

1) تبدیل به بردار: $10 \times 10 \times 5 \rightarrow 500$ (Flatten)
 یا FC: $10 \times 10 \times 5$

ماتریس به بردار تبدیل می شود. روش اول یک بردار محاسبه می شود. روش دوم (به دلیل Flatten) بردار را به ماتریس تبدیل می کند. به همین تفاوت اشاره می کنند.

2) ورودی: 128×128
 تعداد فیلتر: 16
 سایز کانولوشن: 5×5
 stride: 1
 Padding: 2

input size: n , kernel size: 5 , padding: 2

$$\frac{n - k + 2p}{s} + 1 = \frac{128 - 5 + 2 \times 2}{1} + 1 = 128$$

 خروجی: $128 \times 128 \times 16$

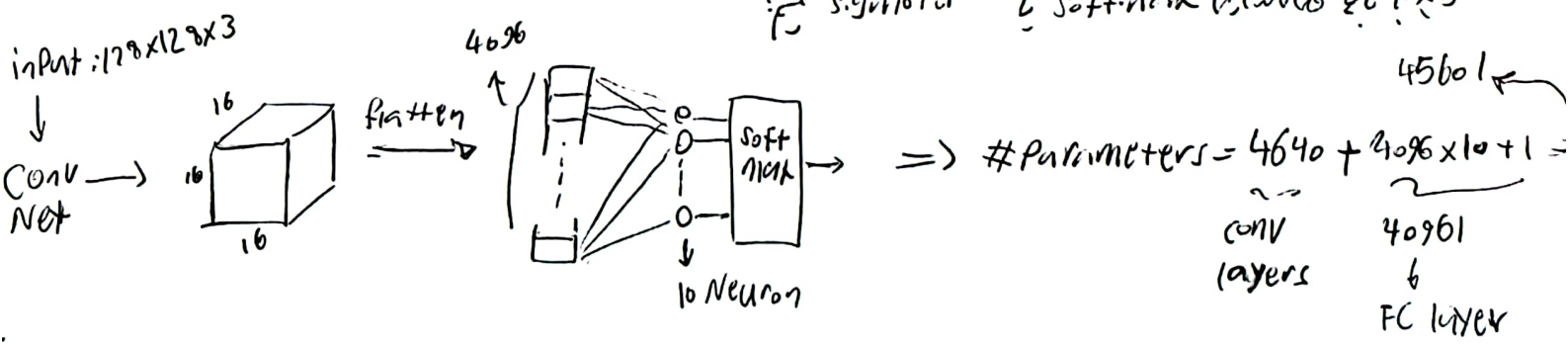
conv $\rightarrow 5 \times 5 \times 16$, stride=1, padding=2
 max Pooling $\rightarrow 2 \times 2$, stride=2
 ReLU

output size = $\frac{\frac{n - 5 + 2 \times 2}{1} + 1 - 2}{2} + 1 = \frac{128 - 2}{2} + 1 = \frac{127}{2} + 1 = \frac{129}{2}$

3) خروجی: $n \times (\frac{1}{2})^3 = \frac{n}{8} = \frac{128}{8} = 16$
 یعنی: $16 \times 16 \times 16$
 یعنی: $16 \times 16 \times 16$

تعداد پارامترها: $(3 \times 3 \times 2 + 1) \times 16 = 448$ (kernel size)
 تعداد پارامترها: $[(3 \times 3 \times 16 + 1) \times 16] \times 2 = 4640$ (max Pooling)
 # Parameters = $4640 + 448 = 4640$

* برای مدل طبقه بندی (Classification) به خروجی لایه کانولوشن اضافه می شود.
 لایه با تابع فعال سازی Softmax یا sigmoid



* Receptive field: در صورتی که در یک شبکه، هر پیکسل از ورودی به یک پیکسل در خروجی مربوط می‌شود.

Receptive Field (در یک شبکه)

$$R'_l = R_{l-1} + (R_{l-1} - 1) \times S_{l-1}$$

در صورتی که Pooling وجود دارد: $R_l = (n_{pooling}) \times R'_l$ \rightarrow در صورتی که Pooling وجود ندارد: $R_l = R'_l$

مثال: $\Rightarrow R_3 = 2 \left(\underset{\substack{\text{pooling} \\ =5}}{K_2} + (\underset{\substack{\text{kernel} \\ =5}}{R_2} - 1) \times \underset{\substack{\text{stride} \\ =1}}{1} \right) = 2(R_2 - 4) \stackrel{(*)}{=} 2(536 + 4) = 1080 \rightarrow \text{Receptive field}$

$R_2 = 2 \left(\underset{\substack{\text{pooling} \\ =5}}{K_2} + (\underset{\substack{\text{kernel} \\ =5}}{R_1} - 1) \times \underset{\substack{\text{stride} \\ =1}}{1} \right) = 2(R_1 + 4) \stackrel{(*)}{=} 2(264 + 4) = 536 \stackrel{(*)}{\leftarrow}$

$R_1 = 2 \left(\underset{\substack{\text{pooling} \\ =5}}{K_1} + (\underset{\substack{\text{kernel} \\ =128}}{R_0} - 1) \times \underset{\substack{\text{stride} \\ =1}}{1} \right) = 2(132) = 264 \stackrel{(*)}{\leftarrow}$

* سوالات مفهومی U-net

3) الف) ویژگی اصلی شبکه U-net نسبت به شبکه‌های کانولوشنی عادی، تعادل بین بزرگ‌ها و کوچک‌ها (کانولوشن در دو سمت) و downsampling و upscaling است.

و Upsampling است و وجود skip connection بین بزرگ‌ها و کوچک‌ها (کانولوشن متناظر در هر دو سمت برای حفظ ویژگی‌های اصلی تصویر است). به افزار U شکل آنتیز به این تعادل گفته شده است.

ب) وجود skip connection در این شبکه‌ها باعث حفظ ویژگی‌های فضا (spatial) (عکس) (رابطه‌ها) سمت downsampling شده و همچنین فرایند error back propagation / رفتن آنتیز به این تعادل وزن‌ها و پارامترهای شبکه تسهیل می‌شود و مانع از Gradient vanishing می‌شود.

ج) وجود skip connection با توجه به اینکه باعث حفظ اطلاعات مکانی تصویر می‌شوند، به دلیل کوچک بودن ریزاست تصاویر پزشکی، توانایی این نوع شبکه‌ها می‌تواند در کاربرد semantic segmentation از اهمیت بالایی برخوردار باشد. در تصویر پزشکی به دلیل شباهت زیاد بافت (texture) (background و foreground)، وجود skip connection می‌تواند با ترکیب ویژگی‌های سطح بالا و سطح پایین تصویر در تارت short skip connection به segment کردن کمک کند.

* سوالات مفهومی DenseNet

الف) در شبکه ResNet در بخش‌هایی از شبکه خروجی‌های کانولوشنی با خروجی‌های کانولوشنی دیگر، از لایه‌های قبل جمع می‌شوند (residual connection).

$$x_l = \text{conv}(x_{l-1}) + x_{l-1}$$

ادامه حل سوالات مربوط به DenseNet بخش الف)

ادامه شبکه DenseNet لایه کانولوشنی از مجرای ۱۶ از ویژگی‌های لایه قبل که به یکدیگر concatenate شده استفاده می‌کنند:

$$x_l = \text{conv}(\left[\text{concat}(x_1, \dots, x_{l-1}) \right])$$

ب) این شبکه با ایجاد skip connection از لایه‌های قبلی به دست می‌آید که اگر این از چندین مسیر عبور کند و منجر به Vanishing Gradient شود. با توجه به این نکته، می‌توانیم شبکه را در ازای کاهش ساینه ویژگی‌ها عمیق‌تر کرده و تعداد پارامترهای شبکه را کاهش داد که به دیگر ارزش‌های این شبکه است.

* سوالات محاسبه‌ای مربوط به U-net

الف) با توجه به اینکه در شبکه U-net در مرحله down Sampling ساینه تصویر از ۴ به ۱۶ می‌شود، پس ساینه تصویر از محقق‌ترین لایه این شبکه بدین‌گونه است:

$$16 \times 16 \leftarrow \frac{256}{24} = \frac{256}{16} = 16$$

ب)

تعداد پارامترهای کانولوشنی در

$$128 \times \left(\overset{977}{3 \times 3 \times 64 + 1} \right) = 73856$$

* سوالات محاسبه‌ای DenseNet:

الف) با توجه به اینکه شبکه DenseNet تمامی فیچرهای لایه‌های قبل از خود را ریاضت می‌کند، پس $128 + 128 = 192$ فیچرهای لایه بعدی خواهد داشت. $256 + 192 = 448$ فیچرهای لایه بعدی خواهد داشت.

ب)

تعداد کانال فیچرهای شبکه

$$\left\{ \begin{array}{l} \text{خروجی لایه 1} : 32 + k \\ \text{خروجی لایه 2} : 32 + 32 + k + k = 64 + 2k \\ \text{خروجی لایه 3} : 32 + \underbrace{(32 + k)}_{\text{لایه 1}} + \underbrace{(64 + 2k)}_{\text{لایه 2}} + k = 128 + 4k = 224 \end{array} \right.$$

الف) در شبکه های کانولوشن معمولی عملیات Grid Sampling ثابت است (در حالی که در کانولوشن های از الگوهای سفارشی همراه با offset یا یاگیری نوک شبکه استفاده میشود. این بدان معناست که شبکه براساس task مشخص شده ^{بهترین} انتخاب sample ها را برای kernel کانولوشن به منظور استخراج بهترین ویژگی ها و افزایش Receptive Field شبکه یاد می بخشد. بنابراین هر فیلتر sample ها کرنل خود را از بخش ها ۱۶ تصویر بر می دارد.

ب) شبکه های deformable به دلیل اینکه محل انتخاب sample ها را در فرآیند یاگیری می فرماید، انعطاف پذیری بیشتری نسبت به شبکه های کانولوشن عادی یا grid sampling دارد. زیرا می تواند ویژگی ها را دقیق به حرکت از تبدیل ها هندسی (geometric Transformations) را با تغییر شکل کرنل خود استخراج کند. هندسی

ج) شبکه های کانولوشن معمولی به دلیل ثابت بودن مساحت کرنل کانولوشن خود صلابت بالایی نسبت به تغییرات هندسی عکس ورودی دارند و به همین خاطر در صورتی که عکس ورودی چرخش زیادی داشته باشد، شبکه قادر به تشخیص و کلاس بندی اشیا داخل تصویر نیست. به عبارتی شبکه های کانولوشن معمولی Shift Invariant هستند.