

Input	۱۲۸ × ۱۲۸ × ۳	(۱) (ب)
Conv 9,32	۱۲۰ × ۱۲۰ × ۳۲	$\lfloor \frac{128-9}{1} + 1 \rfloor = 120$
Maxpool 2	۶۰ × ۶۰ × ۳۲	$\lfloor \frac{120-2}{2} + 1 \rfloor = 60$
Conv 5,64	۵۶ × ۵۶ × ۶۴	$\lfloor \frac{60-5}{1} + 1 \rfloor = 56$
Maxpool 2	۲۸ × ۲۸ × ۶۴	$\lfloor \frac{56-2}{2} + 1 \rfloor = 28$
Conv 5,64	۲۴ × ۲۴ × ۶۴	$\lfloor \frac{28-5}{1} + 1 \rfloor = 24$
Maxpool 2	۱۲ × ۱۲ × ۶۴	$\lfloor \frac{24-2}{2} + 1 \rfloor = 12$
flatten	۹۲۱۶ × ۱	$12 \times 12 \times 64 = 9216$
FC-3	۳ × ۱	$32 \times (9 \times 9 \times 3) + 64 \times (5 \times 5 \times 4) + 64 \times (5 \times 5 \times 4) + 9216 \times 3 = 189024$ $32 + 64 + 64 + 3 = 143 \text{ تعداد لایه ها}$

(ب) قرار دادن تعداد فعال ساز غیر خطی باعث می شود که مدل توانایی یادگیری روابط غیر خطی را نیز داشته باشد. اگر از تعداد فعال ساز خطی استفاده کنیم، هر چند هم که شبکه عمیق باشد، باز هم یک مدل خطی محسوب می شود و نمی تواند روابط غیر خطی را کشف کند. بنیاد این بدان است که قدرت یادگیری شبکه های خطی، از تعداد فعال ساز غیر خطی استفاده می کند.

(ج) مزایای از آنجایی که با استفاده از Pooling اسلاید دردی کمتر می شود، بنابراین می توان خطر overfit شدن مدل را کاهش داد. از طرف دیگر از آنجایی که پارامترهای یادگیری را کاهش می دهد بنابراین هزینه محاسباتی مدل را نیز کاهش می دهد.

معایب: استفاده از Max pooling باعث از دست رفتن اطلاعاتی می شود که ممکن است مفید باشند، زیرا از بین می رانند. ورودی فقط ما شکم آن ما را مجاز می کند و شبکه دیگر ما را در می اندازد.

(د) از آنجایی که آخرین لایه مدل دسته بندی را می ۳ درایه است بنابراین باید از تابع فعال ساز softmax برای استخراج احتمال حرکت از دسته ها استفاده کرد.

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^3 e^{z_j}}$$

$$z = [0.1, 0.5, 0.3]^T \Rightarrow \text{softmax}(z) = \begin{bmatrix} \frac{e^{0.1}}{e^{0.1} + e^{0.5} + e^{0.3}} \\ \frac{e^{0.5}}{e^{0.1} + e^{0.5} + e^{0.3}} \\ \frac{e^{0.3}}{e^{0.1} + e^{0.5} + e^{0.3}} \end{bmatrix} = \begin{bmatrix} 0.127 \\ 0.40 \\ 0.473 \end{bmatrix}$$

۲) الف) از آنجایی که فیلترهای لایه‌های کانولوشنی به صورت تفران به مدای فیلترهای تقویری دردی ایمان می‌شود و این فیلتر دارای پارامترهای مشخص و یکسانی است می‌توان گفت که ویژگی weight sharing وجود دارد. این ویژگی باعث می‌شود که پیدا کردن یک اسکوی مشخص با استفاده از این فیلتر، دایره به ناصیه خاص نباشد و اگر آن اسکوی خاص در نامیدی دیگری از تصویر نیز باشد، می‌توان همان فیلتر آن را مشخصی دارد. به این ویژگی لایه‌های کانولوشنی Translational Equivariance می‌گویند. تابع $f(m)$ مت تابع $g(m)$ equivariant است اگر $f(g(m)) = g(f(m))$ در اینجا می‌توان گفت $f(m)$ تابع کانولوشن و $g(m)$ تابع translation است.

ب) برای اینکه به یک تقویر پردی سیستم مشخص کنیم که مدل مت به کدام سمت‌ها از تصویر پردی می‌شود و آن سمت‌ها را ملاک مشخصی قرار می‌دهد نیاز است که گرایان های لایه آخر (مثل softmax) را به عقب برگردانیم و به کمک این گرایان های حاصل می‌توان فهمید که کدام مدای پردی تأثیر بیشتری در استخراج دارند. حال در شبکه‌های deconvolutional نحوه انتشار به عقب گرایان ما به این صورت است که آن پارامترهای که در برده forward می‌بودند و همچنین آن پارامترهای که در برده انتشار به عقب می‌شدند را همان می‌کنند و تأثیر این رویه‌ها را ندیده می‌گیریم. به این روش $\text{Guided Backpropagation}$ می‌گویند.

ردای که استفاده می‌شود به صورت زیر است: $\text{activation} : f_i^{l+1} = \text{Relu}(f_i^l)$
 $\text{guided backprop} : R_i^l = (f_i^l > 0) \cdot (R_i^{l+1} > 0) \cdot R_i^{l+1}$

در شبکه‌های UP-convolution و خلاصه‌ای شبکه‌ها که برای تقویر پردی گرایان ما را به عقب می‌گردانند، در اینجا سعی می‌شود که یک شبکه جدید آموزش داده شود که به طوری که Φ و این ویژگی‌های حاصل شود، بتواند تصویری پردی را تخمین بزند. در واقع یک شبکه جدید می‌سازد که پردی آن ویژگی‌های استخراج شده (Φ_i) است و پردی آن تصویری پردی (X_i) را پارامترهای این شبکه نیز با استفاده از سیستم کردن تابع هزینه

$$\hat{w} = \arg \min_w \sum_i \|X_i - f(\Phi_i, w)\|_2^2$$

۲) این object detector های تک مرحله ای، تصویر ورودی را ناحیه بندی می کنند و فقط ی اساسی تصویر ورودی یک تعداد بابت Prediction انجام می دهند.

object detector های دو مرحله ای ابتدا تصویر ورودی را ناحیه بندی می کنند و تعدادی درون حریک از نواحی را پیدا می کنند پس در مرحله بعدی، با نته های مرحله قبل را بهبود می دهند و Prediction نهایی را انجام می دهند.

ب) تشخیص دهنده Yolo فقط یک بار تصویر ورودی را می بیند و آن را به شبکه خود وارد می کند و

این را تشخیص می دهد. به همین دلیل نت به پای تشخیص دهنده ما سرعت بیشتری دارد.

ردنی RNN اینها نواحی که محتمل حضور اینا هستند را خارج از دسته بندی آن ها تشخیص می کند. پس

در مرحله بعدی با استفاده از یک شبکه کانولوشنی، ویژگی های حریک از نواحی مرحله قبل را استخراج

می کند و در نهایت با استفاده از یک دسته بندی، دسته مربوط به حریک از نواحی را Predict می کند.

۲) Retina Net یک object detector تک مرحله ای است که با استفاده از دو خاصیت Focal loss و feature pyramid

به دقت بیشتری می رسد. ویژگی FPN به این صورت است که در مرحله پایین به بالا به ویژگی های با رزولوشن کمتر می رسد و

حال با استفاده از یک سری اتصالات کناری و یک شبکه بالا به پایین، ویژگی های لایه های جلوی را با ویژگی های لایه های عقبی

توصیف می کنند به این صورت هر یک سری ویژگی های جدید می رسد و پس این ویژگی های جدید را استفاده می کنند. در واقع نقش

FPN، feature extractor است که به صورت سری ویژگی های لایه های جلوی را با توصیف می کند.

۳) شبکه Inception از تعدادی از بلوک های Inception تشکیل شده است که در هر یک از این بلوک ها

از فیلترهای با سایزهای مختلف استفاده می کنند و در نهایت خروجی اعمال این فیلترها با هم concat

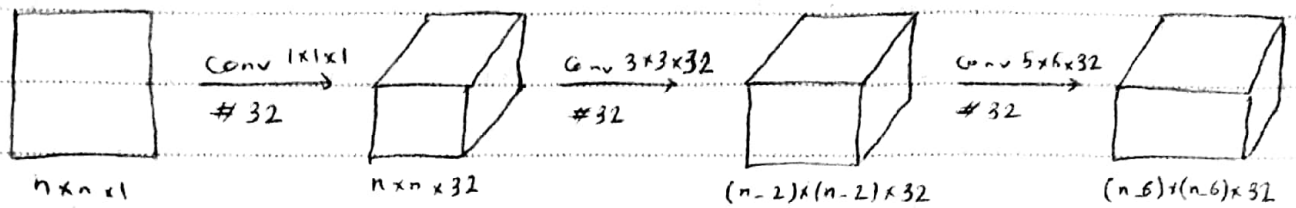
می کنند البته خروجی این فیلترها باید به گونه ای باشند که ابعاد خروجی با دردی برابر باشد یعنی same convolution

که این امر با استفاده از padding درون به درون قابل انجام است.

هدف از اعمال فیلترهای با سایزهای مختلف در هر لایه می تواند این باشد که شبکه را قادر کند که اشیای

نوعی مختلف را در مقیاس های متفاوت را بتواند یاد بگیرد و به این طریق دقت شبکه را

افزایش دهد.



تعداد پارامترها به صورت زیر می شود:

~~$$(1 \times 1 \times 1) + 32 + (3 \times 3 \times 32) + 32 + (5 \times 5 \times 32) + 32 = 1152$$~~

$$32 \left[(1 \times 1 \times 1) + 1 \right] + 32 \left[(3 \times 3 \times 32) + 1 \right] + 32 \left[(5 \times 5 \times 32) + 1 \right] = 34624$$

در شکل دوم که منظره ها یکی به ورودی که یکی کانال دارد اماون حالتند، تعداد پارامترها به صورت زیر می شود:

$$32 \left[(1 \times 1 \times 1) + 1 \right] + 32 \left[(3 \times 3 \times 1) + 1 \right] + 32 \left[(5 \times 5 \times 1) + 1 \right] = 1216$$

با) وقتی که از بین از توابع فعال ساز مانند sigmoid استفاده کنیم، فرای مستقیم در لایه عددی

بین مندریک است. برای که در رله Backprop می ضامیم با استفاده از Chain Rule گرایان ما

ما به عقب برگردانیم، این اعداد بین مندریک در هم ضرب می شوند و شاید این به صورت بخان کوچک می شوند.

تا برای اگر تعداد لایه ها خیلی زیاد باشند، گرایان ما به هیچ وقت به عقب به مندریک می کشند و این باعث

می شود که نتوان پارامتر شبکه را به صورت مناسب به روز رسانی کرد.

نسخه ی Res Net با استفاده از ایسی Skip Connection که مقایسه فرای لایه لایه را به لایه $t+t$ را

انحال می دهد می تواند اثر محو شدن گرایان را کامتی دهد و امکان آید مندی شبکه های عمیقتر را فراهم

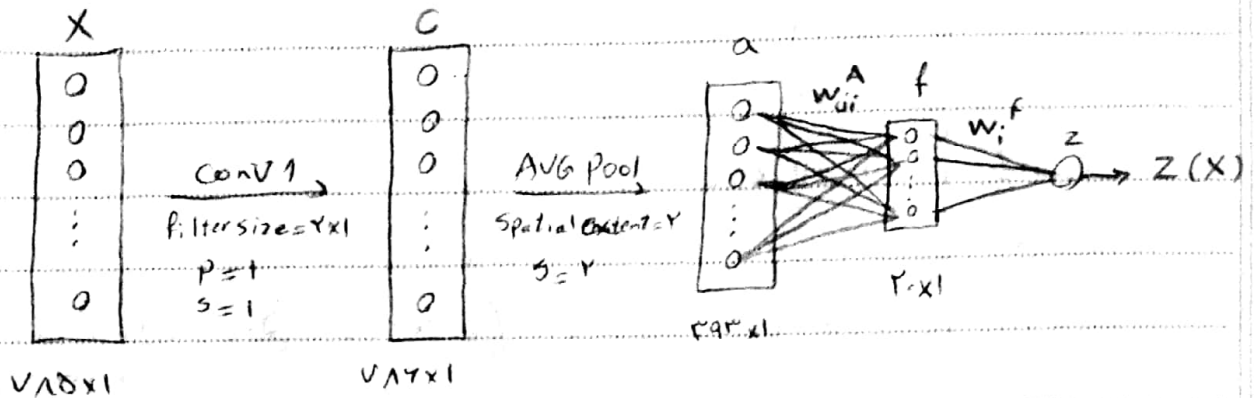
می آید. در رله Backprop نیز گرایان در لایه $t+t$ به لایه لایه منتقل می کنند و این صورت از

محدود کردن گرایان کمتر می شود.

⑤ (الف) شبکه یادار اتصالی که شبکه نیاز به یادگیری دارد برای است:

$$2 + 1 + (4 \times 4 \times 2) + 2 + 2 + 1 = 29.6$$

\downarrow \downarrow \downarrow \downarrow \downarrow
 با ۲ اتصالی، با ۱ با ۲۸ اتصالی لایه FC، با ۲ اتصالی لایه FC، با ۲ اتصالی لایه FC، با ۱



(ب) تابع Loss، پارامتری که شبکه آزمون می‌دهد: $Cost = -y \lg(z(x)) - (1-y) \lg(1-z(x))$

$$\frac{\partial Cost}{\partial w_{ji}^A} = \frac{\partial Cost}{\partial out_{f_i}} \frac{\partial out_{f_i}}{\partial in_{f_i}} \frac{\partial in_{f_i}}{\partial w_{ji}^A}$$

$$\frac{\partial Cost}{\partial out_{f_i}} = \frac{\partial Cost}{\partial out_2} \frac{\partial out_2}{\partial in_2} \frac{\partial in_2}{\partial out_{f_i}} = (z(x) - y) w_i^f$$

$$\frac{\partial Cost}{\partial out_2} = \frac{-y}{z(x)} + \frac{1-y}{1-z(x)} = \frac{-y(1-z(x)) + (1-y)z(x)}{z(x)(1-z(x))} = \frac{z(x) - y}{z(x)(1-z(x))}$$

$$\frac{\partial out_2}{\partial in_2} = z(x)(1-z(x)) \quad \frac{\partial in_2}{\partial out_{f_i}} = w_i^f$$

$$\frac{\partial out_{f_i}}{\partial in_{f_i}} = \begin{cases} 1 & ; in_{f_i} > 0 \\ 0 & ; in_{f_i} < 0 \end{cases} \quad \frac{\partial in_{f_i}}{\partial w_{ji}^A} = a(x^{(j)})$$

$$\text{LL} \Rightarrow \frac{\partial Cost}{\partial w_{ji}^A} = \begin{cases} (z(x) - y) w_i^f a(x^{(j)}) & ; in_{f_i} > 0 \\ 0 & ; in_{f_i} < 0 \end{cases}$$

$$\frac{\partial \text{cost}}{\partial \text{out}_{f_i}} = (z(x) - y) w_i^F$$

در حساب نشان داده شد که

$$\frac{\partial \text{cost}}{\partial a_i} = \sum_{i=1}^r \frac{\partial \text{cost}}{\partial \text{out}_{f_i}} \frac{\partial \text{out}_{f_i}}{\partial \text{in}_{f_i}} \frac{\partial \text{in}_{f_i}}{\partial a_i}$$

حال طبق

$$\frac{\partial \text{out}_{f_i}}{\partial \text{in}_{f_i}} = \begin{cases} 1 & ; \text{in}_{f_i} > 0 \\ 0 & ; \text{in}_{f_i} < 0 \end{cases} \quad \frac{\partial \text{in}_{f_i}}{\partial a_i} = w_{ji}^A$$

$$\Rightarrow \frac{\partial \text{cost}}{\partial a_i} = \begin{cases} \sum_{i=1}^r (z(x) - y) w_i^F w_{ji}^A & ; \text{in}_{f_i} > 0 \\ & ; \text{in}_{f_i} < 0 \end{cases}$$

$$\frac{\partial \text{cost}}{\partial \text{out}_{c_i}} = \frac{\partial \text{cost}}{\partial a_i} \frac{\partial a_i}{\partial \text{out}_{c_i}} = \frac{1}{r} \frac{\partial \text{cost}}{\partial a_i}$$

$$a_i = \frac{\text{out}_{c_i} + \text{out}_{c_r}}{r} \Rightarrow \frac{\partial a_i}{\partial \text{out}_{c_i}} = 1/r$$

$$\frac{\partial \text{cost}}{\partial \text{out}_{c_i}} = \frac{1}{r} \begin{bmatrix} \frac{\partial \text{cost}}{\partial a_1} \\ \frac{\partial \text{cost}}{\partial a_2} \\ \frac{\partial \text{cost}}{\partial a_3} \\ \frac{\partial \text{cost}}{\partial a_4} \\ \vdots \\ \frac{\partial \text{cost}}{\partial a_{n-1}} \\ \frac{\partial \text{cost}}{\partial a_n} \end{bmatrix}_{1 \times n}$$

بنابراین می‌توان گفت که

در نهایت می‌توان نوشت:

$$\frac{\partial \text{cost}}{\partial w_i^c} = \sum_{i=1}^{n+1} \frac{\partial \text{cost}}{\partial \text{out}_{c_i}} \frac{\partial \text{out}_{c_i}}{\partial \text{in}_{c_i}} \frac{\partial \text{in}_{c_i}}{\partial w_i^c}$$

$$\frac{\partial \text{out}_{c_i}}{\partial \text{in}_{c_i}} = \begin{cases} 1 & ; \text{in}_{c_i} > 0 \\ 0 & ; \text{in}_{c_i} < 0 \end{cases}$$

$$\begin{aligned}
 inc_r &= w_1^c x_0 + w_r^c x_1 \Rightarrow \frac{\partial inc_r}{\partial w_1^c} = x_0 \\
 inc_r &= w_1^c x_1 + w_r^c x_r \Rightarrow \frac{\partial inc_r}{\partial w_1^c} = x_1 \\
 &\vdots \\
 inc_{var} &= w_1^c x_{var} + w_r^c x_{var} \Rightarrow \frac{\partial inc_{var}}{\partial w_1^c} = x_{var}
 \end{aligned}$$

0-padding

$$\Leftrightarrow \frac{\partial inc_i}{\partial w_1^c} = x_{i-1}$$

$$\Rightarrow \frac{\partial cost}{\partial w_1^c} = \begin{cases} \sum_{i=1}^{var} \frac{\partial cost}{\partial out_i} x_{i-1} & ; inc_i > 0 \\ 0 & ; inc_i < 0 \end{cases}$$

$$\frac{\partial cost}{\partial w_1^c} = \frac{\partial cost}{\partial out_c} x_0 + \frac{\partial cost}{\partial out_r} x_1 + \dots + \frac{\partial cost}{\partial out_{var}} x_{var}$$

: prob $inc_i > 0$ نيف

$$= \frac{1}{r} \left[\frac{\partial cost}{\partial a_1} (x_0 + x_1) + \frac{\partial cost}{\partial a_r} (x_r + x_r) + \dots + \frac{\partial cost}{\partial a_{var}} (x_{var} + x_{var}) \right]$$

$$\begin{aligned}
 \frac{\partial cost}{\partial w_1^c} &= \frac{1}{r} \left[(x_0 + x_1) \sum_{i=1}^r w_i^F w_{ri}^A + (x_r + x_r) \sum_{i=1}^r w_i^F w_{ri}^A + \dots + \right. \\
 &\quad \left. (x_{var} + x_{var}) \sum_{i=1}^r w_i^F w_{ri}^A \right] \\
 &\quad (Z(x) - y)
 \end{aligned}$$

: prob $inc_i > 0$ نيف

۲) در و لایه محدود دید به میزان $K-1$ از این پیرامون که از آغاز که طول گام های یک است، بدان
مقدار لایه کانولوشن با فیلترهای $K \times K$ ، محدود دید فازی شود با $1 + (K-1)$

ب) اگر اندازه فیلترهای $F \times F$ و اندازه گام استیج D باشد، مثل این است که اندازه فیلتر F
به $1 + (F-1)$ تبدیل شود. حال اگر اندازه ورودی $M \times n$ باشد، آنگاه اندازه
خروجی D خواهد بود:

$$(M - D(F-1) - 1 + 1) \times (N - D(F-1) - 1 + 1) = (M - DF + D) \times (N - DF + D)$$

ج) همان نشان داد که ماتریس A یک ماتریس $D \times D$ خواهد بود که در این مسیر به نظر اول بدست آورده
شده یک D باشد و سایر در این ماتریس منفی هستند. به عنوان مثال اگر $D=2$ و $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$
التماسی تقسیم $A \otimes K$ حاصل شود. با $D=1$ به نظر بدست آورده آن که حتی منفی هستند با
حذف کنیم.

د) Masked convolution در این نوع خاص از فیلتر را معرفی کنیم (Mask) و برای صورت
ی نشان از بعضی از پیکسل ها را در نظر گرفت و این پیکسل ها را استرایم کرد. در این روش
صورت محدود دید فیلتر تغییر می کند و تعداد بعضی از ورودی های آن صفر می شود، یعنی تداوم ارتباط بین پیکسل های
جوار را بر روی بردار می گذارد. برای این روش $dilated convolution$ استفاده می شود و در این
روش برای محدود دید فیلتر را همان بزرگتر کرد و ارتباط بین پیکسل های که با هم فاصله دارند را بهتر
درک کرد.