

(۱) نشان می‌دهم که Hessian تابع یک ماتریس PSD است پس این تابع convex و دارای minimum است
 رابا آلوگورتم یادگیری کانفرنس می‌توانیم به نقطه minimum آن برسیم
 می‌توان به درون w و با extend کردن x در نظر گرفت.

$$z = \sigma(h)$$

$$h = w^T x + w_0 = w^T x$$

$$E(w) = - \sum_i y_i \ln(z_i) + (1 - y_i) \ln(1 - z_i)$$

$$\frac{\partial E(w)}{\partial w} = - \sum \frac{\partial(-)}{\partial w}$$

$$\begin{aligned} \frac{\partial(-)}{\partial w} &= \frac{\partial(y_i \ln(z_i) + (1 - y_i) \ln(1 - z_i))}{\partial w} = y_i \frac{\partial \ln(z_i)}{\partial z_i} * \frac{\partial(z_i)}{\partial h_i} * \frac{\partial h_i}{\partial w} \\ &\quad + (1 - y_i) \frac{\partial \ln(1 - z_i)}{\partial z_i} * \frac{\partial(z_i)}{\partial h_i} * \frac{\partial h_i}{\partial w} \\ &= \frac{\partial(z_i)}{\partial h} * \frac{\partial h}{\partial w} \left[y_i \frac{\partial \ln(z_i)}{\partial z_i} + (1 - y_i) \frac{\partial \ln(1 - z_i)}{\partial z_i} \right] \end{aligned}$$

$$\frac{\partial(z_i)}{\partial h_i} = \sigma(h_i) (\sigma(h_i) - 1) = z_i (z_i - 1)$$

$$\frac{\partial h_i}{\partial w} = x_i$$

$$\Rightarrow \frac{\partial(-)}{\partial w} = x_i (z_i - y_i) \Rightarrow \frac{\partial E(w)}{\partial w} = - \sum_i x_i (\hat{y}_i - y_i) = \sum_i x_i (y_i - \hat{y}_i)$$

$$= X^T (y - \hat{y})$$

$$\begin{cases} \frac{\partial E}{\partial w} = \sum_i x_i (y_i - \hat{y}_i) \\ \frac{\partial E}{\partial b} = \sum_i (y_i - \hat{y}_i) \end{cases}$$

$$H = \frac{\partial}{\partial w} \left[\sum_i x_i (y_i - z_i) \right] = \sum_i x_i \frac{\partial}{\partial w} [y_i - z_i]$$

$$\frac{\partial x_i [y_i - z_i]}{\partial w} = -x_i \frac{\partial z_i}{\partial h_i} * \frac{\partial h_i}{\partial w} = -x_i (z_i) (z_i - 1) x_i^T = \underbrace{(z_i)(1 - z_i)}_{\alpha_i} x_i x_i^T$$

$$0 < z_i = \sigma(h_i) < 1$$

$$0 < 1 - z_i < 1$$

$$\Rightarrow \alpha_i < 1$$

Positive

$$\Rightarrow H = \sum_i \alpha_i x_i x_i^T = X^T \text{diag}(\alpha_i) X = \underbrace{X^T \sqrt{\text{diag}(\alpha_i)}}_{A^T} \underbrace{(\sqrt{\text{diag}(\alpha_i)} X)}_A$$

$$H = A^T A \Rightarrow H \succ 0 \Rightarrow \text{یک نقطه بهینه دارد}$$

$$\text{آلوگورتم یادگیری} \Rightarrow w^{t+1} = w^t - \eta \frac{\partial E}{\partial w}, \quad b^{t+1} = b^t - \eta \frac{\partial E}{\partial b}$$

covariate shift مربوط به تغییر توزیع ورودی ها و node های یک شبکه به ازای نمونه های متفاوت است که می تواند باعث کاهش performance شبکه شود. در حالت کلی ممکن است با آنتروپی شبکه برای داده ها جدید تعمیم پذیری نداشته باشد و این باعث ایجاد مشکلاتی در مدل های شبکه می شود. BN با نرمالیزه کردن داده های miniBatch به سایر مشخص این شکل که متغیر می نماید. در این حالت توزیع نمونه های ورودی شبکه به هم نزدیک اند همچنین در جریان از یک لایه به لایه دیگر نیز تغییر زیادی نمی کند و در نتیجه این روش موجب کاهش اثر covariate shift و مدل های بهتر و سریعتر عمل می شود.

BN مانند اضافه کردن noise به داده ها ورودی است که یک روش regularization مانند dropout است بنابراین به تعمیم پذیری مدل کمک می کند. همچنین با تغییر ورودی ها از overfit جلوگیری می کند و به جریان برداشتن نیز کمک می کند و از بزرگ شدن یا کوچک شدن وزن آن جلوگیری می کند.

$$\frac{\partial L}{\partial x_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{x}_j} \times \frac{\partial \hat{x}_j}{\partial x_i} \quad \Rightarrow \quad \frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial \hat{x}_i} - \frac{1}{n} \sum_{j=1}^n \frac{\partial L}{\partial \hat{x}_j}$$

$$\frac{\partial \hat{x}_j}{\partial x_i} = \begin{cases} 1 - 1/n & i=j \\ -1/n & i \neq j \end{cases} = \delta_{ij} - \frac{1}{n}$$

$$\frac{\partial L}{\partial \hat{x}_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}_j} \times \frac{\partial \hat{y}_j}{\partial \hat{x}_i} \quad \Rightarrow \quad \frac{\partial L}{\partial \hat{x}_i} = \gamma \frac{\partial L}{\partial \hat{y}_i}$$

$$\frac{\partial \hat{y}_j}{\partial \hat{x}_i} = \gamma \delta_{ij}$$

$$\Rightarrow \frac{\partial L}{\partial x_i} = \gamma \frac{\partial L}{\partial \hat{y}_i} - \frac{\gamma}{n} \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}_j}$$

$$\frac{\partial L}{\partial x_i} = ? \quad \Rightarrow \quad \frac{\partial L}{\partial x_i} = 0 \quad n=1$$

ت زمانی که یک داده داریم، تغییرات نسبت به ورودی ها ثابت است چون $\hat{x}=0$

$$n \rightarrow \infty \Rightarrow \gamma \frac{\partial L}{\partial \hat{y}_i}$$

اما هر چه n بزرگتر می شود به سمت صاف می رود، تغییرات نسبت به ورودی

به صورت مستقل محاسبه می گردد و از سایر ورودی ها Batch و نرارتی نمی پذیرد و در نتیجه اثر نرمالیزاسیون روی هر نقطه مستقل از سایر ورودی ها است.

(۳) انت

$$\hat{y}_k = \frac{e^{z_k^{(r)}}}{\sum_d e^{z_d^{(r)}}}$$

$$\frac{\partial \hat{y}_k}{\partial z_i^{(r)}} = \begin{cases} \hat{y}_k (1 - \hat{y}_k) & k=i \\ -\frac{e^{z_k^{(r)}} e^{z_i^{(r)}}}{\left(\sum_{d=1}^k e^{z_d^{(r)}}\right) \left(\sum_{d=1}^k e^{z_d^{(r)}}\right)} = -\hat{y}_k \hat{y}_i & k \neq i \end{cases}$$

$$\frac{\partial L}{\partial z_k^{(r)}} = \frac{\partial}{\partial z_k^{(r)}} \sum_i y_i \log(\hat{y}_i) = \frac{\partial}{\partial z_k^{(r)}} - \sum_i y_i \left[z_i^{(r)} - \log \sum_d e^{z_d^{(r)}} \right] \quad (-)$$

$$= \frac{\partial}{\partial z_k^{(r)}} \sum_{i=1}^k y_i \log \sum_d e^{z_d^{(r)}} = \sum_{i=1}^k y_i \frac{e^{z_k}}{\sum_d e^{z_d^{(r)}}} = \sum_{i=1}^k y_i \hat{y}_k = \hat{y}_k$$

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z^{(1)}} \times \frac{\partial z^{(r)}}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial \hat{a}^{(1)}} \times \frac{\partial \hat{a}^{(1)}}{\partial z^{(1)}} \times \frac{\partial z^{(1)}}{\partial w^{(1)}} \quad \times$$

$$\frac{\partial \hat{a}^{(1)}}{\partial z_i^{(1)}} = \begin{cases} 1 & z_i^{(1)} \geq 0 \\ 0 & z_i^{(1)} < 0 \end{cases} \quad \text{vector}$$

$$\frac{\partial a^{(1)}}{\partial \hat{a}^{(1)}} = \begin{cases} 1 & \text{with probability of } y \\ 0 & \text{with probability of } 1-y \end{cases} \quad \text{vector}$$

$$\frac{\partial z^{(r)}}{\partial a^{(1)}} = w^{(r)}$$

$$\frac{\partial L}{\partial z^{(r)}} = \hat{y}_k$$

بالنسبة
لـ $z^{(r)}$

Jacobian of $R^m \rightarrow R^n$ is $J \in R^{n \times m}$

(1)

$$H(y) = \begin{pmatrix} \frac{\partial^2 y}{\partial u^2} & \frac{\partial^2 y}{\partial u \partial v} & \frac{\partial^2 y}{\partial u \partial z} \\ \frac{\partial^2 y}{\partial v \partial u} & \frac{\partial^2 y}{\partial v^2} & \frac{\partial^2 y}{\partial v \partial z} \\ \frac{\partial^2 y}{\partial z \partial u} & \frac{\partial^2 y}{\partial z \partial v} & \frac{\partial^2 y}{\partial z^2} \end{pmatrix}$$

$$\nabla(y) = \begin{pmatrix} \frac{\partial y}{\partial u} \\ \frac{\partial y}{\partial v} \\ \frac{\partial y}{\partial z} \end{pmatrix}$$

$$F = \begin{pmatrix} f_1(u, v, z) \\ f_r(u, v, z) \\ f_c(u, v, z) \end{pmatrix}$$

$$J(F) = \begin{pmatrix} \frac{\partial f_1}{\partial u} & \frac{\partial f_1}{\partial v} & \frac{\partial f_1}{\partial z} \\ \frac{\partial f_r}{\partial u} & \frac{\partial f_r}{\partial v} & \frac{\partial f_r}{\partial z} \\ \frac{\partial f_c}{\partial u} & \frac{\partial f_c}{\partial v} & \frac{\partial f_c}{\partial z} \end{pmatrix}$$

$$J(\nabla y) = H(y)$$

$$J = \frac{1}{2} (y_d - \sum d_k w_k x_k)^2$$

(5)

$$\frac{\partial J}{\partial w_i} = - (y_d - \sum_{k=1}^n d_k w_k x_k) d_i x_i = \sum_{k=1}^n d_i d_k x_i x_k w_k - y_d d_i x_i$$

$$E \left[\frac{\partial J}{\partial w_i} \right] = E \left[\sum_{\substack{k=1 \\ k \neq i}}^n d_i d_k x_i x_k w_k \right] + E \left[d_i^2 x_i^2 w_i \right] - E \left[y_d d_i x_i \right]$$

$$E[d_i d_k] = E[d_i] E[d_k] = 1$$

$$(1 + \sigma^2) x_i^2 w_i$$

$$- x_i y_d$$

$$\sum_{\substack{k=1 \\ k \neq i}}^n x_i w_k x_k$$

$$= \sum_{\substack{k=1 \\ k \neq i}}^n x_i w_k x_k + (1 + \sigma^2) w_i x_i^2 - x_i y_d = x_i \left[\sum_{k=1}^n w_k x_k + \sigma^2 w_i - y_d \right]$$

$$= - x_i \left[y_d - \sum_{k=1}^n w_k x_k \right] + w_i x_i^2 \sigma^2$$

non-regularized: $J_c = \frac{1}{2} (y_d - \sum w_k x_k)^2$

$$\frac{\partial J_c}{\partial w_i} = - x_i \left[y_d - \sum_{k=1}^n w_k x_k \right]$$

$$\Rightarrow E \left[\frac{\partial J}{\partial w_i} \right] = \frac{\partial J_c}{\partial w_i} + w_i x_i^2 \text{var}(\sigma_i)$$

ماهمی می شود که ضرایب w_k با یک متغیر تصادفی d_k نویزی می شوند و برداریان تابع loss نسبت به ضرایب مانند regularization عمل می کند و مقدار نویز به ضرایب اضافه می گردد

$$f = g(x)$$

ج ۲ اثبات عددی روش نیوتن:

$$f(x^*) = 0, f'(x^*) \neq 0, f'(x)$$

فرض می شود که x_k به x^* میل می کند:

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}; k = 1, 2, \dots$$

$$|x_{k+1} - x^*| \leq M |x_k - x^*|^2 \text{ if } M > \frac{|f''(x^*)|}{2|f'(x^*)|}$$

$$\text{اثبات: } e_k = x_k - x^* \Rightarrow x^* = x_k - e_k$$

$$\text{توسعه تیلور: } f(x) = f(x_k) + (x - x_k)f'(x_k) + R_1$$

$$R_1 = \frac{f''(\xi)}{2!} (x - x_k)^2, x_k < \xi < x^*$$

$$\Rightarrow \underbrace{f(x_k - e_k)}_{f(x^*) = 0} = f(x_k) - e_k f'(x_k) + \frac{(e_k)^2}{2!} f''(\xi)$$

$$\Rightarrow 0 = f(x_k) - e_k f'(x_k) + \frac{e_k^2}{2!} f''(\xi)$$

توجه کنید که x_k به اندازه کافی به x^* نزدیک باشد، چون $f(x^*) = 0$ پس $f(x_k) \neq 0$

$$\Rightarrow 0 = \frac{f(x_k)}{f'(x_k)} - \underbrace{e_k}_{(x_k - x^*)} + \frac{e_k^2}{2!} \frac{f''(\xi)}{f'(x_k)}$$

$$\Rightarrow 0 = \frac{f(x_k)}{f'(x_k)} - (x_k - x^*) + \frac{e_k^2}{2!} \frac{f''(\xi)}{f'(x_k)}$$

$$\xrightarrow{-x_{k+1}} \Rightarrow x_{k+1} - x^* = \frac{(e_k)^2}{2!} \frac{f''(\xi)}{f'(x_k)}$$

$(x_k - x^*)$

$$\Rightarrow x_{k+1} - x^* = (x_k - x^*)^2 \frac{f''(\xi)}{2f'(\eta)}$$

$$\Rightarrow |x_{k+1} - x^*| \leq \frac{|f''(\xi)|}{2|f'(\eta)|} |x_k - x^*|^2$$

بأنوجه به یوستلی f' و f'' ، $f'(x_k)$ به $f'(x^*)$ همگرا می شود و آنکه سین x_k و x^* است نیز به x^* همگرا می شود و $f''(\xi)$ به $f''(x^*)$ همگرا می شود

$$\Rightarrow |x_{k+1} - x^*| \leq M |x_k - x^*|^2, \quad M > \frac{|f''(x^*)|}{2|f'(x^*)|}$$

x^* نقطه ریشه $g(x)$ است چرا که $g'(x^*) = 0$

بأنوجه به نتیجه قبلی:

$$|x_{k+1} - x^*| \leq M |x_k - x^*|^2, \quad M > \frac{|g^{(3)}(x^*)|}{2|g''(x^*)|}$$

$$L(z, y) = - \sum_{i=1}^k y_i \ln \hat{y}_i$$

(1) (v)

$$\frac{\partial L}{\partial z_j} = \frac{\partial}{\partial z_j} \left[\sum_{i=1}^k y_i (z_i - \ln(\sum_{d=1}^k e^{z_d})) \right] = \frac{\partial}{\partial z_j} \left[- \sum_{i=1}^k y_i z_i \right] + \frac{\partial}{\partial z_j} \left[\sum_{i=1}^k y_i \ln(\sum_{d=1}^k e^{z_d}) \right]$$

$$- \frac{\partial}{\partial z_j} \sum_{i=1}^k y_i z_i = - \hat{y}_j, \quad \frac{\partial}{\partial z_j} \left[\sum_{i=1}^k y_i \ln(\sum_{d=1}^k e^{z_d}) \right] = \sum_{i=1}^k y_i \frac{e^{z_j}}{\sum_{d=1}^k e^{z_d}} = \sum_{i=1}^k \hat{y}_i \hat{y}_j$$

$$\Rightarrow \frac{\partial L}{\partial z_j} = \sum_{i=1}^k \hat{y}_i \hat{y}_j - \hat{y}_j = \hat{y}_j - \hat{y}_j \Rightarrow \boxed{\nabla_z L = \hat{y} - y}$$

one hot
همون یه صورت
است

$$H = \nabla_z^2 L = \begin{bmatrix} \frac{\partial^2 L}{\partial z_1 \partial z_1} & \dots & \frac{\partial^2 L}{\partial z_1 \partial z_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 L}{\partial z_k \partial z_1} & \dots & \frac{\partial^2 L}{\partial z_k \partial z_k} \end{bmatrix}, \quad \frac{\partial^2 L}{\partial z_i \partial z_j} = \frac{\partial}{\partial z_j} (\hat{y}_i - \hat{y}_i) = \begin{cases} \hat{y}_i - \hat{y}_i^2 & i=j \\ -\hat{y}_i \hat{y}_j & i \neq j \end{cases} \quad (1)$$

$$\Rightarrow \boxed{H = \text{diag}(\hat{y}_i) - \hat{y} \hat{y}^T}, \quad x^T H x = x^T \text{diag}(\hat{y}_i) x - x^T \hat{y} \hat{y}^T x$$

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

باینری به انگلیسی یعنی احتمال است
و یک عدد مثبت است من به صورت $(\hat{y}_i \sqrt{x_i})^2$ در نظر بگیرم و معنی $\sum_{i=1}^k \hat{y}_i = 1$ چون احتمال است

$$\Rightarrow \text{حسن کوشش کوآرتز} \quad \left(\sum x_i \sqrt{\hat{y}_i} \sqrt{\hat{y}_i} \right)^2 \leq \sum (x_i \sqrt{\hat{y}_i})^2 \leq (\sqrt{\hat{y}_i})^2 = \sum x_i^2 \hat{y}_i \leq \hat{y}_i$$

یعنی چون احتمال است $\sum \hat{y}_i = 1$

$$\Rightarrow \boxed{\sum x_i^2 \hat{y}_i \geq \left(\sum x_i \hat{y}_i \right)^2}$$

$$\Rightarrow x^T H x \geq 0 \quad \forall x \in \mathbb{R}^k \Rightarrow H \geq 0 \quad H \text{ is PSD} \Rightarrow \text{مثبت غیر منفی است}$$

تابع زیان cross entropy نسبت به z یک تابع محدب است