

(1) انت : توابع هزینه : توابعی هستند که با Jacobian شان در شبکه عصبی نسبت به وزنهای شبکه minimize شود. عبارت دیگر، اگر تابع هزینه را با $f(\vec{w}, \vec{x})$ نشان دهیم:

$$\vec{w}^* = \arg \min_{\vec{w}} f(\vec{w}, \vec{x}) \quad (\vec{w}^* \text{ همان بردار وزن بهینه برای شبکه عصبی است.})$$

حقیقت تابع هزینه تقریبی مهمی در update کردن وزنها دارد:

$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} \pm \eta \frac{\partial f(\vec{w}, \vec{x})}{\partial w_{ij}} \quad (\text{learning weights})$$

بنابراین تعریف درست تابع هزینه تأثیر مهمی در training شبکه عصبی دارد. بعنوان مثال، تابع هزینه مورد استفاده در اکثر شبکه‌های عصبی به صورت زیر است:

$$f(\vec{w}, \vec{x}) = \sum_{i \in \text{train}} (\text{sigmoid}(z_i) - \text{actual}_i)^2$$

(z_i خروجی آخرین لایه شبکه عصبی برای sample n_i است.)

• توابع فعالساز (activator): نمایشده مقدار $w^T n$ برای لایه نام به دست آمد، این

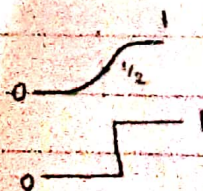
مقدار به یک تابع فعالساز میرود که دامنه‌اش $[0, 1]$ است، پس وارد لایه بعدی میشود.

توابع فعالساز به ما این قابلیت را میدهند که الگوهای غیرخطی بین y و n را تشخیص دهیم.

اگر این توابع نباشند، شبکه عصبی عمیق درگیران خطی خواهد بود:

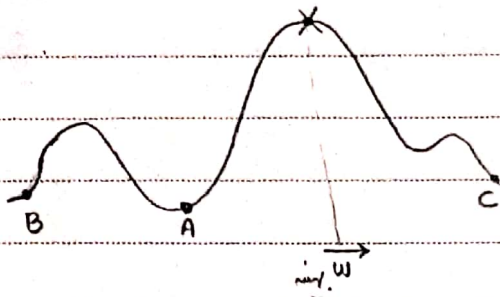


$$\begin{aligned} \text{output} &= w_5 (w_2 n_2 + w_1 n_1) + \\ &w_6 (w_2 n_1 + w_4 n_2) \\ &= \alpha n_1 + \beta n_2 \end{aligned}$$



اگر توابع فعالساز پیوسته باشند، انتقالی به لایه‌ها با مشخص می‌گردد. خود لایه‌ها را مشخص می‌گردد.

(1) ب: فرض کنید توزیع \vec{w} بصورت زیر باشد:



هدف دست یافتن به global optima از طریق gradient descent است. یعنی:

$$\vec{w}_{new} = \vec{w}_{old} - \eta \frac{\partial err}{\partial \vec{w}}$$

حال اینکه ما \vec{w} را با انتخاب \vec{w} (مثلاً میان A و B و C) تأثیرهای زیر را دارد:

(1) اگر \vec{w} نزدیک \vec{w} شروع مطلق باشد، به سرعت به \vec{w} بهینه می‌رسیم و training در زمان کمی انجام می‌شود.

اما در حالتی که چنانچه η مقدار مناسبی داشته باشد، به \vec{w} بهینه می‌رسیم و تعدادی در نتیجه نهایی و زمانها وجود ندارد.

(2) ج: forward عبارت است از ترکیب خطی بردار \vec{w} و خردی شدن آنها در هر لایه propagation.

و عبور دادن از تابع activation یا update نمی‌شود و از لایه input به output می‌رسیم.

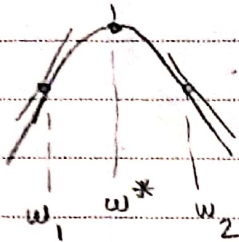
back propagation عبارت است از update کردن \vec{w} که با استفاده از gradient descent

و از لایه آخر به اول باز می‌گردیم.

(1) د: چنانچه η مقدار بسیار کمی داشته باشد، رسیدن به global optima (یا همان train)

بسیار زمانبر خواهد بود. در حالتی که η زیاد باشد، به \vec{w} بهینه می‌رسیم. اگر η مقدار

زیادی هم داشته باشد، ممکن است در local optima گرفتار شویم.



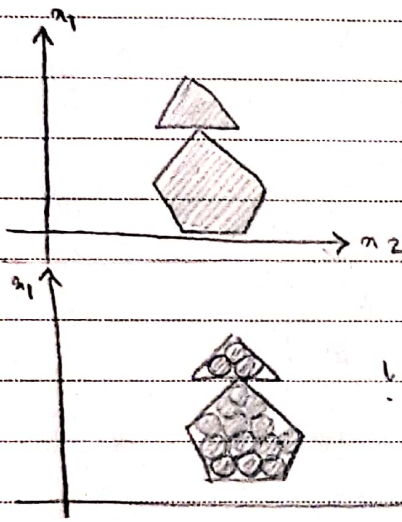
$$w_{new} = w_{old} - \eta \frac{\partial err}{\partial w_{old}}$$

توانستیم که دایره محصوره و خطی نزولی را پیدا کنیم

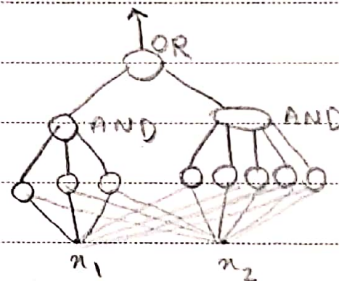
بصورت مداوم حول w^* نوسان می کنیم و علامت w^* را پیدا می کنیم

5) با داشتن لایه های n_1 و n_2 امکان ساختن مرزها و تقسیم پیچیده با تعداد نودان کم فراهم می شود.

در حالتی که نیاز به ساختن مرزها با تعداد لایه کم مریب می شود width (عرض) شبکه بسیار زیاد می شود.



1.5):

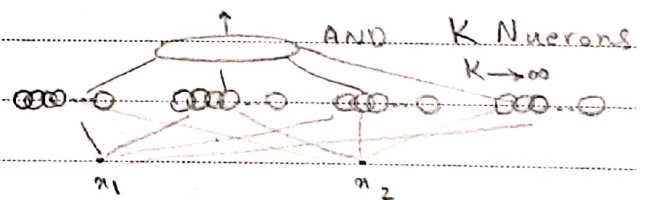


شود. مثلا:

11 Nuerons

2.5): (Polar) (شیرین) یا

دایره



K Nuerons

$K \rightarrow \infty$

دقت شود هر چایره با لایه کم قرار دادن تعداد بسیار زیاد نودان بدست می آید.

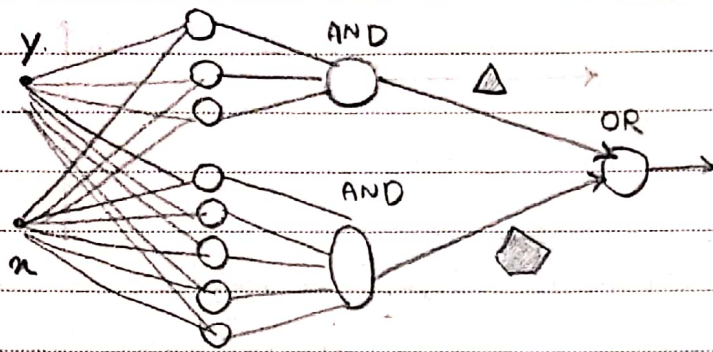
1) Overfitting بزرگترین مشکل train کردن شبکه ها محسوب می شود. از آنجا که ما سعی می کنیم

$$\text{تابع هزینه} = \frac{1}{2} \sum_i [\text{sigmoid}(z_{out_i}) - \text{actual}_i]^2$$

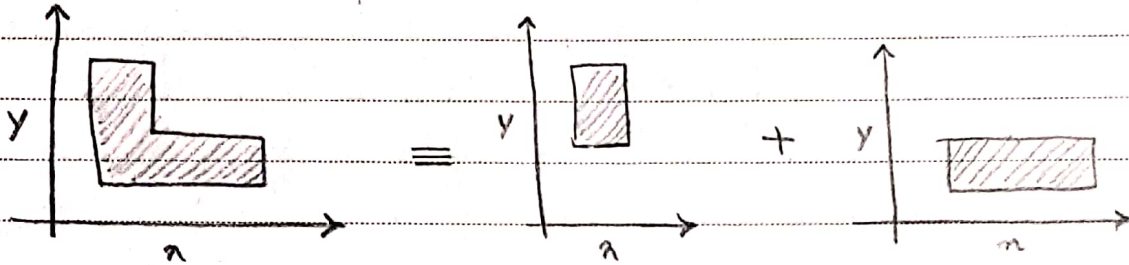
که برعکس، شبکه روی داده های train به دقت بسیار بالایی می رسد که این از علامت overfitting

است. از جمله راه حل هایی که وجود دارد اینست که اجازه حجم در مرحله train شبکه overfit ندهد و

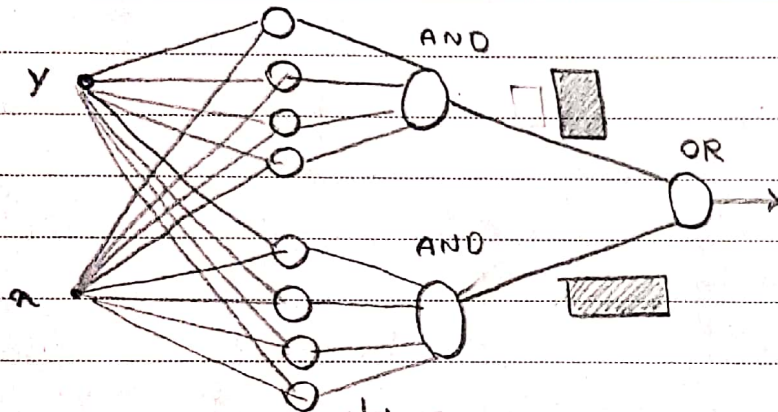
دوینا بهینه شوند، اما در مرحله test تعدادی از دونه ها را مغفول کنیم، یا تعدادی از نودها را حذف کنیم.



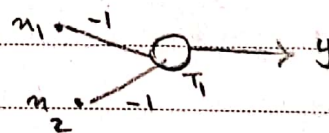
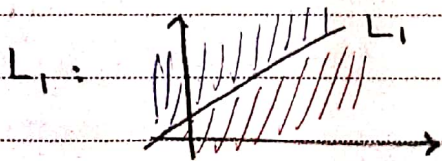
(2) الف: شبکه از 2 لایه
 hidden تشکیل شده به ترتیب
 2 و 8 نود در لایه



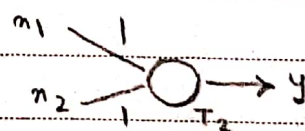
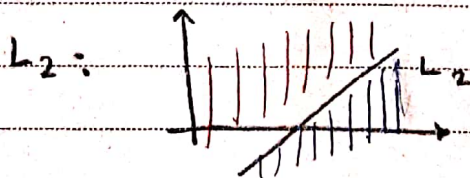
(2) ب:



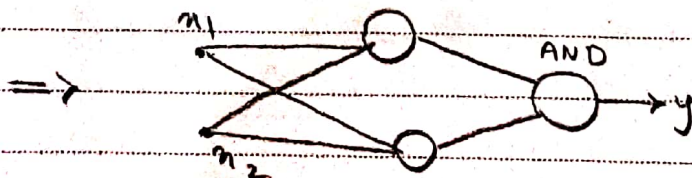
(2) ج: هر کدام از نقاط L_1 و L_2 را میزنان
 با یک perceptron نشان داد و به عبارت دیگر:



$$y = \begin{cases} 1 & -n_1 - n_2 > T_1 \\ 0 & \text{or} \end{cases}$$



$$y = \begin{cases} 1 & n_1 + n_2 > T_2 \\ 0 & \text{or} \end{cases}$$



حاصل AND y_1 و y_2 میشود نمایی لاگت (1)
 و باین نمایی لاگت آبی (0)

(4) الف: چگونگی تعریف تابع relu :

$$\text{relu}(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{و.و.} \end{cases}$$

این تابع حاصل $\sum w_i x_i$ در هر لایه را به عنوان ورودی گرفته و تنها یک شرط روی آن می‌گذارد و به تبع کلیات دیگر این تابع مستقیماً این تابع نیز بسیار آسان است و هنگام update کردن

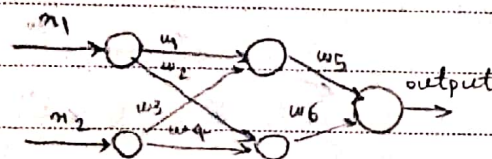
$$\frac{d}{dz} \text{relu}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \end{cases}$$

و تنها با یک محاسبه آسانی داریم:

$$\frac{d}{dz} \tanh(z) = 1 - (\tanh(z))^2$$

$$\frac{d}{dz} \text{sigmoid}(z) = \text{sigmoid}(z) (1 - \text{sigmoid}(z))$$

برای مثال داریم:



$$\frac{\partial \text{error}}{\partial w_5} = \frac{\partial \text{error}}{\partial \text{relu}(w_5 z_1 + w_6 z_2)} \times \frac{\partial \text{relu}(w_5 z_1 + w_6 z_2)}{\partial w_5}$$

(sigmoid, tanh) (در این نقاط مشتق‌ها 1 یا 0 است)

(4) ب: همان شبکه عصبی قبلی را در نظر بگیرید. ضرایب حاصل $|w_5 z_1 + w_6 z_2| > 2$

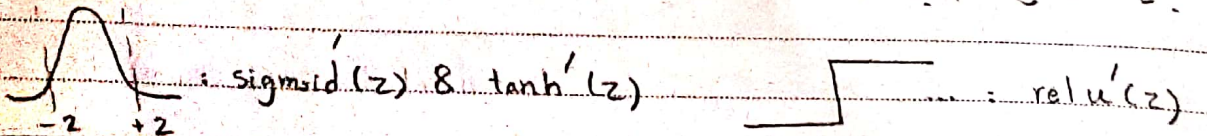
$$\frac{\partial \tanh(w_5 z_1 + w_6 z_2)}{\partial (w_5 z_1 + w_6 z_2)}$$

آنها صفر می‌شود (برای sigmoid هم همینطور است). در نتیجه وزن‌ها

update نمی‌شوند (همان مقدار بقیه را می‌گیرند). به این مشکل vanishing gradient گویند.

دلیل این مشکل $\frac{d}{dz} \tanh(z)$ و $\frac{d}{dz} \text{sigmoid}(z)$ است. اما تابع relu ضرایب $w_i x_i$ متغیر

بسیار بالایی هم بگیرند در آن نقاط هم مشتق مقدار دارد و فرآیند update در آنها اتفاق می‌افتد.



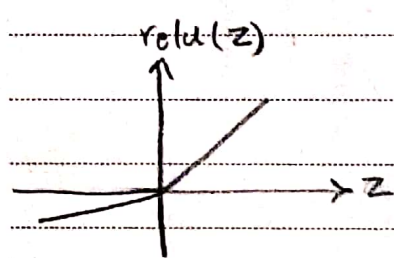
(4) ج: برای classification از sigmoid استفاده می کنیم زیرا bounded است.

(ب) [0, 1] و در cost function داریم:
$$\text{Error}(z) = \frac{1}{2} \sum_i [f(z_i) - \text{output}_i]^2$$

در مسئله کلاسیک باینری خروجی output یا 0 است یا 1، متعلق است که از relu نه

مقادیر بسیار بیشتری از 1 هم می تواند بگیرد استفاده نکنیم. (relu برای لایه آخر regression مناسب است.)

(4) د: تابع
$$\text{leaky_relu}(z) = \begin{cases} z & \text{if } z \geq 0 \\ \alpha z & \text{otherwise} \end{cases}$$



استفاده از leaky-relu فرایندی زیر را دارد:

(1) از ابتدایی که خروجی $\text{relu}(\sum w_i x_i)$ نمایش داده $\sum w_i x_i < 0$ ، صفر است، برخی نودها

فعال نمی شوند و در نتیجه برخی نودها تا آخر صفیری مانند، update نمی شوند. در نتیجه بهتر است

از leaky-relu استفاده کنیم که اطلاعات بیشتری در مورد sample ها می تواند داشته باشیم.

(2) در regression که خروجی مطلوب نمونه نام (output) می تواند متغی هم باشد استفاده از

leaky relu بهتر است.

(5) الف: برای شیر استفاده داریم:

$$H(\text{استفهام}) = -\frac{4}{10} \log\left(\frac{4}{10}\right) - \frac{6}{10} \log\left(\frac{6}{10}\right) = 0.292$$

rest

(ب) برای متغیر فعالیت تنفسی داریم:

$$H(سابقه) = -\frac{3}{4} \log(\frac{3}{4}) - \frac{1}{4} \log(\frac{1}{4}) = +0.244$$

$$H(دارد) = -\frac{1}{4} \log(\frac{1}{4}) - \frac{3}{4} \log(\frac{3}{4}) = +0.244$$

$$H(ندارد) = -0 \log(0) - 1 \log(1) = 0$$

$$\Rightarrow H(تنفسی) = [0.244 \times 4 + 0.244 \times 4 + 0] / 10 = 0.195$$

$$\Rightarrow \text{Info-Gain}(تنفسی) = 0.292 - 0.195 = 0.097$$

متغیر سبب ناری:

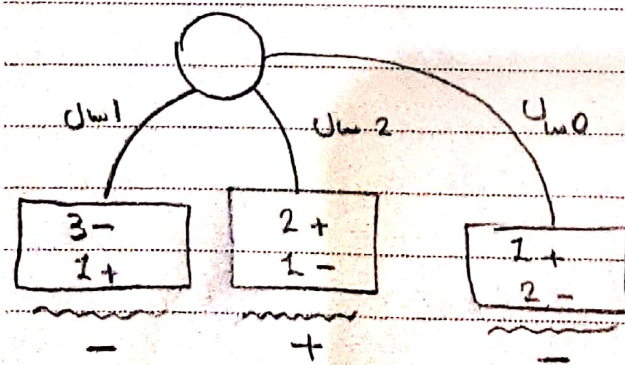
$$H(سبب 1) = -\frac{3}{4} \log(\frac{3}{4}) - \frac{1}{4} \log(\frac{1}{4}) = 0.244$$

$$H(سبب 2) = -\frac{2}{3} \log(\frac{2}{3}) - \frac{1}{3} \log(\frac{1}{3}) = 0.276$$

$$H(سبب 0) = -\frac{1}{3} \log(\frac{1}{3}) - \frac{2}{3} \log(\frac{2}{3}) = 0.276$$

$$\Rightarrow H(سبب ناری) = [4 \times 0.244 + 3 \times 0.276 + 3 \times 0.276] / 10$$

$$\Rightarrow \text{Info-Gain}(سبب ناری) = 0.292 - 0.263 = 0.029$$



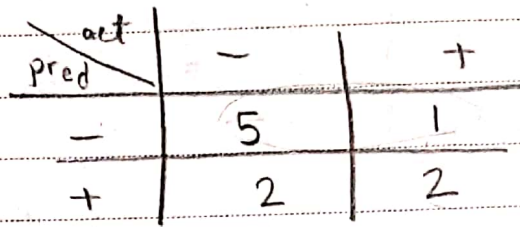
$$\text{accuracy} = \frac{3+2+2}{10} = 0.7$$

$$\text{precision} = \frac{2}{3}$$

$$\text{recall} = \frac{2}{4}$$

$$\text{specificity} = \frac{5}{5+1} = \frac{5}{6}$$

act \ pred	-	+
-	5	2
+	1	2



$$\text{accuracy} = \frac{7}{10}$$

$$\text{recall} = \frac{2}{3}$$

precision = $\frac{2}{4}$

$$\text{specificity} = \frac{5}{7}$$

3 باب: تقریر ادب از لُباسِ ادب (حیدر بالا):

ب: قدری اول از بنا سن اول (دو یا لا):

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & 0 \\ \hline 0 & 1 & 1 & 1 \\ \hline 0 & 1 & 0 & 1 \\ \hline 0 & 1 & 1 & 1 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 0 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 0 & 1 & 1 & 0 & 1 \\ \hline 0 & 1 & 0 & 1 & 0 & 1 \\ \hline 0 & 1 & 1 & 1 & 1 & 1 \\ \hline 0 & 0 & 0 & 1 & 0 & 1 \\ \hline 0 & 1 & 1 & 1 & 1 & 1 \\ \hline \end{array}$$

$\sum w_i x_i + \text{bias}$

1	3
2	6

maxpooling

6

$$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 \end{bmatrix}$$

مقومير ديم از نلاس اول (راست بالا) :

Subject:

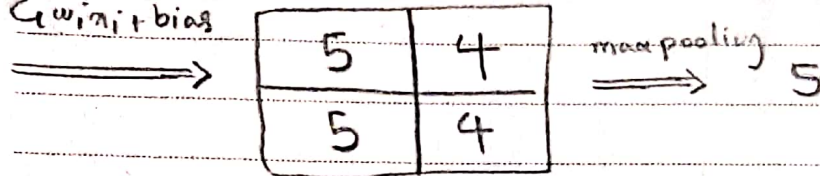
Year:

Month:

Date:

()

$\sum w_i x_i + bias$



0	0	0	0
0	1	0	1
1	1	1	0
0	1	0	0

x

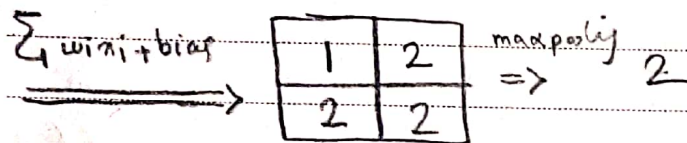
1	1	1
1	0	1
1	1	1

=

تصویر اول از لایه دوم (پایین جبهه):

0	0	0	0	0	0
0	0	0	1	0	1
1	1	1	1	1	0
0	1	0	1	0	1
1	0	1	1	0	0
0	1	0	1	0	0

$\sum w_i x_i + bias$



0	1	0	0
0	1	0	0
0	1	1	0
0	1	0	0

x

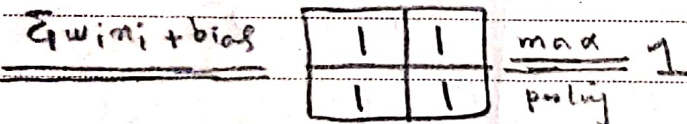
1	1	1
1	0	1
1	1	1

=

تصویر دوم از لایه دوم (پایین جبهه):

0	1	0	1	0	0
0	0	0	1	0	0
0	1	1	1	0	0
0	1	0	1	0	0
0	0	1	1	0	0
0	1	0	1	0	0

$\sum w_i x_i + bias$



ج) میتری که جنبش به ری هر تصویر ای شده، یک حفره میزند بود که توسط تعداد

پیکسل قرمز معامره شده بود (feature لایه 1). حال با convolve و max pooling کردن

بجای هر عکس داشته آیا بخشی از عکس سبز این فیلتر هست؟ (از یابی میشود با فیلتر)

هانفرد که دیده میشود این بجایست برا عکسها لایه اول زیاد است، با براین

با قرار دادن یک threshold می توان عکسها لایه اول دوم را از هم جدا کرد.