سؤال 1-

در ابتدا برای اثبات محدب بودن :

$$E_{w,b} = - \sum_{n} y_n \, ln \, \hat{y}(x_n) + (1-y_n) \, ln(1-\hat{y}(x_n))$$

با فرض اینکه $\hat{y}_n$ لیبل پیش‌بینی‌شده

$$\hat{y}_n = \frac{1}{1+e^{-(w^T x_n + b)}} \longrightarrow \frac{\partial \hat{y}_n}{\partial w} = \hat{y}_n(1-\hat{y}_n)x_n$$

$$\longrightarrow \frac{\partial \hat{y}_n}{\partial b} = \hat{y}_n(1-\hat{y}_n)$$

$$\frac{\partial E_{w,b}}{\partial w} = \sum_{n} \left( \frac{1-y_n}{1-\hat{y}_n} - \frac{y_n}{\hat{y}_n} \right) \frac{\partial \hat{y}_n}{\partial w} = \sum_{n} (\hat{y}_n - y_n)x_n$$

$$\frac{\partial^2 E_{w,b}}{\partial w \, \partial w^T} = \sum_{n} \frac{\partial}{\partial w} \left( (\hat{y}_n - y_n)x_n \right) = \underbrace{\sum_{n} \hat{y}_n (1-\hat{y}_n) x_n x_n^T}_{\substack{\text{هیشن مثبت} \qquad \text{PSD}}}$$

بنابراین ماتریس هسین PSD هست و در نتیجه $E_{w,b}$ محدب هست.

$$\frac{\partial E_{w,b}}{\partial b} = \sum_{n} \left( \frac{1-y_n}{1-\hat{y}_n} - \frac{y_n}{\hat{y}_n} \right) \frac{\partial \hat{y}_n}{\partial b} = \hat{y}_n - y_n$$

ادامه سؤال 1-

$$\frac{\partial^2 E_{w,b}}{\partial b^2} = \sum_n \frac{\partial}{\partial b} \{\hat{y}_n - y_n\} = \underbrace{\sum_n \hat{y}_n (1 - \hat{y}_n)}_{\text{همیشه مثبت}}$$

بنابراین ماتریس هسین PSD هست و تابع محدب هست.

حال چون می‌دانیم که $E_{w,b}$ یک مینیم گلوبال دارد بدین شکل می‌توانیم $b$ و $w$ را در آن نقطه بدست آوریم:

$$w_{n+1} = w_n - \eta_1 \sum_n (\hat{y}_n - y_n) x_n$$

$$b_{n+1} = b_n - \eta_2 \sum_n (\hat{y}_n - y_n)$$

---

سؤال 2-

الف) Covariate Shift به پدیده‌ای می‌گویند که در آن توزیع داده‌های ورودی یک شبکه عصبی در هنگام آموزش تغییر می‌کند و باعث سخت‌تر شدن یادگیری می‌شود. در شبکه‌های عصبی عمدی هر لایه ورودی لایه بعد است و زمانی که وزن‌های هر لایه آپدیت می‌شوند، توزیع ورودی به هر یک از لایه‌ها یعنی تغییر می‌کند. این پدیده سرعت آموزش را کاهش می‌دهد چون شبکه باید هر بار خودش را با توزیع متغیر داده‌ها مطابقت کند.

ادله سؤال 2 الف )

BN برای هر مینی بچ میانگین و واریانس داده‌های ورودی به لایه را
حساب می‌کند و آن را نرمالایز می‌کند که میانگین صفر و واریانس
واحد داشته باشند. به این ترتیب سعی می‌کند که توزیع داده‌ها بهم نریزد.
همچنین بعد از آنکه به میانگین صفر و واریانس واحد رسیدند، BN پارامترهای
قابل آموزشی برای شیفت و اسکیل در نظر می‌گیرد تا اگر میانگین
صفر و واریانس واحد مناسب نبود مدل بتواند آن را تغییر بدهد.

ب )

همین کاهش Covariate shift می‌تواند به جنرالیزیشن شبکه کمک کند.
BN اثرات مشابه با بقیه روش‌های رگولاریزیشن مثل L2 و dropout
دارد و برای همین می‌توانیم آن‌ها را از شبکه حذف کنیم تا مدل ساده‌تر
شود. همچنین می‌توانیم از لرنینگ ریت بالاتری استفاده کنیم.
در نهایت می‌توان گفت که BN مسیر آپتیمیزیشن را نرم‌تر
می‌کند که این می‌تواند به جنرالایز شدن مدل کمک کند.

پ )

$$\mu = \frac{1}{n} \sum_{j=1}^{n} x_j$$

$$\hat{x_i} = x_i - \mu = x_i - \frac{1}{n} \sum_{j=1}^{n} x_j$$

$$L = \sum_{j=1}^{n} L_j$$

$$\longrightarrow \frac{\partial \hat{x}_i}{\partial x_j} = \begin{cases} 1 - \frac{1}{n}, & i = j \\ -\frac{1}{n}, & i \neq j \end{cases} \quad , \quad \begin{array}{l} y_i = \gamma \hat{x}_i + \beta \\ \\ \hookrightarrow \frac{\partial L}{\partial \hat{x}_i} = \gamma \frac{\partial L}{\partial y_i} \end{array}$$

$$\frac{\partial L}{\partial x_i} = \sum_{j=1}^{n} \frac{\partial L}{\partial \hat{x}_j} \times \frac{\partial \hat{x}_j}{\partial x_i} = -\frac{1}{n} \left( \sum_{j=1}^{n} \frac{\partial L}{\partial \hat{x}_j} \right) + \frac{\partial L}{\partial \hat{x}_j}$$

$$= \gamma \left( \frac{\partial L}{\partial y_i} - \frac{1}{n} \sum_{j=1}^{n} \frac{\partial L}{\partial y_j} \right)$$

نکته)

$$n = 1 : \qquad \mu = x_1 \longrightarrow \hat{x}_1 = x_1 - x_1 = 0 \longrightarrow \frac{\partial y_1}{\partial \hat{x}_1} = 0$$

$$\longrightarrow \frac{\partial L}{\partial x_1} = 0$$

$$n \longrightarrow \infty :$$

$\frac{\partial \mu}{\partial x_i} = 0 \longleftarrow$ در این حالت تأثیرهم $x_i$ بر روی $\mu$ بسیار اندک هست

$$\longrightarrow \frac{\partial \hat{x}_i}{\partial x_j} \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases} \longrightarrow \frac{\partial L}{\partial x_i} = \gamma \frac{\partial L}{\partial y_i}$$

٭ نتیجه گرفته می شود که اگر n خیلی بزرگ شود اثر نرمال; پیوند اندین به ١(١١١/.

الف)

$$\hat{y}_k = \text{softmax}(z^{(2)}) = \frac{e^{z_k^{(2)}}}{\sum\limits_{j=1}^{K} e^{z_j^{(2)}}}$$

$$\frac{\partial \hat{y}_k}{\partial z_i^{(2)}} = \frac{\frac{\partial e^{z_k^{(2)}}}{\partial z_i^{(2)}} \sum\limits_{j=1}^{K} e^{z_j^{(2)}} - \left(\sum\limits_{j=1}^{K} \frac{\partial e^{z_j^{(2)}}}{\partial z_i^{(2)}}\right) \times e^{z_k^{(2)}}}{\left(\sum\limits_{j=1}^{K} e^{z_j^{(2)}}\right)^2}$$

$$i = k: \qquad \frac{\partial \hat{y}_k}{\partial z_k^{(2)}} = \frac{e^{z_k^{(2)}} \sum\limits_{j=1}^{K} e^{z_j^{(2)}} - \left(e^{z_k^{(2)}}\right)^2}{\left(\sum\limits_{j=1}^{K} e^{z_j^{(2)}}\right)^2} = \hat{y}_k(1 - \hat{y}_k)$$

$$i \neq k: \qquad \frac{\partial \hat{y}_k}{\partial z_i^{(2)}} = \frac{-e^{z_k^{(2)}} e^{z_i^{(2)}}}{\left(\sum\limits_{j=1}^{K} e^{z_j^{(2)}}\right)^2} = -\hat{y}_k \hat{y}_i$$

ب)

$$L = -\sum_{i=1}^{K} y_i \log(\hat{y}_i) = -\log(\hat{y}_k)$$

$$\longrightarrow \frac{\partial L}{\partial z_i^{(2)}} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i^{(2)}} \qquad\qquad \frac{\partial L}{\partial z_i^{(2)}} = \hat{y}_i - y_i$$

$$\frac{\partial L}{\partial \hat{y}_i} = \begin{cases} \frac{-1}{\hat{y}_k}, & i = k \\ 0, & i \neq k \end{cases} \quad \longrightarrow \quad \frac{\partial L}{\partial z_i^{(2)}} = \begin{cases} \hat{y}_k - 1 & i = k \\ \hat{y}_i & i \neq k \end{cases}$$

$$\frac{\partial L}{\partial z^{(2)}} = \hat{y} - y$$

$$\frac{\partial z^{(2)}}{\partial a^{(1)}} = (w^{(2)})^T \qquad , \qquad \frac{\partial a^{(1)}}{\partial \hat{a}^{(1)}} = \begin{cases} 1, & P = 0.8 \\ 0, & P = 0.2 \end{cases}$$

$$\frac{\partial \hat{a}^{(1)}}{\partial z_i^{(1)}} = \begin{cases} 1, & z_i^{(1)} > 0 \\ 0.01, & z_i^{(1)} \leq 0 \end{cases} \qquad , \qquad \frac{\partial z^{(1)}}{\partial w^{(1)}} = x^T$$

$$\longrightarrow \frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial \hat{a}^{(1)}} \frac{\partial \hat{a}^{(1)}}{\partial z_i^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(1)}}$$

---

سوال -4

$$\nabla y = \begin{bmatrix} \frac{\partial y}{\partial z} \\ \frac{\partial y}{\partial v} \\ \frac{\partial y}{\partial u} \end{bmatrix} \longrightarrow J(\nabla y) = \begin{bmatrix} \frac{\partial^2 y}{\partial u^2} & \frac{\partial^2 y}{\partial u \partial v} & \frac{\partial^2 y}{\partial u \partial z} \\ \frac{\partial^2 y}{\partial v \partial u} & \frac{\partial^2 y}{\partial v^2} & \frac{\partial^2 y}{\partial v \partial z} \\ \frac{\partial^2 y}{\partial z \partial u} & \frac{\partial^2 y}{\partial z \partial v} & \frac{\partial^2 y}{\partial z^2} \end{bmatrix}$$

$$H = \begin{bmatrix} \frac{\partial^2 y}{\partial u^2} & \frac{\partial^2 y}{\partial u \partial v} & \frac{\partial^2 y}{\partial u \partial z} \\ \frac{\partial^2 y}{\partial v \partial u} & \frac{\partial^2 y}{\partial v^2} & \frac{\partial^2 y}{\partial v \partial z} \\ \frac{\partial^2 y}{\partial z \partial u} & \frac{\partial^2 y}{\partial z \partial v} & \frac{\partial^2 y}{\partial z^2} \end{bmatrix} \longrightarrow$$ هسین یا هسین به شرط تقیین

$$\mathcal{F}_1 = \frac{1}{2}\left(\mathcal{Y}_\ell - \sum_{k=1}^{n} \delta_k w_k x_k\right)^2 \quad , \quad \delta_k \sim N(1, \sigma^2)$$

$$\frac{\partial \mathcal{F}_1}{\partial w_i} = \frac{1}{2} \frac{\partial}{\partial v_i}\left(\mathcal{Y}_\ell - \sum_{k=1}^{n} \delta_k w_k x_k\right)^2$$

$$= \left(\mathcal{Y}_\ell - \sum_{k=1}^{n} \delta_k w_k x_k\right) x - \delta_i x_i$$

$$\rightarrow E\left[\frac{\partial \mathcal{F}_1}{\partial w_i}\right] = -x_i \mathcal{Y}_\ell E[\delta_i] + x_i \sum_{k=1}^{n} E[\delta_i \delta_k] w_k x_k$$

$$E[\delta_i \delta_k] = \begin{cases} E[\delta_i]^2 + var(\delta_i) = 1 + \sigma^2 & i = k \\ E[\delta_i] E[\delta_k] = 1 & i \neq k \end{cases}$$

$$\rightarrow E\left[\frac{\partial \mathcal{F}_1}{\partial w_i}\right] = -x_i \mathcal{Y}_\ell + x_i\left((1+\sigma^2) w_i x_i + \sum_{\substack{k=1 \\ k \neq i}}^{n} w_k x_k\right)$$

$$= -x_i\left(\mathcal{Y}_\ell - (1+\sigma^2) w_i x_i - \sum_{\substack{k=1 \\ k \neq i}}^{n} w_k x_k\right)$$

non-regularized $\longrightarrow$ $\tilde{J_i} = \frac{1}{2}(y_d - \sum\limits_{k=1}^{n} w_k x_k)^2$

$\longrightarrow$ $\dfrac{\partial \hat{J_i}}{\partial w_i} = -x_i \times (y_d - w_i x_i - \sum\limits_{\substack{k=1 \\ k \neq i}}^{n} w_k x$

یک ترم $\sigma^2$ کم شده

و در واقع ما یک ترم $-\sigma^2 w_i x_i$ به مشتق اضافه کرده‌ایم و این حالت

این هست که ما یک ترم نویز $w_i$ به اضافه کرده‌ایم.

---

سؤال 6 -

$e_k = x_k - x^* \longrightarrow x^* = x_k - e_k$

$f(x) = g'(x) \longrightarrow x_{k+1} = x_k - \dfrac{g'(x)}{g''(x)} = x_k - \dfrac{f(x)}{f'(x)}$

$f(x^*) = f(x_k - e_k) \simeq f(x_k) - e_k f'(x_k) + \dfrac{e_k^2}{2} f''(\xi_k)$

$x_k < \xi_k < x^*$

$f(x^*) = 0 \longrightarrow f(x_k) - e_k f'(x_k) + \dfrac{e_k^2}{2} f''(\xi_k) = 0$

$$\longrightarrow \quad \frac{f(x_k)}{f'(x_k)} - e_k + \frac{e_k^2}{2f'(x_k)} f'''(\xi_k)$$

$$\longrightarrow \quad x_k - x_{k-1} - x_k + x^* + \frac{(x_k - x^*)^2}{2f'(x_k)} f'''(\xi_k)$$

$$\longrightarrow \quad x_{k-1} - x^* = \frac{(x_k - x^*)^2}{2f'(x_k)} f''(\xi_k) \longrightarrow |x_{k-1} - x^*| \leqslant \frac{(x_k - x^*)^2}{2|f'(x_k)|} |f''(\xi_k)|$$

↳ اگر $k \to \infty$ با توجه به پیوستگی فرض ها داشت :

$$f'(x_k) \to f'(x^*) \quad , \quad x_k \to x^*$$

---

$$\frac{\partial}{\partial z_i} L(z, y) = \frac{\partial}{\partial z_i} \left( -\sum_{k=1}^{K} y_k \log \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}} \right) = \frac{\partial}{\partial z_i} \left( -\sum_{k=1}^{K} y_k z_i + \sum_{k=1}^{K} y_k (\log \sum_{j=1}^{K} e^{z_j}) \right)$$

$$= -\sum_{k=1}^{K} y_k + \sum_{k=1}^{K} y_k \left( \frac{z_i}{\sum_{j=1}^{K} e^{z_j}} \right) = \hat{y}_i - 1 = \hat{y}_i - y_i$$

$$\longrightarrow \nabla_z L(z, y) = \hat{y} - y$$

ب)

$$\frac{\partial^2}{\partial z_i^2} L(z,\hat{y}) = \frac{\partial}{\partial z_i}(\hat{y}_i - 1) = \frac{\partial}{\partial z_i}\left(\frac{e^{z_i}}{\sum\limits_{j=1}^{k} e^{z_j}} - 1\right) = \hat{y}_i(1 - \hat{y}_i)$$

$$\frac{\partial^2}{\partial z_i \partial z_j} L(z,\hat{y}) = \frac{\partial}{\partial z_j}\left(\frac{e^{z_i}}{\sum\limits_{j=1}^{k} e^{z_j}} - 1\right) = -\hat{y}_i \hat{y}_j$$

$$\longrightarrow H = \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & -\hat{y}_1\hat{y}_2 & \cdots & -\hat{y}_1\hat{y}_k \\ -\hat{y}_2\hat{y}_1 & \hat{y}_2(1-\hat{y}_2) & & -\hat{y}_2\hat{y}_k \\ \vdots & & \ddots & \vdots \\ -\hat{y}_k\hat{y}_1 & -\hat{y}_k\hat{y}_2 & \cdots & -\hat{y}_k(1-\hat{y}_k) \end{bmatrix}$$

$$\longrightarrow H = diag(\hat{y}) - \hat{y}\,\hat{y}^T$$

$$\longrightarrow x^T H x = x^T diag(\hat{y})\, x - x^T \hat{y}\hat{y}^T x = \sum_{i=1}^{k}\hat{y}_i x_i^2 - (\hat{y}^T x)^2$$

$$(\hat{y}^T x)^2 = \left(\sum_{i=1}^{k}\hat{y}_i x_i\right)^2$$

$$\sum_{i=1}^{k}\hat{y}_i x_i^2 = \left(\sum_{i=1}^{k}\hat{y}_i\right)\left(\sum_{i=1}^{k}\hat{y}_i x_i^2\right)$$

طبق نامساوی کوشی ـ شوارتز :

$$\left(\sum_{i=1}^{k}\hat{y}_i x_i\right)^2 \leq \left(\sum_{i=1}^{k}\hat{y}_i\right)\left(\sum_{i=1}^{k}\hat{y}_i x_i^2\right) \longrightarrow \sum_{i=1}^{k}\hat{y}_i x_i^2 \geq (\hat{y}^T x)^2$$

$$PSD \longleftarrow x^T H x \geq 0 \quad \longleftarrow$$

چون ماتریس Hessian ، PSD هست تابع ما محدب هست.