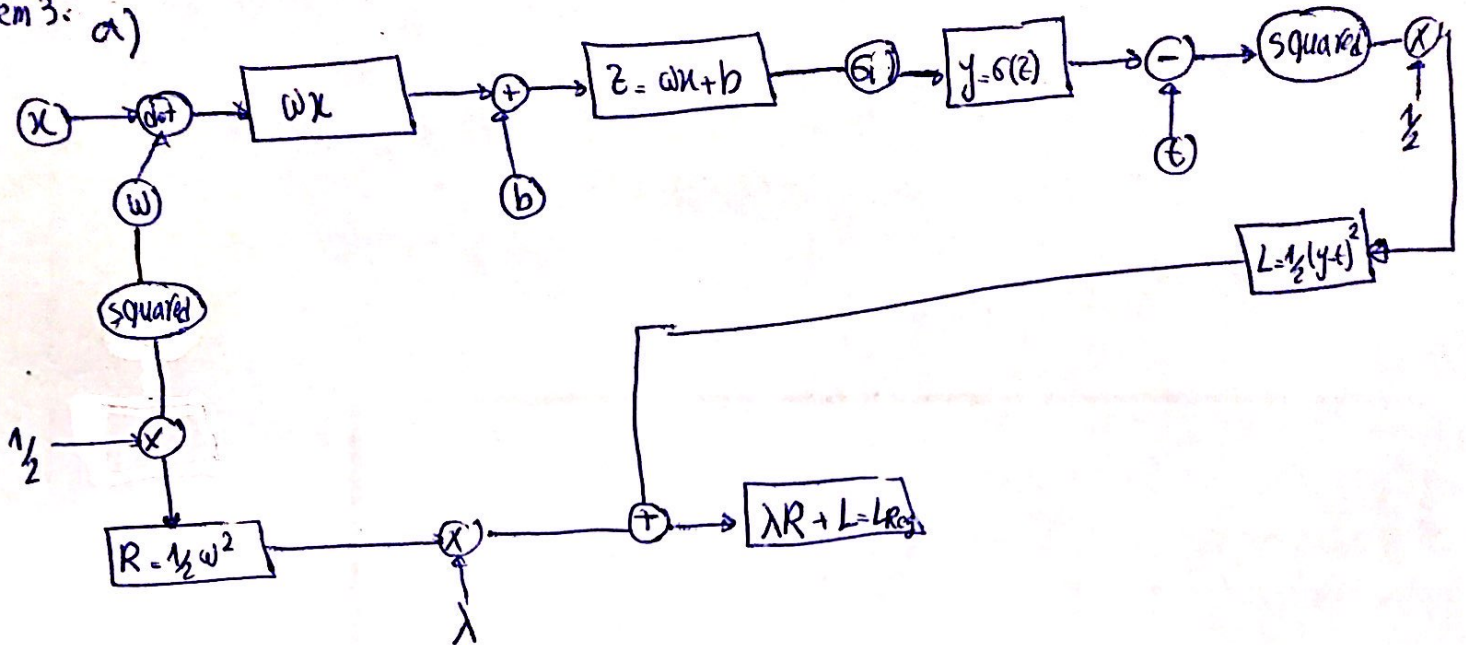
2. Affine Trans. for variability

3. Normalization for consistency

4. Deslanted Training set for improved Generalization

(2)

Problem 3: a)



$$\frac{\partial L_{reg}}{\partial \lambda} = R \quad ; \quad \frac{\partial L_{reg}}{\partial R} = \lambda \quad ; \quad \frac{\partial L_{reg}}{\partial L} = 1$$

$$\frac{\partial L_{reg}}{\partial t} = \frac{\partial L_{reg}}{\partial (y-t)} \cdot \frac{\partial (y-t)}{\partial t} = -1 \cdot (y-t) = t-y$$

$$\frac{\partial L_{reg}}{\partial y} = \frac{\partial L_{reg}}{\partial (y-t)} \cdot \frac{\partial (y-t)}{\partial y} = -1 \cdot (-1) = 1$$

$$\frac{\partial L_{reg}}{\partial z} = \frac{\partial L_{reg}}{\partial y} \cdot \frac{\partial y}{\partial z} = 1 \cdot (y-t) \cdot y(1-y)$$

$$\frac{\partial L_{reg}}{\partial w} = \lambda + \frac{\partial L_{reg}}{\partial z} \cdot \frac{\partial z}{\partial w} = \lambda + (y-t)y(1-y)x$$

$$\frac{\partial L_{reg}}{\partial t} = \frac{\partial L_{reg}}{\partial (y-t)} \cdot \frac{\partial (y-t)}{\partial t} = -1 \cdot (-1) = 1$$

$$\frac{\partial L_{reg}}{\partial x} = \frac{\partial L_{reg}}{\partial z} \cdot \frac{\partial z}{\partial x} = (y-t)y(1-y)w$$

$$z = \sigma(x) \Rightarrow z' = z(1-z)$$

اگر هدف منهای یکی باشند یا تصادفی نباشند، ممکن است در یادگیری نوزن ها دچار مشکل شوید به طوری که از خود رفتارهای متغییری بروز (b)

و نتایج در یادگیری وجود ندارد. و ممکن است در زمانی که مشکل آفرین هستند به دلیل افتار مشابه نوزن ها، مدل مختل شود.

اگر مقدار اولیه وزن ها زیاد باشد، ممکن است activation function در حالت اشباع رفته و دیده vanishing gradient به وجود می آید یا اگر کم باشد،

(3)

c) $\underline{x} = [1, 2]^T$, $\underline{w} = [0.1, 0.3]^T$, $b = 0.2$, $\lambda = 0.1$, $t = 1$, Learning Rate = 0.1

$$\tilde{y} = 0.1 \rightarrow 0.6 + 0.2 = 0.9 \rightarrow y = 0.7109 \rightarrow L = \frac{1}{2} (t - y)^2 = 0.0118$$

$$R = 0.05 \rightarrow L_{req} = 0.0468, \quad \frac{\partial L_{req}}{\partial \underline{w}} = [-0.0434, -0.088]^T, \quad \frac{\partial L_{req}}{\partial b} = -0.0539$$

$$\rightarrow \underline{w} = \underline{w} - \frac{\partial L_{req}}{\partial \underline{w}} \times 0.1 = [0.1049, 0.3089]^T$$

$$b = b - \frac{\partial L_{req}}{\partial b} \times 0.1 = 0.2059$$

Problem 4:

a) $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1}) \rightarrow$ نام Loss
گرایان تابع هدف (f) نسبت به وزن ها در
گام t یعنی آپدیت کردن وزن ها در گام t

• $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t \rightarrow$ (میانگین هدفی گرایان) بر روی momentum برای گام t

• $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \rightarrow$ (میانگین هدفی مربعات) بر روی سرعت برای گام t

• $\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$ تصحیح میانگین گرایان } برای تصحیح بایاس

• $\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$ تصحیح میانگین مربعات گرایان

• $\theta_t \leftarrow \theta_{t-1} - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \rightarrow$ آپدیت کردن وزن ها
با استفاده از گام t

• ϵ رابطی جلوگیری از تقسیم بر صفر می گذارد و α نرخ یادگیری است.

• β_1 و β_2 را معمولاً معادله های گذارده تا الگوریتم m_t و v_t را زیاد به تغییرات لحظه حساس بنماید.

• الگوریتم آدام باعث می شود که ما سریع تر همگرا شویم زیرا سرعت یادگیری ما به طور هوشمند تعیین شود.

b)

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$\beta_1 = 0.9$

$$m_0 = 0$$

$$m_1 = 0.1 g_1 + 0.9(0) = 0.1 g_1$$

$$\hat{m}_1 = g_1$$

$$m_2 = 0.9 [0.1 g_1] + 0.1 g_2 = 0.09 g_1 + 0.1 g_2$$

$$\hat{m}_2 = \frac{0.09 g_1 + 0.1 g_2}{0.1 + 0.09}$$

$$m_3 = 0.9 [0.1 g_2 + 0.09 g_1] + 0.1 g_3 = 0.09 g_2 + 0.081 g_1 + 0.1 g_3$$

در گام‌های نخست m_1 و m_2 بهی مثل m_t ها به صفر با بایس زیادی دارند چون تاثیر $m_0 = 0$ در این گام‌ها زیاد بوده به همین دلیل m_t ها را باید نرمالیز کرد تا تاثیر نقطه شروع در این گام‌ها کمتر بشود. مثلاً $m_1 = 0.1 g_1$ ولی $\hat{m}_1 = g_1$ است که تاثیرش در \hat{m}_1 بیشتر شده [تصمیم با بایس صورت گرفته است] و هر چه جلوتر برویم و t ها بزرگ تر شوند $1 - \beta_1^t$ علی همان یک نشد این کار مشکل با بایس در t ها کوچک را حل می کند.

Problem 5:

$$\omega^T H \omega = L(\omega)$$

$$a) \quad \nabla_{\omega} (\omega^T H \omega) = H \omega + H^T \omega = 2 H \omega \quad \text{و} \quad \nabla_{\omega} L(\omega) = 0 \rightarrow H \omega = 0$$

$$H = H^T$$

$$\text{update rule SGD} \rightarrow \omega_t = \omega_{t-1} - \alpha \nabla_{\omega_t} L \rightarrow \omega_{t-1} - \alpha (2 H \omega_t)$$

$$\rightarrow \omega_t = (I - 2\alpha H) \omega_{t-1} \rightarrow \omega_t = (I - 2\alpha Q \lambda Q^T) \omega_{t-1}$$

$$b) \quad \omega_1 = (I - 2\alpha Q \lambda Q^T) \omega_0 = Q (I - 2\alpha \lambda) Q^T \omega_0$$

$$\omega_2 = (I - 2\alpha Q \lambda Q^T) (I - 2\alpha Q \lambda Q^T) \omega_0 = (I - 2\alpha Q \lambda Q^T)^2 \omega_0$$

$$= Q (I - 2\alpha \lambda) Q^T Q (I - 2\alpha \lambda) Q^T \omega_0 = Q (I - 2\alpha \lambda)^2 \omega_0$$

$$\rightarrow \omega_t = Q (I - 2\alpha \lambda)^t \omega_0$$

(5)

$$c) |1 - 2\alpha\lambda_i| < 1 \rightarrow -1 < 1 - 2\alpha\lambda_i < 1 \rightarrow 0 < 2 - 2\alpha\lambda_i < 2$$

$$\rightarrow 2 > 2\alpha\lambda_i > 0 \rightarrow \frac{2}{\lambda_i} > \alpha > 0 \rightarrow \left(\frac{1}{\lambda_{max}} > \alpha \right)$$

کدام درجهی جاها 2α گرفته و این مقدار به صورت $\frac{2}{\lambda_{max}}$ درآید و همچنین الگوریتم آپدیت

$$\omega_t = \omega_{t-1} - \alpha H \omega_{t-1}$$

$$d) L(\omega) = \omega^T H \omega \quad \nabla_{\omega}(L(\omega)) = 2H\omega \quad \text{و} \quad \nabla_{\omega}^2(L(\omega)) = 2H$$

$$\text{الگوریتم نیوتن} \quad \omega_t = \omega_{t-1} - [\nabla_{\omega}^2(L)]^{-1} \nabla_{\omega}(L(\omega))$$

$$\omega_t = \omega_{t-1} - \frac{1}{2} H^{-1} 2H \omega_{t-1} \rightarrow \omega_t = \omega_{t-1}$$

- این الگوریتم در یک مرحله جواب می‌دهد

مشکل این الگوریتم با α نیز بداند سازی کرد

$$\omega_t = \omega_{t-1} - \alpha \frac{1}{2} H^{-1} 2H \omega_{t-1} \rightarrow \omega_t = (I - \alpha) \omega_{t-1}$$

$$\alpha < 1 \text{ شرط همگرایی}$$

الگوریتم نیوتن به دلیل استفاده از معکوس ماتریس هسیان (H^{-1}) بسیار پرهزینه بوده و برای مسدود با ابعاد بزرگ (e) مشکل بوده و همچنین بعضی مواقع ماتریس H ill condition بوده و حساب کردن معکوس آن مشکل سازی باشد پس روش نیوتن را برای آن استفاده کرد.

Problem 6:

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

$$E \left[\frac{\partial J_1}{\partial \omega_i} \right] = \frac{\partial E[J_1]}{\partial \omega_i}$$

$$E[J_1] = E \left[\frac{1}{2} \left(y_d - \underbrace{\sum_{k=1}^n (\omega_k x_k)}_{J_0} - \sum_{k=1}^n \delta_k x_k \right)^2 \right]$$

$$\Rightarrow E[J_1] = \frac{J_0^2}{2} + \frac{1}{2} E \left[\left(\sum_{k=1}^n \delta_k x_k \right)^2 \right] - J_0 E \left[\sum_{k=1}^n \delta_k x_k \right] \quad \begin{matrix} E[\delta_k] = 0 \\ \text{با فرض استقلال } \delta_k \end{matrix}$$

$$\begin{aligned} \Rightarrow E[J_1] &= \frac{J_0^2}{2} + \frac{1}{2} E \left[\sum_{k=1}^n (\delta_k x_k)^2 \right] + \frac{1}{2} E \left[\sum_{i=1}^n \sum_{j=1}^n \delta_i \delta_j x_i x_j \right] \\ &= \frac{J_0^2}{2} + \frac{1}{2} \sum_{k=1}^n \alpha \omega_k^2 x_k^2 \end{aligned}$$

$$\Rightarrow E[J_1] = \frac{1}{2} \left(y_d - \sum_{k=1}^n \omega_k x_k \right)^2 + \frac{1}{2} \alpha \sum_{k=1}^n \omega_k^2 x_k^2$$

$$\Rightarrow \frac{\partial E[J_1]}{\partial \omega_i} = -\frac{1}{2} x_i + \frac{1}{2} \alpha 2 \omega_i x_i^2 = \alpha \omega_i x_i^2 - 0.5 x_i$$

b)

باید توان به صورت رگولایزن نوشت

$$E[J_1] = \frac{J_0^2}{2} + \frac{\alpha}{2} \sum_{k=1}^n \omega_k^2 x_k^2$$

$$= J_0' + \lambda \sum_{k=1}^n (\omega_k')^2 \rightarrow \text{تنظیم رگولایزن شد}$$

$$\lambda = \frac{\alpha}{2}$$

$$\omega_k' = \omega_k x_k$$

(7)