



Assignment NO.1 Solutions

Neural Networks | Fall 1400 | Dr.Mozayani

Teacher Assistants:

Mohammad Hossein Khojaste

Amir Hossein Aminimehr

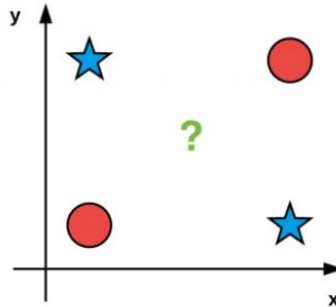
Student name : **Amin Fathi**

Student id : **400722102**

Problem 1.a

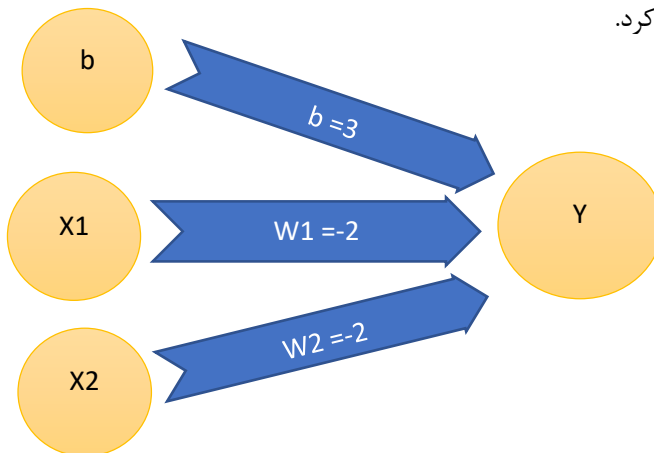
این تابع منطقی درواقع یک تابع XNOR می باشد ، که آن را نمیتوان با یک نورون پیاده سازی کرد ، چرا که برای پیاده سازی آن می بایست نقاط $(0,0)$ و $(1,1)$ را در یک کلاس قرار داد و نقاط $(0,1)$ و $(1,0)$ را در یک کلاس قرار داد که این عمل با کشیدن یک خط در صفحه مختصات ممکن نیست.

x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1



Problem 1.b

این تابع منطقی در واقع یک NAND می باشد ، و میتوان آن را با یک نورون ساده پیاده سازی کرد ، چرا که برای پیاده سازی آن فقط میبایست نقاط $(0,0)$ و $(1,0)$ و $(0,1)$ را توسط یک خط از نقطه $(1,1)$ جدا کرد.



x_1	x_2	y
0	0	1
0	1	1
1	0	1
1	1	0

$w_1 = -2$ و $w_2 = -2$ در نظر میگیریم و مقدار Threshold را هم برابر با ۳ در نظر میگیریم که در آن صورت به ازای هر ۴ حالت معادله برقرار است . (کفایت شرایط زیر برقرار باشد) (مقدار threshold برابر منفی مقدار b است) در واقع باید ضرایبی انتخاب شوند برای w_1 و w_2 که وقتی $w_1 * x_1 + w_2 * x_2$ شد به ازای جفت x_1 ها و x_2 های $(0,0)$ و $(0,1)$ و $(1,0)$ مقدارش از threshold بیشتر شود ، و برای $(1,1)$ از آن کمتر شود

$$w_1 \times 0 + w_2 \times 0 - b \geq 0 \Rightarrow b \geq 0$$

$$w_1 \times 0 + w_2 \times 1 - b \geq 0 \Rightarrow b \geq -w_2$$

$$w_1 \times 1 + w_2 \times 0 - b \geq 0 \Rightarrow b \geq -w_1$$

$$w_1 \times 1 + w_2 \times 1 - b \leq 0 \Rightarrow b \leq -(w_1 + w_2)$$

Problem 2

X1	X2	bias	target
+1	+1	+1	+1
+1	-1	+1	+1
-1	+1	+1	+1
-1	-1	+1	-1

$$Y = b + \sum x_i * w_i$$

$$b(\text{new}) = b(\text{old}) + \eta(d - y)$$

$$w_i(\text{new}) = w_i(\text{old}) + \eta(d - y)x_i$$

$$w_1 = w_2 = b = 0.2$$

$$\eta = 0.1$$

step1:

$$y = 0.2 + 1*0.2 + 1*0.2 = 0.6$$

$$b(\text{new}) = 0.2 + 0.1*(1-0.6) = 0.24$$

$$w_1(\text{new}) = 0.2 + 0.1*(1-0.6)*1 = 0.24$$

$$w_2(\text{new}) = 0.2 + 0.1*(1-0.6)*1 = 0.24$$

$$y = 0.24 + 1*0.24 - 1*0.24 = 0.24$$

$$b(\text{new}) = 0.24 + 0.1*(1-0.24) = 0.316$$

$$w_1(\text{new}) = 0.24 + 0.1*(1-0.24)*1 = 0.316$$

$$w_2(\text{new}) = 0.24 + 0.1*(1-0.24)*-1 = 0.164$$

$$y = 0.316 - 1*0.316 + 1*0.164 = 0.164$$

$$b(\text{new}) = 0.316 + 0.1*(1-0.164) = 0.3996$$

$$w_1(\text{new}) = 0.316 + 0.1*(1-0.164)*-1 = 0.2324$$

$$w_2(\text{new}) = 0.164 + 0.1*(1-0.164)*1 = 0.2476$$

$$y = 0.3996 - 1 \cdot 0.2324 - 1 \cdot 0.2476 = -0.0804$$

$$b(\text{new}) = 0.3443 + 0.1 \cdot (-1 + 0.0804) = 0.30764$$

$$w1(\text{new}) = 0.2324 + 0.1 \cdot (-1 + 0.0804) \cdot -1 = 0.32436$$

$$w2(\text{new}) = 0.2476 + 0.1 \cdot (-1 + 0.0804) \cdot -1 = 0.33956$$

Problem 3.a

اولا، مقادیر خروجی یک پرسپترون به دلیل تابع انتقال (transformation) ، می‌توانند تنها یکی از دو مقدار (۰ یا ۱) را خروجی دهد.

دوما ، پرسپترون فقط برای مجموعه تفکیک پذیری خطی (linear seperable) ها همگراست و می‌تواند کلاس بندی را انجام دهد

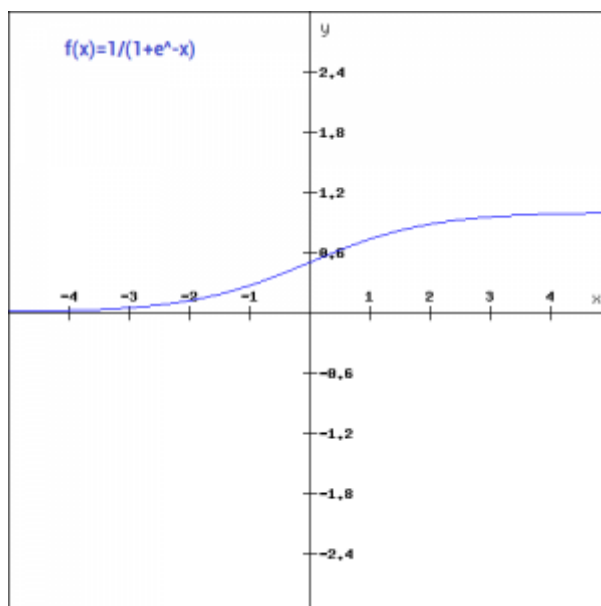
سوما ، پرسپترون توانایی تضمین کیفیت جواب را ندارد .

Problem 3.b

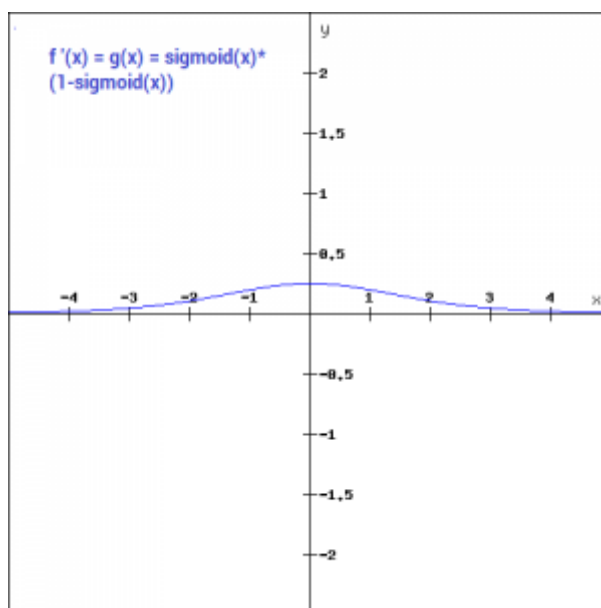
Sigmoid

تابع سیگموئید (Sigmoid function) یک تابع ریاضی است که شکلی مشابه با حرف S در زبان انگلیسی دارد. به طور کلی به توابعی که شکل آن‌ها مشابه با حرف S است تابع‌های سیگموئید یا دارای خمیدگی سیگموئید می‌گویند.

یک تابع غیر خطی است و بنابراین خروجی غیر خطی می‌دهد ، خروجی در حد فاصل ۰ و ۱ است. (مطابق شکل) ، برای کلاس بندی باینری مناسب است .



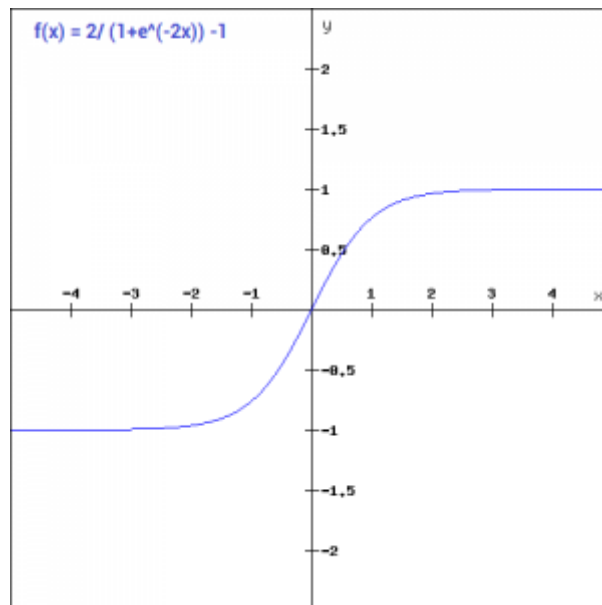
مشتق تابع sigmoid برابر است با شکل زیر :



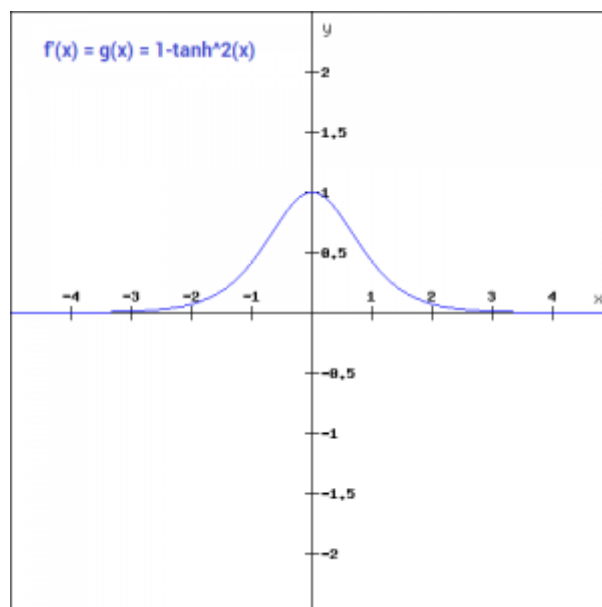
همانطور که در شکل مشاهده می شود مقادیر گرادیان (مشتق) برای محدوده -3 و 3 قابل توجه است، اما نمودار در مناطق دیگر بسیار صاف تر می شود و دارای گرادیان های بسیار کوچک خواهد بود.

Tanh

بسیار شبیه تابع سیگموئید است. تنها تفاوت این است که در اطراف مبدا مختصات متقارن است و محدوده مقادیر در آن از -1 تا 1 است. بنابراین ورودی های لایه های بعدی همیشه علامت یکسانی ندارند. تابع \tanh به این صورت تعریف می شود:



همانند تابع sigmoid این تابع هم پیوسته و در تمام نقاط مشتق پذیر است ، ولی نسبت به سیگنویید از شیب بالاتری برخوردار است.



softmax

اغلب به عنوان ترکیبی از چند سیگموئید توصیف می شود. و می تواند برای مسائل طبقه بندی چند کلاسه استفاده شود. این تابع احتمال تعلق یک نقطه داده به هر کدام از کلاس ها را برمی گرداند. د

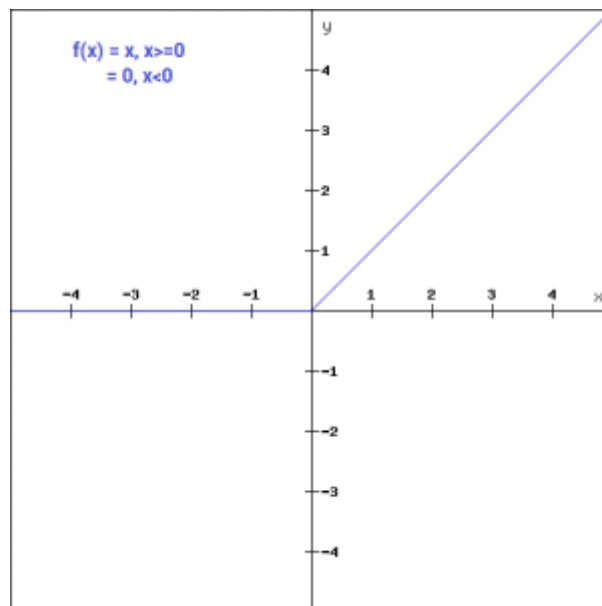
$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

در حین ساختن یک شبکه برای یک مسئله چند کلاسه، لایه خروجی به تعداد کلاس های هدف نرون خواهد داشت. به عنوان مثال اگر ما سه کلاس داریم، سه نرون در لایه خروجی وجود خواهد داشت. فرض کنید خروجی نرون ها را به صورت $[0.75, 0.9, 0.2]$ دریافت کرده اید.

با اعمال تابع softmax بر روی این مقادیر، نتیجه زیر را دریافت خواهیم کرد - $[0.27, 0.31, 0.42]$. اینها نشان دهنده احتمال تعلق نقاط داده (دیتا) به هر کدام از کلاس هاست که همه جمع همه آنها برابر ۱ است.

RELU

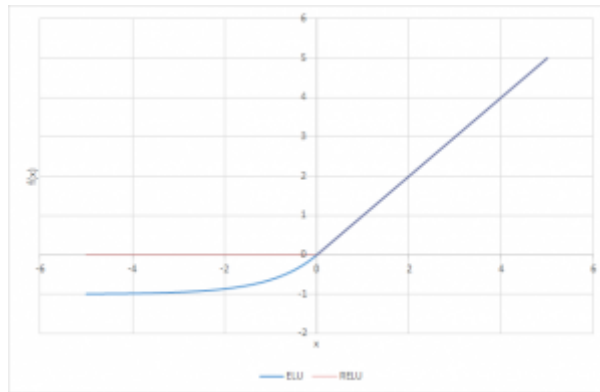
یک تابع فعال سازی غیر خطی است (پر کاربرد ترین) و مزیت آن در این است که همه نرون ها را به طور همزمان فعال نمی کند، در صورتی که ورودی این تابع کوچکتر از صفر باشد، فعال نخواهد بود.



همانطور که در تصویر مشتق RELU مشاهده میشود، مقدار گرادیان این تابع به ازای ورودی های منفی صفر است که این اتفاق در backpropagation میتواند منجر به به روز نشدن تعدادی از وزن ها و در نتیجه غیر فعال شدن تعدادی از نرون ها شود.

ELU

نوعی از ReLU (Rectified Linear Unit) است که شیب قسمت منفی را تغییر می دهد (لگارتیمی می کند)



Problem3.b

معایب SGD

در روش کاهش آماری گرادیان (SGD) برای هر نمونه بروزرسانی پارامترها انجام می شود. تابع هزینه در این روش بسیار ناهموار است و نوسانات شدیدی خواهد داشت زیرا داده ها به طور مکرر به روز می شوند و ممکن است به حداقل های محلی بهتر بپردازند. به همین دلیل برای توابع غیر محدب مناسب نیست

SGD with momentum

مومنتوم اینرسی یک شی را در هنگام حرکت شبیه سازی می کند، یعنی مسیر به روز رسانی قبلی تا حدودی در طول به روز رسانی حفظ می شود، و گرادیان به روز رسانی فعلی برای بهینه سازی جهت به روز رسانی نهایی استفاده می شود. ثبات تا حد معینی افزایش می یابد، بنابراین یادگیری سریعتر است، و همچنین توانایی خلاص شدن از کمینه محلی را دارد.

در فرمول پایین می بینیم که می توانیم تصمیم بگیریم به کدام داده ها وزن بدهیم، که در زمان خاصی بسته به مقدار بتا می آیند.

بنابراین یادگیری از نقاط داده قبلی مفید خواهد بود و جهت خود را حفظ می کند و در نتیجه نویز را حذف می کند.

به گفته محققان، بتای بهینه ۰/۹۵ در نظر گرفته شده است.

Normal Weight Updation Formula in Gradient Descent

$$W_t = W_{t-1} - \eta \frac{dL}{dW_{t-1}}$$

Exponential Weighted Average for momentum

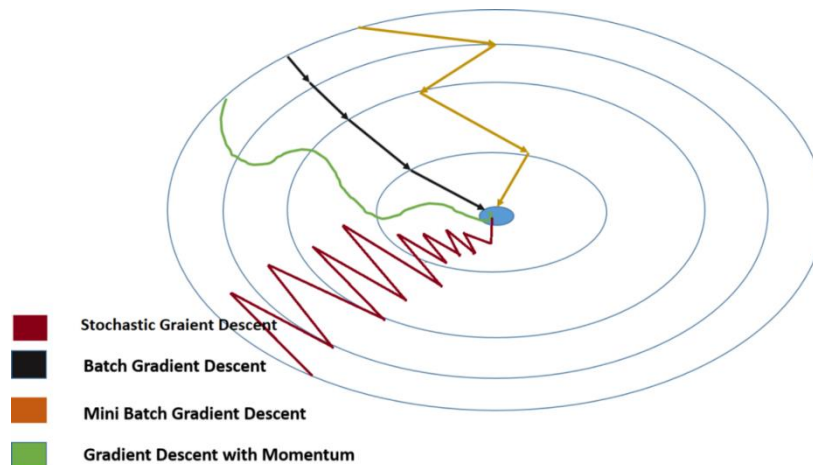
$$W_t = W_{t-1} - \eta Vdw \quad \eta = \text{learning rate}$$

$$b_t = b_{t-1} - \eta Vdb \quad b_t = \text{bias}$$

$$Vdw_t = \beta * Vdw_{t-1} + (1-\beta) * \frac{dL}{dW_{t-1}}$$

$$Vdb_t = \beta * Vdb_{t-1} + (1-\beta) * \frac{dL}{db_{t-1}}$$

If $\beta = 0$ formula will be same as used in gradient descent.



در تصویر بالا مقایسه نحوه رفتار SGD , BGD , SGD with Momentum را مشاهده میکنید

RMS-Prop

یا همان Root Mean Square Prop یک تکنیک بهینه سازی مبتنی بر گرادیان است که در آموزش شبکه های عصبی استفاده می شود جفری هینتون پیشنهاد شد.

با استفاده از میانگین متحرک گرادیان مجذور (moving average of squared

gradients) به نرمال کردن گرادیان می پردازد، این نرمال سازی اندازه ی طول گام را متعادل می کند به معنی که گام را برای

جاهایی که گرادیان در آن زیاد است برای جلوگیری از انفجار کاهش داده و در جایی که گرادیان در آن کم است برای

جلوگیری از ناپدید شدن طول گام را افزایش می دهد. در واقع RMS-Prop از نرخ یادگیری انطباقی adaptive

learning rate به جای نرخ یادگیری ثابت به عنوان یک هایپرپارامتر استفاده می کند. این بدان معناست که نرخ

یادگیری در این روش در طول زمان در حال تغییر است. این بروزرسانی از روابط زیر به دست می آید:

$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$

منابع :

<https://medium.com/analytics-vidhya/a-complete-guide-to-adam-and-rmsprop-optimizer-be83d4e02f75>

<https://towardsdatascience.com/optimizers-c88694e09311>