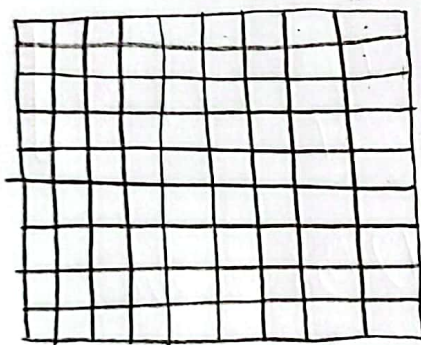
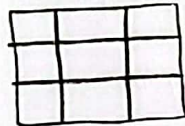
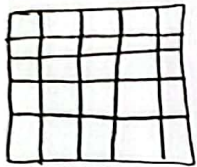


نقشه تئوری درس یادگیری عمیق

سؤال 1 الف

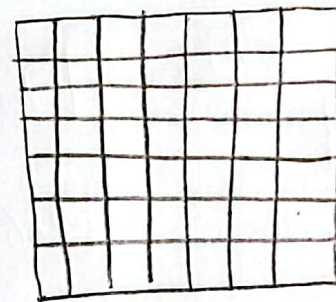
این سؤال مربوط به بحث receptive field می باشد. برای حل این سؤال داریم: در واقع راه حل کلی این است که می گوئیم 1 پیکسل در 1 محاذ 3x3 پیکسل در 1 است و... داریم

یک توریون در لایه چهارم $\leftarrow 3 \times 3$ پیکسل در لایه سوم $\leftarrow 5 \times 5$ پیکسل در لایه دوم



9x9 pixel
در لایه ی صفر

7x7
pixel
لایه اول



بنابر این 81 pixel ✓

سؤال 2 ب

لایه های pooling ایجاد feature map ها را کاهش می دهند و در نتیجه پارامترهای شبکه و محاسبات را کاهش می دهد. همچنین لایه های pooling ویژگی های موجود در یک منطقه از feature map ایجاد شده توسط یک لایه کانولوشن را خلاصه می کند و در نتیجه عملیات بیشتر بر روی ویژگی های خلاصه شده به جای ویژگی های دقیقاً موقعیت یافته (تولید شده توسط لایه کانولوشن) انجام می شود. این باعث می شود که مدل نسبت به تغییرات در موقعیت ویژگی ها در تصویر ورودی مقاوم تر باشد

ج. سائز تصویر ورودی 35×35 و تعداد کانال ها 16
 سائز فیلتر مورد استفاده 3×3 و طول گام 1 و تعداد فیلتر ها برابر 32
 نوع: padding: same

به ازای هر بار اعمال فیلتر یک تصویر 35×35 داریم. اگر 32 فیلتر اعمال کنیم سائز خروجی 8
 $35 \times 35 \times 32$ می شود.

پارامترها $\rightarrow 35 \times 3 \times 3 \times 16 \times 32 + 32 = 4640$
 وزن ها بالاس

د. Input $[227 \times 227 \times 3]$ $\xrightarrow{\text{conv1}}$ $[227 \times 227 \times 256]$ $\xrightarrow{\text{padding: same}}$ output 1
 می دانیم padding: same است داریم:
 پارامتر: $256 \times 5 \times 5 \times 3 + 256 = 19456$

output 1 $[227 \times 227 \times 256]$ $\xrightarrow{\text{conv2}}$ $[227 \times 227 \times 128]$ $\rightarrow \text{output 2}$
 پارامتر: $128 \times 5 \times 5 \times 256 + 128 = 819328$

output 2 $[227 \times 227 \times 128]$ $\xrightarrow{\text{MaxPool}}$ $\frac{227-5}{2} + 1 = 112 \Rightarrow 112 \times 112 \times 128 \rightarrow \text{out 3}$
 پارامتر = 0

Flatten \rightarrow ورودی شده: $112 \times 112 \times 128 = 1605632$

fully connected \rightarrow $512 \rightarrow \text{out 4}$
 $112 \times 112 \times 128 \times 512 + 512 = 822084096$

out 4 512 \rightarrow $100 \rightarrow \text{out 5}$
 $512 \times 100 + 100 = 51300$

هست پارامترها = 822974180

تعداد پارامترهای این شبکه بسیار بالا است. برای آن که پارامترها را کمتر کنیم می توانیم از روش های مختلفی استفاده نماییم. می توانیم از $stride$ های بزرگ تر نیز استفاده کنیم. یکی از راه های ساده این است که $pooling$ را تغییر دهیم. از دیگر مشکلات آن می توان $padding$ به دلیل $stride$ 2×2 در $pooling$ اشاره کرد. در واقع پارامترهای شبکه هنگامی تبدیل $flatten$ زیادی شود و باید به آن جای جای کمی می توان با پارامتری رسید. همچنین در $conv2$ نیز پارامتر زیادی تولید می شود. بهتر است آن چهارم سعی کنیم پارامتر آن کم شود. راه آسان عوض کردن ابعاد $conv$ است که باعث کاهش پارامتری شود اما راه جالب تری می تواند آن باشد که جای $conv1$ و $conv2$ عوض شود داریم:

$$\begin{array}{l} \text{Input} \\ [227 \times 227 \times 3] \end{array} \xrightarrow{\text{conv2}} \begin{cases} \text{out1} : [227 \times 227 \times 128] \\ \text{پارامتر} : 128 \times 5 \times 5 \times 3 + 128 = 9728 \end{cases}$$

$$\begin{array}{l} \text{out1} \\ [227 \times 227 \times 128] \end{array} \xrightarrow{\text{conv1}} \begin{cases} \text{out2} : [227 \times 227 \times 256] \\ \text{پارامتر} : 256 \times 5 \times 5 \times 128 + 256 = 819456 \end{cases}$$

$$\text{out2} \xrightarrow{\text{maxpool } 7 \times 7, \text{ stride } 5 \times 5} \frac{227-7}{5} + 1 = 45 \Rightarrow 45 \times 45 \times 256 \quad \text{out3}$$

$$\text{flatten} \xrightarrow{\text{ورودی شبکه}} 45 \times 45 \times 256 = 1163520$$

$$\xrightarrow{\text{FC}} \begin{cases} \text{out4} = 512 \\ \text{پارامتر} : 512 \times 1163520 + 512 = 595722752 \end{cases}$$

$$\longrightarrow \begin{cases} \text{out5} = 100 \\ \text{پارامتر} = 100 \times 512 + 100 = 51300 \end{cases}$$

$$\text{پارامتر} = 597766756$$

۵ افزایش لایه ها در واقع می توان خاصیت غیر خطی بودن را بیشتر ایجاد کرد
 - اضافه کردن نودون با همان 2 لایه زیرا $f(g(x))$ از $f(x) + g(x)$ خاصیت
 غیر خطی بیشتری ایجاد می کند. بنابراین علی الرغم آن که لایه خاصیت $general$ دارد
 ولی بیشتر دوست داریم تا شبکه عمیق و لاغر داشته باشیم تا عمیق و چاق. مورد
 دیگری که هست از نظر تعداد نودون می باشد. هنگامی که عمیق می شویم به تعداد
 نودون کمتری نسبت به حالتی که چاق می شویم نیاز داریم پس اگر عمیق شویم تعداد
 نودون کمتر و در نتیجه پارامتر کمتر نیاز خواهیم داشت. بنابراین در اکثر مواقع از بیش
 از 2 لایه استفاده می کنیم. (در واقع با عمیق شدن می توان نودون کمتری به ازای پیچیدگی شبکه
 داشت)

۹ برای آن که $vanishing$ را توضیح دهیم می دانیم با اجرای الگوریتم $Error back propagation$
 ممکن است به دلیل کوچکی گرادیان ها هنگامی که لایه های اول می رسم گرادیان کلی برای آپدیت
 آن قدر کم باشد که عملاً آپدیتی صورت نگیرد و در نتیجه شبکه به صورت خیلی خیلی کند
 آپدیت شود و رسیدن به شبکه بهینه سخت شود. این مشکل ناشی از بعضی توابع فعاله
 مثل $sigmoid$ و $tanh$ است و اگر در ناحیهی اشباع قرار بگیرند مشتقشان بسیار به صفر
 نزدیک خواهد بود و اگر چندین لایه این ها در هم ضرب شود (که در EP در لایه های اولیه
 اتفاق می افتد) گرادیان کلی برای آپدیت صفر و شبکه دیگر آپدیت نمی شود در $exploding$
 برخلاف $vanishing$ گرادیان زیاد و زیاد می شود و عملاً باعث واکرایا EP می شود
 آن چه که بین $vanishing$ و $exploding$ تفاوت دارد است که $exploding$ برخلاف
 بر اثر وزن هادر شبکه اتفاق می افتد و نه بر اساس توابع فعال سازی. وزن های لایه های
 پایین تر بیشتر ممکن است برایشان $exploding$ رخ دهد چون گرادیان ها ضرب در اعداد بیشتری می شوند
 پایین تر بیشتر می شود که لایه های پایین به سمت آپادیری می روند. ما می توانیم از این که
 این می تواند باعث شود که لایه های $update$ به $update$ دیگر زیاد است
 مقدار $learning$ داریم یا میسر در $data$ از یک $update$ به $update$ دیگر زیاد است

without zero pad

$$A = \begin{bmatrix} W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 & 0 & 0 & 0 \\ 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 & 0 & 0 \\ 0 & 0 & 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 \\ 0 & 0 & 0 & 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} \end{bmatrix}$$

الف

سؤال 9

in full size

↑
with zero
pad

A =

$$\begin{bmatrix} W_{00} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ W_{01} & W_{00} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & W_{01} & W_{00} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & W_{01} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ W_{10} & 0 & 0 & W_{00} & 0 & 0 & 0 & 0 & 0 & 0 \\ W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 & 0 & 0 & 0 & 0 \\ 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 & 0 & 0 & 0 \\ 0 & 0 & W_{11} & 0 & 0 & W_{01} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & W_{10} & 0 & 0 & W_{00} & 0 & 0 & 0 \\ 0 & 0 & 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 & 0 \\ 0 & 0 & 0 & 0 & W_{11} & W_{10} & 0 & W_{01} & W_{00} & 0 \\ 0 & 0 & 0 & 0 & 0 & W_{11} & 0 & 0 & W_{01} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & W_{10} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & W_{11} & W_{10} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & W_{11} & W_{10} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & W_{11} & W_{10} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & W_{11} \end{bmatrix}$$

سؤال 2

ب

این نوع کانولوشن در واقع در جاهایی استفاده می شود که تعدادی از پیکسل های منطقی از بین رفته است مثل آثار هنری. در این کانولوشن ماسک های دید شده انجام می شود. رابطه ی گون به صورت زیر می باشد

$$X' = W (X \otimes M)$$

یک *spatial separable conv* به صورت ساده یک کرنل را به دو کرنل کوچک تر تقسیم می کند. یک کس معروف این است که یک کرنل 3×3 را به یک 1×3 و یک 3×1 تقسیم کنیم مثل زیر:

$$\begin{bmatrix} 3 & 6 & 9 \\ 4 & 8 & 12 \\ 5 & 10 & 15 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

در واقع به جای انجام یک کانولوشن با ضرب 2 کانولوشن هر یک با 3 ضرب را انجام می دهیم (مجموعه ضرب) تا به اثر یکسان برسیم. با تعداد کمتر ضرب پیچیدگی محاسباتی پایین و شبکه تقابلی این را پیدا خواهد کرد که سریع تر از آن شود.