

تشخیص اشیا با استفاده از شبکه‌های عصبی کانولوشن

حجت رئیسی (نویسنده مسئول)^۱، شیوا بنی هاشمی^۲، مسعود کاوسی^۳، فاطمه حسن نژادبان فرد شیرازی^۴، محمد لهراسبی^۵

^۱ رئیس اداره طراحی و توسعه سامانه‌های هوشمند شهرداری شیراز، شیراز hojjat.raisee@gmail.com

^۲ کارشناس حمل و نقل شهرداری شیراز، شیراز shivabanihashemi@gmail.com

^۳ کارشناس اداره طراحی و توسعه سامانه‌های هوشمند شهرداری شیراز، شیراز masoud@live.com

^۴ کارشناس منطقه شش شهرداری شیراز، شیراز fatemeh.shirazi68@yahoo.com

^۵ کارشناس اداره طراحی و توسعه سامانه‌های هوشمند شهرداری شیراز، شیراز mohammad.lohrasbi69@gmail.com

چکیده

در این مقاله، نشان می‌دهیم که تغییر الگوریتم (طرح‌های محاسباتی با استفاده از شبکه عمیق) منجر به یک حل موثر و قابل توجهی می‌شود که محاسبه طرح در محاسبه شبکه تشخیص، تقریباً بدون هزینه (زمان پردازش اندک) می‌باشد. در انتها، شبکه‌های طرح منطقه جدیدی (RPNS) را معرفی می‌کنیم که لایه‌های کانولوشن را با شبکه‌های تشخیص اشیا جدید و پیشرفته، به اشتراک می‌گذارد. با اشتراک لایه‌ها در زمان آزمایش (زمان تست)، زمان محاسبه طرح‌ها، کم می‌باشد (در حدود ۱۰ میلی ثانیه بر تصویر). مشاهدات ما نشان می‌دهد که ویژگی کانولوشن استفاده شده با استفاده از تشخیص دهندگان مبتنی بر منطقه، مانند R-CNN سریع، می‌تواند در ایجاد طرح‌های منطقه مورد استفاده قرار گیرد. در بالاترین لایه این ویژگی‌های کانولوشن، ما RPNS را با اضافه کردن دو لایه کانولوشن اضافی ایجاد می‌کنیم. این لایه‌ها عبارتند از: لایه‌ای که هر موقعیت کانولوشن را در بردارهای کوتاه تر رمزگذاری می‌کند و لایه‌ای که در هر موقعیت کانولوشن، امتیاز اشیا و مرزهای مربوطه شان را برای طرح‌های منطقه k مربوط به مقیاس‌های متنوع و ضریب منظرهایی (زاویه‌های دید) در آن مکان (موقعیت) به عنوان خروجی تعیین می‌کند یک مقدار معمول است. بنابراین RPNS پیشنهادی نوعی از شبکه کانولوشن کامل است که به خصوص برای ایجاد طرح‌های تشخیص می‌تواند بصورت انتها به انتها آموزش ببیند. به منظور یکپارچه کردن شبکه‌های تشخیص اشیا RPNS با R-CNN سریع، روش آموزش ساده‌ای را پیشنهاد می‌دهیم که بین تنظیم بهینه طرح منطقه و سپس تنظیم بهینه تشخیص اشیا تغییر می‌کند؛ این امر درحالی اتفاق می‌افتد که طرح‌ها را بصورت ثابت حفظ می‌کند. این روش بطور سریع همگرا می‌شود و شبکه‌ای یکنواخت با ویژگی‌هایی که بین هر دو روش به اشتراک گذاشته می‌شود را ایجاد می‌کند.

واژه‌های کلیدی

شبکه عمیق، تشخیص اشیا، شبکه کانولوشن

۱. مقدمه

اخیرا مزایای بدست آمده در تشخیص اشیا بر اساس موفقیت روش‌های طرح‌های منطقه (نظیر مقاله [22]) و شبکه‌های عصبی کانولوشن مبتنی بر منطقه [6] (R-CNNs) بوده است. اگرچه CNN‌های مبتنی بر منطقه بطور گسترده و وسیعی در مقاله [۶] توسعه یافته‌اند، اما هزینه (زمان راه‌اندازی) شان بطور قابل توجهی با استفاده از اشتراک کانولوشن در میان طرح‌ها [۷,۵] کاهش یافته است. پیشرفته‌ترین روش، R-CNN سریع [۵]، به نرخ‌هایی نزدیک به زمان واقعی با استفاده از شبکه‌های بسیار عمیق، (وقتی که زمان گذرانیده شده در طرح‌ها نادیده گرفته می‌شود) دست یافته است.

روش‌های طرح منطقه معمولاً به روش‌های استنتاجی اقتصادی و ویژگی‌های معقول و کم‌هزینه وابسته است. روش جستجوی انتخابی [22] (SS)، یکی از محبوبترین روش‌هایی است که بطور حریصانه‌ای سوپرپیکسل‌های مبتنی بر ویژگی‌های سطح پایین مهندسی شده را استخراج می‌نماید. هنوز هم در مقایسه با شبکه‌های موثر تشخیص، سرعت روش جستجوی انتخابی (SS) به اندازه ۲ ثانیه بر تصویر در بکارگیری زمان پردازش CPU کمتر است. اخیراً باکس‌های لبه [۲۴] شرایط بهینه‌ای را بین کیفیت و سرعت (۰,۲ ثانیه بر تصویر) ارائه می‌دهند. با این وجود، هر گام روش طرح منطقه، زمان پردازشی به اندازه شبکه تشخیص صرف می‌نماید.

به منظور ارزیابی روش پیشنهادی، از پایگاه داده PASCAL VOC [۴] استفاده نموده ایم که نشان می‌دهد روش ترکیبی RPNs با R-CNNs دقت تشخیص بالاتری را نسبت به روش ترکیبی جستجوی تصادفی با R-CNNs سریع ارائه می‌دهد. ضمناً روش این پیشنهادی تقریباً تمام موانع محاسباتی روش SS را در مرحله آزمایش (تست) رفع می‌کند. در روش پیشنهادی، زمان راه‌اندازی موثر برای طرح‌ها تنها برابر با ۱۰ میلی ثانیه است. با استفاده از مدل‌های بسیار عمیق [۱۹]، روش تشخیص پیشنهادی مان‌دارای نرخ فریم ۵ fps (شامل تمام گام‌ها) در GPU بوده که بنابراین از دو نقطه نظر سرعت و دقت، یک سیستم تشخیص اشیای عملی و مناسب محسوب می‌گردد.

۲. سوابق تحقیق

برخی از مقالات اخیر در زمینه روش‌های ارائه شده مبتنی بر شبکه‌های عمیق تعیین مکان یک کلاس خاص یا باکس‌های مرزی کلاس‌های مخصوص بوده است [21,18,3,20]. در روش OverFast [18]، لایه اتصال کامل (fc^1) برای پیش‌بینی مختصات باکس آموزش داده می‌شود. سپس لایه fc به یک لایه کانولوشن برای تشخیص اشیای کلاس مخصوص چندگانه تبدیل می‌شود. مدل‌های چند باکسی [3,20]، طرح‌های منطقه را از شبکه‌ای که آخرین لایه fc را بطور همزمان باکس‌های چندگانه (که برای تشخیص اشیای R-CNN [6] استفاده می‌شود)، پیش‌بینی می‌کند را ایجاد می‌کنند. شبکه پیشنهادی آنها برای یک تصویر تکی یا برش‌های تصویر بزرگ چندتایی نظیر 224×224 بکار می‌رود [20]. در این مقاله به بررسی دقیق روش پیشنهادیمان در مقایسه با روش‌های OverFeat و مالتی‌باکس^۲ می‌پردازیم.

روش محاسبه اشتراکی کانولوشن [18,7,2,5] بطور قابل ملاحظه‌ای در تشخیص بصری دقیق و موثر مورد توجه قرار گرفته است. مقاله Overfeat [18] ویژگی‌های کانولوشن را از تصویر اولیه به منظور طبقه‌بندی، مکانیابی و تشخیص محاسبه می‌کند. روش SPP^۳ [7] بروی ویژگی‌های کانولوشن به اشتراک گذاشته شده به منظور تشخیص بهینه اشیای مبتنی بر منطقه [7,16] و بخش‌بندی معنایی [2] پیشنهاد شده است. روش R-CNN سریع [5] آموزش انتها به انتهای ویژگی‌های کانولوشن به اشتراک گذاشته را مهیا کرده و سرعت و دقت محاسباتی مناسبی را ارائه می‌دهد.

¹ fully-connected

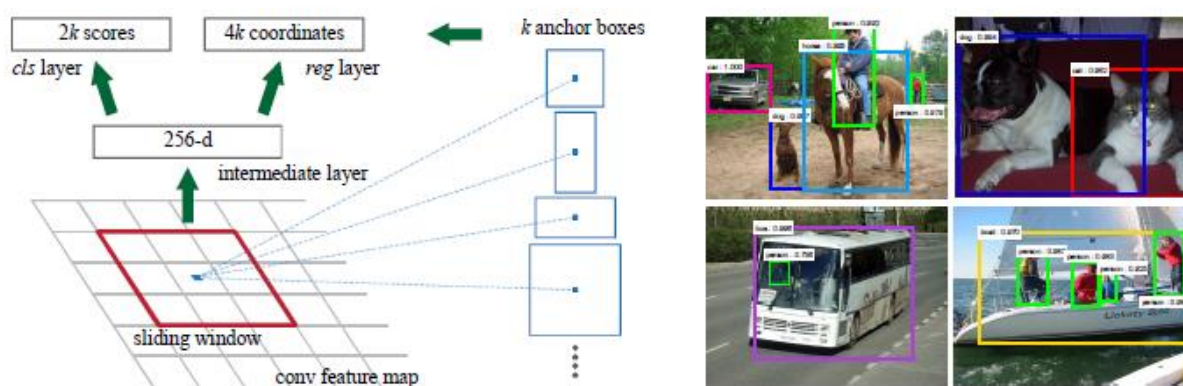
² MultiBox

³ Adaptively-sized pooling

۳. شبکه‌های طرح منطقه

یک شبکه طرح منطقه (RPN) یک تصویر در هر اندازه‌ای را به عنوان ورودی می‌گیرد و مجموعه‌ای از طرح‌های مستطیلی شیء را با امتیازات مربوط به هر شیء به عنوان خروجی ارائه می‌دهد. ما این فرآیند را با استفاده از شبکه کانولوشن کامل [1] که در این بخش به توضیح آن می‌پردازیم را مدلسازی می‌کنیم. از آنجا که خد ف نامحدود ما، اشتراک گذاری محاسبات با شبکه تشخیص اشیای R-CNN سریع می‌باشد، فرض می‌کنیم که هر دو شبکه مجموعه مشترکی از لایه‌های کانولوشن را به اشتراک می‌گذارند. در آزمایشات ما، مدل زیلر و فرگوس^۴ (ZF) [23] را که دارای ۵ لایه قابل اشتراک و مدل سیمونیان و زیسرمن^۵ (VGG) [19] که دارای ۱۳ لایه کانولوشن قابل اشتراک می‌باشند را مورد مطالعه قرار می‌دهیم.

به منظور ایجاد طرح‌های منطقه، شبکه کوچکی بروی ویژگی کانولوشن خروجی توسط آخرین لایه کانولوشن به اشتراک گذاشته شده را تولید می‌کنیم. این شبکه از نوع اتصال کامل^۶ است و بروی پنجره فضایی $n*n$ خروجی توسط آخرین لایه کانولوشن به اشتراک گذاشته شده ایجاد می‌شود. هر پنجره لغزنده^۷ به یک بردار کم‌بُعدتر (ZF برای 256-d و VGG برای 512-d) نسبت داده می‌شود. این بردار در دو لایه اتصال یافته کامل (یک لایه رگرسیون (reg) و یک لایه طبقه بندی (cls)) تغذیه می‌شود. ما در این مقاله از $n=3$ استفاده می‌کنیم؛ توجه شود که میدان پذیرنده موثر بروی تصویر ورودی بسیار بزرگ است (۱۷۱ پیکسل برای ZF و ۲۲۸ پیکسل برای VGG). این شبکه کوچک به عنوان یک موقعیت تکی در شکل ۱ سمت چپ ارائه شده است. توجه شود که از آنجا که شبکه کوچک در یک مُد پنجره لغزنده عمل می‌کند، لایه‌هایی که بطور کامل متصل شده‌اند، در میان تمام مکان‌های فضایی به اشتراک گذاشته می‌شوند. این ساختار معمولاً بصورت $n*n$ لایه کانولوشن که توسط دو لایه $1*1$ کاتولوشنی برای لایه‌های reg و cls تشکیل می‌شوند، بکار برده می‌شود. روش ReLUs [15] برای خروجی $n*n$ لایه کانولوشنی بکار برده می‌شود.



شکل ۱: سمت چپ: شبکه طرح منطقه (RPN). سمت راست: نمونه‌ای از تشخیص با استفاده از طرح‌های RPN بروی داده‌های تست PASCAL VOC 2007. روش مان اشیاء را در طیف گسترده‌ای از مقیاس‌ها و ضریب‌های منظری تشخیص می‌دهد.

⁴ Zeiler and Fergus

⁵ Simonyan and Zisserman

⁶ fully connected

⁷ sliding window

۴. لنگرهای تبدیل ثابت^۸

در مکان هر پنجره لغزنده، بطور همزمان طرح‌های منطقه k را پیش‌بینی می‌کنیم، بنابراین لایه reg دارای خروجی‌های $4k$ است که مختصات باکس k را پیش‌بینی می‌کند. لایه cls امتیازات $2k$ را که احتمال شیء بودن یا نبودن برای هر طرح را پیش‌بینی می‌کند. طرح‌های k نسبت به باکس‌های مرجع k پارامتریک می‌شوند که به آن‌ها لنگر گفته می‌شود. هر لنگر در پنجره لغزنده متمرکز شده و بر اساس مقیاس و ضریب منظر، سازماندهی می‌گردد. ما از ۳ مقیاس و ۳ ضریب منظر استفاده می‌کنیم که منجر به $k=9$ لنگر در هر موقعیت لغزنده می‌شود. برای یک ویژگی کانولوشن با اندازه $W*H$ (معمولا $(2,400)$)، بطور کلی $W H$ لنگر k وجود دارد. نسبت مهمی که روش ما به آن متکی است، "تبدیل ثابت" می‌باشد که هم در لنگرها و هم در توابعی که طرح‌های نسبی را نسبت به لنگرها محاسبه می‌کند، برقرار می‌باشد.

به منظور مقایسه، روش مالتی‌باکس [20] از روش کا-مینز^۹ برای تولید ۸۰۰ لنگر استفاده می‌کند که "تبدیل ثابت" نیستند. اگر هر لنگر، یک شیء را در یک تصویر شناسایی کند، آن طرح و تابع مشابه باید قادر به پیش‌بینی طرح در هر مکان باشد. علاوه بر این، از آنجا که لنگرهای مالتی‌باکس بصورت تبدیل ثابت و غیر متنوعی نیستند، این روش به یک لایه خروجی $800*(4+1)$ بُعدی نیاز دارد؛ این درحالی است که روش پیشنهادی ما نیازمند $9*(4+2)$ بُعد لایه خروجی می‌باشد. لایه‌های پیشنهادی ما دارای پارامترهای کمتر (۲۷ میلیون برای مالتی‌باکس با استفاده از گوگل‌نت^{۱۰} [20] و ۲,۴ میلیون برای RPN با استفاده از VGG-16) می‌باشد، بنابراین دارای ریسک کمتری در زمینه عملکرد مناسب بروی پایگاه داده‌های کوچک نظیر PASCAL VOC می‌باشد.

۵. تابع زیان برای آموزش طرح‌های منطقه

برای آموزش RPNS، یک برچسب کلاس باینری (که بیان‌کننده شیء بودن یا نبودن می‌باشد) به هر لنگر اختصاص می‌دهیم. بدین منظور برچسب مثبت را به دو نوع از لنگرها اختصاص می‌دهیم: (i) لنگر/لنگرهایی با بالاترین همپوشانی^{۱۱} IoU با هر باکس درست؛ یا (ii) هر لنگری که دارای IoU همپوشانی بالاتر از ۰,۷ با هر باکس درست. توجه شود که یک باکس درست تکی ممکن است برچسب‌های مثبتی را به لنگرهای چندتایی اختصاص دهد. همچنین، برچسب منفی را به لنگری غیر مثبتی اختصاص می‌دهیم که نسبت IoU اش برای تمام باکس‌ها کمتر از ۰,۳ باشد. لنگرهایی که نه مثبت هستند و نه منفی، مربوط به اهداف آموزشی نمی‌باشند.

با این تعاریف، تابع هدف را بر اساس تابع زیان چند منظوره در R-CNN سریع [5] کمینه می‌کنیم. تابع زیان برای یک تصویر بصورت زیر تعریف می‌شود:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*). \quad (1)$$

که در این رابطه i بیانگر شاخص لنگر در یک بخش کوچک و p_i بیان‌کننده احتمالی است که یک لنگر نشان‌دهنده یک شیء می‌باشد. در صورتی که برچسب p_i^* برابر با ۱ باشد، لنگر مثبت است و در صورتی که ۰ باشد، لنگر منفی می‌باشد. t_i بیان‌کننده برداری است که ۴ مختصات پارامتریک شده باکس مرزی پیش‌بینی شده می‌باشد و t_i^* باکس مرتبط با لنگر مثبت است. در زیان^{۱۲} طبقه بندی

⁸ Translation-Invariant Anchors

⁹ K-Means

¹⁰ GoogleNet

¹¹ Intersection over-Union

¹² Loss

بندی L_{cls} بیان‌کننده اختلاف دو کلاس (شیء بودن یا نبودن) می‌باشد. برای زیان رگرسیون، از عبارت $L_{reg}(t_i - t_i^*) = R(t_i - t_i^*)$ که R تابع زیان رُباست تعریف شده در [5] است. عبارت $p_i^* L_{reg}$ بیان‌کننده زیان رگرسیون است که تنها برای لنگرهای مثبت ($p_i^* = 1$) بصورت فعال و در غیر اینصورت ($p_i^* = 0$) غیر فعال می‌باشد. خروجی‌های لایه‌های cls و reg به ترتیب شامل $\{t_i\}$ و $\{p_i\}$ می‌باشد. این دو عبارت با N_{cls} و N_{reg} نرمال شده و با وزن λ^3 متعادل می‌شوند.

برای رگرسیون، پارامتریک کردن ۴ مختصات بصورت زیر انجام می‌شود [6]:

$$t_x = (x - x_a)/w_a, \quad t_y = (y - y_a)/h_a, \quad t_w = \log(w/w_a), \quad t_h = \log(h/h_a), \\ t_x^* = (x^* - x_a)/w_a, \quad t_y^* = (y^* - y_a)/h_a, \quad t_w^* = \log(w^*/w_a), \quad t_h^* = \log(h^*/h_a),$$

که x, y, w و h بیان‌کننده دو مختصات مرکز، عرض و ارتفاع باکس می‌باشند. متغیرهای x و x_a به ترتیب برای باکس، باکس لنگر و باکس ground-truth می‌باشند. این موضوع می‌تواند به عنوان رگرسیون باکس مرزی از یک باکس لنگر به نزدیک‌ترین باکس ground-truth در نظر گرفته شود.

با این وجود، روش ما رگرسیون باکس مرزی را با استفاده از روشی متفاوت از روش‌های پیشین مبتنی بر ویژگی انجام می‌دهد [7.5]. در مقالات [7.5]، رگرسیون باکس مرزی بر اساس ویژگی‌های حاصل شده از مناطق سایزبندی شده اختیاری است و وزن‌های رگرسیون توسط تمام این اندازه‌های مناطق به اشتراک گذاشته می‌شوند. در فرمول بندی مان، ویژگی‌های استفاده شده برای رگرسیون دارای اندازه فضایی یکسانی ($n \times n$) بروی ویژگی‌ها هستند. برای در نظرگیری اندازه‌های متغیر، مجموعه‌ای از رگرسیون‌های باکس مرزی k آموزش داده می‌شوند. هر رگرسیون، برای یک مقیاس و یک ضریب منظر (زاویه دید) در نظر گرفته می‌شود و رگرسیون‌های k ، وزنشان را به اشتراک نمی‌گذارند. با این حال، حتی اگر ویژگی‌ها دارای مقیاس/اندازه ثابت نیز باشند، همچنان پیش‌بینی باکس‌ها با اندازه‌های مختلف امکان‌پذیر می‌باشد.

۶. بهینه‌سازی

روش RPN که بطور طبیعی به عنوان یک شبکه کانولوشن کامل بکار برده می‌شود [14]، با استفاده از روش پس انتشار و SGD^{13} می‌تواند بصورت انتها به انتها آموزش ببیند [12]. همچنین می‌توان از استراتژی نمونه برداری "imagecentric" که از [5] تعیین می‌شود، این شبکه را آموزش دهیم. هر بخش کوچک، از تصویر جداگانه‌ای حاصل می‌شود که شامل تعداد زیادی لنگرهای مثبت و منفی می‌باشد. بهینه‌سازی تابع زیان تمام لنگرها نیز امکان‌پذیر است اما بسته به نمونه‌های منفی‌ای دارد که نسبت به نمونه‌های دیگر غالب هستند. بجای آن، ما از ۲۵۶ لنگر نمونه بصورت تصادفی در یک تصویر به منظور محاسبه تابع زیان بخش کوچک استفاده می‌کنیم که لنگرهای مثبت و منفی نمونه دارای نسبت ۱:۱ هستند. اگر مقداری کمتر از ۱۲۸ نمونه مثبت در یک تصویر وجود داشته باشد، بخش کوچکی را با نمونه‌های منفی در نظر می‌گیریم.

حال تمام لایه‌های جدید را بصورت تصادفی با ایجاد وزونی از توزیع گوسین حدود صفر با انحراف استاندارد ۰,۰۱ آماده‌سازی می‌کنیم. تمام لایه‌های دیگر (یعنی لایه‌های کانولوشن به اشتراک گذاشته شده) توسط آموزش مدل برای طبقه بندی ImageNet [17] به

¹³ stochastic gradient descent

عنوان یک عمل استاندارد [6] آماده سازی می شوند. ما تمام لایه های شبکه ZF ، conv3_1 و تا شبکه VGG را به منظور حفظ حافظه تنظیم می کنیم [5]. برای این کار از نرخ آموزش 0.001 برای 60 k بخش کوچک، و 0.001 برای 20 k نمونه کوچک بعدی بروی پایگاه داده PASCAL استفاده می کنیم. همچنین از مومنتوم 0.9 و وزن کاهشی 0.0005 [11] استفاده می کنیم. راه اندازی نرم افزاری ما از Caffe استفاده می کند [10].

7. اشتراک گذاری ویژگی های کانولوشن برای تشخیص اشیا و طرح منطقه

بنابراین، جدا از اینکه چگونه یک شبکه برای تولید طرح منطقه آموزش می بیند، بدون در نظرگیری تشخیص شیء مبتنی بر منطقه، روش CNN این طرح ها را آماده سازی می کند. برای شبکه تشخیص، ابتدا R-CNN سریع [5] را پیاده کرده و سپس الگوریتمی را که لایه های کانولوشن را بین RPN و R-CNN سریع به اشتراک میگذارد، توصیف می نمایم.

هر دو روش RPN و R-CNN سریع بطور مستقل آموزش می بینند و لایه های کانولوشن خودشان را از طریق روش های مختلف، اصلاح می نمایند. بنابراین نیاز به توسعه تکنیک و روشی داریم که بجای اینکه هر دو روش بصورت جداگانه آموزش ببینند، بتواند مجال تبادل لایه های کانولوشن بین هر دو شبکه را مهیا نماید. توجه شود که این کار به اندازه تعریف هر شبکه بصورت مجزا که شامل هر دو روش RPN و R-CNN سریع باشد و سپس با استفاده از روش پس انتشار بهینه سازی گردد، آسان نیست. دلیل این موضوع این است که آموزش R-CNN سریع به طرح های ثابت شیء بستگی دارد. در صورتی که آموزش R-CNN سریع بطور همزمان تغییراتی را ایجاد نماید، این امر منجر به همگرایی مکانیزم طرح خواهد شد. از آنجا که این نقطه بهینه یکی از موضوعات مورد توجه برای تحقیقات آینده است، ما از یک الگوریتم آموزشی 4 مرحله ای برای آموزش ویژگی های به اشتراک گذاشته شده از طریق بهینه سازی جایگزینی استفاده کرده ایم.

در گام اول، همانند آنچه که در بالا توضیح داده شد، RPN را آموزش می دهیم. این شبکه با مدل از پیش آموزش داده شده ImageNet آماده سازی شده و برای استفاده در طرح منطقه بصورت انتها به انتها تنظیم می گردد در گام دوم، شبکه تشخیص را بصورت جداگانه توسط R-CNN سریع با استفاده از طرح های تولید شده توسط گام اول RPN آموزش می دهیم. این شبکه تشخیص نیز توسط مدل از پیش آموزش دیده شده ImageNet آماده سازی می شود. در این مرحله، این دو شبکه لایه های کانولوشن را به اشتراک نمی گذارند. در گام سوم، از شبکه تشخیص دهنده برای آماده سازی آموزش RPN استفاده می کنیم، که لایه های اشتراک گذاشته شده را ثابت و تنها لایه هایی که مختص RPN می باشند را تنظیم بهینه می کند. حال، هر دو شبکه لایه های کانولوشن را به اشتراک می گذارند. در نهایت، با حفظ لایه های کانولوشن به اشتراک گذاشته شده بطور ثابت، لایه های fc روش R-CNN سریع را تنظیم بهینه می کنیم. به همین ترتیب، هر دو شبکه لایه های کانولوشن مشابهی را به اشتراک می گذارند و یک شبکه یکنواخت را تشکیل می دهند.

8. جزئیات پیاده سازی

برای اینکار، هر دو شبکه های تشخیص شیء و طرح منطقه را برای تصاویر مقیاس تکی آموزش و تست می کنیم [7.5]. همچنین مقیاس تصاویر را به نحوی تغییر می دهیم که کوچکترین بخش کناریشان برابر با $s=600$ پیکسل شود [5]. اگرچه ممکن است که استخراج ویژگی چند مقیاسه منجر به بهبود دقت و صحت تشخیص گردد، اما هیچگونه بهبودی در زمینه سرعت استخراج را به همراه نخواهد داشت [5]. همچنین باید توجه داشت که در شبکه های ZF و VGG، گام کلی در آخرین لایه کانولوشن برابر با 16 پیکسل بر اساس تغییر مقیاس تصاویر می باشد، و بنابراین تقریباً برابر با 10 پیکسل بروی تصاویر پایگاه داده PASCAL معمولی (375 * 500 ~) خواهد بود. اگرچه که دقت تشخیص در گام های کمتر بهتر خواهد بود، اما این گام بزرگتر نیز نتایج مطلوبی را به همراه خواهد داشت.

ششمین همایش بین‌المللی افق‌های نوین در مهندسی برق، کامپیوتر و مکانیک

6th International Conference on the New Horizons in Electrical Engineering, Computer and Mechanical

www.mhconf.ir

برای لنگرها، از ۳ مقیاس با مساحت‌های باکس برابر با 512^2 و $128^2, 256^2$ پیکسل و ۳ ضریب منظر (زاویه دید) برابر با ۱:۱، ۱:۲ و ۲:۱ استفاده می‌کنیم. توجه شود الگوریتم پیشنهادیمان اجازه استفاده از باکس‌های لنگری را که بزرگتر از میدان پذیرنده اصلی در زمان پیش‌بینی طرح‌های بزرگتر هستند را می‌دهد. این پیش‌بینی‌ها غیر ممکن بوده و خیلی دقیق نیستند و تنها ممکن است محدوده‌ای از شی را تنها در صورتی که وسط شیء قابل رویت باشد را شناسایی کنند. با این طراحی، راه حل پیشنهادی به ویژگی‌های چند مقیاسه یا پنجره‌های لغزنده چند مقیاسه برای پیش‌بینی مناطق بزرگ نیاز داشته که این امر سبب کاهش قابل توجهی از زمان راه‌اندازی و پردازش می‌گردد. شکل ۱ سمت راست، قابلیت روش پیشنهادیمان را برای طیف گسترده‌ای از مقیاس‌ها و زوایای دید را نشان می‌دهد. جدول زیر اندازه طرح میانگین آموزش داده شده را برای هر لنگر با استفاده از شبکه ZF را نشان می‌دهد. (اعداد برای $s=600$ می‌باشد).

anchor	$128^2, 2:1$	$128^2, 1:1$	$128^2, 1:2$	$256^2, 2:1$	$256^2, 1:1$	$256^2, 1:2$	$512^2, 2:1$	$512^2, 1:1$	$512^2, 1:2$
proposal	188×111	113×114	70×92	416×229	261×284	174×332	768×437	499×501	355×715

باکس‌های لنگر که مرزهای تصویر را برش می‌دهند نیازمند این می‌باشند که به درستی مدیریت شوند. در طول مرحله آموزش، تمام لنگرهای برش مرز را به نحوی نادیده می‌گیریم که منجر به ایجاد زیان در شبکه نشوند. برای یک تصویر معمولی 600×1000 ، معمولاً بطور کلی 20k لنگر $(9 * 40 * 60 \sim)$ وجود خواهد داشت. با صرف نظر از لنگرهای برش مرز، حدوداً 6k لنگر بر تصویر برای آموزش وجود خواهد داشت. اگر لایه‌های خروجی برش‌های مرزی در مرحله آموزش نادیده گرفته شوند، دشواری زیادی در تصحیح خطا بوجود خواهد آمد و مرحله آموزش، همگرا نخواهد شد. در طول مرحله آزمایش نیز، روش RPN کانونلوشن کامل را به کل تصویر اعمال می‌کنیم. این موضوع باکس‌های طرح برش مرزی‌ای را ایجاد خواهد نمود که آن را در مرز تصویر در نظر می‌گیریم.

تعدادی از طرح‌های RPN بشدت بروی یکدیگر همپوشانی ایجاد می‌کنند. به منظور کاهش این افزونگی، از روش NMS¹⁴ بروی طرح‌های منطقه بر اساس امتیازات cls شان استفاده می‌کنیم. همچنین حد آستانه IoU را برای NMS برابر با ۰٫۷ در نظر می‌گیریم که حدود 2k طرح منطقه بر تصویر را به همراه خواهد داشت. همانطور که نشان خواهیم داد، روش NMS به دقت تشخیص نهایی صدمه وارد نخواهد کرد اما منجر به کاهش تعداد طرح‌ها خواهد شد. پس از NMS، از تعداد بالای N طرح منطقه برای تشخیص استفاده می‌کنیم. در ادامه، از روش R-CNN سریع با استفاده از 2k طرح RPN استفاده خواهیم نمود اما تعداد زیادی از طرح‌ها را بطور همزمان مورد ارزیابی قرار می‌دهیم.

منابع

- [1] N. Chavali, H. Agrawal, A. Mahendru, and D. Batra. Object-Proposal Evaluation Protocol is 'Gameable'. arXiv: 1505.05836, 2015.
- [2] J. Dai, K. He, and J. Sun. Convolutional feature masking for joint object and stuff segmentation. In CVPR, 2015.
- [3] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In CVPR, 2014.
- [4] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results, 2007.

¹⁴ Non maximum suppression

- [5] R. Girshick. Fast R-CNN. arXiv:1504.08083, 2015.
- [6] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In CVPR, 2014.
- [7] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In ECCV. 2014.
- [8] J. Hosang, R. Benenson, P. Dollar, and B. Schiele. What makes for effective detection proposals? arXiv:1502.05082, 2015.
- [9] J. Hosang, R. Benenson, and B. Schiele. How good are detection proposals, really? In BMVC, 2014.
- [10] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. arXiv:1408.5093, 2014.
- [11] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.
- [12] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. Neural computation, 1989.
- [13] K. Lenc and A. Vedaldi. R-CNN minus R. arXiv:1506.06981, 2015.
- [14] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In CVPR, 2015.
- [15] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In ICML, 2010.
- [16] S. Ren, K. He, R. Girshick, X. Zhang, and J. Sun. Object detection networks on convolutional feature maps. arXiv:1504.06066, 2015.
- [17] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. arXiv:1409.0575, 2014.
- [18] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In ICLR, 2014.
- [19] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In ICLR, 2015.
- [20] C. Szegedy, S. Reed, D. Erhan, and D. Anguelov. Scalable, high-quality object detection. arXiv:1412.1441v2, 2015.
- [21] C. Szegedy, A. Toshev, and D. Erhan. Deep neural networks for object detection. In NIPS, 2013.
- [22] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. IJCV, 2013.
- [23] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional neural networks. In ECCV, 2014.
- [24] C. L. Zitnick and P. Dollar. Edge boxes: Locating object proposals from edges. In ECCV, 2014.