



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر

گزارش درس شبکه‌های عصبی و یادگیری عمیق

شبکه کانولوشنی نامتراکم

نگارش
مینو دولت‌آبادی

استاد درس
دکتر رضا صفابخش

بهمن ۱۴۰۱

چکیده

علی‌رغم اینکه شبکه‌های عصبی کانولوشنی عمیق عملکرد پیشرفته‌ای را در بسیاری از وظایف بینایی رایانه نشان می‌دهند، ماهیت غنی از پارامتر و محاسبات آن‌ها به طور قابل توجهی مانع استفاده کارآمد از این شبکه‌ها در دستگاه‌های با پهنای باند و توان محدود می‌شود. برای این منظور، در سال‌های اخیر علاقه زیادی در زمینه به حداقل رساندن حافظه و هزینه‌های محاسباتی استنتاج شبکه‌های عصبی مشاهده شده است. یکی از روش‌های مقابله با این مسائل استفاده از شبکه‌های عصبی کانولوشنی نامتراکم است. استفاده از این شبکه‌ها کاملاً کارآمد است زیرا به کمک آن‌ها می‌توان افزودنی در پارامترها را کاهش داد و از طرفی نیازی به اسکن همه پیکسل‌ها یا وکس‌های فضایی نیست و فقط کانولوشن برای عناصر غیر صفر محاسبه می‌شود. لذا در این گزارش به بررسی شبکه‌های عصبی کانولوشنی نامتراکم و تعمیم حاصل از آن‌ها پرداخته می‌شود.

واژه‌های کلیدی:

شبکه عصبی کانولوشنی، مینکوفسکی، کوانتیزاسیون، شبکه عصبی نامتراکم، تشخیص اشیاء سه‌بعدی، شبکه‌های عصبی باقی‌ماندگی

صفحه	فهرست مطالب
۷	فصل اول مقدمه.....
۸	مقدمه.....
۱۰	فصل دوم شبکه‌های عصبی کانولوشنی ساختار یافته به هم پیوسته نامتراکم.....
۱۱	مقدمه.....
۱۱	۱-۲- معرفی.....
۱۲	۲-۲- ساخت یک ماتریس هسته متراکم تشکیل شده.....
۱۴	۳-۲- نتایج.....
۱۶	فصل سوم گسسته‌سازی متمرکز برای شبکه‌های عصبی کانولوشنی نامتراکم.....
۱۷	مقدمه.....
۱۷	۱-۳- روش.....
۱۷	۱-۱-۳- مقدمات: انتقال کوانتیزاسیون.....
۱۸	۳-۱-۲- طراحی تابع کوانتیزاسیون متمرکز.....
۱۹	۳-۱-۳- بهینه‌سازی کوانتیزاسیون متمرکز $Q[\theta]$
۲۱	۴-۱-۳- انتخاب کوانتیزاسیون مناسب.....
۲۲	۲-۳- نتایج.....
۲۴	فصل چهارم شبکه‌های عصبی کانولوشنی فوق نامتراکم.....
۲۵	مقدمه.....
۲۵	۱-۴- هسته کافنا.....
۲۶	۲-۴- پیاده‌سازی عملگر کافنا.....
۲۷	۳-۴- مازول پایه کافنا و کافنا-نت.....
۲۸	۴-۴- مقایسه دو شبکه عصبی تلفن همراه پیشرفته.....
۳۰	فصل پنجم شبکه کانولوشنی نامتراکم سریع.....
۳۱	مقدمه.....
۳۱	۱-۵- پراکندگی شبکه‌ها.....
۳۲	۲-۵- پیاده‌سازی هسته.....
۳۲	۳-۵- بررسی مفهومی.....

فصل ششم شبکه‌های کانولوشنی نامتراکم کانونی برای تشخیص اشیاء سه بعدی ۳۴

مقدمه ۳۵

۱-۶- معرفی ۳۵

۲-۶- مراحل اجرا ۳۷

۳-۶- ادغام کاناکا ۳۸

۱-۳-۶- استخراج ویژگی ۳۸

۲-۳-۶- تراز ویژگی ۳۸

۴-۶- شبکه‌های کاناکا ۳۹

فصل هفتم شبکه‌های عصبی کانولوشنی فضایی-زمانی ۴ بعدی: شبکه‌های عصبی

کانولوشنی مینکوفسکی ۴۰

مقدمه ۴۱

۱-۷- معرفی ۴۱

۲-۷- کانولوشن و تنسور نامتراکم ۴۱

۳-۷- کانولوشن نامتراکم تعمیم یافته ۴۲

۴-۷- مهندسی مینکوفسکی ۴۳

۱-۴-۷- کوانتیزاسیون تنسور نامتراکم ۴۳

۲-۴-۷- ادغام ماکسیمم ۴۵

۳-۴-۷- ادغام مجموع، ادغام میانگین و ادغام سراسری ۴۶

۵-۷- توابع غیر فضایی ۴۶

۶-۷- شبکه‌های عصبی کانولوشنی مینکوفسکی ۴۶

۷-۷- هسته تسراکت و هسته هیبریدی ۴۷

۸-۷- شبکه‌های باقی ماندگی مینکوفسکی ۴۷

۷-۹- CRF-سه گانه ثابت ۴۸

۱۰-۷- تعریف ۴۹

۱۱-۷- استنتاج متغیر ۴۹

۱۲-۷- یادگیری با کانولوشن نامتراکم ۷ بعدی ۵۰

۱۳-۷- نتایج ۵۱

منابع و مراجع ۵۲

فهرست اشکال

صفحه

شکل ۱-۲-۱: IGC - V2	شبکه عصبی کانولوشنی ساختار یافته به هم پیوسته نامتراکم	۱۳
شکل ۲-۲	تأثیر تعداد لایه‌های L	۱۴
شکل ۲-۳	عملکرد شبکه بر اثر تغییر عمق و عرض	۱۵
شکل ۱-۳	توزیع وزن ۸ لایه اول ResNet-18 در ImageNet	۱۸
شکل ۲-۳	فرآیند گام به گام کوانتیزه‌سازی مجدد وزن‌های هرس نشده در block3f/conv1	۲۰
شکل ۳-۳	توزیع وزن لایه block22/conv1 در ResNet-18 نامتراکم آموزش دیده در ImageNet	۲۱
شکل ۴-۳	تأثیر مقادیر مختلف آستانه بر فاصله و سترین	۲۲
شکل ۱-۴	مقایسه هسته کانولوشن معمولی و کافنا	۲۵
شکل ۲-۴	فرآیند عملیاتی انتقال	۲۷
شکل ۱-۵	کانولوشن نامتراکم 1×1 به عنوان SpMM	۳۳
شکل ۳-۵	تأثیر اندازه بلوک بر صحت top-1	۳۳
شکل ۱-۶	چارچوب کانولوشن نامتراکم کانونی و گسترش چند وجهی آن	۳۷
شکل ۱-۷	پیش بینی دوبعدی ابرمکعب در ابعاد مختلف	۴۷
شکل ۲-۷	معماری ResNet18 (سمت چپ) و MinkowskiNet18 (راست)	۴۸

صفحه

فهرست جداول

جدول ۱-۳: دقت، پراکندگی و CR های فشرده‌سازی متمرکز در مدل‌های ImageNet.....	۲۳
جدول ۱-۴: ساختار شبکه	۲۸
جدول ۲-۴: مقایسه دو شبکه عصبی تلفن همراه پیشرفته.....	۲۸
جدول ۱-۷: نتایج تقسیم‌بندی در مجموعه داده ۴ بعدی synthia.....	۵۱
جدول ۲-۷: مقایسه نتایج حاصل از الگوریتم‌های مختلف.....	۵۱

فهرست اختصارات

عنوان اختصاری	عنوان کامل
کانا	کانولوشنی نامتراکم
کاسابنا	کانولوشنی ساختار یافته به هم پیوسته نامتراکم
کافنا	کانولوشنی فوق نامتراکم
کافنا-نت	کانولوشنی فوق نامتراکم-نت
کاناس	کانولوشنی نامتراکم سریع
کاناکا	کانولوشنی نامتراکم قانونی

فصل اول

مقدمه

مقدمه

امروزه شبکه‌های عصبی کانولوشنی عمیق^۱ با اندازه مدل کوچک، هزینه محاسبات پایین و در عین حال دقت بالا به ویژه در دستگاه‌های تلفن همراه بسیار مورد توجه قرار گرفته‌اند. جهت دستیابی به چنین شبکه‌هایی تلاش‌هایی صورت گرفته است:

- فشرده‌سازی شبکه: فشرده‌سازی مدل از پیش آموزش دیده شده با تجزیه ماتریس هسته کانولوشنی یا حذف اتصالات یا کانال‌ها برای حذف افزونگی^۲
- طراحی معماری: طراحی هسته‌های کوچک، هسته‌های نامتراکم یا استفاده از محصول هسته‌هایی که اهمیت بیشتری دارند، جهت نزدیک شدن به یک تک هسته و آموزش شبکه‌ها از ابتدا.

تمرکز این کار مطالعاتی در طراحی معماری با استفاده از محصول هسته‌هایی که اهمیت بیشتری دارند برای ترکیب یک هسته است. برای تحقق این امر دو روند اصلی وجود دارد :

○ ضرب هسته‌های با رتبه پایین (ماتریس) برای تقریب یک هسته با رتبه بالا، به عنوان مثال ماژول های تنگنا^۳ [۱].

○ ضرب ماتریس‌های نامتراکم، که اخیراً تلاش‌های تحقیقاتی را به خود جلب کرده است [۲] [۳].

الگوریتم‌های اخیراً توسعه یافته، مانند کانولوشن گروهی درهم^۴ [۲] و Xception [۳]، یک هسته متراکم را با استفاده از حاصل ضرب دو هسته ساختاریافته - نامتراکم می‌سازند. مشاهده شده است که یکی از این دو هسته را می‌توان بیشتر تقریب زد. با انگیزه مشاهده اینکه کانولوشن‌های موجود در یک کانولوشن گروهی در IGC را می‌توان به همان روش بیشتر تجزیه کرد، ابتداءً فصل دوم کاسابنا (IGC-V2) معرفی می‌شود [۴]. سپس در بخش سوم به معرفی یک استراتژی کوانتیزه‌سازی جدید پرداخته می‌شود [۵]. در بخش چهارم و پنجم به ترتیب دو تعمیم از این شبکه‌ها با نام‌های کافنا و کاناس

¹ Deep convolutional neural networks

² redundancy

³ bottleneck

⁴ interleaved group convolution

معرفی می‌شوند [۶] [۷]. در بخش ششم توضیح داده می‌شود که چگونه این نوع از شبکه‌ها در تشخیص اشیاء سه‌بعدی مفید هستند [۸] و در نهایت در بخش آخر نوع جدیدی از این شبکه‌ها به نام شبکه‌های عصبی کانولوشنی مینکوفسکی معرفی و بررسی می‌شوند [۹].

فصل دوم

شبکه‌های عصبی کانولوشنی ساختار یافته به هم پیوسته نامتراکم^۵

⁵ Interleaved Structured Sparse Convolutional Neural Networks

مقدمه

در این بخش قرار است مسئله طراحی معماری شبکه‌های عصبی کانولوشنی با توجه به حذف افزونگی در هسته‌های کانولوشنی مورد بررسی قرار گیرد. در این بخش کاسابنا (IGC-V2) (شکل ۱-۲) معرفی می‌شود. این شبکه، کانولوشن‌های گروهی که از دو هسته نامتراکم ساختاریافته تشکیل شده است را به محصول هسته‌های نامتراکم ساختار یافته‌تر تعمیم می‌دهد و افزونگی را بیشتر حذف می‌کند.

۲-۱- معرفی

عملیات در یک لایه کانولوشن در شبکه‌های عصبی کانولوشنی به یک عملیات ضرب بردار-ماتریسی در هر مکان متکی است:

$$\mathbf{y} = \mathbf{W}\mathbf{x} \quad (1-1)$$

در این جا ورودی \mathbf{x} متناظر با یک وصله در اطراف مکان کانال‌های ورودی و یک بردار SC_i -بعدی است، به‌طوری که S برابر با اندازه هسته (برای مثال 3×3) و C_i برابر با تعداد کانال‌های ورودی است. خروجی y نیز یک بردار C_o -بعدی است و همچنین متناظر با تعداد کانال‌های خروجی است. \mathbf{w} از هسته‌های C_o تشکیل می‌شود و هر ردیف متناظر با هسته کانولوشن است. در این مقاله جهت افزایش شفافیت فرض می‌شود: $C_o = C_i = C$.

بلوک W_i^g و کانولوشن 1×1 در Xception متراکم هستند و می‌تواند با ضرب ماتریس‌های نامتراکم تشکیل شود. در نتیجه جهت بیشتر از بین بردن افزونگی و صرفه‌جویی در ذخیره‌سازی و زمان چنین فرایندی را می‌توان تعداد بارهای بیشتری تکرار کرد.

$$\begin{aligned} \mathbf{y} &= \mathbf{P}_L \mathbf{W}_L \mathbf{P}_{L-1} \mathbf{W}_{L-1} \dots \mathbf{P}_1 \mathbf{W}_1 \mathbf{x} \\ &= \left(\prod_{l=L}^1 \mathbf{P}_l \mathbf{W}_l \right) \mathbf{x}. \end{aligned} \quad (2-1)$$

در معادله مذکور $P_l W_l$ یک ماتریس نامتراکم است. P_l یک ماتریس جایگشتی^۶ است، و نقش آن مرتب کردن مجدد کانال‌ها به گونه‌ای است که W_l یک ماتریس بلوکی نامتراکم باشد، و متناظر با کانولوشن گروه l ام است، که در آن تعداد کانال‌ها در همه شاخه‌ها یکسان و برابر K_l هستند.

۲-۲- ساخت یک ماتریس هسته متراکم تشکیل شده

در ابتدا شرط زیر که تعمیم یافته کانولوشن‌های گروهی به هم پیوسته است به عنوان یک قاعده برای ساخت کانولوشن‌های گروه l به گونه‌ای که ماتریس هسته کانولوشن تشکیل شده متراکم باشد، مطرح می‌شود.

شرط تعمیم یافته: $\forall m, (W_L \prod_{l=L-1}^m P_l W_l)$ و $(W_{L-1} \prod_{l=L-2}^m P_l W_l)$ متناظر با دو گروه کانولوشنی هستند. دو کانولوشن گروهی فوق را مکمل یکدیگر گویند اگر کانال‌هایی که در یک شاخه در یک کانولوشن گروهی قرار دارند در شاخه‌های مختلف نهفته باشند و از همه شاخه‌های کانولوشن گروه دیگر آمده باشند.

اکنون سوالی که مطرح می‌شود این است که چه زمانی مقدار پارامترها کوچکترین است؟

می‌توان گفت که تعداد پارامترهای کانولوشن گروه L که در معادله (۲-۱) ارائه شده است و در شرط تعمیم یافته صدق می‌کند، پاسخ این سوال است.

تعداد پارامترها در کانولوشن گروه l ام CK_l برای یک کانولوشن 1×1 و CSK_l برای کانولوشن فضایی مانند $s = 3 \times 3$ است. جهت استفاده از پارامترهای کمتر تنها یک کانولوشن فضایی گروهی وجود دارد و بقیه 1×1 هستند. کانولوشن فضایی در هر کانولوشن گروهی نهفته است و بدون تأثیر بر تجزیه و تحلیل، فرض می‌شود که در کانولوشن گروه اول قرار دارد. بنابراین، تعداد کل پارامترهای Q ، تعداد کوچکتری از پارامترها که در ماتریس‌های جایگشت نادیده گرفته شده‌اند، برابر است با:

$$Q = C \sum_{l=2}^L K_l + CSK_1 \quad (3-2)$$

^۶ permutation matrix

با توجه به نابرابری جنسن^۷:

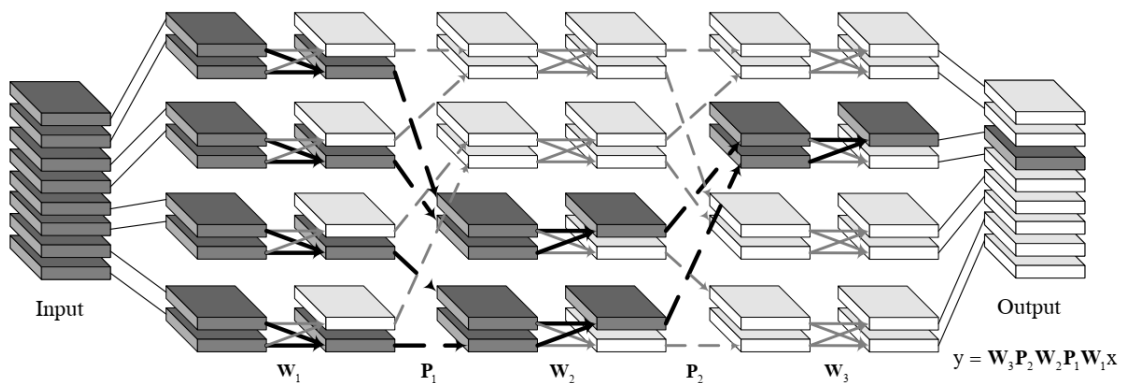
$$\begin{aligned}
 Q &= C \sum_{l=2}^L K_l + CSK_1 \quad (۴-۲) \\
 &\geq CL \left(SK_1 \prod_{l=2}^L K_l \right)^{\frac{1}{L}} \\
 &= CL(SC)^{\frac{1}{L}}
 \end{aligned}$$

در این جا، برابری از خط دوم تا خط سوم به دلیل معادله (۲-۲) برقرار است. برابری در خط دوم برقرار

است، یعنی $Q = CL(SC)^{\frac{1}{L}}$ زمانی که شرط تعادل زیر برآورده شود:

$$SK_1 = K_2 = \dots = K_L \left(= (SC)^{\frac{1}{L}} \right) \quad (۵-۲)$$

علاوه بر این، برای انتخاب L که کمترین مقدار پارامتر را به دست می دهد، یک تحلیل تقریبی با در نظر گرفتن مشتق Q با توجه به L ارائه می شود.

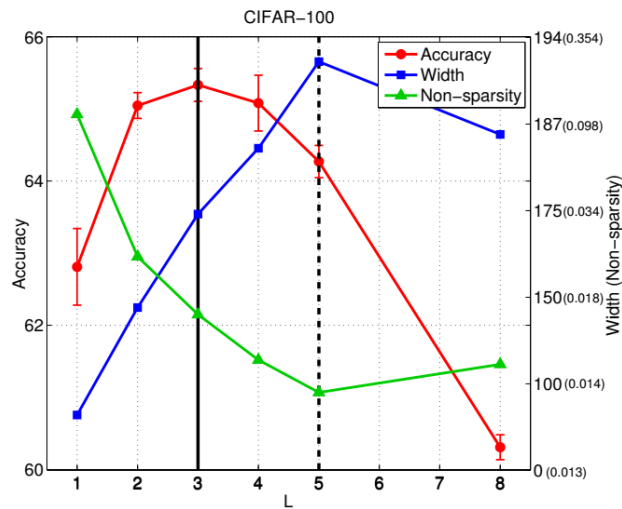


شکل ۲-۱-۲ - IGC: شبکه عصبی کانولوشنی ساختار یافته به هم پیوسته نامتراکم.

⁷ Jensen's inequality

۳-۲- نتایج:

در نمودار زیر نشان داده می‌شود که چگونه تعداد لایه‌های L بر عملکرد CIFAR-100 تأثیر می‌گذارد. حداکثر دقت در مقداری L از به‌دست می‌آید که در آن عرض و درجه عدم پراکندگی به تعادل می‌رسند.

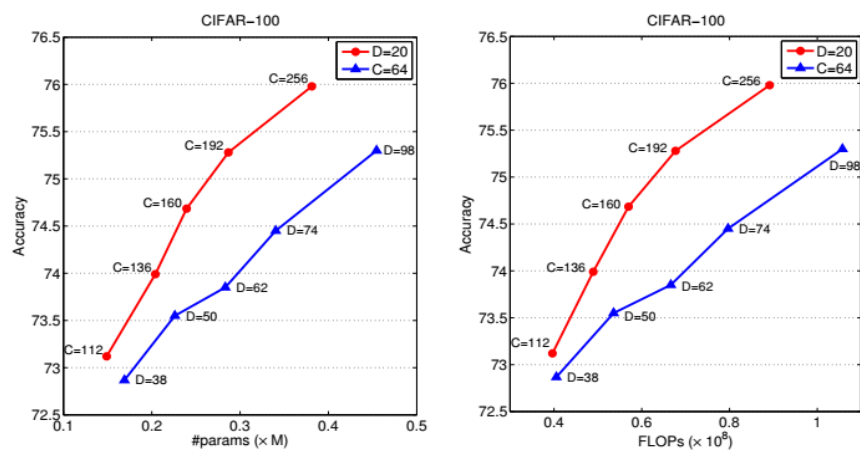


شکل ۲-۱- تأثیر تعداد لایه‌های L .

نمودار زیر چگونگی تغییر عملکرد را زمانی که شبکه ما عمیق‌تر یا گسترده‌تر می‌شود، نشان می‌دهد. در این مقاله از ساختار شبکه $IGC-V2^* (Cx)$ در جدول ۲-۱ استفاده می‌شود و آزمایش‌ها با ابعاد مختلف (D) و با عرض (C) انجام می‌شوند. با توجه به نمودار هم عریض‌تر شدن و هم عمیق‌تر شدن کارایی را افزایش می‌دهد و اثر عریض‌تر شدن بیشتر از عمیق‌تر شدن است.

جدول ۲-۱: معماری استفاده شده شبکه‌هایی که در آزمایش‌ها.

Output size	Xception (Cx)	IGC-V1 (Cx)	IGC-V2 (Cx)	IGC-V2* (Cx)
32×32	$(3 \times 3, x)$	$(3 \times 3, x)$	$(3 \times 3, x)$	$(3 \times 3, 64)$
32×32	$x \times (3 \times 3, 1)$ $(1 \times 1, x) \times B$	$\frac{x}{2} \times (3 \times 3, 2)$ $2 \times (1 \times 1, \frac{x}{2}) \times B$	$x \times (3 \times 3, 1)$ $L-1, x, (1 \times 1, K_{s1}) \times B$	$x \times (3 \times 3, 1)$ $L^*-1, x, (1 \times 1, K) \times B$
16×16	$2x \times (3 \times 3, 1)$ $(1 \times 1, 2x) \times B$	$x \times (3 \times 3, 2)$ $2 \times (1 \times 1, x) \times B$	$2x \times (3 \times 3, 1)$ $L-1, 2x, (1 \times 1, K_{s2}) \times B$	$2x \times (3 \times 3, 1)$ $L^*-1, 2x, (1 \times 1, K) \times B$
8×8	$4x \times (3 \times 3, 1)$ $(1 \times 1, 4x) \times B$	$2x \times (3 \times 3, 2)$ $2 \times (1 \times 1, 2x) \times B$	$4x \times (3 \times 3, 1)$ $L-1, 4x, (1 \times 1, K_{s3}) \times B$	$4x \times (3 \times 3, 1)$ $L^*-1, 4x, (1 \times 1, K) \times B$
1×1	average pool, fc, softmax			
Depth	$3B+2$			



شکل ۳-۰ - عملکرد شبکه بر اثر تغییر عمق و عرض.

فصل سوم

گسسته‌سازی متمرکز برای شبکه‌های عصبی کانولوشنی نامتراکم^۸

^۸ Focused Quantization for Sparse CNNs

مقدمه

در این بخش، به ویژگی‌های آماری شبکه‌های عصبی کانولوشنی پراکنده پرداخته می‌شود و کوانتیزاسیون متمرکز را ارائه می‌دهد، یک استراتژی کوانتیزه‌سازی جدید مبتنی بر توان دو مقادیر، که از توزیع وزن پس از هرس دقیق استفاده می‌کند. روش پیشنهادی به صورت پویا موثرترین نمایش عددی را برای وزن‌ها در لایه‌هایی با پراکندگی‌های مختلف کشف می‌کند که به طور قابل توجهی اندازه مدل را کاهش می‌دهد.

۳-۱- روش

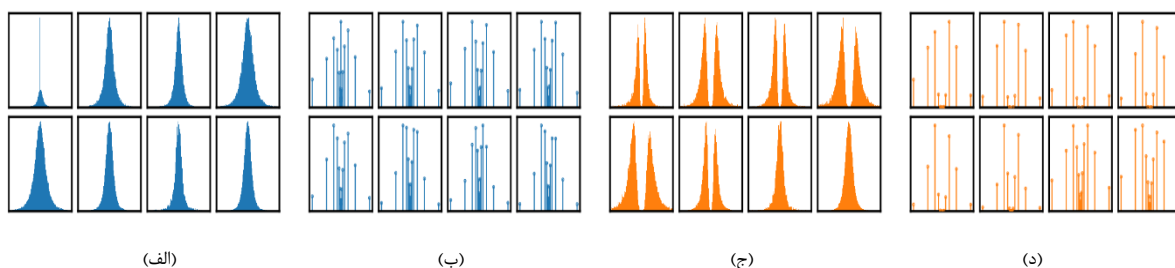
۳-۱-۱- مقدمات: انتقال کوانتیزاسیون^۹

انتقال کوانتیزاسیون یک طرح کوانتیزه‌سازی است که مقادیر وزن را به توان‌های دو یا صفر محدود می‌کند. یک مقدار قابل نمایش برای انتقال کوانتیزاسیون $(k + 2)$ بیتی به صورت زیر است:

$$v = s \cdot 2^{e-b} \quad (۱-۳)$$

که در آن $s = \{-1.0.1\}$ نشان‌دهنده صفر یا علامت مقدار است، e یک عدد صحیح محدود به بازه $[0.2^k - 1]$ و b بایاس است، یک مقدار لایه‌ای ثابت که میزان بزرگی مقدار کوانتیزه را مدیریت می‌کند. در این مقاله از $\hat{\theta} = Q_{n,b}^{shift}[\theta]$ برای نشان دادن انتقال کوانتیزاسیون n بیتی با بایاس b از مقدار وزن θ به نزدیکترین مقدار قابل نمایش $\hat{\theta}$ استفاده شده است. همان‌طور که در شکل ۳-۱ نشان داده شده است، انتقال کوانتیزاسیون در لایه‌های نامتراکم موجب استفاده جزئی از محدوده مقادیر قابل مشاهده می‌شود، یعنی توزیع حاصل پس از کوانتیزه‌سازی $q_{n,b}^{shift}(\theta)$ تقریبی ضعیف از توزیع وزن لایه اصلی $p(\theta|\mathcal{D})$ است، که در آن \mathcal{D} مجموعه داده آموزشی است.

^۹ Shift quantization



شکل ۳-۱- توزیع وزن ۸ لایه اول *ResNet-18* در *ImageNet* (الف) توزیع وزن لایه‌ها، (ج) توزیع وزن لایه‌ها (بدون احتساب صفرها) برای یک نوع نامتراکم. (ب) و (د) به ترتیب توزیع وزن نمودارهای سمت چپ انتقال کوانتیزاسیون ۵ بیتی.

۳-۱-۲- طراحی تابع کوانتیزاسیون متمرکز

به طور شهودی مطلوب است که کوانتیزاسیون بر روی مناطق با احتمال بالا در توزیع وزن در لایه‌های نامتراکم متمرکز شود. در اثر انجام این کار، ما می‌توانیم توزیع وزن‌های کوانتیزه شده را با مقادیر اصلی مطابقت دهیم و در نتیجه همزمان خطاهای گرد کردن کوچک‌تر می‌شوند. کوانتیزاسیون غیرمتمرکز $Q[\theta]$ به طور خاص برای این منظور طراحی شده است و به صورت لایه‌ای اعمال می‌شود. با فرض اینکه $\theta \in \Theta$ مقدار وزن یک لایه کانولوشن باشد، می‌توان $Q[\theta]$ را به صورت زیر تعریف کرد:

$$Q[\theta] = z_{\theta} \alpha \sum_{c \in C} \delta_{c, m_{\theta}} Q_c^{\text{rec}}[\theta]. \text{ where } Q_c^{\text{rec}}[\theta] = Q_{\text{n.b}}^{\text{shift}} \left[\frac{\theta - \mu_c}{\sigma_c} \right] \sigma_c + \mu_c \quad (3-2)$$

در این جا z_{θ} یک ثابت از پیش تعیین شده $\{0.1\}$ دودویی مقدار است که برای نشان دادن اینکه θ هرس شده است یا خیر، و برای تنظیم وزن‌های هرس شده روی 0 استفاده می‌شود. مجموعه مولفه‌های $c \in C$ مکان‌هایی را جهت تمرکز برای اعمال کوانتیزه‌سازی تعیین می‌کند که هر کدام با میانگین μ_c و انحراف معیار σ_c مشخص شده‌اند. دلتای کرونکر^{۱۰} $\delta_{c, m_{\theta}}$ زمانی که $c = m_{\theta}$ برابر با یک یا در غیر

¹⁰ Kronecker delta

این صورت صفر ارزیابی می شود. به عبارت دیگر، ثابت $m_{\theta} \in C$ تعیین می کند که کدام جزء در C برای کوانتیزه کردن θ استفاده شود. در نهایت، $Q_{rec}^c[\theta]$ به صورت محلی مولفه c را با انتقال کوانتیزاسیون، کوانتیزه می کند. علاوه بر این، یک عامل مقیاس پذیری قابل یادگیری از نظر لایه به نام α معرفی شده است که مقدار ۱ را می گیرد، و به طور تجربی دقت را بهبود می بخشد. بنابراین، با تنظیم μ_c و σ_c هر جزء c و یافتن تخصیص مناسب وزن ها به مؤلفه ها، توزیع وزن کوانتیزه شده $q_{\phi}(\theta)$ می تواند بسیار نزدیک با توزیع اصلی باشد، جایی که Φ به عنوان مخفف برای نشان دادن ابرپارامترهای مربوطه مانند μ_c و σ_c استفاده می شود. در بخش زیر توضیح داده خواهد شد که چگونه می توان این ابرپارامترها را به طور موثر بهینه کرد.

۳-۱-۳- بهینه سازی کوانتیزاسیون متمرکز $Q[\theta]$

ابرپارامترهای μ_c و σ_c در کوانتیزه سازی مجدد متمرکز می توانند با اعمال فرآیند دو مرحله ای زیر که به صورت لایه ای است بهینه شوند، ابتدا مناطق با احتمال بالا را شناسایی می کند (نخستین بلوک در شکل ۳-۲)، سپس آن ها را به صورت محلی با انتقال کوانتیزاسیون (بلوک دوم و سوم در شکل ۳-۲) کوانتیزه می کند. در ابتدا متوجه می شویم که توزیع وزن به طور کلی شبیه ترکیبی از توزیع های گوسی است. لذا یافتن یک مدل آمیخته گوسی^{۱۱} $q_{\phi}^{mix}(\theta)$ که توزیع اصلی $p(\theta|\mathcal{D})$ را جهت بهینه سازی دقیق هدف بالا تقریب می زند، موثرتر است:

$$q_{\phi}^{mix}(\theta) = \sum_{c \in C} \lambda_c f(\theta | \mu_c, \sigma_c) \quad (3-3)$$

به طوری که $f(\theta | \mu_c, \sigma_c)$ تابع چگالی احتمال توزیع گاوسی $N(\mu_c, \sigma_c)$ می باشد، λ_c غیرمنفی است و وزن ترکیبی مولفه c ام و $\sum_{c \in C} \lambda_c = 1$ را تعیین می کند. در این مرحله بایستی ابرپارامترهای μ_c ، σ_c و λ_c موجود در Φ را یافت به طوری که $q_{\phi}^{mix}(\theta)$ را به حداکثر می رساند با توجه به اینکه $\theta \sim p(\theta|\mathcal{D})$. این مسئله به عنوان تخمین حداکثر درستنمایی^{۱۲} (MLE) شناخته می شود، و MLE را می توان به طور

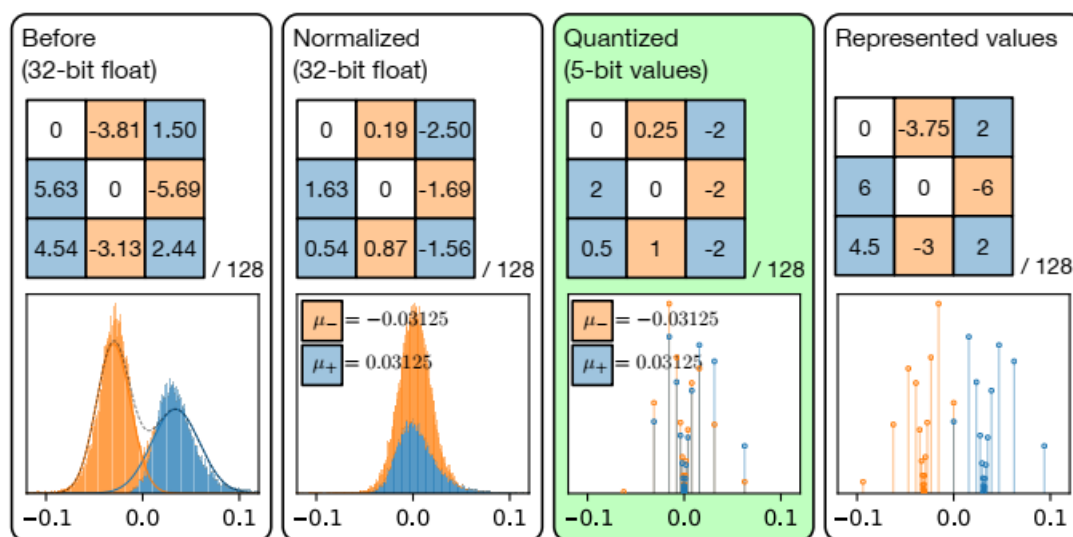
¹¹ gaussian mixture model

¹² maximum likelihood estimate

موثر توسط الگوریتم حداکثرسازی امید ریاضی^{۱۳} (EM) محاسبه کرد. در عمل، استفاده از دو مولفه گاوسی، $C = \{+, -\}$ برای شناسایی مناطق با احتمال بالا در توزیع وزن کافی است. جهت همگرایی سریعتر EM، μ_+ ، σ_+ و μ_- ، σ_- به ترتیب با میانگین و انحراف معیار مقادیر منفی و مثبت در وزن لایه λ_- و λ_+ با $\frac{1}{2}$ مقداردهی اولیه می‌شوند.

سپس m_θ از مدل آمیخته تولید می‌شود، به طوری که به صورت جداگانه مولفه مورد استفاده برای هر وزن را انتخاب می‌کند. برای این کار، m_θ برای هر θ با نمونه‌برداری از یک توزیع طبقه‌ای ارزیابی می‌شود که در آن احتمال اختصاص یک مولفه c به m_θ یعنی $p(m_\theta = c)$ ، $q_\phi^{mix}(\theta)$ است.

در نهایت، ثابت b روی یک مقدار توان دو تنظیم می‌شود، برای اطمینان از اینکه $q_{n,b}^{shift}$ اجازه می‌دهد تا حداکثر نسبت $\frac{1}{2^{n+1}}$ از مقادیر سرریز شود و آنها را به حداکثر مقدار قابل نمایش تقریب می‌زند. در عمل، این انتخاب اکتشافی نتیجه بهتری از سطوح کوانتیزاسیون ارائه شده توسط انتقال کوانتیزاسیون نسبت به عدم اجازه سرریزها می‌دهد. پس از تعیین تمام ابرپارامترهای مربوطه با روشی که در بالا توضیح داده شد، $\hat{\theta} = Q[\theta]$ را می‌توان برای کوانتیزه کردن وزن لایه‌ها ارزیابی کرد.

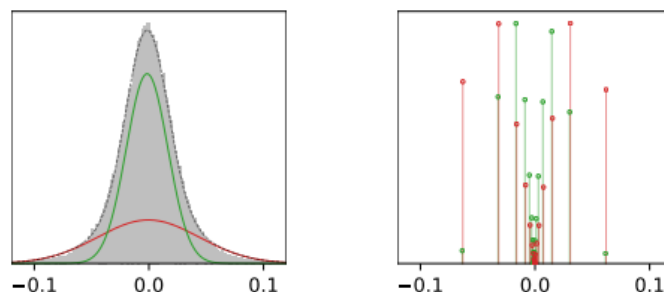


شکل ۳-۲- فرآیند گام به گام کوانتیزه‌سازی مجدد وزن‌های هرس نشده در block3f/conv1

¹³ expectation-maximization

۳-۱-۴- انتخاب کوانتیزاسیون مناسب

توزیع وزن لایه‌های نامتراکم ممکن است همواره دارای چندین ناحیه با احتمال بالا نباشد. به عنوان مثال، برازش یک مدل آمیخته از دو مؤلفه گاوسی بر روی لایه در شکل ۳-۳، اجزای با همپوشانی زیادی را به دست می‌دهد. بنابراین این که از کدام مؤلفه برای کوانتیزه کردن یک مقدار وزن خاص استفاده می‌شود، تأثیر چندانی ندارد. تحت این سناریو، می‌توان به سادگی از انتقال کوانتیزه‌سازی n بیتی $[0]$ $Q_{n,b}^{shift}$ به جای یک n بیت $Q[0]$ استفاده کرد که در داخل خود از یک $(n - 1)$ بیت علامت‌گذاری شده انتقال استفاده می‌کند.



• شکل ۳-۳: توزیع وزن لایه block22/conv1 در ResNet-18 نامتراکم آموزش دیده در

ImageNet

برای تصمیم‌گیری در مورد استفاده از انتقال یا کوانتیزاسیون مجدد، لازم است یک متریک برای مقایسه شباهت بین جفت مؤلفه‌ها معرفی شود. در حالی که واگرایی KL معیاری را برای تشابه ارائه می‌دهد، اما غیر متقارن است لذا برای این منظور نامناسب است. بنابراین در این مقاله پیشنهاد داده شده است که ابتدا توزیع مخلوط را نرمال کنید، سپس از متریک و سرتین^{۱۴} ۲ بین دو مؤلفه گاوسی پس از نرمال‌سازی به عنوان یک معیار تصمیم‌گیری استفاده نمایید، این جداسازی سرتین نامیده می‌شود.

$$\mathcal{W}(c_1, c_2) = \frac{1}{\sigma^2} \left((\mu_{c_1} - \mu_{c_2})^2 + (\sigma_{c_1} - \sigma_{c_2})^2 \right) \quad (4-3)$$

که در آن μ_c و σ_c به ترتیب میانگین و انحراف معیار مؤلفه $c \in \{c_1, c_2\}$ هستند. σ^2 نیز نشان‌دهنده واریانس کل توزیع وزن است. FQ می‌تواند به‌طور تطبیقی برای استفاده از کوانتیزه‌سازی مجدد برای

¹⁴ Wasserstein

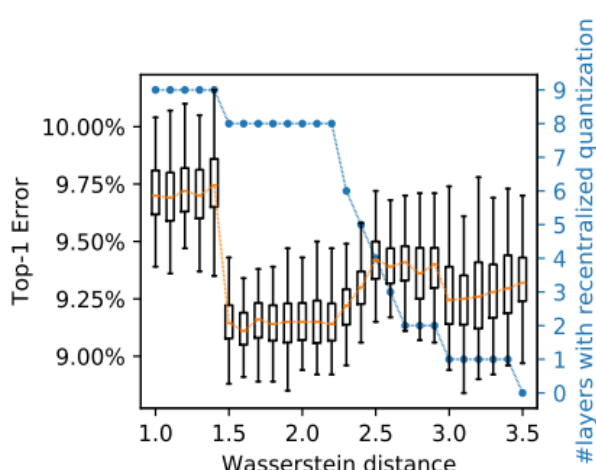
همه لایه‌های نامتراکم انتخاب کند، به جز زمانی که $\mathcal{W}(c_1, c_2) < w_{sep}$ ، و به جای آن از انتقال کوانتیزاسیون استفاده می‌شود.

۵-۲-۳: بهینه‌سازی مدل

برای بهینه‌سازی مدل نامتراکم کوانتیزه‌شده، فرآیند کوانتیزاسیون‌سازی که در بالا توضیح داده شد در آموزش پارامترهای مدل مبتنی بر گرادینان ادغام می‌شود. در ابتدا، ابرپارامترهای μ_c ، σ_c و λ_c برای هر لایه محاسبه می‌شوند، و ماسک انتخاب مولفه m_θ برای هر وزن θ با روش بخش ۳-۲-۳ ایجاد می‌شود. سپس مدل به دست آمده در جایی که حرکت روبه جلو از وزن‌های کوانتیزه $\hat{\theta} = Q[\theta]$ استفاده می‌کند، تنظیم می‌شود، و حرکت رو به عقب پارامترهای وزن ممیز شناور θ را با در نظر گرفتن کوانتیزاسیون به عنوان یک تابع هویت به روزرسانی می‌کند. در طول فرآیند تنظیم دقیق، ابرپارامترهای استفاده شده توسط $Q[\theta]$ با استفاده از توزیع وزن فعلی در هر k دوره بروز می‌شوند.

۳-۲-نتایج

تأثیر مقادیر مختلف آستانه بر فاصله و سترین در نمودار زیر قابل مشاهده است. هرچه آستانه بزرگتر باشد، تعداد لایه‌هایی که از کوانتیزه‌سازی مجدد به جای انتقال کوانتیزه‌سازی استفاده می‌کنند، کمتر می‌شود.



شکل ۳-۴: تأثیر مقادیر مختلف آستانه بر فاصله و سترین

دقت، پراکندگی و CR های فشرده‌سازی متمرکز در مدل‌های ImageNet. مدل‌های پایه قبل از فشرده‌سازی مدل‌های متراکم هستند و از وزن‌های ممیز شناور ۳۲ بیتی استفاده می‌کنند و ۵ بیت و ۷ بیت نشان‌دهنده تعداد بیت‌های استفاده شده توسط وزن‌های جداگانه مدل‌های کوانتیزه شده قبل از رمزگذاری هافمن است.

جدول ۱-۳: دقت، پراکندگی و CR های فشرده‌سازی متمرکز در مدل‌های ImageNet

Model	Top-1	Δ	Top-5	Δ	Sparsity	Size (MB)	CR (\times)
ResNet-18	68.94	—	88.67	—	0.00	46.76	—
Pruned	69.24	0.30	89.05	0.38	74.86	8.31	5.69
5 bits	68.36	-0.58	88.45	-0.22	74.86	2.86	16.33
7 bits	68.57	-0.37	88.53	-0.14	74.86	2.94	15.92
ResNet-50	75.58	—	92.83	—	0.00	93.82	—
Pruned	75.10	-0.48	92.58	-0.25	82.70	11.76	7.98
5 bits	74.86	-0.72	92.59	-0.24	82.70	5.19	18.08
7 bits	74.99	-0.59	92.59	-0.24	82.70	5.22	17.98
MobileNet-V1	70.77	—	89.48	—	0.00	16.84	—
Pruned	70.03	-0.74	89.13	-0.35	33.80	6.89	2.44
7 bits	69.13	-1.64	88.61	-0.87	33.80	2.13	7.90
MobileNet-V2	71.65	—	90.44	—	0.00	13.88	—
Pruned	71.24	-0.41	90.31	-0.13	31.74	5.64	2.46
7 bits	70.05	-1.60	89.55	-0.89	31.74	1.71	8.14

فصل چهارم

شبکه‌های عصبی کانولوشنی فوق نامتراکم^{۱۵}

¹⁵ Super Sparse Convolutional Neural Networks

مقدمه

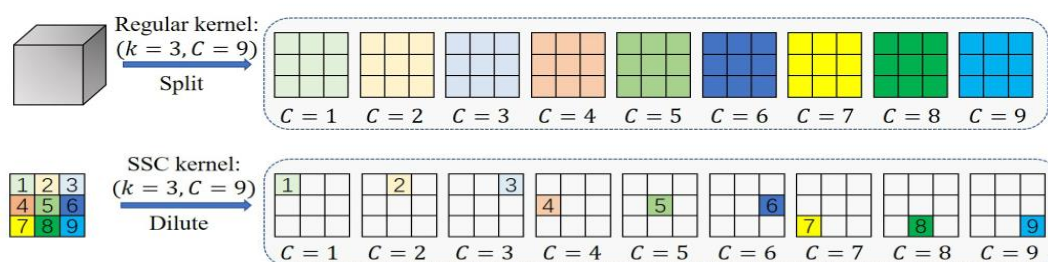
برای ساخت شبکه‌های موبایل کوچک، بدون از دست دادن عملکرد و رسیدگی به مسائل بیش از حد برآزش ناشی از مجموعه داده‌های آموزشی، در این بخش یک هسته کافنا پیشنهاد می‌شود و شبکه متناظر آن کافنا-نت نامیده می‌شود.

۴-۱- هسته کافنا

در شکل ۴-۱، یک هسته کافنا را می‌توان به عنوان تضعیف کردن یک هسته فضایی دو بعدی (به نام هسته اصلی) به یک هسته سه بعدی مشاهده کرد. به این ترتیب، مکان‌های غیرصفر فضایی با هسته اصلی یکسان نگه داشته می‌شوند که می‌تواند ویژگی‌های هندسی عمومی را حفظ کند. و فرآیند تضعیف کردن می‌تواند تفاوت‌های ویژگی‌ها را در طول کانال حفظ کند. فرض کنید هسته کانولوشن معمولی $\mathcal{F} \in \mathbb{R}^{k \times k \times C \times D}$ بعدی \mathcal{F} باشد به طوری که $k \times k$ اندازه هسته فضایی را نشان می‌دهد. C و D به ترتیب به تعداد کانال‌های ورودی و خروجی اشاره دارد. هسته کافنا با $\mathcal{S} \in \mathbb{R}^{k \times k \times C \times D}$ نشان داده می‌شود. لازم ذکر است که \mathcal{S} را می‌توان از تعمیم هسته اصلی $\mathcal{W} \in \mathbb{R}^{k \times k \times C \times D}$ به دست آورد. قابل توجه است که در یک هسته کافنا، $C = k \times k$. بنابراین، تعریف فرمول کافنا می‌تواند بدین صورت باشد:

$$\mathcal{S}_{i,j}^{x,y} = \begin{cases} \mathcal{W}_j^{x,y} & i = x \times k + y \\ 0 & \text{otherwise} \end{cases} \quad (1-4)$$

که x و y نشان‌دهنده موقعیت مکانی هستند و $x, y \in \{0, 1, 2, \dots, k-1\}$. علاوه بر این $i \in \{0, 1, 2, \dots, C-1\}$ و $j \in \{0, 1, 2, \dots, D-1\}$



شکل ۴-۱: مقایسه هسته کانولوشن معمولی و کافنا.

۴-۲- پیاده‌سازی عملگر کافنا

با توجه به یک تانسور^{۱۶} ورودی $\mathcal{I} \in \mathbb{R}^{w \times h \times C}$ که توسط هسته کافنا \mathcal{T} انجام می‌شود، تانسور خروجی مربوطه $\mathcal{O} \in \mathbb{R}^{w \times h \times C}$ را می‌توان به صورت زیر به دست آورد:

$$\mathcal{O}(x, y, q) = \sum_{p=0}^{C-1} \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} \mathcal{T}(i, j, p, q) \mathcal{I}(x + i - \delta_1, y + j - \delta_2, p) \quad (2-4)$$

به طوری که $\delta_1 = \delta_2 = \lfloor k/2 \rfloor$ و $q \in \{0, 1, 2, \dots, D-1\}$. با توجه به معادله (۴-۱)، فرآیند محاسباتی عملیات کافنا را می‌توان به صورت زیر ساده کرد:

$$\mathcal{O}(x, y, q) = \sum_{p=0}^{C-1} \mathcal{T}(i_p, j_p, p, q) \mathcal{I}(x + i_p - \delta_1, y + j_p - \delta_2, p) \quad (3-4)$$

به طوری که (i_p, j_p) تنها مختصات مکانی یک پارامتر غیر صفر در کانال p ام است و $i_p \times k + j_p = p$ از معادله (۳-۴)، واضح است که یک مقدار خروجی در موقعیت مکانی خاص را می‌توان از طریق جمع C بار عملیات ضرب محاسبه کرد. علاوه بر این، از آنجایی که هر صفحه در یک هسته کافنا دارای یک پارامتر است، می‌توان آن را به یک هسته "نقطه‌ای" تبدیل کرد. همچنین، به منظور حفظ روابط موقعیت مکانی در فرآیند محاسبات، از یک هسته انتقال $\mathcal{P} \in \mathbb{R}^{k \times k \times C \times D}$ استفاده می‌کنیم که بر روی نقشه‌های ویژگی برای گرفتن ویژگی‌ها عمل می‌کند. با توجه به موقعیت‌های مکانی غیر صفر هسته‌های کافنا عملیات انتقال در شکل ۴-۲ نشان داده شده است. تعریف این هسته انتقال به صورت زیر نشان داده شده است:

¹⁶ tensor

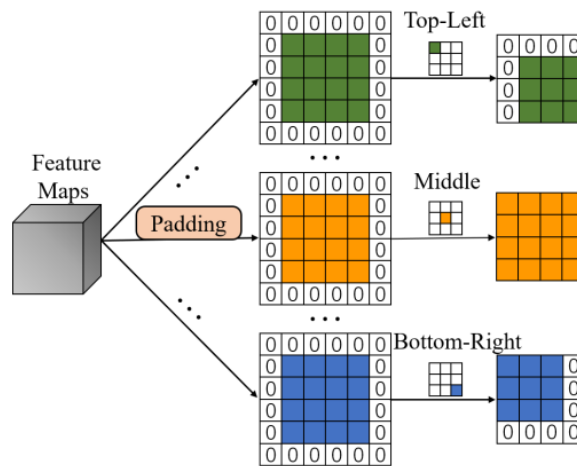
$$\mathcal{P}_{i,j}^{x,y} = \begin{cases} 1. & S_{i,j}^x \neq 0 \\ 0. & \text{otherwise} \end{cases} \quad (4-4)$$

از معادله (۳-۴) و (۴-۴)، عملیات کافنا را می‌توان به صورت زیر فرموله کرد:

$$\mathcal{O}(x, y, q) = \sum_{p=0}^{C-1} \mathcal{T}(p, q) \mathcal{P}(i_p, j_p, p, q) \quad (5-4)$$

$$\mathcal{I}(x + i_p - \delta_1, y + j_p - \delta_2, p)$$

از معادله (۵-۴)، عملیات کافنا را می‌توان با عملیات کانولوشنی انتقال و نقطه‌ای پیاده‌سازی کرد، که می‌تواند پارامترها و فلاپ‌های محاسباتی بسیار بیشتری را ذخیره کند.



شکل ۴-۲: فرآیند عملیاتی انتقال.

۴-۳- مازول پایه کافنا و کافنا-نت^{۱۷}

در ابتدا یک لایه کانولوشن «نقطه‌ای» جهت نگاشت نقشه‌های ویژگی ورودی به یک فضای ابعادی مورد نیاز استفاده می‌شود، که همچنین می‌تواند ویژگی‌های خروجی مازول قبلی را از گروه‌های مختلف ترکیب کند. کانال‌های نقشه ویژگی پیش‌بینی شده به دست آمده M هستند. سپس، دو لایه کافنا با یک عملیات

¹⁷ SSC-Nets

"کانال مختلط"^{۱۸} در میان آن‌ها مجهز می‌شوند. قبل از هر عملیات کافنا، یک لایه کانولوشنال «گروهی نقطه‌ای» دیگر استفاده می‌شود تا هر ورودی کافنا را وادار کند تا در حد امکان ویژگی‌های هندسی مشابهی داشته باشد. در نهایت، نگاشت هویت نیز در ماژول کافنا اتخاذ شده است. ساختار جزئی ماژول پایه کافنا-نت در جدول ۴-۱ نشان داده شده است.

جدول ۴-۱: ساختار شبکه

Stage	Operation	Groups	Channels/Group
Conv1	1×1	1	M
Conv2	1×1	G	C
1^{st} SSC	<i>shift</i>	G	C
	1×1	G	C
Shuffle	<i>shuffle</i>	C	G
Conv3	1×1	C	G
2^{nd} SSC	<i>shift</i>	G	C
	1×1	G	C
	1×1	C	G

۴-۴- مقایسه دو شبکه عصبی تلفن همراه پیشرفته

در نهایت، به منظور بررسی عملکرد کافنا-نت، کافنا-نت‌های مختلف با مدل مدرن تلفن همراه معرفی شده در مقاله ۱ مقایسه می‌شوند. در CIFAR، همان‌طور که در جدول ۴-۲ نشان داده شده است که کافنا-نت می‌تواند عملکرد بسیار بهتری در CIFAR-10 ایجاد کند. در مورد CIFAR 100، از آنجایی که کافنا به طور کامل توسط کانولوشن‌های نقطه‌ای در عمل پیاده‌سازی می‌شود، کافنا-نت همواره بهترین دقت را به دست نمی‌آورد، اما می‌تواند الزامات دقت عمومی را برآورده کند.

جدول ۴-۲: مقایسه دو شبکه عصبی تلفن همراه پیشرفته.

روش	پارامترها (M)	CIFAR-10	CIFAR 100
-----	---------------	----------	-----------

¹⁸ shuffle-channel

IGCV2*-C416 (Xie et al. 2018)	0.7	94.51	77.05
SSC-Net-6-9	1.2	95.14	75.99
SSC-Net-3-18	2.1	94.55	77.13
SSC-Net-4-18	2.8	94.95	77.67

فصل پنجم

شبکه کانولوشنی نامتراکم سریع^{۱۹}

^{۱۹} Fast Sparse ConvNets

مقدمه

تنک‌شدگی وزن معمولاً منجر به مدل‌های کوچک‌تر و کارآمدتر از نظر محاسباتی (از نظر تعداد عملیات‌های ممیز شناور) می‌شود، اما اغلب به عنوان ابزاری عملی برای تسریع مدل‌ها نادیده گرفته می‌شود، زیرا این تصور نادرست وجود دارد که عملیات نامتراکم نمی‌تواند به اندازه کافی سریع باشد تا به سرعت‌های واقعی در طول استنتاج دست یابد. برای پرداختن به این تصور اشتباه، در این مقاله، هسته‌های سریعی برای ضرب ماتریس نامتراکم - متراکم²⁰ (SpMM) معرفی می‌شود که به طور خاص شتاب شبکه‌های عصبی نامتراکم را هدف قرار می‌دهد.

۵-۱- پراکندگی شبکه‌ها

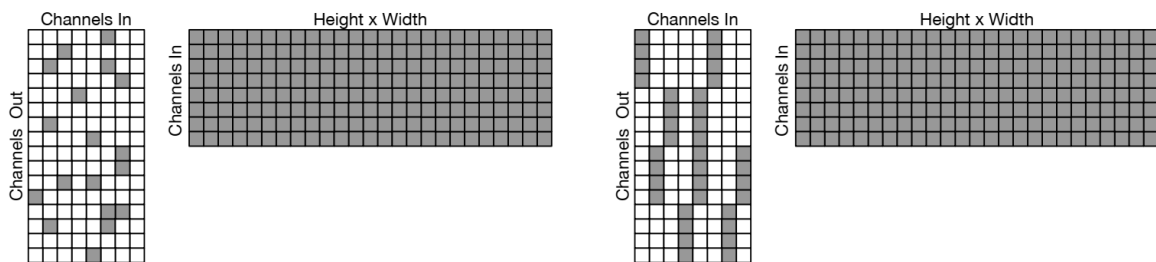
در این مقاله روی مجموعه داده ImageNet با تقویت استاندارد آموزش می‌دهیم و دقت‌های بالای ۱ را در مجموعه اعتبارسنجی نمونه ارائه شده ۵۰ کیلویی گزارش می‌کنیم. برای پراکندگی شبکه‌ها از تکنیک هرس تدریجی [۴۸] استفاده می‌کنیم.

ما اولین پیچیدگی کامل را در ابتدای هر سه شبکه هرس نمی‌کنیم. سهم کلی آن در شمارش پارامترها، تعداد FLOP و زمان اجرا اندک است و نیازی به معرفی یک عملیات پراکنده جدید ندارد. در عوض، ما یک هسته کانولوشنال متراکم را پیاده‌سازی می‌کنیم که تصویر را در طرح استاندارد HWC به عنوان ورودی می‌گیرد و طرح CHW مصرف شده توسط عملیات پراکنده در بقیه را خروجی می‌دهد

²⁰ Sparse Matrix-Dense Matrix

۵-۲- پیاده‌سازی هسته

نمودار کانولوشن 1×1 به عنوان SpMM در شکل ۵-۱ مشاهده می‌شود. طرح ارائه شده در این مقاله مستلزم این است که تنسورهای فعال‌سازی در قالب CHW ذخیره شوند، برخلاف کتابخانه‌های استنتاج سیار متراکم که خروجی‌های آن‌ها به صورت HWC هستند



شکل ۵-۱: کانولوشن نامتراکم 1×1 به عنوان SpMM. سمت چپ: پراکندگی بدون ساختار (با اندازه بلوک 1). سمت راست: اندازه بلوک کانال خروجی 4.

سه بینش کلیدی وجود دارد که موجب عملکرد بالای هسته‌ها می‌شود:

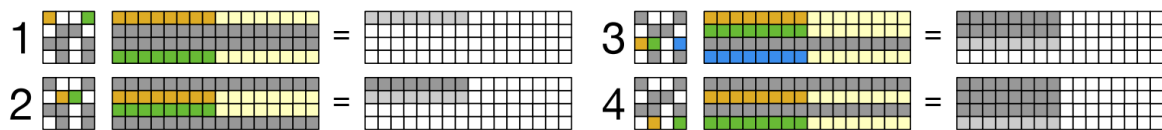
۱. در حالی که ماتریس وزن نامتراکم است، ماتریس فعال‌سازی متراکم باشد. این بدان معنی است که می‌توان بارها بردار را از ماتریس فعال‌سازی اجرا کرد و چندین موقعیت مکانی را به طور همزمان پردازش نمود.

۲. با پردازش ماتریس به ترتیب درست، می‌توان مقادیری را که به صورت تصادفی در حافظه پنهان L1 به آن‌ها دسترسی پیدا می‌کنند، نگه داشت، که دسترسی تصادفی از آن‌ها سریع و در زمان ثابت است.

۳. هنگامی که تعداد کانال‌های ورودی به اندازه کافی کم باشد، واکنشی اولیه از فعال‌سازی‌ها می‌تواند از دست دادن حافظه پنهان را کاهش دهد.

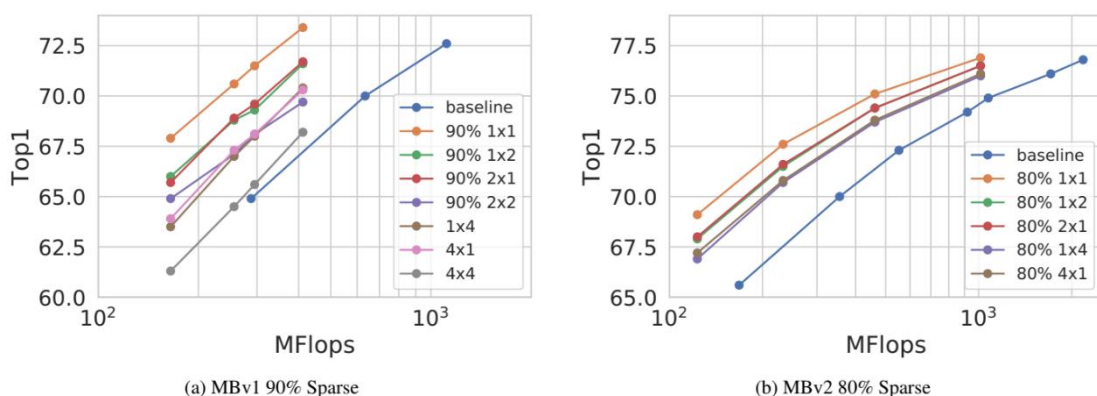
۵-۳- بررسی مفهومی

شکل ۵-۲ الگوهای خواندن و نوشتن حافظه در چند مرحله از هسته را نشان می‌دهد. شکل، 8 عنصر را نشان می‌دهد که همزمان برای تجسم پردازش می‌شوند، اما وجود ۱۶ عنصر برای یک اجرای واقعی



شکل ۵-۲: خواندن و نوشتن الگوریتم توسط تجسم حافظه.

طبیعی تر است زیرا مربوط به یک خط از حافظه است. حلقه بیرونی روی ستون‌ها و حلقه داخلی روی ردیف‌ها قرار دارد. این مسئله اجازه می‌دهد تا هر نوار از ۱۶ موقعیت مکانی در فعال سازی‌ها در حافظه پنهان L1 باقی بماند تا زمانی که دیگر مورد نیاز نباشد. در شکل ۵-۲ مراحل ۱ و ۲ حافظه پنهان را آماده می‌کند، در حالی که مراحل بعدی ۳ و ۴ تمام مقادیر سمت راست را از حافظه پنهان L1 بارگیری می‌کند. علاوه بر بردار سازی در بعد HW، بهره‌گیری از مقادیر اندک ساختار، در ماتریس وزن می‌تواند با افزایش استفاده مجدد از داده‌ها پس از بارگذاری مقادیر، افزایش عملکرد قابل توجهی را ارائه دهد. محدود کردن الگوی پراکندگی به طوری که چندین کلنال خروجی یا ورودی همه یک الگوی صفر / غیر صفر را به اشتراک می‌گذارند "بلوک‌ها" را در ماتریس وزن ایجاد می‌کند (سمت راست شکل ۵-۲). بلوک‌ها در بعد کانال خروجی امکان استفاده مجدد از داده‌ها را بیشتر از بلوک‌های بعد کانال ورودی می‌دهند. آزمایش‌ها (شکل ۵-۳) نشان می‌دهند که هر کدام از گزینه‌ها تأثیر یکسانی بر دقت دارد، بنابراین در این مقاله، مسدود کردن کانال خروجی با اندازه‌های ۲ و ۴ اجرا می‌شود. نامگذاری برای هسته‌ها این گونه است که عرض برداری مکانی آن‌ها به دنبال اندازه بلوک کانال خروجی ارائه داده می‌شود - 2×16 به معنای 16 پیکسل و پردازش 2 کانال خروجی در حلقه داخلی می‌باشد.



شکل ۵-۳: تأثیر اندازه بلوک بر صحت $top-1$.

فصل ششم

شبکه‌های کانولوشنی نامتراکم کانونی برای تشخیص اشیاء سه بعدی^{۲۱}

²¹ Focal Sparse Convolutional Networks for 3D Object Detection

مقدمه

داده‌های نامتراکم سه‌بعدی غیریکنواخت، مانند ابرهای نقطه‌ای^{۲۲} یا وکسل‌ها در موقعیت‌های فضایی مختلف، به روش‌های متفاوتی به وظیفه تشخیص اشیاء سه‌بعدی کمک می‌کنند. در این مقاله، دو ماژول جدید برای افزایش قابلیت شبکه‌های نامتراکم معرفی می‌شود، که هر دو مبتنی بر یادگیری ویژگی نامتراکم با پیش‌بینی اهمیت موقعیت هستند. در این مقاله برای نخستین بار، نشان داده شده است که پراکندگی فضایی قابل یادگیری در کانولوشن نامتراکم برای تشخیص اشیاء سه‌بعدی پیچیده ضروری است.

۶-۱- معرفی

یک چالش کلیدی در تشخیص اشیاء سه‌بعدی، یادگیری بازنمایی‌های مؤثر از داده‌های هندسی سه‌بعدی بدون ساختار و نامتراکم مانند ابرهای نقطه ای است.

صرف نظر از کانولوشن نامتراکم عادی یا زیرخمینه^{۲۳}، موقعیت‌های خروجی P_{out} در تمام $p \in P_{in}$ ثابت می‌باشد که نامطلوب است. در مقابل، در این مقاله تعیین تطبیقی اندازه‌های نامتراکم یا میدان پذیرنده به روشی دقیق انجام می‌شود. موقعیت‌های خروجی آزادسازی^{۲۴} می‌شوند تا P_{out} به صورت پویا توسط ویژگی‌های نامتراکم تعیین شود. این فرآیند پیشنهادی در شکل ۶-۱ نشان داده شده است. در فرمول، موقعیت‌های خروجی P_{out} به اجتماعی از همه موقعیت‌های مهم با ناحیه منبسط و سایر موقعیت‌های بی‌اهمیت تعمیم می‌یابند. نواحی منبسط تغییر شکل پذیر و پویا نسبت به موقعیت‌های ورودی هستند. معادله زیر حاصل می‌شود:

$$P_{out} = \left(\bigcup_{p \in P_{in}} P(p, K_{in}^d(p)) \right) \cup P_{in/im} \quad (۶-۱)$$

²²point clouds

²³ submanifold

²⁴ relax

این فرآیند به سه مرحله تقسیم می‌شود: (الف) پیش‌بینی اهمیت مکعبی، (ب) انتخاب ورودی مهم، و (ج) تولید شکل خروجی پویا

- پیش‌بینی اهمیت مکعبی:

یک نقشه اهمیت مکعبی I_p شامل اهمیت ویژگی‌های نامزد خروجی در اطراف ویژگی ورودی در موقعیت p است. هر نقشه اهمیت مکعبی همان شکل K^d را با وزن هسته کانولوشن پردازشی اصلی دارد، به عنوان مثال، $K^3 = 3 \times 3 \times 3$ با اندازه هسته ۳. این نقشه توسط یک کانولوشن نامتراکم زیر خمینه اضافی با تابع سیگموئید پیش‌بینی می‌شود. مراحل آخر به نقشه‌های اهمیت مکعبی پیش‌بینی شده بستگی دارد.

- انتخاب ورودی مهم:

P_{im} زیر مجموعه ای از P_{in} است. این شامل موقعیت‌های ویژگی‌های ورودی نسبتاً مهم است. می‌توان P_{im} را به صورت زیر انتخاب کرد:

$$P_{im} = \{p \mid I_0^p \geq \tau, p \in P_{in}\} \quad (۲-۶)$$

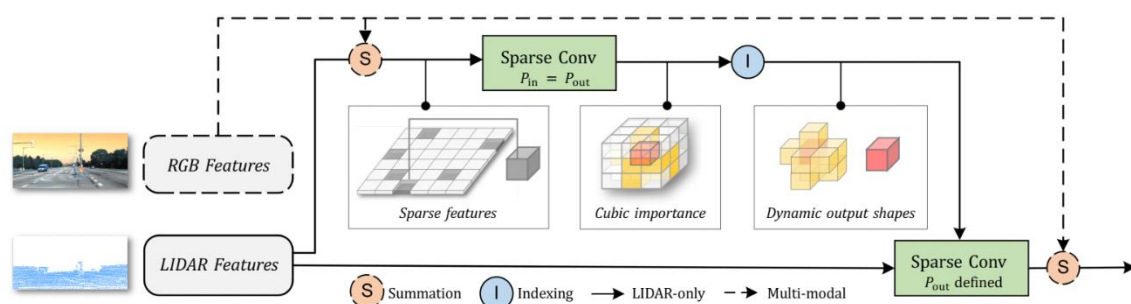
که در آن I_0^p مرکز نقشه اهمیت مکعبی در موقعیت p است. و τ یک آستانه از پیش تعریف شده است. هنگامی که τ به ترتیب ۰ یا ۱ باشد، فرمول ما به کانولوشن نامتراکم عادی یا زیرخمینه تبدیل می‌شود.

- تولید شکل خروجی پویا:

ویژگی‌های P_{im} به یک شکل پویا منبسط می‌شوند. خروجی اطراف p توسط شکل خروجی پویا $K_{im}^d(p)$ تعیین می‌شود. توجه داشته باشید که اشکال خروجی قابل تغییر شکل در داخل اتساع اصلی بدون انحرافات هرس می‌شوند. مشابه با فرمول (۲-۶) به صورت زیر محاسبه می‌شود:

$$K_{im}^d(p) = \{k \mid p + k \in P_{im}, I_k^p \geq \tau, k \in K^d\}$$

برای ویژگی‌های بی‌اهمیت باقی‌مانده، موقعیت‌های خروجی آن‌ها به عنوان ورودی، یعنی زیرمینه ثابت می‌شوند. حذف مستقیم آن‌ها یا استفاده از روشی کاملاً پویا بدون حفظ آن‌ها، روند آموزش را ناپایدار می‌کند.



شکل ۶-۱: چارچوب کانولوشن نامتراکم کانونی و گسترش چند وجهی آن.

۶-۲- مراحل اجرا

در تشخیص اشیاء سه بعدی، اشیاء پیش‌زمینه اطلاعات ارزشمندتری می‌باشند. بر اساس این موضوع، تابع زیان^{۲۵} به عنوان یک تابع زیان هدف برای نظارت بر پیش‌بینی اهمیت اعمال در نظر گرفته می‌شود. اهداف عینی برای مراکز و کسل‌های ویژگی در جعبه‌های داده مرجع سه بعدی ساخته می‌شود. وزن زیان آن به عنوان 1 برای کلیه ماژول‌ها نگه داشته می‌شود.

نظارت اضافی از ضرب نقشه‌های اهمیت مکعبی پیش‌بینی شده در ویژگی‌های خروجی به عنوان توجه حاصل می‌شود. این امر باعث می‌شود که شاخه پیش‌بینی اهمیت به طور طبیعی متمایز شود. به طور تجربی نشان داده شده است که این شیوه توجه برای کلاس‌های کوچک، به عنوان مثال، عابر پیاده و دوچرخه‌سوار در دیتاست KITTI مفید است.

²⁵ loss

۶-۳-ادغام^{۲۶} کاناکا

یک نسخه چند وجهی از کانولوشن نامتراکم کانونی، همان طور که در شکل (۶-۱) نشان داده شده است، در نظر گرفته می شود. در این بخش ویژگی های RGB از تصاویر استخراج می شوند و ویژگی های لایدار^{۲۷} با آن ها تراز می شوند. ویژگی های استخراج شده با ویژگی های نامتراکم ورودی و خروجی مهم در کانولوشن نامتراکم کانونی ترکیب می شوند.

۶-۳-۱-استخراج ویژگی

ماژول ادغام سبک وزن است و شامل یک لایه conv-bn-relu و یک لایه max-pooling است. این عملیات وضوح تصویر ورودی را $\frac{1}{4}$ می کند. به دنبال آن ۳ لایه conv-bn-relu با اتصال باقی ماندگی^{۲۸} وجود دارد. سپس شماره کانال کاهش می یابد تا با ویژگی های نامتراکم، با یک لایه از پرسپترون چند لایه سازگار باشد.

۶-۳-۲-تراز ویژگی

یکی از مشکلات رایج در حین ادغام، ناهمبستگی در نگاشت از فضای سه بعدی به فضای دو بعدی است. داده های ابر نقطه ای معمولاً با تبدیل و تقویت پردازش می شوند. تبدیلات شامل معکوس کردن، مقیاس مجدد، چرخش و انتقال است. تقویت معمولی نمونه برداری از داده مرجع و کپی کردن اشیاء از صحنه های دیگر است. برای تبدیل های معکوس، مختصات ویژگی های نامتراکم را با پارامترهای تبدیل ثبت شده معکوس می کنیم. برای نمونه برداری از داده مرجع، اشیاء دو بعدی مربوطه را روی تصاویر کپی می کنیم. به جای استفاده از یک مدل تقسیم بندی اضافی، برای ساده سازی، اشیاء مستقیماً در جعبه های مرزی برش داده می شود.

مراحل ادغام.

²⁶Fusion

²⁷ LIDAR

²⁸ residual

ویژگی‌های RGB تراز شده مستقیماً با ویژگی‌های نامتراکم ترکیب می‌شوند، زیرا شماره‌های کانال یکسانی دارند. ویژگی‌های RGB تراز شده دو بار در این ماژول با ویژگی‌های نامتراکم ترکیب می‌شوند. ابتدا با ویژگی‌های ورودی برای پیش‌بینی اهمیت مکعب ترکیب می‌شود. سپس ویژگی‌های RGB فقط با ویژگی‌های نامتراکم خروجی مهم ترکیب می‌شود. به طور کلی، لایه‌های چند حالتی از نظر پارامترها و استراتژی‌های ادغام سبک وزن هستند. آن‌ها به طور مشترک با آشکارسازها آموزش می‌بینند. این روش یک راه‌حل کارآمد و اقتصادی برای ماژول ادغام در تشخیص اشیاء سه بعدی ارائه می‌دهد.

۶-۴- شبکه‌های کاناکا

کاناکا و گسترش چندوجهی آن می‌توانند به راحتی جایگزین همتایان خود در شبکه‌های ماز^{۲۹} آشکارسازهای سه بعدی شوند. در طول آموزش، از تنظیمات اولیه یا نرخ یادگیری خاصی برای ماژول های معرفی شده استفاده نمی‌شود. شاخه پیش‌بینی اهمیت به وسیله پس انتشار از طریق تابع ضرب توجه و از تابع زیان هدف آموزش داده می‌شود. شبکه‌های ماز در آشکارسازهای شی سه بعدی معمولاً از یک لایه ساقه و چهار مرحله تشکیل شده‌اند. هر مرحله، به جز مرحله اول، شامل یک کانولوشن منظم نامتراکم با زیرنمونه‌برداری و دو بلوک زیرخمینه است. در مرحله اول، یک یا دو لایه کانولوشن نامتراکم وجود دارد. به طور پیش فرض، هر کانولوشن نامتراکم با نرمال‌سازی دسته‌ای و فعال‌سازی ReLU دنبال می‌شود.

²⁹ backbone network

فصل هفتم شبکه‌های عصبی کانولوشنی فضایی-زمانی ۴ بعدی^{۳۰}: شبکه‌های عصبی کانولوشنی مینکوفسکی^{۳۱}

³⁰ 4D Spatio-Temporal ConvNets

³¹ Minkowski Convolutional Neural Networks

مقدمه

در بسیاری از برنامه‌های رباتیک و VR/AR، ویدیوهای سه بعدی منابع ورودی به راحتی در دسترس هستند (توالی از تصاویر عمقی یا اسکن‌های لایدار). با این حال، در بسیاری از موارد، فیلم‌های سه بعدی فریم به فریم یا از طریق شبکه‌های دو بعدی یا الگوریتم‌های ادراک سه بعدی پردازش می‌شوند. در این بخش، شبکه‌های عصبی کانولوشنی ۴ بعدی برای ادراک فضایی-زمانی پیشنهاد می‌شود که می‌تواند مستقیماً چنین ویدئوهای سه بعدی را با استفاده از پیچش‌های با ابعاد بالا پردازش کند.

۷-۱- معرفی

ادراک فضایی-زمانی ۴ بعدی اساساً به ادراک سه بعدی به عنوان برشی از ۴ بعد در امتداد بعد زمانی یک اسکن سه بعدی نیاز دارد. با این حال، در درجه اول بایستی ادراک سه بعدی به ویژه تقسیم‌بندی سه بعدی با استفاده از شبکه‌های عصبی پوشش داده شود.

۷-۲- کانولوشن و تانسور نامتراکم

در داده‌های گفتار، متن یا تصویر سنتی، ویژگی‌ها به طور فشرده استخراج می‌شوند. با این حال، برای اسکن‌های سه بعدی، چنین نمایش متراکمی ناکارآمد است، زیرا بیشتر فضا خالی است. در عوض، می‌توان فضای غیر خالی را تحت عنوان مختصات آن و ویژگی مرتبط ذخیره کرد. این نمایش یک بسط N بعدی از یک ماتریس نامتراکم است. به طور خلاصه، می‌توان مجموعه‌ای از مختصات ۴ بعدی را به صورت $C = \{(x_i, y_i, z_i, t_i)\}_i$ یا به عنوان یک ماتریس C و ویژگی‌های مرتبط $F = \{f_i\}_i$ یا به عنوان ماتریس F نشان داد. لذا یک تانسور نامتراکم را می‌توان به صورت زیر نوشت:

$$C = \begin{bmatrix} x_1 & y_1 & z_1 & t_1 & b_1 \\ & & \vdots & & \\ x_N & y_N & z_N & t_N & b_N \end{bmatrix}, F = \begin{bmatrix} f_1^T \\ \vdots \\ f_N^T \end{bmatrix} \quad (1-7)$$

که در آن b_i شاخص‌های دسته‌ای مختصات و f_i یک بردار است.

۷-۳- کانولوشن نامتراکم تعمیم یافته:

کانولوشن نامتراکم تعمیم یافته نه تنها همه کانولوشن های نامتراکم بلکه کانولوشن های متراکم معمولی را نیز در بر می گیرد. فرض کنید $\mathbf{x}_u^{\text{in}} \in \mathbb{R}^{N^{\text{in}}}$ یک بردار ویژگی ورودی در یک فضای D بعدی در \mathbb{R}^D (یک مختصات D بعدی باشد)، و وزن هسته کانولوشنی $\mathbf{W} \in \mathbb{R}^{K^D \times N^{\text{out}} \times N^{\text{in}}}$ باشد. سپس وزن ها با ماتریس های K^D به اندازه $N^{\text{out}} \times N^{\text{in}}$ تحت عنوان W_i برای $|i| = K^D$ به وزن های فضایی تقسیم می شوند. در نهایت، کانولوشن متراکم معمولی در بعد D به صورت زیر است:

$$\mathbf{x}_u^{\text{out}} = \sum_{i \in \mathcal{V}^D(K)} W_i \mathbf{x}_{u+i}^{\text{in}} \text{ for } \mathbf{u} \in \mathbb{Z}^D \quad (2-7)$$

به طوری که $\mathcal{V}^D(K)$ لیستی از انحرافات در ابرمکعب D بعدی است که در مرکز مبدا قرار دارد. به عنوان مثال، $\mathcal{V}^1(3) = \{-1, 0, 1\}$ کانولوشن نامتراکم تعمیم یافته در معادله (۷-۳) تعداد پارامترهای معادله (۲-۷) را کم می کند.

$$\mathbf{x}_u^{\text{out}} = \sum_{i \in \mathcal{N}^D(\mathbf{u}, \mathcal{C}^{\text{in}})} W_i \mathbf{x}_{u+i}^{\text{in}} \text{ for } \mathbf{u} \in \mathcal{C}^{\text{out}} \quad (3-7)$$

که در آن \mathcal{N}^D مجموعه ای از انحرافات است که شکل یک هسته را تعریف می کند و $\mathcal{N}^D(\mathbf{u}, \mathcal{C}^{\text{in}}) = \{i | \mathbf{u} + i \in \mathcal{C}^{\text{in}}, i \in \mathcal{N}^D\}$ به عنوان مجموعه ای از انحرافات از مرکز فعلی، \mathbf{u} ، که در \mathcal{C}^{in} وجود دارد، شناخته می شود. \mathcal{C}^{in} و \mathcal{C}^{out} مختصات ورودی و خروجی از پیش تعریف شده تنسورهای نامتراکم هستند. توجه داشته باشید که مختصات ورودی و خروجی لزوماً معادل نیستند. شکل هسته کانولوشن به طور دلخواه با \mathcal{N}^D تعریف می شود. این تعمیم بسیاری از موارد خاص مانند هسته های کانولوشنی منبسط و هسته های ابرمکعبی را در بر می گیرد. یکی دیگر از موارد خاص این حالت است که $\mathcal{C}^{\text{in}} = \mathcal{C}^{\text{out}}$ و $\mathcal{N}^D = \mathcal{V}^D(K)$ در این حالت "کانولوشن زیرخمینه نامتراکم" است. اگر $\mathcal{C}^{\text{in}} = \mathcal{C}^{\text{out}} = \mathbb{Z}^D$ و $\mathcal{N}^D = \mathcal{V}^D(K)$ ، کانولوشن نامتراکم تعمیم یافته معادل کانولوشن متراکم است. برای کانولوشن های گام به گام، $\mathcal{C}^{\text{in}} \neq \mathcal{C}^{\text{out}}$.

۷-۴-مهندسی مینکوفسکی

۷-۴-۱-کوانتیزاسیون تنسور نامتراکم

اولین مرحله در شبکه عصبی کانولوشنی نامتراکم، پردازش داده‌ها برای تولید یک تنسور نامتراکم است که ورودی را به مختصات منحصر به فرد و ویژگی‌های مرتبط تبدیل می‌کند. در الگوریتم ۱، تابع GPU برای این فرآیند فهرست می‌شود. به طور خاص، برای بخش‌بندی معنایی، یک برچسب برای هر جفت ویژگی مختصات ورودی ایجاد می‌کند. اگر بیش از یک برچسب معنایی مختلف در یک وکسل وجود داشته باشد، در طول آموزش با علامت‌گذاری آن با IGNORE_LABEL، این وکسل نادیده گرفته می‌شود. ابتدا، همه مختصات به کلیدهای هش تبدیل می‌شود و جفت برچسب هش-کلید منحصر به فرد برای حذف برخورد پیدا می‌شود. توجه داشته باشید که SortByKey، UniqueByKey، و ReduceByKey همگی توابع استاندارد کتابخانه Thrust هستند. تابع کاهش $\leq f((l_x, i_x), (l_y, i_y))$

Algorithm 1 GPU Sparse Tensor Quantization

Inputs: coordinates $C_p \in \mathbb{R}^{N \times D}$, features $F_p \in \mathbb{R}^{N \times N_f}$,
target labels $\mathbf{l} \in \mathbb{Z}_+^N$, quantization step size v_l
 $C'_p \leftarrow \text{floor}(C_p / v_l)$
 $\mathbf{k} \leftarrow \text{hash}(C'_p)$, $\mathbf{i} \leftarrow \text{sequence}(N)$
 $((\mathbf{i}', \mathbf{l}'), \mathbf{k}') \leftarrow \text{SortByKey}((\mathbf{i}, \mathbf{l}), \text{key}=\mathbf{k})$
 $(\mathbf{i}'', (\mathbf{k}'', \mathbf{l}'')) \leftarrow \text{UniqueByKey}(\mathbf{i}', \text{key}=(\mathbf{k}', \mathbf{l}'))$
 $(\mathbf{l}''', \mathbf{i}''') \leftarrow \text{ReduceByKey}((\mathbf{l}'', \mathbf{i}'''), \text{key}=\mathbf{k}'', \text{fn}=f)$
return $C'_p[\mathbf{i}''', :], F_p[\mathbf{i}''', :], \mathbf{l}'''$

(IGNORE_LABEL, i_x) جفت کلیدهای برچسب را می‌گیرد و برچسب نادیده گرفتن را برمی‌گرداند زیرا حداقل دو جفت کلید برچسب در یک کلید به معنای برخورد برچسب است. یک نسخه CPU نیز به طور مشابه کار می‌کند با این تفاوت که تمام کاهش‌ها و مرتب‌سازی‌ها به صورت سریال پردازش می‌شوند.

مرحله بعدی در خط لوله، تولید مختصات خروجی \mathcal{C}^{out} با توجه به مختصات ورودی \mathcal{C}^{in} (معادل ۷-۳) است. هنگامی که در شبکه‌های عصبی معمولی استفاده می‌شود، این فرآیند فقط به اندازه گام لایه کانولوشنی (یا ادغام)، مختصات ورودی، و اندازه گام تنسور نامتراکم ورودی (حداقل فاصله بین مختصات) نیاز دارد. علاوه بر این، همچنین از تنظیم پویا یک مختصات خروجی دلخواه \mathcal{C}^{out} برای کانولوشن نامتراکم تعمیم‌یافته پشتیبانی می‌شود.

در مرحله بعد، برای ادغام ورودی‌ها با یک هسته، به یک نگاشت نیاز است تا مشخص شود کدام ورودی‌ها بر روی خروجی‌ها تأثیر می‌گذارند. این نگاشت نقشه‌های هسته نامیده می‌شود و آن‌ها به عنوان جفت لیست‌هایی از شاخص‌های ورودی و شاخص‌های خروجی، $M = \{(I_i, O_i)\}_i$ برای هر $i \in \mathcal{N}^D$ تعریف می‌شوند. در نهایت، با توجه به مختصات ورودی و خروجی، نقشه هسته، و وزن هسته W_i ، می‌توان کانولوشن نامتراکم تعمیم‌یافته را با تکرار در هر یک از انحرافات $i \in \mathcal{N}^D$ (الگوریتم ۲) محاسبه کرد.

Algorithm 2 Generalized Sparse Convolution

Require: Kernel weights \mathbf{W} , input features F^i , output feature placeholder F^o , convolution mapping \mathbf{M} ,

- 1: $F^o \leftarrow \mathbf{0}$ // set to 0
 - 2: **for all** $W_i, (I_i, O_i) \in (\mathbf{W}, \mathbf{M})$ **do**
 - 3: $F_{\text{tmp}} \leftarrow W_i[F_{I_i[1]}^i, F_{I_i[2]}^i, \dots, F_{I_i[n]}^i]$ // (cu)BLAS
 - 4: $F_{\text{tmp}} \leftarrow F_{\text{tmp}} + [F_{O_i[1]}^o, F_{O_i[2]}^o, \dots, F_{O_i[n]}^o]$
 - 5: $[F_{O_i[1]}^o, F_{O_i[2]}^o, \dots, F_{O_i[n]}^o] \leftarrow F_{\text{tmp}}$
 - 6: **end for**
-

به طوریکه $I[n]$ و $O[n]$ به ترتیب عنصر n ام فهرست شاخص‌های I و O و F_n^i و F_n^o نیز به ترتیب بردارهای ورودی و خروجی n ام هستند. کانولوشن نامتراکم تعمیم‌یافته جابجا شده (دکانولوشن^{۳۲}) به طور مشابه کار می‌کند با این تفاوت که نقش مختصات ورودی و خروجی معکوس است.

³² deconvolution

۷-۴-۲-ادغام ماکسیمم ۳۳

بر خلاف تنسورهای متراکم، در تنسورهای نامتراکم، تعداد ویژگی‌های ورودی در هر خروجی متفاوت است. بنابراین، این یک پیاده‌سازی غیر ضروری برای ادغام ایجاد می‌کند. فرض کنید I و O برداری باشند

Algorithm 3 GPU Sparse Tensor MaxPooling

Input: input feature F , output mapping O
 $(I', O') \leftarrow \text{SortByKey}(I, \text{key}=O)$
 $S \leftarrow \text{Sequence}(\text{length}(O'))$
 $S', O'' \leftarrow \text{ReduceByKey}(S, \text{key}=O', \text{fn}=f)$
return $\text{MaxPoolKernel}(S', I', O'', F)$

که همه $\{I_i\}_i$ و $\{O_i\}_i$ را به ترتیب برای $i \in \mathcal{N}^D$ به هم الحاق کرده است. ابتدا تعداد ورودی‌ها در هر مختصات خروجی و شاخص‌های آن ورودی‌ها را پیدا کنید. الگوریتم ۳ ویژگی‌های ورودی را که به همان مختصات خروجی نگاشت می‌شوند، کاهش می‌دهد. $\text{Sequence}(n)$ دنباله‌ای از اعداد صحیح از 0 تا $n - 1$ و تابع کاهش $f((k_1, v_1), (k_2, v_2)) = \min(v_1, v_2)$ را تولید می‌کند که حداقل مقدار جفت کلید-ارزش را برمی‌گرداند. MaxPoolKernel یک هسته CUDA سفارشی است که ویژگی‌ها را با استفاده از S' کاهش می‌دهد که شامل شاخص ابتدایی I و شاخص‌های خروجی مربوطه O است.

³³ Max Pooling

۷-۴-۳- ادغام مجموع، ادغام میانگین و ادغام سراسری

ادغام میانگین و ادغام سراسری، میانگین ویژگی‌های ورودی را برای هر مختصات خروجی محاسبه می‌کند. الگوریتم این کار بدین صورت است:

Algorithm 4 GPU Sparse Tensor AvgPooling

Input: mapping $\mathbf{M} = (\mathbf{I}, \mathbf{O})$, features F , one vector $\mathbf{1}$
 $S_M = \text{coo2csr}(\text{row}=\mathbf{O}, \text{col}=\mathbf{I}, \text{val}=\mathbf{1})$
 $F' = \text{cuspars_csrmm}(S_M, F)$
 $N = \text{cuspars_csrmv}(S_M, \mathbf{1})$
return F'/N

۷-۵- توابع غیر فضایی

برای توابعی که نیازی به اطلاعات مکانی (مختصات) ندارند، می‌توان توابع را مستقیماً روی ویژگی‌های F اعمال کرد. برای مثال، غیر خطی‌ها به اطلاعات مکانی مانند ReLU نیاز ندارند.

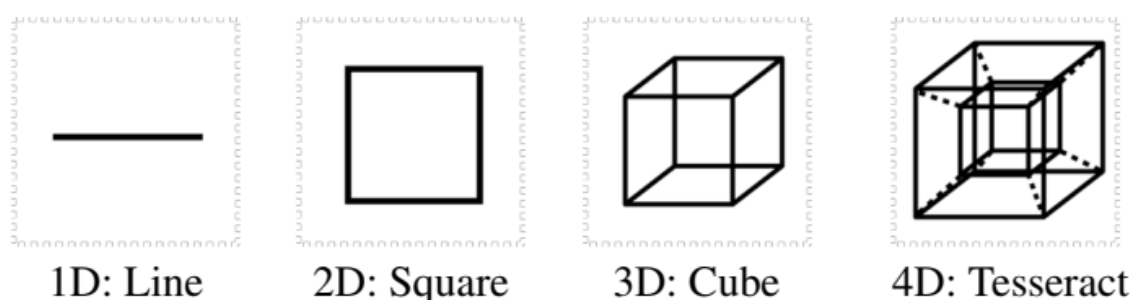
۷-۶- شبکه‌های عصبی کانولوشنی مینکوفسکی

از ویژگی خاصی از کانولوشن نامتراکم تعمیم‌یافته استفاده می‌شود و اشکال هسته غیرمتعارف پیشنهاد می‌گردد که با تعمیم بهتر، حافظه و محاسبات را ذخیره می‌کند. از طرفی جهت سازگاری مکانی-زمانی، یک میدان تصادفی شرطی با ابعاد بالا (در فضای ۷ بعدی فضا-زمان-رنگ) پیشنهاد می‌شود که می‌تواند ثبات را اعمال کند و شبکه پایه و میدان تصادفی شرطی را از مبدا به مقصد^{۳۴} آموزش دهد.

³⁴³⁴ end-to-end

۷-۷- هسته تسراکت^{۳۵} و هسته هیبریدی^{۳۶}

مساحت سطح داده‌های ۳ بعدی به صورت خطی نسبت به زمان و به صورت درجه دوم نسبت به وضوح فضایی افزایش می‌یابد. با این حال، اگر از یک ابرمکعب ۴ بعدی یا یک تسراکت (شکل ۷-۱) برای هسته‌های کانولوشن استفاده شود، افزایش تصاعدی در تعداد پارامترها به احتمال زیاد منجر به پارامترسازی بیش از حد، برازش بیش از حد، و همچنین هزینه محاسباتی و مصرف حافظه بالا می‌شود.



شکل ۷-۱: پیش‌بینی دوبعدی ابرمکعب در ابعاد مختلف.

لذا، یک هسته ترکیبی پیشنهاد می‌شود که از شکل هسته دلخواه کانولوشن نامتراکم تعمیم‌یافته، \mathcal{N}^D استفاده می‌کند.

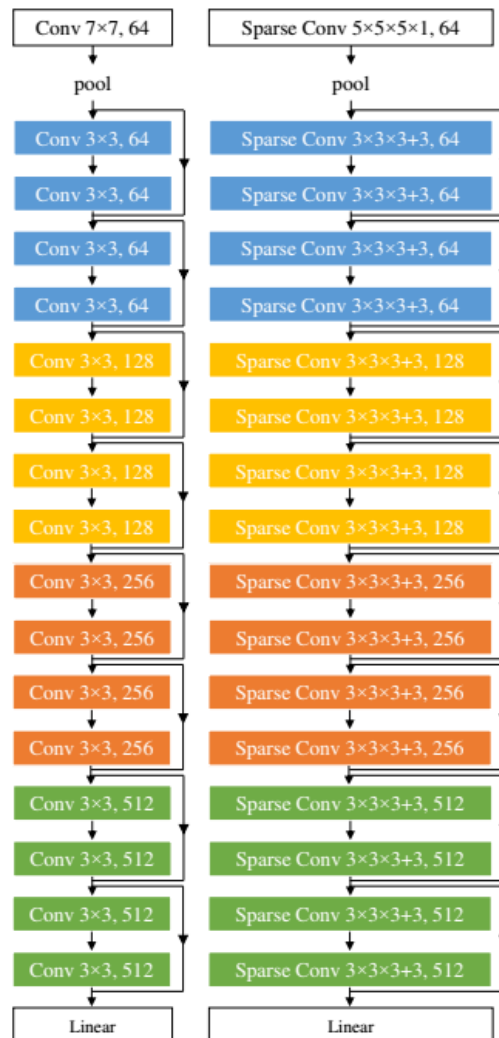
۷-۸- شبکه‌های باقی‌ماندگی مینکوفسکی^{۳۷}

کانولوشن نامتراکم تعمیم‌یافته این اجازه را می‌دهد تا گام‌ها و اشکال هسته به طور دلخواه تعریف شوند. لذا برای لایه اول، به جای کانولوشن 7×7 دوبعدی از کانولوشن نامتراکم تعمیم‌یافته $5 \times 5 \times 5 \times 1$ استفاده می‌شود. و برای مابقی، از طرح اصلی پیروی می‌شود و نسخه چهار بعدی ResNet18 در شکل ۷-۲ تجسم می‌شود.

³⁵ Tesseract

³⁶ Hybrid

³⁷ Residual Minkowski



شکل ۷-۲: معماری ResNet18 (سمت چپ) و MinkowskiNet18 (راست).

۷-۹- CRF-سه گانه ثابت

پیش‌بینی‌های شبکه مینکوفسکی برای مراحل زمانی مختلف لزوماً در سراسر محور زمانی سازگار نیستند. برای شفاف‌تر کردن چنین سازگاری و بهبود پیش‌بینی‌ها، یک میدان تصادفی شرطی با یک هسته ثابت تعریف شده در یک فضای سه‌جانبه پیشنهاد می‌شود. فضای سه‌ضلعی از فضای ۳ بعدی، زمان ۱ بعدی و

فضای رنگی ۳ بعدی تشکیل شده است. این گسترش فضای دوطرفه در پردازش تصویر است. فضای رنگی اجازه می‌دهد تا نقاطی با رنگ‌های مختلف که از نظر فضایی مجاور هستند (مثلاً روی یک مرز) در فضای رنگی از هم دور باشند.

۷-۱۰-تعریف

فرض کنید یک گره CRF^3 در فضای ۷-بعدی (فضا-زمان-کروما) x_i باشد. از اجزای بیرونی دوربین برای تبدیل مختصات فضایی یک گره x_i استفاده می‌شود تا در سیستم مختصات جهانی تعریف شود تا نقاط استاتیک حتی زمانی که ناظر حرکت می‌کند، در همان مختصات باقی بمانند.

برای هر گره x_i پتانسیل یگانه را به صورت $\phi_u(x_i)$ و پتانسیل جفت به صورت $\phi_p(x_i, x_j)$ تعریف می‌شود به طوری که x_j همسایه x_i یعنی $\mathcal{N}^7(x_i)$ است. فیلد تصادفی شرطی نهایی به صورت تعریف می‌شود:

$$P(\mathbf{X}) = \frac{1}{Z} \exp \sum_i \left(\phi_u(x_i) + \sum_{j \in \mathcal{N}^7(x_i)} \phi_p(x_i, x_j) \right) \quad (۴-۷)$$

که در آن Z تابع افراز است. X مجموعه‌ای از تمام گره‌ها است. و ϕ_p بایستی شرط ثلث بودن $\phi_p(\mathbf{u}, \mathbf{v}) = \phi_p(\mathbf{u} + \tau_u, \mathbf{v} + \tau_v)$ را برای هر $\tau_u, \tau_v \in \mathbb{R}^D$ برآورده کند.

۷-۱۱-استنتاج متغیر

مسئله بهینه‌سازی $\arg \max_X P(X)$ غیرقابل حل است. بنابراین، از استنتاج تغییرات برای به حداقل رساندن واگرایی بین $P(X)$ بهینه و توزیع تقریبی $Q(X)$ استفاده می‌شود. به طور ویژه، از تقریب میدان

میانگین $Q = \prod_i Q_i(x_i)$ به عنوان راه حل شکل بسته استفاده می‌شود. Q یک حداکثر محلی است اگر و فقط اگر:

$$Q_i(x_i) = \frac{1}{Z_i} \exp \mathbf{E}_{\mathbf{x}_{-i} \sim Q_{-i}} \left[\phi_u(x_i) + \sum_{j \in \mathcal{N}^7(x_i)} \phi_p(x_i, x_j) \right] \quad (5-7)$$

Q_{-i} و \mathbf{x}_{-i} همه گره‌ها یا متغیرها را به جز متغیر i ام نشان می‌دهند. معادله نقطه ثابت نهایی معادله زیر است:

$$Q_i^+(x_i) = \frac{1}{Z_i} \exp \left\{ \phi_u(x_i) + \sum_{j \in \mathcal{N}^7(x_i)} \sum_{x_j} \phi_p(x_i, x_j) Q_j(x_j) \right\} \quad (6-7)$$

۷-۱۲- یادگیری با کانولوشن نامتراکم ۷ بعدی

جالب توجه است که جمع وزنی $\phi_p(x_i, x_j) Q_j(x_j)$ در معادله (۶-۷) معادل یک کانولوشن نامتراکم تعمیم‌یافته در فضای ۷ بعدی است زیرا ϕ_p ثابت است و لبه‌ها را می‌توان با استفاده از \mathcal{N}^7 تعریف کرد. الگوریتم نهایی در الگوریتم ۵ آمده است.

Algorithm 5 Variational Inference of TS-CRF

Require: Input: Logit scores ϕ_u for all x_i ; associated coordinate C_i , color F_i , time T_i
 $Q^0(X) = \exp \phi_u(X)$, $C_{\text{crf}} = [C, F, T]$
for n from 1 to N **do**
 $\tilde{Q}^n = \text{SparseConvolution}((C_{\text{crf}}, Q^{n-1}), \text{kernel}=\phi_p)$
 $Q^n = \text{Softmax}(\phi_u + \tilde{Q}^n)$
end for
return Q^N

در نهایت، از ϕ_u به عنوان پیش‌بینی‌های یک شبکه ۴ بعدی مینکوفسکی استفاده می‌شود و هر دو ϕ_u و ϕ_p با استفاده از یک شبکه ۴ بعدی و یک شبکه ۷ بعدی مینکوفسکی با استفاده از معادله زیر آموزش داده می‌شود:

$$\frac{\partial L}{\partial \phi_p} = \sum_n^N \frac{\partial L}{\partial Q^{n+}} \frac{\partial Q^{n+}}{\partial \phi_p} \cdot \frac{\partial L}{\partial \phi_u} = \sum_n^N \frac{\partial L}{\partial Q^{n+}} \frac{\partial Q^{n+}}{\partial \phi_u} \quad (7-7)$$

۷-۱۳- نتایج

نتایج تقسیم‌بندی در مجموعه داده ۴ بعدی synthia در جدول ۷-۱ آورده شده است:

جدول ۷-۱: نتایج تقسیم‌بندی در مجموعه داده ۴ بعدی synthia

Method	mIOU	mAcc
3D MinkNet20	76.24	89.31
3D MinkNet20 + TA	77.03	89.20
4D Tesseract MinkNet20	75.34	89.27
4D MinkNet20	77.46	88.013
4D MinkNet20 + TS-CRF	78.30	90.23
4D MinkNet32 + TS-CRF	78.67	90.51

مقایسه الگوریتم‌های متفاوت نیز در جدول ۷-۲ آورده شده است:

جدول ۷-۲: مقایسه نتایج حاصل از الگوریتم‌های مختلف

Method	mIOU	mAcc
PointNet [23]	41.09	48.98
SparseUNet [9]	41.72	64.62
SegCloud [31]	48.92	57.35
TangentConv [30]	52.8	60.7
3D RNN [33]	53.4	71.3
PointCNN [16]	57.26	63.86
SuperpointGraph [15]	58.04	66.5
MinkowskiNet20	62.60	69.62
MinkowskiNet32	65.35	71.71

منابع و مراجع

- [1] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770-778. 2016.
- [2] Zhang, Ting, Guo-Jun Qi, Bin Xiao, and Jingdong Wang. "Interleaved group convolutions." In Proceedings of the IEEE international conference on computer vision, pp. 4373-4382. 2017.
- [3] Chollet, François. "Xception: Deep learning with depthwise separable convolutions." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1251-1258. 2017.
- [4] Xie, Guotian, Jingdong Wang, Ting Zhang, Jianhuang Lai, Richang Hong, and Guo-Jun Qi. "Interleaved structured sparse convolutional neural networks." In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8847-8856. 2018.
- [5] Zhao, Yiren, Xitong Gao, Daniel Bates, Robert Mullins, and Cheng-Zhong Xu. "Focused quantization for sparse cnns." *Advances in Neural Information Processing Systems* 32 (2019).
- [6] Lu, Yao, Guangming Lu, Bob Zhang, Yuanrong Xu, and Jinxing Li. "Super sparse convolutional neural networks." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 4440-4447. 2019.
- [7] Elsen, Erich, Marat Dukhan, Trevor Gale, and Karen Simonyan. "Fast sparse convnets." In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 14629-14638. 2020.
- [8] Chen, Yukang, Yanwei Li, Xiangyu Zhang, Jian Sun, and Jiaya Jia. "Focal Sparse Convolutional Networks for 3D Object Detection." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5428-5437. 2022.
- [9] Choy, Christopher, JunYoung Gwak, and Silvio Savarese. "4d spatio-temporal convnets: Minkowski convolutional neural networks." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3075-3084. 2019.



