



تمرین ۳ درس یادگیری ماشین
آریا جلالی
۹۸۱۰۵۶۶۵

۱ Comparison between some models

۱.۱

درخت تصمیم در بسیاری از موارد روی داده‌ی آموزش overfit می‌کند و طبق Bias-variance tradeoff دارای بایاس کم و واریانس بالا هستند. طبق رابطه‌ی بایاس می‌توانیم بنویسیم

$$Bias = (f - E[f'])^2$$

است. همانطور که مشخص است، اگر ما دفعات متعددی درخت تصمیم آموزش دهیم و میانگین آن‌ها را به عنوان classifier نهایی قرار دهیم، خطای ما به مراتب روی داده‌ی تست کاهش پیدا می‌کند.

یکی از روش‌هایی که باعث می‌شود واریانس ما کمتر شود (و خطا روی داده‌ی تست نیز کاهش پیدا کند) استفاده از تعدادی درخت تصمیم با correlation کم است. برای کم کردن هرچه بیشتر این ارتباط می‌توانیم feature‌های متفاوت یا داده‌های متفاوت برای آموزش به هر درخت بدهیم و در نهایت بین نتایج بدست آمده رای‌گیری داشته باشیم. در واقع استفاده از Random Forests با استفاده کردن از چندین درخت تصمیم جلوی overfit شدن که یکی از خصوصیات بارز درخت تصمیم است را می‌گیرد.

۲.۱

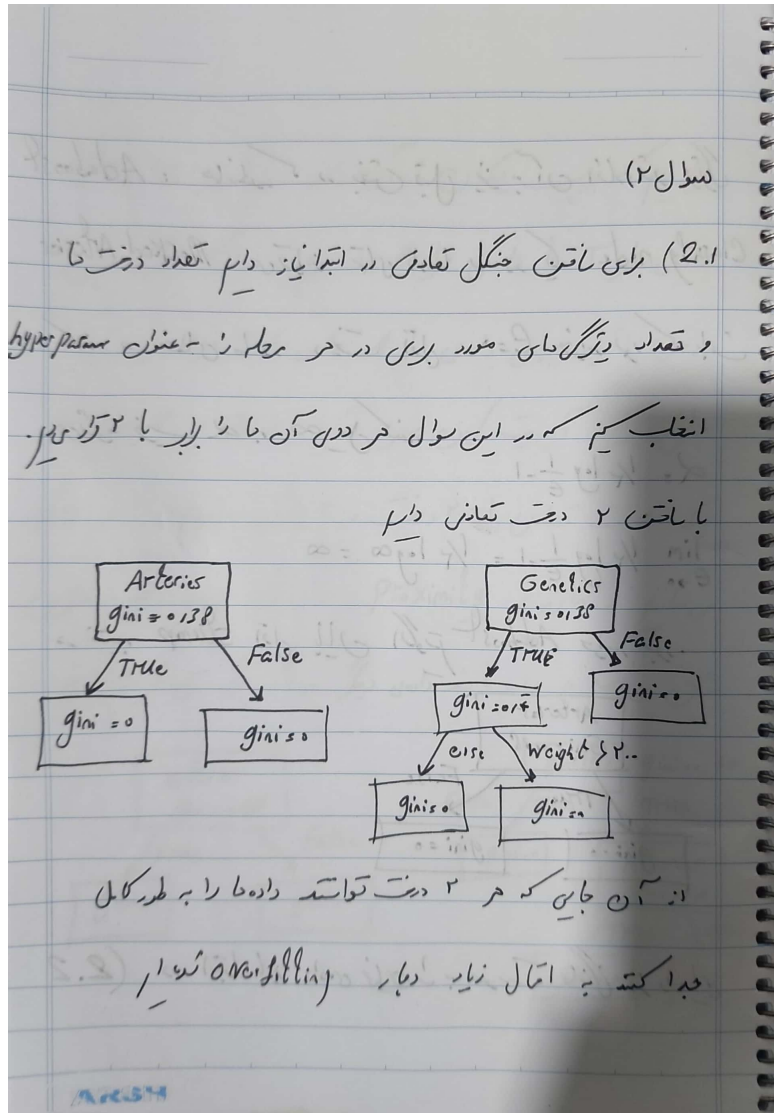
XGBoost در مقایسه با Random Forests عملکرد بهتری در دیتاست‌های unbalanced از خود نشان می‌دهد، زیرا اگر در گام اول نتواند به درستی نتیجه را پیش‌بینی کند، در گام‌های بعدی توجه و اهمیت بیشتری به آن کلاس می‌دهد.

XGBoost سرعت بالاتری نسبت به AdaBoost دارد و حساسیت کمتری نسبت به نویز موجود در دیتاست از خود نشان می‌دهد.

XGBoost در مقابل Gradient Boosting همانند بخش قبل سریعتر است و با استفاده از loss های $l1$ و $l2$ احتمال overfit شدن را می‌گیرد و قدرت generalization بیشتری از خود نسبت به Gradient Boosting نشان می‌دهد.

Fitting a model ۲

۱.۲



Adaboost: همانند که در بحث قبل نیز به آن اشاره شد درخت

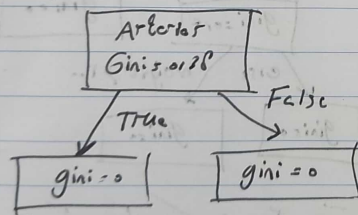
Mocked Arborescences می توانیم آنها را در یک Classify node

کنند و خطای ما در درخت اول $\epsilon = 0$ خواهد بود که بافت

می شود زیرا α به ∞ میل کند.
 $\alpha = \frac{1}{\epsilon} \log \frac{1}{\epsilon} - 1$

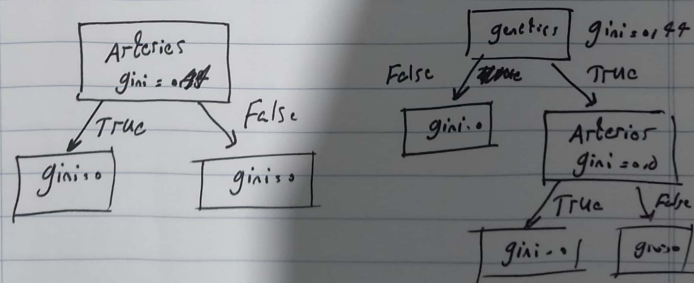
$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \log \frac{1}{\epsilon} - 1 = \frac{1}{\epsilon} \log \infty = \infty$$

در نتیجه stump اول بیان الگوریتم Adaboost خواهد بود.



2.2 در ابتدا باید داده‌های نمای را به صورت عایقی از داده‌ها

فردم که به یکسان با آن دانش در این دیتا تنها
 دانش ۳ با دانش ۴ در بهر متغیر است. پس
 حدس اولیه برای دانش ۴ دانش ۳ خواهد بود.
 حال برای آسپت کردن حدس اولیه با استفاده از ۳ دانش ادی
 یک درخت تصانی با ۲ دخت و ۲ برگ زدیم در هر مرحله
 می‌سازیم و مارتس proximity را تشکیل می‌دهیم. این مارتس
 نشان می‌دهد هر داده با چند دانش دیگر در یک برگ قرار گرفته است



در هر ۲ دخت دهن سطر 4 تنها با دهن سطر ۳ در
یک برگ را میگیرد، آئینت نمیخورد، الکتر ما درین
مرحله به پایان میرسد.

$$X_4 = X_7 = \text{Yes} - 167 - \text{No} - \text{Yes}$$

۳ Gradient Boost

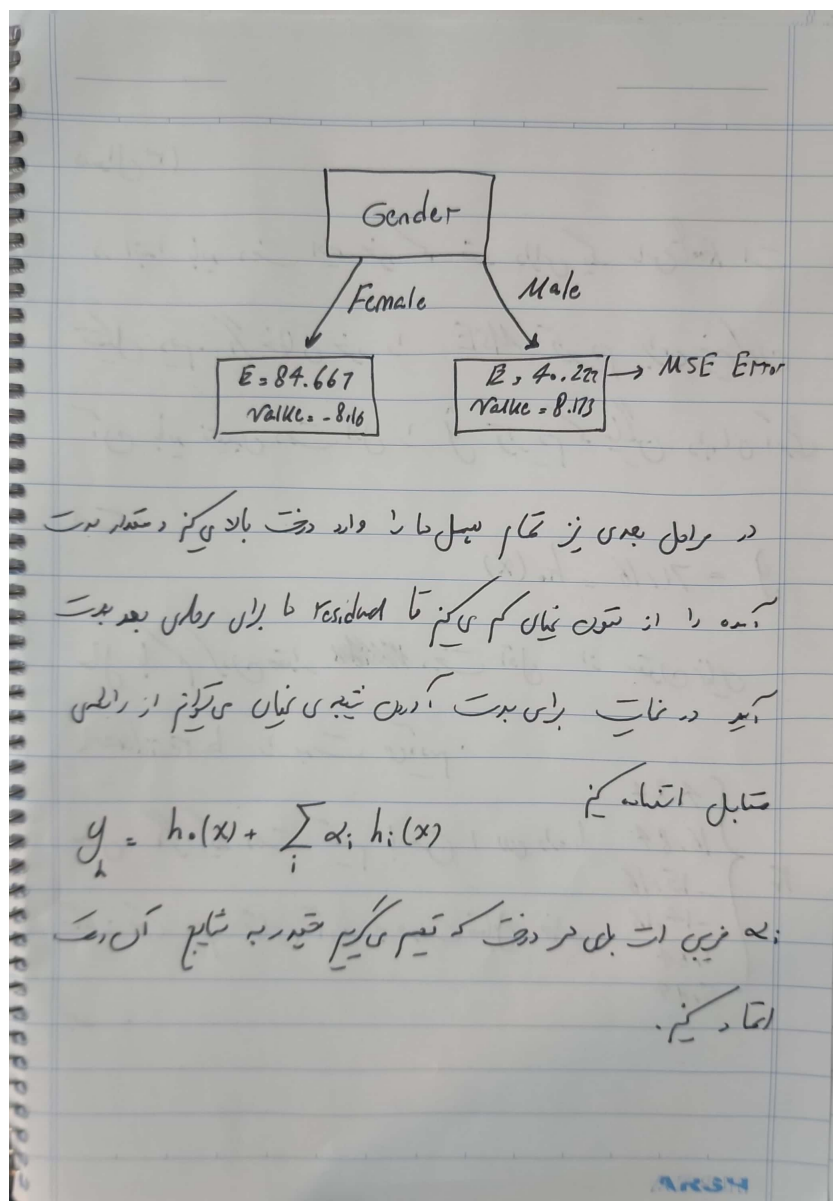
سوال ۳)

در ابتدا باید درخت اولیه خود که فقط دارای یک راس R_{root} است
تسکیل دهیم. اگر خطای خود را MSE قرار دهیم برای تسکیل کردن
آن باید نتایج درخت اول را \hat{y} قرار دهیم که میانی داده های آموزش
است.

$$\hat{y} = 71,16 = h_0(x)$$

حال با کم کردن مقدار ~~MSE~~ درخت اول از ستون نمان
residual را بدست می آوریم.

$$\begin{cases} 4,84 \\ 16,84 \\ -15,16 \\ -14,16 \\ 1,84 \\ 5,84 \end{cases} \quad \begin{array}{l} \text{حال اگر یک درخت تقسیم با محور ۱ پس داده ها} \\ \text{آموزش دهیم نتایج آن به صورت مقابل خواهد بود} \end{array}$$



۴ Descriptive Questions

۱.۴

Exploding Gradient زمانی رخ می‌دهد که گرادیان ارورها به صورت نمایی بزرگ می‌شوند و این باعث می‌شود در مرحله‌ی آپدیت وزن مقادیر w_i بسیار افزایش یا کاهش پیدا کنند و مدل ما unstable شود و نتواند از داده‌ی آموزشی چیز معنی‌داری یاد بگیرد. برای حل این مشکل می‌توانیم از روش Weight Initialization استفاده کنیم. به این صورت که مقادیر اولیه وزن‌ها را به صورت رندوم انتخاب نمی‌کنیم و سعی می‌کنیم مقادیر داده شده از یک

توزیع نرمال با یک سری پارامتر خاص پیروی کنند تا وزن‌ها در بازه‌ی محدودی قرار بگیرند و مشکل Exploding Gradient با احتمال کمتری رخ دهد. یکی دیگر از روش‌های حل این مشکل استفاده از Gradient Clipping است. در این روش اگر مقدار گرادیان بزرگتر از یک بازه‌ای باشد مقدار آن برابر با کران بالای بازه‌ی انتخاب شده قرار داده می‌شود و همین عمل برای گرادیان‌ها کوچکتر از بازه و کران پایین انجام می‌شود. با این روش گرادیان‌ها در محدوده‌ی خاصی قرار می‌گیرند و به صورت نامحدود رشد نمی‌کنند.

۲.۴

اضافه کردن لایه‌های بیشتر به مدل باعث پیچیده‌تر شدن آن می‌شود و می‌توانیم بایاس را کاهش دهیم. ولی اگر پیچیدگی مدل به مراتب بیشتر از حد نیاز برای حل سوال باشد، دچار overfitting می‌شویم و مدل ما واریانس بالایی خواهد داشت. از طرفی اضافه کردن لایه‌های زیاد باعث افزایش زمان train و به وجود آمدن مشکلات متعددی مانند Exploding Gradient و Vanishing Gradient می‌شود.

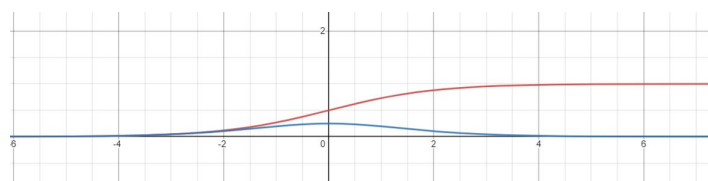
۳.۴

در صورتی که داده‌ی آموزشی ما از قبل به صورت تصادفی چیده نشده باشد، ممکن است ارتباطی بین داده‌های نزدیک به هم باشد و مدل ما سعی می‌کند با یاد گرفتن این ارتباط خطا را کاهش دهد. در صورتی که هدف ما یاد گرفتن توزیع آماری مدل به صورت کلی است. این bias در داده‌ی آموزشی باعث می‌شود مدل به جای یاد گرفتن ویژگی‌ها و روابط توصیف کننده‌ی توزیع داده، روابط بین داده‌ها را یاد بگیرد و دچار overfitting شدید روی داده‌ی آموزشی شود.

۴.۴

$$(1 + e^{-x})\sigma = 1 \Rightarrow -e^{-x}\sigma + (1 + e^{-x})\frac{d\sigma}{dx} = 0$$

$$\frac{d\sigma}{dx} = \sigma \cdot \frac{e^{-x}}{(1 + e^{-x})} = \sigma \cdot \frac{(1 + e^{-x}) - 1}{(1 + e^{-x})} = \sigma \cdot \left[1 - \frac{1}{(1 + e^{-x})}\right] = \sigma \cdot (1 - \sigma)$$



شکل ۱: شکل تابع و مشتق سیگموئید

همانطور که از شکل بالا مشخص است در صورتی که مقدار وزن یا مقدار ورودی به تابع بسیار بزرگ باشد، گرادیان آن تقریباً برابر با ۰ خواهد بود و دچار مشکل Vanishing Gradient می‌شویم و مدل در تغییر وزن‌ها دچار مشکل می‌شود و فرایند Train به اتمام نمی‌رسد یا زمان بسیار زیادی طول خواهد کشید. یکی از راه‌های حل این مشکل تغییر تابع activation به یک تابع دیگر مانند RELU است.

۵.۴

در این حالت ۲ مشکل ممکن است رخ دهد. اگر همانند صورت سوال فرض کنیم مقادیر بزرگتر از ۵.۰ برابر با کلاس ۱ قرار داده می‌شوند و مقادیر ۵.۰ کمتر برابر با کلاس ۰، مشکلی از لحاظ پیشبینی همواره یک کلاس نخواهیم داشت (در صورتی که مقادیر ۵.۰ و بزرگتر به عنوان کلاس ۰ پیشبینی می‌شدند به دلیل نامنفی بودن خروجی تابع RELU همواره پیشبینی مدل کلاس ۱ بود). ولی مشکل اساسی این مدل این است که اگر اعضای کلاس ۱ به اشتباه به کلاس ۰ نسبت داده شوند به دلیل منفی بودن نتیجه و صفر بودن مشتق تابع $\sigma(ReLU(x))$ در نقاط منفی گرادیان برابر با ۰ می‌بود و مدل از این پیشبینی اشتباه چیزی یاد نمی‌گرفت.

۶.۴

در صورتی که تابع activation نداشته باشیم، خروجی هر لایه را به صورت یک ضرب ماتریسی به صورت مقابل نشان داد

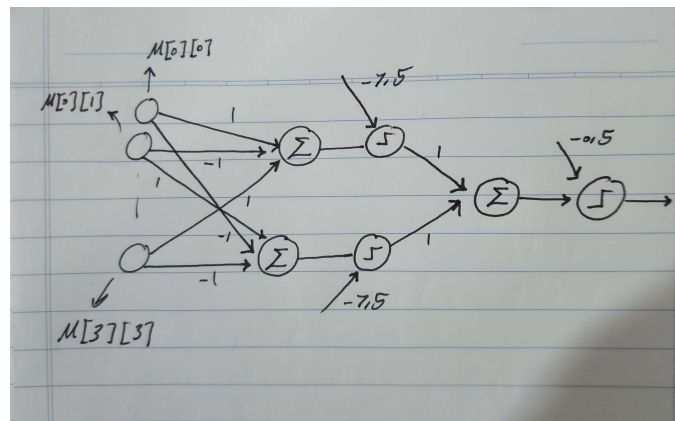
$$A^n = W^n(A^{n-1}) + b^{n-1} = W^n(W^{n-1}A^{n-2} + b^{n-1}) + b^{n-2} = \dots = W'(X) + b^0$$

$$W'(X) = W^n(W^{n-1}(W^{n-2}(\dots) + b^{n-3}) + b^{n-2})$$

با ضرب متوالی ماتریس‌ها به این نتیجه می‌رسیم که خروجی هر لایه ترکیب خطی‌ای از ورودی‌ها است و مدل ما در واقع یک Regressor است.

Intuitive Questions ۵

۱.۵



شکل ۲: مدل طراحی شده برای تشخیص شکل شطرنجی

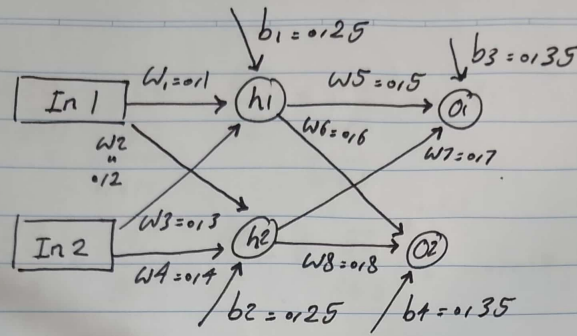
در مدل بالا به پیکسل‌هایی که باید در حالت مدنظر ما ۱ باشند وزن ۱ و به پیکسل‌های ۰ وزن -۱ را می‌دهیم. در این صورت حداکثر مقدار جمع برابر با ۸ خواهد بود، زیرا اگر پیکسل دیگری روشن باشد حداقل یکی از خروجی کم می‌شود. در نهایت این مقدار را با بایاس -۵.۷ جمع می‌کنیم و اگر شکل صفحه همانند صفحه‌ی نشان داده شده باشد مقدار ورودی به تابع پله برابر با ۵.۰ خواهد بود و در غیر اینصورت مقداری منفی خواهد بود. عکس همین کار برای حالتی که جای پیکسل‌های سیاه و سفید شود انجام شده است و در نهایت نتیجه‌ی خروجی‌ها باهم OR شده‌اند. در نهایت اگر صفحه‌ی مدنظر تشخیص داده شود خروجی مدل ما برابر با ۱ و در غیر اینصورت برابر با ۰ خواهد بود.

۲.۵

MLP ها رابطه‌ی بین پیکسل‌های اطراف را در نظر نمی‌گیرند و هر پیکسل را مستقل از پیکسل‌های اطراف در نظر می‌گیرند (از دانش prior استفاده‌ای نمی‌شود). از طرفی مدل‌های MLP حساس به مکان object در تصویر هستند و در اکثر اوقات دچار overfit روی داده‌ی آموزشی می‌شوند. این مدل‌ها دارای Invariance to translation نیستند و تنها داده‌های جدید را می‌توانند تشخیص دهند که در داده‌ی آموزشی نمونه‌ای از آن در همان نقطه وجود داشته باشد.

مشکل دیگر استفاده از MLP و در نظر نگرفتن رابطه‌ی بین پیکسل‌ها وجود پارامترهای فراوان در عکس‌ها است و این باعث کاهش سرعت یادگیری می‌شود. از طرفی منابع پردازشی زیادی برای آموزش نیاز است.

Computational Question 6



$$\text{output } h_1 = \sigma(0.1 + 0.18 + 0.18) \approx 0.601$$

$$\text{output } h_2 = \sigma(0.2 + 0.2 + 0.25) \approx 0.615$$

$$O_1 = \sigma(0.5 \cdot 0.601 + 0.7 \cdot 0.615 + 0.35) \approx 0.747$$

$$O_2 = \sigma(0.6 \cdot 0.601 + 0.8 \cdot 0.615 + 0.35) \approx 0.769$$

$$E = \frac{1}{2} \sum_i (t_i - o_i)^2 \approx 0.259$$

$$\frac{\partial E}{\partial O_1} = (O_1 - t_1) = 0.697, \quad \frac{\partial E}{\partial O_2} = (O_2 - t_2) = -0.181$$

$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial O_1} \frac{\partial O_1}{\partial w_5} = 0.1079$$

$$0.697 \cdot (0.747) \cdot (0.253) \cdot (0.601)$$

$$0.147, 0.1253$$

$$\frac{\partial E}{\partial w_7} = \frac{\partial E}{\partial a_1} \times \frac{\partial a_1}{\partial z} \times \frac{\partial z}{\partial w_7} = 0.181$$

0.697 0.615

$$\frac{\partial E}{\partial w_6} = \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z} \times \frac{\partial z}{\partial w_6} = -0.1019$$

-0.181 0.1769 0.1231 0.601

$$\frac{\partial E}{\partial w_8} = \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z} \times \frac{\partial z}{\partial w_8} = -0.102$$

$$\frac{\partial E}{\partial b_3} = \frac{\partial E}{\partial a_1} \times \frac{\partial a_1}{\partial z} \times \frac{\partial z}{\partial b_3} = 0.132$$

$$\frac{\partial E}{\partial b_4} = \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z} \times \frac{\partial z}{\partial b_4} = -0.1032$$

☆ برای محاسبه گراییان: به کمک یک یا از به گراییان یا محاسبه فردی:

دایر h, h, h

$$\frac{\partial E}{\partial x_5} = \frac{\partial E}{\partial x_7} = 0.132, \quad \frac{\partial E}{\partial x_6} = \frac{\partial E}{\partial x_8} = -0.1032$$

$$\frac{\partial E}{\partial w_1} = \underbrace{\frac{\partial E}{\partial x_5}}_{0.132} \frac{\partial x_5}{\partial z} \frac{\partial z}{\partial w_1} + \frac{\partial E}{\partial x_6} \frac{\partial x_6}{\partial z} \frac{\partial z}{\partial w_1} = 0.001$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial x_7} \frac{\partial x_7}{\partial z} \frac{\partial z}{\partial w_2} + \frac{\partial E}{\partial x_8} \frac{\partial x_8}{\partial z} \frac{\partial z}{\partial w_2} = 0.002$$

$$\frac{\partial E}{\partial w_3} = \frac{\partial E}{\partial x_5} \frac{\partial x_5}{\partial z} \frac{\partial z}{\partial w_3} + \frac{\partial E}{\partial x_6} \frac{\partial x_6}{\partial z} \frac{\partial z}{\partial w_3} = 0.006$$

$$\frac{\partial E}{\partial w_4} = \frac{\partial E}{\partial x_7} \frac{\partial x_7}{\partial z} \frac{\partial z}{\partial w_4} + \frac{\partial E}{\partial x_8} \frac{\partial x_8}{\partial z} \frac{\partial z}{\partial w_4} = 0.008$$

$$\frac{\partial E}{\partial b_1} = \frac{\partial E}{\partial x_5} \frac{\partial x_5}{\partial z} \frac{\partial z}{\partial b_1} + \frac{\partial E}{\partial x_6} \frac{\partial x_6}{\partial z} = 0.011$$

$$\frac{\partial E}{\partial b_2} = \frac{\partial E}{\partial x_7} \frac{\partial x_7}{\partial z} + \frac{\partial E}{\partial x_8} \frac{\partial x_8}{\partial z} = 0.016$$

پس از آنکه کد در دسترس داریم

$$w'_1 = 0.099, w'_2 = 0.294, w'_3 = 0.1421, w'_4 = 0.619$$

$$w'_2 = 0.198, w'_4 = 0.392, w'_6 = 0.619, w'_8 = 0.82$$

$$b'_1 = 0.239, b'_2 = 0.234, b'_3 = 0.218, b'_4 = 0.382$$