

سوال ۳ الف

قسمت اول: بطور کلی رگولارایز کردن بایاس، مقدار آن را نزدیک به صفر می کند. حال رگرسیون لاجستیک را در نظر بگیرید. با صفر شدن بایاس، ابرصفحه مرز تصمیم گیری حتما از مبداء فضای ویژگی می گذرد که این واریانس مدل را کاهش داده و می تواند پیشبینی مدل را با خطای زیادی مواجه کند. حال یک شبکه عصبی دو لایه را در نظر بگیرید. در این حالت صفر کردن بایاس ها باعث می شود، برای نمونه هایی که مقدار بردار ویژگی شان نزدیک به صفر است، مدل خروجی نزدیک به صفر بدهد، زیرا خود وزن ها نیز رگولارایز می شوند. بطور کلی این اتفاق می تواند برای شبکه های عمیق تر نیز رخ دهد ولی با افزایش عمق شبکه تاثیر صفر شدن بایاس ها در کاهش واریانس مدل کمتر شده و عملا در شبکه های عمیق و کانولوشنی بایاس ها در واریانس مدل تاثیر ناچیزی دارند. هرچند رگولارایز کردن آن ها می تواند تاثیری منفی بگذارد اما بطور کلی در شبکه های عمیق رگولارایز کردن یا نکردن بایاس معمولا تاثیر چندانی در عملکرد شبکه ندارد. با توجه به این نکات و اینکه از تحلیل شبکه های سطحی می فهمیم که رگولارایز کردن بایاس می تواند تاثیر منفی داشته باشد (در صورتی که اگر آن را رگولارایز نکنیم، اگر با توجه به داده ها بهینه باشد، خودش به صفر می رسد)، لذا در شبکه های عصبی معمولا بایاس رگولارایز نمی شود.

مشتق رگولارایز L1 برای یک وزن $\text{sign}(w)$ می شود. رابطه بروزرسانی وزن به صورت زیر نوشته می شود:

$$w := w - \eta \text{sign}(w) - \frac{\partial \mathcal{L}}{\partial w}$$

قسمت دوم: فارغ از تاثیر گرادیان هزینه، در این رابطه اگر وزن مثبت باشد در مرحله بعدی وزن به اندازه نرخ یادگیری به سمت صفر سوق داده می شود و اگر منفی باشد دوباره به همان اندازه به سمت صفر می رود. این موضوع که این اتفاق با یک مقدار مشخص برای هر به روز رسانی اتفاق می افتد، باعث می شود تعداد زیادی از وزن ها صفر شده و مدل بدست آمده اصطلاحا به جواب تنکی برسد. البته شهودهای هندسی نیز برای این موضوع قابل ارائه است.

سوال ۳ ج

قسمت اول: بچ نرمالیزیشن سرعت پردازش هر بچ را کاهش می دهد ولی چون می تواند باعث شود مدل زودتر همگرا شود، در مجموع فرایند آموزش را سرعت می بخشد.

قسمت دوم: اگر ساینز بچ کوچک باشد، میانگین و واریانسی که در بچ نرمالیزیشن بدست می آید نویز بالایی خواهند داشت و نمایانگر توزیع واقعی نخواهند بود و در زمان آزمون هم می تواند پیشبینی شبکه را دچار مشکل کند.

سوال ۴ الف

خروجی شبکه ۰/۵ می شود و با یک بار به روز رسانی وزن ها می توانند مثبت یا منفی شوند و نمی توان بطور قطعی اظهار نظر کرد.

Q1. a

$$w_j := w_j - \alpha \frac{\partial \mathcal{L}}{\partial w_j}$$
$$b := b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{1}{2n} \sum_i \frac{\partial \mathcal{L}}{\partial w_j} (y^{(i)} - \hat{y}^{(i)})^2 = \frac{1}{2n} \sum_i 2x_j (\hat{y}^{(i)} - y^{(i)})$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{1}{2n} \sum_i 2(\hat{y}^{(i)} - y^{(i)})$$

Q1. b

$$\text{Set } X := \begin{bmatrix} 1 & -x^{(1)T} \\ \vdots & \vdots \\ 1 & -x^{(n)T} \end{bmatrix}_{n \times (d+1)}, \quad w := \begin{bmatrix} b \\ w \\ 1 \end{bmatrix}_{(d+1) \times 1}, \quad Y := \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}_{n \times 1}$$

$$\text{then } \mathcal{L}(w, b) = \frac{1}{2n} \|Y - Xw\|_2^2$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w} = \frac{-2}{2n} X^T (Y - Xw) = 0 \Rightarrow \boxed{w = (X^T X)^{-1} X^T Y}$$

Q1. c

Gradient Descent = $O(nm)$

closed form = $O(n^3 + n^2d)$

Q2.a

حالت اینکه در backprop قاعده زنجیره ای
گرادیان در عمق فرب می گوزد، اگر مقدار گرادیان
کوچک باشد (< 1) صورت نمایی گرادیان
نمایی به صفر میل می کند و باعث می شود وزن
مربوطه آپدیت نشود. \leftarrow Vanishing grad

$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \Rightarrow \sigma'(x) < 1$
لذا ضرب گرادیان خروجی را سگموشی می تواند مختبر به این مشکل شود.

$$\text{Relu}'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

لذا برای ReLU این مشکل پیش می آید چون برای (x_1, x_2, \dots)
نمایی گرادیان را کوچه می کند. (x_1, x_2, \dots) می شود و صورت

Q2.c

مقدار بزرگ وزن می باشد مقدار ورودی به
تابع فعال ساز سگموشی خروجی بزرگ باشد
که در نتیجه گرادیان بزرگتری کوچه شده و در نهایت
باعث مشکل Vanishing grad می شود.

Q3.6

$$\mathcal{L}'(\omega) = \frac{1}{2n} \sum_i \left(y^{(i)} - (x^{(i)} + \delta^{(i)})^T \omega \right)^2$$

where $\delta^{(i)} \sim \mathcal{N}(0, \sigma^2 I)$.

$$\begin{aligned} &= \frac{1}{2n} \sum_i \left(y^{(i)} - x^{(i)T} \omega - \delta^{(i)T} \omega \right)^2 \\ &= \frac{1}{2n} \sum_i \left(y^{(i)} - x^{(i)T} \omega \right)^2 - 2 \left(y^{(i)} - x^{(i)T} \omega \right) \delta^{(i)T} \omega \\ &\quad + \left(\delta^{(i)T} \omega \right)^2 \end{aligned}$$

$$= \underbrace{\frac{1}{2n} \sum_i \left(y^{(i)} - x^{(i)T} \omega \right)^2}_{\mathcal{L}(\omega)} + \frac{1}{2n} \left[\sum_i -2 \left(y^{(i)} - x^{(i)T} \omega \right) \delta^{(i)T} \omega + \sum_i \left(\delta^{(i)T} \omega \right)^2 \right]$$

$$\Rightarrow \mathbb{E}_{\delta \sim \mathcal{N}} [\mathcal{L}'(\omega)] = \mathbb{E}_{\delta \sim \mathcal{N}} [\mathcal{L}(\omega)] + \frac{1}{2n} \sum_i \mathbb{E}_{\delta \sim \mathcal{N}} [\dots]$$

we have that $\begin{cases} \mathbb{E}_{\delta \sim \mathcal{N}} [-2(y^{(i)} - x^{(i)T} \omega) \delta^{(i)T} \omega] \\ \quad = -2(y^{(i)} - x^{(i)T} \omega) \underbrace{\mathbb{E}_{\delta \sim \mathcal{N}} [\delta^{(i)T} \omega]}_{=0} \\ \mathbb{E}_{\delta \sim \mathcal{N}} [(\delta^{(i)T} \omega)^2] = \sigma^2 \|\omega\|_2^2 \end{cases}$

$$\Rightarrow \boxed{\mathbb{E}_{\delta \sim \mathcal{N}} [\mathcal{L}'(\omega)] = \mathcal{L}(\omega) + \sigma^2 \|\omega\|_2^2}$$

Q3.d

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla_{\theta} J(\theta_t)$$

$$\Rightarrow m_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \nabla_{\theta} J(\theta_i)$$

$$\Rightarrow E[m_t] = E\left[(1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \nabla_{\theta} J(\theta_i)\right]$$

$$= (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} E[\nabla_{\theta} J(\theta_i)]$$

Since $\stackrel{iid}{=} (1 - \beta_1) E[\nabla_{\theta} J(\theta_t)] \sum_{i=1}^t \beta_1^{t-i}$

$$= E[\nabla_{\theta} J(\theta_t)] (1 - \beta_1) \frac{\beta_1^{t-1}}{\beta_1 - 1}$$

$$= E[\nabla_{\theta} J(\theta_t)] (1 - \beta_1^t)$$

لذا به ازای مقدار β_1 نزدیک به 1 مقدار $1 - \beta_1^t$ برای t بزرگ

اول به صفر میل می کند (تقریباً برابر است) و در نتیجه مقدار اخیر m_t تقریباً

تقریباً برابر می شود لذا ما m_t از $\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$ استفاده می شود که:

$$E[\hat{m}_t] = E[\nabla_{\theta} J(\theta_t)] \quad \text{که همان چیزی است که می خواهیم}$$

Q4.6

این شبیه word2vec است که بردار آن را می‌خواهیم پس.
فقط باید که $W \in \mathbb{R}^{d \times H}$ و $W' \in \mathbb{R}^{H \times d}$ و یکی بردار و دیگری
است که فرض می‌کنیم m آن یکی است.

داریم:

$$\frac{\partial \mathcal{L}}{\partial u_j} = \hat{y}_j - \mathbb{I}\{j=m\} = 0$$

$$\frac{\partial \mathcal{L}}{\partial W'_{ij}} = \sum_{k=1}^d \frac{\partial \mathcal{L}}{\partial u_k} \frac{\partial u_k}{\partial W'_{ij}} = \frac{\partial \mathcal{L}}{\partial u_j} \frac{\partial u_j}{\partial W'_{ij}}$$

چون W'_{ij} فقط در u_j حضور دارد.

$$\frac{\partial u_j}{\partial W'_{ij}} = \sum_{k=1}^d W_{ik} x_k = 0 \rightarrow \frac{\partial \mathcal{L}}{\partial W'_{ij}} = 0$$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}} = \sum_{k=1}^d \frac{\partial \mathcal{L}}{\partial u_k} \frac{\partial u_k}{\partial W_{ij}} = \sum_{k=1}^d (\hat{y}_k - \mathbb{I}\{k=m\}) \frac{\partial u_k}{\partial W_{ij}}$$

$$\frac{\partial u_k}{\partial W_{ij}} = W'_{jk} x_i = 0 \rightarrow \frac{\partial \mathcal{L}}{\partial W_{ij}} = 0$$

پس از آن این است که قابل می‌باشد.

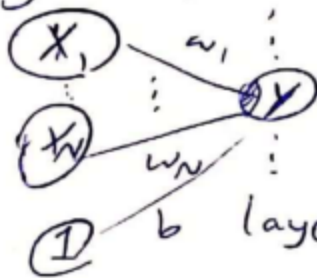
Xavier Initialization

if the weights are too small, then the variance of the input signal starts diminishing as it passes through layers of the network. this finally becomes so small that won't be useful. (causing Vanishing Gradients).

Now if the weights are too large, then the variance of the input data tends to rapidly increase with each passing layer; finally becoming too large to be useful anymore. (causing exploding gradients).

So the main idea is that we want the variance of input data not to change when it passes through the network layers.

layer [L]



$$y = g(z)$$

$$z = w_1 x_1 + \dots + w_N x_N + b$$

we want the $\text{var}(z) = \text{var}(x_i)$ for all $i \in \{1, \dots, N\}$.

$$\text{Var}(z) = \text{Var}(w_1 x_1) + \dots + \text{Var}(w_N x_N) + \text{Var}(b)$$

assume $w_i \sim \text{Gaussian}(0, \sigma)$ for all $i \in \{1, \dots, N\}$

also assume that $E[x_i] = 0$

then

$$\text{Var}(z) = N \text{Var}(w_k) \text{Var}(x_k) \quad \text{for any } k.$$

(Layer [L])

$$\Rightarrow \text{Var}(w_k) = \frac{1}{N} \quad \text{for all } k \in \{1, \dots, N\}$$

↑
(number of input layer units)

Xavier Initialization: Initialize weights mean=0 & Var = $\frac{1}{N}$ from Gaussian of

So in Xavier Initialization, we will initialize like so:

$$W^{[L]} \sim \text{Gaussian}\left(0, \frac{1}{N^{[L]}}\right)$$

← $N^{[L]}$
units at layer [L-1]