



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



یادگیری ماشین

تمرین شماره ۳

نام و نام خانوادگی

امیرحسین چمی

دی ۱۴۰۴

۴..... فهرست شکل ها

۵..... فهرست جدول ها

۶..... چکیده

۷..... پرسش ۱ – PERCEPTRON

۷..... ۱-۱. رسم نمودار امتیاز B بر حسب امتیاز A

۷..... ۲-۱. اجرای یک حلقه از الگوریتم پرسپترون

۸..... ۳-۱. وزن ها و جداسازی داده ها

۸..... ۴-۱. توانایی جداسازی الگوریتم پرسپترون

۸..... ۱-۴-۱. حالت a

۸..... ۲-۴-۱. حالت b

۸..... ۳-۴-۱. حالت c

۹..... پرسش ۲ – NEURAL NETWORKS: REPRESENTATION

۹..... ۱-۲. نمودارهای A

۹..... ۲-۲. نمودارهای B

۱۰..... ۳-۲. نمودار C

پرسش ۳ – NEURAL NETWORKS: FORWARD AND BACKWARD

۱۱..... PROPAGATION

۱۱..... ۱-۳. ابعاد هر پارامتر و وکتورایز کردن شبکه

۱۳..... ۲-۳. محاسبه $\partial J / \partial y_i$ و $\partial J / \partial y$

۱۳..... ۳-۳. محاسبه $\partial y_i / \partial z_2$

۱۴..... ۴-۳. محاسبه $\partial z_2 / \partial a_1$

۱۴..... ۵-۳. محاسبه $\partial a_1 / \partial z_1$

۱۵..... ۶-۳. محاسبه $\partial z_1 / \partial W_1$

۱۶..... ۷-۳. محاسبه $\partial J / \partial W_1$

۱۷..... پرسش ۴ – HESSIAN OF A TWO-LAYER NETWORK

۱۷	۱-۴. هر دو وزن در لایه دوم
۱۸	۲-۴. هر دو وزن در لایه اول
۲۰	۳-۴. هر وزن از یکی از لایه ها
۲۱	پرسش ۵- LINEAR DISCRIMINANT ANALYSIS (LDA)
۲۳	پرسش ۶- FEATURE CONDITIONING AND PCA
۲۵	پرسش ۷- FEATURE CONDITIONING
	COMPUTATIONAL EXERCISE: FLOWER IMAGE CLASSIFICATION
۲۸	
۲۸	۱-۸. پیش پردازش داده ها
۲۸	۲-۸. بخش اول: آموزش از صفر
۲۸	۱-۲-۱. مدل خطی
۲۹	۲-۲-۱. شبکه با یک لایه پنهان و <i>ReLU</i>
۲۹	۳-۲-۱. شبکه‌ی بهبود یافته
۳۰	۳-۸. بخش دوم: یادگیری انتقالی
۳۰	۱-۳-۱. مدل <i>AlexNet</i>
۳۱	۲-۳-۱. مدل <i>VGG ۱۶</i>
۳۲	۴-۸. مقایسه و تحلیل نهایی
۳۲	۵-۸. نتیجه گیری
۳۳	مراجع

فهرست شکل‌ها

- شکل ۱. نمودار داده‌ها..... ۷
- شکل ۲. نمودار اول b..... ۹
- شکل ۳. نمودار دوم b..... ۱۰
- شکل ۴. نمودار اول c..... ۱۰
- شکل ۵. نمودار دوم c..... ۱۰
- شکل ۶. نمودار CCR الگوریتم انتخاب پیشرونده..... ۲۶
- شکل ۷. نمودار CCR الگوریتم حذف پسرونده..... ۲۷
- شکل ۸. نمودار دقت و loss مدل خطی..... ۲۹
- شکل ۹. نمودار دقت و loss شبکه با یک لایه پنهان و ReLU..... ۲۹
- شکل ۱۰. نمودار دقت و loss شبکه بهبود یافته..... ۳۰
- شکل ۱۱. نمودار دقت و loss مدل AlexNet..... ۳۱
- شکل ۱۲. نمودار دقت و loss مدل VGG^{۱۶}..... ۳۱

فهرست جدول‌ها

جدول ۱. جدول مقایسه مدل‌ها ۳۲

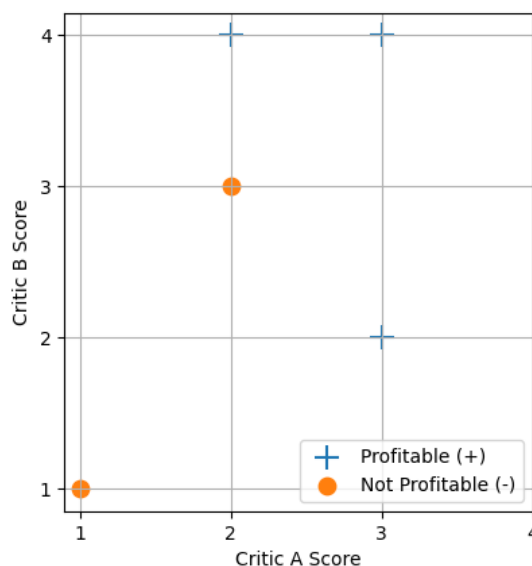
هدف از این مجموعه تمرین، بررسی توانایی، محدودیت‌ها و رفتار یادگیری مدل‌های مختلف یادگیری ماشین از دیدگاه تئوری و عملی است. در بخش‌های ابتدایی، با تحلیل مسئله‌ی پرسپترون، مفهوم جداسازی خطی و وابستگی قدرت مدل به نوع ویژگی‌ها مورد مطالعه قرار می‌گیرد و نشان داده می‌شود که حتی الگوریتم‌های یادگیری ساده نیز در صورت نامناسب بودن نمایش داده، قادر به حل برخی مسائل نخواهند بود. همچنین در ادامه، با بررسی شبکه‌های عصبی با توابع فعال‌سازی ReLU، ظرفیت نمایش (Representational Capacity) شبکه‌ها و نقش تعداد نورون‌ها و لایه‌های میانی در پیاده‌سازی توابع قطعه‌ای خطی تحلیل می‌شود.

در بخش‌های بعدی، تمرین به تحلیل دقیق‌تر شبکه‌های عصبی از منظر پس‌انتشار خطا، مشتق‌گیری و ساختار گرادینان و هسین می‌پردازد و در کنار آن، مفاهیم آماری مهمی مانند LDA و PCA برای درک بهتر نقش توزیع داده و کاهش بُعد بررسی می‌شوند. در نهایت، با اجرای عملی روش‌های Feature Selection شامل Forward Selection و Backward Elimination بر روی یک دسته‌بند Naive Bayes، تأثیر انتخاب ویژگی‌ها بر عملکرد مدل به صورت تجربی ارزیابی می‌شود. به نظرم این تمرین تقریباً دید جامعی از ارتباط میان پیچیدگی مدل، نمایش داده و عملکرد عملی در مسائل یادگیری ماشین ارائه می‌دهد.

پرشی ۱ – Perceptron

۱-۱. رسم نمودار امتیاز B بر حسب امتیاز A

با استفاده از یک کد پایتون نمودار آن را به صورت زیر رسم کردم. علامت مثبت برای فیلم های سودده و منفی برای فیلم های زیانده هستند. (سوال کد استفاده شده را خواسته به همین دلیل ارائه نمی شود).



شکل ۱. نمودار داده ها

واضا این داده ها را نمی توان با استفاده از یک خط جداسازی کرد.

۱-۲. اجرای یک حلقه از الگوریتم پرسپترون

مرحله به مرحله پیش می رویم.

$$score = -1 \cdot 1 + 0 \cdot 1 + 0 \cdot 1 = -1 \quad (1)$$

پیش بینی منفی شد و فیلم اول واقعا زیانده بوده پس پیشبینی درست است و به روز رسانی وزن ها نیازی نیست.

$$score = -1 \cdot 1 + 0 \cdot 3 + 0 \cdot 2 = -1 \quad (2)$$

پیش بینی منفی شد و فیلم دوم سودده بوده پس پیشبینی درست نیست و باید وزن ها به روز رسانی شود.

$$W_{new} = W_{old} + yX = [-1, 0, 0] + [1, 3, 2] = [0, 3, 2] \quad (3)$$

حال به سراغ نمونه ی بعدی می رویم.

$$score = 0 + 3 \cdot 2 + 2 \cdot 4 = 14 \quad (4)$$

پیش بینی مثبت شد و فیلم سوم واقعا سودده بوده پس پیش بینی درست است و نیازی به آپدیت وزن ها نیست.

$$score = 0 + 3.3 + 2.4 = 17 \quad (5)$$

پیش بینی مثبت شد و فیلم چهارم واقعا سودده بوده پس پیش بینی درست است و نیازی به آپدیت وزن ها نیست.

$$score = 0 + 3.2 + 2.3 = 12 \quad (6)$$

پیش بینی مثبت شد و فیلم پنجم واقعا زیانده بوده پس پیش بینی درست نیست و باید وزن ها به روز رسانی شود.

$$W_{new} = W_{old} + y X = [0, 3, 2] - [1, 2, 3] = [-1, 1, -1] \quad (7)$$

پس وزن نهایی بعد از یک دور $[-1, 1, -1]$ است.

۳-۱. وزن ها و جداسازی داده ها

در بخش اول دیدیم که داده ها جداپذیر خطی نیستند. وزن ها در این الگوریتم با یک دور اجرا نتوانستند نقاط را جدا کنند ولی حتی اگر بینهایت دور هم الگوریتم را اجرا کنیم، نمی توانند این داده ها را به درستی جدا کنند.

۴-۱. توانایی جداسازی الگوریتم پرسپترون

۴-۱-۱. حالت a

در این حالت مرز تصمیم برابر $A + B = 8$ است. از آنجا که این مرز تصمیم با یک خط مشخص می شود پس این الگوریتم در این حالت موفق است.

۴-۱-۲. حالت b

این حالت یک ناحیه مربعی را با نقاط $(2, 2)$ ، $(2, 3)$ ، $(3, 2)$ و $(3, 3)$ مشخص می کند و یک خط نمی تواند داخل و بیرون مرتبط را جدا کند پس این الگوریتم در این حالت موفق نیست.

۴-۱-۳. حالت c

در این حالت مرز تصمیم برابر $A - B = 0$ است. از آنجا که این مرز تصمیم با یک خط مشخص می شود پس الگوریتم پرسپترون در این حالت موفق است.

پرسش ۲ – Neural Networks: Representation

منطق حل این سوال واضح و برای تمام شبکه ها و توابع داده شده یکسان است. هر تابع ReLU حداکثر یک شکست را در خط می تواند مدلسازی کند.

شبکه های G_1, G_2, H_1 و H_2 تابع ReLU ندارند پس فقط خط ها را می توانند مدلسازی کنند (یاد بگیرند).

شبکه های G_3, G_4 و G_5 یک ReLU دارند بنابراین حداکثر یک شکست را مدلسازی می کنند (یاد بگیرند).

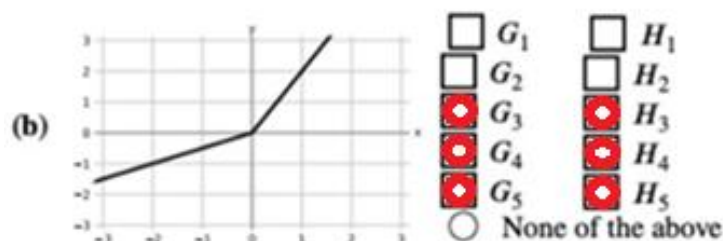
شبکه های H_3, H_4 و H_5 چند ReLU موازی دارند (به دلیل لایه های مخفی) بنابراین می توانند حداکثر چند شکست را مدلسازی کنند (یاد بگیرند). بر همین اساس می توانیم تکلیف تمامی نمودارها را مشخص کنیم.

۱-۲. نمودارهای a

با توجه به اینکه هر دوی این نمودارها خط صاف بدون شکستگی هستند، به راحتی می توان با استفاده از همهی این شبکه ها این دو نمودار را یاد گرفت.

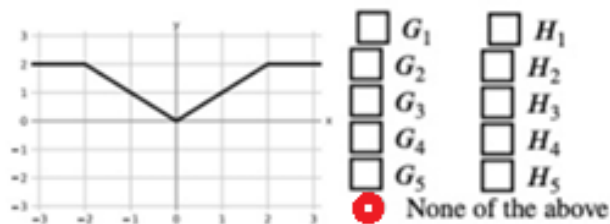
۲-۲. نمودارهای b

نمودار اول تنها شامل یک شکستگی است.



شکل ۲. نمودار اول b

این نمودار صرفاً یک شکستگی دارد پس با استفاده از شبکه های مشخص شده در تصویر (۱) قابل یادگیری است. نمودار دوم شامل سه شکستگی است.

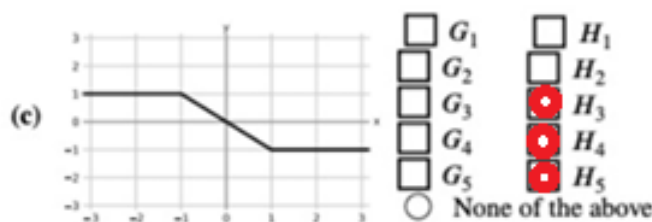


شکل ۳. نمودار دوم b

تمامی شبکه‌هایی که داریم نهایتاً می‌توانند دو تا شکستگی را یاد بگیرند پس با استفاده از این شبکه‌ها نمی‌توان این نمودار را یاد گرفت.

۳-۲. نمودار c

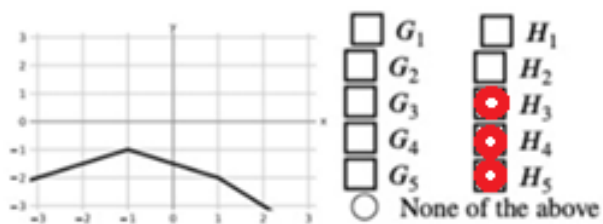
نمودار اول دو شکستگی دارد.



شکل ۴. نمودار اول c

از آنجا که دو شکستگی داریم، تنها با استفاده از شبکه‌های مشخص شده می‌توان در شکل (۳) می‌توان آن را یاد گرفت.

نمودار دوم دو شکستگی دارد.



شکل ۵. نمودار دوم c

از آنجا که دو شکستگی داریم، تنها با استفاده از شبکه‌های مشخص شده می‌توان در شکل (۳) می‌توان آن را یاد گرفت.

پرسش ۳ – Neural Networks: Forward and Backward Propagation

۱-۳. ابعاد هر پارامتر و وکتورایز کردن شبکه

طبق صورت سوال می دانیم که ابعاد z_1 به صورت $1 \times D_{a_1}$ است. بر همین اساس طبق $z_1 = W_1 x^{(i)} + b_1$ می فهمیم که ابعاد $W_1 x^{(i)}$ باید برابر $1 \times D_{a_1}$ باشد و از آنجا که ابعاد $x^{(i)}$ برابر $1 \times D_x$ است پس ابعاد W_1 باید $D_x \times D_{a_1}$ باشد. در این صورت ابعاد در ضرب ماتریسی همخوانی خواهند داشت. بایاس این رابطه هم باید مطابق خروجی باشد در نتیجه b_1 بعد $1 \times D_{a_1}$ را دارد.

از آنجا که تابع $ReLU$ به صورت elementwise عمل می کند پس ابعاد a_1 هم $1 \times D_{a_1}$ می باشد.

برای پارامتر های بعدی ابتدا باید توجه کنیم که $\hat{y}^{(i)}$ یک اسکالر است (زیرا خروجی یک تابع لاجیستیک است) و همچنین تابع $loss$ هم binary cross-Entropy می باشد. بنابراین تعداد نرون های خروجی صرفا یکی است. پس z_2 باید 1×1 باشد. پس در رابطه $z_2 = W_2 a_1 + b_2$ پارامتر b_2 و حاصل ضرب $W_2 a_1$ باید 1×1 باشد پس از آنجا که a_1 ابعاد $1 \times D_{a_1}$ دارد پس W_2 باید $1 \times D_{a_1}$ باشد.

حالا می خواهیم شبکه را وکتورایز کنیم. ابتدا ورودی های x و خروجی های y را بررسی می کنیم. در صورت سوال دیدیم که $x^{(i)}$ ابعاد $1 \times D_x$ دارد. حال ماتریس X را به صورت زیر تعریف می کنیم.

$$X = [x^{(1)}, x^{(2)}, \dots, x^{(m)}] \quad (۸)$$

یعنی هر داده ی $x^{(i)}$ را به صورت ستونی در یک ماتریس چیدیم. ابعاد X برابر $D_x \times m$ می شود. (Y هم در انتها بررسی خواهد شد).

در ادامه هر رابطه را بازنویسی کرده و ابعاد را مشخص می کنیم. در رابطه $Z_1 = W_1 X + B_1$ جمله ی $W_1 X$ و پارامتر Z_1 ابعاد $D_{a_1} \times m$ دارند و به صورت زیر هستند.

$$W_1 X = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ W_1 x^{(1)} & W_1 x^{(2)} & \dots & W_1 x^{(m)} \\ \vdots & \vdots & \dots & \vdots \end{bmatrix} \quad (۹)$$

$$Z_1 = [z_1^{(1)} \quad z_1^{(2)} \quad \dots \quad z_1^{(m)}] \quad (10)$$

در واقع ضرایب W_1 همان ابعاد قبلی یعنی $D_{a_1} \times D_x$ را دارد و به همین شکل در تک تک ستون های X ضرب می شود. همچنین بایاس برای تمامی نمونه ها یکسان است. پس B_1 به صورت زیر است.

$$B_1 = [b_1 \quad b_1 \quad \dots \quad b_1] \quad (11)$$

ابعاد B_1 باید برابر با جمله ی اول باشد، در نتیجه برابر $D_{a_1} \times m$ است. البته توجه کنید این یک راه استاندارد برای نمایش B_1 است. در متلب و کتابخانه های پایتون مثل *numpy* برای جلوگیری از مصرف بی دلیل رم و افزایش بهینگی از روشی به اسم *broadcasting* استفاده می شود. یعنی به جای اینکه b_1 را در یک ماتریس تکرار کنیم، صرفاً خود b_1 را به کار گرفته و در محاسبات *numpy* مشکل ابعاد را به صورت خودکار حل می کند. یعنی تقریباً به صورت زیر عمل می کند. (در مورد *broadcasting* در *numpy* اینجا [۱] بخوانید)

$$\mathbf{1}_m = [1 \quad 1 \quad \dots \quad 1] \quad (12)$$

$$B_1 = b_1 \mathbf{1}_m \quad (8)$$

در این حالت ابعاد B_1 همان $D_{a_1} \times m$ و ابعاد b_1 همچنان $D_{a_1} \times 1$ است.

اتفاق مشابهی در مورد رابطه ی $Z_2 = W_2 A_1 + B_2$ می افتد. یعنی:

$$A_1 = [a_1^{(1)} \quad a_1^{(2)} \quad \dots \quad a_1^{(m)}] \quad , \quad A_1 \in \mathbb{R}^{D_{a_1} \times m} \quad (9)$$

جمله ی $W_2 A_1$ ابعاد $1 \times m$ دارد که به صورت زیر است.

$$W_2 A_1 = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ W_2 a_1^{(1)} & W_2 a_1^{(2)} & \dots & W_2 a_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad (10)$$

در واقع ضرایب W_2 همان ابعاد قبلی یعنی $1 \times D_{a_1}$ را دارد و به همین شکل در تک تک ستون های A_1 ضرب می شود. همچنین بایاس برای تمامی نمونه ها یکسان است. پس B_2 به صورت زیر است.

$$B_2 = [b_2 \quad b_2 \quad \dots \quad b_2] \quad (11)$$

ابعاد B_2 باید برابر با جمله ی اول باشد، در نتیجه برابر $1 \times m$ است و ابعاد b_2 همچنان 1×1 است. (می توان مشابه روابط (۱۲) و (۱۳) نیز B_2 را نشان داد.) نهایتاً Z_2 و Y به صورت زیر هستند.

$$Z_2 = [z_2^{(1)} \quad z_2^{(2)} \quad \dots \quad z_2^{(m)}] \quad , \quad Z_2 \in \mathbb{R}^{1 \times m} \quad (12)$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] = \sigma(Z_2) \quad (13)$$

$$= \left[\sigma(z_2^{(1)}) \quad \sigma(z_2^{(2)}) \quad \dots \quad \sigma(z_2^{(m)}) \right] \quad , \quad Y \in \mathbb{R}^{1 \times m}$$

بر این اساس Y هم $m \times 1$ خواهد بود؛ زیرا $\hat{y}^{(i)}$ یک اسکالر است.

۲-۳. محاسبه $\frac{\partial J}{\partial \hat{y}^{(i)}}$ و $\frac{\partial J}{\partial \hat{y}}$

طبق صورت سوال میدانیم که:

$$L^{(i)} = y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \quad (14)$$

و تابع هزینه به صورت زیر است.

$$J = -\frac{1}{m} \sum_{i=1}^m L^{(i)} \quad (20)$$

واضح است که صرفاً جمله‌ی i -ام به $\hat{y}^{(i)}$ وابسته بوده و بقیه کنار می‌روند.

$$\frac{\partial J}{\partial \hat{y}^{(i)}} = -\frac{1}{m} \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} \quad (21)$$

حال از $L^{(i)}$ نسبت به $\hat{y}^{(i)}$ مشتق می‌گیریم.

$$\begin{aligned} \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} &= \frac{\partial}{\partial \hat{y}^{(i)}} (y^{(i)} \log \hat{y}^{(i)}) + \frac{\partial}{\partial \hat{y}^{(i)}} ((1 - y^{(i)}) \log(1 - \hat{y}^{(i)})) \\ &= \frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{1 - y^{(i)}}{1 - \hat{y}^{(i)}} \end{aligned} \quad (15)$$

پس خواهیم داشت:

$$\delta_1^{(i)} = \frac{\partial J}{\partial \hat{y}^{(i)}} = -\frac{1}{m} \left(\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{1 - y^{(i)}}{1 - \hat{y}^{(i)}} \right) \quad (16)$$

حال می‌خواهیم $\frac{\partial J}{\partial \hat{y}}$ را، یعنی مشتق نسبت به معادل وکتورایز شده محاسبه کنیم.

$$\frac{\partial J}{\partial \hat{y}} = \left[\frac{\partial J}{\partial \hat{y}^{(1)}} \quad \frac{\partial J}{\partial \hat{y}^{(2)}} \quad \dots \quad \frac{\partial J}{\partial \hat{y}^{(m)}} \right] = [\delta^{(1)} \quad \delta^{(2)} \quad \dots \quad \delta^{(m)}] \quad (17)$$

با توجه به آنچه در محاسبه‌ی $\frac{\partial J}{\partial \hat{y}^{(i)}}$ دیدیم، به صورت زیر خواهد بود.

$$\frac{\partial J}{\partial \hat{y}} = -\frac{1}{m} \left(\frac{y}{\hat{y}} - \frac{1 - y}{1 - \hat{y}} \right) \quad (18)$$

۳-۳. محاسبه $\frac{\partial \hat{y}^{(i)}}{\partial z_2^{(i)}}$

دیدیم که $\hat{y}^{(i)} = \sigma(z_r)$ و $\sigma(z)$ تابع سیگموئید است. ابتدا باید مشتق تابع سیگموئید حساب شود. مشتق تابع سیگموئید معروف است اما با این حال از ابتدا محاسبه می کنیم. تابع سیگموئید:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (19)$$

محاسبه‌ی مشتق:

$$\frac{d\sigma(z)}{dz} = -(1 + e^{-z})^{-2} \cdot (-e^{-z}) = \frac{e^{-z}}{(1 + e^{-z})^2} \quad (20)$$

می توان به صورت زیر نیز این مشتق را بازنویسی کرد که مشابه فرمول معروف مشتق تابع سیگموئید است.

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z)) = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right) \quad (21)$$

حال از $\hat{y}^{(i)}$ و $z_r^{(i)}$ استفاده می کنیم.

$$\frac{\partial \hat{y}^{(i)}}{\partial z_r^{(i)}} = \sigma(z_r^{(i)}) (1 - \sigma(z_r^{(i)})) \quad (22)$$

از آنجاکه $\hat{y}^{(i)} = \sigma(z_r^{(i)})$ پس به صورت زیر بازنویسی می کنیم.

$$\delta_r^{(i)} = \frac{\partial \hat{y}^{(i)}}{\partial z_r^{(i)}} = \hat{y}^{(i)} (1 - \hat{y}^{(i)}) \quad (30)$$

۳-۴. محاسبه‌ی $\frac{\partial z_2}{\partial a_1}$

با توجه به $z_r = W_r a_1 + b_r$ مشتق آن به صورت زیر است.

$$\delta_r^{(i)} = \frac{\partial z_r}{\partial a_1} = W_r \quad (31)$$

۳-۵. محاسبه‌ی $\frac{\partial a_1}{\partial z_1}$

می دانیم:

$$a_1 = \text{ReLU}(z_1) = \begin{cases} z_1 & z_1 > 0 \\ 0 & z_1 \leq 0 \end{cases} \quad (32)$$

می دانیم که ابعاد z_1 برابر $1 \times D_{a_1}$ است. بر این اساس می توان به صورت زیر هم آن را نوشت.

$$(k = 1, \dots, D_{a_1})$$

$$ReLU(z_{\backslash,k}) = \begin{cases} z_{\backslash,k} & z_{\backslash,k} > 0 \\ 0 & z_{\backslash,k} \leq 0 \end{cases} \quad (33)$$

پس مشتق a_{\backslash} نسبت به z_{\backslash} به صورت *element-wise* به این شکل خواهد بود.

$$\frac{\partial a_{\backslash,k}}{\partial z_{\backslash,k}} = \begin{cases} 1 & z_{\backslash,k} > 0 \\ 0 & z_{\backslash,k} \leq 0 \end{cases} \quad (34)$$

اگر بخواهیم این مشتق را به صورت برداری/ ماتریسی بنویسیم خواهیم داشت:

$$\delta_{\backslash}^{(i)} = \frac{\partial a_{\backslash}}{\partial z_{\backslash}} = \text{diag}(\mathbb{1}_{z_{\backslash}^{(i)} > 0}) \quad (35)$$

ولی در حالت کلی به این شکل است.

$$\delta_{\backslash}^{(i)} = \frac{\partial a_{\backslash}}{\partial z_{\backslash}} = \frac{\partial}{\partial z_{\backslash}} ReLU(z_{\backslash}) = \mathbb{1}_{z_{\backslash} > 0}. \quad (36)$$

۳-۶. محاسبه $\frac{\partial z_1}{\partial W_1}$

ابتدا توجه کنید که W_1 یک ماتریس و z_1 یک بردار است. پاسخ این مشتق سه بعد خواهد داشت یعنی انتظار یک تانسور داریم. با یک مثال می توان توضیح داد؛ خروجی k -ام و یک وزن رندوم مثلا $W_{\backslash,ij}$ را در نظر بگیرید. (k اندیس هر درایه از z_1 و i, j اندیس های هر درایه از ماتریس W_1 هستند).

مقدار $\frac{\partial z_{\backslash,k}}{\partial W_{\backslash,ij}}$ تنها در یک حالت شامل وزن $W_{\backslash,ij}$ خواهد بود و آن هم در صورتی است که $k = i$ باشد در غیر این صورت مشتق صفر خواهد بود. پس در کل:

اگر $k = i$ ، مشتق $W_{\backslash,kj} x_j$ نسبت به $W_{\backslash,kj}$ برابر x_j خواهد بود و اگر $k \neq i$ ، مشتق صفر است. به همین دلیل است که این مشتق یک تانسور خواهد بود.

$$\delta_{\Delta}^{(i)} = \frac{\partial z_{\backslash,k}}{\partial W_{\backslash,ij}} = \begin{cases} x_j & k = i \\ 0 & k \neq i \end{cases} \quad (37)$$

این مشتق به همین شکل حساب شد؛ اما می توان آن را راحت تر و واضح تر با استفاده از تابع دلتای کرونکر نشان داد. این تابع به صورت زیر است:

$$\delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad (38)$$

در مورد این تابع و خواص آن می توانید در [۲] بخوانید. با استفاده از این تابع می توان مشتق را به صورت زیر بازنویسی کرد.

$$\delta_{\delta}^{(i)} = \left(\frac{\partial z_1}{\partial W_1} \right)_{kij} = \delta_{ki} x_j \quad (39)$$

توجه کنید؛ شاید در ابتدا اینجور به نظر برسد که صرفاً در نظر گرفتن حالت $k = i$ کافی است و اگر بر این اساس پیش برویم به یک ماتریس (۲-D) برسیم. ولی چنین کاری اشتباه است. می‌توان بر اساس مفهوم مشتق آن را توضیح داد. وقتی می‌خواهیم $\frac{\partial z_1}{\partial W_1}$ را حساب کنیم یعنی می‌خواهیم ببینیم تغییر هر درایه از W_1 چه اثری روی هر عنصر از z_1 خواهد داشت، اگر صرفاً $k = i$ را در نظر بگیریم یعنی تاثیر تعداد زیادی از درایه‌های W_1 را در نظر نگرفته و بررسی نکرده ایم.

۷-۳. محاسبه‌ی $\frac{\partial J}{\partial W_1}$

با توجه به بخش‌های قبل به راحتی می‌توان این مشتق را به صورت زیر حساب کرد.

$$\frac{\partial J}{\partial W_1} = \sum_{i=1}^m \frac{\partial J}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z_r^{(i)}} \frac{\partial z_r^{(i)}}{\partial a_1^{(i)}} \frac{\partial a_1^{(i)}}{\partial z_1^{(i)}} \frac{\partial z_1^{(i)}}{\partial W_1} = \sum_{i=1}^m \delta_1^{(i)} \delta_r^{(i)} \delta_r^{(i)} \delta_r^{(i)} \delta_{\delta}^{(i)} \quad (23)$$

از آنجا که هر کدام از این مشتق‌ها طولانی بوده و نوشتن اصل روابط صرفاً باعث طولانی شدن رابطه شده و کمک خاصی به درک مطلب نمی‌کرد، صرفاً با استفاده علامت گذاری این مشتق نوشته شد.

یک راه دیگر هم برای محاسبه‌ی این مشتق وجود دارد. با توجه به J مشتق آن به صورت زیر است.

$$\frac{\partial J}{\partial W_1} = - \frac{1}{m} \sum_{i=1}^m \frac{\partial L^{(i)}}{\partial W_1} \quad (41)$$

با استفاده از قاعده‌ی زنجیره‌ای می‌توان این مشتق را به صورت زیر نوشت.

$$\frac{\partial L^{(i)}}{\partial W_1} = \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z_r^{(i)}} \frac{\partial z_r^{(i)}}{\partial a_1^{(i)}} \frac{\partial a_1^{(i)}}{\partial z_1^{(i)}} \frac{\partial z_1^{(i)}}{\partial W_1} \quad (42)$$

توجه کنید که در این مشتق زنجیره‌ای از $\frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}}$ استفاده شده و در بخش‌های قبلی این مشتق را حساب نکردیم. به همین دلیل صرفاً به رابطه‌ی (۴۰) اکتفا کرده و روش دوم محاسبه را ادامه نمی‌دهیم.

پرسش ۴ – Hessian of a Two-Layer Network

در این سوال می خواهیم ماتریس هسین یک شبکه دو لایه را محاسبه کنیم. روابط به صورت زیر هستند.

$$a_j = \sum_i w_{ji}^{(1)} x_i, \quad z_j = h(a_j) \quad (43)$$

$$a_k = \sum_j w_{kj}^{(2)} z_j$$

و E_n خطای داده‌ی n است. بر این اساس مشتق‌های E_n را با علائم زیر نشان می دهیم.

$$\delta_k \equiv \frac{\partial E_n}{\partial a_k}, \quad M_{kk'} \equiv \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} \quad (44)$$

توجه کنید که در این اثبات از i و i' برای مشخص کردن ورودی‌ها، j و j' برای مشخص کردن لایه‌های مخفی و از k و k' برای مشخص کردن خروجی‌ها استفاده شده است.

۴-۱. هر دو وزن در لایه دوم

باید مشتق دوم زیر را حساب کرد که برای محاسبه آن از رابطه‌ی مشتق زنجیره‌ای استفاده می کنیم.

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(2)} \partial w_{k'j'}^{(2)}} = \frac{\partial}{\partial w_{kj}^{(2)}} \left(\frac{\partial E_n}{\partial w_{k'j'}^{(2)}} \right) \quad (45)$$

با محاسبه مشتق داخلی شروع می کنیم.

$$\frac{\partial E_n}{\partial w_{k'j'}^{(2)}} = \frac{\partial E_n}{\partial a_{k'}} * \frac{\partial a_{k'}}{\partial w_{k'j'}^{(2)}} = \delta_{k'} * z_{j'} \quad (46)$$

حالا مشتق دوم این عبارت را محاسبه می کنیم.

$$\frac{\partial}{\partial w_{kj}^{(r)}} \left(\frac{\partial E_n}{\partial w_{k'j'}^{(r)}} \right) = \frac{\partial(\delta_{k'} * z_{j'})}{\partial w_{kj}^{(r)}} = \frac{\partial(\delta_{k'} * z_{j'})}{\partial a_k} * \frac{\partial a_k}{\partial w_{kj}^{(r)}} \quad (47)$$

$$= z_j * \frac{\partial(\delta_{k'} * z_{j'})}{\partial a_k} = z_j \left[z_{j'} * \frac{\partial \delta_{k'}}{\partial a_k} + \delta_{k'} * \frac{\partial z_{j'}}{\partial a_k} \right]$$

عبارت $\frac{\partial z_{j'}}{\partial a_k}$ برابر صفر است. پس

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(r)} \partial w_{k'j'}^{(r)}} = z_j z_{j'} M_{kk'} \quad (48)$$

به این صورت مشتق ها با استفاده از قاعده‌ی زنجیره ای حساب شد.

۲-۴. هر دو وزن در لایه اول

باید مشتق دوم زیر را حساب کرد که برای محاسبه آن از رابطه‌ی مشتق زنجیره ای استفاده می کنیم.

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} = \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{j'i'}^{(1)}} \right) \quad (49)$$

با محاسبه مشتق داخلی شروع می کنیم. می توانیم این کار را با تقسیم به دو حالت $j = j'$ و $j \neq j'$ انجام دهیم. در این صورت اثبات واضح تر خواهد بود.

فرض کنید $j \neq j'$:

$$\begin{aligned} \frac{\partial E_n}{\partial w_{j'i'}^{(1)}} &= \sum_{k'} \frac{\partial E_n}{\partial a_{k'}} \frac{\partial a_{k'}}{\partial w_{j'i'}^{(1)}} = \sum_{k'} \frac{\partial E_n}{\partial a_{k'}} \frac{\partial a_{k'}}{\partial z_{j'}} \frac{\partial z_{j'}}{\partial a_{j'}} \frac{\partial a_{j'}}{\partial w_{j'i'}^{(1)}} = \\ &= \sum_{k'} \frac{\partial E_n}{\partial a_{k'}} w_{k'j'}^{(r)} h'(a_{j'}) x_{i'} \end{aligned} \quad (50)$$

عبارات $w_{k'j'}^{(r)}$ به سادگی از وابستگی $a_{k'}$ به وزن های لایه‌ی دوم $(w_{k'j'}^{(r)})$ ، رابطه‌ی z_j و نهایتاً a_j با استفاده از مشتق زنجیره ای می آید. حالا مشتق دوم این عبارت را محاسبه می کنیم.

$$\begin{aligned}
\frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{j'i'}^{(1)}} \right) &= \frac{\partial}{\partial w_{ji}^{(1)}} \left(\sum_{k'} \frac{\partial E_n}{\partial a_{k'}} w_{k'j'}^{(r)} h'(a_{j'}) x_{i'} \right) \quad (51) \\
&= \sum_{k'} h'(a_{j'}) x_{i'} \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j'}^{(2)} \right) \\
&= \sum_{k'} h'(a_{j'}) x_{i'} \sum_k \frac{\partial}{\partial a_k} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j'}^{(2)} \right) \frac{\partial a_k}{\partial w_{ji}^{(1)}} \\
&= \sum_{k'} h'(a_{j'}) x_{i'} \sum_k \frac{\partial}{\partial a_k} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j'}^{(2)} \right) (w_{kj}^{(2)} h'(a_j) x_i) \\
&= \sum_{k'} h'(a_{j'}) x_{i'} \sum_k M_{kk'} w_{k'j'}^{(r)} (w_{kj}^{(r)} h'(a_j) x_i)
\end{aligned}$$

حال با خارج کردن پارامترهایی که به سامیشتن و ایندکس آن ارتباطی ندارند می توانیم ساده سازی کنیم.

$$\frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{j'i'}^{(1)}} \right) = x_{i'} x_i h'(a_{j'}) h'(a_j) \sum_{k'} \sum_k w_{k'j'}^{(r)} w_{kj}^{(r)} M_{kk'} \quad (52)$$

حالا حالتی را فرض کنید که $j = j'$ ، فرایند تا حدودی شبیه به حالت قبلی خواهد بود. مشتق داخلی در این حالت به صورت زیر است.

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \sum_{k'} \frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} h'(a_j) x_{i'} \quad (24)$$

حالا مشتق دوم این عبارت را محاسبه می کنیم.

$$\begin{aligned}
\frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{j'i'}^{(1)}} \right) &= \frac{\partial}{\partial w_{ji}^{(1)}} \left(\sum_{k'} \frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} h'(a_j) x_{i'} \right) = \sum_{k'} \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} h'(a_j) x_{i'} \right) \quad (25) \\
&= x_{i'} \sum_{k'} \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} \right) h'(a_j) + x_{i'} \sum_{k'} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} \right) \frac{\partial h'(a_j)}{\partial w_{ji}^{(1)}} \\
&= x_{i'} x_i h'(a_j) h'(a_j) \sum_{k'} \sum_k w_{k'j}^{(r)} w_{kj}^{(r)} M_{kk'} + x_{i'} \sum_{k'} \left(\frac{\partial E_n}{\partial a_{k'}} w_{k'j}^{(r)} \right) \frac{\partial h'(a_j)}{\partial w_{ji}^{(1)}} \\
&= x_{i'} x_i h'(a_j) h'(a_j) \sum_{k'} \sum_k w_{k'j}^{(r)} w_{kj}^{(r)} M_{kk'} + h''(a_j) x_i x_{i'} \sum_{k'} \delta_{k'} w_{k'j}^{(r)}
\end{aligned}$$

(توجه کنید برای اینکه هر مساوی در یک خط قرار گرفته و روابط شکل قابل درک تری داشته باشند مجبور شدم سائز فونت را کوچک تر کنم. لطفا زوم کنید.)

حالا که در هر دو حالت این مشتق را بررسی کردیم، با جمع روابط (52) و (54) به رابطه ی زیر خواهیم رسید.

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} = x_{i'} x_i \left(h'(a_{j'}) h'(a_j) \sum_{k'} \sum_k w_{k'j'}^{(2)} w_{kj}^{(2)} M_{kk'} + h''(a_j) I_{jj'} \sum_{k'} \delta_{k'} w_{k'j}^{(2)} \right) \quad (26)$$

به این صورت مشتق ها با استفاده از قاعده‌ی زنجیره ای حساب شد.

۳-۴. هر وزن از یکی از لایه ها

باید مشتق دوم زیر را حساب کرد که برای محاسبه آن از رابطه‌ی مشتق زنجیره ای استفاده می کنیم.

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{kj'}^{(2)}} = \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{kj'}^{(2)}} \right) \quad (56)$$

با محاسبه مشتق داخلی شروع می کنیم. به طور مشابه قسمت قبل با تقسیم به دو حالت $j \neq j'$ و $j = j'$ پیش می رویم. در این صورت اثبات واضح تر خواهد بود. با حالت $j \neq j'$ شروع می کنیم.

$$\frac{\partial E_n}{\partial w_{kj'}^{(2)}} = \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial w_{kj'}^{(2)}} = \delta_k z_{j'} \quad (57)$$

حالا مشتق این عبارت را نسبت به وزن های لایه اول حساب می کنیم.

$$\begin{aligned} \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{kj'}^{(2)}} \right) &= \frac{\partial (\delta_k z_{j'})}{\partial w_{ji}^{(1)}} = z_{j'} \sum_{k'} \frac{\partial \delta_k}{\partial a_{k'}} \frac{\partial a_{k'}}{\partial w_{ji}^{(1)}} \\ &= z_{j'} \sum_{k'} M_{kk'} \left(w_{k'j}^{(2)} h'(a_j) x_i \right) \\ &= x_i h'(a_j) z_{j'} \sum_{k'} M_{kk'} w_{k'j}^{(2)} \end{aligned} \quad (58)$$

حالا حالتی را فرض کنید که $j = j'$ ، فرایند تا حدودی شبیه به حالت قبلی خواهد بود. مشتق داخلی در این حالت به صورت زیر است.

$$\frac{\partial E_n}{\partial w_{kj}^{(2)}} = \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial w_{kj}^{(2)}} = \delta_k z_j \quad (59)$$

حالا مشتق این عبارت را نسبت به وزن های لایه اول حساب می کنیم.

$$\begin{aligned} \frac{\partial}{\partial w_{ji}^{(1)}} \left(\frac{\partial E_n}{\partial w_{kj}^{(2)}} \right) &= \frac{\partial (\delta_k z_j)}{\partial w_{ji}^{(1)}} = z_j \frac{\partial \delta_k}{\partial w_{ji}^{(1)}} + \delta_k \frac{\partial z_j}{\partial w_{ji}^{(1)}} \\ &= x_i h'(a_j) z_j \sum_{k'} M_{kk'} w_{k'j}^{(2)} + \delta_k h'(a_j) x_i \end{aligned} \quad (60)$$

حال که هر دو حالت را بررسی کردیم، با جمع روابط (58) و (60) به رابطه‌ی زیر خواهیم رسید.

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{kj}^{(2)}} = x_i h'(a_j) \left(\delta_k I_{jj'} + z_{j'} \sum_{k'} w_{k'j}^{(2)} M_{kk'} \right) \quad (61)$$

به این صورت تمامی درایه‌های ماتریس هسین یک شبکه دو لایه محاسبه شد.

پرسش ۵ – Linear Discriminant Analysis (LDA)

در این سوال، نمایشی دیگری از تابع هزینه‌ی LDA را اثبات خواهیم کرد.

فرض شده که دو دیتاست به صورت زیر داریم.

$$Y_1 = \{y_i\}_{i=1}^{n_1}, \quad Y_2 = \{y_j\}_{j=1}^{n_2} \quad (62)$$

حال از رابطه‌ی اصلی شروع می‌کنیم.

$$\begin{aligned} J &= \frac{1}{n_1 n_2} \sum_{y_i} \sum_{y_j} (y_i - y_j)^2 = \frac{1}{n_1 n_2} \sum_{y_i} \sum_{y_j} (y_i^2 + y_j^2 - 2y_i y_j) \\ &= \frac{1}{n_1 n_2} \left[\sum_{y_i} \sum_{y_j} y_i^2 + \sum_{y_i} \sum_{y_j} y_j^2 - 2 \sum_{y_i} \sum_{y_j} y_i y_j \right] \end{aligned} \quad (63)$$

در رابطه‌ی (63)، جملات $\sum_{y_i} \sum_{y_j} y_i^2$ و $\sum_{y_i} \sum_{y_j} y_j^2$ به دلیل جمع شدن روی اندیسی که در جمع وجود ندارد، به ترتیب n_2 بار و n_1 بار تکرار می‌شوند. همچنین $\sum_{y_i} \sum_{y_j} y_i y_j$ برابر حاصل ضرب میانگین هاست. پس بر این اساس رابطه‌ی (63) به صورت زیر در خواهد آمد.

$$J = \frac{1}{n_1 n_2} \left[n_2 \sum_{y_i} y_i^2 + n_1 \sum_{y_j} y_j^2 - 2m_1 m_2 \right] = \frac{1}{n_1} \sum_{y_i} y_i^2 + \frac{1}{n_2} \sum_{y_j} y_j^2 - 2m_1 m_2 \quad (27)$$

حالا مقادیر \bar{m}_1^2 و \bar{m}_2^2 را به رابطه‌ی بدست آمده اضافه و کم می کنیم.

$$\begin{aligned} J &= \frac{1}{n_1} \sum_{y_i} y_i^2 + \frac{1}{n_2} \sum_{y_j} y_j^2 - \bar{m}_1 \bar{m}_2 + \bar{m}_1^2 + \bar{m}_2^2 - \bar{m}_1^2 \quad (65) \\ &= (m_1 - m_2)^2 + \frac{1}{n_1} \sum_{y_i} y_i^2 + \frac{1}{n_2} \sum_{y_j} y_j^2 + m_1^2 + m_2^2 \\ &\quad - \bar{m}_1^2 - \bar{m}_2^2 \end{aligned}$$

توجه کنید هر میانگین نمونه‌ای به صورت $m = \frac{1}{n} \sum_i y_i$ بدست آمده است. پس روابط \bar{m}_1^2 و \bar{m}_2^2 را باز می کنیم و با استفاده از تعریف میانگین نمونه ای بازنویسی می کنیم.

$$\begin{aligned} J &= (m_1 - m_2)^2 + \frac{1}{n_1} \sum_{y_i} y_i^2 + \frac{1}{n_2} \sum_{y_j} y_j^2 + m_1^2 + m_2^2 - \bar{m}_1 \frac{1}{n_1} \sum_{y_i} y_i \\ &\quad - \bar{m}_2 \frac{1}{n_2} \sum_{y_j} y_j \quad (66) \end{aligned}$$

حال می توانیم s_1^2 و s_2^2 را بسازیم.

$$\begin{aligned} J &= (m_1 - m_2)^2 + \frac{1}{n_1} \sum_{y_i} (y_i^2 + m_1^2 - \bar{m}_1 y_i) + \frac{1}{n_2} \sum_{y_j} (y_j^2 + m_2^2 - \bar{m}_2 y_j) \quad (67) \\ &= (m_1 - m_2)^2 + \frac{1}{n_1} \sum_{y_i} (y_i - m_1)^2 + \frac{1}{n_2} \sum_{y_j} (y_j - m_2)^2 \\ &= (m_1 - m_2)^2 + \frac{1}{n_1} s_1^2 + \frac{1}{n_2} s_2^2 \end{aligned}$$

رابطه‌ی مورد نظر اثبات شد.

پرسی ۶ – Feature Conditioning and PCA

در این سوال با استفاده از روش PCA یکی از روابط پایه ای و البته مهم را اثبات می کنیم.

فرض کنید دیتاست $\{x_n\}_{n=1}^N$ را داشته باشیم که در آن هر x_n یک بردار D بعدی است و می خواهیم با استفاده از روش PCA این داده ها را به فضای یک بعدی ببریم. بردار w را در این فضای D بعدی فرض کنید که جهت داده ها را در این فضا مشخص می کند. از آنجا که صرفاً جهت این بردار برای ما اهمیت دارد و نه اندازه ای آن، پس فرض می کنیم این بردار واحد باشد. ($w^T w = 1$) این فرض بدون از دست رفتن کلیت مسئله بوده و تاثیری بر آنچه می خواهیم نتیجه بگیریم ندارد.

هر داده ای x_n با استفاده از بردار w به یک نقطه به صورت $w^T x_n$ تبدیل می شود. میانگین داده های پراجکت شده (تبدیل شده) برابر $w^T \bar{X}$ است که در آن \bar{X} میانگین نمونه ای به صورت زیر است.

$$\bar{X} = \frac{1}{N} \sum_{n=1}^N x_n \quad (68)$$

واریانس این داده های پراجکت شده به صورت زیر است.

$$\frac{1}{N} \sum_{n=1}^N (w^T x_n - w^T \bar{X})^2 = w^T S w \quad (69)$$

در رابطه‌ی (۶۹) پارامتر S ماتریس کوواریانس داده هاست.

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{X})(x_n - \bar{X})^T \quad (70)$$

حال یک مسئله بهینه سازی داریم که باید واریانس $w^T S w$ را ماکزیمم کنیم. اما باید توجه داشت برای اینکه در این ماکزیمم سازی اندازه w به سمت بینهایت میل نکند، $(\|w\| \rightarrow \infty)$ ، باید فرض کنیم که بردار مورد نظر همان طور که قبلا هم اشاره شد حتما واحد باشد. $(w^T w = 1)$

حال با این شرط می توان با استفاده از روش لاگرانژ این ماکزیمم سازی را انجام داد. ضریب λ_1 را برای شرط در نظر می گیریم.

$$\mathcal{L} = w^T S w + \lambda_1 (1 - w^T w) \quad (71)$$

حال نسبت به w از این رابطه مشتق می گیریم و مساوی صفر قرار می دهیم.

$$S w - \lambda_1 w = 0 \quad (72)$$

و خواهیم داشت:

$$S w = \lambda_1 w \quad (73)$$

رابطه (۷۳) دقیقا همان معادله‌ی مقدار ویژه و بردار ویژه است.

پرسش ۷ – Feature Conditioning

برای پیاده سازی کد این سوال از [۳] و [۴] استفاده کردم. الگوریتم های انتخاب پیشرونده و حذف پسرونده اجرا شد و نتیجه هر کدام به صورت زیر است.

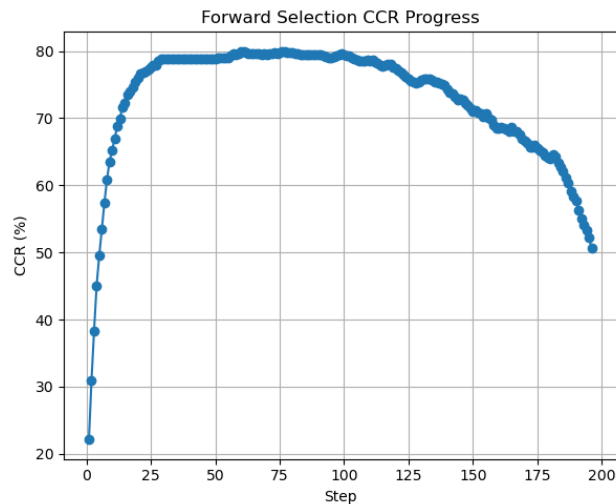
انتخاب پیشرونده:

===== BEST FORWARD SELECTION RESULT =====

Features: [۳۵, ۱۰۲, ۹۱, ۱۰۶, ۱۳۱, ۸۹, ۴۹, ۱۰۴, ۴۷, ۶۵, ۱۳۲, ۶۳, ۷۴, ۱۳۳, ۵۲, ۱۶۲, ۱۳۵, ۷۶, ۷۷, ۷۸, ۱۴۴, ۱۴۹, ۱۱۶, ۶۱, ۱۱۸, ۶۴, ۹۰, ۷۵, ۱۴۸, ۰, ۱, ۲, ۳, ۱۲, ۱۳, ۱۴, ۱۵, ۲۷, ۲۸, ۸۴, ۹۸, ۱۵۴, ۱۶۷, ۱۶۸, ۱۸۱, ۱۸۲, ۱۸۳, ۱۹۴, ۱۹۵, ۱۲۰, ۵۱, ۱۳۴, ۱۴۷, ۸۰, ۱۰۱, ۱۰۳, ۴۶, ۶۰, ۷۹, ۴۸, ۱۴۶]

Best CCR: ۷۹,۹۶٪

نمودار تغییر CCR نسبت به پیشروی الگوریتم (هر مرحله) به صورت زیر است.



شکل ۶. نمودار CCR الگوریتم انتخاب پیشرونده

همانگونه که انتظار داشتیم این نمودار تا جایی به صورت صعودی بوده و پس از آن انتخاب ویژگی های بیشتر باعث کاهش CCR شده است.

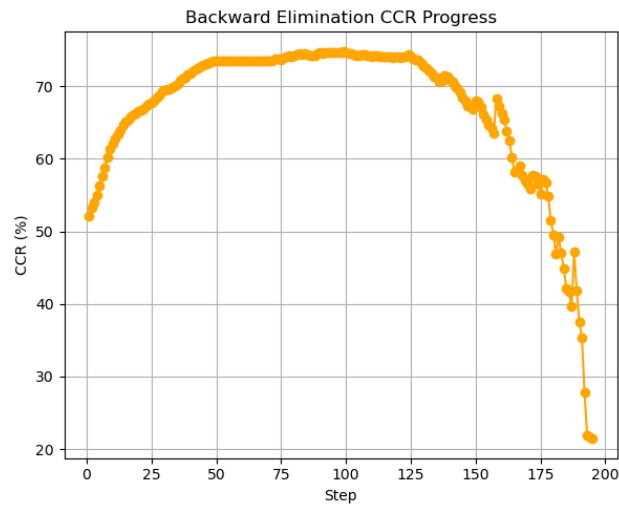
حذف پسرونده:

===== BEST BACKWARD ELIMINATION RESULT =====

Features: [۷, ۹, ۱۱, ۱۸, ۲۰, ۲۲, ۲۳, ۲۵, ۲۶, ۳۱, ۳۴, ۳۶, ۳۷, ۳۸, ۳۹, ۴۰, ۴۴, ۴۵, ۴۶, ۴۷, ۴۸, ۴۹, ۵۰, ۵۲, ۵۳, ۵۵, ۶۰, ۶۱, ۶۳, ۶۴, ۶۶, ۶۸, ۶۹, ۷۰, ۷۱, ۷۴, ۷۵, ۷۷, ۷۹, ۸۰, ۸۲, ۸۷, ۸۸, ۸۹, ۹۰, ۹۱, ۹۲, ۹۳, ۹۵, ۹۶, ۹۷, ۱۰۱, ۱۰۲, ۱۰۳, ۱۰۴, ۱۰۵, ۱۰۶, ۱۱۰, ۱۱۶, ۱۱۸, ۱۲۰, ۱۲۱, ۱۲۶, ۱۳۰, ۱۳۱, ۱۳۲, ۱۳۴, ۱۳۵, ۱۳۶, ۱۳۸, ۱۴۰, ۱۴۴, ۱۴۵, ۱۴۸, ۱۴۹, ۱۵۱, ۱۵۶, ۱۵۷, ۱۵۸, ۱۵۹, ۱۶۱, ۱۶۲, ۱۷۰, ۱۷۵, ۱۷۶, ۱۷۷, ۱۷۹, ۱۸۰, ۱۸۵, ۱۸۶, ۱۸۷, ۱۸۸, ۱۸۹, ۱۹۰, ۱۹۱, ۱۹۲, ۱۹۳]

Best CCR: ۷۴,۸۰٪

نمودار تغییر CCR نسبت به پیشروی الگوریتم (هر مرحله) به صورت زیر است.



شکل ۷. نمودار CCR الگوریتم حذف پسرونده

همانگونه که انتظار داشتیم این نمودار تا جایی به صورت صعودی بوده و پس از آن حذف ویژگی های بیشتر باعث کاهش CCR شده است.

پیش‌پرسش ۸ – Computational Exercise: Flower Image Classification

هدف این آزمایش، بررسی و مقایسه‌ی عملکرد مدل‌های مختلف یادگیری عمیق در مسئله‌ی طبقه‌بندی تصاویر گل‌ها است. در این تمرین، ابتدا چند مدل به صورت آموزش از صفر پیاده‌سازی و ارزیابی شدند و سپس از یادگیری انتقالی با استفاده از شبکه‌های از پیش آموزش دیده استفاده گردید.

۸-۱. پیش پردازش داده‌ها

تمام تصاویر به اندازه‌ی 224×224 تغییر مقیاس داده شدند. برای سازگاری با مدل‌های از پیش آموزش دیده، از نرمال‌سازی استاندارد ImageNet استفاده شد:

$$\text{Mean} = [0.485, 0.456, 0.406]$$

$$\text{Std} = [0.229, 0.224, 0.225]$$

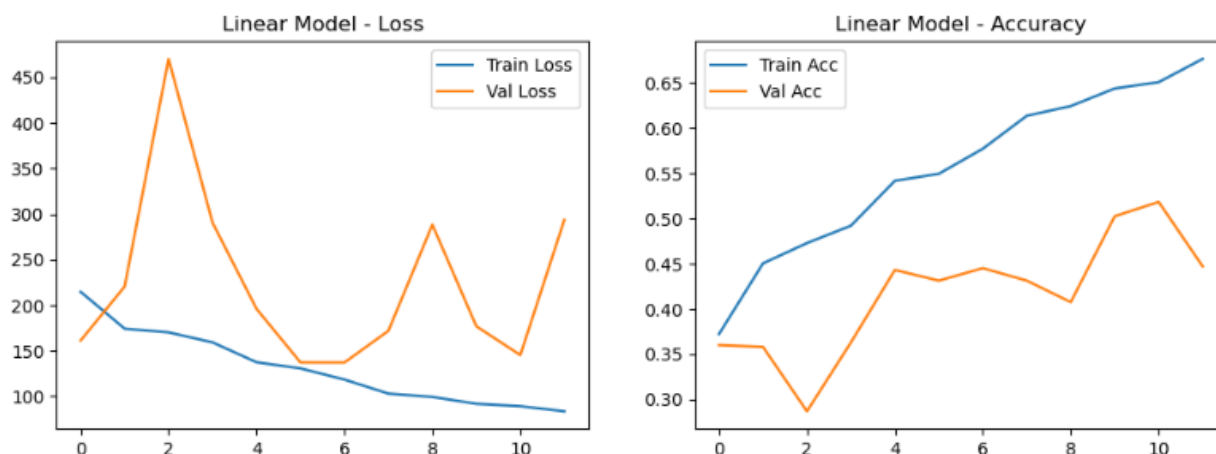
این پیش‌پردازش باعث پایداری بیشتر آموزش و بهبود همگرایی مدل‌ها شد.

۸-۲. بخش اول: آموزش از صفر

در این بخش، سه مدل مختلف که همگی بدون استفاده از وزن‌های از پیش آموزش دیده و به صورت کامل از ابتدا آموزش داده شدند، مورد بررسی قرار گرفتند. هدف از این بخش، تحلیل تاثیر افزایش ظرفیت مدل و استفاده از غیرخطی سازی و تکنیک‌های منظم سازی بر عملکرد نهایی شبکه بود. تمامی مدل‌ها با استفاده از حلقه‌ی آموزش دستی و مکانیزم Early Stopping آموزش داده شدند تا از بیش برازش جلوگیری شود.

۸-۲-۱. مدل خطی

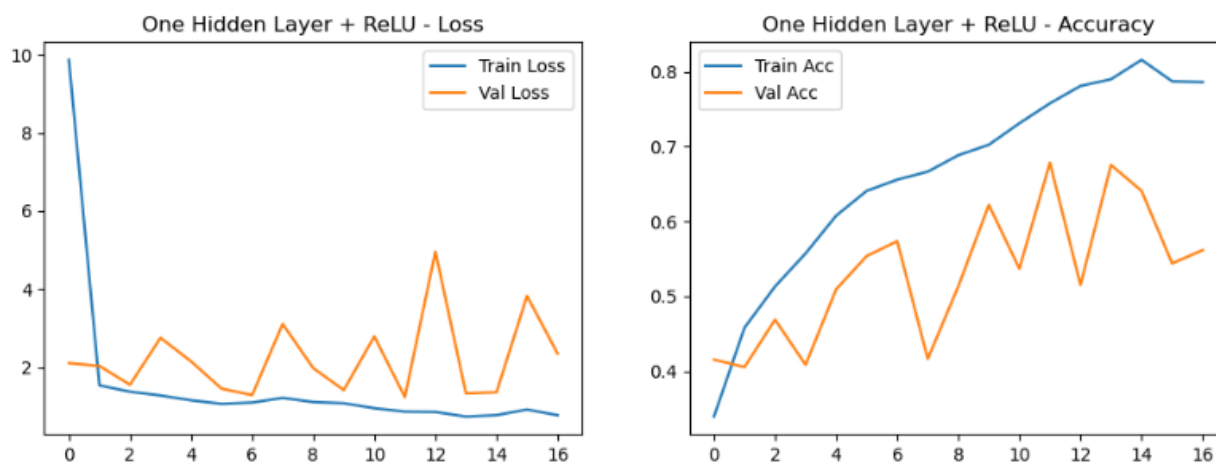
اولین مدل مورد بررسی یک مدل خطی ساده بدون استفاده از هرگونه تابع فعال ساز غیرخطی بود. این مدل تنها قادر به یادگیری مرزهای تصمیم خطی است و در نتیجه توانایی محدودی در استخراج الگوهای پیچیده‌ی موجود در تصاویر دارد. نتایج نشان داد که این مدل به دقت اعتبارسنجی حدود ۵۲ درصد دست یافته است که بیانگر underfitting و ناکافی بودن ظرفیت مدل برای این مسئله است.



شکل ۸. نمودار دقت و **loss** مدل خطی

۸-۲-۲. شبکه با یک لایه پنهان و ReLU

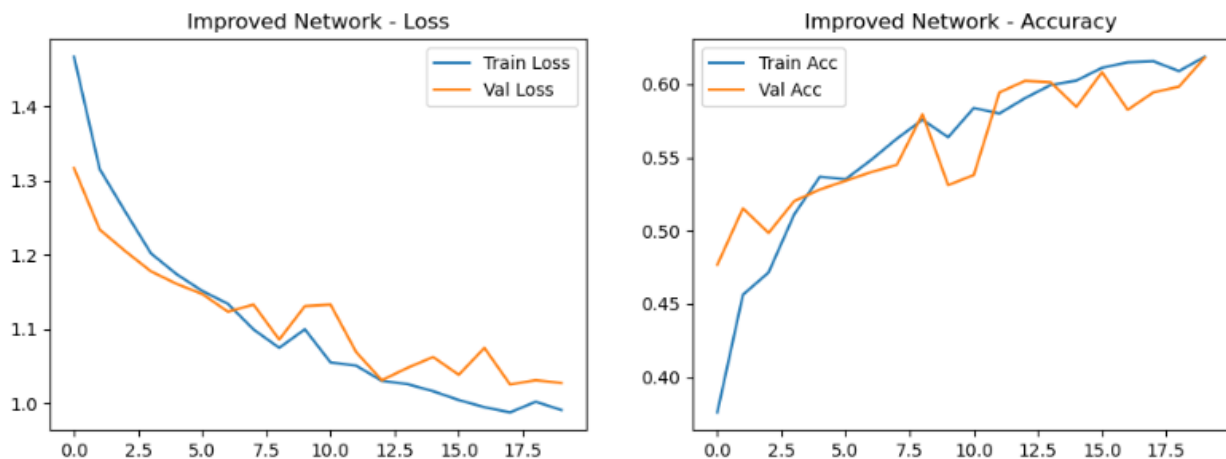
در مدل دوم، با اضافه کردن یک لایه‌ی پنهان و استفاده از تابع فعال ساز ReLU، امکان یادگیری روابط غیرخطی فراهم شد. این تغییر منجر به بهبود قابل توجه عملکرد مدل گردید و دقت اعتبارسنجی به حدود ۶۸ درصد افزایش یافت. این مدل بهترین عملکرد را در میان مدل‌های آموزش از صفر ارائه داد و به وضوح نشان داد که استفاده از غیرخطی‌سازی نقش کلیدی در افزایش قدرت بیان شبکه‌های عصبی دارد.



شکل ۹. نمودار دقت و **loss** شبکه با یک لایه پنهان و ReLU

۸-۲-۳. شبکه‌ی بهبود یافته

مدل سوم به عنوان شبکه‌ی بهبود یافته شامل تکنیک‌هایی نظیر Batch Normalization، Dropout و regularization بود. این شبکه رفتار پایدارتری در طول آموزش از خود نشان داد و اختلاف بین دقت آموزش و اعتبارسنجی در آن کاهش یافت. با این حال، دقت اعتبارسنجی این مدل در حدود ۶۲ درصد باقی ماند که کمی کمتر از مدل قبلی است. این نتیجه نشان می‌دهد که اعمال regularization قوی در دیتاست‌های نسبتاً کوچک می‌تواند منجر به محافظه کاری بیش از حد مدل و کاهش عملکرد نهایی شود.



شکل ۱۰. نمودار دقت و loss شبکه بهبود یافته

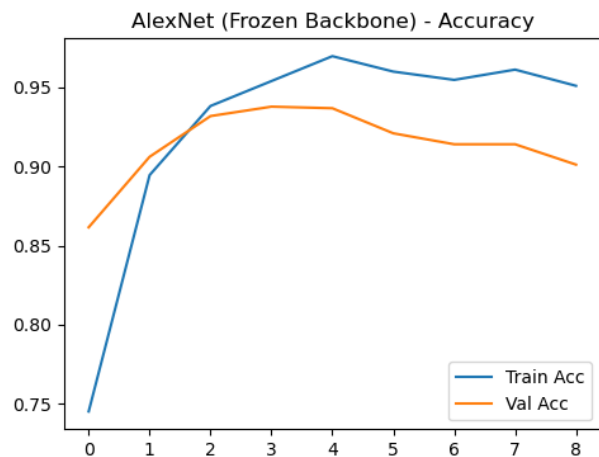
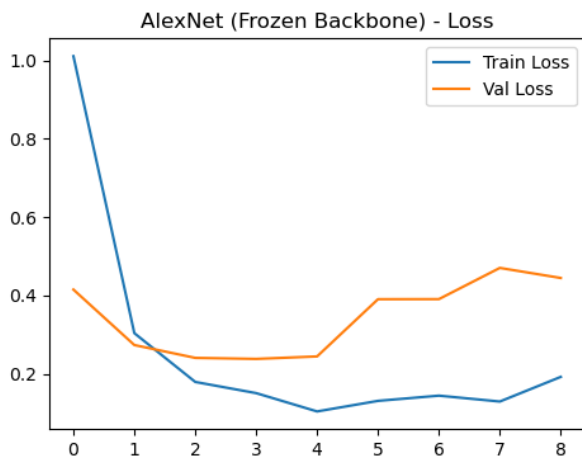
۸-۳. بخش دوم: یادگیری انتقالی

در این بخش از پروژه، از یادگیری انتقالی با استفاده از شبکه‌های AlexNet و VGG^{۱۶} که قبلاً روی دیتاست ImageNet آموزش داده شده‌اند استفاده شد. در هر دو مدل، لایه‌های استخراج ویژگی فریز شدند و تنها لایه‌های نهایی طبقه‌بندی متناسب با تعداد کلاس‌های مسئله آموزش داده شدند. انتظار می‌رود که این رویکرد به دلیل بهره‌گیری از ویژگی‌های سطح بالایی از پیش آموخته شده، عملکرد بهتری نسبت به آموزش کامل از صفر ارائه دهد.

۸-۳-۱. مدل AlexNet

مدل AlexNet با یک classifier جدید روی ویژگی‌های استخراج شده آموزش داده شد. نتایج نشان داد که این مدل از همان epoch‌های ابتدایی به دقت بالایی دست پیدا می‌کند. دقت اعتبارسنجی در epoch اول حدود ۸۶٪ بود و با ادامه‌ی آموزش افزایش یافت. بهترین دقت اعتبارسنجی این مدل در حدود ۹۳/۸٪ به دست آمد که نشان‌دهنده‌ی عملکرد بسیار قوی AlexNet در این مسئله است.

پس از این نقطه، با افزایش تعداد epoch‌ها، دقت اعتبارسنجی به تدریج کاهش یافت، در حالی که دقت آموزش همچنان بالا باقی ماند. این رفتار نشان‌دهنده‌ی شروع بیش‌برازش است و تایید می‌کند که استفاده از early stopping برای انتخاب بهترین مدل ضروری بوده است.

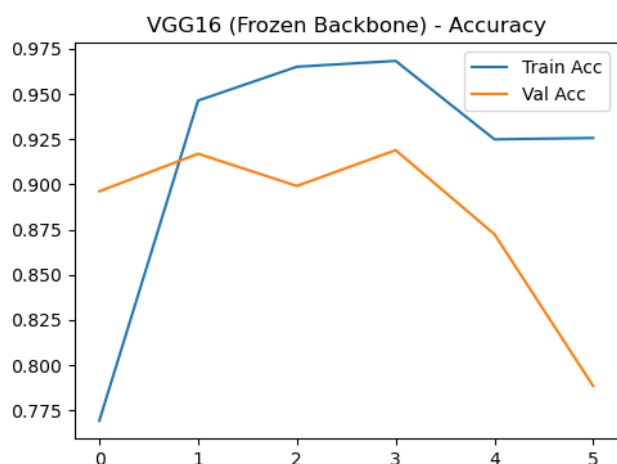
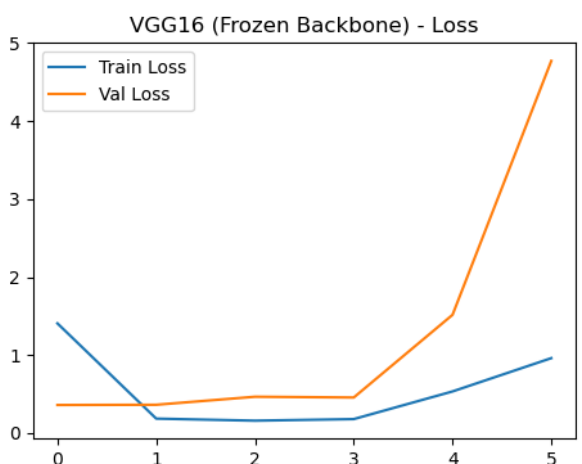


شکل ۱۱. نمودار دقت و loss مدل AlexNet

۸-۳-۲. مدل VGG۱۶

در ادامه، مدل VGG۱۶ با رویکرد مشابه مورد استفاده قرار گرفت. این شبکه به دلیل عمق بیشتر، ظرفیت بالاتری برای استخراج ویژگی‌های پیچیده دارد. نتایج نشان داد که VGG۱۶ نیز در epochهای ابتدایی به دقت بالایی دست می‌یابد و بهترین دقت اعتبارسنجی آن حدود ۹۱/۹٪ ثبت شد.

با این حال، با ادامه‌ی آموزش، کاهش محسوسی در دقت اعتبارسنجی مشاهده شد، به‌طوری که در epochهای پایانی افت شدیدی در عملکرد رخ داد. این موضوع نشان می‌دهد که VGG۱۶ نسبت به AlexNet سریع‌تر دچار overfitting شده است، که می‌تواند به دلیل تعداد پارامترهای بسیار بیشتر این شبکه و محدود بودن حجم دیتاست باشد.



شکل ۱۲. نمودار دقت و loss مدل VGG۱۶

۸-۴. مقایسه و تحلیل نهایی

مقایسه‌ی نتایج به دست آمده نشان می‌دهد که استفاده از یادگیری انتقالی منجر به بهبود بسیار چشمگیر عملکرد مدل‌ها نسبت به آموزش از صفر شده است. در حالی که بهترین مدل در Part I دقتی در حدود ۶۸٪ داشت، مدل AlexNet توانست دقت اعتبارسنجی نزدیک به ۹۴٪ را ارائه دهد.

در مقایسه‌ی دو مدل انتقالی، AlexNet عملکرد پایدارتر و بهتری نسبت به VGG^{۱۶} نشان داد. اگرچه VGG^{۱۶} شبکه‌ای عمیق‌تر است، اما در این مسئله خاص و با حجم داده‌ی محدود، دچار بیش برآزش سریع‌تری شد. این نتایج نشان می‌دهد که انتخاب مدل مناسب در یادگیری انتقالی تنها به عمق شبکه وابسته نیست و باید با توجه به اندازه و ماهیت دیتاست انجام شود.

جدول ۱. جدول مقایسه مدل‌ها

مدل	روش آموزش	بهترین دقت اعتبار سنجی (درصد)
خطی	آموزش از ابتدا	۵۲
یک لایه پنهان با ReLU	آموزش از ابتدا	۶۷.۹
شبکه بهبود یافته	آموزش از ابتدا	۶۲
AlexNet	یادگیری انتقالی	۹۳.۸
VGG ^{۱۶}	یادگیری انتقالی	۹۱.۹

۸-۵. نتیجه گیری

در این تمرین نشان داده شد که مدل‌های ساده آموزش دیده از صفر، اگر چه قادر به یادگیری الگوهای اولیه هستند، اما عملکرد آن‌ها در مقایسه با مدل‌های مبتنی بر یادگیری انتقالی به مراتب ضعیف‌تر است. افزودن غیرخطی سازی باعث بهبود قابل توجه عملکرد شد، اما بیشترین جهش دقت با استفاده از شبکه‌های از پیش آموزش دیده به دست آمد.

به طور کلی، نتایج این آزمایش تایید می‌کند که یادگیری انتقالی یک راهکار بسیار موثر برای مسائل یادگیری ماشین، به ویژه در شرایطی با داده‌ی محدود، محسوب می‌شود و می‌تواند بدون نیاز به آموزش کامل شبکه از ابتدا، دقت بسیار بالاتری ارائه دهد.

- [١] Available: [ادرون خطی]. "broadcasting",
<https://numpy.org/doc/stable/user/basics.broadcasting.html>
- [٢] Available: [ادرون خطی]. "Kronecker delta," wikipedia"
https://en.wikipedia.org/wiki/Kronecker_delta
- [٣] Available: [https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html) [ادرون خطی]. "GaussianNB," scikit"
[learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html)
- [٤] Available: [https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html) [ادرون خطی]. "accuracy_score," scikit"
[learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html)