

بنام خدا  
تمرین سری سوم یادگیری ماشین (شبکه‌های عصبی)

مریم رضوانی  
شماره دانشجویی: ۹۹۲۱۱۶۰۰۱۹

(1)

**الف) گزینه A؛** تکنیک dropout یک روش منظم سازی است که تقریباً معادل آموزش تعداد زیادی شبکه عصبی با معماری های مختلف به صورت موازی است. در این روش، در هر مرحله آموزش، برخی از خروجی های لایه ها به صورت تصادفی صفر می شوند یا «حذف می شوند». این کار باعث می شود که شبکه به نوروں های خاصی وابسته نشود و تمام نوروں ها بتوانند بهتر تعمیم دهند.

**Bagging:** این تکنیک شامل آموزش چندین شبکه عصبی مستقل با داده های زیرمجموعه ای از داده های آموزشی است. سپس، پیشبینی های این شبکه ها با هم ترکیب می شوند تا یک پیشبینی نهایی ایجاد کنند. این تکنیک باعث می شود که شبکه ها از تنوع بیشتری برخوردار شوند و کمتر به نویز داده ها حساس باشند

**بنابراین، می توان گفت که تکنیک dropout مشابه تکنیک bagging است،** با این تفاوت که در dropout، فقط یک شبکه عصبی آموزش داده می شود که در هر مرحله، بخشی از نوروں های آن حذف می شوند. این کار باعث می شود که شبکه عصبی مانند یک مجموعه از شبکه های عصبی با معماری های مختلف رفتار کند. این روش دارای مزایایی مانند سرعت بالا، حافظه کمتر و پیاده سازی آسان است

**ب)** تفاوت اصلی آنها در این است که batch gradient descent از تمام داده های آموزشی برای محاسبه گرادیان در هر مرحله استفاده می کند، در حالی که stochastic gradient descent فقط از یک نمونه یا یک زیرمجموعه از داده های آموزشی برای محاسبه گرادیان در هر مرحله استفاده می کند.

این تفاوت باعث می شود که batch gradient descent و stochastic gradient descent دارای مزایا و معایب مختلفی باشند. برخی از این مزایا و معایب عبارتند از:

- batch gradient descent دقیق تر و پایدارتر است، زیرا از تمام داده های آموزشی برای محاسبه گرادیان استفاده می کند و به سمت کمینه محلی یا سراسری تابع هزینه حرکت می کند. Stochastic gradient descent ناپایدارتر و نویزی تر است، زیرا ممکن است در هر مرحله به سمت های مختلفی حرکت کند.

- Stochastic gradient descent سریع تر و کم حافظه تر است، زیرا فقط نیاز به محاسبه گرادیان برای یک نمونه از داده های آموزشی در هر مرحله دارد و نیازی به ذخیره سازی تمام داده های آموزشی ندارد. batch gradient descent کندتر و حافظه برتر است، زیرا نیاز به محاسبه گرادیان برای تمام داده های آموزشی در هر مرحله دارد و نیاز به ذخیره سازی تمام داده های آموزشی دارد.

- Stochastic gradient descent مناسب تر برای بهینه سازی توابع هزینه غیر محدب است، زیرا می تواند از کمینه های محلی خارج شود و به سمت کمینه سراسری حرکت کند. batch gradient descent مناسب تر برای بهینه سازی توابع هزینه محدب است، زیرا می تواند به سرعت به کمینه محلی یا سراسری برسد.

ج) مزیت استفاده از چندین لایه در یک شبکه عصبی چندلایه این است که شبکه می‌تواند روابط پیچیده و غیرخطی بین ورودی‌ها و خروجی‌ها را یاد بگیرد و داده‌هایی را متمایز کند که به صورت خطی قابل تفکیک نیستند. برای مثال، شبکه عصبی چندلایه می‌تواند مسئله XOR را حل کند که پرسپترون تک لایه نمی‌تواند.

استفاده از چندین لایه در یک شبکه عصبی چندلایه می‌تواند به شبکه امکان دهد تا از ویژگی‌های سطح بالاتر و انتزاعی‌تری از داده‌ها استفاده کند. برای مثال، در یک شبکه عصبی چندلایه که برای شناسایی تصاویر آموزش داده شده است، لایه‌های اولیه می‌توانند ویژگی‌های ساده مانند لبه‌ها و رنگ‌ها را تشخیص دهند، لایه‌های میانی می‌توانند ویژگی‌های پیچیده‌تر مانند شکل‌ها و الگوها را تشخیص دهند، و لایه‌های آخر می‌توانند ویژگی‌های بسیار انتزاعی مانند چهره‌ها و اشیاء را تشخیص دهند.

استفاده از چندین لایه در یک شبکه عصبی چندلایه می‌تواند به شبکه کمک کند تا از مشکلاتی مانند محو‌گرادیان و انفجار گرادیان جلوگیری کند.

استفاده از چندین لایه در یک شبکه عصبی چندلایه می‌تواند به شبکه امکان دهد تا از مزایای روش‌های مختلف یادگیری استفاده کند.

اما استفاده از چندین لایه در یک شبکه عصبی چندلایه ممکن است آموزش شبکه را سریع‌تر نکند. بلکه ممکن است آموزش شبکه را کندتر و پیچیده‌تر کند. زیرا با افزایش تعداد لایه‌ها، تعداد پارامترها و محاسبات شبکه نیز افزایش می‌یابد و نیاز به بیشترین داده‌ها و منابع محاسباتی دارد. بنابراین، باید تعادل مناسبی بین تعداد لایه‌ها و کیفیت شبکه ایجاد کرد.

د) تفاوت بین یک Feedforward Network سنتی و یک Recurrent Network در جهت انتقال داده‌ها، حافظه نوروها، پیچیدگی محاسباتی و مناسب بودن برای مسائل مختلف است.

یک Recurrent Network یک نوع شبکه عصبی مصنوعی است که در آن ورودی‌ها به صورت دنباله‌ای از مقادیر متغیر در طول زمان انتقال می‌یابند. (وابستگی زمانی لحاظ می‌شود). این شبکه از یک لایه ورودی، یک یا چند لایه پنهان و یک لایه خروجی تشکیل شده است. هر لایه شامل چندین نرون است که با وزن‌ها و بایاس‌هایی به لایه بعدی متصل هستند. هر نرون یک تابع فعال‌سازی را روی جمع وزن‌دار ورودی‌های خود و خروجی‌های خود از مرحله قبلی اعمال می‌کند و خروجی خود را به لایه بعدی و به خودش در مرحله بعدی می‌فرستد. این شبکه می‌تواند برای یادگیری روابط بین دنباله‌های ورودی و خروجی در مسائل مختلف مانند ترجمه زبان، تبدیل گفتار به متن و کنترل رباتیک استفاده شود.

یک Feedforward Network سنتی یک نوع شبکه عصبی مصنوعی است که در آن داده‌ها یا ورودی‌ها فقط در یک جهت انتقال می‌یابند. این شبکه از یک لایه ورودی، یک یا چند لایه پنهان و یک لایه خروجی تشکیل شده است. هر لایه شامل چندین نرون است که با وزن‌ها و بایاس‌هایی به لایه بعدی متصل هستند. هر نرون یک تابع فعال‌سازی را روی جمع وزن‌دار ورودی‌های خود اعمال می‌کند و خروجی خود را به لایه بعدی می‌فرستد. این شبکه می‌تواند برای یادگیری روابط بین ورودی‌ها و خروجی‌ها در مسائل مختلف مانند شناسایی الگو، تشخیص گفتار و تشخیص حروف استفاده شود.

ه) **نرخ یادگیری** مقداری است که نشان می‌دهد که چقدر وزن‌ها و بایاس‌های شبکه در هر تکرار به‌روزرسانی می‌شوند. نرخ یادگیری باید به اندازه کافی بزرگ باشد تا شبکه بتواند به سرعت یاد بگیرد و به نقطه کمینه تابع هزینه برسد. اما اگر نرخ یادگیری بیش از حد بزرگ باشد، ممکن است شبکه ناپایدار شود و از نقطه کمینه عبور کند. برعکس، اگر نرخ یادگیری بیش از حد کوچک باشد، ممکن است شبکه به آهستگی یاد بگیرد و در نقاط کمینه محلی گیر کند. بنابراین، انتخاب نرخ یادگیری مناسب برای بهینه‌سازی شبکه عصبی بسیار مهم است.

**ممنتم** یک پارامتر است که نشان می‌دهد که چقدر شبکه از گرادیان‌های قبلی برای به‌روزرسانی وزن‌ها و بایاس‌ها استفاده می‌کند. ممنتم باعث می‌شود که شبکه در جهتی که گرادیان‌ها همگرا هستند، سرعت بیشتری داشته باشد و از نوسانات ناخواسته جلوگیری کند. ممنتم می‌تواند به شبکه کمک کند تا از نقاط کمینه محلی خارج شود و به نقطه کمینه سراسری

نزدیک‌تر شود. ممنتیم باید به اندازه کافی بزرگ باشد تا شبکه بتواند از گرادیان‌های قبلی استفاده کند و از نقاط کمینه محلی عبور کند. اما اگر ممنتیم بیش از حد بزرگ باشد، ممکن است شبکه نتواند در نقطه کمینه متوقف شود و از آن فراتر رود. بنابراین، انتخاب ممنتیم مناسب برای بهینه‌سازی شبکه عصبی بسیار مهم است.

و) تفاوت بین *multiclass classification* و *multiple classification* در تعداد کلاس‌هایی که هر نمونه می‌تواند به آنها تعلق داشته باشد، است. بعنوان مثال: در *multiclass classification*، هر نمونه به یک و تنها یک کلاس از چندین کلاس موجود اختصاص می‌یابد. برای مثال، اگر داده‌ها شامل سه کلاس گل‌های آیریس باشند، هر نمونه فقط می‌تواند به یکی از کلاس‌های *setosa*، *versicolor* یا *virginica* تعلق داشته باشد و در *multiple classification*، هر نمونه می‌تواند به چندین کلاس از چندین کلاس موجود اختصاص یابد. برای مثال، اگر داده‌ها شامل چندین کلاس برچسب مانند رنگ، شکل، اندازه و جنس باشند، هر نمونه می‌تواند به چندین برچسب متفاوت تعلق داشته باشد.

ن) **محو گرادیان:** همانطور که الگوریتم پس انتشار از لایه خروجی به سمت لایه ورودی به سمت پایین (یا به عقب) پیش می‌رود، گرادیان‌ها اغلب کوچکتر و کوچکتر می‌شوند و به صفر نزدیک می‌شوند که در نهایت وزن لایه‌های اولیه یا پایین تقریباً بدون تغییر باقی می‌ماند. در نتیجه، گرادیان کاهشی هرگز به حد مطلوب همگرا نمی‌شود. این به عنوان مشکل محو گرادیان (*Vanish gradient*) شناخته می‌شود.

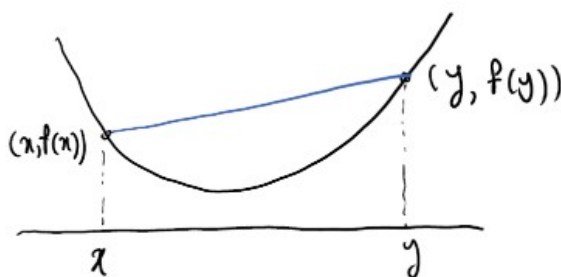
**انفجار گرادیان:** برعکس، در برخی موارد، با پیشرفت الگوریتم پس انتشار، گرادیان‌ها بزرگتر و بزرگتر می‌شوند. این به نوبه خود باعث به‌روز رسانی وزن بسیار بزرگ می‌شود و باعث واگرایی گرادیان کاهشی می‌شود. این به عنوان مشکل گرادیان انفجاری (*Exploding gradient*) شناخته می‌شود.

ی) عبارت درست است؛ زیرا در یک مسأله دسته‌بندی چند کلاسه، لایه خروجی شبکه باید به اندازه کلاس‌های موجود در داده‌ها نورون داشته باشد. هر نورون در لایه خروجی می‌تواند احتمال تعلق یک نمونه به یک کلاس را محاسبه کند. و به ما بستگی دارد که مسأله را چند کلاسه طراحی کنیم برای هر کلاس یک نورون در لایه خروجی قرار می‌دهیم.

## 2) تعریف تابع محدب:

**Definition 1.** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if its domain is a convex set and for all  $x, y$  in its domain, and all  $\lambda \in [0, 1]$ , we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$



مجموعه محدب مجموعه‌ای است که در آن هر پاره خطی که دو نقطه از مجموعه را به هم متصل می‌کند نیز به طور کامل در مجموعه قرار می‌گیرد. این بدان معنی است که مجموعه‌های محدب ساده هستند و ساختار مشخصی دارند.

یک سطح تصمیم خطی را در نظر بگیرید که با یک ابر صفحه تعریف شده با معادله  $w^T x + w_0 = 0$ ، که در آن  $w$  بردار وزن و  $w_0$  عبارت بایاس است، در نظر بگیرید. فرض کنید  $x_1$  و  $x_2$  دو نقطه در یک کلاس باشند، و اجازه دهید  $x_3$  یک میانگین وزنی  $x_1$  و  $x_2$  باشد:

$$x_3 = \lambda x_1 + (1 - \lambda)x_2, \text{ where } 0 \leq \lambda \leq 1$$

با جایگزینی  $x_3$  در معادله سطح تصمیم، می‌بینیم:

$$w^T x_3 + w_0 = \lambda w^T x_1 + (1 - \lambda)w^T x_2 + w_0 = 0$$

از آنجایی که  $w^T x_1$  و  $w^T x_2$  هر دو برابر 0 هستند، داریم:

$$\lambda w^T x_1 + (1 - \lambda)w^T x_2 = 0$$

این به این معنی است که  $x_3$  نیز روی سطح تصمیم قرار دارد، و بنابراین، هر نقطه دیگری در پاره خطی که  $x_1$  و  $x_2$  را به هم متصل می‌کند نیز باید روی سطح تصمیم قرار گیرد. این ثابت می‌کند که سطوح تصمیم‌گیری خطی محدب هستند.

شبکه‌های عصبی و جداسازی ناحیه محدب نورون‌ها، واحدهای محاسباتی اساسی شبکه‌های عصبی، همچنین می‌توانند برای جداسازی مناطق محدب در فضای ورودی استفاده شوند. این به این دلیل است که نورون‌ها می‌توانند انواع توابع غیرخطی را پیاده‌سازی کنند که می‌توانند فضای ورودی را به فضایی با ابعاد بالاتر که در آن سطح تصمیم‌گیری خطی است، ترسیم کنند.

نورون‌ها می‌توانند نواحی محدب را با استفاده از سطوح تصمیم خطی یا غیرخطی که فضای ورودی را به مناطق مختلف تقسیم می‌کنند، جدا کنند. نورون‌ها می‌توانند سطوح تصمیم‌گیری خطی یا غیرخطی را با استفاده از توابع فعال‌سازی مختلف ایجاد کنند که مجموع وزنی ورودی‌های آنها را به خروجی تبدیل می‌کند. به عنوان مثال، یک نورون با یک تابع فعال‌سازی خطی می‌تواند یک سطح تصمیم‌گیری خطی، مانند یک ابر صفحه، که محدب است ایجاد کند. یک نورون با تابع فعال‌سازی سیگموئید می‌تواند یک سطح تصمیم‌گیری غیرخطی، مانند منحنی سیگموئید، که همچنین محدب است، ایجاد کند. با این حال، یک نورون با تابع فعال‌سازی درجه دوم می‌تواند یک سطح تصمیم غیرخطی، مانند سهمی، که محدب نیست ایجاد کند.

برای نشان دادن این موضوع، یک شبکه عصبی ساده با یک نورون و یک تابع فعال‌سازی سیگموئید را در نظر بگیرید. نورون را می‌توان با معادله نشان داد:

$$f(x) = \sigma(wx + b), \text{ که } \sigma \text{ تابع سیگموئید و } w \text{ و } b \text{ وزن‌ها و بایاس نورون هستند.}$$

تابع سیگموئید به صورت زیر تعریف می‌شود:

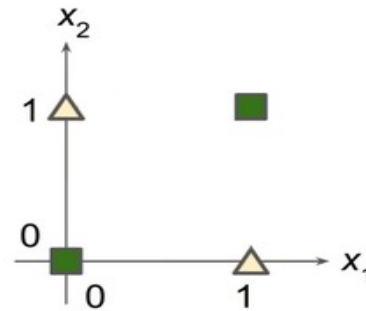
$$\sigma(x) = 1 / (1 + e^{-x})$$

تابع سیگموئید مقدار ورودی  $x$  را به مقداری بین 0 و 1 تقسیم می‌کند. این بدان معناست که نورون می‌تواند فضای ورودی را به دو ناحیه جدا کند: یکی که  $f(x) < 0.5$  و دیگری  $f(x) > 0.5$ . اگر بردار وزن  $w$  با دقت انتخاب شود، سطح تصمیم تشکیل شده توسط نورون می‌تواند محدب باشد. این را می‌توان با آموزش نورون با مجموعه‌ای از نقاط داده که نشان دهنده دو کلاسی است که باید از هم جدا شوند به دست آورد. فرآیند آموزش، وزن‌های  $w$  نورون را تنظیم می‌کند تا خطا بین برچسب‌های کلاس پیش‌بینی شده و برچسب‌های کلاس واقعی نقاط داده به حداقل برسد. به طور خلاصه، سطوح تصمیم‌گیری محدب هستند و می‌توان از نورون‌ها برای تقریب آن‌ها استفاده کرد. این باعث می‌شود نورون‌ها برای کارهایی مانند طبقه‌بندی و خوشه‌بندی، که اغلب شامل جداسازی داده‌ها به کلاس‌های مختلف است، مناسب باشند.

## جدول XOR

Input 1	Input 2	Output
0	0	0
0	1	1
1	0	1
1	1	0

## نمودار XOR

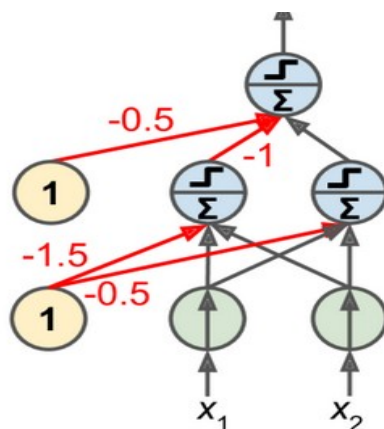


طبق نمودار رسم شده می بینیم که داده ها بصورت خطی قابل تفکیک نیست و برای جداسازی آن ها نیاز به (MLP) پرسپترون چندلایه داریم؛ برای حل این مشکل، یک لایه اضافی به پرسپترون خود اضافه می کنیم. ما این لایه اضافی را لایه پنهان می نامیم.

لایه اول ورودی های ما دارای ۲ نورون یعنی  $x_1, x_2$  هستند که یک نورون اضافی بایاس هم در هر مرحله داریم. لایه دوم که لایه پنهان نامیده می شود شامل ۲ نورون  $y_1, y_2$  است که برای تبدیل مسأله XOR به یک تابع خطی اضافه شده اند. که کار AND و OR که با تابع خطی قابل حل است را انجام می دهند. پس در این مسأله ما ۳ لایه (ورودی (۳ نورون)، پنهان (۳ نورون) و خروجی (۱ نورون) داریم.

$$y_j = f\left(\sum_{i=0}^m x_i w_{i,j}\right)$$

معادله ۱: معادله خروجی پرسپترون



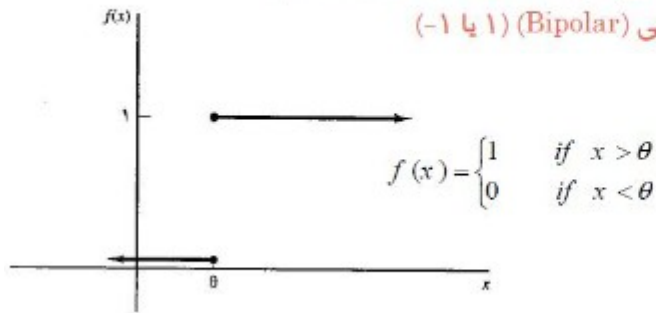
شکل ۴: یک شبکه عصبی mlp که می تواند مسئله XOR را حل کند.

اتصالات خاکستری در شکل وزن  $w=1$  دارند.  $y_1$  یک گیت AND و  $y_2$  یک گیت OR است. محاسبه شکل بالا را در زیر می توان مشاهده کرد. که خروجی در هر مرحله توسط تابع Step مشخص می شود.

## تابع پله‌ای دودویی (Step Function)

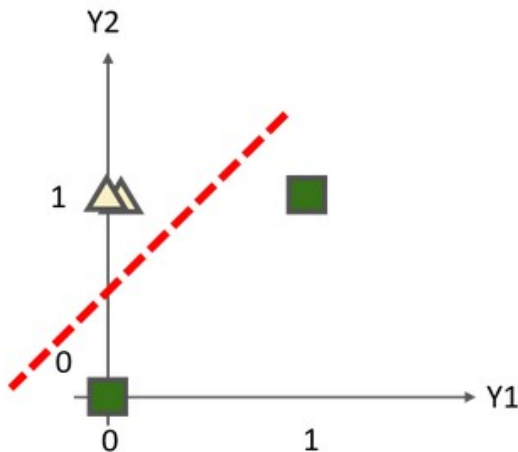
○ تابع آستانه (Threshold Function) یا تابع هویساید (Heaviside Function)

○ خروجی = سیگنال دودویی (۱ یا ۰) یا دوقطبی (۱ یا -۱)



X1	X2	Y1	Y2	Output
0	0	Step(-1.5*1 + 0*1 + 0*1) = 0	Step(-0.5*1 + 0*1 + 0*1) = 0	Step((-0.5*1 - 1*0 + 1*0) = 0
0	1	Step(-1.5*1 + 0*1 + 1*1) = 0	Step(-0.5*1 + 0*1 + 1*1) = 1	Step((-0.5*1 - 1*0 + 1*1) = 1
1	0	Step(-1.5*1 + 1*1 + 0*1) = 0	Step(-0.5*1 + 1*1 + 0*1) = 1	Step((-0.5*1 - 1*0 + 1*0) = 1
1	1	Step(-1.5*1 + 1*1 + 1*1) = 1	Step(-0.5*1 + 1*1 + 1*1) = 1	Step(-0.5*1 - 1*1 + 1*1) = 0

یعنی دو نورون Y1 و Y2 خروجی‌ها را به فضای جدیدی بردند، طوری‌که تفکیک داده‌ها راحت است. در شکل زیر تفکیک داده‌ها مطابق خروجی جدید مشاهده می‌شود:



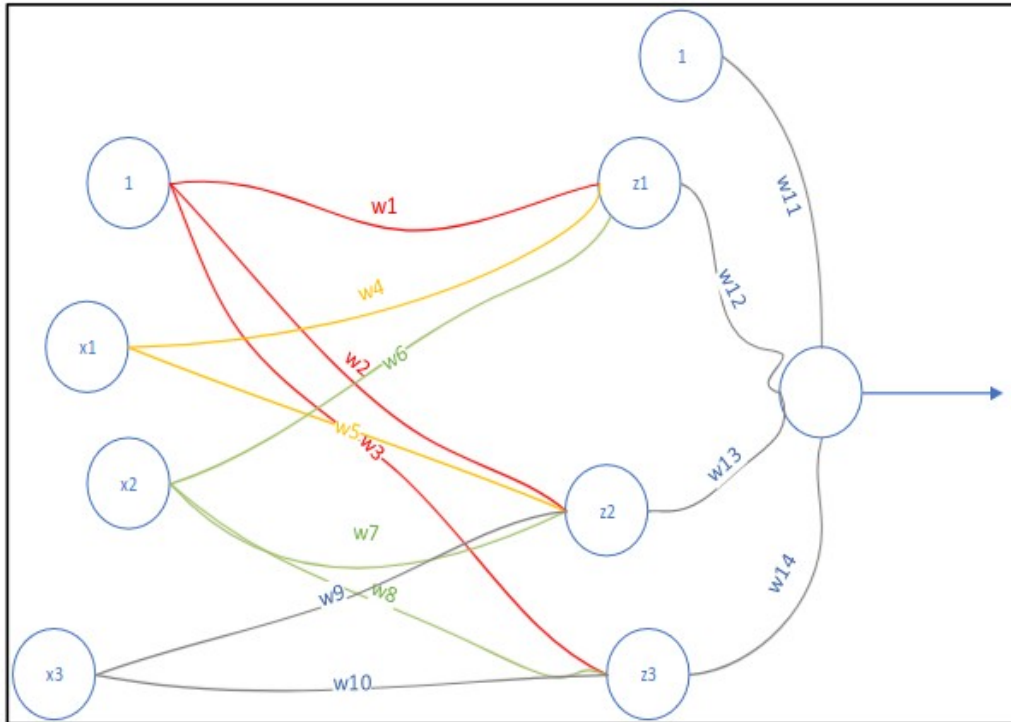
شکل ۵: تغییر ورودی‌ها توسط دو نورون Y1 و Y2

(4)

برای پاسخ به این سوال، باید به این نکته توجه کنیم که گرادیان تابع خطا نسبت به وزن‌ها، به صورت جزئی از مشتق‌های خطا نسبت به خروجی‌های هر لایه تشکیل شده است. اگر همه وزن‌های اولیه یکسان باشند، خروجی‌های هر لایه نیز یکسان خواهند بود. این باعث می‌شود که مشتق‌های خطا نسبت به خروجی‌ها نیز یکسان باشند. در نتیجه،

گرایان تابع خطا نسبت به وزن‌ها نیز یکسان خواهند بود. این یعنی همه وزن‌ها با یک مقدار یکسان به‌روزرسانی می‌شوند و هیچ تفاوتی بین آن‌ها ایجاد نمی‌شود. این باعث می‌شود که شبکه عصبی قادر به یادگیری ویژگی‌های مختلف ورودی‌ها نباشد و در نتیجه همگرا نشود.

(5)



**Problem 5)**  $\sum w_i x_i + b$

$Z_1 = w_1 + x_1 w_4 + x_2 w_6 \rightarrow f(Z_1) = a(Z_1) \leftarrow$  تابع فعل

$Z_2 = w_2 + x_1 w_5 + x_2 w_7 + x_3 w_9 \rightarrow f(Z_2) = a(Z_2) \leftarrow$  تابع فعل

$Z_3 = w_3 + x_2 w_8 + x_3 w_{10} \rightarrow f(Z_3) = a(Z_3) \leftarrow$  تابع فعل

$0 = w_{11} + (Z_1 \cdot w_{12}) + Z_2 \cdot w_{13} + Z_3 \cdot w_{14} \rightarrow f(0) = \frac{1}{1+e^0} = 0.5$

(فرهنگی شبکه عصبی دارد شبکه عصبی است)

فرهنگی در این Input دارد یکبار

1)  $x_1, x_2, x_3$  به  $z_1, z_2, z_3$  و  $z_1, z_2, z_3$  به  $\Sigma$

2)  $x_1, x_2, x_3$  به  $z_1, z_2, z_3$  و  $z_1, z_2, z_3$  به  $\Sigma$

3)  $x_1, x_2, x_3$  به  $z_1, z_2, z_3$  و  $z_1, z_2, z_3$  به  $\Sigma$

دقیقاً برعکس به روش شبکه NOR که با استفاده از شبکه AND و OR ساخته می‌شود.

این شبکه هم برای یک لایه مخفی به 3 نورونیکه‌هاست یک لایه تعمیم می‌دهد.

P4PCO

شبکه عصبی با هر تعداد لایه، بدون توابع غیرخطی، شبیه یک شبکه عصبی تک لایه است. بله می‌توان این شبکه چندلایه را به یک تک لایه تبدیل نمود. به صورت زیر عمل می‌کنیم:

با توجه به مفروضات گفته شده، میتوان این شبکه را با یک شبکه تک لایه جایگزین کرد. برای این کار، باید از رابطه زیر استفاده کنیم:

$$f(z) = \text{sigmoid}(z) = 1 / (1 + e^{-z})$$

که در آن sigmoid تابع سیگموید و  $z$  خروجی هر نورون است. اگر فرض کنیم که ورودی شبکه یک بردار  $x$  باشد و وزن های بین لایه ورودی و لایه پنهان را با  $W_1$  و بایاسهای لایه پنهان را با  $b_1$  نشان دهیم، می‌توان خروجی لایه پنهان را به شکل زیر محاسبه کرد:

$$W_1 = w_4, w_5, w_6, w_7, w_8, w_9, w_{10}$$

$$b_1 = w_1, w_2, w_3$$

$$z_1 = z_1, z_2, z_3$$

$$z_1 = f(W_1 x + b_1)$$

که در آن  $z_1$  یک بردار سه بعدی است که هر عنصر آن نشان دهنده خروجی یک نورون پنهان است. سپس، اگر وزن های بین لایه پنهان و لایه خروجی را با  $W_2$  و بایاس های لایه خروجی را با  $b_2$  نشان دهیم، می‌توان خروجی شبکه را به شکل زیر محاسبه کرد:

$$W_2 = w_{12}, w_{13}, w_{14}$$

$$b_2 = w_{11}$$

$$y = \text{sigmoid}(W_2 z_1 + b_2)$$

که در آن  $y$  یک بردار یک بعدی است که نشان دهنده خروجی شبکه است. حال، اگر بخواهیم این شبکه را با یک شبکه تک لایه جایگزین کنیم، باید از رابطه زیر استفاده کنیم:

$$y = \text{sigmoid}(W x + b)$$

که در آن  $W$  و  $b$  وزن ها و بایاس های شبکه تک لایه هستند. برای محاسبه این وزن ها و بایاس ها، باید از رابطه زیر استفاده کنیم:

$$W = W_2 f(W_1)$$

$$b = W_2 f(b_1) + b_2$$



که در آن  $f(W_1)$  و  $f(b_1)$  بردارهایی هستند که هر عنصر آن ها برابر با مقدار تابع سیگموئید روی هر عنصر ماتریس  $W_1$  و بردار  $b_1$  است. برای مثال، اگر فرض کنیم که وزن ها و بایاس های شبکه با یک لایه پنهان به شکل زیر باشند:

$$W_1 = \begin{bmatrix} 0.1 & -0.2 & 0.3 \\ 0.4 & -0.5 & 0.6 \\ 0.7 & -0.8 & 0.9 \end{bmatrix}$$

$$b_1 = [0.1, 0.2, 0.3]$$

$$W_2 = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix}$$

$$b_2 = [0.1]$$

می توان وزن ها و بایاس های شبکه تک لایه را به شکل زیر محاسبه کرد:

$$f(W_1) = \begin{bmatrix} 0.524 & 0.450 & 0.574 \\ 0.598 & 0.378 & 0.645 \\ 0.668 & 0.310 & 0.710 \end{bmatrix}$$

$$f(b_1) = [0.524, 0.549, 0.574]$$

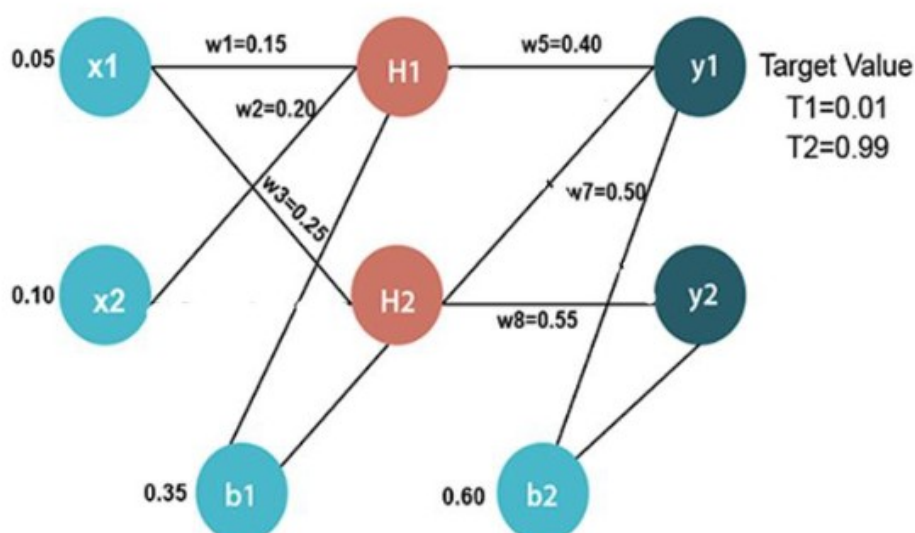
$$W = W_2 f(W_1) = [0.186, 0.133, 0.201]$$

$$b = W_2 f(b_1) + b_2 = [0.214]$$

بنابراین، معماری پیشنهادی برای شبکه تک لایه به شکل زیر است:

$$y = \text{sigmoid}([0.186, 0.133, 0.201] x + [0.214])$$

(6)



$x_1 = 0.5$

$w_1 = 0.15$

$w_5 = 0.140$

(6)

$x_2 = 0.10$

$w_2 = 0.120$

$w_7 = 0.150$

Target Value

$b_1 = 0.135$

$w_3 = 0.125$

$w_8 = 0.155$

$T_1 = 0.1$

$b_2 = 0.160$

$T_2 = 0.199$

Problem 6)

activation function is Sigmoid =  $\frac{1}{1+e^{-x}}$ 

$H_1 = x_1 w_1 + x_2 w_2 + b_1$

$$\text{out } H_1 = \frac{1}{1+e^{-H_1}}$$

forward pass

$$(a_1) \leftarrow H_1 = x_1 w_1 + x_2 w_2 + b_1$$

$$= (0.5 \times 0.15) + (0.1 \times 0.12) + 0.135 = 0.137 \rightarrow \sigma(H_1) = 0.59 = a_1$$

$$(a_2) \leftarrow H_2 = x_1 w_3 + x_2 w_4 + b_1 \rightarrow$$

$$= (0.5 \times 0.125) + (0.1 \times 0.1) + 0.135 = 0.1362 \rightarrow \sigma(H_2) = 0.58 = a_2$$

$$y_1 = a_1 w_5 + a_2 w_7 + b_2$$

$$= (0.59 \times 0.14) + (0.58 \times 0.15) + 0.16 = 0.112 \rightarrow \sigma(y_1) = 0.75 = (z_1)$$

$$y_2 = (a_1 w_6) + (a_2 w_8) + b_2 =$$

$$0 + (0.59 \times 0.155) + 0.16 = 0.1919 \rightarrow \sigma(y_2) = 0.71 = (z_2)$$

$$E_{\text{total}} = \sum \frac{1}{2} \left( \underset{\text{target}}{t} - \underset{\text{output}}{o} \right)^2 = \frac{1}{2} (T_1 - Z_1)^2 + \frac{1}{2} (T_2 - Z_2)^2$$

$$E_{\text{total}} = \frac{1}{2} (0.1 - 0.75)^2 + \frac{1}{2} (0.199 - 0.71)^2 = 0.273 + 0.1312$$

$$E_{\text{total}} = \frac{1}{2} (0.547) + \frac{1}{2} (0.78) = 0.273 + 0.392 = 0.665$$



دریغ دیر بی انتظار!

$$E = \frac{1}{2} (t_i - o_i)^2$$

$$E_1 = \frac{1}{2} (t_1 - y_1)^2 = \frac{1}{2} (1.1 - 1.75)^2 = \frac{1.547}{2} = 1.273$$

$$E_2 = \frac{1}{2} (t_2 - y_2)^2 = \frac{1}{2} (1.99 - 1.71)^2 = \frac{1.078}{2} = 1.0392$$

$$\frac{\partial E_1}{\partial w_{ji}} = \underbrace{\frac{\partial E}{\partial o_j}}_{-(t_i - o_j)} \times \underbrace{\frac{\partial o_j}{\partial \text{net } j}}_{o_j(1 - o_j)} \times \underbrace{\frac{\partial \text{net } j}{\partial w_{ji}}}_{x_{ji}}$$

## Backpropagation:

Consider  $w_5$ 

$$\text{Error at } w_5 = \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial z_1} \times \frac{\partial z_1}{\partial y_1} \times \frac{\partial y_1}{\partial w_5} \times \alpha_{ji}$$

$$\frac{\partial E_{\text{total}}}{\partial z_1} = \underbrace{\frac{\partial E_{\text{total}}}{\partial z_1}}_{-(T_1 - z_1)} = -(T_1 - z_1)$$

$$\textcircled{1} 1.74 \times 1.18 \times 1.59 = 1.08 \rightarrow \frac{\partial E_{\text{total}}}{\partial w_5} = 1.082 \quad \text{میزان تغییر وزن } w_5$$

$$\text{updating } w_5: w_5 = w_5 - \eta \times \frac{\partial E_{\text{total}}}{\partial w_5} = 1.4 - 0.1(1.082) = 1.391 \approx 1.4$$

in the same way  $w_7, w_8$ 

$$\frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial z_1} \times \frac{\partial z_1}{\partial y_1} \times \frac{\partial y_1}{\partial w_7} = 1.74 \times 1.18 \times 1.58 = 1.077$$

$$w_7 = w_7 - \eta \frac{\partial E_{\text{total}}}{\partial w_7} = 1.5 - 1.1(1.077) = 1.49 = 1.5$$

$$w_8 = \frac{\partial E_{\text{total}}}{\partial w_8} = \frac{\partial E}{\partial z_2} \times \frac{\partial z_2}{\partial y_2} \times \frac{\partial y_2}{\partial w_8}$$

$$= -(T_2 - z_2) \cdot z_2(1 - z_2) \cdot \alpha_2 = 1.28 \cdot 1.2 \cdot 1.58 =$$

$$w_8 = w_8 - \eta \frac{\partial E_{\text{total}}}{\partial w_8} = 1.55 - 1.1(1.032) = 1.54 \quad \textcircled{1.032}$$

Now at Hidden layer updating  $w_1, w_2, w_3$

$$\textcircled{1} \frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{total}}}{\partial a_1} \times \frac{\partial a_1}{\partial H_1} \times \frac{\partial H_1}{\partial w_1} \quad w_1 = 0.5$$

$$\frac{\partial E_{\text{total}}}{\partial a_1} = \frac{\partial E_1}{\partial a_1} + \frac{\partial E_2}{\partial a_1}$$

$$\frac{\partial E_1}{\partial a_1} = \frac{\partial E_1}{\partial y_1} \times \frac{\partial y_1}{\partial a_1} \rightarrow w_5 = 14$$

$$\frac{\partial E_1}{\partial a_1} = 113 \times 14 = 1555$$

$$\frac{\partial E_1}{\partial a_1} = 113 \times 14 = 1555$$

$$\frac{\partial E_1}{\partial y_1} = \frac{\partial E_1}{\partial z_1} \times \frac{\partial z_1}{\partial y_1} = 174 \times 118 = 113$$

$$\frac{\partial E_2}{\partial a_1} = \frac{\partial E_2}{\partial y_2} \times \frac{\partial y_2}{\partial a_1} = -0.19$$

$$\frac{\partial E_{\text{total}}}{\partial a_1} = 1555 + (-0.19) = 1554.81$$

$$\star \frac{\partial a_1}{\partial H_1} = a_1(1-a_1) = 0.59(1-0.59) = 0.24$$

$$H_1 = w_1 x_1 + w_2 x_2 + b_1$$

$$\frac{\partial H_1}{\partial w_1} = x_1 = 0.5$$

$$\textcircled{2} \frac{\partial E_{\text{total}}}{\partial w_1} = 1554.81 \times 0.24 \times 0.5 = 196.58$$

$$\text{updating } w_1 = w_1 - \eta \frac{\partial E_{\text{total}}}{\partial w_1} = 0.5 - 0.1(196.58) = 0.30342$$

$$w_2 = w_2 - \eta \frac{\partial E_{\text{total}}}{\partial w_2} = 0.2 - 0.1(196.58) = 0.00342$$

$$1554.81 \times 0.24 \times 0.1 = 37.31544$$

$$\frac{\partial E_{total}}{\partial w_3} = \left( \frac{\partial E_{total}}{\partial a_2} \right) \times \left( \frac{\partial a_2}{\partial H_2} \right) \times \left( \frac{\partial H_2}{\partial w_3} \right) \quad w_3 = .125$$

$$\frac{\partial E_{total}}{\partial a_2} = \frac{\partial E_1}{\partial a_2} + \frac{\partial E_2}{\partial a_2} + \dots$$

$$\frac{\partial E_1}{\partial a_2} = \frac{\partial E_1}{\partial y_2} \times \frac{\partial y_2}{\partial a_2}$$

$$\frac{\partial E_2}{\partial a_2} = \frac{\partial E_2}{\partial y_2} \times \frac{\partial y_2}{\partial a_2} =$$

$$\frac{\partial a_2}{\partial H_2} = a_2(1-a_2) = .58(1-.58) = .12$$

$$\frac{\partial E_{total}}{\partial w_3} = 0.1 \times .12 \times .125 = .01$$

$$\text{updating } w_3 = w_3 - \eta \left( \frac{\partial E_{total}}{\partial w_3} \right) = .125 - .1(.01) = .1249$$


---