



تمرین چهارم

پاسخ مسئله ۱.

(الف)

- موارد خواسته شده در صورت سوال، همگی در جدول زیر قرار گرفته‌اند.

ویژگی	شبکه تمام متصل	شبکه پیچشی
تعداد واحدهای خروجی	I	$(W - K + 1) \times (H - K + 1) \times I$
تعداد اتصالات	$H \times W \times I \times J$	$K \times K \times J \times (W - K + 1) \times (H - K + 1) \times I$
تعداد وزن‌های قابل یادگیری	$H \times W \times I \times J$	$K \times K \times J \times I$

در ادامه توضیحات دقیق در ارتباط با نحوه محاسبه این مقادیر آورده شده است.

۱. تعداد واحدهای خروجی:

شبکه تمام متصل: در یک شبکه تمام متصل، تعداد واحدهای خروجی برابر با تعداد نوروهای موجود در لایه خروجی است که با I نمایش داده می‌شود. این مقدار به ابعاد ورودی $(H \times W)$ وابسته نیست زیرا تمام ورودی‌ها به تمام خروجی‌ها متصل هستند.

شبکه پیچشی: در یک شبکه پیچشی، هر کرنل به صورت محلی عمل می‌کند و تعداد موقعیت‌های ممکن برای جایگذاری کرنل در ورودی به صورت زیر محاسبه می‌شود:

$$(W - K + 1) \times (H - K + 1)$$

است. بنابراین تعداد کل واحدهای خروجی با در نظر گرفتن تعداد فیلترها (I) برابر است با:

$$(W - K + 1) \times (H - K + 1) \times I$$

۲. تعداد اتصالات:

شبکه تمام متصل: در این شبکه، هر واحد ورودی $(H \times W)$ به هر واحد خروجی (I) متصل است و این اتصالات برای تمام کانال‌های ورودی (J) برقرار است. بنابراین تعداد کل اتصالات برابر است با:

$$H \times W \times I \times J$$

شبکه پیچشی: در شبکه پیچشی، هر کرنل $(K \times K)$ در هر موقعیت ممکن از ورودی اعمال می‌شود. تعداد این موقعیت‌ها برابر با تعداد واحدهای خروجی است:

$$(W - K + 1) \times (H - K + 1)$$

هر کرنل برای تمام کانال‌های ورودی (J) اعمال می‌شود و خروجی هر کرنل به تعداد فیلترها (I) افزوده می‌شود. بنابراین تعداد کل اتصالات برابر است با:

$$K \times K \times J \times (W - K + 1) \times (H - K + 1) \times I$$

۳. تعداد وزن‌های قابل یادگیری:

شبکه تمام‌متصل: برای هر اتصال بین واحدهای ورودی و خروجی، یک وزن مستقل وجود دارد. چون تعداد کل اتصالات برابر با $H \times W \times I \times J$ است، تعداد وزن‌های قابل یادگیری نیز دقیقاً برابر با تعداد اتصالات خواهد بود:

$$H \times W \times I \times J$$

شبکه پیچشی: در شبکه پیچشی، هر کرنل $(K \times K)$ برای تمام کانال‌های ورودی (J) و فیلترهای خروجی (I) مشترک است. بنابراین تعداد کل وزن‌های قابل یادگیری برابر است با:

$$K \times K \times J \times I$$

- 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)

- **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 نرمال‌سازی شوند. این کار برای حفظ همبستگی مکانی داده‌ها ضروری است. به همین دلیل، نرمال‌سازی برای هر کانال به صورت مستقل و برای تمامی مکان‌های فضایی به طور مشترک اعمال می‌شود.

نرمال‌سازی دسته‌ای در شبکه‌های پیچشی، پایداری گرادین‌ها را افزایش داده و یادگیری را کارآمدتر می‌کند. این ویژگی به خصوص در شبکه‌های عمیق‌تر اهمیت بیشتری دارد زیرا وابستگی محلی داده‌ها حفظ می‌شود و اختلالی در این ارتباطات ایجاد نمی‌گردد.

- در لایه‌های پیچشی، حذف بایاس مشکلی در عملکرد شبکه ایجاد نمی‌کند. هنگام استفاده از نرمال‌سازی دسته‌ای، خروجی لایه پیچشی ابتدا نرمال می‌شود و سپس دو پارامتر قابل یادگیری β و γ به خروجی اضافه می‌شوند. در این فرآیند، نقش β معادل بایاس است و بایاس قبلی نادیده گرفته می‌شود. از این رو، می‌توان بایاس را بدون تأثیر منفی حذف کرد.

- اگر وزن‌های لایه‌ها (W) یا ورودی‌ها در عددی ثابت مانند α ضرب شوند، این عمل تنها مقادیر اولیه خروجی را تغییر می‌دهد. نرمال‌سازی دسته‌ای با ثابت نگه داشتن میانگین و واریانس مقادیر ورودی، این تغییر را خنثی می‌کند. بنابراین، عملکرد شبکه در طول آزمایش تغییر نخواهد کرد، زیرا لایه نرمال‌سازی این مقیاس را به صورت خودکار تنظیم می‌کند.

- یکی از سناریوهایی که نرمال‌سازی دسته‌ای ممکن است عملکرد ضعیف‌تری داشته باشد، زمانی است که توزیع داده‌های ورودی به مرور زمان تغییر می‌کند. به عنوان مثال، اگر یک شبکه برای شناسایی اشیاء در شرایط روز و شب آموزش دیده باشد و آماره‌های نرمال‌سازی دسته‌ای بر اساس داده‌های روز محاسبه شده باشند، ممکن است در زمان آزمایش روی داده‌های شب، ویژگی‌های مهم تصویر مختوش شوند و دقت کاهش یابد. در چنین مواردی، روش‌های جایگزینی مانند نرمال‌سازی لایه‌ای (Layer Normalization) که به جای مینی‌بچ، نرمال‌سازی را برای هر لایه و نمونه به صورت مستقل انجام می‌دهد، می‌تواند مناسب‌تر باشد. همچنین، روش‌های دیگری از قبیل نرمال‌سازی گروهی (Group Normalization) یا نرمال‌سازی وزن‌ها (Weight Normalization) می‌توانند وابستگی به اندازه دسته را کاهش دهند و عملکرد پایداری در شرایط مختلف ارائه دهند.

(ج)

- ابعاد تصویر ورودی 256×256 است. در هر لایه از انکدر، 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.

- تعداد فیلترهای انکدر به صورت ۵۱۲, ۵۱۲, ۲۵۶, ۱۲۸, ۶۴ است. برای لایه کانولوشن دوم با 3×3 کرنل‌ها:

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

برای لایه کانولوشن دوم انکدر:

$$73856 = 128 + 128 \times (64 \times 3 \times 3) = \text{تعداد پارامترها}$$

بنابراین تعداد پارامترهای این لایه ۷۳۸۵۶ است.

د) شبکه‌های Up-Convolutional: این شبکه‌ها که گاهی به نام Transposed Convolution شناخته می‌شوند، ابزارهایی برای افزایش ابعاد فضایی نقشه‌های ویژگی هستند. در مقاله‌ی "Inverting Visual Representations with Convolutional Networks" (مانند HOG, SIFT, و خروجی لایه‌های مختلف AlexNet) به کار گرفته شده‌اند. در این روش، یک شبکه Up-Convolutional با یادگیری معکوس نگاشت ویژگی‌ها، تصویر اولیه را تخمین می‌زند. این بازسازی‌ها نشان می‌دهد 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 "Apple" یا "Tree") 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 در AlexNet نشان می‌دهد که حتی اطلاعاتی مانند رنگ‌ها از مقادیر پیش‌بینی کلاس‌ها نیز قابل استخراج هستند.

شبکه‌های De-Convolutional: در مقاله‌ی "Striving for Simplicity: The All Convolutional Net"، از شبکه‌های De-Convolutional به‌عنوان 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 موارد زیر است:

۱. شروع از یک فعال‌سازی خاص در لایه بالا و صفر کردن سایر مقادیر.
 ۲. بازگرداندن این فعال‌سازی‌ها به ورودی تصویر، به طوری که نواحی تصویر که بیشترین تأثیر را روی آن نورون دارند، بازسازی شوند.
 ۳. استفاده از روش‌های مختلف برای بازگشت از لایه‌های ReLU و پیچشی مانند Guided Backpropagation که نویز کمتری 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. با این بازسازی‌ها می‌توان 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 را 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.

پاسخ مسئله‌ی ۲.

دقت کنید که می‌توان نوشت:

$$\begin{aligned}\frac{\partial L}{\partial W_j} &= \sum_{i=1}^{N-K} \frac{\partial L}{\partial Z_i} \times \frac{\partial Z_i}{\partial W_j} \\ &= \sum_{i=1}^{N-K} \frac{\partial L}{\partial Z_i} \times \frac{\partial \left(\sum_{k=1}^{K-1} W_k X_{i+k} \right)}{\partial W_j} \\ &= \sum_{i=1}^{N-K} \frac{\partial L}{\partial Z_i} X_{j+i}\end{aligned}$$

بنابراین با توجه به عبارت به دست آمده، می‌توان نتیجه گرفت که $\frac{\partial L}{\partial W_j}$ در واقع برابر مولفه j ام بردار حاصل از اعمال فیلتر کانولوشنی $F = \left[\frac{\partial L}{\partial Z_1} \quad \dots \quad \frac{\partial L}{\partial Z_{N-K}} \right]$ بر روی بردار X می‌باشد. به عبارت دیگر طبق به تعریف بیان شده در صورت سوال، می‌توان نوشت:

$$\frac{\partial L}{\partial W} = F * X \iff \frac{\partial L}{\partial W_j} = \sum_{i=1}^{N-K} F_i X_{j+i}$$

پاسخ مسئله‌ی ۳.

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

لایه اول:

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

لایه دوم:

$$P_{i,j}^{(2)} = \sum_{k=1}^{\P} \sum_{w=\bullet}^{\wedge} \sum_{h=\bullet}^{\wedge} W_{w+\wedge, h+\wedge, k}^{(2)} P_{i+w, j+h, k}^{(1)}$$

لایه سوم:

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

لایه چهارم:

$$h^{(4)} = \begin{bmatrix} P_{\wedge, \wedge}^{(3)} & P_{\wedge, \P}^{(3)} & P_{\P, \wedge}^{(3)} & P_{\P, \P}^{(3)} \end{bmatrix}$$

لایه پنجم:

$$z = \sum_{i=1}^{\P} W_i^{(fc)} h_i^{(4)}$$

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

(الف)

۱. دقت کنید که می توان نوشت:

$$\frac{\partial L}{\partial h_k^{(4)}} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial h_k^{(4)}} = \frac{\partial L}{\partial z} W_k^{(fc)}$$

از طرفی طبق تعریف داریم:

$$P_{i,j}^{(3)} = h_{\P i + j - \P}^{(4)} \quad (1 \leq i, j \leq \P)$$

پس به دست می آید:

$$\frac{\partial h_k^{(4)}}{\partial P_{i,j}^{(3)}} = \mathbb{I}(k = \P i + j - \P)$$

و در نتیجه خواهیم داشت:

$$\begin{aligned} \frac{\partial L}{\partial P_{i,j}^{(3)}} &= \sum_{k=1}^{\P} \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial h_k^{(4)}} \times \frac{\partial h_k^{(4)}}{\partial P_{i,j}^{(3)}} \\ &= \frac{\partial L}{\partial z} W_{\P i + j - \P}^{(fc)} \end{aligned}$$

۲. دقت کنید که می توان نوشت:

$$\frac{\partial L}{P_{i,j,k}^{(\imath)}} = \sum_{m=\imath}^{\mathfrak{f}} \sum_{n=\imath}^{\mathfrak{f}} \frac{\partial L}{P_{m,n}^{(\mathfrak{r})}} \times \frac{\partial P_{m,n}^{(\mathfrak{r})}}{P_{i,j,k}^{(\imath)}}$$

$$\frac{\partial L}{P_{m,n}^{(\mathfrak{r})}} = \sum_{p=\imath}^{\mathfrak{r}} \sum_{q=\imath}^{\mathfrak{r}} \frac{\partial L}{P_{p,q}^{(\mathfrak{r})}} \times \frac{\partial P_{p,q}^{(\mathfrak{r})}}{P_{m,n}^{(\mathfrak{r})}}$$

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

$$\frac{\partial P_{p,q}^{(\mathfrak{r})}}{P_{m,n}^{(\mathfrak{r})}} = \frac{\imath}{\mathfrak{f}} \mathbb{I}((m = \mathfrak{r}p \vee m = \mathfrak{r}p - \imath) \wedge (n = \mathfrak{r}q \vee n = \mathfrak{r}q - \imath))$$

$$\frac{\partial P_{m,n}^{(\mathfrak{r})}}{P_{i,j,k}^{(\imath)}} = W_{m-i,n-j,k}^{(\mathfrak{r})} \mathbb{I}((m = i \vee m = i + \imath) \wedge (n = j \vee n = j + \imath))$$

همچنین طبق قسمت قبل می دانیم:

$$\frac{\partial L}{P_{p,q}^{(\mathfrak{r})}} = \frac{\partial L}{\partial z} W_{\mathfrak{r}p+q-\mathfrak{r}}^{(fc)}$$

بنابراین جواب نهایی به شکل زیر به دست می آید:

$$\begin{aligned} \frac{\partial L}{P_{i,j,k}^{(\imath)}} &= \sum_{m=i}^{i+\imath} \sum_{n=j}^{j+\imath} \frac{\partial L}{P_{m,n}^{(\mathfrak{r})}} W_{m-i,n-j,k}^{(\mathfrak{r})} \\ &= \frac{\imath}{\mathfrak{f}} \sum_{m=i}^{i+\imath} \sum_{n=j}^{j+\imath} \frac{\partial L}{\partial z} W_{\mathfrak{r}[\frac{m}{\mathfrak{r}}] + [\frac{n}{\mathfrak{r}}] - \mathfrak{r}}^{(fc)} W_{m-i,n-j,k}^{(\mathfrak{r})} \end{aligned}$$

ب) دقت کنید که می توان نوشت:

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

حال با جایگذاری $\frac{\partial L}{P_{i,j,k}^{(\imath)}}$ به دست آمده از قسمت قبل، نتیجه می شود:

$$\frac{\partial L}{\partial W_{\imath,\imath,k}^{(\imath)}} = \frac{\imath}{\mathfrak{f}} \sum_{m=i}^{i+\imath} \sum_{n=j}^{j+\imath} \frac{\partial L}{\partial z} W_{\mathfrak{r}[\frac{m}{\mathfrak{r}}] + [\frac{n}{\mathfrak{r}}] - \mathfrak{r}}^{(fc)} W_{m-i,n-j,k}^{(\mathfrak{r})} X_{i,j,k}$$

ج) دقت کنید که می توان نوشت:

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

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

$$P_{m,n}^{(\gamma)} = \sum_{k=1}^{\gamma} \sum_{w=0}^1 \sum_{h=0}^1 W_{w+1,h+1,k}^{(\gamma)} P_{m+w,n+h,k}^{(1)}$$

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

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

و در غیر این صورت داریم:

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

حال با جایگذاری $\frac{\partial P_{m,n}^{(\gamma)}}{\partial W_{i,j,k}^{(\gamma)}}$ و مشتق خواسته شده به شکل زیر محاسبه می شود:

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

پاسخ مسئله‌ی ۴.

(الف)

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

– **Dense Connections در DenseNet**: 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.

– **Residual Connections در ResNet**: در ResNet، اتصال بین لایه‌ها از طریق جمع (Summation) ویژگی‌ها انجام می‌شود. این روش گرادیان‌ها را مستقیماً به لایه‌های قبلی بازمی‌گرداند 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 (Shortcut Connections) به لایه‌های اولیه بازمی‌گردد.
پارامترهای کمتری نیاز دارد زیرا ویژگی‌های تکراری مجدداً یاد 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 موثرتری از مشکل vanishing gradient جلوگیری می‌کند.	احتمال بیشتری برای vanishing gradient دارد (هرچند کاهش یافته است).

جدول ۱: مقایسه بین DenseNet و ResNet

- 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 و بالعکس بدون هیچ ضعیفی 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.

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

$$\text{تعداد لایه‌ها} \times k + \text{ورودی اولیه} = \text{تعداد کانال‌های خروجی}$$

$$\text{کانال} = ۱۰۴ = ۳۲ + ۲۴ \times ۳ = \text{تعداد کانال‌های خروجی}$$

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

(ب)

• حالت اول: کانولوشن عادی (5×5)

- در این حالت، از ۳۲ فیلتر 5×5 استفاده می‌شود.
- هر فیلتر با ورودی $192 \times 28 \times 28$ اعمال شده و خروجی تولید می‌کند.
- تعداد کل عملیات محاسباتی به صورت زیر محاسبه می‌شود:

$$\text{عملیات} = 120422400 = (5 \times 5) \times 32 \times 192 \times 28 \times 28$$

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

- در این حالت، ابتدا کانال‌های ورودی به ۱۶ 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:

$$\text{عملیات} = 2407424 = (1 \times 1) \times 16 \times 192 \times 28 \times 28$$

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

$$\text{عملیات} = 10648576 = (5 \times 5) \times 16 \times 32 \times 192 \times 28 \times 28$$

* جمع کل:

$$\text{عملیات} = 13056000 = 2407424 + 10648576$$

درصد کاهش پیچیدگی:

- کاهش پیچیدگی با استفاده از فرمول زیر محاسبه می‌شود:

$$\frac{120422400 - 13056000}{120422400} \times 100 \approx 89.16\%$$

مشاهده می‌شود که پیچیدگی محاسباتی مدل حدوداً ۹۰٪ کاهش پیدا کرده است که به شدت قابل توجه می‌باشد.

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

$$28 \times 28 \times 192 \times 64 \times (1 \times 1) = 9625600$$

برای فیلتر 1×1 با ۶۴ کانال.

$$28 \times 28 \times 192 \times 96 \times (1 \times 1) = 14438400$$

برای فیلتر 1×1 با ۹۶ کانال.

$$28 \times 28 \times 192 \times 16 \times (1 \times 1) = 2407424$$

برای فیلتر 1×1 با ۱۶ کانال.

$$28 \times 28 \times 192 \times (3 \times 3) = 1354752$$

برای 3×3 Max Pooling.

$$28 \times 28 \times 96 \times 128 \times (3 \times 3) = 27704576$$

برای فیلتر 3×3 با ۱۲۸ کانال.

$$28 \times 28 \times 16 \times 32 \times (5 \times 5) = 10648576$$

برای فیلتر 5×5 با ۳۲ کانال.

$$28 \times 28 \times 192 \times 32 \times (1 \times 1) = 5949440$$

برای فیلتر 1×1 با ۳۲ کانال.

جمع کل عملیات ها:

$$9625600 + 14438400 + 2407424 + 1354752$$

$$+ 27704576 + 10648576 + 5949440 = 72128368 \text{ عملیات}$$

دلایل استفاده از فیلترها با اندازه های مختلف :

- پردازش چند مقیاسی: 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.

- معماری GoogLeNet با معرفی طبقه‌بندهای کمکی (Auxiliary Classifiers) تلاش می‌کند تا مشکل vanishing gradient را کاهش داده و پایداری بهینه‌سازی را بهبود بخشد. این طبقه‌بندها در لایه‌های میانی شبکه، مانند Inception (4a) و Inception (4d)، قرار گرفته‌اند. وظیفه آن‌ها کمک به انتشار گرادیان در طول شبکه به سمت لایه‌های پایین‌تر است. این ویژگی باعث می‌شود که لایه‌های اولیه شبکه، که گرادیان ضعیف‌تری دریافت می‌کنند، در فرآیند آموزش تأثیر بیشتری بپذیرند.

طبقه‌بندهای کمکی در مرحله آموزش به عنوان خروجی‌های موقت عمل می‌کنند و از طریق اضافه کردن ضرر (Loss) مرتبط با آن‌ها به ضرر کل شبکه، به یادگیری بهتر شبکه کمک می‌کنند. به این ترتیب، این طبقه‌بندها به عنوان یک نوع regularization عمل می‌کنند و از بیش‌برازش (Overfitting) جلوگیری می‌کنند. اما در زمان تست، این طبقه‌بندها حذف می‌شوند تا عملکرد شبکه تنها بر اساس خروجی نهایی سنجیده شود.

این رویکرد علاوه بر کاهش vanishing gradient، به بهبود جریان گرادیان در شبکه کمک می‌کند و پایداری بهینه‌سازی را افزایش می‌دهد. با توجه به عمیق بودن معماری GoogLeNet، این ویژگی باعث می‌شود که مدل بتواند اطلاعات مفیدتری از داده‌های آموزشی استخراج کند و ویژگی‌های پیچیده‌تر را با دقت بالاتری یاد بگیرد.

با این حال، محدودیت‌های این طراحی نیز قابل توجه است. تنظیم دقیق ضرر طبقه‌بندهای کمکی و وزن‌دهی آن‌ها در یادگیری شبکه نیازمند تلاش زیادی است. همچنین، در معماری‌های بسیار عمیق‌تر یا شبکه‌های مدرن‌تر، ممکن است این روش کارایی مشابهی نداشته باشد. به همین دلیل، استفاده از این تکنیک در برخی کاربردهای جدیدتر یادگیری عمیق محدود شده است.

ج

- شبکه‌های کانولوشن معمولی از یک ساختار شبکه نمونه‌گیری ثابت استفاده می‌کنند. به عنوان مثال، در یک هسته 3×3 ، 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 به صورت 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:

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

معادله کانولوشن Deformable با اضافه کردن آفست‌ها:

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

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 به صورت محلی و مستقل یاد گرفته می‌شوند.

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

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

که در آن $G(q, p)$ کرنل درونیابی 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.

پاسخ مسئله‌ی ۵.

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

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

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

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

که در آن:

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

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

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) 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 می باشد. همچنین فرض می کنیم که شماره گذاری سطر ها و ستون های ماتریس از صفر شروع می شوند.

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

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

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

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

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

$$K' = K \otimes A$$

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 این همان خواسته مسئله است.

ج) در ابتدا متغیر r_l را برابر تعداد پیکسل های قابل مشاهده توسط عنصر (i, j) خروجی در لایه l ام شبکه قرار دهید. دقت کنید که به وضوح و با توجه به تعریف، $r_3 = 1$ می باشد. در این قسمت، هدف محاسبه r (لایه صفرم همان لایه ورودی شبکه است) می باشد.

حال به این نکته توجه کنید که طبق نتیجه به دست آمده در قسمت قبل، اعمال کانولوشن گسترش یافته با فیلتر $K_{F \times F}$ و گام D ، معادل اعمال کانولوشن عادی با فیلتر $K'_{(FD-D+1) \times (FD-D+1)}$ است. به همین دلیل، از اینجا به بعد فرض می کنیم که تمامی عملیات های کانولوشن موجود در شبکه عادی هستند و با این فرض مسئله را حل می کنیم. به وضوح رابطه بازگشتی زیر برای r_l ها برقرار است:

$$r_l = r_{l+1} + (F_l D_l - D_l + 1) - 1 = r_{l+1} + D_l (F_l - 1)$$

چرا که با اعمال کانولوشن و عبور از لایه l به لایه $l + 1$ ، تعداد پیکسل های قابل مشاهده توسط عنصر خروجی به میزان "بعد فیلتر منهای یک" عدد افزایش می یابد (منهای یک به منظور در نظر نگرفتن خود پیکسل است). حال با توجه به رابطه بازگشتی بالا و مقادیر زیر، r را به دست می آوریم:

$$F_0 = w - 1, \quad F_1 = w, \quad F_2 = w + 1 \quad D_0 = d - 1, \quad D_1 = d, \quad D_2 = d + 1$$

خواهیم داشت:

$$\begin{aligned} r_3 = 1 &\implies r_2 = 1 + w(d + 1) \\ &\implies r_1 = 1 + w(d + 1) + (w - 1)d \\ &\implies r_0 = 1 + w(d + 1) + (w - 1)d + (w - 2)(d - 1) \end{aligned}$$

و این همان خواسته مسئله است.

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

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

محدودیت‌ها:

۱. **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.

۲. **عمق زیاد شبکه:** برای افزایش دامنه دید، نیاز به تعداد زیادی لایه وجود دارد که این موضوع می‌تواند محاسبات را پیچیده کرده و مشکلاتی مانند 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:

۱. **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.

۲. **حفظ ساختار محلی و جلوگیری از دسترسی غیرمجاز:** ترکیب کانولوشن ماسک شده با کانولوشن گسترش یافته، امکان حفظ دقت 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.

۳. **کاهش عمق شبکه و پیچیدگی محاسباتی:** با گسترش دامنه دید در هر لایه از طریق کانولوشن گسترش یافته، نیاز به افزودن لایه‌های متعدد برای افزایش دامنه دید کاهش می‌یابد. این امر منجر به کاهش پیچیدگی محاسباتی و افزایش کارایی مدل می‌شود.

۴. **افزایش انعطاف‌پذیری مدل:** کانولوشن گسترش یافته به مدل اجازه می‌دهد تا تعادلی بهینه بین دقت محلی و اطلاعات کلی برقرار کند، که این امر برای درک بهتر ویژگی‌های مختلف تصویر ضروری است.

پاسخ مسئله‌ی ۶.

(الف) در ابتدا به این نکته دقت کنید که چون مشتق تابع ReLU تنها دو حالت دارد:

$$f'(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

هرگاه حتی در یک مرحله از لایه‌های بعد از i مشتق صفر باشد، حاصل ضرب کل به صفر منجر می‌شود. به منظور مشاهده‌ی بیشترین مقدار ممکن برای $|\delta_i|$ در شرایط ایده‌آل فرض می‌کنیم که در تمام مراحل مشتق ReLU برابر ۱ است. در این حالت:

$$\begin{aligned} |\delta_i| = \|\delta_i\|_2 &= \left\| \delta_L \left(\prod_{k=i+1}^L W_k \right) \right\|_2 \\ &\leq \|\delta_L\|_2 \prod_{k=i+1}^L \|W_k\|_2 \\ &= |\delta_L| \prod_{k=i+1}^L \sigma_{\max}(W_k) \end{aligned}$$

در نتیجه خواهیم داشت:

$$\begin{aligned} \lim_{L \rightarrow \infty} |\delta_i| &\leq \lim_{L \rightarrow \infty} |\delta_L| \prod_{k=i+1}^L \sigma_{\max}(W_k) \\ &= |\delta_L| \lim_{L \rightarrow \infty} \prod_{k=i+1}^L \sigma_{\max}(W_k) \end{aligned}$$

از آنجا که طبق فرض مسئله $\sigma_{\max}(W_k)$ (بزرگ‌ترین مقدار تکین ماتریس W_k) برای هر k کوچک‌تر از یک است، زمانی که $L \rightarrow \infty$ عبارت $\prod_{k=i+1}^L \sigma_{\max}(W_k)$ برابر حاصلضرب بی‌شمار عدد کوچک‌تر از یک می‌شود و به صفر میل می‌کند و در نتیجه خواهیم داشت:

$$\lim_{L \rightarrow \infty} \prod_{k=i+1}^L \sigma_{\max}(W_k) = 0 \implies \lim_{L \rightarrow \infty} |\delta_i| = 0$$

با توجه به نتیجه فوق و برای مثال تا ابتدای شبکه، داریم:

$$|\delta_0| \leq |\delta_L| \prod_{k=1}^{100} \|W_k\| \leq |\delta_L| (0.9)^{100}$$

با توجه به اینکه $(0.9)^{100} \approx 2.6561 \times 10^{-5}$ (تقریباً بسیار کوچک)، می‌توان دید که $|\delta_0|$ بسیار کوچک شده و گرادین در لایه‌های اولیه شبکه تقریباً به صفر نزدیک می‌شود. این امر، پدیده‌ی vanishing gradient را به خوبی نشان می‌دهد.

ب) برای مقداردهی اولیه (Initialization) از روش Kaiming Initialization یا He Initialization استفاده می‌کنیم. در این روش، برای یک توزیع نرمال، واریانس وزن‌ها برابر است با:

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

که n_{in} تعداد ورودی‌های موثر بر هر Kernel است. اگر فرض کنیم لایه کانولوشن ما دارای کانال ورودی ۳ و کرنل با ابعاد 8×8 باشد، آنگاه:

$$n_{\text{in}} = 3 \times 8 \times 8 = 192.$$

در نتیجه:

$$\text{Var}(w) = \frac{2}{192} = \frac{1}{96} \approx 0.0104$$

انحراف معیار σ به صورت:

$$\sigma = \sqrt{\text{Var}(w)} = \sqrt{0.0104} \approx 0.102$$

پس کرنل‌های این لایه از توزیع نرمال:

$$w \sim \mathcal{N}(0, 0.102^2)$$

نمونه‌برداری می‌شوند.

ج) مقایسه‌ی ReLU و Sigmoid از منظر گرادیان:

- تابع Sigmoid:

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

برای ورودی‌های بزرگ مثبت، $\text{Sigmoid}(x) \approx 1$ و مشتق آن بسیار کوچک می‌شود. همچنین برای ورودی‌های بزرگ منفی $\text{Sigmoid}(x) \approx 0$ و مشتق آن باز هم ناچیز است. بنابراین گرادیان در لایه‌های ابتدایی شبکه که ورودی‌شان به نواحی اشباع منتقل می‌شود، بسیار کوچک خواهد شد:

$$\frac{d}{dx} \text{Sigmoid}(x) = \text{Sigmoid}(x)(1 - \text{Sigmoid}(x))$$

که در حالت اشباع به مقدار نزدیک صفر میل می‌کند.

- تابع ReLU:

$$\text{ReLU}(x) = \max(0, x)$$

در نواحی مثبت، مشتق ReLU ثابت و برابر ۱ است:

$$\frac{d}{dx} \text{ReLU}(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

این ویژگی مانع از محوشدن گرادیان می‌شود و به یادگیری سریع‌تر و مؤثرتر شبکه‌های عمیق کمک می‌کند.

د) مشکل Dead Neurons (نورون‌های مرده) در ReLU وقتی رخ می‌دهد که ورودی نورون همواره منفی باشد، به طوری که خروجی ReLU آن همیشه ۰ باقی مانده و گرادیانی برای به‌روزرسانی وزن‌ها دریافت نکند. در نتیجه نورون برای همیشه غیرفعال می‌شود.

توابع Leaky ReLU، PReLU (Parametric ReLU) و ELU جهت رفع این مشکل معرفی شده‌اند:

- Leaky ReLU:

$$f(x) = \begin{cases} x & x > 0 \\ \alpha x & x \leq 0 \end{cases}$$

α یک مقدار کوچک ثابت (مثلاً ۰/۰۱) است. این تابع حتی در ناحیه منفی، خروجی منفی غیرصفر تولید می‌کند و مشتق آن:

$$f'(x) = \begin{cases} 1 & x > 0 \\ \alpha & x \leq 0 \end{cases}$$

- (Parametric ReLU) PReLU:

$$f(x) = \begin{cases} x & x > 0 \\ ax & x \leq 0 \end{cases}$$

که a پارامتری قابل یادگیری است و مشتق:

$$f'(x) = \begin{cases} 1 & x > 0 \\ a & x \leq 0 \end{cases}$$

این تابع انعطاف‌پذیری بیشتری نسبت به Leaky ReLU دارد.

- (Exponential Linear Unit) ELU:

$$f(x) = \begin{cases} x & x > 0 \\ \alpha(e^x - 1) & x \leq 0 \end{cases}$$

مشتق آن:

$$f'(x) = \begin{cases} 1 & x > 0 \\ \alpha e^x & x \leq 0 \end{cases}$$

این توابع با ایجاد خروجی‌های منفی و در نتیجه گرادیان غیرصفر در ناحیه منفی، مانع از مردن نورون‌ها می‌شوند.

البته دقت کنید که این توابع هر یک مزایا و معایبی دارند:

- Leaky ReLU: پیاده‌سازی ساده ولی α ثابت ممکن است بهینه نباشد.
- PReLU: انعطاف‌پذیری بالاتر با قابلیت یادگیری α ، ولی افزایش پارامترها ممکن است خطر overfitting را بالا ببرد.
- ELU: همگرایی سریع‌تر و میانگین خروجی نزدیک به صفر، اما محاسبات نسبت به ReLU پیچیده‌تر و هزینه‌برتر است.

ه) در شبکه‌های دارای بلاک‌های باقی‌مانده (Residual Blocks)، خروجی بلاک به صورت:

$$y_l = F(x_l) + x_l$$

تعریف می‌شود. در هنگام محاسبه گرادیان تابع هزینه J نسبت به ورودی بلاک:

$$\frac{\partial J}{\partial x_l} = \frac{\partial J}{\partial y_l} (F'(x_l) + 1)$$

وجود عبارت $(+1)$ که از مسیر میانبر (Skip Connection) ناشی می‌شود، یک جریان گرادیان مستقیم را فراهم می‌کند. حتی اگر $F'(x_l)$ کوچک باشد، عبارت $+1$ مانع از نابودی گرادیان می‌شود و به حفظ مقدار گرادیان در لایه‌های اولیه کمک می‌کند.

در شبکه‌های بدون Residual Blocks:

$$y_l = F(x_l) \implies \frac{\partial J}{\partial x_l} = \frac{\partial J}{\partial y_l} F'(x_l)$$

اگر $F'(x_l)$ کوچک باشد (مثلاً در شبکه‌هایی که از Sigmoid استفاده می‌کنند و ورودی در ناحیه اشباع است)، گرادیان به سرعت در عمق شبکه از بین می‌رود (پدیده vanishing gradient). اما در صورت وجود Residual Blocks، با توجه به جمله $+1$ ، نحوه‌ی انتقال گرادیان به گونه‌ای تغییر می‌یابد که از محو شدن زودهنگام آن جلوگیری شود و شبکه‌های بسیار عمیق بتوانند مؤثرتر آموزش ببینند.