

به نام خدا

یادگیری عمیق، تکلیف دوم

مهدی کافی ۹۹۲۱۰۷۵۳

## مسئله ۱. Network Design

(آ) فرض کنید که شبکه ای با ساختار زیر داریم، خروجی هر لایه و تعداد پارامترهای آن را محاسبه کنید.

$$\begin{aligned} \text{Input} &= 224 * 224 * 3 \rightarrow 1 \\ 1 &\rightarrow \text{Conv}[3 * 3 * 64, \text{stride} = 1, \text{pad} = 1] \rightarrow 2 \\ 2 &\rightarrow \text{ReLU} \rightarrow 3 \\ 3 &\rightarrow \text{Pooling}[2 * 2 * 32, \text{stride} = 2] \rightarrow 4 \\ 4 &\rightarrow \text{FC}[10 - \text{class}] \rightarrow 5 \end{aligned}$$

\* اعداد ۱ تا ۴ برای مشخص کردن خروجی و ورودی های لایه ها استفاده شده اند.

(ب) در صورتیکه بخواهیم فقط با یک لایه ی تمام متصل (FC) تصویری با ابعاد ورودی و خروجی مشابه داشته باشیم، تعداد پارامترهای لایه ی تمام متصل را محاسبه کنید

(ج) تفاوت Deformable Convolution با Standard Convolution در چیست و چه مزایایی نسبت به حالت استاندارد دارد؟

الف) می دانیم که اگر تصویری با ابعاد  $(h * w * c)$  داشته باشیم و کرنلی با ابعاد  $(k_1, k_2, f)$  و با گام  $s$  و پدینگ با  $p$  بر روی آن اعمال کنیم. ابعاد خروجی به صورت زیر خواهد بود.

$$\left( \left\lfloor \frac{h - k_1 + 2p}{s} \right\rfloor + 1, \left\lfloor \frac{w - k_2 + 2p}{s} \right\rfloor + 1, f \right)$$

و اگر لایه max pooling با ابعاد  $(k_1, k_2)$  با گام  $s$  داشته باشیم. ابعاد خروجی به صورت زیر خواهد بود.

$$\left( \left\lfloor \frac{w - k_1}{s} \right\rfloor + 1, \left\lfloor \frac{h - k_2}{s} \right\rfloor + 1, c \right)$$

بنابراین خروجی لایه اول به صورت زیر خواهد بود.

$$1 \rightarrow \left( \left\lfloor \frac{224 - 3 + 2 \times 1}{1} \right\rfloor + 1, \left\lfloor \frac{224 - 3 + 2 \times 1}{1} \right\rfloor + 1, 64 \right) = (224, 224, 64)$$

خروجی لایه دوم و سوم نیز به صورت زیر محاسبه می شود.

$$2 \rightarrow (224, 224, 64)$$

$$3 \rightarrow (112, 112, 64)$$

خروجی لایه آخر نیز یک وکتور به ابعاد روبرو خواهد بود.

$$4 \rightarrow (10, 1)$$

برای محاسبه تعداد پارامترها نیز به صورت زیر عمل می‌کنیم که یک لایه convolutional با ابعاد  $(k_1, k_2, f)$  تعداد  $k_1 * k_2 * f$  پارامتر خواهد داشت و همینطور هنگامیکه می‌خواهیم تصویر را به لایه Fully Connected ببریم نیاز داریم ماتریس وزنی با ابعاد  $(fc, h * w * c)$  به علاوه تعداد  $fc$  تا bias داشته باشیم. بنابراین تعداد پارامترهای هر لایه نیز به صورت زیر محاسبه می‌شود.

$$1 \rightarrow 3 \times 3 \times 64 = 576$$

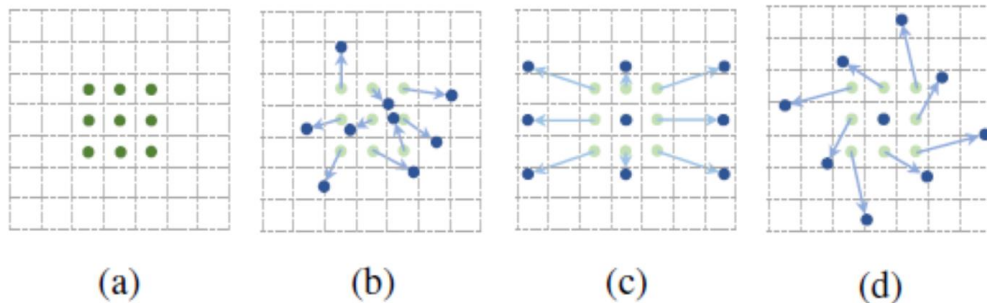
$$2 \rightarrow \text{no parameters}$$

$$3 \rightarrow \text{no parameters}$$

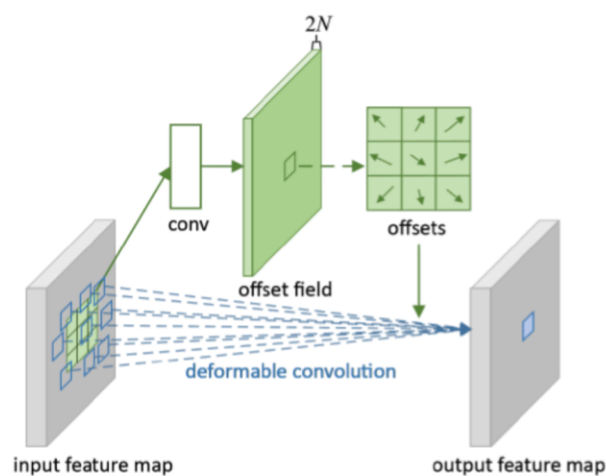
$$4 \rightarrow (112 \times 112 \times 64 \times 10) + 10 = 8028170$$

(ب) در صورتیکه بخواهیم با لایه تمام متصل تصویر ورودی را به خروجی مذکور ببریم به  $(224 * 224 * 3 * 10) + 10 = 1505290$  پارامتر نیاز خواهیم داشت.

(ج) هنگام استفاده از شبکه‌های convolutional مشکلی داریم و آن این است که ممکن است تصاویر نگاشت‌های هندسی داشته باشند نسبت به حالت استاندارد خودشان و برای مواجهه با این مشکل می‌توانیم به طور مثال داده‌ها را Augment کنیم و این تغییرات هندسی را تا حد ممکن روی داده‌ها ایجاد کنیم. مشکل این روش این است که این تغییراتی که ایجاد می‌کنیم محدود و ثابت هستند و تعمیم پذیری کمی دارند. حال می‌توانیم از Deformable Convolutionها استفاده کنیم. این لایه‌های کانولوشنی به این صورت عمل می‌کنند که لزوماً نقاط کنار هم در یک کرنل را نگاه نمی‌کنند مانند Standard Convolutionها بلکه می‌توانند نقاط را با تغییرات هندسی و یا بدون قاعده خاصی نگاه کنند و این تغییر در اعمال نقاط کرنل به صورت Adaptive صورت می‌گیرد. نشان داده شده است که این روش پارامتر و سربار محاسباتی کمی تولید می‌کند و در عوض دقت شبکه را به نسبت قابل توجهی افزایش می‌دهد.



(a) Conventional Convolution, (b) Deformable Convolution, (c) Special Case of Deformable Convolution with Scaling, (d) Special Case of Deformable Convolution with Rotation



Deformable Convolution

Regular convolution

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n)$$

Deformable convolution

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n + \Delta p_n)$$

where  $\Delta p_n$  is generated by a sibling branch of regular convolution

به طور مثال برای Object Detection اگر object مورد نظر کوچک و یا بزرگ باشد به طور مثال یک ماشین در فاصله‌ای دور و یا یکی اتوبوس بزرگ در فاصله‌ای نزدیک، این روش به راحتی آن‌ها را شناسایی می‌کند.



Examples: three levels of 3×3 deformable filters for three activation units (green points) on the background (left), a small object (middle), and a large object (right)

## مسئله‌ی ۲. Transposed Convolution

کانولوشن بین ورودی  $X$  و فیلتر  $W$  به صورت زیر است:

$$X = \begin{bmatrix} x_{(.,.)} & x_{(.,1)} & x_{(.,2)} \\ x_{(1,.)} & x_{(1,1)} & x_{(1,2)} \\ x_{(2,.)} & x_{(2,1)} & x_{(2,2)} \end{bmatrix} \quad W = \begin{bmatrix} w_{(.,.)} & w_{(.,1)} \\ w_{(1,.)} & w_{(1,1)} \end{bmatrix}$$

می‌توان عملیات کانولوشن را به صورت ضرب ماتریسی نوشت که ورودی و خروجی را به صورت یک بردار و فیلتر را به صورت یک ماتریس در نظر گرفت. ورودی  $X$  را به صورت بردار زیر نمایش می‌دهیم:

$$X = \begin{bmatrix} x_{(.,.)} & x_{(.,1)} & \dots & x_{(2,2)} \end{bmatrix}^T$$

(آ) عملیات کانولوشن بالا با  $S = 1$  را به صورت ضرب ماتریسی  $Y = AX$  بنویسید.

(ب) با استفاده از نمایش ماتریسی بالا می‌توان گرادیان پس انتشار نسبت به ورودی  $X$ ، را به صورت

$$\frac{\partial L}{\partial X} = A^T \frac{\partial L}{\partial Y}$$

نمایش داد.

عملیات transposed convolution را می‌توان مشابه عملیات گرادیان پس انتشار نسبت به ورودی در نظر گرفت و می‌توان این عملیات را به صورت کانولوشن مستقیم با ماتریس  $A^T$  به عنوان فیلتر دانست.

فرض کنید  $X$  خروجی یک عملیات کانولوشن با فیلتر  $W$  و  $\text{stride} = 2$  است. خروجی transposed convolution را با ورودی‌های زیر محاسبه کنید.

$$X = \begin{bmatrix} 2 & 1 \\ 3 & 6 \end{bmatrix} \quad W = \begin{bmatrix} 1 & 2 \\ 4 & 1 \end{bmatrix}$$

آ) برای انجام کانولوشن به صورت ضرب ماتریسی، به صورت زیر عمل می‌کنیم.

$$\begin{bmatrix} w_{(0,0)} & w_{(0,1)} & 0 & w_{(1,0)} & w_{(1,1)} & 0 & 0 & 0 & 0 \\ 0 & w_{(0,0)} & w_{(0,1)} & 0 & w_{(1,0)} & w_{(1,1)} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{(0,0)} & w_{(0,1)} & 0 & w_{(1,0)} & w_{(1,1)} & 0 \\ 0 & 0 & 0 & 0 & w_{(0,0)} & w_{(0,1)} & 0 & w_{(1,0)} & w_{(1,1)} \end{bmatrix} \times \begin{bmatrix} x_{(0,0)} \\ x_{(0,1)} \\ x_{(0,2)} \\ x_{(1,0)} \\ x_{(1,1)} \\ x_{(1,2)} \\ x_{(2,0)} \\ x_{(2,1)} \\ x_{(2,2)} \end{bmatrix}$$

$$= \begin{bmatrix} w_{(0,0)}x_{(0,0)} + w_{(0,1)}x_{(0,1)} + w_{(1,0)}x_{(1,0)} + w_{(1,1)}x_{(1,1)} \\ w_{(0,0)}x_{(0,1)} + w_{(0,1)}x_{(0,2)} + w_{(1,0)}x_{(1,1)} + w_{(1,1)}x_{(1,2)} \\ w_{(0,0)}x_{(1,0)} + w_{(0,1)}x_{(1,1)} + w_{(1,0)}x_{(2,0)} + w_{(1,1)}x_{(2,1)} \\ w_{(0,0)}x_{(1,1)} + w_{(0,1)}x_{(1,2)} + w_{(1,0)}x_{(2,1)} + w_{(1,1)}x_{(2,2)} \end{bmatrix}$$

ب) ماتریس  $W$  که در ورودی اولیه ضرب شده‌است و خروجی  $X$  را تولید کرده‌است به صورت زیر بوده‌است.

$$W = \begin{bmatrix} 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 \end{bmatrix}$$

بنابراین برای عملیات **transposed convolution** که منجر به **up-sampling** می‌شود، نیاز داریم که این  $W^T$  را در بردار حاصل از

ماتریس خروجی  $X$  به صورت زیر ضرب کنیم.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 4 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 2 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 4 \\ 0 & 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 2 \\ 1 \\ 3 \\ 6 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 1 \\ 2 \\ 8 \\ 2 \\ 4 \\ 1 \\ 3 \\ 6 \\ 6 \\ 12 \\ 12 \\ 3 \\ 24 \\ 6 \end{bmatrix} \rightarrow \text{up-sampled} = \begin{bmatrix} 2 & 4 & 1 & 2 \\ 8 & 2 & 4 & 1 \\ 3 & 6 & 6 & 12 \\ 12 & 3 & 24 & 6 \end{bmatrix}$$

### مسئله‌ی ۳. Object Detection

یکی از کاربردهای شبکه‌های کانولوشن، آشکارسازی اشیا است. مقالات مرتبط را مطالعه کرده و به سوالات زیر پاسخ دهید.

(آ) روش‌های YOLO و RCNN و RFCN را بررسی کرده و با یکدیگر مقایسه کنید.

(ب) در آشکارسازی‌های real-time استفاده از کدام یک را پیشنهاد می‌کنید.

(ج) لایه Spatial Pyramid Pooling را مختصراً توضیح داده و مزایای استفاده از آن را، در آشکارسازی تصاویر بیان کنید.

(آ) روش RCNN در ابتدا ۲۰۰۰ مکان از تصویر را انتخاب می‌کند و سپس مراحل بعدی را فقط با این ۲۰۰۰ مکان ادامه می‌دهد و نه تعداد زیادی مکان در تصویر. برای انتخاب ۲۰۰۰ مکان در ابتدا مکان‌هایی را انتخاب می‌کند و سپس این مکان‌ها را با یکدیگر ترکیب می‌کند تا در نهایت به ۲۰۰۰ مکان مورد نظر برسد و سپس این مکان‌ها را برای کلاس بندی به شبکه Convolutional وارد می‌کند. این روش نسبت به روش‌های غیر کانولوشنی مزیت‌های بسیاری دارد اما برای کاربردهای real-time بسیار کند است و همینطور گذاشتن عدد ثابت ۲۰۰۰ برای تمام کاربردها مناسب نخواهد بود.

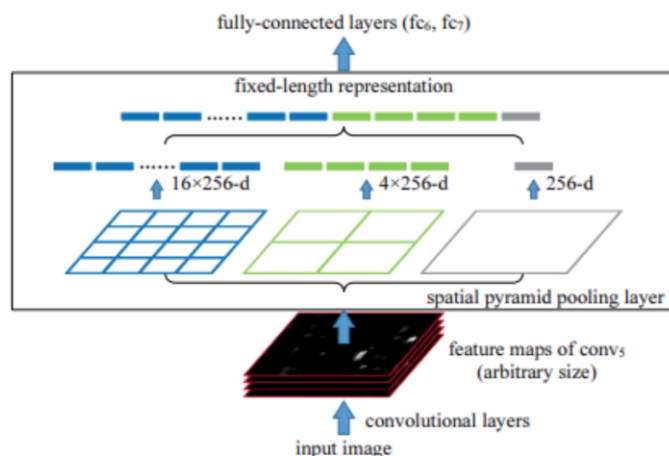
برخلاف روش قبلی به هر مکان یک بخش از شبکه را صدها بار اعمال می‌کند، روش RFCN، روشی دقیق و بهینه است که محاسبات را بر روی تمام تصویر پخش می‌کند. این روش یک جدول امتیاز برای مکان‌ها ایجاد می‌کند و برای محاسبه این امتیاز یک trade-off بین مستقل از نگاشت در کلاس بندی تصویر و وابسته به نگاشت در پیدا کردن اشیا استفاده می‌کند.

روش You Only Look Once یا به اختصار YOLO به این صورت عمل می‌کند که در هنگام تست به تمام تصویر به یکباره نگاه می‌کند و تخمین آن بر اساس مفهوم کلی تصویر است و همینطور برخلاف مدل‌هایی نظیر RCNN این مدل تخمین را با یک شبکه انجام می‌دهد. این

روش یک **grid** بر روی تصویر اعمال می‌کند و اگر مرکز یک شی در یکی از خانه‌های **grid** بیافتد. آن خانه مسئول تشخیص آن شی می‌شود. هر خانه **grid** دارای ۵ **bouding box** است و برای این **bouding box**ها امتیاز اطمینان محاسبه می‌کند این امتیاز میزان اطمینان از وجود یک شی در **bouding box** را مشخص می‌کند. برای هر **bouding box** یک کلاس نیز تخمین می‌زند. ترکیب امتیاز اطمینان و کلاس، احتمال وجود یک شی از کلاس خاص را در **bouding box** مشخص می‌کند.

(ب) برای کاربردهای **real-time** معمولا روش **YOLO** ترجیح داده می‌شود.

(ج) هنگام استفاده از شبکه‌های **Covolutional** معمولا به این صورت عمل می‌کنیم که پس از چند لایه **Convolutional** ویژگی‌های استخراج شده را در قالب یک بردار در می‌آوریم و به لایه **Fully Connected** می‌دهیم که این وکتور می‌تواند سایزهای متفاوتی داشته باشد. روش **Spatial Pyramid Pooling** به این صورت عمل می‌کند که در سطوح متفاوتی عمل **Pooling** را انجام می‌دهد. به طور مثال در تصویر زیر سه سطح **Pooling** داریم که در سطح اول که خاکستری است فقط یک **bin** داریم که از کل تصویر در لایه‌های متفاوت آن **max-pool** می‌گیرد بنابراین با داشتن ۲۵۶ لایه ویژگی، ۲۵۶ عدد حاصل از **max-pooling** خواهیم داشت. در سطح بعدی که سبز رنگ است ۴ **bin** داریم و بنابراین در هر یک از **bin**ها عمل **max-pooling** انجام می‌شود و  $4 \times 256$  عدد خواهیم داشت. سپس این اعداد را در کنار یکدیگر قرار می‌دهیم و بردار حاصل را به لایه تمام متصل می‌دهیم. این روش اطلاعات مکانی را در جایگاه خودشان نگه می‌دارد و همینطور ما را از یکسان کردن سایز تصاویر بی‌نیاز می‌کند.



Spatial Pyramid Pooling (Credits: [Paper](#))