



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

یادگیری ماشین

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین: ایزدی

مهلت ارسال نهایی: ۲۳ آذر

تمرین چهارم

مهلت ارسال امتیازی: ۱۶ آذر

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روزهای مشخص شده است.
- در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر مجاز داشته باشید و در مجموع حداکثر ۱۵ روز تأخیر مجاز خواهید داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت جداگانه محاسبه می‌شود و مجموع تأخیر هر دو نباید بیشتر از ۱۵ روز شود. پس از اتمام زمان مجاز، دو روز اضافی برای آپلود غیرمجاز در نظر گرفته شده است که در این بازه، به ازای هر ساعت تأخیر، ۲ درصد از نمره نهایی تمرین کسر خواهد شد.
- اگر بخش عملی یا تئوری تمرین را قبل از مهلت ارسال امتیازی آپلود کنید، ۲۰ درصد نمره اضافی به آن بخش تعلق خواهد گرفت و پس از آن، ویدیویی تحت عنوان راهنمایی برای حل تمرین منتشر خواهد شد.
- حتماً تمرین‌ها را بر اساس موارد ذکر شده درک شده در صورت سوالات حل کنید. در صورت وجود هرگونه ابهام، آن را در صفحه تمرین در سایت کوئرا مطرح کنید و به پاسخ‌هایی که از سوی دستیار آموزشی مربوطه ارائه می‌شود، توجه کنید.
- در صورت هم‌فکری و یا استفاده از منابع، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال را ذکر کنید.
- فایل پاسخ‌های سوالات نظری را در قالب یک فایل pdf به فرمت HW4_T_[STD_ID].pdf آماده کنید و برای سوالات عملی، هر یک را در یک فایل zip جداگانه قرار دهید. فایل مربوط به نوتبوک i ام را به فرمت HW4_P[i]_[STD_ID].zip نام‌گذاری کرده و هرکدام را به صورت جداگانه آپلود کنید

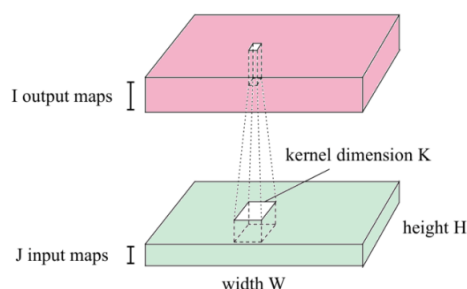
گردآورندگان تمرین: امیرحسین ایزدی، محمدحسین شالچیان، امیرعزتی، علی بختیاری

سوالات نظری (۱۰۰+۱۰ نمره)

۱. (۲۲ نمره) به سوالات زیر پاسخ دهید.

الف) دو لایه متوالی از شبکه ژرفی را در نظر بگیرید به طوری که لایه ورودی دارای I و لایه خروجی دارای I ویژگی نگاشت باشد. اتصال میان این دو لایه می‌تواند یا به صورت اتصال تمام‌متصل یا به صورت پیچشی باشد. برای هر دو حالت، ویژگی نگاشت لایه ورودی دارای ابعاد $W \times H$ است و کرنل استفاده شده در لایه پیچشی دارای ابعاد $K \times K$ است. نمای کلی این دو لایه در شکل ۲ قابل مشاهده هست.

- تعداد واحدهای خروجی، اتصالات و وزن‌های قابل یادگیری را در صورت امکان محاسبه؛ گزارش و با یکدیگر مقایسه کنید. (از بایاس صرف نظر کنید)



شکل ۱: لایه ورودی با رنگ سبز و خروجی با رنگ قرمز مشخص شده است.

پاسخ:

	FC layer	CONV layer
#units output	WHI	WHI
#weights	W^*H^*IJ	K^*IJ
#connections	W^*H^*IJ	WHK^*IJ

(آ) از محاسبات قبلی می‌توان نتیجه گرفت که شبکه‌های عمیق کاملاً پیچشی نسبت به شبکه‌های کاملاً متصل عموماً نیازمند داده آموزشی کمتری برای آموزش هستند؟ آیا این نتیجه‌گیری لزوماً صحیح است و مرتبط است یا خیر؟

پاسخ:

- به طور کلی، شبکه‌های عمیق پیچشی نسبت به شبکه‌های کاملاً متصل نیازمند داده آموزشی کمتری برای آموزش هستند. دلیل این امر آن است که شبکه‌های پیچشی از اشتراک وزن‌ها و کاهش تعداد پارامترهای قابل یادگیری استفاده می‌کنند. این ویژگی‌ها باعث می‌شود که شبکه‌های پیچشی بهره‌وری بالاتری در استفاده از داده آموزشی داشته باشند و نیاز به داده کمتر برای دستیابی به همان سطح از دقت داشته باشند.
- با این حال، این نتیجه‌گیری لزوماً همیشه صحیح نیست و به عوامل متعددی بستگی دارد. به عنوان مثال، اگر داده ورودی پیچیدگی بالایی داشته باشد و یا ارتباطات غیرمحلی مهم باشند، شبکه‌های کاملاً متصل ممکن است عملکرد بهتری داشته باشند. علاوه بر این، معماری شبکه، اندازه داده آموزشی، و نحوه پردازش داده‌ها نیز می‌توانند تأثیر قابل توجهی بر نیاز به داده آموزشی داشته باشند.

ب) شبکه CNN ای را در نظر بگیرید که از بلاک‌هایی به فرم زیر استفاده می‌کند:

(ConvLayer) \rightarrow (BatchNorm) \rightarrow (Activation)

- با مطالعه این [مقاله](#) نحوه انجام نرمال سازی بچ در شبکه های تماماً متصل و شبکه های پیچشی را با یکدیگر مقایسه کنید

پاسخ:

شبکه‌های تمام‌متصل: در شبکه‌های تمام‌متصل، ورودی‌ها به صورت بردارهای یک‌بعدی و مسطح شده ارائه می‌شوند. در این ساختار، نرمال‌سازی دسته‌ای (Batch Normalization) به طور مستقل برای هر ویژگی اعمال می‌شود، به طوری که میانگین مقادیر ورودی به هر نورون صفر و واریانس آن‌ها یک می‌شود. این امر باعث می‌شود که تغییرات داده‌های ورودی در طول آموزش کاهش یافته، و فرآیند یادگیری تسریع شود.

همچنین، به دلیل اینکه وابستگی آماری بین ویژگی‌ها در شبکه‌های تمام‌متصل کمتر است، نرمال‌سازی به صورت مستقل انجام می‌شود. این عملیات موجب پایداری گرادینت‌ها در طول یادگیری شده و به مدل کمک می‌کند تا با نرخ یادگیری بالاتر و پارامترهای اولیه کمتر حساس، به سرعت بهینه شود.

شبکه‌های پیچشی: در شبکه‌های پیچشی، ورودی‌ها معمولاً داده‌های دوبعدی یا تصاویر هستند که دارای ویژگی‌های محلی مهم می‌باشند. نرمال‌سازی دسته‌ای در این شبکه‌ها به گونه‌ای طراحی شده

که تمامی مقادیر یک نقشه ویژگی (Feature Map) به صورت مشابه نرمال سازی شوند. این کار برای حفظ همبستگی مکانی داده ها ضروری است. به همین دلیل، نرمال سازی برای هر کانال به صورت مستقل و برای تمامی مکان های فضایی به طور مشترک اعمال می شود. نرمال سازی دسته ای در شبکه های پیچشی، پایداری گرادین ها را افزایش داده و یادگیری را کارآمدتر می کند. این ویژگی به خصوص در شبکه های عمیق تر اهمیت بیشتری دارد زیرا وابستگی محلی داده ها حفظ می شود و اختلالی در این ارتباطات ایجاد نمی گردد.

- آیا حذف بایاس b از لایه کانولوشن در کارکرد این شبکه اختلالی ایجاد می کند؟ چرا؟

پاسخ:

- حذف بایاس در لایه هایی که از نرمال سازی دسته ای استفاده می کنند معمولاً اختلالی در شبکه ایجاد نمی کند. دلیل این امر این است که لایه نرمال سازی دسته ای شامل پارامترهای قابل یادگیری (مقیاس γ و شیفت β) است که می توانند اثر نبود بایاس را جبران کنند.

- **دلیل:** مرحله تبدیل خطی (Affine Transformation) در نرمال سازی دسته ای (پس از نرمال سازی) انعطاف پذیری لازم را برای جبران تغییرات فراهم می کند. بنابراین، وجود بایاس ضروری نیست.

- همچنین فرض کنید شبکه را آموزش داده ایم؛ آیا ضرب کردن وزن ها در یک عدد مانند α در زمان آزمایش (Inference)، عملکرد شبکه را تغییر می دهد؟ ضرب کردن این ضریب در تمام درایه های ورودی شبکه چگونه؟

پاسخ:

تأثیر ضرب کردن وزن ها در حین inference:
اگر وزن ها را در α ضرب کنیم:

$$x_{in} \Rightarrow \text{Conv Layer} \Rightarrow \alpha x_{out} \Rightarrow \text{Batch Norm} \Rightarrow \alpha \gamma \left(\frac{\alpha x_{out} - \alpha \mu}{|\alpha| \sigma} \right) + \alpha \beta$$

$$= \alpha \left(\gamma \text{sign}(\alpha) \left(\frac{x_{out} - \mu}{\sigma} \right) + \beta \right) = \alpha y$$

حال اگر تابع فعال سازی غیرخطی باشد در این صورت مقدار $f(\alpha y)$ رابطه مشخصی با $f(y)$ ندارد، بنابراین عملکرد مدل تحت تأثیر قرار می گیرد و رفتار مشابه قبل نخواهد داشت.

تأثیر ضرب کردن ورودی ها در حین inference:
داریم:

$$\alpha x_{in} \Rightarrow \text{Conv Layer} \Rightarrow \alpha x_{out} \Rightarrow \text{Batch Norm} \Rightarrow \gamma \left(\frac{\alpha x_{out} - \alpha \mu}{|\alpha| \sigma} \right) + \beta$$

$$= \text{sign}(\alpha) \left(\gamma \left(\frac{x_{out} - \mu}{\sigma} \right) + \beta \right)$$

می توان دید که در صورت مثبت بودن α تغییری در عملکرد مدل ایجاد نمی شود.

- یک سناریو مرتبط به حوزه پردازش تصویر را شرح دهید که در آن نرمال سازی دسته ای یا بیچ ممکن است در آن کارایی کمتری داشته باشد و یا حتی نتیجه در معکوس دهد. در چنین مواردی چه تکنیک های جایگزینی می توانند مورد توجه قرار گیرند؟

پاسخ:

- **سناریو:** در حوزه‌هایی مانند تقویت تصاویر گم‌نور، یا تحلیل تصاویر پزشکی، که شدت پیکسل‌های ورودی توزیع‌های نامتعادلی دارند، نرمال‌سازی دستگاه ممکن است باعث نتایج ضعیف شود.

* به عنوان مثال، اگر مقادیر پیکسل‌ها توزیع به شدت نامتعالل داشته باشند، نرمال‌سازی ممکن است جزئیات ظریفی که برای تصمیم‌گیری مهم هستند را از بین ببرد.

استراتژی‌ها برای رفع این مشکلات:

- استفاده از نرمال‌سازی نمونه‌ای (Instance Normalization) به جای نرمال‌سازی دستگاهی در وظایف خاص (مانند تحلیل سبک).

- استفاده از نرمال‌سازی گروهی (Group Normalization) یا نرمال‌سازی لایه‌ای (Layer Normalization) که وابسته به آمار دسته نیستند.

* تنظیم پارامترهای اولیه γ و β برای انطباق بهتر با داده‌های خاص دامنه.

* آموزش مجدد شبکه روی داده‌هایی که توزیع مشابهی با مجموعه آزمایش دارند، برای تطبیق مناسب لایه‌های نرمال‌سازی.

ج) معماری Unet را در نظر بگیرید:

- تصور کنید که ابعاد تصویر ورودی ما برای این شبکه 256×256 می‌باشد. حال فرض می‌شود که در این معماری هر لایه در آنکدر ابعاد را به نصف کاهش می‌دهد و در دیکدر دو برابر می‌کند. در پایین‌ترین لایه (عمیق‌ترین لایه) این معماری، فضای ویژگی ما چند پیکسل خواهد داشت؟

پاسخ:

ابعاد تصویر ورودی 256×256 است. در هر لایه از آنکدر، ابعاد به نصف کاهش می‌یابد. بنابراین، ابعاد در هر لایه به صورت زیر خواهد بود:

- لایه اول: 256×256 .

- لایه دوم: 128×128 .

- لایه سوم: 64×64 .

- لایه چهارم: 32×32 .

- لایه پنجم (عمیق‌ترین لایه): 16×16 .

بنابراین در پایین‌ترین لایه (عمیق‌ترین لایه)، فضای ویژگی ما ابعاد 16×16 خواهد داشت.

- فرض کنید آنکودر دارای لایه‌هایی با $256, 128, 64$ و 512 فیلتر است. اگر هر لایه کانولوشن از کرنل‌های 3×3 استفاده کند، تعداد پارامترهای لایه کانولوشن دوم آنکدر را محاسبه کنید.

پاسخ:

تعداد فیلترهای آنکدر به صورت $512, 512, 256, 128, 64$ است. برای لایه کانولوشن دوم با 3×3 کرنل‌ها: **فرمول کلی:**

$$\text{تعداد پارامترها} = (3 \times 3) \times \text{ورودی} \times \text{خروجی} + \text{بایاس}$$

برای لایه کانولوشن دوم:

$$128 \times 64 \times (3 \times 3) + 128 = 73,856$$

بنابراین تعداد پارامترهای این لایه ۷۳,۸۵۶ است.

د) تفسیرپذیری شبکه‌های عصبی از اهمیت بالایی برخوردار است. هنگام دسته‌بندی تصاویر، علاقه‌مندیم بدانیم کدام بخش‌های تصویر در دسته‌بندی، تأثیر بیشتری داشته‌اند. در **مقاله** ایده‌ی شبکه‌های de-convolutional و در **مقاله** ایده‌ی شبکه‌های up-convolutional مطرح شده است. با بررسی این مقالات، توضیح دهید هر کدام از دو روش به چه صورت منجر به تفسیرپذیری شبکه کانولوشنی می‌شوند.

پاسخ:

شبکه‌های Up-Convolutional: این شبکه‌ها که گاهی به نام Transposed Convolution شناخته می‌شوند، ابزاری برای افزایش ابعاد فضایی نقشه‌های ویژگی هستند. در مقاله‌ی "Inverting Visual Representations"، این شبکه‌ها برای بازسازی تصاویر از نمایش‌های ویژگی (SIFT، HOG) و خروجی لایه‌های مختلف AlexNet) به کار گرفته شدند. در این روش، یک شبکه Up-Convolutional با یادگیری معکوس نگاشت ویژگی‌ها، تصویر اولیه را تخمین می‌زند. این بازسازی‌ها نشان می‌دهند که حتی در لایه‌های عمیق شبکه:

- رنگ‌ها و کانتورهای کلی اشیاء حفظ می‌شوند.
- اطلاعات محلی دقیق‌تر با پیشرفت به لایه‌های بالاتر شبکه کاهش می‌یابد.
- ویژگی‌های سطح بالا مانند دسته‌بندی‌های مفهومی (مانند "Tree" یا "Apple") به صورت کلی حفظ می‌شوند.

این بازسازی‌ها کمک می‌کنند تا درک شود که چگونه شبکه‌های پیچشی اطلاعات را حفظ یا حذف می‌کنند. برای مثال، بازسازی از لایه‌های بالاتر مانند FC8 نشان می‌دهد که حتی اطلاعاتی مانند رنگ‌ها و مقادیر پیشین کلاس‌ها نیز قابل استخراج هستند.

شبکه‌های De-Convolutional: در مقاله‌ی "Striving for Simplicity: The All Convolutional Network"، از شبکه‌های De-Convolutional به عنوان ابزاری برای تحلیل ویژگی‌های یادگرفته‌شده استفاده شده است. این شبکه‌ها با معکوس کردن جریان داده، نواحی تصویر که به‌طور خاص باعث فعال‌سازی نورون‌های خاص در لایه‌های بالاتر می‌شوند را مشخص می‌کنند. فرآیند این شبکه‌ها شامل مراحل زیر است:

- i. شروع به فعال‌سازی خاص در لایه‌ی بالا و صفر کردن سایر مقادیر.
 - ii. بازگرداندن این فعال‌سازی‌ها به ورودی تصویر، به طوری که نواحی تصویر که بیشترین تأثیر را روی آن نورون دارند، بازسازی شوند.
 - iii. استفاده از روش‌های مختلف برای بازگشت از لایه‌های ReLU و پیچشی مانند Guided Backpropagation که نویز کمی تولید می‌کند و مناطق مؤثرتر را برجسته می‌سازد.
- نتیجه این فرآیند نشان می‌دهد که:

- در لایه‌های پایین‌تر، نواحی فعال‌سازی بسیار محلی و مرتبط با الگوهای ساده (مانند لبه‌ها) هستند.
- در لایه‌های بالاتر، نواحی فعال‌سازی به الگوهای پیچیده‌تر و انتزاعی‌تر (مانند اشیاء کامل یا دسته‌های مفهومی) گسترش می‌یابند.

چگونگی کمک به تفسیرپذیری: هر دو روش به صورت مکمل به تفسیرپذیری شبکه کمک می‌کنند:

- **شبکه‌های Up-Convolutional:** این شبکه‌ها بازسازی تصویر اولیه را از نمایش ویژگی ممکن می‌کنند. بازسازی‌ها می‌توانند فهمید که هر لایه از شبکه چه اطلاعاتی را حفظ می‌کند. مثلاً حفظ رنگ و کانتورها نشان‌دهنده اهمیت این اطلاعات در فرآیند تصمیم‌گیری شبکه است.

- **شبکه‌های De-Convolutional:** این شبکه‌ها با مشخص کردن نواحی حساس به هر نورون، کمک می‌کنند تا نقش هر ویژگی در پیش‌بینی‌های شبکه درک شود. به‌طور خاص، این روش نشان می‌دهد که چه بخش‌هایی از تصویر بیشترین تأثیر در فعال‌سازی لایه‌های بالاتر دارند. به‌طور کلی، شبکه‌های Up-Convolutional برای بازسازی تصویر کلی و نمایش اطلاعات موجود در ویژگی‌ها مناسب هستند، در حالی که شبکه‌های De-Convolutional روی تحلیل مناطق بحرانی تمرکز دارند که به تصمیم‌گیری شبکه کمک می‌کنند.

۲. (۸ نمره) فرض کنید برداری به طول N دارید و قصد دارید یک لایه کانولوشن یک بعدی روی آن اعمال کنید. حاصل اعمال یک لایه کانولوشن را از طریق رابطه‌ی زیر به‌دست آورید:

$$Z = W * X \longrightarrow Z_i = \sum_{j=0}^{K-1} W_j X_{i+j}$$

به دست می‌آوریم که K اندازه‌ی فیلتر را نشان می‌دهد. اگر مقدار $\frac{\partial L}{\partial Z_i}$ برای تمامی مقادیر i بدانیم، رابطه‌ی مربوط به $\frac{\partial L}{\partial W_j}$ را به طور دقیق پیدا کنید. نشان دهید این رابطه، عملاً معادل اعمال یک فیلتر کانولوشن است.

پاسخ:

یک مرحله forward عملیات کانولوشن ۱ بعدی را می‌توان به صورت زیر نوشت:

$$Z_i = \sum_{j=0}^{K-1} W_j X_{i+j}, \quad \text{for } i \in [0, N-K]$$

برای محاسبه گرادیان $\frac{\partial L}{\partial W_j}$ ، می‌توانیم قانون مشتق زنجیری را اعمال کنیم:

$$\frac{\partial L}{\partial W_j} = \sum_{i=1}^N \frac{\partial L}{\partial Z_i} \cdot \frac{\partial Z_i}{\partial W_j}$$

ما می‌دانیم که $\frac{\partial L}{\partial Z_i}$ داده شده است، و می‌توانیم $\frac{\partial Z_i}{\partial W_j}$ را به صورت زیر استخراج کنیم:

$$\frac{\partial Z_i}{\partial W_j} = X_{i+j}$$

با جای‌گذاری این معادله در معادله قانون زنجیری، خواهیم داشت:

$$\frac{\partial L}{\partial W_j} = \sum_{i=1}^N \frac{\partial L}{\partial Z_i} \cdot X_{i+j}$$

این معادله یک عملیات کانولوشن را نشان می‌دهد، جایی که گرادیان $\frac{\partial L}{\partial Z_i}$ با ورودی X_i جابه‌جا شده با موقعیت‌های j ضرب می‌شود. بنابراین، فرمول در واقع نتیجه یک عملیات کانولوشن است. به طور خلاصه:

$$\frac{\partial L}{\partial W_j} = \left(\frac{\partial L}{\partial Z} \right) * X_{\text{by shifted } j}$$

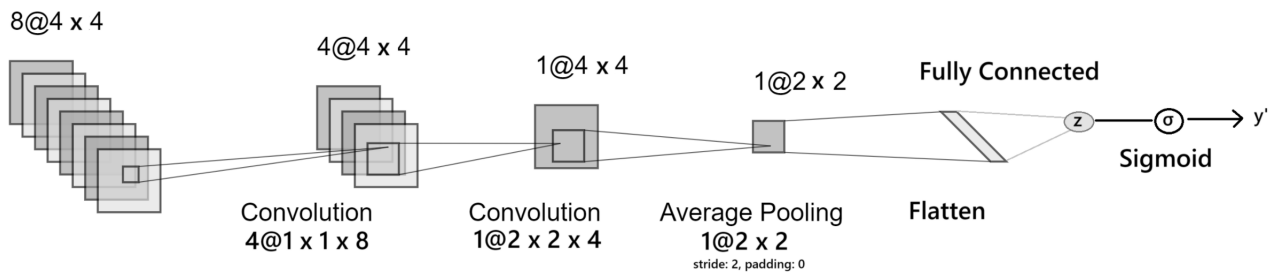
یا به صورت کلی تر:

$$\frac{\partial L}{\partial W} = \left(\frac{\partial L}{\partial Z} \right) * X$$

یعنی کافیت گرایان upstream را یکبار با ورودی لایه کانولوشن، convolve کنیم تا گرایان نسبت به کرنل کانولوشن به دست آید.

۳. (۱۵ نمره)

یکی از اساسی ترین اجزا در یک شبکه ژرف، انتشار به عقب یا Backpropagation است. شبکه داده شده در شکل ۲ را در نظر بگیرید. در این شکل، وزن های لایه اول با $W^{(1)}$ و وزن های لایه دوم با $W^{(2)}$ و وزن های لایه تمام متصل با $W^{(fc)}$ نمایش داده شده اند. همچنین، $P^{(1)}$ ، $P^{(2)}$ و $P^{(3)}$ به ترتیب خروجی های لایه اول، دوم و سوم شبکه هستند.



شکل ۲: معماری شبکه

- ورودی شبکه: یک تانسور با ابعاد $4 \times 4 \times 8$
- لایه اول: لایه کانولوشن با ۴ فیلتر $1 \times 1 \times 8$ که خروجی $P^{(1)}$ با ابعاد $4 \times 4 \times 4$ تولید می کند.
- لایه دوم: لایه کانولوشن با ۱ فیلتر $2 \times 2 \times 4$ که خروجی $P^{(2)}$ با ابعاد $4 \times 4 \times 1$ تولید می کند.
- لایه سوم: لایه Avg. Pooling با کرنل 2×2 و گام ۲، که خروجی $P^{(3)}$ با ابعاد $2 \times 2 \times 1$ ایجاد می کند.
- لایه چهارم: لایه Flatten که $P^{(3)}$ را به یک بردار v با ابعاد 4×1 تبدیل می کند.
- لایه پنجم: لایه تمام متصل با تابع فعال ساز Sigmoid که خروجی نهایی را تولید می کند.

وزن های فیلترهای کانولوشن در لایه های اول و دوم ($W^{(1)}$ و $W^{(2)}$) و وزن های لایه تمام متصل ($W^{(fc)}$) از جمله پارامترهای قابل آموزش شبکه هستند. هدف شبکه، بهینه سازی این وزن ها برای کمینه سازی یک تابع زیان L است.

وظایف:

(آ) با استفاده از قاعده مشتق زنجیره ای و بر اساس مشتق $\frac{\partial L}{\partial z}$ ، عبارت های زیر را محاسبه کنید:

- $\frac{\partial L}{\partial P_{i,j}^{(3)}}$: گرایان زیان نسبت به خروجی لایه Average Pooling.
- $\frac{\partial L}{\partial P_{i,j,k}^{(1)}}$: گرایان زیان نسبت به خروجی لایه اول کانولوشن.

پاسخ:

ابتدا با توجه به معماری شبکه ی عصبی معرفی شده، مقادیر مختلف در لایه های مختلف را به دست آوریم.

لایه اول

$$P_{i,j,t}^{(1)} = \sum_{k=1}^{\Lambda} W_{1,1,k,t}^{(1)} X_{i,j,k}$$

لایه دوم

$$P_{i,j}^{(2)} = \sum_{k=1}^{\mathfrak{F}} \sum_{w=1}^{\mathfrak{Y}} \sum_{h=1}^{\mathfrak{Y}} W_{w,h,k}^{(2)} P_{i+w-1,j+h-1,k}^{(1)}$$

لایه سوم

$$P_{i,j}^{(3)} = \frac{1}{\mathfrak{F}} \sum_{w=\mathfrak{Y}i-1}^{\mathfrak{Y}i} \sum_{h=\mathfrak{Y}j-1}^{\mathfrak{Y}j} P_{w,h}^{(2)}$$

لایه چهارم

$$h^{(\mathfrak{F})} = P_{1,1}^{(3)}, P_{1,\mathfrak{Y}}^{(3)}, P_{\mathfrak{Y},1}^{(3)}, P_{\mathfrak{Y},\mathfrak{Y}}^{(3)}$$

لایه پنجم

$$z = \sum_{i=1}^{\mathfrak{F}} W_i^{(\text{fc})} h_i^{(\mathfrak{F})}$$

$$\hat{y} = \sigma(z)$$

i.

دقت کنید که می‌توان نوشت :

$$\frac{\partial L}{\partial h_k^{(\mathfrak{F})}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\hat{y}}{z} \times \frac{\partial z}{\partial h_k^{(\mathfrak{F})}} = \frac{\partial L}{\partial \hat{y}} \times \hat{y}(1 - \hat{y}) \times W_k^{(\text{fc})}$$

طبق تعریف داریم:

$$P_{i,j}^{(3)} = h_{\mathfrak{Y}i+j-\mathfrak{Y}}^{(\mathfrak{F})} \quad (1 \leq i, j \leq \mathfrak{Y})$$

پس به دست می‌آید:

$$\frac{\partial h_k^{(\mathfrak{F})}}{\partial P_{i,j}^{(3)}} = I(k = \mathfrak{Y}i + j - \mathfrak{Y})$$

در نتیجه خواهیم داشت:

$$\frac{\partial L}{\partial P_{i,j}^{(3)}} = \sum_{k=1}^{\mathfrak{F}} \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial h_k^{(\mathfrak{F})}} \times \frac{\partial h_k^{(\mathfrak{F})}}{\partial P_{i,j}^{(3)}} = \frac{\partial L}{\partial \hat{y}} \times \hat{y}(1 - \hat{y}) \times W_{\mathfrak{Y}i+j-\mathfrak{Y}}^{(\text{fc})}$$

ii.

می توان نوشت :

$$\frac{\partial L}{\partial P_{i,j,k}^{(1)}} = \sum_{m=1}^F \sum_{n=1}^F \frac{\partial L}{\partial P_{m,n}^{(2)}} \times \frac{\partial P_{m,n}^{(2)}}{\partial P_{i,j,k}^{(1)}}$$

ابتدا داریم:

$$\frac{\partial L}{\partial P_{m,n}^{(2)}} = \sum_{p=1}^2 \sum_{q=1}^2 \frac{\partial L}{\partial P_{p,q}^{(3)}} \times \frac{\partial P_{p,q}^{(3)}}{\partial P_{m,n}^{(2)}}$$

به توجه به آنچه آوردیم:

$$\frac{\partial P_{p,q}^{(3)}}{\partial P_{m,n}^{(2)}} = \frac{1}{F} I((m = 2p \vee m = 2p - 1) \wedge (n = 2q \vee n = 2q - 1))$$

$$\frac{\partial P_{m,n}^{(2)}}{\partial P_{i,j,k}^{(1)}} = W_{m-i,n-j,k}^{(2)} I((m = i \vee m = i + 1) \wedge (n = j \vee n = j + 1))$$

بنابراین خواهیم داشت:

$$\frac{\partial L}{\partial P_{i,j,k}^{(1)}} = \frac{1}{F} \times \hat{y}(1 - \hat{y}) \times \sum_{m=i}^{i+1} \sum_{n=j}^{j+1} W_{\lceil m/2 \rceil + \lceil n/2 \rceil - 2}^{(fc)} W_{m-i,n-j,k}^{(2)}$$

(ب) گرادیان زیان نسبت به وزن $W_{1,1,k}^{(1)}$ یکی از فیلترهای لایه اول کانولوشن را محاسبه کنید:

$$\frac{\partial L}{\partial W_{1,1,k}^{(1)}}$$

پاسخ:

حال با توجه به رابطه ای که به دست آوردیم داریم :

$$\frac{\partial L}{\partial W_{1,1,k}^{(1)}} = \frac{\partial L}{\partial P_{i,j,k}^{(1)}} \times \frac{\partial P_{i,j,k}^{(1)}}{\partial W_{1,1,k}^{(1)}} = \frac{\partial L}{\partial P_{i,j,k}^{(1)}} X_{i,j,k}$$

با جایگذاری رابطه قبلی داریم :

$$\frac{\partial L}{\partial W_{1,1,k}^{(1)}} = \frac{1}{F} \times \hat{y}(1 - \hat{y}) \times \sum_{m=i}^{i+1} \sum_{n=j}^{j+1} W_{\lceil m/2 \rceil + \lceil n/2 \rceil - 2}^{(fc)} W_{m-i,n-j,k}^{(2)} X_{i,j,k}$$

(ج) عبارت گرادیان $\frac{\partial L}{\partial W_{i,j,k}^{(2)}}$ ، یکی از وزنهای فیلتر کانولوشن لایه دوم را به دست آورید.

پاسخ:

$$\frac{\partial L}{\partial W_{i,j,k}^{(y)}} = \sum_{m=1}^f \sum_{n=1}^f \frac{\partial L}{\partial P_{m,n}^{(y)}} \times \frac{\partial P_{m,n}^{(y)}}{\partial W_{i,j,k}^{(y)}}$$

با توجه به آنچه که در ابتدا به دست آوردیم، داریم:

$$P_{m,n}^{(y)} = \sum_{k=1}^f \sum_{w=1}^y \sum_{h=1}^y W_{w,h,k}^{(y)} P_{m+w-1, n+h-1, k}^{(1)}$$

بنابراین اگر $m+i$ و $j+n$ در محدوده لایه دوم قرار بگیرند، خواهیم داشت:

$$\frac{\partial P_{m,n}^{(y)}}{\partial W_{i,j,k}^{(y)}} = P_{m+i, n+j, k}^{(1)}$$

و در غیر این صورت داریم:

$$\frac{\partial P_{m,n}^{(y)}}{\partial W_{i,j,k}^{(y)}} = 0$$

حال با جایگذاری مشتق خواسته شده به شکل زیر محاسبه می‌شود:

$$\frac{\partial L}{\partial W_{i,j,k}^{(y)}} = \frac{1}{f} \hat{y}(1 - \hat{y}) \sum_{m=1}^f \sum_{n=1}^f W_{\lceil m/y \rceil + \lceil n/y \rceil - y}^{(fc)} P_{m+i-1, n+j-1, k}^{(1)}$$

۴. (۲۲ نمره) در این سوال قصد داریم برخی گونه‌های مختلف شبکه‌های پیچشی و همچنین معماری‌های مبتنی بر CNN معرفی شده در دهه اخیر را مورد بررسی و تحلیل قرار دهیم.

الف) شبکه DenseNet

- تفاوت اصلی DenseNet's dense connections و ResNet's residual connections را بیان کنید. در مورد هر کدام نیز، ابتدا توضیح مختصری بدهید.

پاسخ:

تفاوت اصلی دو مدل ذکر شده را می‌توان در نحوه ارتباط لایه‌های مختلف شبکه، به شکل زیر بیان نمود:

- DenseNet در Dense Connections: در DenseNet، هر لایه به تمام لایه‌های قبلی متصل است و ویژگی‌ها از طریق الحاق (Concatenation) منتقل می‌شوند. این امر باعث استفاده مجدد از ویژگی‌ها و افزایش بازدهی مدل می‌شود.
- ResNet در Residual Connections: در ResNet، اتصال بین لایه‌ها از طریق جمع ویژگی‌ها (Summation) انجام می‌شود. این روش گرادینان‌ها را مستقیماً به لایه‌های قبلی باز می‌گرداند اما در مقایسه با DenseNet ویژگی‌های کمتری را در طول شبکه حفظ می‌کند.

DenseNet	ResNet
خروجی هر لایه به تمام لایه‌های بعدی الحاق می‌شود.	خروجی هر لایه با خروجی لایه قبلی جمع می‌شود.
گرادیان مستقیماً از لایه خروجی به تمام لایه‌ها منتقل می‌شود.	گرادیان از طریق مسیرهای کوتاه (Short-cut Connections) به لایه‌های اولیه بازمی‌گردد.
پارامترهای کمتری نیاز دارد زیرا ویژگی‌های تکراری مجدداً در نظر گرفته نمی‌شوند.	پارامترهای بیشتری به دلیل جمع ویژگی‌ها در هر لایه نیاز دارد.
هر لایه از تمام ویژگی‌های قبلی بهره می‌برد که به یادگیری بهتر کمک می‌کند.	ممکن است با افزایش عمق، برخی لایه‌ها اطلاعات کمتری بیاموزند.
به طور مؤثرتری از مشکل vanishing gradient جلوگیری می‌کند.	احتمال بیشتری برای vanishing gradient دارد (هرچند کاهش یافته است).

- توضیح دهید که DenseNet چگونه مشکل vanishing gradient را کاهش می‌دهد و مزیت محاسباتی آن چیست؟

پاسخ:

در DenseNet، هر لایه به طور مستقیم به تمام لایه‌های قبلی متصل است. این ارتباط مستقیم باعث می‌شود جریان گرادیان از لایه‌های ابتدایی به لایه‌های پایانی بدون کاهش ضعیف انجام شود. به همین دلیل مشکل vanishing gradient در این معماری کاهش می‌یابد. همچنین اتصال متراکم باعث بهبود انتشار اطلاعات و کاهش نیاز به تعداد زیاد پارامترها می‌شود.

- با در نظر گرفتن نرخ رشد k در DenseNet، اگر هر لایه k ویژگی نگاشت جدید تولید کند و ورودی یک dense block دارای ۳۲ کانال باشد، اگر $k = ۲۴$ باشد، لایه سوم در بلوک چند کانال خروجی خواهد داشت؟

پاسخ:

اگر تعداد کانال‌های ورودی اولیه ۳۲ باشد و بلوک دارای ۳ لایه باشد:

تعداد کانال‌های خروجی = ورودی اولیه + $k \times x = ۳۲ + ۲۴ \times ۳$

$$۳۲ + ۲۴ \times ۳ = ۳۲ + ۷۲ = ۱۰۴$$

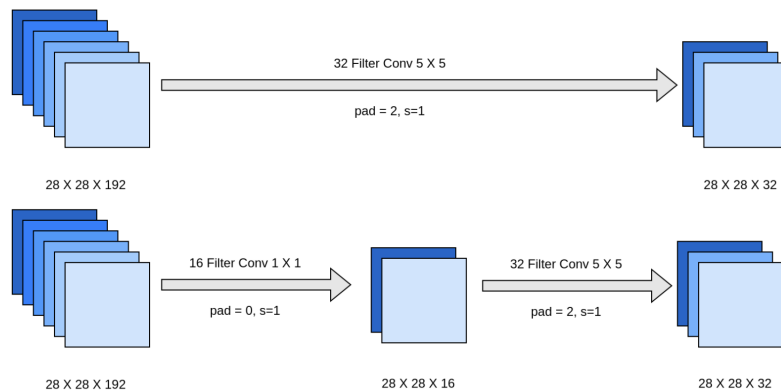
$۱۰۴ =$ تعداد کانال‌های خروجی

این فرمول نشان‌دهنده افزایش تدریجی کانال‌ها در هر بلوک با توجه به نرخ رشد k است.

ب) شبکه GoogleNet

پیمانه پیدایش با هدف کاهش پیچیدگی محاسباتی به وجود آمد. قصد داریم ابتدا با پیمانه پیدایش و سپس شبکه GoogleNet آشنا شویم.

- در شکل ۳ یک لایه کانولوشن به دو صورت نشان داده شده است. ورودی هر دو حالت $28 \times 28 \times 192$ و خروجی هر دو $28 \times 28 \times 32$ است. با محاسبه تعداد کل عملیات‌های انجام شده در هر حالت، پیچیدگی محاسباتی دو حالت نشان داده شده را مقایسه کنید. مشخص کنید که افزودن فیلتر کانولوشن 1×1 ، پیچیدگی محاسباتی را چند درصد کاهش یا افزایش داده است؟



شکل ۳: لایه کانولوشن معمولی (بالا) و لایه کانولوشن با فیلتر 1×1 (پایین)

پاسخ:

- حالت اول: کانولوشن عادی (5×5)

- * در این حالت، از 32 فیلتر 5×5 استفاده می‌شود.
 - * هر فیلتر با ورودی $28 \times 28 \times 192$ اعمال شده و خروجی تولید می‌کند.
 - * تعداد کل عملیات محاسباتی به صورت زیر محاسبه می‌شود:
- $$28 \times 28 \times 192 \times 32 \times (5 \times 5) = 120,422,400 \text{ عملیات}$$

- حالت دوم: کانولوشن با کاهش ابعاد (1×1)

- * در این حالت، ابتدا کانال‌های ورودی به 16 کانال کاهش می‌یابد. سپس کانولوشن 5×5 اعمال می‌شود.
- * تعداد عملیات محاسباتی در دو مرحله انجام می‌شود:
- . مرحله اول: کاهش ابعاد با کانولوشن (1×1)
- تعداد عملیات:

$$28 \times 28 \times 192 \times 16 \times (1 \times 1) = 2,407,424$$

- . مرحله دوم: کانولوشن 5×5 روی کانال‌های کاهش یافته
- تعداد عملیات:

$$28 \times 28 \times 16 \times 32 \times (5 \times 5) = 10,648,576$$

* جمع کل:

$$2,407,424 + 10,648,576 = 130,056,000 \text{ عملیات}$$

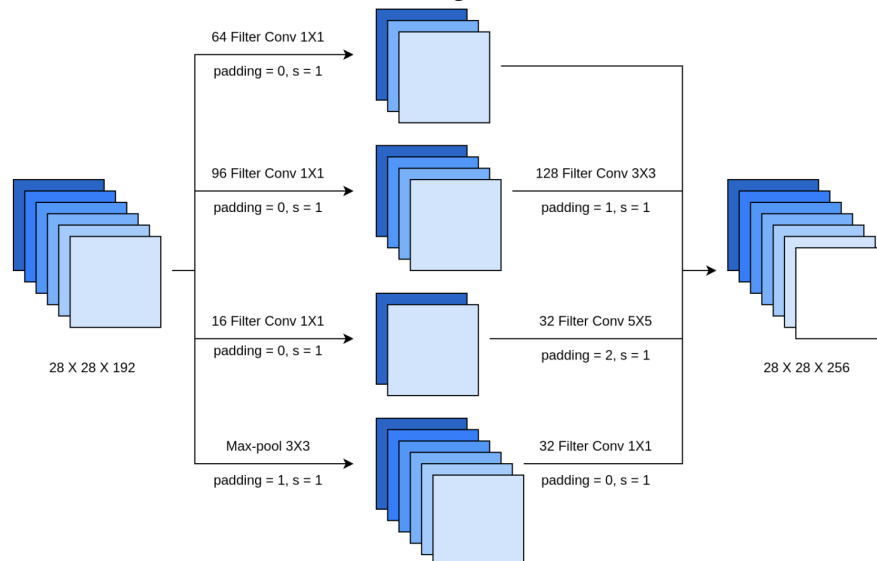
- درصد کاهش پیچیدگی:

- * کاهش پیچیدگی با استفاده از فرمول زیر محاسبه می‌شود:

$$\frac{120,422,400 - 130,056,000}{120,422,400} \times 100 \approx 89.7\%$$

مشاهده می‌شود که پیچیدگی محاسباتی مدل حدوداً 90% کاهش پیدا کرده است که به شدت قابل توجه می‌باشد.

- در شکل ۴ یک پیمانه پیدایش به صورت کامل نشان داده شده است. ترکیب متوالی این پیمانه‌ها با هم، شبکه GoogleNet را تشکیل می‌دهند. توضیح دهید که دلیل استفاده از فیلترهایی با اندازه متفاوت (برای مثال 1×1 ، 3×3 و 5×5) و حتی فیلترهای با انواع متفاوت (استفاده از Max pooling در کنار کانولوشن) چیست؟ تعداد کل محاسبات در این شبکه را بدست آورید.



شکل ۴: شمای کلی پیمانه پیدایش

پاسخ:

محاسبه تعداد عملیات برای هر نوع فیلتر در شبکه به صورت زیر انجام می‌شود:

- برای فیلتر 1×1 با ۶۴ کانال:
 $28 \times 28 \times 192 \times 64 \times (1 \times 1) = 9,625,600$
- برای فیلتر 1×1 با ۹۶ کانال:
 $28 \times 28 \times 192 \times 96 \times (1 \times 1) = 14,438,400$
- برای فیلتر 1×1 با ۱۶ کانال:
 $28 \times 28 \times 192 \times 16 \times (1 \times 1) = 2,407,424$
- برای Max Pooling 3×3 :
 $28 \times 28 \times 192 \times (3 \times 3) = 1,354,752$
- برای فیلتر 3×3 با ۱۲۸ کانال:
 $28 \times 28 \times 96 \times 128 \times (3 \times 3) = 27,704,576$
- برای فیلتر 5×5 با ۳۲ کانال:
 $28 \times 28 \times 16 \times 32 \times (5 \times 5) = 10,648,576$
- برای فیلتر 1×1 با ۳۲ کانال:
 $28 \times 28 \times 192 \times 32 \times (1 \times 1) = 5,949,440$

جمع کل عملیات‌ها:

$$9,625,600 + 14,438,400 + 2,407,424 + 1,354,752 + 27,704,576 + 10,648,576 + 5,949,440 = 72,128,368$$

دلایل استفاده از فیلترها با اندازه‌های مختلف:

- پردازش چند مقیاس: فیلترهای 3×3 و 5×5 ویژگی‌های بزرگ‌تر و پیچیده‌تر را استخراج می‌کنند، در حالی که فیلتر 1×1 ویژگی‌های محلی را بررسی می‌کند و تعداد کانال‌ها را کاهش می‌دهد.

- **کاهش ابعاد:** کانولوشن 1×1 تعداد کانال‌های ورودی را کاهش داده و پیچیدگی محاسباتی فیلترهای بزرگ‌تر را بهینه می‌کند.
- **ترکیب اطلاعات:** مسیرهای موازی اطلاعات متنوع را استخراج کرده و در خروجی ترکیب می‌کنند.

• آیا این معماری نسبت به vanishing gradient مقاوم است؟ همچنین توضیح دهید که auxiliary classifier موجود در این معماری چگونه به بهبود جریان گرادیان و پایداری بهینه‌سازی در شبکه کمک می‌کند. در نهایت، محدودیت‌های احتمالی این طراحی در کاربردهای مدرن یادگیری عمیق را ارزیابی کنید.

پاسخ:

گوگل نت (GoogLeNet)، که به عنوان بخشی از خانواده معماری‌های Inception در سال ۲۰۱۴ معرفی شد، طراحی شبکه‌های عصبی کانولوشنی را با تمرکز بر کارایی، ماژولار بودن و مقیاس‌پذیری متحول کرد.

۱. مقاومت در برابر مشکل ناپدید شدن گرادیان (Vanishing Gradient)

مشکل ناپدید شدن گرادیان زمانی رخ می‌دهد که گرادیان‌ها هنگام بازپراکنش (backpropagation) در شبکه‌های عمیق به‌طور نمایی کاهش می‌یابند. این مشکل مانع از آموزش مؤثر شبکه‌های بسیار عمیق می‌شود.

گوگل نت به‌طور غیرمستقیم این مشکل را از طریق چندین ویژگی معماری حل می‌کند:

- **ماژولار بودن در عمق شبکه:** به‌جای انباشتن لایه‌های کانولوشنی یکپارچه، گوگل نت از ماژول‌های Inception استفاده می‌کند. این ماژول‌ها بلوک‌هایی ماژولار هستند که فیلترهایی با اندازه‌های مختلف را به‌صورت موازی ترکیب می‌کنند. این طراحی به شبکه اجازه می‌دهد عمیق‌تر باشد بدون اینکه مشکل ناپدید شدن گرادیان را تشدید کند، چرا که هر ماژول بر استخراج ویژگی‌ها در سطوح مختلف تمرکز دارد.

- **نرمال‌سازی دسته‌ای (Batch Normalization):** اگرچه در نسخه اولیه گوگل نت وجود نداشت (در Inception-v2 معرفی شد)، نرمال‌سازی دسته‌ای، فعال‌سازی‌های میانی را نرمال کرده و تغییرات ناگهانی داخلی را کاهش می‌دهد، که به پایداری گرادیان کمک می‌کند.

- **کلاس‌یفایرهای کمکی (Auxiliary Classifiers):** کلاس‌یفایرهای کمکی لایه‌های softmax میانی هستند که در بخش‌های مختلف شبکه قرار داده می‌شوند. این کلاس‌یفایرها به‌عنوان "تقویت‌کننده‌های گرادیان" عمل می‌کنند و سیگنال‌های گرادیان اضافی را به لایه‌های قبلی در طی بازپراکنش ارسال می‌کنند. این ویژگی به کاهش مشکل ناپدید شدن گرادیان کمک می‌کند.

۲. کلاس‌یفایرهای کمکی در گوگل نت

کلاس‌یفایرهای کمکی لایه‌های خروجی میانی هستند که برای کمک به فرآیند آموزش اضافه شده‌اند. این کلاس‌یفایرها سیگنال‌های خطای اضافی را به لایه‌های قبلی ارسال کرده و جریان گرادیان را بهبود می‌بخشند.

اثر بخشی:

- **بهبود جریان گرادیان:** کلاس‌یفایرهای کمکی با تزریق سیگنال‌های گرادیان در نقاط میانی، به کاهش ناپدید شدن گرادیان کمک می‌کنند و تضمین می‌کنند که لایه‌های اولیه به‌طور مؤثری آموزش ببینند.

- **منظم‌سازی (Regularization):** با مجبور کردن لایه‌های میانی به تولید ویژگی‌های معنادار، آموزش شبکه را تثبیت کرده و تعمیم‌پذیری آن را بهبود می‌بخشند.

محدودیت‌ها:

- **راه حل ناقص:** این کلاس‌یفایرها مشکل ناپدید شدن گرادیان را تا حدی کاهش می‌دهند، اما به‌طور کامل حل نمی‌کنند، به‌ویژه در شبکه‌های بسیار عمیق.
- **استفاده موقتی:** تأثیر این کلاس‌یفایرها فقط در طی آموزش است و در فرآیند استنتاج (Inference) نقشی ندارند.
- **جایگزین‌های مدرن:** تکنیک‌هایی مانند اتصالات باقی‌مانده (Residual Connections) در ResNet جایگزین کلاس‌یفایرهای کمی شده‌اند، زیرا راه‌حل‌های قوی‌تری برای جریان گرادیان ارائه می‌دهند.

با توجه به توضیحات بالا می‌توان نتیجه گرفت کلاس‌یفایرهای کمی در گوگل‌نت تا حدی در جلوگیری از ناپدید شدن گرادیان مؤثر بودند، به‌ویژه در شبکه‌هایی با عمق مشابه. اما امروزه با ظهور روش‌های برتر در معماری‌های مدرن مانند اتصالات باقی‌مانده، این تکنیک به‌طور گسترده‌ای منسوخ شده است.

ج) شبکه‌های کانولوشنی Deformable

- تفاوت شبکه‌های کانولوشنی عادی و شبکه‌های کانولوشنی Deformable را از نظر grid sampling مقایسه کنید. همچنین شبکه‌های Deformable چگونه می‌توانند نسبت به Geometric transformation در تصاویر انعطاف‌پذیر باشند؟
- مفهوم offset در این شبکه‌ها به چه معناست و چگونه محاسبه می‌شود؟

پاسخ:

شبکه‌های کانولوشن معمولی از یک ساختار شبکه نمونه‌گیری ثابت استفاده می‌کنند. به‌عنوان مثال، در یک هسته $K \times K$ ، نقاط نمونه‌گیری همیشه در موقعیت‌های از پیش تعریف‌شده قرار دارند. این ساختار ثابت باعث می‌شود که این شبکه‌ها در مواجهه با تغییرات هندسی انعطاف‌پذیری محدودی داشته باشند.

در مقابل، شبکه‌های Deformable با اضافه کردن آفست‌های قابل یادگیری به نقاط شبکه نمونه‌گیری، این مشکل را حل می‌کنند. این آفست‌ها، که در طی فرآیند آموزش و بر اساس ویژگی‌های محلی یاد گرفته می‌شوند، به شبکه اجازه می‌دهند تا به‌صورت پویا و محلی تغییرات هندسی را مدل‌سازی کند. شبکه‌های Deformable قادر به مدل‌سازی تغییرات هندسی پیچیده مانند تغییر مقیاس، چرخش، و تغییر نسبت ابعاد هستند. این ویژگی به دلیل وجود آفست‌های متغیر و قابل یادگیری است که امکان نمونه‌گیری از مکان‌های غیرثابت را فراهم می‌کنند.

معادله کانولوشن معمولی:

$$y(p_{\cdot}) = \sum_{p_n \in R} w(p_n) \cdot x(p_{\cdot} + p_n)$$

معادله کانولوشن Deformable با اضافه کردن آفست‌ها:

$$y(p_{\cdot}) = \sum_{p_n \in R} w(p_n) \cdot x(p_{\cdot} + p_n + \Delta p_n)$$

در اینجا Δp_n آفست‌هایی هستند که در طی فرآیند یادگیری تعیین می‌شوند.

آفست‌ها به کمک یک لایه کانولوشن اضافی محاسبه می‌شوند که ورودی آن همان ویژگی‌های نقشه ورودی است. خروجی این لایه، یک میدان آفست است که در همان وضوح نقشه ورودی قرار دارد. این آفست‌ها برای هر مکان به‌صورت محلی و مستقل یاد گرفته می‌شوند.

برای مکان‌های کسری که به دلیل آفست ایجاد می‌شوند، از درونیایی دوجمله‌ای استفاده می‌شود:

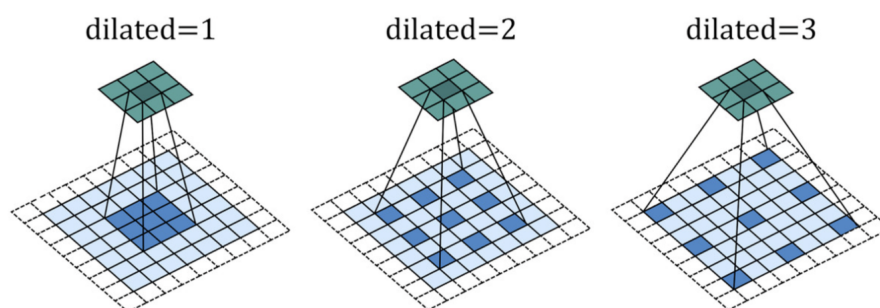
$$x(p) = \sum_q G(q, p) \cdot x(q)$$

که در آن $G(q, p)$ وزن درونیایی دوجمله‌ای است و p یک مکان کسری است.

مزایای شبکه‌های Deformable:

- **انعطاف‌پذیری:** این شبکه‌ها می‌توانند تغییرات هندسی پیچیده را به صورت پویا مدل‌سازی کنند.
- **کارایی بالا:** به دلیل معماری سبک و ساده، اضافه کردن این ماژول‌ها سربار محاسباتی کمی دارد.
- **بهبود عملکرد:** این شبکه‌ها در وظایف پیچیده‌ای مانند تشخیص اشیاء و بخش‌بندی معنایی عملکرد بهتری دارند.

۵. (۲۱ نمره) در شبکه‌های پیچشی به صورت متداول از لایه‌های کانولوشن ساده استفاده می‌شود که با آن آشنا هستید. نوع دیگری از لایه‌ها که می‌توان از آنان در شبکه‌های پیچشی استفاده نمود، لایه کانولوشن گسترش یافته یا متسع است. در شکل ۵ تصویر شهودی از فیلتر کانولوشن گسترش یافته ارائه شده است، این فیلترها میان خانه‌هایی که فیلتر با استفاده از اطلاعات آن‌ها لایه بعد را محاسبه می‌کند فاصله می‌اندازند یا به بیانی دیگر در زمان اعمال فیلتر و انجام عملیات ضرب کانولوشن، بر روی ورودی با گام (dilated) بزرگتری حرکت می‌کنیم، توجه کنید طول گام مفهومی متفاوت نسبت به طول گام (stride) در لایه‌های شبکه کانولوشن دارد.



شکل ۵: شهودی از کانولوشن گسترش یافته با گام‌های متفاوت

همانطور که در شکل ۱ نیز مشخص است این روش، یک روش کم هزینه برای افزایش محدوده دید شبکه‌های پیچشی است. کانولوشن گسترش یافته بصورت فرم بسته ریاضی زیر تعریف می‌شود.

$$(K \star_D I)(i, j) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} K(m, n) I(i + Dm, j + Dn)$$

الف) برای ورودی $I \in \mathbb{R}^{M \times N}$ و کرنل $K \in \mathbb{R}^{F \times F}$ نشان دهید خروجی عملگر کانولوشن گسترش یافته دارای ابعاد $(M - DF + D) \times (N - DF + D)$ است.

پاسخ:

طبق فرض صورت سوال، کرنل اولیه دارای ابعاد $F \times F$ می‌باشد. در هر ردیف و ستون این کرنل، تعداد فضاهای خالی برابر با $F - 1$ است و هر فضای خالی اندازه $D - 1$ را اشغال می‌کند. بنابراین، ابعاد کرنل پس از اعمال گسترش به صورت زیر محاسبه می‌شود:

$$\text{Dimensions Kernel New} = F + (F - 1)(D - 1) = FD - D + 1$$

این نشان می‌دهد که کرنل جدید دارای ابعاد $(FD - D + 1) \times (FD - D + 1)$ خواهد بود. برای محاسبه ابعاد خروجی کانولوشن، از فرمول زیر استفاده می‌کنیم:

$$\text{Dimensions Output} = (N - FD + D) \times (M - FD + D)$$

که در آن:

$$N' = N - (FD - D) + 1 = N - FD + D$$

$$M' = M - (FD - D) + 1 = M - FD + D$$

بنابراین، اگر ماتریس ورودی دارای ابعاد $N \times M$ باشد، ابعاد خروجی پس از اعمال کرنل گسترش یافته به صورت $(N - FD + D) \times (M - FD + D)$ خواهد بود.

ب) نشان دهید کانولوشن گسترش یافته معادل کانولوشن با کرنل متسع شده K' است. ماتریس A را مشخص کنید. $K' = K \otimes A$ ، عملگر \otimes کرونکر محصول product kronecker است.

پاسخ:

مطابق صورت سوال، تعریف ریاضی کانولوشن گسترش یافته به شکل زیر است:

$$(K *_D I)(i, j) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} K(m, n) \cdot I(i + Dm, j + Dn)$$

که در اینجا D نشان دهنده گام گسترش و I ورودی تصویر می‌باشد. همچنین فرض می‌کنیم که شماره گذاری سطرها و ستون‌های ماتریس از صفر شروع می‌شوند.

ضرب کرونکر بین دو ماتریس M و N به شکل زیر تعریف می‌کنیم:

$$(M \otimes N)_{ij} = M_{ij} \cdot N$$

این بدان معناست که هر عنصر M_{ij} ماتریس M در کل ماتریس N ضرب می‌شود و نتیجه حاصل یک ماتریس بزرگتر خواهد بود. برای تعریف کانولوشن گسترش یافته با استفاده از ضرب کرونکر، ابتدا ماتریس $A_{D \times D}$ را به صورت زیر تعریف می‌کنیم:

$$A_{i,j} = \begin{cases} 1 & i = j = 0 \\ 0 & \text{otherwise} \end{cases}$$

سپس کرنل گسترش یافته K' را به صورت زیر محاسبه می‌کنیم:

$$K' = K \otimes A$$

بنابراین، عناصر K' به صورت زیر خواهند بود:

$$K'_{i,j} = \begin{cases} K_{\frac{i}{D}, \frac{j}{D}} & \frac{i}{D}, \frac{j}{D} \in \mathbb{Z} \\ 0 & \text{otherwise} \end{cases}$$

خروجی کانولوشن با کرنل گسترش یافته K' به صورت زیر محاسبه می شود:

$$\begin{aligned} (K' * I)(i, j) &= \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} K'(m, n) \cdot I(i + m, j + n) \\ &= \sum_{p=-\infty}^{\infty} \sum_{q=-\infty}^{\infty} K(p, q) \cdot I(i + pD, j + qD) \\ &= \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} K(m, n) \cdot I(i + Dm, j + Dn) = (K *_{\mathcal{D}} I)(i, j) \end{aligned}$$

که همان خواسته مسئله است.

ج) فرض کنید یک لایه ورودی با ابعاد $M \times N$ به یک شبکه عصبی با سه لایه کانولوشن گسترش یافته داده می شود. در این شبکه، هر لایه کانولوشن شامل تعدادی فیلتر است که بر روی قسمت های مختلف ورودی حرکت می کند و ویژگی های گوناگون آن را استخراج می کند. هدف ما بررسی تعداد پیکسل های قابل مشاهده توسط یک عنصر خاص مانند (i, j) از خروجی است. به عبارت دیگر، محدوده ای از ورودی که این عنصر را تحت تأثیر قرار می دهد را در صورت وجود امکان محاسبه به صورت پارامتری بیابید.

ابعاد فیلترها به ترتیب از چپ به راست برابر هستند با:

$$(w - 1) \times (w - 1), \quad w \times w, \quad (w + 1) \times (w + 1)$$

پارامتر گسترش این فیلترها به ترتیب از چپ به راست برابر است با:

$$d - 1, \quad d, \quad d + 1$$

گستره دید در شبکه های عصبی کانولوشنی

۱. گستره دید در یک بعد

۱.۱ تنظیمات و تعاریف

ما L لایه کانولوشنی داریم که با $\ell = 1, 2, \dots, L$ شماره گذاری شده اند. برای هر لایه ℓ ، موارد زیر تعریف می شوند:

- اندازه کرنل: $k_{\ell} \in \mathbb{N}$.
- اتساع: $d_{\ell} \in \mathbb{N}$.
- گام: $s_{\ell} \in \mathbb{N}$.

فرض می کنیم Padding اضافی وجود ندارد. (یا فقط بخش هایی را که واقعاً ورودی را مشاهده می کنند در نظر بگیریم). دو مقدار در لایه ℓ تعریف می شوند:

- R_ℓ : اندازه گستره دید در ورودی اصلی برای یک موقعیت خروجی در لایه ℓ .
 - jump_ℓ : میزان حرکت (به واحد اندیس‌های ورودی) وقتی که یک گام در خروجی لایه ℓ حرکت کنیم.
- برای لایه پایه $\ell = 0$:

$$R_0 = 1, \quad \text{jump}_0 = 1.$$

۲.۱ رابطه بازگشتی لایه به لایه

برای هر لایه $\ell \geq 1$:

$$R_\ell = R_{\ell-1} + (k_\ell - 1) \cdot d_\ell \cdot \text{jump}_{\ell-1}, \quad (1)$$

$$\text{jump}_\ell = \text{jump}_{\ell-1} \cdot s_\ell. \quad (2)$$

توضیح شهودی

- **رشد گستره دید:** اگر یک موقعیت خروجی در لایه ℓ را در نظر بگیریم، آن خروجی با کانولوشن یک فیلتر با اندازه k_ℓ و اتساع d_ℓ در خروجی لایه $\ell - 1$ محاسبه می‌شود. گستره دید با مقدار $(k_\ell - 1) \cdot d_\ell$ گسترش می‌یابد.
- **اندازه گام:** جابه‌جایی یک موقعیت در خروجی لایه ℓ معادل جابه‌جایی s_ℓ موقعیت در خروجی لایه $\ell - 1$ است.

۳.۱ روابط کلی

با استفاده از روابط بازگشتی:

$$\text{jump}_0 = 1,$$

$$\text{jump}_\ell = \text{jump}_{\ell-1} \cdot s_\ell,$$

$$R_0 = 1,$$

$$R_\ell = R_{\ell-1} + (k_\ell - 1) \cdot d_\ell \cdot \text{jump}_{\ell-1}.$$

۲. بسط به دوبعدی

۱.۲ روابط در ابعاد افقی و عمودی

برای کانولوشن دوبعدی، موارد زیر تعریف می‌شوند:

• اندازه کرنل: (k_ℓ^H, k_ℓ^W) .

• اتساع: (d_ℓ^H, d_ℓ^W) .

• گام: (s_ℓ^H, s_ℓ^W) .

روابط برای هر بُعد $\alpha \in \{H, W\}$:

$$R_\ell^{(\alpha)} = R_{\ell-1}^{(\alpha)} + (k_\ell^{(\alpha)} - 1) \cdot d_\ell^{(\alpha)} \cdot \text{jump}_{\ell-1}^{(\alpha)},$$

$$\text{jump}_\ell^{(\alpha)} = \text{jump}_{\ell-1}^{(\alpha)} \cdot s_\ell^{(\alpha)}.$$

۲.۲ کرنل‌های مربعی و گام و اتساع یکسان

اگر $s_\ell^H = s_\ell^W = s_\ell$ ، $d_\ell^H = d_\ell^W = d_\ell$ ، $k_\ell^H = k_\ell^W = k_\ell$ و گستره دید در هر دو بُعد برابر خواهد بود:

$$R_\ell^{(H)} = R_\ell^{(W)} = R_\ell^{(1D)}.$$

بنابراین گستره دید نهایی دوبعدی برابر است با:

$$R_\ell^{(2D)} = (R_\ell^{(1D)})^2.$$

پاسخ:

محاسبه گستره دید برای شبکه سه لایه

(آ) لایه ۱:

$$R_1 = 1 + (k_1 - 1) \cdot d_1 = 1 + (w - 2)(d - 1).$$

(ب) لایه ۲:

$$R_2 = R_1 + (k_2 - 1) \cdot d_2 \cdot s = [1 + (w - 2)(d - 1)] + (w - 1)ds.$$

(ج) لایه ۳:

$$R_3 = R_2 + (k_3 - 1) \cdot d_3 \cdot s^2 = [1 + (w - 2)(d - 1) + (w - 1)ds] + w(d + 1)s^2.$$

فرمول نهایی گستره دید در یک بُعد

$$R_3 = 1 + (w - 2)(d - 1) + (w - 1)ds + w(d + 1)s^2.$$

گستره دید نهایی در دوبعد

برای کرنل‌های مربعی:

$$R_3^{(2D)} = (R_3)^2.$$

(د) در مورد Masked Convolution، کاربرد و محدودیت‌های آن تحقیق کنید و ضمن توضیح مختصری در این باره پاسخ دهید چگونه می‌توان با استفاده از کانولوشن گسترش یافته محدودیت Masked Convolution را بهبود بخشید؟

پاسخ:

کانولوشن ماسک شده (Masked Convolution) در حوزه یادگیری عمیق، به ویژه در مدل‌هایی مانند Pix-eCNN که برای پیش‌بینی توالی‌ها و مدل‌های خودتوضیح‌دهنده طراحی شده‌اند، کاربرد دارد. در این روش، فیلترها با استفاده از یک ماسک (Mask) محدود می‌شوند تا تنها از اطلاعات یک یا نقاط خاصی بهره‌برداری کنند. به عنوان مثال:

- در پردازش تصویر، کانولوشن ماسک شده تضمین می‌کند که مقدار هر پیکسل تنها به پیکسل‌های قبلی وابسته باشد.
- در مسائل توالی‌ای مانند پیش‌بینی سری‌های زمانی، این تکنیک از دسترسی به داده‌های آینده جلوگیری می‌کند.

محدودیت‌ها:

- **دامنه دید محدود:** ماسک‌گذاری باعث می‌شود که هر لایه تنها به بخشی محدود از ورودی دسترسی داشته باشد، که ممکن است برای برخی کاربردها ناکافی باشد.
- **عمق زیاد شبکه:** با افزایش دامنه دید، نیاز به تعداد زیادی لایه وجود دارد که این موضوع می‌تواند محاسبات را پیچیده کرده و ممکن است باعث ناپدید شدن گرادیان (Vanishing Gradient) را ایجاد کند.

کانولوشن گسترش یافته (Dilated Convolution) با افزایش فاصله بین نقاط نمونه‌برداری، دامنه دید را بدون افزایش تعداد پارامترها یا تعداد لایه‌ها گسترش می‌دهد. این تکنیک به ویژه در ترکیب با کانولوشن ماسک شده مزایای قابل توجهی دارد:

مزایا:

- **گسترش دامنه دید بدون افزایش پارامترها:** با استفاده از کانولوشن گسترش یافته، می‌توان دامنه دید را به طور قابل توجهی افزایش داد بدون نیاز به افزودن کرنل‌های بزرگ‌تر یا لایه‌های بیشتر.
- **حفظ ساختار محلی و جلوگیری از دسترسی غیرمجاز:** ترکیب کانولوشن ماسک شده با کانولوشن گسترش یافته، امکان حفظ نواحی محلی و در عین حال گسترش دامنه دید را فراهم می‌آورد. ماسک اعمال شده اطمینان حاصل می‌کند که تنها به اطلاعات مجاز دسترسی داشته باشیم، در حالی که کانولوشن گسترش یافته اجازه می‌دهد تا مدل اطلاعات گسترده‌تری را در هر لایه پردازش کند.
- **کاهش عمق شبکه و پیچیدگی محاسباتی:** با گسترش دامنه دید در هر لایه از طریق کانولوشن گسترش یافته، نیاز به وجود لایه‌های متعدد برای افزایش عمق دید کاهش می‌یابد. این امر منجر به کاهش پیچیدگی محاسباتی و افزایش کارایی می‌شود.
- **افزایش انعطاف‌پذیری مدل:** کانولوشن گسترش یافته به مدل امکان می‌دهد تا تعادل بهینه بین دقت محلی و اطلاعات کلی برقرار کند که این امر برای درک بهتر ویژگی‌های مختلف تصویر ضروری است.

۶. (۲۲ نمره) در این سوال به بررسی پدیده vanishing gradient یا ”گرادیان ناپدید شونده” می‌پردازیم.

الف) (امتیازی) پیش از ابداع روش‌های مبتنی بر شبکه‌های عصبی کانولوشنی، feedforward neural network ها قادر به استخراج ویژگی‌ها از تصاویر به طور مستقل نبودند و معمولاً استخراج ویژگی‌ها با روش‌هایی غیر از یادگیری عمیق انجام می‌شد. از طرفی، می‌دانیم که عمیق‌تر کردن شبکه عصبی می‌تواند باعث شود مدل ویژگی‌های پیچیده‌تری را بیاموزد، اما عمیق‌تر کردن شبکه بدون استفاده از تکنیک‌های خاص مشکلاتی نیز به وجود می‌آورد، یکی از مشکلات اصلی که در این فرآیند بروز می‌کند، پدیده vanishing gradient است.

فرض کنید یک شبکه عصبی با تعداد زیادی لایه داریم. می‌دانیم که گرادیان وزن‌های لایه i از طریق عبارت زیر محاسبه می‌شود (که در آن، δ_L مشتق تابع ضرر نسبت به خروجی لایه آخر، W_k وزن لایه k ، f تابع فعال‌سازی و z ورودی تابع فعال‌سازی است).

$$\frac{\partial J}{\partial W^i} = (\delta^i)^T \times \frac{\partial z^i}{\partial W^i} = (\delta^i)^T \times a^{i-1}$$

$$\delta_i = \delta_L \left(\prod_{k=i+1}^L W_k \cdot f'_{k-1}(z_{k-1}) \right)$$

فرض کنید بزرگ‌ترین singular value (مقدار تکین) برای همه ماتریس‌های وزن W_k کوچک‌تر از یک باشد، با فرض این‌که $\delta_L = M$ و تابع فعال‌سازی ما از نوع ReLU باشد، در این صورت نشان دهید که اگر $L \rightarrow \infty$ در این صورت، $|\delta_i| \rightarrow 0$.

همچنین فرض کنید بزرگ‌ترین singular value همه ماتریس‌های وزن برابر 0.9 باشد، برای $L = 100$ با شرایطی که بیان شد یک حد بالا برای $|\delta_i|$ در صورتی که $i = 0$ باشد؛ پیدا کنید.

پاسخ:

اثبات $|\delta_i| \rightarrow 0$ وقتی $L \rightarrow \infty$

گرادیان δ_i به صورت زیر تعریف شده است:

$$\delta_i = \delta_L \prod_{k=i+1}^L W_k \cdot f'(z_{k-1}).$$

با فرض اینکه:

- بزرگ‌ترین مقدار تکین هر ماتریس وزن $\sigma_{\max}(W_k) < 1$.
- مشتق تابع فعال‌سازی ReLU $\|f'(z)\| \leq 1$.

طبق قاعده مثلثی:

$$\|\delta_i\| \leq \|\delta_L\| \cdot \prod_{k=i+1}^L \|W_k\| \cdot \|f'(z_{k-1})\|.$$

از آنجا که $\|W_k\| = \sigma_{\max}(W_k) < 1$ و نتیجه می‌گیریم:

$$\|\delta_i\| \leq \|\delta_L\| \cdot \prod_{k=i+1}^L \sigma_{\max}(W_k).$$

با فرض اینکه $\sigma_{\max}(W_k) = \sigma < 1$ برای همه k ، داریم:

$$\|\delta_i\| \leq \|\delta_L\| \cdot \sigma^{L-i}.$$

وقتی $L \rightarrow \infty$ ، $\sigma^{L-i} \rightarrow 0$ زیرا $\sigma < 1$. بنابراین:

$$\|\delta_i\| \rightarrow 0.$$

حد بالا برای $|\delta_i|$ وقتی $L = 100$ و $\sigma_{\max} = 0.9$ و $i = 0$

برای $L = 100$ و $\sigma = 0.9$ ، داریم:

$$\|\delta_i\| \leq \|\delta_L\| \cdot \sigma^{L-i}.$$

برای مثال:

(آ) اگر $i = 0$:

$$\|\delta_0\| \leq \|\delta_L\| \cdot (0.9)^{100-0} = \|\delta_L\| \cdot (0.9)^{100}.$$

با محاسبه $(0.9)^{50}$:

$$(0.9)^{50} \approx 0.000027$$

بنابراین:

$$\|\delta_{50}\| \leq \|\delta_L\| \cdot 0.000027$$

ب) یکی از راهکارها برای جلوگیری از vanishing gradient استفاده از یک initialization مناسب است. فرض کنید یک لایه کانولوشن با طول و عرض 8×8 داریم که یک لایه نهفته با تعداد کانال ۳ را به یک لایه نهفته با تعداد کانال ۵ متصل می‌کند. اگر از **Kaiming initialization** با توزیع نرمال استفاده کنیم، پارامترهای توزیع نرمال را برای مقداردهی اولیه کرنل‌های این لایه به دست آورید. **راهنمایی:** تعداد واحدهای ورودی به هر کرنل، حاصل ضرب تعداد کانال‌های ورودی در مساحت کرنل است.

پاسخ:

برای Kaiming Initialization با فرض تابع فعال‌سازی ReLU، وزن‌های هر کرنل از توزیع نرمال زیر نمونه‌برداری می‌شوند:

$$W \sim \mathcal{N}(0, \text{Var}(W)).$$

واریانس این توزیع به صورت زیر محاسبه می‌شود:

$$\text{Var}(W) = \frac{2}{n_{\text{in}}},$$

که در آن n_{in} تعداد واحدهای ورودی به هر کرنل است.

محاسبه n_{in} :

تعداد واحدهای ورودی به هر کرنل برابر است با:

$$n_{\text{in}} = (\text{تعداد کانال‌های ورودی}) \times (\text{مساحت کرنل}) = 3 \times (8 \times 8) = 192.$$

توزیع نرمال:

با جایگذاری $n_{\text{in}} = 192$ در فرمول واریانس:

$$\text{Var}(W) = \frac{2}{192} = 0.0104.$$

انحراف معیار توزیع برابر است با:

$$\text{Std}(W) = \sqrt{\text{Var}(W)} = \sqrt{0.0104} \approx 0.102.$$

نتیجه نهایی:

برای هر کرنل، وزن‌ها از توزیع نرمال با مشخصات زیر نمونه‌برداری می‌شوند:

$$W \sim \mathcal{N}(0, 0.0104).$$

توضیح تکمیلی: در صورتی که در پارامترهای توزیع نرمال، به جای واریانس ($\text{Var} = 0.0104$)، انحراف معیار ($\text{Std} = 0.102$) ذکر شود نیز پاسخ معتبر است.

ج) در اوایل توسعه‌ی شبکه‌های عصبی، تابع فعال‌سازی Sigmoid برای بسیاری از شبکه‌ها انتخاب می‌شد. این تابع، که مقادیر ورودی را بین صفر و یک نگاشت می‌کند، مناسب به نظر می‌رسید. اما به مرور زمان و با افزایش عمق شبکه‌ها، معلوم شد که استفاده از تابع Sigmoid نیز می‌تواند باعث بروز پدیده vanishing gradient شود، یکی از دلایلی که معماری‌های Modern Convolutional Neural Network در ابتدا نمی‌توانستند عملکرد کافی داشته باشند عدم امکان عمیق‌سازی آن‌ها به اندازه کافی بود.

در سال ۲۰۱۲ معماری AlexNet معرفی شد که نسبت به معماری مشابه قبلی خود LeNet پیشرفت بسیار قابل توجهی داشت، یکی از تفاوت‌های این معماری عمق بیشتر و استفاده از تابع فعال‌سازی ReLU به جای Tanh / Sigmoid بود.

با توجه به این موضوع، پاسخ دهید: چرا استفاده از تابع ReLU به جای Sigmoid به جلوگیری از مشکل vanishing gradient کمک می‌کند؟

پاسخ:

استفاده از تابع ReLU به جای Sigmoid به دلایل زیر به کاهش پدیده vanishing gradient کمک می‌کند:

- **مشکل گرادیان کوچک در Sigmoid:** در تابع Sigmoid، اگر مقدار ورودی x کمی از ۰ دور شود (چه مثبت چه منفی)، مقدار گرادیان بسیار کوچک می‌شود:

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

این مسئله باعث می‌شود که گرادیان‌ها هنگام انتشار به لایه‌های ابتدایی تقریباً صفر شوند. در مقابل، ReLU در محدوده $x > 0$ گرادیانی برابر ۱ دارد و از این کاهش جلوگیری می‌کند.

- **عدم اشباع در ReLU:** برخلاف Sigmoid که خروجی آن بین ۰ و ۱ محصور است و به اشباع منجر می‌شود، ReLU هیچ محدودیتی برای مقادیر مثبت ندارد. این ویژگی کمک می‌کند گرادیان در لایه‌های عمیق کاهش نیابد و یادگیری بهتر انجام شود.

د) اگرچه ReLU مزایای زیادی ارائه می‌دهد، یکی از معایب آن مسئله‌ی Dead Neurons است، توضیح دهید این مشکل چیست. Leaky ReLU، Parametric ReLU و ELU توابعی هستند که برای رفع مشکلات تابع ReLU طراحی شده‌اند، با استفاده از معادله و معادله مشتق این توابع بیان کنید هر کدام چه معایب و مزایایی نسبت به تابع ReLU دارند.

پاسخ:

اگر ورودی یک نورون مقدار منفی داشته باشد، خروجی ReLU صفر می‌شود و گرادینان نیز صفر خواهد بود. این موضوع باعث می‌شود نورون دیگر به‌روزرسانی نشود و به اصطلاح ”مرده” بماند.

توابع جایگزین ReLU:

۱. Leaky ReLU:

$$f(x) = \begin{cases} x & x \geq 0 \\ \alpha x & x < 0 \end{cases}, \quad f'(x) = \begin{cases} 1 & x \geq 0 \\ \alpha & x < 0 \end{cases}$$

مزایا: جلوگیری از Dead Neurons با دادن یک شیب کوچک ($\alpha > 0$) به مقادیر منفی.
معایب: انتخاب مقدار α ثابت ممکن است برای همه مسائل بهینه نباشد.

۲. Parametric ReLU (PReLU):

$$f(x) = \begin{cases} x & x \geq 0 \\ \alpha x & x < 0 \end{cases}, \quad f'(x) = \begin{cases} 1 & x \geq 0 \\ \alpha & x < 0 \end{cases}$$

α پارامتری قابل یادگیری است.
مزایا: انعطاف‌پذیری بیشتر با یادگیری مقدار بهینه α .
معایب: پیچیدگی بیشتر در آموزش و خطر بیش‌برازش.

۳. Exponential Linear Unit (ELU):

$$f(x) = \begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}, \quad f'(x) = \begin{cases} 1 & x \geq 0 \\ \alpha e^x & x < 0 \end{cases}$$

مزایا: کاهش Dead Neurons و میانگین خروجی نزدیک به صفر برای سرعت‌دهی به همگرایی.
معایب: محاسبات بیشتر به دلیل استفاده از تابع نمایی.

ه) یک شبکه با ۲ بلاک residual را در نظر بگیرید. نشان دهید چگونه این بلاک‌ها در شبکه residual کمک می‌کند تا از مشکل vanishing gradient جلوگیری شود؟ این شبکه را با شبکه‌های نرمال بدون residual blocks مقایسه کنید. برای سادگی در نوشتار، تابعی که در هر بلاک روی ورودی اعمال می‌شود را با F نشان دهید که شامل activation function هم می‌شود.

پاسخ:

بلاک‌های residual با افزودن ورودی (x) به خروجی یک لایه ($\mathcal{F}(x)$) به حل مشکل vanishing gradient کمک می‌کنند. ساختار ریاضی شبکه دو بلاکی residual به صورت زیر است:

$$\begin{aligned}y_1 &= \mathcal{F}_1(x) + x, \\ y_2 &= \mathcal{F}_2(y_1) + y_1.\end{aligned}$$

برای محاسبه گرادیان‌ها:

$$\frac{\partial y_2}{\partial w_2} = \mathcal{F}'_2(y_1).$$

$$\frac{\partial y_2}{\partial w_1} = (\mathcal{F}'_2(y_1) + 1) \cdot \mathcal{F}'_1(x),$$

در شبکه نرمال:

$$\frac{\partial y_2}{\partial w_2} = \mathcal{F}'_2(y_1).$$

$$\frac{\partial y_2}{\partial w_1} = \mathcal{F}'_2(y_1) \cdot \mathcal{F}'_1(x),$$

در شبکه residual، وجود جمله اضافی $+1$ در گرادیان باعث می‌شود که جریان گرادیان بدون کاهش شدید از یک بلاک به بلاک دیگر منتقل شود. اما در شبکه‌های معمولی، گرادیان تنها به توابع تبدیل \mathcal{F}_i وابسته است و در اثر ضرب‌های متوالی، با افزایش عمق شبکه به شدت کاهش می‌یابد، که این مسئله مشکل vanishing gradient را ایجاد می‌کند، اما در لایه‌های residual تا حد خوبی از این مشکل جلوگیری می‌شود.