



نام و نام خانوادگی:

رضا قربانی پاچی

شماره دانشجویی:

403206565

تمرین چهارم درس یادگیری ماشین

د سوال ۱

for convolution: $h_{\text{total}} = (h_{\text{conv}} \times \text{filter})$
 $= (k \times k \times k \times 1) \times h_{\text{conv}}$
 $h_{\text{conv}} = 16 \times 16 \times 16$

حسبه های dense تعداد یا مترقی سویی وای آموزش حسبه به CNN نیاز دارند و همچنین حسبه های CNN می توانند حسبه های داده های مورد نیاز وای آموزش dense حسبه به CNN به قدرت حسبه باشند. اما با توجه به وظیفه ای که داریم جواب این سوال همیشه این است که CNN همیشه با تعداد داده کمتری حسبه به dense عمل می کند زیرا حسبه های CNN حسبه های داده را به صورت خودکار استخراج می کنند.

(ب)

2- در بسیاری از بیابان‌های همگامی که از ۵۰۰۰ سال پیش می‌گویند، برای بار بار این دریا را می‌توان مشاهده کرد. زیرا این دریا در گذشته دریا بوده و اکنون خشک شده است. این دریا را می‌توان در استان خراسان مشاهده کرد.

3- در مرحله استنتاج، یا از اهرهای میانی و در آنش در زمان آهسته ثابت هستند. در این مرحله مراکز تبدیل از اهری ثابت به دوری دگرگونی‌های درونی انجام می‌دهند. این تبدیل ~~بهره~~ از جزیی می‌توان با فرض‌های کمپلوس ادا کرد و در نتیجه هیچ برابر با حسابی امانه و ~~است~~ استثنائی در عملکرد ایجاد می‌شود. در عمل، استفاده از مرحله استنتاج همچنانکه در یادگیری می‌گردد. زیرا فرجه‌ی هسته دست به تغییرات انکساری (توزیع و درجی) می‌دهد و به همین دلیل در

4- در طراحی هاست در دسترس بودن داده‌های مهم، اندازه mini-batch های بسیار کوچک، یا نرخ‌های یادگیری بالا، مشکل (همه) ، کارایی مورد انتظار را کاهش می‌دهد یا حتی نتیجه معکوس دهد، روش‌های دیگری وجود دارند.

تکنیک‌های تغییر: Layer norm, weight normalization, Group Normalization, Instance Normalization و

Batch Renorm---

(ع)

1- در معماری encoder، ابعاد فضایی به نصف کاهش می‌یابد. ساختار U-net، 4 مرحله downsampling در encoder دارد:

input: 256×256

layer 1: 128×128

layer 2: 64×64

layer 3: 32×32

layer 4: $16 \times 16 \rightarrow$ در این مرحله ابعاد فضایی به نصف کاهش می‌یابد و در آن کارایی می‌شود 16×16 که به یک چهارم می‌رسد

$$\begin{aligned} \text{تعداد کانالها} &= (1 + \text{تعداد کانالها} \times \text{تعداد کانالها} \times \text{تعداد کانالها}) \times \text{تعداد کانالها} \\ &= (1 + 3 \times 3 \times 64) \times 128 = 73,856 \end{aligned}$$

-2

-4

de-convolution (Deconv)

در این روش، از حوضه‌های ویژگی به ابعاد فضایی بزرگ‌تر می‌رویم. شبکه شروع کرده و با استفاده از یک فرآیند وارون سازی (شماره یک) که در این عملیات کانولوشن، غیر فعال سازی (Pooling) و معکوس در فضایی ورودی بازسازی می‌شود. این تصویر نشان می‌دهد که فیلتر با ابعاد دقیقاً معکوس فیلترهای ورودی تصویر اصلی را برمی‌گرداند.

U-convolution (transposed convolution)

در یک معماری Encoder-decoder، در بخش encoder، ویژگی‌ها فشرده شده و بخش encoder با استفاده از لایه‌های U-conv و ویژگی‌های فشرده به ابعاد تصویر ورودی برمی‌گردد. این بازسازی تصویر، به ابعاد تصویر مورد نیاز می‌رسد و معنی می‌دهد که این روش، وگرنه آن‌ها را در ابعاد اصلی تصویر بازسازی می‌کند.

سوال (2)

سوال (2):

$$\frac{\partial L}{\partial \omega_j} = \sum_{i=0}^{N-k} \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial \omega_j}$$

X_{i+j}

$$\rightarrow \sum_{i=0}^{N-k} \frac{\partial L}{\partial z_i} \times \omega_j X_{i+j}$$

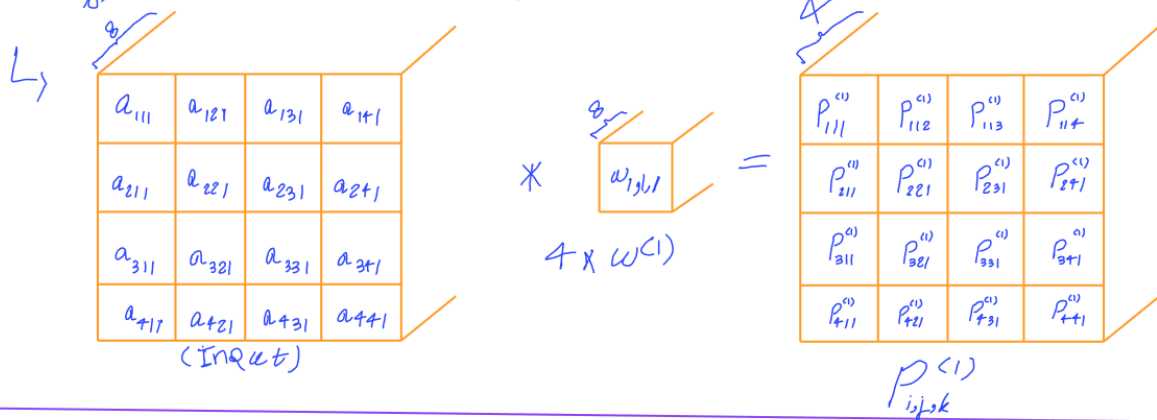
$$\hookrightarrow \frac{dL}{d\omega} = \frac{\partial L}{\partial z} * X$$

در واقع داریم X و مشتق گرادیان داریم

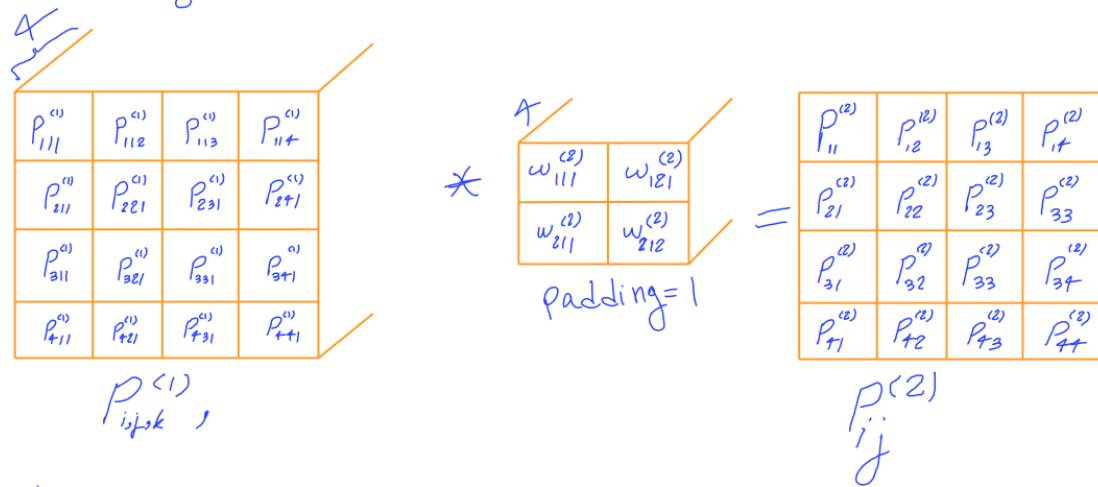
$\frac{\partial L}{\partial z}$ - کارهای مشتق می کنیم

Feed Forward

$$\text{layer 1: } P_{ijk}^{(1)} = \text{conv}(\text{input} \times w^{(1)}) = \sum_{k=1}^8 a_{ijk} \times w_{ijk}$$



$$\text{layer 2: } P_{ij}^{(2)} = \text{conv}(P_{ijk}^{(1)}, w^{(2)})$$

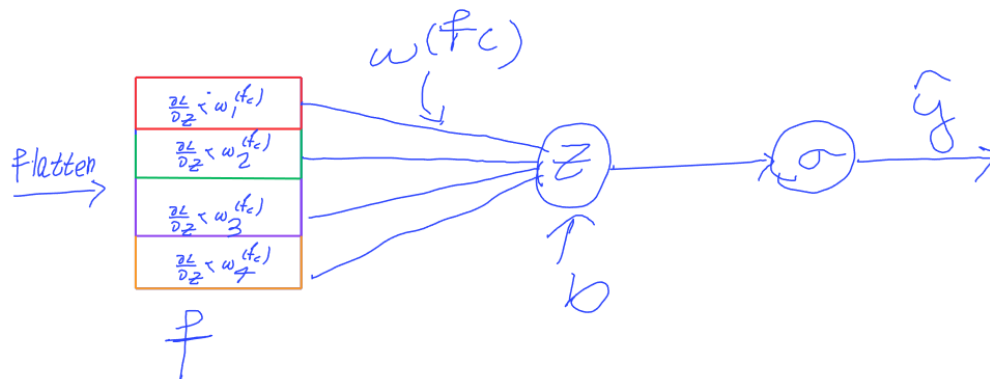


$$P_{ij}^{(2)} = \sum_{j=1}^4 \sum_{k=1}^4 \left[P_{(i,j,k)}^{(1)} \times w_{(1,1,k)}^{(2)} + P_{(i,j,k)}^{(1)} \times w_{(1,2,k)}^{(2)} + P_{(i,j,k)}^{(1)} \times w_{(1,3,k)}^{(2)} + P_{(i,j,k)}^{(1)} \times w_{(1,4,k)}^{(2)} \right]$$

$$\text{layer 3: } P_{ij}^{(3)} = \text{avg-pool}_{2 \times 2}(P_{i,j}^{(2)})$$



fc layer



$$z = f \times (w^{(fc)})^T + b$$

$$\frac{\partial \mathcal{L}}{\partial f} = \frac{\partial \mathcal{L}}{\partial z} \times \frac{\partial z}{\partial p^{(3)}_{i,j}} = \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)} \rightarrow \frac{\partial \mathcal{L}}{\partial p^{(3)}_{i,j}} = \frac{\partial \mathcal{L}}{\partial z} \cdot \text{reshape}(p^{(3)}_{i,j})$$

① (↑)

$$z = (w^{(fc)})^T f + b \rightarrow \frac{\partial z}{\partial f} = w^{(fc)}$$

$$\frac{\partial \mathcal{L}}{\partial f} = \begin{array}{|c|} \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 \\ \hline \end{array} \xrightarrow[\substack{\text{Reshape} \\ p^{(3)}_{i,j}}]{\text{Reshape}} \begin{array}{|c|c|} \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 \\ \hline \end{array}$$

②

$$\text{Padded} \left(\frac{1}{4} \begin{array}{|c|c|c|c|} \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 \\ \hline \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_3 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_4 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_1 & \frac{\partial \mathcal{L}}{\partial z} \times w^{(fc)}_2 \\ \hline \end{array} \right) \xrightarrow[\text{convolution}]{180^\circ \text{ rotate}} \left(\begin{array}{|c|c|} \hline w^{(2)}_{111} & w^{(2)}_{121} \\ \hline w^{(2)}_{211} & w^{(2)}_{212} \\ \hline \end{array} \right)$$

padding=1

$w^{(2)}_{i,j,k}$

$\frac{\partial \mathcal{L}}{\partial p^{(2)}_{i,j}}$

(ب)

$$\frac{\partial L}{\partial w_{i,j,k}^{(1)}} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial p_{ij}^{(3)}} \times \frac{\partial p_{ij}^{(3)}}{\partial p_{ij,k}^{(2)}} \times \frac{\partial p_{ij,k}^{(2)}}{p_{i,j,k}^{(1)}} \times \frac{\partial p_{ij,k}^{(1)}}{\partial w_{ij,k}^{(1)}}$$

در نقش آ (ا) مناسب نیست

$$\frac{\partial L}{\partial w_{i,j,k}^{(1)}} = \text{conv} \left(\frac{\partial L}{\partial p_{i,j,k}^{(1)}}, \text{input} \right)$$

(ج)

$$\frac{\partial L}{\partial w_{i,j,k}^{(2)}} = \text{conv} \left(\frac{\partial L}{\partial p_{i,j,k}^{(2)}}, p_{i,j,k}^{(1)} \right)$$

0	0	0	0	0
$p_{111}^{(1)}$	$p_{112}^{(1)}$	$p_{113}^{(1)}$	$p_{114}^{(1)}$	0
$p_{211}^{(1)}$	$p_{212}^{(1)}$	$p_{213}^{(1)}$	$p_{214}^{(1)}$	0
$p_{311}^{(1)}$	$p_{312}^{(1)}$	$p_{313}^{(1)}$	$p_{314}^{(1)}$	0
$p_{411}^{(1)}$	$p_{412}^{(1)}$	$p_{413}^{(1)}$	$p_{414}^{(1)}$	0

$p^{(1)}, \text{padding}=1$

$\times \frac{1}{4}$

$\frac{\partial L}{\partial p_{111}^{(2)}}$	$\frac{\partial L}{\partial p_{112}^{(2)}}$	$\frac{\partial L}{\partial p_{113}^{(2)}}$	$\frac{\partial L}{\partial p_{114}^{(2)}}$
$\frac{\partial L}{\partial p_{211}^{(2)}}$	$\frac{\partial L}{\partial p_{212}^{(2)}}$	$\frac{\partial L}{\partial p_{213}^{(2)}}$	$\frac{\partial L}{\partial p_{214}^{(2)}}$
$\frac{\partial L}{\partial p_{311}^{(2)}}$	$\frac{\partial L}{\partial p_{312}^{(2)}}$	$\frac{\partial L}{\partial p_{313}^{(2)}}$	$\frac{\partial L}{\partial p_{314}^{(2)}}$
$\frac{\partial L}{\partial p_{411}^{(2)}}$	$\frac{\partial L}{\partial p_{412}^{(2)}}$	$\frac{\partial L}{\partial p_{413}^{(2)}}$	$\frac{\partial L}{\partial p_{414}^{(2)}}$

$\frac{\partial L}{\partial p_{i,j,k}^{(2)}}$

$=$

$\frac{\partial L}{\partial w_{111}^{(2)}}$	$\frac{\partial L}{\partial w_{112}^{(2)}}$
$\frac{\partial L}{\partial w_{211}^{(2)}}$	$\frac{\partial L}{\partial w_{212}^{(2)}}$

سوال 4)

پاسخ بخش (الف):

تفاوت اصلی Residual Connections و Dense Connections :

- در ResNet خروجی هر لایه به لایه بعدی اضافه می شود (skip connection) و تنها یک مسیر مستقیم رو به جلو ایجاد می کند.

- در DenseNet، هر لایه خروجی خود را به تمام لایه های بعدی متصل می کند (به صورت الحاق یا concatenation)، بنابراین هر لایه مجموعه ای از تمام ویژگی های استخراج شده توسط لایه های قبلی را به عنوان ورودی دریافت می کند.

به بیان ساده: در ResNet یک میان بر (short-cut) بین لایه های غیرمتوالی وجود دارد، اما در DenseNet هر لایه به طور متراکم به تمامی لایه های قبل از خود مرتبط است.

چگونگی کاهش مشکل vanishing gradient در DenseNet و مزیت آن:

اتصالات متراکم DenseNet باعث می شوند که گرادیان ها هنگام پس انتشار، بتوانند مسیرهای متعدد و کوتاه تری برای برگشت به لایه های اولیه داشته باشند. این امر سبب می شود گرادیان ها بدون محو شدن، به راحتی به عقب منتقل شوند. در نتیجه:

- یادگیری پایدارتر و عمیق تر صورت می گیرد.

- استفاده کاراتر از پارامترها و بهبود کیفیت ویژگی های استخراج شده امکان پذیر می شود.

محاسبه تعداد کانال خروجی در Dense Block با نرخ رشد: k

در DenseNet، اگر نرخ رشد (growth rate) برابر k باشد، هر لایه k ویژگی جدید تولید کرده و به کانال های موجود می افزاید.

- اگر ورودی بلوک 32 کانال داشته باشد و $k = 24$ باشد، پس از 3 لایه (هر لایه 24 کانال جدید) تعداد کانال های خروجی به صورت زیر است:

$$24 \times 3 + 32 = 32 + 72 = 104$$

پاسخ بخش (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 3، یک لایه کانولوشن به دو صورت نشان داده شده است. ورودی هر دو حالت $192 \times 28 \times 28$ و خروجی هر دو نیز $32 \times 28 \times 28$ 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 (به عنوان پیمانه) با چند درصد کاهش یا افزایش در محاسبات همراه است؟

محاسبات حالت معمولی (بدون پیمانه):

• ابعاد ورودی: $192 \times 28 \times 28$

• فیلتر: 5×5

• تعداد فیلتر خروجی: 32

تعداد عملیات:

$$(28 \times 28) \times 192 \times 32 \times (5 \times 5) = 784 \times 192 \times 32 \times 25 = 120,422,400$$

محاسبات با پیمانه (ابتدا کانولوشن 1×1 برای کاهش کانال):

1. کانولوشن 1×1 از 192 کانال به 16 کانال

تعداد عملیات:

$$16 \times 192 \times (28 \times 28) = 2,408,448$$

2. کانولوشن 5×5 با ورودی 16 کانال و خروجی 32 کانال

تعداد عملیات:

$$(28 \times 28) \times 16 \times 32 \times (5 \times 5) = 784 \times 16 \times 32 \times 25 = 10,035,200$$

مجموع عملیات در حالت پیمانه: $12,443,648 = 10,035,200 + 2,408,448$

مقایسه:

بدون پیمانه: حدود 120 میلیون عملیات

با پیمانه: حدود 12.4 میلیون عملیات

کاهش تقریباً 90٪ در پیچیدگی محاسباتی اتفاق می افتد.

توضیح ایده پیمانه‌های پیدایش (Inception) در GoogleNet

در معماری GoogleNet، ایده به کارگیری پیمانه‌های پیدایش (Inception Modules) برای کاهش پیچیدگی محاسباتی و در عین حال استخراج ویژگی‌های چند-مقیاسی (Multi-scale) از ورودی مطرح شد. در این پیمانه‌ها، به جای استفاده از یک لایه کانولوشن بزرگ با تعداد فیلتر زیاد، از چند مسیر موازی با فیلترهای مختلف (1×1 ، 3×3 ، 5×5) و حتی یک مسیر شامل Max Pooling استفاده می‌شود. هر مسیر روی همان ورودی عمل می‌کند و در نهایت خروجی‌ها در جهت عمق (Channel) با هم الحاق (Concatenate) می‌شوند.

دلایل استفاده از این ترکیب چندمسیری عبارتند از:

استخراج ویژگی‌های چند-مقیاسی: فیلترهای 1×1 ، 3×3 و 5×5 هر کدام محدوده فضایی متفاوتی از ویژگی‌ها را استخراج می‌کنند. فیلترهای بزرگ‌تر (مثلاً 5×5) به بافت‌های گسترده‌تر حساس‌اند، در حالی که فیلترهای کوچک‌تر (1×1 یا 3×3) جزئیات موضعی‌تر را بررسی می‌کنند.

کاهش محاسبات با 1×1 Convolution: قبل از اعمال فیلترهای بزرگ‌تر، ابتدا از کانولوشن 1×1 برای کاهش تعداد کانال‌ها استفاده می‌شود. این کار تعداد ضرب و جمع‌های مورد نیاز برای کانولوشن‌های بعدی (مثلاً 3×3 یا 5×5) را به شدت کاهش می‌دهد.

تنوع در استخراج ویژگی‌ها: استفاده از Max Pooling در کنار کانولوشن‌ها، نوعی ویژگی مجزاست که تنوع ارائه اطلاعات به لایه بعدی را افزایش می‌دهد.

تعداد محاسباتی برای یک پیمانه پیدایش:

فرض کنید ورودی یک پیمانه $28 \times 28 \times 192$ است. این پیمانه چند شاخه (Branch) دارد:

شاخه اول:

کانولوشن 1×1 با 64 فیلتر:

$$28 \times 28 \times 192 \times 64 = 9,633,792$$

شاخه دوم:

ابتدا کانولوشن 1×1 با 96 فیلتر:

$$28 \times 28 \times 192 \times 96 = 14,450,688$$

سپس کانولوشن 3×3 با 128 فیلتر بر خروجی قبلی (ورودی حالا $28 \times 28 \times 96$):

$$28 \times 28 \times 96 \times 128 \times 3 \times 3 = 86,704,128$$

شاخه سوم:

ابتدا کانولوشن 1×1 با 16 فیلتر:

$$28 \times 28 \times 192 \times 16 = 2,408,448$$

سپس کانولوشن 5×5 با 32 فیلتر:

$$28 \times 28 \times 16 \times 32 \times 5 \times 5 = 10,035,200$$

شاخه چهارم (Max Pool و 1×1)

پس از 3×3 Max Pool و سپس کانولوشن 1×1 با 32 فیلتر روی $28 \times 28 \times 192$:

$$28 \times 28 \times 192 \times 32 = 4,816,896$$

هر شاخه پس از کاهش کانال‌ها، اعمال فیلتر بزرگ‌تر را بسیار ارزان‌تر می‌کند.

ترکیب نتایج شاخه‌های مختلف (1×1) ، 3×3 ، 5×5 و Max Pool یک نمایش غنی و چندمقیاسی از ویژگی‌ها را با هزینه محاسباتی کمتر نسبت به استفاده مستقیم از فیلترهای بزرگ به دست می‌دهد.

در نهایت، محاسبات کل یک پیمانه پیدایش با جمع عملیات تمام شاخه‌ها محاسبه شده و معمولاً بسیار کمتر از حالتی است که مستقیماً از کانولوشن‌های بزرگ و بدون کاهش کانال استفاده شود. این تعادل میان تنوع فیلترها و کاهش محاسبات، از موفقیت‌های کلیدی GoogleNet در زمان معرفی آن بوده است.

سوال سوم)

معماری GoogleNet به دلیل استفاده از پیمانه‌های چندمسیره و کانولوشن‌های 1×1 در میان مسیرها، نسبت به vanishing gradient مقاوم‌تر از معماری‌های عمیق پیشین است. افزون بر این، حضور auxiliary classifier ها (خروجی‌های کمکی میانی) کمک می‌کند تا گرادینان از میانه شبکه نیز به عقب منتشر شده و لایه‌های اولیه بهتر آموزش ببینند. این امر جریان گرادینان را بهبود داده و پایداری یادگیری را افزایش می‌دهد. با این حال، در معماری‌ها و تکنیک‌های مدرن که مسئله vanishing gradient کمتر مطرح است، نیاز به auxiliary classifier ها کمتر احساس می‌شود.

پاسخ بخش (ج):

افزودن فرمول‌ها برای تشریح تفاوت کانولوشن عادی و Deformable:

در کانولوشن عادی، خروجی در نقطه (x, y) با استفاده از جمع وزن‌دار مقادیر ورودی از یک گرید منظم محاسبه می‌شود. فرض کنید هسته کانولوشن (Kernel) به اندازه $F \times F$ باشد. فرمول کانولوشن عادی را می‌توان به شکل زیر بیان کرد:

$$G(x, y) = \sum_{(i,j) \in \text{Kernel}} w_{i,j} \cdot F(x + i, y + j)$$

در این معادله، (i, j) محدوده‌ای از اندیس‌ها است که محل نمونه‌برداری هسته کانولوشن از تصویر (یا نقشه ویژگی) F را مشخص می‌کند. این هسته یک گرید منظم است: برای هر موقعیت (x, y) ، نمونه‌برداری دقیقاً $(x + i, y + j)$ انجام می‌شود و هیچ تغییری در مکان این نقاط رخ نمی‌دهد.

اما در کانولوشن تغییرپذیر (Deformable Convolution)، به جای این که از یک گرید منظم ثابت استفاده شود، مکان نمونه‌برداری می‌تواند با استفاده از آفست‌هایی که شبکه در حین آموزش یاد می‌گیرد، تغییر کند. فرمول Deformable Convolution به صورت زیر است:

$$G(x, y) = \sum_{(i,j) \in \text{Kernel}} w_{i,j} \cdot F(x + i + \Delta x_{i,j}, y + j + \Delta y_{i,j})$$

در این رابطه، $\Delta x_{i,j}$ و $\Delta y_{i,j}$ آفست‌هایی هستند که موقعیت نمونه‌برداری را از یک گرید منظم به یک گرید تغییرپذیر تبدیل می‌کنند. این آفست‌ها توسط شبکه (معمولاً یک زیرشبکه کوچک) محاسبه می‌شوند و در زمان آموزش یاد گرفته می‌شوند.

انعطاف‌پذیری نسبت به Geometric Transformation :

شبکه‌های Deformable می‌توانند به طور خودکار الگوهای پیچیده در تصاویر را با اعمال تغییر در محل نمونه‌برداری فیلترها دنبال کنند. مثلاً اگر شیء مورد نظر در تصویر کمی کشیده یا کج شود، یا شکل آن اندکی تغییر کند، فیلترهای Deformable می‌توانند نقاط نمونه‌برداری خود را بازتنظیم کنند تا الگوی مورد نظر را همچنان استخراج نمایند. این خاصیت باعث می‌شود این شبکه‌ها نسبت به تصاویر با تغییرات هندسی غیرصلب (Non-rigid Transformations) انعطاف‌پذیرتر باشند.

مفهوم Offset در Deformable Convolution :

هر آفست از دو مقدار $\Delta x_{i,j}$ و $\Delta y_{i,j}$ تشکیل می‌شود که به مختصات اصلی موقعیت (i, j) هسته کانولوشن اضافه می‌شوند. به عبارت دیگر، به جای نمونه‌برداری از پیکسل $(x + i, y + j)$ در تصویر یا نقشه ویژگی، فیلتر از $(x + i + \Delta x_{i,j}, y + j + \Delta y_{i,j})$ نمونه‌برداری می‌کند. این جابه‌جایی‌های قابل یادگیری، هسته کانولوشن را قادر می‌سازند تا خود را با الگوهای هندسی پیچیده‌تر در تصویر تطبیق دهند.

نحوه محاسبه offset :

آفست‌ها مستقیماً از داده ورودی و از طریق یک زیرشبکه (زیرمدل) کوچک محاسبه می‌شوند. مراحل معمول محاسبه آن‌ها به این صورت است:

1. زیرشبکه تولید آفست:

معمولاً در پیاده‌سازی Deformable Convolution، یک لایه کانولوشن اضافی (با خروجی 2×2 تعداد نقاط هسته) کانال روی همان ویژگی‌های ورودی اعمال می‌شود.

○ اگر هسته کانولوشن $F \times F$ باشد، این زیرشبکه باید $2F^2$ مقدار خروجی دهد (2 مقدار برای هر نقطه نمونه‌برداری: Δx و Δy).

○ این لایه کانولوشن اغلب هسته 3×3 و پدینگ مناسب دارد تا ابعاد فضایی حفظ شده و آفست‌ها برای هر مکان محاسبه شوند.

2. یادگیری در خلال آموزش:

وزن‌های این لایه کانولوشن که آفست‌ها را تولید می‌کند، طی فرایند آموزش با استفاده از backpropagation و بر مبنای معیار خطای نهایی شبکه (مثلاً خطای طبقه‌بندی یا تشخیص) به‌روزرسانی می‌شوند. در نتیجه، شبکه می‌آموزد که کدام جابه‌جایی‌ها (offsets) برای بهبود استخراج ویژگی از داده‌ها مفیدند. به طور تجربی، این آفست‌ها معمولاً یاد می‌گیرند روی مناطق مهم یا ساختارهای خاص در تصویر متمرکز شوند و فیلتر را در جهاتی تغییر دهند که حساسیت به شکل، زاویه، یا مکان اجزا در تصویر بیشتر شود.

سوال 5

پاسخ بخش (الف):

در کانولوشن‌های گسترشی (Dilated Convolution)، به جای نمونه‌برداری متوالی و مجاور پیکسل‌ها در یک ناحیه ثابت (مثلاً $k \times k$)، فاصله‌ای میان خانه‌های نمونه‌برداری وجود دارد. این فاصله با پارامتر گسترش (Dilation) کنترل می‌شود. زمانی که مقدار گسترش $D > 1$ باشد، فیلتر با «پردن» روی پیکسل‌ها (با فاصله‌های منظم) اعمال می‌شود و بنابراین یک «ناحیه مؤثر» بزرگ‌تر را بدون افزایش تعداد پارامترها پوشش می‌دهد.

مفهوم:

- در کانولوشن عادی (بدون گسترش)، اگر هسته یک فیلتر 3×3 داشته باشیم، این هسته مستقیماً روی یک بلوک 3×3 از پیکسل‌ها اعمال می‌شود.
- در کانولوشن گسترشی با ضریب D ، هسته همان 3×3 است، اما میان هر دو پیکسل افقی و عمودی $D-1$ پیکسل جهش وجود دارد. مثلاً با $D=2$ ، فیلتر از مختصات (i,j) ، سپس $(i,j+2)$ ، $(i+2,j)$ ، ... نمونه‌برداری می‌کند. این باعث می‌شود ناحیه‌ای به اندازه بزرگ‌تر $(3+(3-1)*(D-1))$ برای هسته 3×3 (تحت پوشش قرار گیرد، بدون این که تعداد پارامترهای فیلتر افزایش یابد.

مزایا:

- محدوده دید (Receptive Field) شبکه با هزینه کم افزایش می‌یابد. به عبارت دیگر، شبکه می‌تواند اطلاعات فضایی در مقیاس بزرگ‌تر را ببیند، بدون اینکه مجبور باشیم لایه‌های بیشتری اضافه کنیم یا اندازه کرنل را افزایش دهیم.
- برای وظایف بینایی کامپیوتری مانند درک صحنه‌های پیچیده (مثلاً در تقسیم‌بندی معنایی تصاویر)، این امر مفید است زیرا شبکه می‌تواند با همان تعداد پارامتر، بافت و ساختارهای بزرگ‌تری را تحلیل کند.

فرمول ریاضی:

فرمول کانولوشن گسترشی (Dilated Convolution) به صورت زیر است:

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

در این رابطه:

- I ورودی با ابعاد $M \times N$ است.
- K هسته کانولوشن با اندازه $F \times F$ است.
- D پارامتر گسترش است.

- وقتی $D=1$ باشد، همان کانولوشن عادی را داریم.
 - وقتی $D>1$ باشد، فواصل میان پیکسل‌های ورودی که فیلتر روی آن‌ها اعمال می‌شود، افزایش می‌یابد و محدوده مؤثر فیلتر برابر خواهد بود با $(F-1)(D-1)+F$.
 - $Height(rows) = M - (kernel_size - 1) = M - (DF - D + 1 - 1) = M - DF + D$
 - $Width(columns) = N - (kernel_size - 1) = N - (DF - D + 1 - 1) = N - DF + D$
- نتیجه خروجی کانولوشن گسترشی نیز با ابعاد $(M - DF + D)(N - DF + D)$ خواهد بود که نشان می‌دهد با افزایش D ، بدون اضافه کردن پارامترهای بیشتر، می‌توان محدوده فضایی تحت پوشش فیلتر را افزایش داد.

پاسخ بخش (ب):

برای نشان دادن اینکه کانولوشن گسترشی معادل کانولوشن با کرنل متسع شده (K') است، نیاز داریم مفهوم محصول کرونکر (Kronecker Product) را معرفی کنیم. محصول کرونکر یک روش ریاضی برای توسعه ماتریس‌ها است که در اینجا برای گسترش کرنل K استفاده می‌شود.

تعریف Kronecker Product:

اگر k یک ماتریس $n \times m$ باشد و A یک ماتریس $q \times p$ ، Kronecker Product آنها ($A \otimes K$) یک ماتریس با ابعاد $(q \cdot n) \times (p \cdot m)$ است که به صورت زیر تعریف می‌شود:

$$K' = A \otimes K$$

If \mathbf{A} is an $m \times n$ matrix and \mathbf{B} is a $p \times q$ matrix, then the Kronecker product $\mathbf{A} \otimes \mathbf{B}$ is the $pm \times qn$ block matrix:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix},$$

more explicitly:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}b_{11} & a_{11}b_{12} & \cdots & a_{11}b_{1q} & \cdots & \cdots & a_{1n}b_{11} & a_{1n}b_{12} & \cdots & a_{1n}b_{1q} \\ a_{11}b_{21} & a_{11}b_{22} & \cdots & a_{11}b_{2q} & \cdots & \cdots & a_{1n}b_{21} & a_{1n}b_{22} & \cdots & a_{1n}b_{2q} \\ \vdots & \vdots & \ddots & \vdots & & & \vdots & \vdots & \ddots & \vdots \\ a_{11}b_{p1} & a_{11}b_{p2} & \cdots & a_{11}b_{pq} & \cdots & \cdots & a_{1n}b_{p1} & a_{1n}b_{p2} & \cdots & a_{1n}b_{pq} \\ \vdots & \vdots & & \vdots & \ddots & & \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots & & \ddots & \vdots & \vdots & & \vdots \\ a_{m1}b_{11} & a_{m1}b_{12} & \cdots & a_{m1}b_{1q} & \cdots & \cdots & a_{mn}b_{11} & a_{mn}b_{12} & \cdots & a_{mn}b_{1q} \\ a_{m1}b_{21} & a_{m1}b_{22} & \cdots & a_{m1}b_{2q} & \cdots & \cdots & a_{mn}b_{21} & a_{mn}b_{22} & \cdots & a_{mn}b_{2q} \\ \vdots & \vdots & \ddots & \vdots & & & \vdots & \vdots & \ddots & \vdots \\ a_{m1}b_{p1} & a_{m1}b_{p2} & \cdots & a_{m1}b_{pq} & \cdots & \cdots & a_{mn}b_{p1} & a_{mn}b_{p2} & \cdots & a_{mn}b_{pq} \end{bmatrix}$$

Photo from Wikipedia

در اینجا A یک ماتریس است که برای گسترش مکان‌های نمونه‌برداری استفاده می‌شود. اگر A به عنوان ماتریس گسترش با تک 1 در مکان مشخص و بقیه عناصر صفر باشد، محصول کرونکر بین K و A باعث ایجاد یک کرنل جدید K' می‌شود که گسترش یافته است و مکان‌های نمونه‌برداری $D-1$ فاصله دارند.

کرنل گسترش یافته: (K')

برای گسترش کرنل K با استفاده از پارامتر گسترش D :

1. ماتریس A را به صورت یک ماتریس هویت $D \times D$ در نظر بگیرید که فاصله‌های صفر اضافی بین المان‌های اصلی دارد.

2. Kronecker Product $K' = A \otimes K$ اعمال می‌شود، که کرنل گسترشی را ایجاد می‌کند. این کرنل K' حالا فاصله $D-1$ بین نقاط اصلی گرید نمونه‌برداری دارد.

به عنوان مثال A را ماتریسی 2×2 و D را نیز 2 در نظر می‌گیریم حاصل Kronecker Product به صورت زیر می‌شود:

x	y
z	r

 \otimes

1	0
0	0

 $=$

x	0	y	0
0	0	0	0
z	0	r	0
0	0	0	0

پاسخ بخش (ج):

محاسبه محدوده تأثیر: (Receptive Field)

در کانولوشن، هر فیلتر تأثیر خود را بر محدوده‌ای از پیکسل‌های ورودی اعمال می‌کند. این محدوده، که به آن **Receptive Field** گفته می‌شود، وابسته به ابعاد فیلتر ($w \times w$) و میزان گسترش (Dilation Factor) در هر لایه است.

محاسبه R_1 :

$$R_1 = (w - 1) + (w - 1 - 1)(d - 1 - 1)$$

$$R_1 = (w - 1) + (w - 2)(d - 2)$$

$$R_1 = wd - w - 2d + 3$$

محاسبه R_2 :

$$R_2 = R_1 + [w + (w - 1)(d - 1) - 1]$$

محاسبه R_3 :

$$R_3 = R_2 + [(w + 1) + (w + 1 - 1)(d + 1 - 1) - 1]$$

$$R_3 = 3(wd - d + 1)$$

(د)

Masked Convolution چیست؟

Masked Convolution روشی است که در آن وزن‌های فیلتر کانولوشن به صورت انتخابی ماسک (Mask) می‌شوند تا دسترسی به برخی از پیکسل‌های ورودی محدود شود. این روش معمولاً در مسائل مربوط به پیش‌بینی‌های ترتیبی یا داده‌های وابسته به زمان (مانند مدل‌های autoregressive) به کار می‌رود. ماسک‌ها می‌توانند به شکلی طراحی شوند که تضمین کنند فقط اطلاعات از گذشته یا از نقاط خاص ورودی استفاده شود.

کاربردها:

1. مدل‌های Autoregressive:

- در شبکه‌هایی مانند PixelCNN، از Masked Convolution برای تولید تصویر پیکسل به پیکسل استفاده می‌شود. ماسک تضمین می‌کند که پیکسل فعلی فقط به پیکسل‌های قبلی خود وابسته است.

2. مدل‌های زمان‌سری:

- در تحلیل داده‌های زمانی، Masked Convolution می‌تواند تضمین کند که مقدار پیش‌بینی شده در زمان t تنها به داده‌های قبل از t وابسته است.

3. پردازش زبان طبیعی (NLP):

- در معماری‌های مبتنی بر کانولوشن برای NLP، Masked Convolution محدودیت ترتیب زمانی را رعایت می‌کند تا پیش‌بینی فقط به کلمات قبلی متکی باشد.

محدودیت‌های Masked Convolution:

1. کاهش کارایی محاسباتی:

ماسک کردن وزن‌های فیلترها ممکن است باعث افزایش پیچیدگی محاسبات و کاهش بهره‌وری GPU شود.

2. سختی در طراحی ماسک‌ها:

طراحی مناسب ماسک‌ها به صورتی که محدودیت‌های مسئله را رعایت کند، می‌تواند چالش‌برانگیز باشد.

3. کاهش اطلاعات در دسترس:

ماسک کردن می‌تواند اطلاعات مفید موجود در نواحی خاصی از ورودی را محدود کند، که ممکن است منجر به کاهش دقت مدل شود.

بهبود محدودیت‌های Masked Convolution با کانولوشن گسترشی:

کانولوشن گسترشی (Dilated Convolution) می‌تواند این محدودیت‌ها را کاهش دهد:

1. افزایش محدوده دید (Receptive Field):

کانولوشن گسترشی بدون نیاز به ماسک کردن نقاط اضافی، می‌تواند نواحی بزرگ‌تری از ورودی را تحت پوشش قرار دهد، که به مدل اجازه می‌دهد اطلاعات بیشتری را از گذشته استخراج کند.

2. حفظ ترتیب زمانی:

با تنظیم صحیح پارامتر گسترش (Dilation Factor)، می‌توان ترتیب زمانی یا ترتیبی داده‌ها را بدون نیاز به طراحی ماسک پیچیده رعایت کرد.

3. افزایش کارایی محاسباتی:

برخلاف Masked Convolution که مستلزم حذف وزن‌ها و نمونه‌برداری‌های غیرمؤثر است، کانولوشن گسترشی کارایی بیشتری داشته و از عملیات معمول کانولوشن بهره می‌برد.

الف)

$$\delta^{(i)} = \frac{\partial F}{\partial z^{(i)}} = \underbrace{f'(z^{(i)})}_{\text{مشتق فعال}} \underbrace{\omega^{(i+1)} \delta^{(i+1)}}_{\text{محاسبه خروجی}}$$

$$\xrightarrow{\text{جمع}} \delta^{(i)} = \sum_{k=i+1}^L (\omega^{(k)} f'(z^{(k)})) \delta^{(k)}$$

فرض $\left\{ \begin{array}{l} \text{دگرگونی مقدارها} \rightarrow \text{all } \delta \text{ max}(\omega^{(k)}) \\ \text{مشتق فعال} \rightarrow \text{at=Relu} \end{array} \right.$

$$\|\delta^{(i)}\| \leq \sum_{k=i+1}^L \underbrace{\sigma_{\max}(\omega^{(k)})}_{*} \|\delta^{(k)}\|$$

* اگر تعداد کمتر از باشد باز یاد شدن تعداد لایه به منجر می شود یعنی gradient vanishing

مترادف خاص

$$\sigma_{\max}(\omega^{(k)}) = 0.9 \quad L = 100$$

$$\|\delta^{(i)}\| = (0.9)^{100} \|\delta^{(L)}\|$$

$$\rightarrow \sim 2.87 \times 10^{-5}$$

اگر نزدیکترین وزن به افراط باشد
تعدد لایه ها صاف می شود

$$\sigma = \sqrt{\frac{2}{\text{بارهای ورودی}}} \quad \text{بارهای ورودی} = 3 \times 8 \times 8 = 192$$

$$\sigma = \sqrt{\frac{2}{192}} = 0.102 \rightarrow \mathcal{N}(0, \sqrt{\frac{2}{192}})$$

ج) تابع sigmoid به دلیل فرجه می شود (یعنی 0 و 1) و مشتق کوچکی در نواحی اشباع مشکل vanishing gradient ایجاد می کند و یادگیری شبکه های عمیق را دشوار می کند تابع ReLU با داشتن گرادینت ثابت یکی مقادیر مثبت و نبود ناحیه اشباع، این مشکل را برطرف نموده و یادگیری کارآمدتری را در شبکه های عمیق ممکن می سازد

(7) مشکل Dead neurons زمانی رخ میدهد که ورودن یک نورون هیچ سیگنالی نمی‌دهد و محدودیات آن صفر می‌شود.
در نتیجه نورون در محول آموزش غیر فعال باقی می‌ماند و هیچ گاه بروز سیگنالی نمی‌شود.

توانع جلوگیری

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

اجازه می‌دهد که ادیان در نواحی منفی نیز منفی، صفر نشود. بنابراین از مشکل dead neurons جلوگیری می‌کنند.
اسباب و چالش‌ها دارند ممکن است بی‌عملکردی را تأثیر بگذارد.

- Parametric ReLU \rightarrow مشابه قبلی است اما α در $x \leq 0$ در x قابل یادگیری است.

امکانات دیگری نسبت به Leaky ReLU و توانایی تعیین مقدار α با داده‌ها.

- ELU (Exponential Linear Unit)

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

کاهش میانگین خروجی لایه‌ها و بهبود سرعت همگرایی، α که ادیان در نواحی منفی صفر نمی‌شود.
اما پیچیدگی محاسباتی بیشتر نسبت به ReLU و Leaky ReLU دارد.

$$x^{(i+1)} = F(x^{(i)}) \xrightarrow{\text{فرونی مدی RB}} x^{(i+1)} = F(x^{(i)}) + x^{(i)} \rightarrow \text{فرونی با RB} \quad (5)$$

مزایای شبکه‌های عمیق:

- 1- جلوگیری از vanishing gradient: زیرا اتصال مستقیم فرونی باعث می‌شود این مشکل را نداشته باشیم
- 2- بدلیل داشتن مشکل vanishing gradient امکان این وجود دارد که شبکه‌های عمیق تر داشته باشیم
- 3- شبکه می‌تواند سریعتر همگرا شود

* در این شبکه‌ها، skip connections اجازه می‌دهد گره‌های بدون تغییر یا تضعیف شدید به لایه‌های قبلی منتقل شوند. این امر یادگیری شبکه‌های بسیار عمیق را امکان پذیر می‌کند