

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان
دانشکده علوم ریاضی

نگاهی به یادگیری عمیق از دیدگاه ریاضیات کاربردی

پروژه کارشناسی

پریمه الصفی

شماره دانشجویی: ۹۹۱۸۷۴۳

استاد راهنما

دکتر رضا مختاری

۱۴۰۳

فهرست مطالب

۱	چکیده
فصل اول: مقدمه‌ای بر یادگیری عمیق و شبکه‌های عصبی مصنوعی	
۱-۱	معماری و اصول پایه‌ای شبکه‌های عصبی مصنوعی
۱-۱-۱	روش‌های پیشرفته دسته‌بندی تصاویر با استفاده از شبکه‌های عصبی
۲-۱	بررسی جامع توابع فعال‌سازی در یادگیری عمیق
۳-۱	تحلیل ساختار شبکه‌های عصبی و کاربردهای آن
۴-۱	تابع هزینه
فصل دوم: روش‌های بهینه‌سازی و بهبود عملکرد شبکه‌های عصبی	
۱-۲	بهبود عملکرد شبکه‌های عصبی چند لایه با تکنیک‌های بهینه‌سازی
۱-۱-۲	نمادگذاری
۲-۱-۲	استراتژی‌های بهینه‌سازی و کاهش خطا در شبکه‌های عصبی
۳-۱-۲	تکنیک‌های بهینه‌سازی با استفاده از روش گرادینان تصادفی
۲-۲	روش پس‌انتشار
فصل سوم: کاربردهای عملی شبکه‌های عصبی و ابزارهای پیاده‌سازی	
۱-۳	شبکه‌های عصبی پیچشی: اصول، ساختار و کاربردها
۲-۳	پیشگیری از بیش‌برازش در پردازش تصاویر
۱-۲-۳	بیش‌برازش
۳-۳	تحلیل عمیق توابع هزینه و فعال‌سازی در شبکه‌های عصبی
۴-۳	بهینه‌سازی دسته‌بندی تصاویر
۱-۴-۳	پیاده‌سازی و ارزیابی مدل‌های دسته‌بندی تصویر
۲-۴-۳	تحلیل نتایج و تاثیر حذف تصادفی در بهبود دقت مدل
۵-۳	شبیه‌سازی و ارزیابی مدل‌های شبکه عصبی

۳۶	۳-۶ جمع بندی
۳۹	مراجع

فهرست تصاویر

۳	۱-۱ ماتریس RGB
۴	۲-۱ تابع سیگموئید
۵	۳-۱ معماری شبکه-ساختار نورون
۵	۴-۱ نقاط در فضای \mathbb{R}^2 . دایره‌ها نشان‌دهنده نقاط دسته A و ضربدرها نشان‌دهنده نقاط دسته B هستند.
۶	۵-۱ نمایی از یک شبکه عصبی با چهار لایه
۸	۶-۱ خروجی شبکه عصبی بر روی داده‌های شکل ۱-۴
۱۰	۱-۲ نمایی از یک شبکه عصبی با پنج لایه
۲۰	۲-۲ تابع هزینه در مقابل تعداد تکرار
۲۹	۱-۳ نمای کلی معماری شبکه
۳۱	۲-۳ آموزش شبکه بدون حذف تصادفی
۳۲	۳-۳ آموزش شبکه با حذف تصادفی
۳۳	۴-۳ ماتریس در هم‌ریختگی بدون حذف تصادفی
۳۳	۵-۳ ماتریس در هم‌ریختگی با حذف تصادفی
۳۴	۶-۳ دقت آموزش شبکه شبیه‌سازی شده بدون حذف تصادفی
۳۴	۷-۳ خطای شبکه شبیه‌سازی شده با حذف تصادفی
۳۵	۸-۳ دقت آموزش شبکه شبیه‌سازی شده با حذف تصادفی
۳۵	۹-۳ خطای شبکه شبیه‌سازی شده با حذف تصادفی
۳۵	۱۰-۳ ماتریس در هم‌ریختگی مدل شبیه‌سازی شده بدون حذف تصادفی
۳۶	۱۱-۳ ماتریس در هم‌ریختگی مدل شبیه‌سازی شده با حذف تصادفی

چکیده

این گزارش به بررسی جامع روش‌ها و تکنیک‌های مورد استفاده در شبکه‌های عصبی مصنوعی و یادگیری عمیق می‌پردازد. در بخش ابتدایی، به ساختار شبکه‌های عصبی مصنوعی (ANN) و مفاهیم پایه‌ای همچون توابع فعال‌سازی، شبکه‌های عصبی چندلایه و شبکه‌های عصبی پیچشی (CNN) پرداخته شده است. همچنین، روش‌های بهینه‌سازی شامل پس‌انتشار خطا، گرادیان نزولی و روش‌های بهینه‌سازی تصادفی مانند «گرادیان نزول تصادفی» (SGD) و «حذف تصادفی» (Dropout) به منظور بهبود عملکرد و کاهش بیش‌برازش معرفی می‌شوند. در ادامه، کاربردهای عملی این تکنیک‌ها در زمینه‌های مختلف از جمله پردازش تصویر، دسته‌بندی داده‌های پیچیده، و یادگیری ویژگی‌های غیرخطی در داده‌های حجیم بررسی می‌شود. ابزارهای نرم‌افزاری نظیر MatConvNet که برای پیاده‌سازی شبکه‌های عمیق و تسریع فرایند یادگیری مورد استفاده قرار می‌گیرند، معرفی شده‌اند. گزارش همچنین به تحلیل دقیق روش‌های بهینه‌سازی پارامترها با استفاده از توابع هزینه مختلف مانند «تابع هزینه نرم بیشینه» (Softmax Loss) و «تابع هزینه درجه دوم» می‌پردازد. نتایج به‌دست‌آمده نشان می‌دهند که استفاده از شبکه‌های عصبی پیچشی برای مسائل پردازش تصویر، همراه با تکنیک‌های بهینه‌سازی مناسب، می‌تواند دقت و کارایی مدل‌ها را به‌طور قابل توجهی بهبود بخشد.

فصل ۱

مقدمه‌ای بر یادگیری عمیق و شبکه‌های عصبی مصنوعی

یادگیری عمیق^۱ شاخه‌ای از یادگیری ماشین^۲ است که با بهره‌گیری از شبکه‌های عصبی عمیق^۳ و تعداد زیادی لایه‌های مخفی، قادر به استخراج ویژگی‌ها و الگوهای پیچیده از داده‌های حجیم و غیرساخت‌یافته هست. این تکنیک به شبکه‌های عصبی امکان می‌دهد تا با تحلیل لایه به لایه داده‌ها، روابط پنهان و غیرمشخص را کشف کنند.

۱-۱ معماری و اصول پایه‌ای شبکه‌های عصبی مصنوعی

شبکه‌های عصبی مصنوعی^۴ ساختارهای محاسباتی الهام گرفته از شبکه عصبی مغز انسان هستند. این شبکه‌ها شامل مجموعه‌ای از نودها یا نورون‌ها هستند که به صورت لایه‌ای سازماندهی شده‌اند. هر نورون ورودی‌هایی را دریافت و با وزن‌های مشخص ترکیب می‌کند، سپس یک تابع فعال‌سازی را اعمال کرده و خروجی را به نورون‌های بعدی ارسال می‌کند. [۱]

۱-۱-۱ روش‌های پیشرفته دسته‌بندی تصاویر با استفاده از شبکه‌های عصبی

فرض کنید می‌خواهیم یک شبکه عصبی طراحی کنیم که تصاویر گربه را از تصاویر غیر گربه تشخیص دهد، به طوری که اگر تصویر ورودی یک گربه باشد، خروجی ۱ و در غیر این صورت ۰ باشد. تصاویر معمولاً به صورت

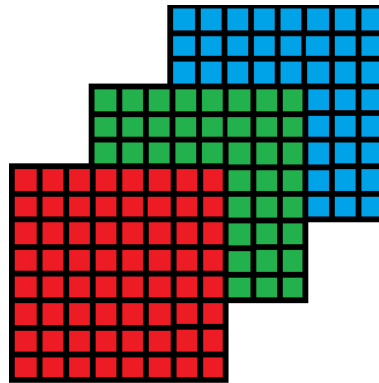
¹Deep Learning

²Machine Learning

³Neural Networks

⁴Artificial Neural Networks

RGB^۱ هستند ۱-۱ و به همین دلیل به شکل بردار یا ماتریس به شبکه وارد می‌شوند. [۲]



شکل ۱-۱: ماتریس RGB

برای آموزش چنین شبکه‌ای، مهم‌ترین مسئله، دقت شبکه در فاز آموزشی و قابلیت تعمیم آن به داده‌های جدید است. در این فرایند، ورودی‌ها که به صورت ماتریس و بردار هستند، باید به اعدادی باینری (۰، ۱) تبدیل شوند. برای این کار از توابع فعال‌سازی مختلفی استفاده می‌شود.

۲-۱ بررسی جامع توابع فعال‌سازی در یادگیری عمیق

یکی از مهم‌ترین توابع فعال‌سازی^۲، تابع سیگموئید^۳ است که به صورت

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (1-1)$$

این تابع را می‌توان به عنوان نسخه‌ای از یک تابع پله‌ای^۴ در نظر گرفت. به این معنا که رفتار یک نورون در مغز را تقلید می‌کند؛ به این صورت که وقتی ورودی به اندازه کافی بزرگ باشد، فعال می‌شود (خروجی ۱) و در غیر این صورت غیرفعال می‌ماند (خروجی ۰).

تابع سیگموئید خاصیت مفیدی دارد که مشتق آن به شکل ساده‌ای به صورت

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)). \quad (2-1)$$

برای ساده‌سازی بیشتر، تابع سیگموئید را به صورت برداری تعریف می‌کنیم. اگر $\mathbf{z} \in \mathbb{R}^m$ باشد، و $\sigma : \mathbb{R}^m \rightarrow \mathbb{R}^m$ تابع سیگموئید با اعمال به اجزای هر مؤلفه تعریف می‌شود، به طوری که

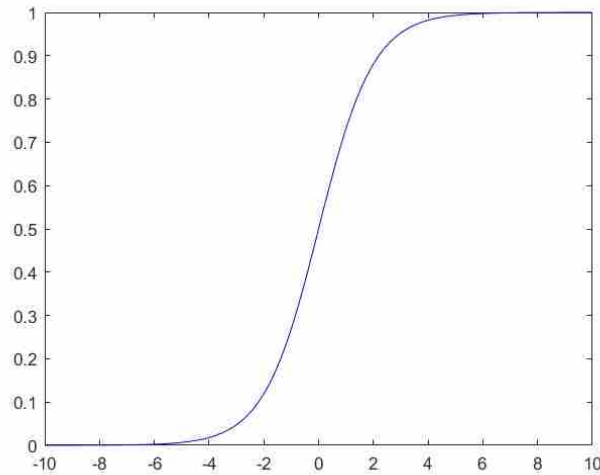
¹Red, Green and Blue

²Activation Functions

³Sigmoid Function

⁴Step Function

$$(\sigma(\mathbf{z}))_i = \sigma(z_i). \quad (3-1)$$



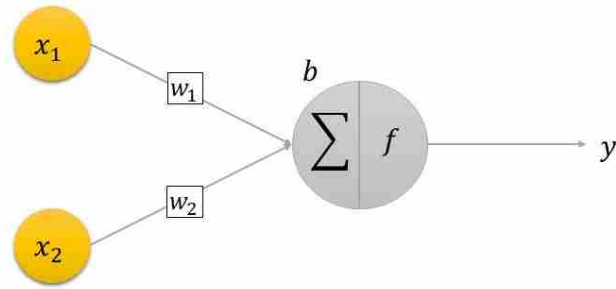
شکل ۱-۲: تابع سیگموئید

۳-۱ تحلیل ساختار شبکه‌های عصبی و کاربردهای آن

با نمادگذاری گفته شده در قسمت قبل، می‌توانیم لایه‌هایی از نورون‌ها را ایجاد کنیم. در هر لایه، هر نورون یک عدد حقیقی تولید می‌کند که به نورون‌های لایه بعدی ارسال می‌شود. در لایه بعدی، هر نورون ترکیب وزنی مختص به خود از این مقادیر را تشکیل داده، بایاس خود را اضافه کرده و تابع سیگموئید را اعمال می‌کند. اگر اعداد حقیقی تولید شده توسط نورون‌ها در یک لایه در یک بردار a جمع‌آوری شوند، بردار خروجی‌ها از لایه بعدی به صورت (مانند شکل ۱-۳)

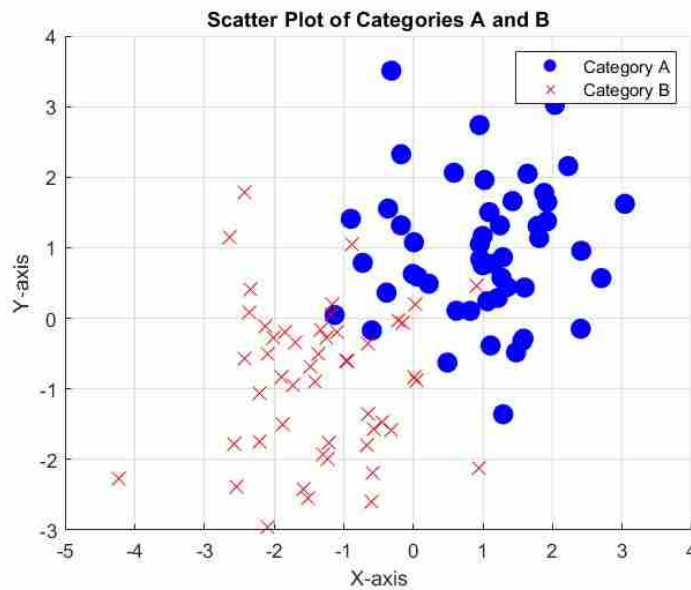
$$\sigma(Wa + b). \quad (4-1)$$

در اینجا، W یک ماتریس و b یک بردار است. W حاوی وزن‌ها و b حاوی بایاس‌ها است. تعداد ستون‌های W با تعداد نورون‌هایی که بردار a را در لایه قبلی تولید کرده‌اند، مطابقت دارد. تعداد ردیف‌های W و مؤلفه‌های b با تعداد نورون‌های لایه فعلی مطابقت دارد. فرض کنید می‌خواهیم دسته‌بندی روی یک مجموعه داده انجام دهیم. به عنوان مثال، مجموعه داده درباره دو منطقه حفاری است که به دودسته مناطق A و



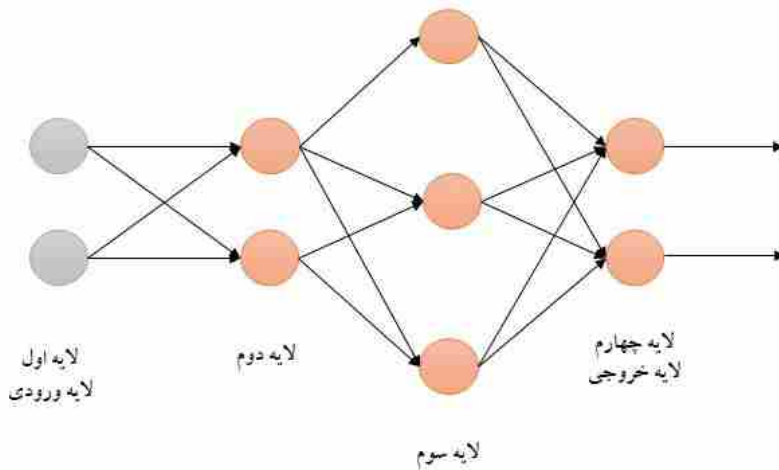
شکل ۱-۳: معماری شبکه- ساختار نرون

مناطق B دسته‌بندی شده‌اند. پس از برازش، خروجی به صورت نمودار نشان داده می‌شود، به طوری که مناطق A با دایره و مناطق B با ضربدر مشخص شده‌اند.



شکل ۱-۴: نقاط در فضای \mathbb{R}^2 . دایره‌ها نشان‌دهنده نقاط دسته A و ضربدرها نشان‌دهنده نقاط دسته B هستند.

شکل ۱-۵ یک شبکه عصبی مصنوعی با چهار لایه را نشان می‌دهد. این شکل از شبکه را برای مسئله‌ای که در شکل ۱-۴ تعریف شده است، به کار خواهیم گرفت. در شبکه نشان‌داده‌شده در شکل ۱-۵، لایه اول (ورودی) با دو دایره نشان‌داده‌شده است، زیرا نقاط داده ورودی ما دارای دو مؤلفه هستند. لایه دوم دارای دو دایره توپر است که نشان می‌دهد دو نرون در حال استفاده هستند. پیکان‌های بین لایه ۱ و لایه ۲ نشان می‌دهند که هر دو مؤلفه داده‌های ورودی به دو نرون در لایه ۲ منتقل می‌شوند. از آنجاکه داده‌های ورودی به صورت $x \in \mathbb{R}^2$ است، وزن‌ها و بایاس‌های لایه ۲ به ترتیب با ماتریس $\mathbf{W}^{[2]} \in \mathbb{R}^{2 \times 2}$ و بردار $\mathbf{b}^{[2]} \in \mathbb{R}^2$ نشان داده



شکل ۱-۵: نمایی از یک شبکه عصبی با چهار لایه

می‌شوند. خروجی از لایه ۲ به صورت

$$\sigma(\mathbf{W}^{[2]}\mathbf{x} + \mathbf{b}^{[2]}) \in \mathbb{R}^2, \quad (5-1)$$

لایه ۳ دارای سه نورون است که هر کدام ورودی در \mathbb{R}^2 دریافت می‌کنند، بنابراین وزن‌ها و بایاس‌های لایه ۳ به ترتیب با ماتریس $\mathbf{W}^{[3]} \in \mathbb{R}^{3 \times 2}$ و بردار $\mathbf{b}^{[3]} \in \mathbb{R}^3$ نشان داده می‌شوند. خروجی از لایه ۳ به صورت

$$\sigma(\mathbf{W}^{[3]}\sigma(\mathbf{W}^{[2]}\mathbf{x} + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) \in \mathbb{R}^3, \quad (6-1)$$

لایه چهارم (خروجی) دارای دو نورون است که هر کدام ورودی در \mathbb{R}^3 دریافت می‌کنند، بنابراین وزن‌ها و بایاس‌های این لایه به ترتیب با ماتریس $\mathbf{W}^{[4]} \in \mathbb{R}^{2 \times 3}$ و بردار $\mathbf{b}^{[4]} \in \mathbb{R}^2$ نشان داده می‌شوند. خروجی از لایه ۴، و در نتیجه از کل شبکه، به صورت

$$F(\mathbf{x}) = \sigma(\mathbf{W}^{[4]}\sigma(\mathbf{W}^{[3]}\sigma(\mathbf{W}^{[2]}\mathbf{x} + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) + \mathbf{b}^{[4]}) \in \mathbb{R}^2. \quad (7-1)$$

عبارت ۷-۱ تابع $F: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ را بر اساس ۲۳ پارامترش تعریف می‌کند. ما نیاز داریم که برای نقاط داده در دسته A به $[1, 0]^T$ و برای نقاط داده در دسته B به $[0, 1]^T$ نزدیک باشد. سپس، با داشتن یک نقطه جدید $x \in \mathbb{R}^2$ ، منطقی است که آن را بر اساس بزرگ‌ترین مؤلفه $F(x)$ دسته‌بندی کنیم؛ به این صورت که اگر $F_1(x) > F_2(x)$ باشد، دسته A و اگر $F_1(x) < F_2(x)$ باشد، دسته B، دارد.

۴-۱ تابع هزینه

نیازمندی بیان شده در F می‌تواند از طریق یک تابع هزینه مشخص شود. با نشان دادن ده نقطه داده در شکل ۴-۱ با $\{x^{(i)}\}_{i=1}^{10}$ ، از $y(x^{(i)})$ برای خروجی هدف استفاده می‌کنیم؛ به این صورت که:

$$y(x^{(i)}) = \begin{cases} [1, 0]^T & \text{اگر } x^{(i)} \text{ متعلق به دسته A باشد} \\ [0, 1]^T & \text{اگر } x^{(i)} \text{ متعلق به دسته B باشد} \end{cases} \quad (۸-۱)$$

بنابراین تابع هزینه به صورت

$$\text{cost}(W^{[2]}, W^{[3]}, W^{[4]}, b^{[2]}, b^{[3]}, b^{[4]}) = \frac{1}{10} \sum_{i=1}^{10} \frac{1}{2} \|y(x^{(i)}) - F(x^{(i)})\|_2^2. \quad (۹-۱)$$

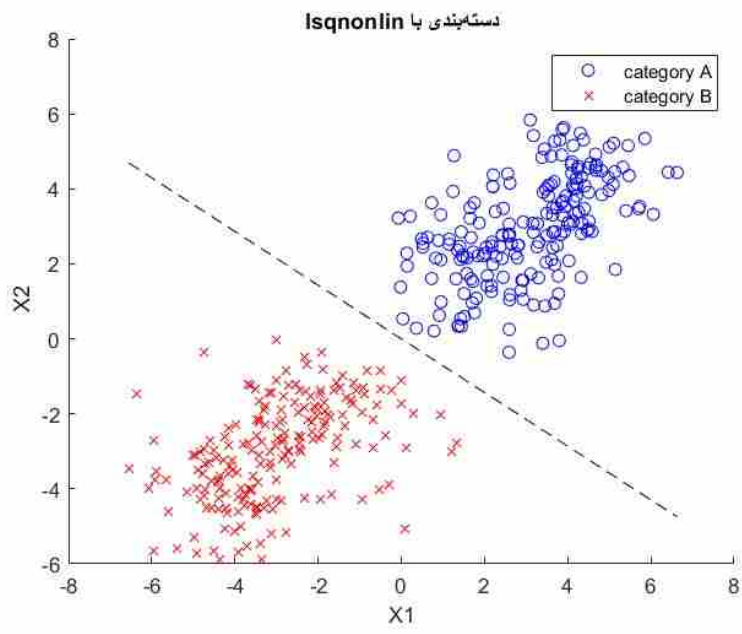
در اینجا، ضریب $\frac{1}{2}$ برای راحتی لحاظ شده است؛ زمانی که شروع به مشتق‌گیری می‌کنیم، این امر مسائل را ساده‌تر می‌کند. تأکید می‌کنیم که تابع هزینه، تابعی از وزن‌ها و بایاس‌ها است.

آموزش شبکه

انتخاب وزن‌ها و بایاس‌ها به گونه‌ای که تابع هزینه را به حداقل برساند، به عنوان آموزش شبکه شناخته می‌شود. توجه داریم که در اصل، تغییر مقیاس یک تابع هدف، مسئله بهینه‌سازی را تغییر نمی‌دهد، به این معنی که اگر تابع هزینه را به عنوان مثال به $100 \cdot \text{cost}$ یا $\text{cost}/30$ تغییر دهیم، باید به همان نقطه کمینه برسیم، بنابراین ضرایب $\frac{1}{2}$ و $\frac{1}{10}$ در ۹-۱ نباید بر نتیجه تأثیری داشته باشند.

بهینه‌سازی پارامترها

برای تحلیل داده‌های موجود در شکل ۴-۱، ما از جعبه‌ابزار بهینه‌سازی MATLAB برای بهینه‌سازی تابع هزینه ۹-۱ استفاده کردیم که این تابع به ۲۳ پارامتر مختلف شامل ماتریس‌های وزن $W^{[2]}, W^{[3]}, W^{[4]}$ و بردارهای بایاس $b^{[2]}, b^{[3]}, b^{[4]}$ وابسته است. به طور خاص، از الگوریتم «lsqnonlin» که یک روش برای حل مسائل بهینه‌سازی با کمینه کردن خطاهای مربعی به طور غیرخطی است، بهره بردیم. در نتیجه، شبکه آموزش دیده، که در شکل ۶-۱ نشان داده شده است، مرز تفکیکی را ایجاد می‌کند به طوری که $F_1(x) > F_2(x)$ برای دسته A و $F_1(x) < F_2(x)$ برای دسته B.



شکل ۱-۶: خروجی شبکه عصبی بر روی داده های شکل ۱-۴

فصل ۲

روش‌های بهینه‌سازی و بهبود عملکرد شبکه‌های عصبی

۱-۲ بهبود عملکرد شبکه‌های عصبی چند لایه با تکنیک‌های بهینه‌سازی

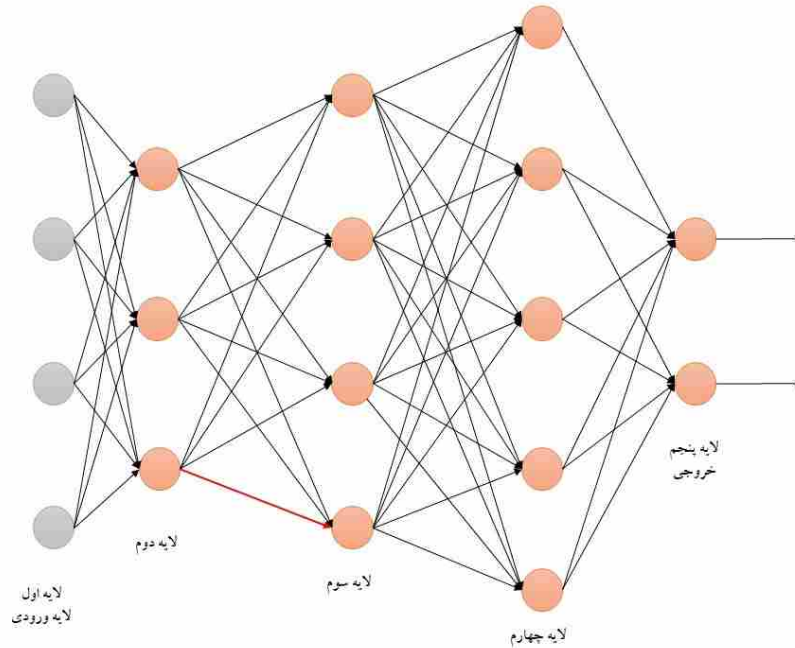
مثال چهار لایه‌ای در فصل ۱ مفهوم نورون‌ها را معرفی کرد، که با استفاده از تابع سیگموئید نشان داده می‌شوند و در لایه‌های مختلف عمل می‌کنند. در یک لایه عمومی، هر نورون ورودی یکسانی دریافت می‌کند و یک مقدار حقیقی تولید می‌کند که به هر نورون در لایه بعدی ارسال می‌شود. دو لایه استثنایی وجود دارد. در لایه ورودی، «لایه قبلی» وجود ندارد و هر نورون بردار ورودی را دریافت می‌کند. در لایه خروجی، «لایه بعدی» وجود ندارد و این نورون‌ها خروجی کلی را فراهم می‌کنند. لایه‌های بین این دو لایه، «لایه‌های مخفی» نامیده می‌شوند که محاسبات میانی را انجام می‌دهند. یادگیری عمیق به معنای استفاده از تعداد زیادی لایه‌های مخفی است.

۱-۱-۲ نمادگذاری

در این بخش، فرم کلی نمادگذاری را به تفصیل بیان می‌کنیم. فرض می‌کنیم شبکه دارای L لایه است، به طوری که لایه‌های ۱ و L به ترتیب لایه‌های ورودی و خروجی هستند. فرض کنید لایه l ، برای $l = 1, 2, 3, \dots, L$ ، شامل n_l نورون است. بنابراین n_1 بُعد داده‌های ورودی است و به طور کلی، شبکه از R^{n_1} به R^{n_L} نگاشت می‌کند. ما از $W^{[l]} \in \mathbb{R}^{n_l \times n_{l-1}}$ برای نشان دادن ماتریس وزن‌ها در لایه l استفاده می‌کنیم. به طور دقیق‌تر، $w_{jk}^{[l]}$ وزنی است که نورون j در لایه l به خروجی نورون k در لایه $l-1$ اعمال می‌کند. به همین ترتیب، بردار $b^{[l]} \in \mathbb{R}^{n_l}$ بایاس‌ها برای لایه l است، به طوری که نورون j در لایه l از بایاس $b_j^{[l]}$ استفاده می‌کند.

در شکل ۲-۱، مثالی با $L = 5$ لایه ارائه می‌دهیم. در این مثال، $n_1 = 4, n_2 = 3, n_3 = 4, n_4 = 5, n_5 = 2$

$$W^{[2]} \in \mathbb{R}^{3 \times 4}, \quad W^{[3]} \in \mathbb{R}^{4 \times 3}, \quad W^{[4]} \in \mathbb{R}^{5 \times 4}, \quad W^{[5]} \in \mathbb{R}^{2 \times 5}, \quad b^{[2]} \in \mathbb{R}^3, \quad b^{[3]} \in \mathbb{R}^4, \\ b^{[4]} \in \mathbb{R}^5, \quad b^{[5]} \in \mathbb{R}^2.$$



شکل ۱-۲: نمایی از یک شبکه عصبی با پنج لایه

با داشتن یک ورودی $\mathbf{x} \in \mathbb{R}^{n_1}$ ، می‌توانیم به طور مرتب عملکرد شبکه را با توجه به اینکه $a_j^{[l]}$ خروجی یا فعال‌سازی نورون j در لایه l را نشان می‌دهد، خلاصه کنیم؛ بنابراین

$$a^{[1]} = \mathbf{x} \in \mathbb{R}^{n_1}, \quad (1-2)$$

$$a^{[l]} = \sigma(W^{[l]}a^{[l-1]} + b^{[l]}) \in \mathbb{R}^{n_l} \quad \text{for } l = 2, 3, \dots, L. \quad (2-2)$$

روابط فوق، نحوه عملکرد الگوریتم را در پردازش ورودی از طریق شبکه به تولید یک خروجی $a^{[L]} \in \mathbb{R}^{n_L}$ نشان می‌دهند.

۲-۱-۲ استراتژی‌های بهینه‌سازی و کاهش خطا در شبکه‌های عصبی

فرض کنید که N داده یا نقاط آموزشی در \mathbb{R}^{n_1} داریم، یعنی $\{x^{(i)}\}_{i=1}^N$ که برای آن‌ها خروجی‌های هدف $\{y(x^{(i)})\}_{i=1}^N$ در \mathbb{R}^{n_L} داده شده‌اند. با تعمیم رابطه ۱-۹، تابع هزینه به شکل زیر است:

$$\text{Cost} = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|y(x^{(i)}) - a^{[L]}(x^{(i)})\|_2^2. \quad (3-2)$$

۳-۱-۲ تکنیک‌های بهینه‌سازی با استفاده از روش گرادیان تصادفی

به طور کلی وزن‌ها و بایاس‌ها به صورت ماتریس‌ها و بردارها هستند، اما در این مرحله راحت‌تر است که تصور کنیم آن‌ها به عنوان یک بردار واحد ذخیره می‌شوند که آن را \mathbf{p} می‌نامیم. مثال شکل ۱-۵ شامل مجموعاً ۲۳ وزن و بایاس است. بنابراین، در آن حالت، $\mathbf{p} \in \mathbb{R}^{23}$ است. به طور کلی، فرض خواهیم کرد که $\mathbf{p} \in \mathbb{R}^s$ و تابع هزینه در ۳-۲ را به صورت $\text{Cost}(\mathbf{p})$ می‌نویسیم بنابراین $\text{Cost} : \mathbb{R}^s \rightarrow \mathbb{R}$.

اکنون یک روش کلاسیک در بهینه‌سازی معرفی می‌کنیم که اغلب به عنوان نزول گرادیان^۱ شناخته می‌شود. این روش به صورت تکراری عمل می‌کند و دنباله‌ای از بردارها در \mathbb{R}^s را محاسبه می‌کند، با هدف همگرا شدن به یک بردار که تابع هزینه را به حداقل می‌رساند. فرض کنید که بردار فعلی ما \mathbf{p} است؛ چگونه باید یک تغییر، $\Delta\mathbf{p}$ ، انتخاب کنیم به طوری که بردار بعدی، $\mathbf{p} + \Delta\mathbf{p}$ ، بهبود را نشان دهد؟ اگر $\Delta\mathbf{p}$ کوچک باشد، با نادیده‌گیری جملات مرتبه $\|\Delta\mathbf{p}\|^2$ ، بسط سری تیلور به صورت

$$\text{Cost}(\mathbf{p} + \Delta\mathbf{p}) \approx \text{Cost}(\mathbf{p}) + \sum_{r=1}^s \frac{\partial \text{Cost}(\mathbf{p})}{\partial p_r} \Delta p_r. \quad (4-2)$$

در اینجا $\frac{\partial \text{Cost}(\mathbf{p})}{\partial p_r}$ مشتق جزئی تابع هزینه نسبت به r امین پارامتر را نشان می‌دهد. برای راحتی، اجازه می‌دهیم که $\nabla \text{Cost}(\mathbf{p}) \in \mathbb{R}^s$ بردار مشتقات جزئی را نشان دهد که به عنوان گرادیان شناخته می‌شود، به طوری که

$$(\nabla \text{Cost}(\mathbf{p}))_r = \frac{\partial \text{Cost}(\mathbf{p})}{\partial p_r},$$

سپس معادله ۴-۲ به شکل زیر تبدیل می‌شود

$$\text{Cost}(\mathbf{p} + \Delta\mathbf{p}) \approx \text{Cost}(\mathbf{p}) + \nabla \text{Cost}(\mathbf{p})^T \Delta\mathbf{p}. \quad (5-2)$$

هدف ما کاهش مقدار تابع هزینه است. رابطه ۲-۵ ایده انتخاب $\Delta\mathbf{p}$ را برای منفی کردن $\nabla \text{Cost}(\mathbf{p})^T \Delta\mathbf{p}$ تا حد ممکن انگیزه می‌دهد. ما می‌توانیم این مسئله را با استفاده از نابرابری کوشی - شوارتز حل کنیم که

¹Gradient descent

بیان می‌کند برای هر $f, g \in \mathbb{R}^s$ ، داریم $\|f^T g\| \leq \|f\|_2 \|g\|_2$ ؛ بنابراین منفی‌ترین مقداری که $f^T g$ می‌تواند باشد، $-\|f\|_2 \|g\|_2$ است که زمانی رخ می‌دهد که $f = -g$ ؛ بنابراین، بر اساس ۲-۵، باید Δp را در جهت $-\nabla \text{Cost}(p)$ انتخاب کنیم. با در نظر داشتن اینکه ۲-۵ یک تقریب است که تنها برای Δp کوچک مرتبط است، ما خود را به یک قدم کوچک در آن جهت محدود خواهیم کرد. این منجر به به‌روزرسانی زیر می‌شود

$$p \rightarrow p - \eta \nabla \text{Cost}(p). \quad (6-2)$$

اینجا η یک اندازه گام کوچک است که در این زمینه به‌عنوان نرخ یادگیری^۱ شناخته می‌شود. این معادله روش نزول تندترین^۲ را تعریف می‌کند. ما یک بردار اولیه انتخاب می‌کنیم و با استفاده از ۲-۶ تکرار می‌کنیم تا زمانی که یک معیار توقف برآورده شود یا تعداد تکرارها از بودجه محاسباتی فراتر رود. تابع هزینه ما ۲-۳ شامل مجموعی از جملات فردی است که بر روی داده‌های آموزشی اجرا می‌شود؛ بنابراین مشتق جزئی $\nabla \text{Cost}(p)$ جمعی از مشتقات جزئی فردی بر روی داده‌های آموزشی است.

$$C_{x^{(i)}} = \frac{1}{2} \|y(x^{(i)}) - a^{[L]}(x^{(i)})\|_2^2, \quad (7-2)$$

سپس، از ۲-۳،

$$\nabla \text{Cost}(p) = \frac{1}{N} \sum_{i=1}^N \nabla C_{x^{(i)}}(p). \quad (8-2)$$

هنگامی که تعداد زیادی پارامتر و نقاط آموزشی داریم، محاسبه بردار گرادیان ۲-۸ در هر تکرار از روش نزول گرادیان ۲-۶ می‌تواند به‌شدت هزینه‌بر باشد. جایگزین بسیار ارزان‌تری این است که میانگین گرادیان‌های فردی را بر روی تمام نقاط آموزشی با گرادیان در یک نقطه آموزشی تصادفی انتخاب‌شده جایگزین کنیم. این منجر به ساده‌ترین شکل از روشی می‌شود که به آن روش گرادیان تصادفی^۳ می‌گویند که یک گام از آن را به‌صورت زیر خلاصه می‌شود:

۱. یک عدد صحیح i را به صورت تصادفی و یکنواخت از مجموعه $\{1, 2, 3, \dots, N\}$ انتخاب کنید.

۲. به‌روزرسانی کنید

$$p \rightarrow p - \eta \nabla C_{x^{(i)}}(p). \quad (9-2)$$

¹ Learning rate

² Steepest descent

³ Stochastic gradient

به زبان ساده، در هر مرحله از روش گرادیان تصادفی، یک نقطه آموزشی از میان تمامی داده‌ها به صورت تصادفی انتخاب می‌شود و از آن به عنوان نماینده کل مجموعه داده‌ها استفاده می‌شود. با پیشرفت هر مرحله، مدل داده‌های بیشتری را مشاهده می‌کند، به همین دلیل انتظار می‌رود که این روش، حتی اگر در هر مرحله کاهش هزینه زیادی به همراه نداشته باشد، در طول زمان به طور کلی مؤثر باشد.

نکته مهم این است که حتی با انتخاب یک مقدار بسیار کوچک برای نرخ یادگیری (η)، تضمینی وجود ندارد که تابع هزینه کلی کاهش یابد زیرا به جای میانگین‌گیری از تمام داده‌ها، از یک نمونه استفاده شده است؛ بنابراین، با وجود اینکه عبارت «نزول گرادیان تصادفی» به طور گسترده‌ای استفاده می‌شود، ما ترجیح می‌دهیم از عبارت «گرادیان تصادفی» استفاده کنیم تا دقت بیشتری در بیان داشته باشیم.

نسخه‌ای از روش گرادیان تصادفی که در رابطه ۲-۹ معرفی کردیم، ساده‌ترین شکل از یک طیف گسترده از روش‌های موجود است. در این روش، شاخص i با استفاده از نمونه‌گیری با جایگزینی انتخاب می‌شود—به این معنا که پس از استفاده از یک نقطه آموزشی، این نقطه به مجموعه آموزشی بازگردانده می‌شود و در گام بعدی به اندازه هر نقطه دیگر احتمال انتخاب شدن دارد.

یکی از جایگزین‌های موجود، نمونه‌گیری بدون جایگزینی است که در آن، نقاط آموزشی به ترتیب تصادفی و بدون تکرار انتخاب می‌شوند. انجام N گام به این شیوه، به عنوان تکمیل یک «دوره»^۱ شناخته می‌شود و می‌تواند به صورت زیر خلاصه شود:

۱. اعداد $\{1, 2, 3, \dots, N\}$ را به ترتیب جدیدی $\{k_1, k_2, k_3, \dots, k_N\}$ مرتب کنید.

۲. برای $i = 1$ تا N ، به روزرسانی کنید

$$p \rightarrow p - \eta \nabla C_{x_{\{k_i\}}}(p). \quad (10-2)$$

اگر روش گرادیان تصادفی را به عنوان تخمینی از میانگین روی تمام نقاط آموزشی در رابطه ۲-۸ با یک نمونه تکی در نظر بگیریم، منطقی است که یک راه میانه را نیز در نظر بگیریم که در آن از میانگین یک نمونه کوچک استفاده شود. برای برخی $m \ll N$ می‌توانیم گام‌هایی از فرم زیر برداریم:

۱. m عدد صحیح $\{k_1, k_2, k_3, \dots, k_m\}$ را به صورت تصادفی و یکنواخت از $\{1, 2, 3, \dots, N\}$ انتخاب کنید.

۲. به روزرسانی کنید

$$p \rightarrow p - \eta \frac{1}{m} \sum_{i=1}^m \nabla C_{x_{\{k_i\}}}(p). \quad (11-2)$$

¹epoch

در این روش، مجموعه‌ای از داده‌ها به نام $\{x^{(k_i)}\}_{i=1}^m$ به‌عنوان یک minibatch یا دسته کوچک شناخته می‌شود. یک روش دیگر، بدون جایگزینی داده‌ها، این است که فرض کنیم تعداد کل داده‌ها N برابر با Km است که در آن k تعداد دسته‌های کوچک^۱ است. در این روش، مجموعه داده‌ها به‌صورت تصادفی به k دسته کوچک مجزا تقسیم می‌شود و سپس به ترتیب بین این دسته‌ها حرکت می‌شود.

از آنجا که روش گرادیان تصادفی معمولاً برای محاسبات بسیار بزرگ استفاده می‌شود، انتخاب‌های الگوریتمی مانند اندازه دسته‌های کوچک و نحوه تصادفی‌سازی داده‌ها، بر اساس نیازهای سیستم‌های محاسباتی با عملکرد بالا تعیین می‌شود. همچنین، این امکان وجود دارد که در طول آموزش، اندازه دسته‌ها و دیگر پارامترها مانند نرخ یادگیری را به‌صورت دینامیکی تغییر دهیم تا فرایند همگرایی سریع‌تر انجام شود. در واقع، داشتن یک الگوریتم مناسب برای تعیین نرخ یادگیری، یکی از بخش‌های کلیدی برای اجرای مؤثر محاسبات است.

۲-۲ روش پس‌انتشار

اکنون در موقعیتی هستیم که می‌توانیم روش گرادیان تصادفی را برای آموزش یک شبکه عصبی مصنوعی به کار بگیریم. بردار عمومی پارامترها، \mathbf{p} که در بخش قبل استفاده شد، اکنون به ورودی‌های ماتریس‌های وزن و بردارهای بایاس تغییر می‌کند. وظیفه ما محاسبه مشتقات جزئی تابع هزینه نسبت به هر $w_{jk}^{[l]}$ و $b_j^{[l]}$ است.

روش گرادیان تصادفی از ساختار تابع هزینه بهره می‌برد: زیرا ۲-۳ ترکیبی خطی از عبارات فردی است که بر روی داده‌های آموزشی اجرا می‌شود، همین امر در مورد مشتقات جزئی آن نیز صدق می‌کند؛ بنابراین توجه خود را بر محاسبه آن مشتقات جزئی فردی متمرکز می‌کنیم.

برای یک نقطه آموزشی ثابت، $C_{x^{(i)}}$ در ۲-۷ به‌عنوان تابعی از وزن‌ها و بایاس‌ها در نظر گرفته می‌شود؛ بنابراین می‌توانیم وابستگی به $x^{(i)}$ را حذف کرده و به‌سادگی بنویسیم

$$C = \frac{1}{2} \|\mathbf{y} - \mathbf{a}^{[L]}\|_2^2. \quad (12-2)$$

به یاد می‌آوریم از ۲-۲ که $\mathbf{a}^{[L]}$ خروجی شبکه عصبی مصنوعی است. وابستگی C به وزن‌ها و بایاس‌ها تنها از طریق $\mathbf{a}^{[L]}$ حاصل می‌شود. برای استخراج عبارات مفید برای مشتقات جزئی، معرفی دو مجموعه متغیر دیگر سودمند است. ابتدا داریم

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]} \in \mathbb{R}^{n_l} \quad \text{for } l = 2, 3, \dots, L, \quad (13-2)$$

¹Minibatch

$z_j^{[l]}$ به عنوان ورودی وزن دار برای نورون j در لایه l اشاره می‌کند. رابطه بنیادی ۲-۲ که اطلاعات را در شبکه منتشر می‌کند، سپس

$$\mathbf{a}^{[l]} = \sigma(\mathbf{z}^{[l]}) \quad \text{for } l = 2, 3, \dots, L, \quad (14-2)$$

آنگاه، $\delta_j^{[l]} \in \mathbb{R}^{n_l}$ را به صورت

$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} \quad \text{for } 1 \leq j \leq n_l \quad \text{and } 2 \leq l \leq L. \quad (15-2)$$

عبارت بالا که اغلب به عنوان «خطای نورون j » در لایه l نامیده می‌شود، یک کمیت میانی است که برای تحلیل و محاسبات به کار می‌رود. با این حال، باید توجه داشت که استفاده از اصطلاح «خطا» در اینجا ممکن است کمی گمراه کننده باشد. در یک لایه میانی، مشخص نیست که تا چه حد می‌توان هر نورون را برای اختلافات در خروجی نهایی مسئول دانست. همچنین، در لایه خروجی L ، عبارت $\delta_j^{[L]}$ مستقیماً این اختلافات را اندازه گیری نمی‌کند. دلیل اینکه به $\delta_j^{[L]}$ در اینجا به عنوان «خطا» اشاره می‌شود، این است که تابع هزینه فقط زمانی می‌تواند به حداقل برسد که تمام مشتقات جزئی آن صفر باشند؛ بنابراین، داشتن $\delta_j^{[L]} = 0$ می‌تواند هدف مناسبی باشد. به طور کلی می‌توان گفت که $\delta_j^{[L]}$ در واقع حساسیت تابع هزینه نسبت به ورودی وزن دار نورون j در لایه l را اندازه گیری می‌کند.

در این مرحله، ما همچنین نیاز داریم که ضرب هادامار^۱ یا ضرب مؤلفه به مؤلفه دو بردار^۲ را تعریف کنیم. اگر $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ، آنگاه $\mathbf{x} \circ \mathbf{y} \in \mathbb{R}^n$ به صورت $(\mathbf{x} \circ \mathbf{y})_i = x_i y_i$ تعریف می‌شود. به زبان ساده، ضرب هادامار با ضرب جفت به جفت مؤلفه‌های متناظر تشکیل می‌شود.

لم ۱۰۲: داریم

$$\delta^{[L]} = \sigma'(z^{[L]}) \circ (\mathbf{a}^{[L]} - \mathbf{y}), \quad (16-2)$$

$$\delta^{[l]} = \sigma'(z^{[l]}) \circ (W^{[l+1]})^T \delta^{[l+1]} \quad \text{for } 2 \leq l \leq L-1, \quad (17-2)$$

$$\frac{\partial C}{\partial b_j^{[l]}} = \delta_j^{[l]} \quad \text{for } 2 \leq l \leq L, \quad (18-2)$$

¹Hadamard

²componentwise

$$\frac{\partial C}{\partial b_j^{[l]}} = \delta_j^{[l]} \quad \text{for } 2 \leq l \leq L. \quad (19-2)$$

اثبات: ابتدا با اثبات رابطه ۲-۱۶ شروع می‌کنیم. رابطه ۲-۱۴ با $l = L$ نشان می‌دهد که $a_j^{[L]}$ و $z_j^{[L]}$ توسط $a^{[L]} = \sigma(z^{[L]})$ به هم متصل هستند و از این رو

$$\frac{\partial a_j^{[L]}}{\partial z_j^{[L]}} = \sigma'(z_j^{[L]}),$$

همچنین، از رابطه ۲-۱۲ داریم

$$\frac{\partial C}{\partial a_j^{[L]}} = \frac{\partial}{\partial a_j^{[L]}} \frac{1}{2} \sum_{k=1}^{n_L} (y_k - a_k^{[L]})^2 = -(y_j - a_j^{[L]}),$$

بنابراین، با استفاده از قاعده زنجیره‌ای

$$\delta_j^{[L]} = \frac{\partial C}{\partial z_j^{[L]}} = \frac{\partial C}{\partial a_j^{[L]}} \frac{\partial a_j^{[L]}}{\partial z_j^{[L]}} = (a_j^{[L]} - y_j) \sigma'(z_j^{[L]}),$$

که این شکل جزئی از رابطه ۲-۱۶ است. برای نشان دادن رابطه ۲-۱۷، از قاعده زنجیره‌ای استفاده می‌کنیم تا

از $z_j^{[l]}$ به $\{z_k^{[l+1]}\}_{k=1}^{n_{l+1}}$ تبدیل کنیم. با اعمال قاعده زنجیره‌ای و استفاده از تعریف ۲-۱۵ داریم

$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \sum_{k=1}^{n_{l+1}} \frac{\partial C}{\partial z_k^{[l+1]}} \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} = \sum_{k=1}^{n_{l+1}} \delta_k^{[l+1]} \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}}. \quad (20-2)$$

حال، از رابطه ۲-۱۳ می‌دانیم که $z_k^{[l+1]}$ و $z_j^{[l]}$ به صورت زیر به هم متصل هستند

$$z_k^{[l+1]} = \sum_{s=1}^{n_l} w_{ks}^{[l+1]} \sigma(z_s^{[l]}) + b_k^{[l+1]},$$

بنابراین،

$$\frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} = w_{kj}^{[l+1]} \sigma'(z_j^{[l]}),$$

در رابطه ۲-۲۰ داریم

$$\delta_j^{[l]} = \sum_{k=1}^{n_{l+1}} \delta_k^{[l+1]} w_{kj}^{[l+1]} \sigma'(z_j^{[l]}),$$

که می‌توان آن را به صورت زیر بازنویسی کرد

$$\delta_j^{[l]} = \sigma'(z_j^{[l]}) \left((\mathbf{W}^{[l+1]})^T \delta^{[l+1]} \right)_j.$$

این شکل جزئی از رابطه ۱۷-۲ است. برای نشان دادن رابطه ۱۸-۲، توجه داریم که از روابط ۱۳-۲ و ۱۴-۲ می‌توان فهمید که $z_j^{[l]}$ به $b_j^{[l]}$ به صورت زیر متصل است

$$z_j^{[l]} = (\mathbf{W}^{[l]} \sigma(\mathbf{z}^{[l-1]}))_j + b_j^{[l]},$$

چون $\mathbf{z}^{[l-1]}$ به $b_j^{[l]}$ وابسته نیست،

$$\frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = 1,$$

سپس با استفاده از قاعده‌ی زنجیره‌ای،

$$\frac{\partial C}{\partial b_j^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} \frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} = \delta_j^{[l]},$$

با استفاده از تعریف ۲-۱۵. این نتیجه رابطه ۱۸-۲ را به دست می‌دهد. در نهایت، برای به دست آوردن رابطه ۲-۱۹، با نسخه‌ی جزء به جزء رابطه ۲-۱۵ شروع می‌کنیم:

$$z_j^{[l]} = \sum_{k=1}^{n_{l-1}} w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]},$$

که نتیجه می‌دهد

$$\frac{\partial z_j^{[l]}}{\partial w_{jk}^{[l]}} = a_k^{[l-1]}, \quad (21-2)$$

که به صورت مستقل از j است. و

$$\frac{\partial z_s^{[l]}}{\partial w_{jk}^{[l]}} = 0 \quad s \neq j. \quad (22-2)$$

به زبان ساده، روابط ۲۱-۲ و ۲۲-۲ به این دلیل حاصل می‌شوند که نورون j -ام در لایه l فقط از وزن‌های ردیف j -ام ماتریس $\mathbf{W}^{[l]}$ استفاده می‌کند و این وزن‌ها را به صورت خطی اعمال می‌کند. سپس، با استفاده از قاعده‌ی زنجیره‌ای، روابط ۲۱-۲ و ۲۲-۲ به صورت

$$\frac{\partial C}{\partial w_{jk}^{[l]}} = \sum_{s=1}^{n_l} \frac{\partial C}{\partial z_s^{[l]}} \frac{\partial z_s^{[l]}}{\partial w_{jk}^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} \frac{\partial z_j^{[l]}}{\partial w_{jk}^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} a_k^{[l-1]} = \delta_j^{[l]} a_k^{[l-1]},$$

□ که در مرحله‌ی آخر از تعریف $\delta_j^{[l]}$ در رابطه ۲-۱۵ استفاده شده است.

بسیاری از جنبه‌های لم ۱.۲ نیاز به توجه دارند. از روابط ۱-۲، ۱۳-۲ و ۱۴-۲ می‌توان یادآوری کرد که خروجی $\mathbf{a}^{[L]}$ را می‌توان از طریق یک عبور رو به جلو از شبکه محاسبه کرد، به ترتیب محاسبه‌ی

$$\mathbf{a}^{[1]}, \mathbf{z}^{[2]}, \mathbf{a}^{[2]}, \mathbf{z}^{[3]}, \dots, \mathbf{a}^{[L]}.$$

پس از انجام این کار، از رابطه‌ی ۲-۱۶ می‌بینیم که $\delta^{[L]}$ به طور مستقیم در دسترس است. سپس، با استفاده از رابطه‌ی ۲-۱۷، $\delta^{[2]}, \dots, \delta^{[l-2]}, \delta^{[l-1]}$ می‌تواند در عبور رو به عقب محاسبه شود. از روابط ۲-۱۸ و ۲-۱۹ می‌توانیم به مشتقات جزئی دسترسی پیدا کنیم. محاسبه‌ی گرادیان‌ها به این روش به نام پس‌انتشار^۱ شناخته می‌شود.

درک عمیق‌تر روابط پیش‌انتشار و مشتقات جزئی

برای درک بیشتر فرمول‌های پیش‌انتشار ۲-۱۸ و ۲-۱۹ در لم ۱.۲، مفید است که تعریف بنیادی مشتق جزئی را به یاد بیاوریم. مشتق جزئی $\frac{\partial C}{\partial w_{jk}^{[l]}}$ نشان می‌دهد که چگونه تابع هزینه C با تغییرات کوچک در وزن $w_{jk}^{[l]}$ تغییر می‌کند. به عنوان مثال، شکل ۲-۱ وزن $w_{43}^{[3]}$ را برجسته می‌کند. تغییر در این وزن تأثیری بر خروجی لایه‌های قبلی ندارد، بنابراین برای محاسبه $\frac{\partial C}{\partial w_{43}^{[3]}}$ نیازی به دانستن مشتقات جزئی در لایه‌های قبلی نداریم. اما باید بتوان $\frac{\partial C}{\partial w_{43}^{[3]}}$ را به کمک مشتقات جزئی در لایه‌های بعدی بیان کرد به طور دقیق‌تر، فعال‌سازی که به نورون چهارم در لایه ۳ تغذیه می‌شود $z_4^{[3]}$ است و طبق تعریف، $\delta_4^{[3]}$ حساسیت C نسبت به این ورودی را اندازه‌گیری می‌کند. ورودی به این نورون به صورت ثابت $+ w_{43}^{[3]} a_3^{[2]}$ است، بنابراین منطقی است که

$$\frac{\partial C}{\partial w_{43}^{[3]}} = \delta_4^{[3]} a_3^{[2]}, \quad (23-2)$$

به طور مشابه، از نظر بایاس، ثابت $+ b_4^{[3]}$ به نورون تغذیه می‌شود، که توضیح می‌دهد چرا

$$\frac{\partial C}{\partial b_4^{[3]}} = \delta_4^{[3]} \times 1. \quad (24-2)$$

برای اجتناب از استفاده از ماتریس‌های قطری در روابط ۲-۱۶ و ۲-۱۷، می‌توانیم از نماد ضرب هادامارد استفاده کنیم. اجازه دهید $D^{[l]} \in \mathbb{R}^{n_l \times n_l}$ ماتریس قطری باشد که عنصر (i, i) آن به صورت $\sigma'(z_i^{[l]})$ تعریف می‌شود. بنابراین داریم

$$\delta^{[L]} = D^{[L]}(a^{[L]} - y), \quad (25-2)$$

و

$$\delta^{[l]} = D^{[l]}(W^{[l+1]})^T \delta^{[l+1]}, \quad (26-2)$$

که می‌توان این را به صورت زیر گسترش داد

$$\delta^{[l]} = D^{[l]}(W^{[l+1]})^T D^{[l+1]}(W^{[l+2]})^T \dots D^{[L-1]}(W^{[L]})^T D^{[L]}(a^{[L]} - y). \quad (27-2)$$

¹Backward pass

همچنین از رابطه ۱-۲ به یاد داریم که محاسبه $\sigma'(z)$ ساده است.

رابطه ۲-۱۸ نشان می‌دهد که $\delta^{[l]}$ دقیقاً برابر با گرادیان تابع هزینه نسبت به بایاس‌ها در لایه l است. اگر به عنوان مؤلفه (j, k) در یک ماتریس از مشتقات جزئی در لایه l در نظر گرفته شود، آنگاه رابطه ۲-۱۹ نشان می‌دهد که این ماتریس حاصل ضرب خارجی $\delta^{[l]}(a^{[l-1]})^T \in \mathbb{R}^{n_l \times n_{l-1}}$ است. با ترکیب این موارد، می‌توانیم شبه کد زیر را برای الگوریتمی که شبکه‌ای را با استفاده از تعداد ثابتی از تکرارهای گرادیان نزولی تصادفی N_{iter} آموزش می‌دهد، بنویسیم. برای سادگی، نسخه ساده ۲-۹ را در نظر می‌گیریم که در آن نمونه‌ها به صورت با جایگزینی انتخاب می‌شوند. برای هر نقطه آموزشی، یک عبور روبه‌جلو از شبکه انجام می‌دهیم تا فعال‌سازی‌ها، ورودی‌های وزنی و خروجی کلی $a^{[L]}$ را ارزیابی کنیم. سپس یک عبور روبه‌عقب انجام می‌دهیم تا خطاها و به‌روزرسانی‌ها را محاسبه کنیم.

با این روند، برای هر نقطه آموزشی، ابتدا شبکه را در جهت پیشرو^۱ ارزیابی می‌کنیم و سپس خطاها را با استفاده از پیش انتشار^۲ محاسبه می‌کنیم. این روش به‌روزرسانی وزن‌ها و بایاس‌ها را برای بهبود عملکرد شبکه در پیش‌بینی داده‌های آموزشی ممکن می‌سازد.

مثال عملی

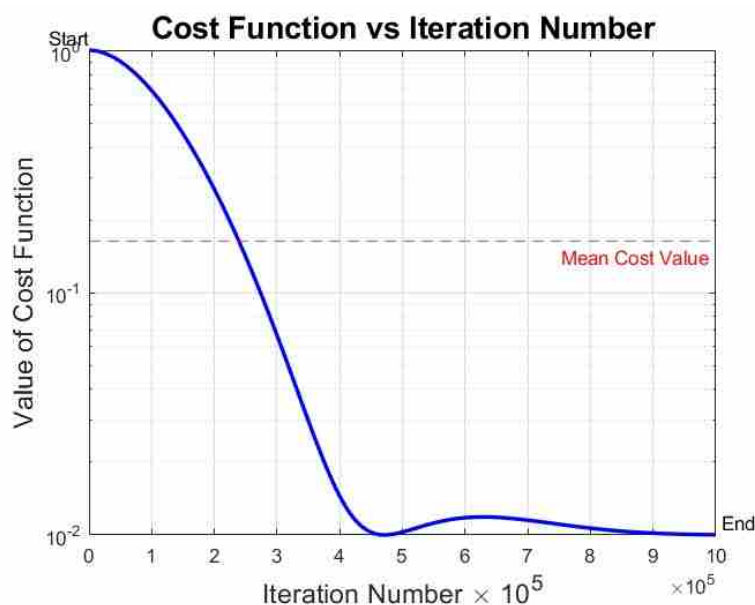
در این بخش، یک مثال عملی از پیش انتشار و روش گرادیان نزولی تصادفی را ارائه می‌دهیم. <<Listing1>> نشان می‌دهد که چگونه یک شبکه عصبی با ساختار نشان‌داده‌شده در شکل ۱-۵ می‌تواند بر روی داده‌های موجود در شکل ۱-۴ مورد استفاده قرار گیرد. تعداد لایه‌ها به طور ثابت تعیین شده و مراحل عبور پیشرو و پس رو خطبه‌خط اجرا شده‌اند. از آنجا که وزن‌ها و بایاس‌ها در هر لایه ابعاد یکسانی ندارند، ذخیره آن‌ها در یک آرایه سه‌بعدی مناسب نیست. می‌توانیم از یک آرایه سلولی یا ساختاری استفاده کنیم و سپس عملیات پیشرو و پس رو را در حلقه‌های `for` پیاده‌سازی کنیم. تابع netbp در <<Listing1>> شامل یک تابع تودرتو به نام cost است که نسخه‌ای مقیاس شده از هزینه در رابطه ۱-۹ را محاسبه می‌کند. از آنجا که این تابع تودرتو است، به متغیرهای موجود در تابع اصلی، به‌ویژه داده‌های آموزشی، دسترسی دارد. لازم به ذکر است که تابع تو در تو cost به طور مستقیم در مراحل پیشرو و پس رو استفاده نمی‌شود. این تابع در هر تکرار از روش گرادیان نزولی تصادفی فراخوانی می‌شود تا پیشرفت آموزش را پایش کنیم.

کد <<Listing2>> تابع activate را نشان می‌دهد که توسط netbp استفاده می‌شود و تابع سیگموئید را به صورت برداری اعمال می‌کند. در ابتدای netbp، داده‌های آموزشی و مقادیر هدف y را همان‌طور که در رابطه

¹Forward pass

²Backward pass

۱-۸ تعریف شده‌اند، تنظیم می‌کنیم. سپس تمام وزن‌ها و بایاس‌ها را با استفاده از مولد اعداد تصادفی نرمال randn مقاداردهی اولیه می‌کنیم. برای سادگی، نرخ یادگیری ثابت $\eta = 0.05$ را تنظیم کرده و تعداد ثابتی از تکرارها را با مقدار $N_{iter} = 1e6$ انجام می‌دهیم. ما از روش تکرار گرادیان تصادفی ساده که در پایان بخش ۵ خلاصه شده است، استفاده می‌کنیم. در اینجا، دستور randi(10) یک عدد صحیح را به طور یکنواخت و مستقل از بین ۱ تا ۱۰ انتخاب می‌کند. پس از ذخیره‌سازی مقدار تابع هزینه در هر تکرار، از دستور semilogy برای نمایش پیشرفت تکرار استفاده می‌کنیم. در این آزمایش، حدس اولیه ما برای وزن‌ها و بایاس‌ها مقدار ۳.۵ را برای تابع هزینه به دست آورد. پس از 10^6 مرحله گرادیان تصادفی، این مقدار به 7.3×10^{-4} کاهش یافت. شکل ۲-۲ نمودار semilogy را نشان می‌دهد و می‌بینیم که کاهش هزینه یک‌نواخت نیست.



شکل ۲-۲: تابع هزینه در مقابل تعداد تکرار

این مثال MATLAB نشان می‌دهد که چگونه می‌توان یک شبکه عصبی ساده را با استفاده از پیش انتشار و پس انتشار آموزش داد. نتایج نشان می‌دهند که این روش گرادیان نزولی تصادفی می‌تواند به تدریج خطای شبکه را کاهش داده و به بهبود عملکرد پیش‌بینی کمک کند. این روند کاهش هزینه در شکل ۲-۲ نمایش داده شده است و نشان می‌دهد که اگرچه فرایند کاهش هزینه یک‌نواخت نیست، اما با گذر زمان و تعداد تکرارهای زیاد، بهبود قابل توجهی حاصل می‌شود.

```

1
2 function netbp
3 % NETBP - Neural network using backpropagation for training
4 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5
6 %%%%%%%%% Initial Setup %%%%%%%%%
7 % Input data
8 x1 = [0.1, 0.3, 0.1, 0.6, 0.4, 0.6, 0.5, 0.9, 0.4, 0.7];
9 x2 = [0.1, 0.4, 0.5, 0.9, 0.2, 0.3, 0.6, 0.2, 0.4, 0.6];
10
11 % Output labels
12 y = [ones(1, 5), zeros(1, 5); zeros(1, 5), ones(1, 5)];
13
14 % Initialize weights and biases
15 rng(5000);
16 W2 = initializeWeights(2, 2);
17 W3 = initializeWeights(3, 2);
18 W4 = initializeWeights(2, 3);
19 b2 = initializeWeights(2, 1);
20 b3 = initializeWeights(3, 1);
21 b4 = initializeWeights(2, 1);
22
23 % Training settings
24 eta = 0.05; % Learning rate
25 Niter = 1e6; % Number of iterations
26 savecost = zeros(Niter, 1); % Cost function value at each iteration
27
28 % Training loop
29 for counter = 1:Niter
30 k = randi(10); % Choose a random training sample
31 x = [x1(k); x2(k)];
32
33 % Forward pass
34 [a2, a3, a4] = forwardPass(x, W2, W3, W4, b2, b3, b4);
35
36 % Backward pass and weight update
37 [W2, W3, W4, b2, b3, b4] = backPropagation(x, y(:, k), a2, a3, a4,
38 W2, W3, W4, b2, b3, b4, eta);
39
40 % Compute and save cost function value
41 savecost(counter) = computeCost(x1, x2, y, W2, W3, W4, b2, b3, b4);
42 end
43
44 % Plot cost function decay
45 semilogy(1:1e4:Niter, savecost(1:1e4:Niter));
46 xlabel('Iteration');
47 ylabel('Cost');
48 title('Cost Function Decay');
49 grid on;
50
51 %%%%%%%%% Helper Functions %%%%%%%%%
52
53 % Function to initialize weights
54 function W = initializeWeights(rows, cols)
55 W = 0.5 * randn(rows, cols);
56 end

```

```

57 % Function for forward pass
58 function [a2, a3, a4] = forwardPass(x, W2, W3, W4, b2, b3, b4)
59 a2 = activate(x, W2, b2);
60 a3 = activate(a2, W3, b3);
61 a4 = activate(a3, W4, b4);
62 end
63
64 % Function for backpropagation and updating weights
65 function [W2, W3, W4, b2, b3, b4] = backPropagation(x, y, a2, a3,
66     a4, W2, W3, W4, b2, b3, b4, eta)
67 delta4 = a4 .* (1 - a4) .* (a4 - y);
68 delta3 = a3 .* (1 - a3) .* (W4' * delta4);
69 delta2 = a2 .* (1 - a2) .* (W3' * delta3);
70
71 W4 = W4 - eta * delta4 * a3';
72 W3 = W3 - eta * delta3 * a2';
73 W2 = W2 - eta * delta2 * x';
74 b4 = b4 - eta * delta4;
75 b3 = b3 - eta * delta3;
76 b2 = b2 - eta * delta2;
77 end
78
79 % Function to compute cost
80 function costval = computeCost(x1, x2, y, W2, W3, W4, b2, b3, b4)
81 costvec = zeros(10, 1);
82 for i = 1:10
83     x = [x1(i); x2(i)];
84     [a2, a3, a4] = forwardPass(x, W2, W3, W4, b2, b3, b4);
85     costvec(i) = norm(y(:, i) - a4, 2);
86 end
87 costval = norm(costvec, 2)^2;
88 end
89
90 % Activation function (sigmoid)
91 function a = activate(x, W, b)
92 a = 1 ./ (1 + exp(-(W * x + b)));
93 end
94
95 % Second MATLAB Code Block

```

Listing 1: Neural Network with Backpropagation

```
1 function y = activate(x, W, b)
2 % ACTIVATE - Computes the sigmoid function output
3 %
4 % This function takes the input vector x and computes the sigmoid
5 % output by applying weights and biases.
6 %
7 % Inputs:
8 %   x - Input vector (n x 1)
9 %   W - Weight matrix (m x n)
10 %   b - Bias vector (m x 1)
11 %
12 % Output:
13 %   y - Sigmoid output vector (m x 1)
14 %
15 % Computation:
16 %   y = 1 ./ (1 + exp(-(W * x + b)));
17
18
19 % Compute sigmoid output
20 z = W * x + b; % Linear combination of inputs
21 y = 1 ./ (1 + exp(-z)); % Apply sigmoid function to the linear
22 % combination
end
```

Listing 2: Activation Function

فصل ۳

کاربردهای عملی شبکه‌های عصبی و ابزارهای پیاده‌سازی

در این بخش به یک مسئله کاربردی‌تر می‌پردازیم که به ما امکان می‌دهد قابلیت‌های رویکرد یادگیری عمیق را به صورت مؤثرتری به نمایش بگذاریم. برای این منظور، از جعبه‌ابزار MatConvNet استفاده می‌شود که برای ارائه بلوک‌های اصلی یادگیری عمیق به صورت دستورات ساده MATLAB طراحی شده است؛ بنابراین این ابزار، محیطی ایده‌آل برای نمونه‌سازی سریع و کاربردهای آموزشی به شمار می‌آید. علاوه بر این، پشتیبانی از GPU ها همچنین MatConvNet را برای محاسبات در مقیاس بزرگ کارآمد می‌سازد و شبکه‌های از پیش آموزش‌دیده نیز به سادگی قابل استفاده هستند. استفاده از MatConvNet در یک مسئله در مقیاس بزرگ همچنین این فرصت را فراهم می‌کند تا مفاهیم بیشتری را که به محاسبات عملی را مورد بررسی قرار دهیم.

۱-۳ شبکه‌های عصبی پیچشی: اصول، ساختار و کاربردها

جعبه ابزار MatConvNet از یک دسته ویژه از شبکه‌های عصبی مصنوعی به نام شبکه‌های عصبی پیچشی^۱ (CNNs) استفاده می‌کند که به یک ابزار استاندارد در کاربردهای بینایی کامپیوتر تبدیل شده‌اند. به عنوان مثال، یک تصویر رنگی که از 200×200 پیکسل تشکیل شده و هر پیکسل دارای اجزای قرمز، سبز و آبی است را در نظر بگیرید. این به یک بردار ورودی با ابعاد

$$n_1 = 200 \times 200 \times 3 = 120,000$$

^۱Convolutional Neural Network

و در نتیجه یک ماتریس وزن $W^{[2]}$ در سطح ۲ با $120,000$ ستون منجر می‌شود. اگر وزن‌ها و بایاس‌ها را به‌طور کلی در نظر بگیریم، این رویکرد به‌وضوح غیرممکن است. CNN‌ها با محدود کردن مقادیر مجاز در ماتریس‌های وزن و بردارهای بایاس، این مسئله را به‌طور مؤثری مدیریت می‌کنند. به‌جای یک تبدیل خطی کامل، CNN‌ها یک هسته یا فیلتر خطی کوچک را به‌طور مکرر در بخش‌های مختلف داده‌های ورودی خود اعمال می‌کنند. در واقع، ماتریس‌های وزن مورد استفاده در CNN‌ها بسیار پراکنده و با ساختار خاص هستند.

برای درک اینکه چرا این رویکرد ممکن است مفید باشد، در نظر بگیرید که یک بردار ورودی در \mathbb{R}^6 را با

ماتریس

$$\begin{bmatrix} 1 & -1 & & & & \\ & 1 & -1 & & & \\ & & 1 & -1 & & \\ & & & 1 & -1 & \\ & & & & 1 & -1 \\ & & & & & 1 \end{bmatrix} \in \mathbb{R}^{5 \times 6} \quad (1-3)$$

ضرب کنیم. این کار یک بردار در \mathbb{R}^5 تولید می‌کند که از تفاوت‌های بین مقادیر مجاور تشکیل شده است. در این مورد، از یک فیلتر $[1, -1]$ و گام طول یک استفاده می‌شود (فیلتر بعد از هر بار استفاده به یک مکان جلو می‌رود). تعمیم‌های مناسب این ماتریس به حالت بردارهای ورودی ناشی از تصاویر دو بعدی می‌تواند برای شناسایی لبه‌ها در تصویر استفاده شود (که در صورت وجود تغییر ناگهانی در مقادیر پیکسل‌های مجاور، یک مقدار بزرگ را تولید می‌کند). حرکت دادن یک فیلتر در سراسر یک تصویر می‌تواند ویژگی‌های دیگری مانند انواع خاصی از منحنی‌ها یا لکه‌های هم‌رنگ را نیز آشکار کند؛ بنابراین، با مشخص کردن اندازه فیلتر و طول گام، می‌توانیم اجازه دهیم که فرایند آموزش وزن‌های فیلتر را به‌عنوان ابزاری برای استخراج ساختارهای مفید یاد بگیرد.

اصطلاح «پیچشی»^۱ به این دلیل به کار می‌رود که تبدیل‌های خطی مرتبط با این شبکه‌ها را می‌توان به‌صورت

پیچشی نوشت. در حالت یک‌بعدی، کانولوشن بردار $x \in \mathbb{R}^p$ با فیلتر $g_1, g_2, \dots, g_{p-2}, g_{p-1}$ دارای مؤلفه

k -ام به‌صورت زیر

$$y_k = \sum_{n=1}^p x_n g_{k-n}. \quad (2-3)$$

مثال موجود در ۳-۱ مربوط به فیلتری با $g_0 = 1$ ، $g_{-1} = -1$ ، و سایر $g_k = 0$ است. در حالتی که

^۱Convolutiona

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} a & b & c & d & & & & & & \\ & a & b & c & d & & & & & \\ & & a & b & c & d & & & & \\ & & & a & b & c & d & & & \\ & & & & a & b & c & d & & \\ & & & & & a & b & c & d & \\ & & & & & & a & b & c & d \\ & & & & & & & a & b & c & d \\ & & & & & & & & & & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \quad (3-3)$$

ما از فیلتری با چهار وزن a, b, c, d استفاده می‌کنیم که دارای طول گام دو است. به دلیل عدم تطابق طول بردار ورودی x با طول فیلتر، یک مقدار صفر به آن اضافه کرده‌ایم.

۳-۲ پیشگیری از بیش‌برازش در پردازش تصاویر

در عمل، داده‌های تصویری به طور معمول به‌عنوان یک تانسور^۱ سه‌بعدی در نظر گرفته می‌شوند به طوری که هر پیکسل دارای دو مختصات فضایی و یک مقدار قرمز/سبز/آبی است. از این دیدگاه، فیلتر به شکل یک تانسور کوچک است که به طور متوالی به بخش‌های مختلف تانسور ورودی اعمال می‌شود و عملیات پیچشی مربوطه به‌صورت چندبعدی انجام می‌گیرد. از دید محاسباتی، یکی از مزایای کلیدی شبکه‌های عصبی پیچشی این است که محاسبه حاصل ضرب‌های ماتریس - بردار که در مراحل پیشرو و پس‌رو در شبکه دخیل هستند، می‌تواند به طور بسیار کارآمدی با استفاده از تکنیک‌های تبدیل سریع انجام شود.

یک لایه پیچشی معمولاً با یک لایه ادغام^۲ دنبال می‌شود که بعد را با نگاشت مناطق کوچک از پیکسل‌ها به اعداد منفرد کاهش می‌دهد. به‌عنوان مثال، وقتی این مناطق کوچک به‌صورت مربع‌هایی از چهار پیکسل همسایه در یک تصویر دوبعدی انتخاب می‌شوند، یک لایه بیشینه ادغام^۳ هر مجموعه از چهار پیکسل را به ترتیب با بیشترین یا میانگین مقدار آن‌ها جایگزین می‌کند.

۳-۲-۱ بیش‌برازش

بیش‌برازش^۴ زمانی رخ می‌دهد که یک شبکه آموزش دیده روی داده‌های موجود عملکرد بسیار دقیقی دارد، اما نمی‌تواند به خوبی به داده‌های جدید تعمیم یابد. به‌طور کلی، این بدان معنی است که فرایند برازش بیش از حد روی «نویز»های غیرمهم و غیرنمایشی در داده‌های موجود متمرکز شده است. روش‌های بسیاری برای مقابله با بیش‌برازش پیشنهاد شده است که برخی از آن‌ها می‌توانند به طور هم‌زمان استفاده شوند. یک تکنیک مفید این است که داده‌های موجود را به دو گروه مجزا تقسیم کنید.

• **داده‌های آموزشی:**^۵ داده‌های آموزشی برای تعریف تابع هزینه‌ای که مسئله بهینه‌سازی را مشخص می‌کند،

¹Tensor

²pooling

³Max Pooling

⁴Over-fitting

⁵Training Data

استفاده می‌شوند. این داده‌ها فرآیند یادگیری را هدایت می‌کنند، که در طی آن وزن‌ها به‌طور تکراری به‌روزرسانی می‌شوند.

• **داده‌های اعتبارسنجی:**^۱ این داده‌ها در فرآیند بهینه‌سازی استفاده نمی‌شوند (آن‌ها هیچ تأثیری در نحوه به‌روزرسانی وزن‌ها در هر مرحله ندارند). ما از داده‌های اعتبارسنجی فقط برای ارزیابی عملکرد شبکه فعلی استفاده می‌کنیم. در هر مرحله از بهینه‌سازی، می‌توانیم تابع هزینه مربوط به داده‌های اعتبارسنجی را ارزیابی کنیم. این عملکرد مجموعه فعلی از وزن‌های آموزش‌دیده روی داده‌های نادیده را اندازه‌گیری می‌کند.

به‌طور شهودی، بیش‌برازش به وضعیتی اشاره دارد که فرآیند بهینه‌سازی در حال کاهش تابع هزینه خود (ارائهٔ برازش بهتر به داده‌های آموزشی) است، اما تابع هزینه برای خطای اعتبارسنجی دیگر در حال کاهش نیست (بنابراین عملکرد در داده‌های نادیده بهبود نمی‌یابد)؛ بنابراین، منطقی است که آموزش را در مرحله‌ای متوقف کنیم که هیچ بهبودی در داده‌های اعتبارسنجی مشاهده نمی‌شود.

یک روش محبوب دیگر برای مقابله با بیش‌برازش، حذف تصادفی و مستقل نورون‌ها در طول فاز آموزش است. به‌عنوان مثال، در هر مرحله از روش گرادینان تصادفی، می‌توانیم هر نورون را با احتمال p حذف کنیم و شبکه باقی‌مانده را آموزش دهیم. در پایان فرآیند، به دلیل اینکه وزن‌ها و بایاس‌ها بر روی این شبکه‌های کوچک‌تر آموزش‌دیده‌اند، می‌توانیم هر یک را برای استفاده در شبکه با اندازه کامل، در عامل p ضرب کنیم. این تکنیک که به نام حذف تصادفی^۲ شناخته می‌شود، این روش به‌نوعی تفسیر شهودی اشاره دارد «ما با ایجاد میانگینی از بسیاری از شبکه‌های آموزش‌دیده، به اجماعی دست می‌یابیم که از هر یک از شبکه‌های فردی به‌تنهایی قابل‌اعتمادتر است».

۳-۳ تحلیل عمیق توابع هزینه و فعال‌سازی در شبکه‌های عصبی

در فصل‌های قبل، ما از توابع فعال‌سازی به شکل سیگموئید ۱-۱ و تابع هزینهٔ درجه دوم ۲-۳ استفاده کردیم. انتخاب‌های متداول دیگری نیز وجود دارد و عملکرد نسبی آن‌ها بستگی به کاربرد خاص دارد. در زمینهٔ طبقه‌بندی تصویر، استفاده از یک واحد خطی یک‌سو شده^۳ معمول است که به‌صورت زیر تعریف می‌شود:

$$\sigma(x) = \begin{cases} 0 & \text{for } x \leq 0, \\ x & \text{for } x > 0. \end{cases} \quad (۴-۳)$$

^۱Validation Data

^۲Dropout

^۳Rectified Linear Unit (ReLU)

۴-۳ بهینه‌سازی دسته‌بندی تصاویر

اگر داده‌های آموزشی ما $\{x^{(i)}\}_{i=1}^N$ شامل K دسته‌بندی مختلف با برچسب‌های مشخص باشد، برچسب مربوط به داده $x^{(i)}$ را با $l_i \in \{1, 2, \dots, K\}$ نمایش می‌دهیم. به‌جای استفاده از تابع هزینه درجه دوم ۲-۳، می‌توانیم از روش بیشینه هموار^۱ به شکل زیر استفاده کنیم. خروجی $v^{(i)} = a^{[L]}(x^{(i)})$ از شبکه به صورت یک بردار در \mathbb{R}^K است، به طوری که مؤلفه j -ام از این بردار بزرگ‌تر است زمانی که تصویر متعلق به دسته j باشد

$$(v^{(i)})_s \mapsto \frac{e^{v_s^{(i)}}}{\sum_{j=1}^K e^{v_j^{(i)}}.$$

این کار، مؤلفه‌های بزرگ را تقویت کرده و یک بردار از وزن‌های مثبت تولید می‌کند که مجموع آن‌ها برابر با یک است و می‌توان آن را به عنوان احتمال تفسیر کرد. اکنون هدف ما این است که مقدار بیشینه هموار برای نقطه آموزشی $x^{(i)}$ در مؤلفه l_i (که همان برچسب صحیح است) تا حد ممکن به عدد یک نزدیک باشد. به‌جای استفاده از تابع هزینه درجه دوم، از یک معیار لگاریتمی استفاده می‌کنیم و به تابع هزینه زیر می‌رسیم

$$-\sum_{i=1}^N \log \left(\frac{e^{v_{l_i}^{(i)}}}{\sum_{j=1}^K e^{v_j^{(i)}} \right). \quad (۵-۳)$$

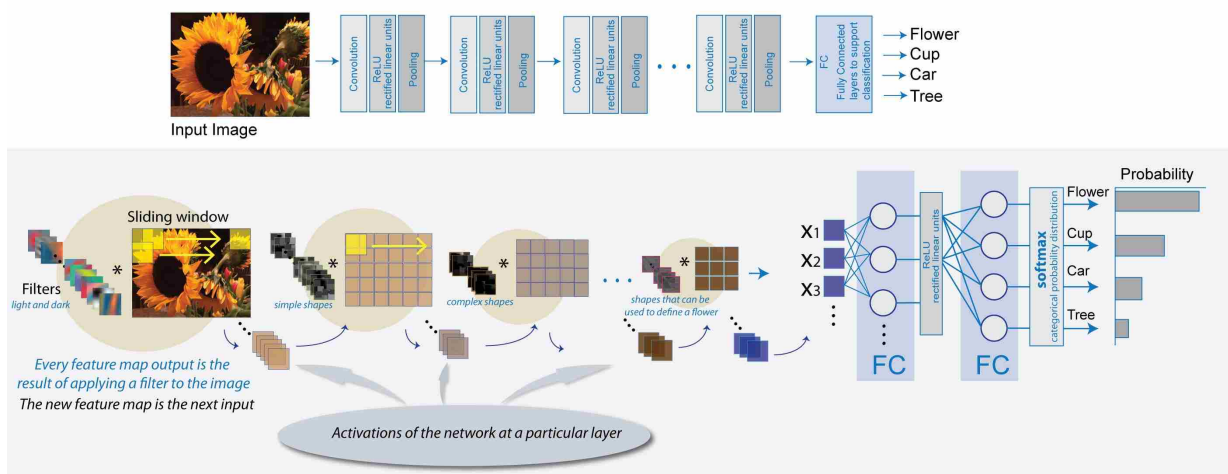
۳-۴-۱ پیاده‌سازی و ارزیابی مدل‌های دسته‌بندی تصویر

اکنون نتایج یک کار یادگیری نظارت شده در دسته‌بندی تصویر را نشان می‌دهیم. ما مجموعه‌ای از تصاویر را در نظر می‌گیریم که هر کدام دقیقاً در یکی از ده دسته: هواپیما، اتومبیل، پرنده، گربه، گوزن، سگ، قورباغه، اسب، کشتی، کامیون قرار می‌گیرند. ما از داده‌های برچسب‌دار مجموعه CIFAR-۱۰ استفاده می‌کنیم. تصاویر کوچک هستند و دارای ۳۲ در ۳۲ پیکسل هستند که هر کدام دارای یک مؤلفه قرمز، سبز و آبی هستند؛ بنابراین یک قطعه داده آموزشی شامل $3 \times 32 \times 32 = 3072$ مقدار است. ما از یک مجموعه آموزشی شامل ۵۰,۰۰۰ تصویر و ۱۰,۰۰۰ تصویر دیگر به عنوان مجموعه اعتبارسنجی استفاده می‌کنیم. پس از اتمام بهینه‌سازی و آموزش شبکه، عملکرد آن را روی مجموعه جدیدی از ۱۰,۰۰۰ تصویر تست، با ۱,۰۰۰ تصویر از هر دسته، ارزیابی می‌کنیم. توجه داشته باشید که استراتژی‌های اعتبارسنجی متقاطع پیشرفته‌تری می‌توانند مورد استفاده قرار گیرند، جایی که همان مجموعه داده کلی به طور مکرر به مجموعه‌های مجزا برای آموزش، اعتبارسنجی و تست تقسیم می‌شود. بر اساس معماری مورد استفاده در کدهای مربوط به جعبه‌ابزار MatConvNet، شبکه‌ای را راه‌اندازی

^۱ Softmax

می‌کنیم که لایه‌های آن به پنج بلوک تقسیم می‌شوند. در اینجا ابعاد ورودی‌ها/خروجی‌ها و وزن‌ها را به صورت فشرده در قالب تانسور توصیف می‌کنیم.

بلوک ۱ شامل یک لایه پیچشی است که با یک لایه ادغام و فعال‌سازی دنبال می‌شود. این تبدیل ورودی اصلی $32 \times 32 \times 3$ را به بعد $16 \times 16 \times 32$ تبدیل می‌کند. به طور دقیق‌تر، لایه پیچشی از فیلترهای 5×5 استفاده می‌کند که همچنین در سه کانال رنگی اسکن می‌کنند. ۳۲ فیلتر مختلف وجود دارد، بنابراین به طور کلی وزن‌ها می‌توانند در یک آرایه $32 \times 32 \times 3 \times 5 \times 5$ نمایش داده شوند. خروجی از هر فیلتر ممکن است به عنوان یک نقشه ویژگی توصیف شود. فیلترها با یک طول گام واحد^۱ اعمال می‌شوند. به این ترتیب، هر یک از ۳۲ نقشه ویژگی دارای ابعاد 32×32 است. سپس بیشینه ادغام با استفاده از طول گام دو به هر نقشه ویژگی اعمال می‌شود. این ابعاد نقشه‌های ویژگی را به 16×16 کاهش می‌دهد. سپس از فعال‌سازی $ReLU$ استفاده می‌شود. [۳]



شکل ۳-۱: نمای کلی معماری شبکه

بلوک ۲: در این بلوک ابتدا عملیات پیچشی اعمال می‌شود، سپس تابع فعال‌سازی و در نهایت یک لایه ادغام‌کننده. این کار ابعاد را به $8 \times 8 \times 32$ کاهش می‌دهد. به طور دقیق‌تر، ما از ۳۲ فیلتر استفاده می‌کنیم که هر کدام 5×5 هستند و در تمام ۳۲ نقشه ویژگی اسکن می‌کنند؛ بنابراین وزن‌ها را می‌توان به صورت یک تانسور $32 \times 32 \times 5 \times 5$ در نظر گرفت. طول گام برابر یک است، بنابراین نقشه‌های ویژگی ۳۲ تایی بعد از کانولوشن همچنان ابعاد 16×16 دارند. پس از اعمال تابع فعال‌سازی $ReLU$ ، یک لایه میانگین ادغام با گام دو اعمال می‌شود که ابعاد هر کدام از نقشه‌های ویژگی را به 8×8 کاهش می‌دهد.

بلوک ۳: این بلوک ابتدا یک لایه پیچشی و سپس تابع فعال‌سازی و عملیات ادغام را اعمال می‌کند که ابعاد

¹ Stride Length

را به $4 \times 4 \times 64$ کاهش می‌دهد. به طور دقیق‌تر، ۶۴ فیلتر به کار می‌روند. هر فیلتر 5×5 است و در تمام ۳۲ نقشه ویژگی اسکن می‌کند؛ بنابراین وزن‌ها را می‌توان به صورت یک تانسور $5 \times 5 \times 32 \times 64$ در نظر گرفت. طول گام برابر یک است و نقشه‌های ویژگی حاصل ابعاد 8×8 دارند. پس از فعال‌سازی *ReLU*، یک لایه میانگین ادغام با گام دو اعمال می‌شود که ابعاد هر کدام از ۶۴ نقشه ویژگی را به 4×4 کاهش می‌دهد.

بلوک ۴: این بلوک فقط از پیچشی و سپس فعال‌سازی استفاده می‌کند، و ابعاد را به $1 \times 1 \times 64$ کاهش می‌دهد. به طور دقیق‌تر، از ۶۴ فیلتر استفاده می‌شود. هر فیلتر 4×4 است و در ۶۴ نقشه ویژگی اعمال می‌شود؛ بنابراین وزن‌ها را می‌توان به صورت یک تانسور $4 \times 4 \times 64 \times 64$ در نظر گرفت و هر فیلتر یک عدد واحد تولید می‌کند.

بلوک ۵: این بلوک شامل عملیات پیچشی نیست. از یک ماتریس وزن عمومی (به طور کامل متصل) از نوعی که در بخش‌های ۲ تا ۶ بررسی شد، استفاده می‌کند تا خروجی با ابعاد $1 \times 1 \times 10$ را بدهد. این معادل یک ماتریس وزن با ابعاد 10×64 است. در نهایت، تابع فعال‌سازی بیشینه هموار هر کدام از ده مؤلفه خروجی را به بازه $[0, 1]$ تبدیل می‌کند. شکل ۱.۷ نمای کلی معماری شبکه را به صورت بصری نشان می‌دهد.

خروجی شبکه یک بردار شامل ده عدد حقیقی است. تابع هزینه برای مسئله بهینه‌سازی به صورت لاگ‌لاس با سافت مکس و $K = 10$ در نظر گرفته شده است. برای بهینه‌سازی از روش گرادیان تصادفی با momentum استفاده می‌کنیم که در آن «میانگین متحرک» جهت‌های گرادیان فعلی و گذشته محاسبه می‌شود. داده‌ها به صورت دسته‌های کوچک^۱ با اندازه ۱۰۰ مورد استفاده قرار می‌گیرند و تعداد ثابتی از ۴۵ دوره (*epoch*) برای آموزش تعیین شده است. نرخ یادگیری به صورت پیش فرض برای هر دوره به این شکل تعریف می‌شود: $\eta = 0.05$ برای ۳۰ دوره اول، $\eta = 0.0005$ برای ۱۰ دوره بعدی، و $\eta = 0.0005$ برای ۵ دوره نهایی. اجرای این شبکه روی یک پردازنده گرافیکی Tesla C۲۰۷۵ در حالت دقت تک، ۴۵ دوره^۲ را در کمتر از چهار ساعت تکمیل می‌کند. به عنوان آزمایشی اضافی، ما همچنین شبکه را با استفاده از تکنیک حذف تصادفی^۳ آموزش می‌دهیم. در اینجا، در هر مرحله گرادیان تصادفی، خروجی هر نورون با یک احتمال مستقل به صفر بازنشانی می‌شود.

۳-۴-۲ تحلیل نتایج و تاثیر حذف تصادفی در بهبود دقت مدل

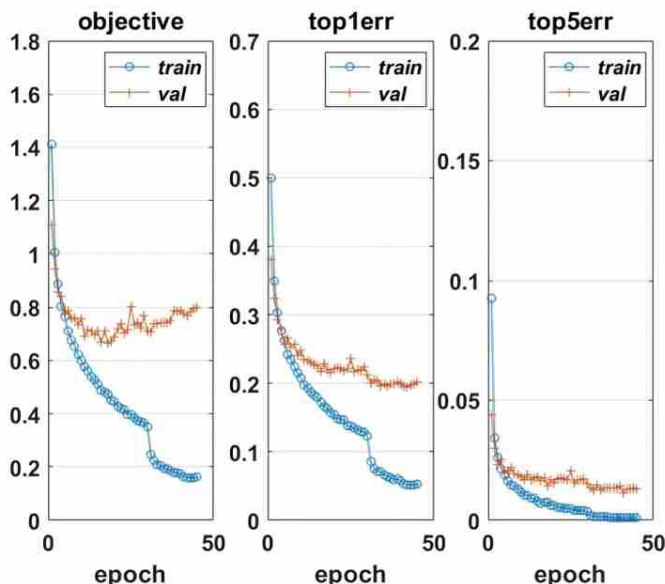
در شکل ۳-۲، فرایند آموزش را در حالت بدون حذف تصادفی نشان می‌دهیم. در نمودار سمت چپ، دایره‌ها نشان‌دهنده کاهش تابع هدف ۳-۵ پس از هر یک از ۴۵ دوره آموزشی هستند. همچنین از ضربدرها برای نمایش مقدار تابع هدف بر روی داده‌های اعتبارسنجی استفاده می‌شود. (این اندازه‌گیری‌ها به صورت میانگین

¹Minibatches

²Epoch

³Dropout

بر روی بچ‌هایی که دوره آموزش را تشکیل می‌دهند، محاسبه شده‌اند. توجه داشته باشید که وزن‌ها پس از هر بچ به‌روزرسانی می‌شوند.) هدف کلی ما طبقه‌بندی تصاویر به یکی از ده دسته است، و نمودار میانی در شکل ۲-۳ درصد خطاهای طبقه‌بندی را با بالاترین احتمال نشان می‌دهد. به‌طور مشابه، نمودار سمت راست درصد مواردی را نشان می‌دهد که دسته صحیح در میان پنج دسته برتر نیست.

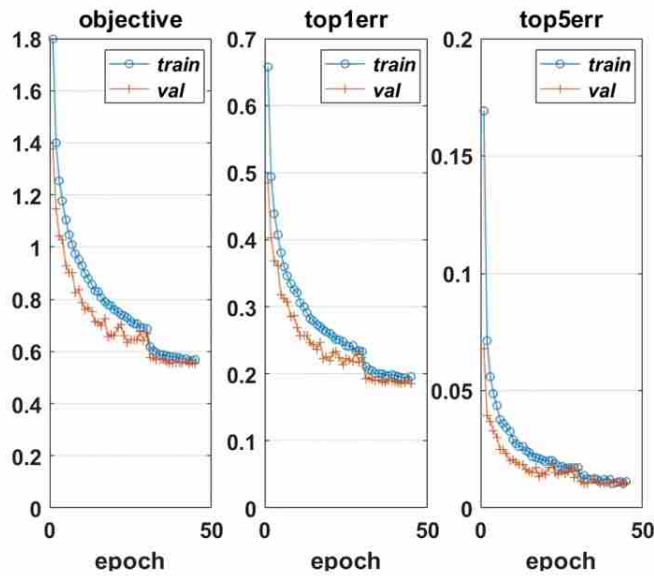


شکل ۲-۳: آموزش شبکه بدون حذف تصادفی

از شکل ۲-۳ مشاهده می‌کنیم که خطای اعتبارسنجی در مرحله‌ای به حالت ثابت می‌رسد، درحالی‌که روش گرادینت تصادفی همچنان کاهش‌های قابل‌توجهی در خطای آموزش ایجاد می‌کند. این نشان‌دهنده بیش‌برازش است یعنی شبکه در حال یادگیری جزئیات دقیق داده‌های آموزشی است که به بهبود عملکرد در مواجهه با داده‌های دیده نشده کمک نخواهد کرد.

شکل ۳-۳ نتایج مشابه را در حالتی که از حذف تصادفی استفاده شده است، نشان می‌دهد. در این حالت، خطاهای آموزشی به‌طور قابل‌توجهی بزرگ‌تر از شکل ۲-۳ هستند، و خطاهای اعتبارسنجی دارای بزرگی مشابهی هستند. با این حال، دو ویژگی کلیدی در حالت حذف تصادفی عبارت‌اند از: (الف) خطای اعتبارسنجی کمتر از خطای آموزشی است و (ب) خطای اعتبارسنجی به همراه خطای آموزشی به کاهش ادامه می‌دهد. این موارد نشان می‌دهند که فرایند بهینه‌سازی در حال استخراج اطلاعات مفید در طول تمام دوره‌های آموزش است.

شکل ۴-۳ عملکرد شبکه آموزش دیده بدون حذف تصادفی (پس از ۴۵ دوره) را به‌صورت یک ماتریس سردرگمی خلاصه می‌کند. مقدار صحیح در ورودی عمومی (i, j) تعداد مواردی را نشان می‌دهد که شبکه دسته i را برای تصویری از دسته j پیش‌بینی کرده است. عناصر خارج از قطر اصلی نشان‌دهنده خطاهای طبقه‌بندی هستند.



شکل ۳-۳: آموزش شبکه با حذف تصادفی

برای مثال، عنصر (1, 1) برابر با ۸۱۴ تعداد تصاویر هواپیما را که به درستی به عنوان هواپیما طبقه‌بندی شده‌اند، نشان می‌دهد، و عنصر (1, 2) برابر با ۲۱ تعداد تصاویر خودرو را که به اشتباه به عنوان هواپیما طبقه‌بندی شده‌اند، نشان می‌دهد. در زیر هر عدد صحیح درصد مربوطه، به یک رقم اعشار گرد شده، داده شده است، با توجه به اینکه داده‌های آزمایشی از هر دسته شامل ۱۰۰۰ تصویر است. سطر اضافی با برجسب «همه» ورودی‌های هر ستون را خلاصه می‌کند. به عنوان مثال، مقدار ۸۱/۴٪ در ستون اول سطر نهایی به این دلیل است که ۸۱۴ از ۱۰۰۰ تصویر هواپیما به درستی طبقه‌بندی شده‌اند. در زیر این، مقدار ۱۸/۶٪ به این دلیل است که ۱۸۶ از این تصاویر هواپیما به اشتباه طبقه‌بندی شده‌اند. ستون نهایی ماتریس، با برجسب «همه»، هر سطر را خلاصه می‌کند. به عنوان مثال، مقدار ۸۲/۴٪ در ستون نهایی سطر اول به این دلیل است که ۹۸۸ تصویر توسط شبکه به عنوان هواپیما طبقه‌بندی شده‌اند که از این تعداد ۸۱۴ طبقه‌بندی صحیح بوده‌اند. در زیر این، مقدار ۱۷/۶٪ به این دلیل است که باقی‌مانده ۱۷۴ از این ۹۸۸ طبقه‌بندی هواپیما نادرست بوده‌اند. در نهایت، ورودی‌های گوشه پایین سمت راست همه دسته‌ها را خلاصه می‌کنند. می‌بینیم که ۸۰/۱٪ از تمام تصاویر به درستی طبقه‌بندی شده‌اند (و بنابراین ۱۹/۹٪ به اشتباه طبقه‌بندی شده‌اند).

شکل ۳-۵ نتایج مربوط به حالتی که در آن از حذف تصادفی استفاده شده است را نشان می‌دهد. مشاهده می‌کنیم که استفاده از حذف تصادفی به طور کلی عملکرد را بهبود بخشیده است و به ویژه نرخ موفقیت کلی را از ۸۰/۱٪ به ۸۱/۱٪ افزایش داده است. حذف تصادفی در نه دسته از ده دسته مقادیر بیشتری در عناصر قطری ماتریس سردرگمی به وجود آورده است.

Confusion Matrix

Output Class	airplane	auto	bird	cat	deer	dog	frog	horse	ship	truck	all
airplane	814 8.1%	21 0.2%	39 0.4%	14 0.1%	15 0.1%	9 0.1%	5 0.1%	5 0.1%	42 0.4%	24 0.2%	82.4%
auto	13 0.1%	869 8.7%	2 0.0%	3 0.0%	1 0.0%	2 0.0%	5 0.1%	2 0.0%	23 0.2%	60 0.6%	88.7%
bird	45 0.4%	7 0.1%	747 7.5%	44 0.4%	43 0.4%	39 0.4%	40 0.4%	20 0.2%	7 0.1%	4 0.0%	75.0%
cat	17 0.2%	4 0.0%	49 0.5%	618 6.2%	44 0.4%	139 1.4%	44 0.4%	37 0.4%	8 0.1%	6 0.1%	64.0%
deer	12 0.1%	2 0.0%	48 0.5%	57 0.6%	795 8.0%	39 0.4%	20 0.2%	42 0.4%	6 0.1%	0 0.0%	77.9%
dog	6 0.1%	2 0.0%	43 0.4%	146 1.5%	31 0.3%	706 7.1%	19 0.2%	39 0.4%	5 0.1%	1 0.0%	70.7%
frog	7 0.1%	7 0.1%	41 0.4%	63 0.6%	28 0.3%	21 0.2%	855 8.6%	2 0.0%	8 0.1%	3 0.0%	82.6%
horse	10 0.1%	5 0.1%	23 0.2%	31 0.3%	37 0.4%	37 0.4%	5 0.1%	844 8.4%	2 0.0%	7 0.1%	84.3%
ship	44 0.4%	25 0.3%	4 0.0%	11 0.1%	3 0.0%	2 0.0%	6 0.1%	3 0.0%	883 8.8%	13 0.1%	88.8%
truck	32 0.3%	58 0.6%	4 0.0%	13 0.1%	3 0.0%	6 0.1%	1 0.0%	6 0.1%	16 0.2%	882 8.8%	86.4%
all	81.4%	86.9%	74.7%	61.8%	79.5%	70.6%	85.5%	84.4%	88.3%	88.2%	80.1%
	18.6%	13.1%	25.3%	38.2%	20.5%	29.4%	14.5%	15.6%	11.7%	11.8%	19.9%
	airplane	auto	bird	cat	deer	dog	frog	horse	ship	truck	all

Target Class

شکل ۳-۴: ماتریس در هم‌ریختگی بدون حذف تصادفی

Confusion Matrix

Output Class	airplane	auto	bird	cat	deer	dog	frog	horse	ship	truck	all
airplane	847 8.5%	10 0.1%	51 0.5%	17 0.2%	18 0.2%	10 0.1%	9 0.1%	14 0.1%	41 0.4%	19 0.2%	81.8%
auto	15 0.1%	905 9.0%	1 0.0%	4 0.0%	1 0.0%	2 0.0%	3 0.0%	2 0.0%	9 0.1%	56 0.6%	90.7%
bird	27 0.3%	2 0.0%	714 7.1%	52 0.5%	33 0.3%	40 0.4%	27 0.3%	28 0.3%	2 0.0%	2 0.0%	77.0%
cat	14 0.1%	1 0.0%	35 0.4%	618 6.2%	36 0.4%	154 1.5%	28 0.3%	22 0.2%	6 0.1%	11 0.1%	66.8%
deer	13 0.1%	1 0.0%	66 0.7%	56 0.6%	823 8.2%	49 0.5%	21 0.2%	58 0.6%	2 0.0%	3 0.0%	75.4%
dog	1 0.0%	1 0.0%	45 0.4%	135 1.4%	10 0.1%	684 6.8%	9 0.1%	29 0.3%	0 0.0%	0 0.0%	74.8%
frog	5 0.1%	6 0.1%	53 0.5%	68 0.7%	39 0.4%	23 0.2%	893 8.9%	9 0.1%	8 0.1%	5 0.1%	80.5%
horse	9 0.1%	0 0.0%	24 0.2%	21 0.2%	33 0.3%	32 0.3%	2 0.0%	828 8.3%	2 0.0%	5 0.1%	86.6%
ship	41 0.4%	19 0.2%	3 0.0%	14 0.1%	2 0.0%	2 0.0%	6 0.1%	3 0.0%	915 9.2%	19 0.2%	89.4%
truck	28 0.3%	55 0.5%	8 0.1%	15 0.1%	5 0.1%	4 0.0%	2 0.0%	7 0.1%	15 0.1%	880 8.8%	86.4%
all	84.7%	90.5%	71.4%	61.8%	82.3%	68.4%	89.3%	82.8%	91.5%	88.0%	81.1%
	15.3%	9.5%	28.6%	38.2%	17.7%	31.6%	10.7%	17.2%	8.5%	12.0%	18.9%
	airplane	auto	bird	cat	deer	dog	frog	horse	ship	truck	all

Target Class

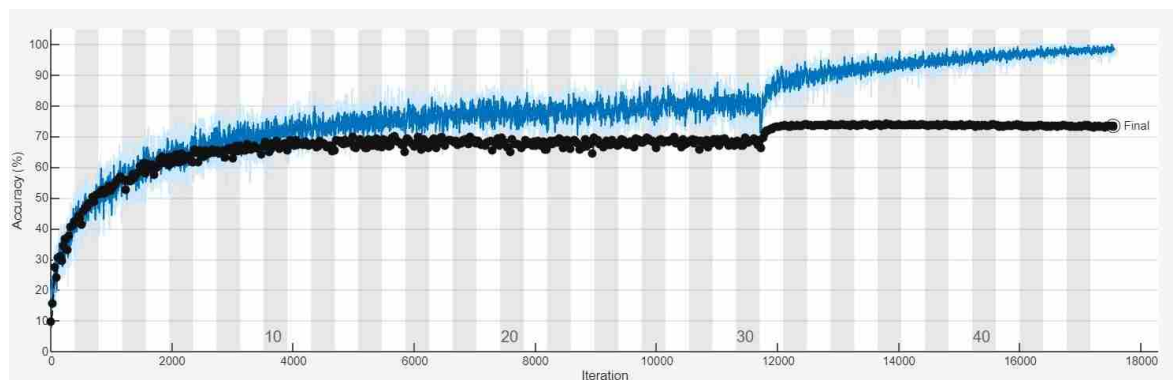
شکل ۳-۵: ماتریس در هم‌ریختگی با حذف تصادفی

۵-۳ شبیه‌سازی و ارزیابی مدل‌های شبکه عصبی

در ادامه به شبیه‌سازی شبکه مورد نظر با استفاده از کدهای <<Listing3>> پرداخته و نتایج حاصل از آن را ارائه می‌دهیم. نمودار ۳-۶ دقت^۱ کد شبیه‌سازی شده را نشان می‌دهد؛ خط آبی نشان‌دهنده دقت مدل در طول آموزش

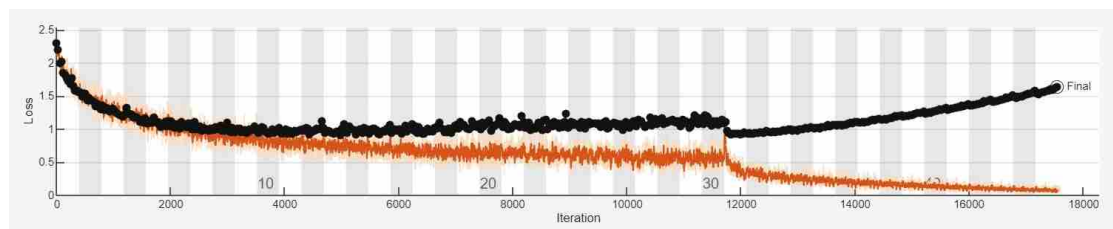
¹Accuracy

است و خط مشکی نشان‌دهنده دقت مدل در مجموعه داده‌های اعتبارسنجی است. دقت مدل در ابتدا افزایشی بوده و در نهایت به دقتی حدود ۷۳ درصد می‌رسد همچنین دقت اعتبارسنجی نیز به تریج روند صعودی داشته و در نهایت به مقدار ثابتی نزدیک شده است. نمودار ۳-۷ میزان خطای شبکه^۱ را نشان می‌دهد. خط نارنجی



شکل ۳-۶: دقت آموزش شبکه شبیه سازی شده بدون حذف تصادفی

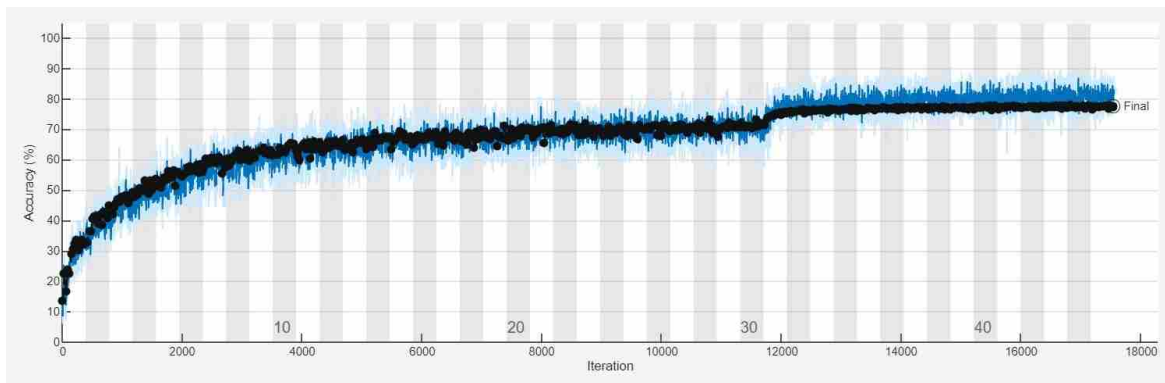
نشان‌دهنده خطای مدل در طول آموزش و خط مشکی نشان‌دهنده خطای مدل در مجموعه داده‌های اعتباری است. خطای مدل به تدریج کاهش یافته و به سطح پایداری نزدیک می‌شود؛ اما در مدل اعتبارسنجی به طور ناگهانی افزایش می‌یابد و می‌تواند نشانه‌ای از بیش برآزش باشد. به همین ترتیب، دو نمودار ۳-۸، ۳-۹ نشان‌دهنده



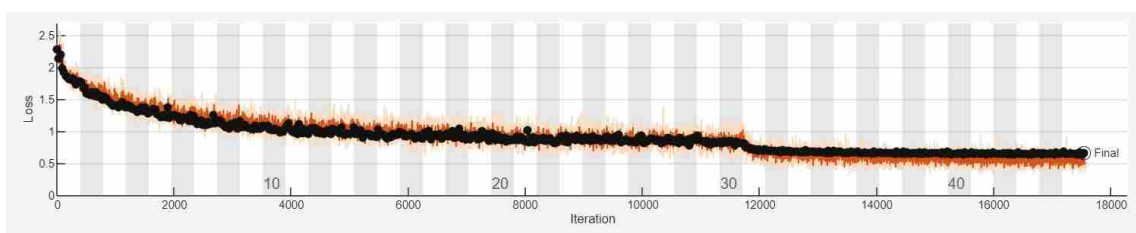
شکل ۳-۷: خطای شبکه شبیه سازی شده با حذف تصادفی

دقت و خطای مدل با اعمال حذف تصادفی در شبکه هستند. در این حالت، دقت اعتبارسنجی به بیش از ۷۷ درصد رسیده است که نشان‌دهنده بهبود عملکرد مدل در مقایسه با مدل اولیه است. دو نمودار ۳-۱۰، ۳-۱۱ ماتریس درهم‌ریختگی در دو مدل ذکر شده را نشان می‌دهد.

¹Loss



شکل ۳-۸: دقت آموزش شبکه شبیه سازی شده با حذف تصادفی



شکل ۳-۹: خطای شبکه شبیه سازی شده با حذف تصادفی

Confusion Matrix - Without Dropout

0	774	9	49	28	13	7	14	12	57	37	77.4%	22.6%
1	14	852	8	6	2	5	11		28	74	85.2%	14.8%
2	55	7	621	82	55	66	67	28	7	12	62.1%	37.9%
3	22	4	55	543	46	178	75	35	20	22	54.3%	45.7%
4	19		67	70	695	44	43	52	7	3	69.5%	30.5%
5	11	5	37	164	39	655	31	41	7	10	65.5%	34.5%
6	12	4	42	62	32	24	803	3	11	7	80.3%	19.7%
7	18	7	30	49	46	64	13	752	5	16	75.2%	24.8%
8	62	20	9	18	4	7	6	3	844	27	84.4%	15.6%
9	32	55	6	22	6	12	6	7	38	816	81.6%	18.4%

76.0%	88.5%	67.2%	52.0%	74.1%	61.7%	75.1%	80.6%	82.4%	79.7%
24.0%	11.5%	32.8%	48.0%	25.9%	38.3%	24.9%	19.4%	17.6%	20.3%
0	1	2	3	4	5	6	7	8	9

Predicted Class

شکل ۳-۱۰: ماتریس در هم‌ریختگی مدل شبیه سازی شده بدون حذف تصادفی

Confusion Matrix - With Dropout

True Class \ Predicted Class	0	1	2	3	4	5	6	7	8	9	Accuracy	Loss
0	815	15	47	26	18	1	11	11	35	21	81.5%	18.5%
1	9	889	4	3	2	4	7	2	25	55	88.9%	11.1%
2	60	3	650	52	91	47	75	14	7	1	65.0%	35.0%
3	16	5	55	592	78	136	85	16	8	9	59.2%	40.8%
4	10	2	57	35	792	18	56	22	7	1	79.2%	20.8%
5	7		45	166	53	655	39	31	4		65.5%	34.5%
6	8	2	28	35	22	11	889	3	2		88.9%	11.1%
7	15	1	33	49	86	41	6	763	2	4	76.3%	23.7%
8	52	9	13	12	8	3	8	4	879	12	87.9%	12.1%
9	36	61	5	9	6	3	13	11	19	837	83.7%	16.3%

79.3%	90.1%	69.4%	60.5%	68.6%	71.3%	74.8%	87.0%	89.0%	89.0%
20.7%	9.9%	30.6%	39.5%	31.5%	28.7%	25.2%	13.0%	11.0%	11.0%

شکل ۳-۱۱: ماتریس در هم‌ریختگی مدل شبیه سازی شده با حذف تصادفی

۳-۶ جمع‌بندی

در این گزارش، به بررسی مفاهیم پایه‌ای و تکنیک‌های پیشرفته در حوزه شبکه‌های عصبی مصنوعی و یادگیری عمیق پرداخته شد. روش‌های مختلف بهینه‌سازی و ساختارهای متنوع شبکه‌های عصبی از جمله شبکه‌های چندلایه (MLP) و شبکه‌های عصبی پیچشی (CNN) به تفصیل مورد بحث قرار گرفتند. همچنین، تکنیک‌های بهبود عملکرد مانند پس انتشار خطا، گرادیان نزولی و حذف تصادفی که به کاهش بیش برآزش و افزایش دقت مدل‌ها کمک می‌کنند، معرفی شدند.

در بخش‌های مختلف گزارش، کاربردهای عملی این تکنیک‌ها در مسائل پیچیده‌ای نظیر پردازش تصویر و یادگیری ویژگی‌های غیرخطی مورد بررسی قرار گرفت. استفاده از ابزارهایی مانند MatConvNet نشان داد که می‌توان با بهره‌گیری از کتابخانه‌های مناسب، فرایند آموزش مدل‌ها را تسریع و بهینه کرد.

نتایج نشان می‌دهد که استفاده هوشمندانه از تکنیک‌های بهینه‌سازی و تنظیمات مختلف شبکه‌های عصبی می‌تواند تأثیر چشمگیری در افزایش دقت و کارایی مدل‌ها داشته باشد. با این حال، برای دستیابی به بهترین عملکرد، نیازمند تنظیم دقیق پارامترها و روش‌های بهینه‌سازی هستیم که این امر بستگی به نوع داده‌ها و مسئله مورد نظر دارد.

به‌طور کلی، این گزارش نشان می‌دهد که یادگیری عمیق و شبکه‌های عصبی با وجود چالش‌های موجود، ابزارهای قدرتمندی برای حل مسائل پیچیده در حوزه‌های مختلف هستند و با پیشرفت‌های مداوم در الگوریتم‌ها و سخت‌افزارها، انتظار می‌رود که کاربردهای آن‌ها به طور قابل توجهی گسترش یابد.

```

1 % Load CIFAR-10 dataset
2 cifar10Data = fullfile(tempdir, 'cifar10');
3 cifar10URL = 'https://www.cs.toronto.edu/~kriz/cifar-10-matlab.tar.
   gz';
4
5 % Download and extract CIFAR-10 dataset
6 if ~exist(cifar10Data, 'dir')
7 disp('Downloading CIFAR-10 dataset...');
8 downloadFolder = tempdir;
9 cifar10TarFile = fullfile(downloadFolder, 'cifar-10-matlab.tar.gz')
   ;
10 websave(cifar10TarFile, cifar10URL);
11 untar(cifar10TarFile, cifar10Data);
12 end
13
14 % Load CIFAR-10 data
15 cifar10Data = fullfile(cifar10Data, 'cifar-10-batches-mat');
16 [trainImages, trainLabels, valImages, valLabels] = loadCIFAR10(
   cifar10Data);
17
18 % Define neural network layers without Dropout
19 layers_no_dropout = [
20 imageInputLayer([32 32 3], 'Name', 'input')
21 convolution2dLayer(5, 32, 'Padding', 2, 'Stride', 1, 'Name', 'conv1
   ')
22 batchNormalizationLayer('Name', 'bn1')
23 reluLayer('Name', 'relu1')
24 maxPooling2dLayer(3, 'Stride', 2, 'Padding', [0 1 0 1], 'Name', '
   pool1')
25 convolution2dLayer(5, 32, 'Padding', 2, 'Stride', 1, 'Name', 'conv2
   ')
26 reluLayer('Name', 'relu2')
27 averagePooling2dLayer(3, 'Stride', 2, 'Padding', [0 1 0 1], 'Name',
   'pool2')
28 convolution2dLayer(5, 64, 'Padding', 2, 'Stride', 1, 'Name', 'conv3
   ')
29 reluLayer('Name', 'relu3')
30 averagePooling2dLayer(3, 'Stride', 2, 'Padding', [0 1 0 1], 'Name',
   'pool3')
31 convolution2dLayer(4, 64, 'Padding', 0, 'Stride', 1, 'Name', 'conv4
   ')
32 reluLayer('Name', 'relu4')
33 convolution2dLayer(1, 10, 'Padding', 0, 'Stride', 1, 'Name', 'conv5
   ')
34 softmaxLayer('Name', 'softmax')
35 classificationLayer('Name', 'classification')];
36
37 % Set training options
38 options = trainingOptions('sgdm', ...
39 'InitialLearnRate', 0.1, ...
40 'LearnRateSchedule', 'piecewise', ...
41 'LearnRateDropPeriod', 30, ...
42 'LearnRateDropFactor', 0.1, ...
43 'MaxEpochs', 45, ...
44 'ValidationData', {valImages, valLabels}, ...
45 'ValidationFrequency', 30, ...
46 'Verbose', false, ...

```

```

47     'Plots', 'training-progress');
48
49     % Train network without Dropout
50     net_no_dropout = trainNetwork(trainImages, trainLabels,
51                                 layers_no_dropout, options);
52
53     % Function to load CIFAR-10 data
54     function [trainImages, trainLabels, valImages, valLabels] =
55             loadCIFAR10(cifar10Data)
56     [trainImages, trainLabels] = loadBatch(fullfile(cifar10Data, '
57             data_batch_1.mat'));
58     for i = 2:5
59         [images, labels] = loadBatch(fullfile(cifar10Data, ['data_batch_'
60             num2str(i) '.mat']));
61     trainImages = cat(4, trainImages, images);
62     trainLabels = [trainLabels; labels];
63     end
64     [valImages, valLabels] = loadBatch(fullfile(cifar10Data, '
65             test_batch.mat'));
66     end
67
68     function [images, labels] = loadBatch(filename)
69     load(filename, 'data', 'labels');
70     images = reshape(data, 32, 32, 3, []);
71     images = permute(images, [2 1 3 4]);
72     labels = categorical(labels);
73     end

```

Listing 3: MATLAB code for loading and training neural networks

مراجع

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <https://www.deeplearningbook.org>
- [2] MathWorks. *8x8 RGB LED Matrix*. 2024. https://de.mathworks.com/help/simulink/supportpkg/raspberrypi_ref/8x8rgbledmatrix.html Accessed: 2024-09-04
- [3] MathWorks. *Introduction to Convolutional Neural Networks*. 2024. <https://de.mathworks.com/help/deeplearning/ug/introduction-to-convolutional-neural-networks.html> Accessed: 2024-09-04
- [4] MathWorks. *MATLAB Code for CNN on CIFAR-10*. 2024. https://viewer.mathworks.com/?viewer=plain_code&url=https%3A%2F%2Ffr.mathworks.com%2Fmatlabcentral%2Fmlc-downloads%2Fdownloads%2Fe5aef38d-4a80-11e4-9553-005056977bd0%2Fc36e6b81-a2f6-b65f-4ecb-8fa3c6764b50%2Ffiles%2Fexamples%2Fcifar%2Fcnn_cifar_init.m&embed=web Accessed: 2024-09-04
- [5] Catherine F. Higham and Desmond J. Higham. *Deep Learning: An Introduction for Applied Mathematicians*. SIAM Review, 61(4):860-891, 2019.