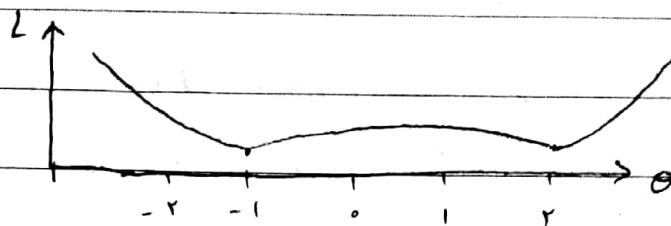


$$L = \frac{1}{N} \sum_{i=1}^N |y_i - f_{\theta}(x_i)|, \quad f_{\theta}(x_i) = x_i(\theta^2 - \theta) \quad \text{① (الف)}$$

$$D = \{(3, 2), (2, 0), (1, 2)\}$$

$$L = \frac{1}{3} [|4 - 3(\theta^2 - \theta)| + |0 - 2(\theta^2 - \theta)| + |2 - (\theta^2 - \theta)|]$$

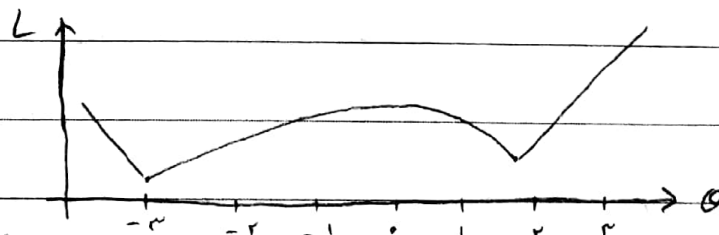


ملاحظه شود که مقدار L دارای دو کینه محلی که برابر می هستند می باشد. / به نظر می آید استاندارد از gd به یک سیستم محلی می رسیم که در اینجا هر دو سیستم محلی برابر می باشند.

$$L = \frac{1}{N} \sum_{i=1}^N |y_i - f_{\theta}(x_i)|, \quad f_{\theta}(x_i) = \ln(1 + e^{x_i \theta}) \quad \text{ب)}$$

$$D = \{(3, 4), (-1, 3), (1, 0)\}$$

$$L = \frac{1}{3} [|4 - \ln(1 + e^{3\theta})| + |3 - \ln(1 + e^{-\theta})| + |0 - \ln(1 + e^{\theta})|]$$



ملاحظه شود که مقدار L دارای دو کینه محلی که یکی از آن ها برابر است با کینه دیگر است. با تقدم به وزن دمی اولیه ممکن است به کینه محلی که برابر است با کینه دیگر رسیم.

ج) معمولاً مقدار زیادی کینه محلی وجود دارد که برابر می هستند و ممکن است به این کینه های محلی

همراه رسیم. برای که کردن افعال رفتار این پدیده می توان از تکنیک های زیر بهره برد:

مقدار دمی اولیه وزن مابین مدت فواصل، تنظیم مقیاس یادآور $Regularization$ ، $Momentum$ ، $learning rate$

$Batch Normalization$ ، روش های Ensemble

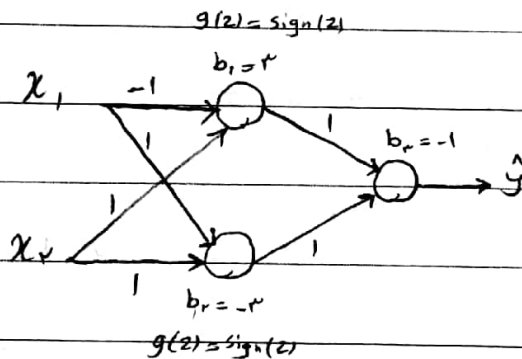
(۲) آ) با ترمال کردن داده‌های ورودی، بازه مقایسه ویژگی‌ها به یک بازه یکسان در منفی
 معمولاً [ارواح] نگاشت می‌شوند. این بازه مقایسه یکسان برای ویژگی‌ها باعث حلاله
 که الگوریتم بهینه‌سازی بتواند راحت‌تر و سریع‌تر به بهینه‌ترین حالتی برای حل مسئله
 در واقع با ترمال کردن داده‌ها، قدم‌ها در gradient descent مقدار کمتر می‌شود
 و می‌توان از مقادیر learning rate های بزرگتر استفاده کرد که باعث حرکتی سریع‌تر
 مدل می‌شود. علاوه بر این ترمال کردن داده‌ها، مقایسه گره‌های درخت را نیز کوچک
 نگه می‌دارد و به این طریق ~~از~~ انقباض / پایداری مدل گره‌های درخت را گامی می‌باید.
 همچنین ترمال کردن داده‌ها می‌تواند باعث بهبود دقت نهایی مدل نیز شود، زیرا با این کار
 اهمیت / وزن های یکسان برای تمامی ویژگی‌ها در نظر گرفته می‌شود و برای صورت یک
 ویژگی صرفاً چون مقایسه بزرگتری دارد، نمی‌تواند باعث اهمیت و اثر دادن بیشتر مدل به آن شود.

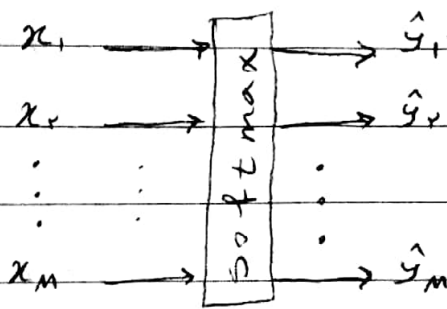
ب) استفاده از منظم‌سازها که معمولاً به صورت اضافه کردن نرم به دار وزن‌ها به تابع هزینه است،
 معمولاً باعث می‌شود که مقایسه وزن‌ها خیلی بزرگ نشود. از این رو به این روش‌ها
 weight decay نیز می‌گویند. این جدگیری از افزایش مقادیر وزن‌ها باعث می‌شود که مدل
 زیادی به یک یا چند ویژگی اهمیت ندهد و به این طریق از overfit شدن مدل
 جدگیری می‌شود. البته روش‌های دیگری مانند $\text{Batch Normalization}$ و dropout نیز
 اثر منظم‌ساز دارند. به طور کلی استفاده از منظم‌سازها باعث افزایش و بهبود عملکرد مدل روی
 داده‌های که در زمان آموزش آن‌ها را ندیده می‌شود. در نهایت منظم‌سازها باعث می‌شود که
 وزن ویژگی‌های که خیلی مهم هستند مهم‌تر شوند و به این طریق generalizability مدل افزایش می‌یابد.
 ج) هر مقدار لایه لایه نرم اما با تابع فعال‌ساز خطی، هیچ قدرت اضافه کردن نیست به یک لایه خطی ندارد. زیرا
 هیچ non linearity به مدل اضافه نشده است. شبکه دوم اما چون non linearity دارد قدرت مدل کردن
 را بعد از خطی را به قدری که اول دارد، به این قدرت مدل کردن شبکه دوم بیشتر است. آنگاه می‌تواند
 اول ممکن است می‌تواند شبکه دوم بیشتر باشد ولی هیچ برتری نسبت به مدل‌های خطی ندارد.

(۳) (آ) مشخص است که هیچ رمز تقسیم‌گیری فعلی که خودی یک بدل فعل است، نمی‌تواند داده‌های موجود در شکل را به‌طور کامل در فزب دسته‌بندی کند.

با رمز تقسیم‌گیری که به صورت $X_2 = |X_1 - 3|$ باشد می‌تواند تمامی داده‌ها را به خودی دسته‌بندی کند. بنابراین در لایه اول نیاز به دو نورون داریم که حرکت از آنها وظیفه بدل‌سازی حرکت از بین فعل‌های این رمز تقسیم را دارند. در لایه آخر هم نیاز به یک گیت AND داریم که مشخص کند یک نمونه داده بالای ورودی بین فعل باشد. جای آنکه در لایه آخر از گیت AND استفاده کنیم لازم است که تابع فعال ساز در لایه اول تابع $\text{sign}(z)$ باشد تا بتوانیم خودی لایه اول را به صورت باینری تفسیر کنیم.

$$\left. \begin{aligned} \text{معادله بین فعل اول: } -X_1 + X_2 + 3 &= 0 \\ \text{معادله بین فعل دوم: } X_1 - X_2 - 3 &= 0 \end{aligned} \right\}$$





(f)

$$L(\hat{y}, y) = - \sum_{i=1}^m y_i \log \hat{y}_i \quad \text{if true label} = t \rightarrow L(\hat{y}, y) = - \log \hat{y}_t$$

$$\hat{y}_t = \frac{e^{x_t}}{\sum_{m=1}^m e^{x_m}}$$

$$\frac{\partial L}{\partial x_j} = \frac{\partial L}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial x_j}$$

$$L = -\log \hat{y}_t \Rightarrow \frac{\partial L}{\partial \hat{y}_t} = -\frac{1}{\hat{y}_t}$$

$$\text{if } j \neq t \Rightarrow \frac{\partial \hat{y}_t}{\partial x_j} = -\frac{e^{x_j} e^{x_t}}{(\sum_{m=1}^m e^{x_m})^2} = -\hat{y}_j \hat{y}_t$$

$$\text{if } j = t \Rightarrow \frac{\partial \hat{y}_t}{\partial x_j} = \frac{e^{x_t} \sum_{m=1}^m e^{x_m} - e^{2x_t}}{(\sum_{m=1}^m e^{x_m})^2} = \hat{y}_j (1 - \hat{y}_j)$$

$$\Rightarrow \frac{\partial L}{\partial x_j} = \begin{cases} \hat{y}_j & ; j \neq t \\ \hat{y}_j - 1 & ; j = t \end{cases} \Rightarrow \frac{\partial L}{\partial x_j} = \hat{y}_j - y_j$$

$$\begin{matrix} x_1 & 0 \\ x_2 & 0 \\ \vdots & \vdots \\ x_n & 0 \end{matrix} \quad \begin{matrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{matrix}$$

m نود

⑤ اگر تعداد نودهای لایه پنهان را m فرض کنیم

بنابراین از آنجایی که $x^{(i)} \in \mathbb{R}^n$ آنگاه باید

$$W_1 \in \mathbb{R}^{m \times n}$$

همچنین از آنجایی که مسئله رگرسیون باید در

لایه خروجی یک نود داشته باشیم بنابراین باید

$$W_2 \in \mathbb{R}^{1 \times m}$$

از آنجایی که هر نود یک بایاس دارد بنابراین $b_1 \in \mathbb{R}^{m \times 1}$ و $b_2 \in \mathbb{R}$

ب) ابعاد وزن ها و بایاس به تعداد داده های ورودی شبکه وابسته است بنابراین ابعاد آن ها تغییر می کنند.

$$L = \frac{1}{m} \left[\sum_{i=1}^m -y^{(i)} \ln \hat{y}^{(i)} - (1-y^{(i)}) \ln (1-\hat{y}^{(i)}) \right] \quad (2)$$

$$\frac{\partial L}{\partial \hat{y}^{(i)}} = \frac{1}{m} \left[\sum_{i=1}^m -\frac{y^{(i)}}{\hat{y}^{(i)}} + \frac{1-y^{(i)}}{1-\hat{y}^{(i)}} \right] = \frac{1}{m} \left[\sum_{i=1}^m \frac{\hat{y}^{(i)} - y^{(i)}}{\hat{y}^{(i)}(1-\hat{y}^{(i)})} \right]$$

$$\frac{\partial \hat{y}^{(i)}}{\partial z_r} = g(z_r)(1-g(z_r)) = \hat{y}^{(i)}(1-\hat{y}^{(i)})$$

$$\frac{\partial z_r}{\partial a_1} = w_r$$

$$\frac{\partial a_r}{\partial z_1} = \begin{cases} 1 & ; z_1 > 0 \\ 0 & ; z_1 < 0 \end{cases}$$

$$\frac{\partial z_1}{\partial w_1} = x^{(i)}$$

⑥) از آزادی که منظم ساز به ما می‌دهد، به سبب L_2 ، $weight\ decay$ نیز گفته می‌شود.

$$L = \frac{1}{N} \|Xw - y\|^2 + \lambda \|w\|^2 \quad (ب)$$

$$\Rightarrow \frac{\partial L}{\partial w} = \frac{2}{N} X^T (Xw - y) + 2\lambda w = 0$$

$$\Rightarrow X^T X w - X^T y + \lambda N w = 0$$

$$\Rightarrow w(X^T X + \lambda N I) = X^T y \Rightarrow w = (X^T X + \lambda N I)^{-1} X^T y$$

ج) ابتدا ثابت می‌کنیم ماتریس $X^T X$ ماتریس $positive\ semi\ definite$ است. (و PSD بین

ماتریس M وای است: M positive semi definit $\Leftrightarrow v^T M v \geq 0$ for all v

$$v^T X^T X v = (Xv)^T (Xv) = \|Xv\|^2 \geq 0 \quad \text{حال داریم}$$

بنابراین $X^T X$ ماتریس PSD است یعنی $eigen\ value$ های آن بزرگتر مساوی صفر هستند.

حال ثابت می‌کنیم ماتریس $\lambda N I$ ماتریس $positive\ definite$ است. شرط P.D بین ماتریس

M وای است: M positive definite $\Leftrightarrow v^T M v > 0$ for all $v \neq 0$

$$v^T \lambda N I v = \lambda N \|v\|^2 > 0 \quad v \neq 0 \quad \text{حال داریم}$$

بنابراین ماتریس $X^T X$ PSD است، و ماتریس $\lambda N I$ ماتریس PD است. می‌دانیم جمع یک ماتریس

PSD با یک ماتریس PD واری یک ماتریس PD است. از این رو

بنابراین ماتریس $X^T X + \lambda N I$ یک ماتریس PD است یعنی تمام $eigen\ value$ های آن بزرگتر از صفر

هستند، بنابراین این ماتریس $full\ rank$ و معکوس پذیر است.

نتیجه:

نتیجه: