

Multi-Level Ternary Quantization for Improving Sparsity and Computation in Embedded Deep Neural Networks

Hosna Manavi Mofrad¹, Seyed Ali Ansarmohammadi², Mostafa Ersali Salehi Nasab^{3*}

¹ Electrical and Computer Engineering, University of Tehran, Tehran, Iran

² Electrical and Computer Engineering, University of Tehran, Tehran, Iran

³ Electrical and Computer Engineering, University of Tehran, Tehran, Iran

Received: 03 January 2023, Revised: 15 February 2023, Accepted: 25 May 2023

Paper type: Research

Abstract

Deep neural networks (DNNs) have achieved great interest due to their success in various applications. However, the computation complexity and memory size are considered to be the main obstacles for implementing such models on embedded devices with limited memory and computational resources. Network compression techniques can overcome these challenges. Quantization and pruning methods are the most important compression techniques among them. One of the famous quantization methods in DNNs is the multi-level binary quantization, which not only exploits simple bit-wise logical operations, but also reduces the accuracy gap between binary neural networks and full precision DNNs. Since, multi-level binary can't represent the zero value, this quantization does not take advantage of sparsity. On the other hand, it has been shown that DNNs are sparse, and by pruning the parameters of the DNNs, the amount of data storage in memory is reduced while computation speedup is also achieved. In this paper, we propose a pruning and quantization-aware training method for multi-level ternary quantization that takes advantage of both multi-level quantization and data sparsity. In addition to increasing the accuracy of the network compared to the binary multi-level networks, it gives the network the ability to be sparse. To save memory size and computation complexity, we increase the sparsity in the quantized network by pruning until the accuracy loss is negligible. The results show that the potential speedup of computation for our model at the bit and word-level sparsity can be increased by 15x and 45x compared to the basic multi-level binary networks.

Keywords: Deep Neural Networks, Multi-Level Ternary Quantization, Sparse Neural Network, Pruning, Embedded Devices.

* Corresponding Author's email: mersali@ut.ac.ir

چندی سازی غیریکنواخت سه حالتی جهت بهبود تنگی و محاسبات شبکه های عصبی عمیق در کاربردهای نهفته

حسنا معنوی مفرد^۱، سید علی انصار محمدی^۲، مصطفی ارسالی صالحی نسب^{۳*}

^۱ دانشجو مقطع ارشد، دانشکده مهندسی برق و کامپیوتر، دانشکدگان فنی، دانشگاه تهران، تهران، ایران

^۲ دانشجو مقطع دکتری، دانشکده مهندسی برق و کامپیوتر، دانشکدگان فنی، دانشگاه تهران، تهران، ایران

^۳ استادیار، دانشکده مهندسی برق و کامپیوتر، دانشکدگان فنی، دانشگاه تهران، تهران، ایران

تاریخ دریافت: ۱۴۰۱/۱۰/۱۳ تاریخ بازبینی: ۱۴۰۱/۱۱/۲۶ تاریخ پذیرش: ۱۴۰۲/۰۳/۰۴

نوع مقاله: پژوهشی

چکیده

شبکه های عصبی عمیق به دلیل موفقیت در کاربردهای مختلف، به جذابیت فوق العاده ای دست یافته اند. اما پیچیدگی محاسبات و حجم حافظه از موانع اصلی برای پیاده سازی آن ها در بسیاری از دستگاه های نهفته تلقی می شود. از مهم ترین روش های بهینه سازی که در سال های اخیر برای برطرف نمودن این موانع ارائه شده، می توان به چندی سازی و هرس کردن اشاره کرد. یکی از روش های معروف چندی سازی، استفاده از نمایش اعداد غیریکنواخت دو حالتی است که علاوه بر بهره بردن از محاسبات بیتی، افت صحت شبکه های دو حالتی را در مقایسه با شبکه های دقت کامل کاهش می دهد. اما به دلیل نداشتن قابلیت نمایش عدد صفر در آن ها، مزایای تنگی داده ها را از دست می دهند. از طرفی، شبکه های عصبی عمیق به صورت ذاتی تنک هستند و با تنک کردن پارامترهای شبکه عصبی عمیق، حجم داده ها در حافظه کاهش می یابد و همچنین به کمک روش هایی می توان انجام محاسبات را تسریع کرد. در این مقاله می خواهیم هم از مزایای چندی سازی غیریکنواخت و هم از تنگی داده ها بهره ببریم. برای این منظور چندی سازی غیریکنواخت سه حالتی برای نمایش اعداد ارائه می دهیم که علاوه بر افزایش صحت شبکه نسبت به شبکه غیریکنواخت دو حالتی، قابلیت هرس کردن را به شبکه می دهد. سپس میزان تنگی در شبکه چندی شده را با استفاده از هرس کردن افزایش می دهیم. نتایج نشان می دهد که تسریع بالقوه شبکه ما در سطح بیت و کلمه می تواند به ترتیب ۱۵ و ۴۵ برابر نسبت به شبکه غیریکنواخت دو حالتی پایه افزایش یابد.

کلیدواژگان: شبکه های عصبی عمیق، چندی سازی غیریکنواخت سه حالتی، شبکه عصبی تنک، هرس کردن، دستگاه های نهفته.

* رایانامه نویسنده مسؤول: mersali@ut.ac.ir

۱- مقدمه

شوند تا بتوانند روی این دستگاه‌ها گسترش یابند [۵]. سطوح بهینه‌سازی شبکه‌های عصبی عمیق جهت کاهش حافظه و حجم محاسبات شامل موارد زیر است: (۱) طراحی شبکه‌های بهینه مانند MobileNets [۶]، SqueezeNet [۷]، (۲) استفاده از روش‌های فشرده‌سازی شامل چندی‌سازی^۹ [۸-۱۰] و هرس کردن^{۱۰} [۱۱-۱۳]، (۳) بهینه‌سازی در سطح الگوریتم برای سرعت بخشیدن به محاسبات مانند ضرب عمومی ماتریس^{۱۱} یا روش winograd [۱۴] و (۴) استفاده از سخت‌افزار اختصاصی برای تسریع جریان داده‌ها و انجام محاسبات به صورت موازی [۱۵].

چندی‌سازی و هرس کردن دو روش مؤثر در فشرده‌سازی شبکه‌های عصبی عمیق است. در روش چندی‌سازی سعی بر این است که پارامترها را (وزن‌ها^{۱۲} و فعال‌سازها^{۱۳}) که در حالت عادی به صورت ممیز شناور ۳۲ بیتی هستند، با یک سیستم نمایش اعداد با تعداد بیت کمتر نشان دهند تا منجر به هزینه محاسباتی و مصرف حافظه کمتر شود. به عنوان اولین تلاش‌ها در این زمینه، سیستم اعداد ممیز ثابت چندی شده [۱۶-۱۸] مطرح شد که سعی می‌کرد شبکه را با تعداد بیت کمتر طراحی کند و برای کاهش صحت از دست رفته، چندی‌سازی همراه با یادگیری شبکه را ارائه داده اند. چالشی که در این سیستم اعداد حتی بعد از چندی‌سازی وجود دارد این است که هزینه محاسبات و میزان حافظه مصرفی هنوز قابل توجه است.

برای برطرف کردن این مسئله سعی شد که دقت اعداد با شدت بیشتری کاهش پیدا کند به این صورت که وزن‌ها و فعال‌سازها با حتی یک بیت نمایش داده شوند. به این مدل چندی‌سازی، دو حالتی^{۱۴} [۹، ۱۹] می‌گویند که شبکه‌های عصبی دو حالتی را ایجاد می‌کنند. در چندی‌سازی دو حالتی پارامترهای شبکه با یک مقدار ۱ یا ۰ تخمین زده می‌شوند که با یک بیت قابل نمایش است. در نتیجه حافظه لازم برای ذخیره‌سازی پارامترهای مدل و مقادیر میانی به شدت کاهش می‌یابد. از طرفی در شبکه‌های عصبی دو حالتی، یک عملیات کانولوشن را می‌توان با روش Xnor_PopCount که شامل عملیات بیتی Xnor و عملیات ساده شمارش بیت است، انجام داد. به دلیل استفاده از عملیات ساده بیتی، در شبکه‌های عصبی دو حالتی پیچیدگی محاسبات نیز کاهش خواهد یافت. اما

در سال‌های اخیر تحقیقات در زمینه هوش مصنوعی، به‌ویژه یادگیری ماشین، به‌طور چشم‌گیری افزایش پیدا کرده است. یادگیری ماشین^۱ از زیرمجموعه‌های هوش مصنوعی^۲ است که هدف آن هوشمند کردن رایانه‌هاست بدون اینکه مستقیماً به آن‌ها یاد بدهیم چطور رفتار کنند. یادگیری ماشین شامل الگوریتم‌های متعددی است که پرکاربردترین آن‌ها الگوریتم‌های برگرفته از مغز انسان هستند؛ یادگیری عمیق، در این دسته قرار دارند، در بیشتر روش‌های یادگیری عمیق از معماری شبکه عصبی مصنوعی^۳ استفاده می‌شود، به همین دلیل مدل‌های یادگیری عمیق اغلب تحت عنوان شبکه‌های عصبی عمیق شناخته می‌شوند. اصطلاح عمیق معمولاً به تعداد لایه‌های پنهان در شبکه عصبی اشاره دارد. شبکه‌های عصبی سنتی فقط شامل ۲-۳ لایه مخفی هستند، در حالی که شبکه‌های عمیق می‌توانند تا ۱۵۰ لایه داشته باشند. یکی از رایج‌ترین انواع شبکه‌های عصبی عمیق، تحت عنوان شبکه‌های عصبی کانولوشنی^۴ شناخته می‌شود. این شبکه‌ها توانایی‌های خارق‌العاده‌ای در کاربردهای پیچیده مانند طبقه‌بندی تصویر^۵، تشخیص اشیاء^۶ و قطعه‌بندی معنایی^۷ تصاویر نشان داده‌اند [۱].

اگرچه طراحی‌های اخیر شبکه‌های عصبی کانولوشنی با میلیاردها پارامتر قابلیت‌هایی در سطح فرا انسانی دارند، با این حال، اندازه بزرگ مدل و هزینه بالای محاسبات برای بسیاری از برنامه‌ها کاربردی به‌ویژه با هدف پیاده‌سازی بر روی دستگاه‌های نهفته که محدودیت حافظه، منابع محاسباتی و توان دارند، همچنان مانع بزرگی محسوب می‌شود [۲]. به‌طور مثال شبکه vgg16 با مجموعه داده ImageNet به ۵۵۲ مگابایت پارامتر و ۳۰٫۸ گیگافلاپس محاسبات برای هر تصویر نیاز دارد. این میزان از حافظه و محاسبات برای گسترش دادن شبکه‌های عصبی عمیق روی دستگاه‌های با محدودیت منابع دشوار است [۳]. در حال حاضر، توجه به سمت گسترش دادن شبکه‌های عصبی عمیق روی دستگاه‌های نهفته^۸ جلب شده است [۴].

با توجه به محدودیت‌های دستگاه‌های نهفته، لازم است که شبکه‌های عصبی از نظر حافظه‌ی مصرفی و حجم محاسبات فشرده

⁸ Embedded device

⁹ Quantization

¹⁰ Pruning

¹¹ General matrix multiplication (GEMM)

¹² Weights

¹³ Activations

¹⁴ Binary

¹ Machine learning

² Artificial Intelligence

³ Artificial neural network

⁴ Convolutional neural network

⁵ Image classification

⁶ Object detection

⁷ Semantic segmentation

می‌شود، بنابراین میزان حافظه مورد نیاز کاهش پیدا می‌کند. از طرفی به این دلیل که ضرب صفر در هر عدد، صفر می‌شود بنابراین می‌توان از این نوع محاسبات در شبکه کانولوشنی جلوگیری کرد که این روش باعث کم شدن پیچیدگی محاسبات نیز خواهد شد [۲۵-۲۸]. برای فشرده‌سازی حداکثری شبکه از آموزش شبکه همگام با چندی‌سازی و هرس کردن بهره برده اند.

یکی از مزایای تنکی در وزن‌ها، شانس استفاده از کدگذاری جهت فشرده‌سازی وزن‌ها در حافظه‌ی خارجی تراشه هست که هر چه تنکی بیشتر باشد تاثیر عملکرد فشرده‌سازی بیشتر می‌شود و از هزینه‌های دسترسی بالای انتقال داده از حافظه‌های خارجی به داخل تراشه می‌کاهد. الگوریتم‌های مختلفی برای فشرده‌سازی وزن‌های غیرصفر وجود دارد که چهار مورد از مهمترین آن‌ها به نام فشرده‌سازی سطری^۳ [۲۹]، فشرده‌سازی ستونی^۴ [۳۰، ۳۱]، فشرده‌سازی با توجه به ابعاد ورودی^۵ [۳۲-۳۴] و کدبندی طول اجرا^۶ [۳۵] هستند.

در این مقاله قصد داریم روش چندی‌سازی ارایه دهیم که از مزایای تنکی بهره برد و هم‌زمان بتوانیم از تنکی ذاتی شبکه‌های عصبی جهت فشرده‌سازی حافظه و جلوگیری از محاسبات بی‌تاثیر استفاده کنیم و تا حد امکان از کاهش صحت شبکه جلوگیری کنیم. لذا در مرحله‌ی اول، با استفاده از چندی‌سازی غیریکنواخت سه‌حالتی امکان بهبود صحت شبکه و نمایش مقادیر صفر را ارایه می‌کنیم و سپس با استفاده از اعمال تنکی در پارامترهای وزن‌های شبکه همگام با چندی‌سازی و عملیات آموزش شبکه، حجم پارامترها و حجم محاسبات در اجرای لایه‌های شبکه را کاهش می‌دهیم. در پایان با ارائه یک معماری پایه جهت مقایسه عملیات چندی‌سازی غیریکنواخت دو‌حالتی با سه‌حالتی و با امکان حذف عملیات غیرموثر ضرب، حداکثر تسریع بالقوه محاسبات با استفاده از سیستم نمایش اعداد غیریکنواخت سه‌حالتی پیشنهادی را در سطح کلمه و بیت نسبت به دو‌حالتی نشان خواهیم داد.

نوآوری‌های این مقاله به شرح زیر است:

- استفاده از چندی‌سازی غیریکنواخت سه‌حالتی برای رفع مشکل تنکی در شبکه‌های عصبی عمیق دو‌حالتی
- آموزش شبکه غیریکنواخت سه‌حالتی با آگاهی هم‌زمان از چندی‌سازی و هرس
- افزایش تنکی وزن‌ها در سطح کلمه به کمک هرس کردن با

چالشی که در این شبکه‌های دو‌حالتی وجود دارد افت صحت^۱ است [۲۰]. همین امر باعث ظهور روش دیگری از چندی‌سازی به نام سه‌حالتی^۲ [۱۰، ۲۱، ۲۲] شد که در این حالت هر پارامتر را با (۱ و ۰) نشان می‌دهیم. در مقایسه با نمایش دو‌حالتی که برای نمایش هر پارامتر یک بیت کافی بود، در این نمایش برای نمایش هر یک از مقادیر (۱ و ۰) به ۲ بیت نیاز داریم. دلیل اصلی پیدایش این نوع چندی‌سازی پر کردن فاصله صحت بین مدل دو‌حالتی و ۳۲ بیت ممیز شناور بود. هرچند صحت شبکه‌های عصبی عمیق سه‌حالتی از دو‌حالتی بهتر شد ولی با وجود دو برابر شدن حجم پارامترهای شبکه هنوز با صحت شبکه‌های ۳۲ بیتی ممیز شناور فاصله‌ی قابل توجه‌ای داشتند.

مقاله‌های زیادی جهت جبران صحت از دست رفته‌ی شبکه‌های عصبی عمیق دو‌حالتی چاپ شده‌است، از جمله: روش‌های مبتنی بر افزایش ابعاد عرضی شبکه (افزایش تعداد فیلترهای هر لایه) [۲۳]. هرچند این روش‌ها نتایج قابل ملاحظه‌ای در جهت بهبود صحت شبکه داشته‌اند ولی در این روش‌ها کم‌کم حجم بالای پارامترها آزاردهنده هست و همچنین نیاز است تا شبکه با فرآیندهای جدید طراحی و آموزش داده شود. برای برطرف کردن فاصله صحت بین شبکه دو‌حالتی و شبکه اصلی (۳۲ بیتی ممیز شناور) روش غیریکنواخت دو‌حالتی [۲، ۲۴] ارائه گردید که در این حالت به جای نمایش یک عدد ۳۲ بیتی ممیز شناور، از چند رقم دو‌حالتی (دو‌حالتی ۲، ۳ و ... رقمی) برای نمایش استفاده می‌شود. با توجه به غیریکنواخت بودن سطوح چندی‌سازی در این روش و همچنین استفاده از چند رقم دو‌حالتی، می‌توان از این روش به عنوان چندی‌سازی غیریکنواخت چندرقمی دو‌حالتی نام برد ولی به اختصار در این مقاله از چندی‌سازی غیریکنواخت دو‌حالتی استفاده شده است. مقاله‌های [۲، ۲۴] نشان داده‌اند که با دو‌حالتی کردن وزن‌ها و فعال‌سازها به ۳ یا ۴ رقم، صحت شبکه قابل قبول و نزدیک به صحت شبکه در حالت دقت کامل است. لذا با این نوع چندی‌سازی هم از خواص عملیات بیتی در محاسبات دو‌حالتی استفاده شده‌است و هم افت صحت شبکه قابل چشم‌پوشی هست.

روش موثر دیگر در فشرده‌سازی شبکه‌های عصبی کانولوشنی، روش هرس کردن پارامترهای صفر شبکه است. درواقع با هرس کردن پارامترهای نزدیک صفر که تأثیر خیلی کمی در صحت شبکه دارند، حذف می‌شوند. در این حالت تعداد پارامترهای شبکه موردنظر کمتر

⁴ Compressed Sparse Column (CSC)

⁵ Compressed Image Siz (CIS)

⁶ Run Length Coding (RLC)

¹ Accuracy

² Ternary

³ Compressed Sparse Row (CSR)

۳۷]. در واقع شبکه‌های عصبی کانولوشنی به صورت ذاتی شبکه‌هایی تنک هستند. در چندی‌سازی ممیز ثابت با کمتر از ۸ بیت که به عنوان یکی از مرسوم‌ترین و محبوب‌ترین چندی‌سازی‌ها مطرح شده، مقاله‌های [۳۶-۳۸] نشان داده‌اند که با از دست دادن مقداری قابل چشم‌پوشی از صحت، وزن‌های شبکه‌های عصبی عمیق می‌توانند، تا ۷۰ درصد در سطح کلمه و تا ۹۲ درصد در سطح بیت صفر داشته باشند. این خود انگیزه‌ای شده است تا از این تنکی شبکه‌ها جهت کاهش حجم پارامترها در حافظه‌های خارجی تراشه و همچنین کاهش محاسبات استفاده کنیم.

از طرف دیگر روش‌های چندی‌سازی دو حالتی مطرح شدند تا پارامترها را فقط با دو مقدار ۱- و یا ۱ با تنها یک بیت نمایش دهند و باعث کاهش هرچه بیشتر حجم داده‌ها در حافظه و ساده کردن عملیات کانولوشنی با استفاده از عملیات بیتی و شمارش شوند. ولی به دلیل اینکه فاصله بین صحت شبکه‌های دو حالتی نسبت به شبکه‌های ممیز ثابت زیاد است شبکه‌های غیریکنواخت دو حالتی مطرح شدند. این شبکه‌ها علاوه بر در نظر گرفتن سطوح چندی‌سازی غیریکنواخت، با استفاده از افزایش سطوح چندی‌سازی باعث افزایش صحت شبکه‌ی چندی شده‌اند. به عنوان مثال برای مجموعه داده‌ی ImageNet مقاله‌ی [۲] با استفاده از ۸ یا ۱۶ سطح چندی غیریکنواخت به صحت نزدیک به شبکه با دقت کامل رسیده است که برای نمایش این ۸ یا ۱۶ سطح چندی به ۳ یا ۴ بیت در سخت‌افزار نیاز هست.

همان‌طور که پیش‌تر اشاره شد شبکه‌های کانولوشنی به طور ذاتی قابلیت تنکی زیادی در پارامترهای وزن دارند. در مقابل شبکه‌های چندی دو حالتی غیریکنواخت در کنار بهره بردن از محاسبات سبک بیتی و همچنین جبران صحت از دست‌رفته و نزدیک شدن به صحت شبکه با دقت کامل، شبکه‌های محبوب و مناسبی برای دستگاه‌های نهفته می‌توانند باشند. اما در شبکه‌های چندی دو حالتی غیریکنواخت چون برای نمایش هر رقم تنها از (۱ و ۰-) استفاده می‌شود، عدد صفر قابل‌نمایش نیست. بنابراین در شبکه‌های غیریکنواخت دو حالتی نمی‌توان از مزایای تنکی در سطح کلمه برای کاهش محاسبات و حافظه بهره برد.

ایده این مقاله از همین نکته به ذهن ما رسید، که روش چندی‌سازی ارائه کنیم که علاوه بر بهره بردن از محاسبات دو حالتی جهت کاهش حافظه و انرژی، امکان نشان دادن اعداد صفر را داشته باشد تا بتوانیم از مزایای تنکی شبکه‌های عصبی استفاده کنیم و تا جایی که صحت شبکه آسیب نبیند بتوانیم شبکه را هرس کنیم. بنابراین روش استفاده از سیستم اعداد غیریکنواخت سه حالتی

حفظ صحت در شبکه غیریکنواخت سه حالتی

- جبران سربار حافظه ناشی از روش چندی‌سازی غیریکنواخت سه حالتی به کمک فشرده‌سازی وزن‌های تنک شده
- ارائه معماری مناسب برای حذف عملیات غیر موثر ضرب در مقدار صفر جهت کاهش حجم محاسبات شبکه‌های غیریکنواخت سه حالتی

ساختار این مقاله به شرح زیر سازماندهی شده است: در بخش دو به انگیزه‌ای که باعث شد این تحقیق انجام شود خواهیم پرداخت و در بخش سوم به بیان کارهای که در پژوهش‌های پیشین در این حوزه انجام شده می‌پردازیم. ایده مقاله در بخش چهارم که شامل نحوه‌ی چندی‌سازی شبکه و تنک کردن شبکه است ارائه می‌شود. همچنین نحوه‌ی فشرده‌سازی حافظه و کاهش محاسبات با استفاده از واحدهای محاسباتی پایه معرفی می‌شود. در نهایت نتایج را در بخش پنجم بررسی و تحلیل می‌کنیم و در انتها نیز نتیجه‌گیری و مراجع به ترتیب در بخش‌های ششم و هفتم قرار گرفته است.

۲- انگیزه پژوهش

یکی از روش‌های رایج برای کاهش محاسبات شبکه‌های عصبی عمیق در راستای قابل پیاده‌سازی شدن آن‌ها برای یک دستگاه نهفته، استفاده از تنکی پارامترهای شبکه است. به این معنی که مقادیر خیلی کوچک را به صفر تقریب می‌زنند و به این ترتیب سعی می‌کنند تعداد پارامترهای با مقدار صفر را افزایش دهند. انگیزه‌ی اصلی از ایجاد تنکی در شبکه‌ها، در واقع امکان حذف عملیات ضرب در صفر در محاسبات شبکه است. این مقادیر صفر می‌توانند برای وزن‌ها یا ورودی‌ها و یا خروجی‌های یک نرون باشند که معمولاً به آن تنکی در سطح کلمه می‌گویند. علاوه بر سطح کلمه، در سطح بیت، حذف عملیات بی اثر برای بیت‌های صفر هم می‌تواند سبب کاهش عملیات محاسباتی شود که به آن تنکی در سطح بیت گفته می‌شود. مقاله‌های [۳۶، ۳۷] برای تنکی در سطح بیت ۵ الی ۹ برابر بهبود در سرعت و کارایی انرژی را گزارش کرده‌اند. در واقع با توجه به این شواهد، امید داریم که با تنک کردن شبکه به محدودیت‌های ناشی از پیاده‌سازی شبکه‌های عصبی کانولوشنی روی دستگاه‌های نهفته غلبه کنیم.

در کارهای پیشین که در زمینه تنکی شبکه‌های عصبی انجام شده، نتایج به دست آمده حاکی از آن است که شبکه‌های عصبی چه در حالت ممیز شناور چه در حالتی که به صورت ممیز ثابت چندی شده اند، حاوی تعداد زیادی پارامتر صفر و یا نزدیک به صفر در سطح کلمه و بیت (هم در وزن‌ها و هم در فعال‌سازها) هستند [۲۹، ۳۶،

کاهش صحت ناشی از اعمال این روش‌ها غلبه کنند نیاز هست که شبکه را دوباره آموزش دهند یا از هم‌زمانی آموزش شبکه با توجه به چندی‌سازی یا تنکی شبکه استفاده شود.

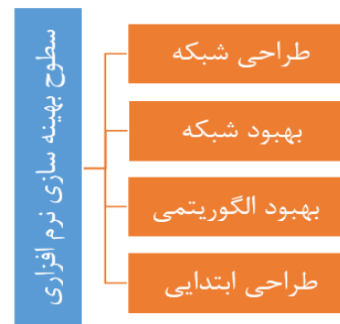
چندی‌سازی پارامترهای شبکه یک روش برای فشرده‌سازی حافظه شبکه‌های عصبی عمیق، کاهش هزینه سخت‌افزار و هزینه محاسبات با استفاده از عملیات منطقی با پهنای بیت کم است. وقتی که وزن‌ها و فعال‌سازها به‌درستی کوانتیزه شوند، می‌توان عملیات کانولوشن را به‌طور مؤثر توسط عملگرهای بیتی محاسبه کرد [۲، ۴۴]. این روش با رویکردهایی مانند شبکه‌های عصبی دو حالتی [۲]، شبکه‌های سه حالتی [۱۰] معرفی شد که بین صحت شبکه‌های عصبی عمیق چندی شده و نمونه‌های کاملاً دقیق آن‌ها فاصله وجود داشت، خصوصاً وقتی شبکه با پهنای بیت بسیار کم مانند یک بیت یا دو بیت چندی شود. برای رفع این شکاف Lqnet [۲] و RebNet [۲۴] چندی‌سازی غیریکنواخت دو حالتی برای وزن‌ها و فعال‌سازها پیشنهاد دادند که صحت شبکه آن‌ها با صحت شبکه ممیز شناور تقریباً برابر است.

علاوه بر روش‌های چندی‌سازی، بهره بردن از تنکی بیت‌های وزن و فعال‌ساز و جلوگیری از عملیات بی‌اثر، یک روش کارآمد برای کاهش هزینه محاسبات و سربار حافظه است. اساساً، در شبکه‌های عصبی عمیق دو سطح تنکی وجود دارد: سطح کلمه و سطح بیت. اکثر شتاب‌دهنده‌های مبتنی بر تنکی که اخیراً پیشنهاد شده‌اند، تنکی سطح کلمه را هدف قرار داده‌اند که مبتنی به از بین بردن ضریب‌های صفر عملوند و ذخیره وزن‌ها به‌صورت فشرده هستند. چندین معماری برای بهره بردن از تنک بودن بیت‌ها در داخل هر کلمه پیشنهاد شده است، از جمله صرفه‌جویی در انرژی لازم برای محاسبات، برای مقادیر صفر فعال‌سازها [۲۹، ۳۰]، یا وزن‌ها [۳۲] یا هر دو [۳۱، ۳۸]. لاکونیک [۳۷] تنکی سطح بیت را هدف قرار داده و یک شتاب‌دهنده سخت‌افزاری برای بهره بردن از تنکی در وزن و فعال‌ساز پیشنهاد داده است. برای این منظور، Laconic هر ورودی را به فهرستی از توان‌های دو علامت‌دار یا عباراتی با رمزگذاری بوث تبدیل می‌کند و سپس عبارات را به‌صورت سریال ضرب و جمع می‌کند. هر چند هر عملیات ضرب وزن در فعال‌ساز به سیکل‌های متعددی نیاز دارد، اما عملکرد کلی این روش بهتر است. درحالی‌که معماری‌های قبلی عمدتاً بر بهره بردن از تنکی در پارامترهای با دقت بالا (ممیز شناور) یا پهنای بیت بالا متمرکز بودند، این کار به دنبال یک شتاب‌دهنده است که به‌طور مؤثر هم از نظر تنکی در سطح کلمه و هم در سطح بیت، هر دو پارامتر، وزن

جهت چندی‌سازی شبکه را ارائه خواهیم داد که در بخش ایده مقاله به تفصیل درباره آن صحبت خواهیم کرد.

۳- کارهای پیشین

با توجه به محدودیت‌های دستگاه‌های نهفته، یعنی قدرت محاسباتی محدود، ظرفیت کم حافظه و محدودیت مصرف توان، چندین هدف برای بهینه‌سازی وجود دارد. برای مثال، ممکن است کسی تصمیم بگیرد که صحت شبکه عصبی را قربانی کند تا مصرف حافظه، یا میزان پیچیدگی محاسبات را کاهش دهد. بسته به هدف بهینه‌سازی، شبکه‌های عصبی طیف وسیعی از فرصت‌های بهینه‌سازی را ارائه می‌دهند. همان‌طور که در شکل ۱ نشان داده شده است [۵]، این فرصت‌ها را به چند دسته کلی تقسیم می‌کنیم:



شکل ۱. سطوح بهینه‌سازی شبکه‌های عصبی

سطح طراحی شبکه: بهینه‌سازی در سطح طراحی شبکه مجموعه‌ای از روش‌هایی است که ساختار شبکه را قبل یا در حین آموزش برای بهبود تأخیر یا هزینه استنتاج شبکه تنظیم می‌کند. نمونه‌هایی از این موارد MobileNet-V1/V2 [۳۹]، SqueezeNet [۷] و ShuffleNet [۴۰] هستند که به صورت دستی طراحی شده‌اند. مجموعه‌ای از کارهای جدیدتر، جستجوی معماری عصبی^۱ را به‌عنوان روشی برای کاهش دخالت انسان برای طراحی چنین معماری‌های پیچیده و زمان‌بری معرفی کردند. نمونه‌هایی از این موارد MNASnet [۴۱]، FbNet [۴۲] و Lemonade [۴۳] هستند که از NAS آگاه از سخت‌افزار از طریق یادگیری تقویتی، الگوریتم‌های تکاملی یا روش‌های مبتنی بر گرادینت برای کشف ساختارهای شبکه عصبی با دقت خوب و عملکرد بالا استفاده کردند.

سطح بهبود شبکه: روش‌های بهبود در سطح شبکه‌های عصبی بر اصلاح نحوه نمایش شبکه با بهره بردن از انواع داده‌های با دقت پایین‌تر (چندی‌سازی) و یا استفاده از تنکی وزن‌ها و فعال‌سازها (هرس کردن) مبتنی هستند. در بهبود شبکه، اغلب برای اینکه به

^۱ Neural architecture search (NAS)

نمایش اعداد دو حالتی و غیریکنواخت دو حالتی و نحوه‌ی محاسبات ضرب در آن بررسی می‌شود. سپس در بخش بعدی سیستم نمایش اعداد غیریکنواخت سه حالتی پیشنهادی ارائه می‌شود و در مرحله بعد عملیات چندی‌سازی و تنکی درحین یادگیری شبکه انجام می‌شود تا با استفاده از حداقل بیت و حداکثر تنکی در پارامترهای وزن، حداقل افت را در صحت شبکه داشته باشیم. در گام بعدی جهت کاهش حجم موردنیاز برای ذخیره‌سازی پارامترهای تنک شده‌ی شبکه در حافظه‌ی خارجی تراشه، الگوریتم‌های فشرده‌سازی مورد بررسی قرار می‌گیرد. در پایان، معماری واحدهای محاسباتی پایه‌ای برای نمایش اعداد سه حالتی و دو حالتی ارائه شده است تا بتوانیم نتایج حاصل از روش‌های چندی‌سازی به همراه تنکی در سیستم نمایش اعداد دو حالتی و سه حالتی غیریکنواخت مختلف را در سطح نتایج سخت‌افزاری موردبررسی قرار دهیم.

۴-۱- ادبیات و مفاهیم پایه‌ای

۴-۱-۱- شبکه‌های عصبی عمیق دو حالتی

همان‌طور که در بخش مقدمه نیز بیان شد، یکی از شبکه‌های عصبی عمیق معروف، شبکه‌های عصبی کانولوشنی است، که از آن‌ها به‌طور گسترده در تشخیص اشیاء، کلاس‌بندی تصاویر و ... استفاده می‌شود. این شبکه‌ها از چندین لایه شامل کانولوشن، جمع‌بندی^۲ و تمام‌اتصال^۳ تشکیل شده‌اند. در این شبکه‌ها به‌منظور استخراج ویژگی، از لایه‌های کانولوشن و جمع‌بندی استفاده می‌شود و لایه تمام‌اتصال در لایه آخر شبکه برای کلاس‌بندی خروجی‌ها طراحی شده است.

درواقع در این مدل شبکه‌های عمیق، بیشترین بار محاسباتی در لایه‌های کانولوشنی و تمام‌اتصال است، علت آن این است که این شبکه‌ها حاوی میلیون‌ها وزن هستند و ضرب این میزان وزن در فعال‌سازی آن‌ها منجر به هزینه بالا در حافظه و محاسبات پیچیده و سنگین می‌شود. این هزینه‌های بالای حافظه و محاسبات مانع بزرگی برای پیاده‌سازی شبکه‌های عصبی کانولوشنی بر روی دستگاه‌های نهفته با محدودیت منابع تلقی می‌شود.

برای غلبه به موانع موجود برای پیاده‌سازی CNNها بر روی دستگاه‌های نهفته، روش‌هایی برای فشرده‌سازی این شبکه‌ها وجود دارد که یکی از این روش‌ها، چندی‌سازی است. درواقع هدف اصلی چندی‌سازی این است که نوع نمایش وزن‌ها و فعال‌سازی‌ها را به‌جای ۳۲ بیت ممیز شناور، با چند بیت محدود نشان دهد. مقاله‌های

و فعال‌سازی را کنترل کند که منجر به عملیات منطقی ساده با پهنای بیت پایین شده است.

سطح بهبود الگوریتم: در این بخش روی بهینه‌سازی عملیات هر لایه تمرکز می‌شود به‌عنوان مثال روش‌های مختلفی برای انجام کانولوشن وجود دارد یکی از آن‌ها روش مستقیم^۱ است در این روش کانولوشن را با ۶ حلقه تودرتو پیاده‌سازی می‌کنند [۴۵]. روش winogard [۴۶] یک روش دیگر برای بهبود الگوریتمی به‌منظور کاهش هزینه‌های محاسباتی عملیات کانولوشن در شبکه‌های عصبی کانولوشنی است. این الگوریتم با تبدیل عملیات سنگین و پیچیده مانند عملیات ضرب به چندین عملیات ساده مانند جمع باعث کاهش پیچیدگی کلی عملیات کانولوشن می‌شود. پیاده‌سازی این الگوریتم‌ها برای سیستم‌های نهفته با توان کم مناسب است، زیرا منابع و بودجه توان بسیار محدود است. در مقابل، آن‌ها هزینه بیشتری در مصرف حافظه و دقت دارند [۱۴].

سطح طراحی اولیه: در این قسمت به پایین‌ترین سطح بهبود نرم‌افزاری یک شبکه عصبی خواهیم پرداخت. بهینه‌سازی در این سطح می‌تواند از تغییر چیدمان داده در حافظه [۴۷]، تغییر نحوه استفاده از دستورالعمل‌های برداری برای پردازش لایه‌ها [۴۸] تا نوشتن پیاده‌سازی اسمبلی هسته‌ها برای رسیدن به عملکرد بهتر پردازنده‌های خاص متفاوت باشد.

با هدف بهره بردن از مزایای تنکی و چندی‌کردن در شبکه‌های عصبی عمیق، در این مقاله می‌خواهیم برای اولین بار یک روش چندی‌سازی غیریکنواخت سه حالتی همراه با تنکی را ارائه دهیم. با استفاده از این روش علاوه بر حفظ صحت شبکه از محاسبات سبک‌بیتی به همراه حذف عملیات ضرب بی‌تاثیر حداکثر بهره را جهت فشرده‌سازی شبکه می‌بریم. انتظار داریم کارایی این شبکه نسبت به حالت‌هایی که تنها از روش چندی‌کردن و یا فقط روش تنکی استفاده می‌کنند بهتر باشد. درواقع انتظار می‌رود مصرف حافظه، پیچیدگی عملیات و انرژی و توان مصرفی نسبت به حالتی که تنها از یکی از روش‌ها استفاده می‌کند بیشتر کاهش یابد و شبکه ایجادشده به محدودیت‌های موجود در دستگاه‌های نهفته غلبه کند.

۴- چندی‌سازی غیریکنواخت سه حالتی

در این بخش از مقاله در ابتدا مقدمه‌ای از ساختار لایه‌های کانولوشنی دو حالتی و مبانی عملکرد چندی‌سازی جهت کاهش تعداد بیت موردنیاز برای پارامترهای شبکه بیان می‌شود. در ادامه

³ Fully connected

¹ Direct

² Pooling

عملیات Xnor_PopCount در معادله ۳ نشان داده شده است، در این معادله \vec{w} و \vec{a} به ترتیب نشان‌دهنده بردار ممیز شناور وزن و فعال‌ساز، \vec{B}_a و \vec{B}_w بیانگر بردار علامت وزن و فعال‌ساز، γ_a و γ_w نیز ضریب مقیاس برای وزن‌ها و فعال‌سازها و بردارهای $Enc(\vec{B}_w)$ و $Enc(\vec{B}_a)$ نیز بردارهای کدگذاری شده از روی بردار علامت هستند. در واقع همان‌طور که در معادله ۳ نشان داده شده است، برای ضرب وزن‌ها و فعال‌سازها کافی است ضریب مقیاس وزن‌ها و فعال‌سازها از آن بردارها جدا شود و آن‌ها به‌تنهایی ضرب شوند و سپس بردارهای علامت وزن‌ها و فعال‌سازها کدگذاری شود ($\rightarrow -1$ $\rightarrow +1, 0$) و درنهایت بین بردارهای کدگذاری شده، عملیات بیتی ساده انجام پذیرد.

$$\begin{aligned} \text{dot}(\vec{w}, \vec{a}) &= \gamma_w \gamma_a \text{dot}(\vec{B}_w, \vec{B}_a) = \\ & \gamma_w \gamma_a \text{XnorPopcount}(Enc(\vec{B}_w), Enc(\vec{B}_a)) \end{aligned} \quad (3)$$

عملیات Xnor_PopCount که در معادله ۳ به آن اشاره شد، در شکل ۲ با ذکر مثال توضیح داده شده است. در این مثال قرار است یک بردار وزن در بردار فعال‌ساز نظیرش که هر دو شامل چهار المان تک‌بیتی هستند ضرب شوند. در شکل ۲ (الف) بردار وزن‌ها و فعال‌سازها، بردارهای چندی شده دو حالتی به همراه ضریب مقیاس هستند که برای انجام عملیات کانولوشن بر روی آنها نیاز است که هر کدام از وزن‌ها در فعال‌ساز نظیرش ضرب شود و درنهایت با یکدیگر جمع شوند. اما به دلیل اینکه ضریب مقیاس وزن‌ها و فعال‌سازها یکسان هستند، روش ساده‌تر، جدا کردن ضریب مقیاس وزن‌ها و فعال‌سازها همانند شکل ۲ (ب) است، که درنهایت خروجی جمع در ضرایب مقیاس ضرب می‌شوند. در مدل (ب) به دلیل اینکه نیازی نیست در هر بار ضرب وزن در فعال‌ساز نظیرش، ضرایب مقیاس ۳۲ بیتی نیز ضرب شوند و بعد از ضرب مقادیر دو حالتی، در خروجی حاصله ضرایب مقیاس را اعمال می‌کنیم، پس بنابراین عملیات ضرب مدل (ب) نسبت به مدل (الف) بهتر است، اما هنوز عملیات ضرب بین وزن‌ها و فعال‌سازها پابرجاست. روشی که پیچیدگی محاسبات ضرب در شبکه‌های دو حالتی را کاهش می‌دهند روش Xnor_PopCount است (شکل ۲ (پ))، به این صورت که بردارهای علامت که حاوی ضرایب مقیاس نیستند را کد می‌کند و بیت‌های وزن و فعال‌ساز نظیرش را Xnor کرده و سپس عملیات PopCount را به کمک فرمول 2P-N روی خروجی Xnor اعمال خواهیم کرد. به این‌صورت که دو برابر تعداد یک‌ها (p) در بردار خروجی را منهای N یا اندازه بردار خروجی (تعداد بیت‌های بردار خروجی) می‌کنیم و در نهایت ضرایب مقیاس وزن‌ها و فعال‌سازها را در خروجی ضرب خواهیم کرد.

زیادی از نمایش اعداد ممیز ثابت، با عرض ۱۶، ۸ و حتی ۴ بیت جهت چندی‌سازی شبکه‌های عصبی عمیق استفاده کرده‌اند و حداکثر حدود یک درصد افت در صحت شبکه داشته‌اند. از طرفی شبکه‌های عصبی عمیق دو حالتی جهت کاهش حداکثری حجم پارامترها و حجم محاسبات معرفی شده‌اند. در چندی‌سازی شبکه‌های عصبی عمیق دو حالتی هر پارامتر از وزن‌ها و فعال‌سازها تنها با یک بیت نمایش داده می‌شود. یکی از ساده‌ترین روش‌های چندی‌سازی دو حالتی بر روی وزن‌ها، با استفاده از تابع Sign است (معادله ۱)، که در آن w و w^b نشان‌دهنده وزن‌های با دقت کامل و وزن‌های دو حالتی می‌باشند.

$$w^b = \text{Sign}(w) = \begin{cases} +1 & w \geq 0, \\ -1 & w \leq 0. \end{cases} \quad (1)$$

البته برای جبران قسمتی از افت صحت شبکه ناشی از چندی‌سازی دو حالتی، معمولاً در نمایش اعداد دو حالتی از یک ضریب مقیاس استفاده می‌شود که این ضریب مقیاس می‌تواند برای گروهی از پارامترها مشترک باشد [۲، ۲۴]. به‌عنوان مثال برای وزن‌های هر فیلتر در یک لایه‌ی کانولوشنی، از یک ضریب مقیاس مستقل و یکسان استفاده می‌شود. لذا مقدار عددی پارامترهای چندی شده w^q به‌صورت معادله ۲ نمایش داده می‌شود.

$$w^q = \gamma * \text{Sign}(w) = \begin{cases} +\gamma & w \geq 0, \\ -\gamma & w \leq 0. \end{cases} \quad (2)$$

به دلیل آن که تعداد وزن‌ها در شبکه‌های کانولوشنی بسیار زیاد است، با چندی‌کردن وزن‌ها حجم داده‌های ذخیره‌شده در حافظه به‌شدت کاهش پیدا می‌کند. در واقع اگر از n بیت برای چندی‌سازی استفاده کنیم، میزان فشردگی شبکه در مقایسه با شبکه ممیز شناور ۳۲ بیتی، $\frac{32}{n}$ خواهد شد. به‌طور مثال زمانی که چندی‌سازی به کمک معادله ۲ باشد، یعنی هر عدد ۳۲ بیتی ممیز شناور به یک عدد یک بیتی دو حالتی چندی شود، در این صورت مقدار n، برابر با یک خواهد بود و حجم وزن‌های شبکه چندی شده ۳۲ برابر نسبت به شبکه اصلی فشرده شده است. چندی‌سازی علاوه بر اینکه حجم حافظه را کم می‌کند عملیات محاسباتی را نیز ساده‌تر خواهد کرد. یعنی به‌جای آن که وزن‌ها و فعال‌سازها را با عملیات ممیز شناور در هم ضرب کنیم، در صورتی که چندی‌سازی از نوع ممیز ثابت باشد، این عملیات با پیچیدگی به‌مراتب کمتر انجام خواهد شد. به‌علاوه اگر عملیات چندی‌سازی دو حالتی بر روی وزن‌ها و فعال‌سازها پیاده‌سازی شود، به‌جای عملیات سنگین ضرب، از عملیات بیتی که شامل عملیات Xnor و PopCount هست استفاده می‌شود که با این روش نیز می‌توان سنگینی و پیچیدگی محاسبات لایه کانولوشن را به‌طور چشمگیری کاهش داد.

برای اعداد صحیح بدون علامت نشان داده می‌شود. فرض می‌شود مقدار اسکالر q با عرض k رقم وجود دارد این دو مقدار را با استفاده از دو بردار ضریب مقیاس (γ) و دو حالتی (b) همانند معادله ۴ می‌توان نمایش داد.

$$\gamma = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ 2^{k-1} \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix}, b_i \in \{0,1\} \quad (4)$$

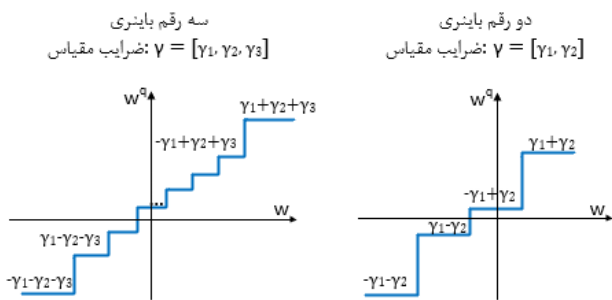
$$\Rightarrow q = \gamma^T b$$

مقادیر قابل نمایش در این حالت به صورت یکنواخت و با فاصله‌ی یکسان از یکدیگر قرار می‌گیرند. در نمایش اعداد غیریکنواخت دو حالتی بردار ضریب مقیاس لزوماً مقادیر توان دو نیست و می‌تواند هر مقداری داشته باشد. لذا برای مقدار اسکالر q که به صورت k رقم دو حالتی هستند نحوه‌ی نمایش اعداد به صورت معادله ۵ می‌شود.

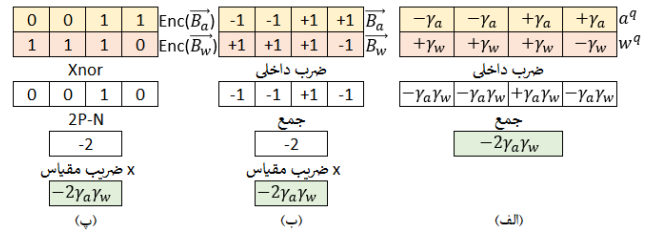
$$\gamma = \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_k \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix}, b_i \in \{-1,1\}, \quad (5)$$

$$\gamma_i \in Z \rightarrow q = \gamma^T b = \gamma_1 b_1 + \gamma_2 b_2 + \dots + \gamma_k b_k$$

این سبب می‌شود تا سطوح مقادیر قابل نمایش با k رقم به صورت یکنواخت نباشد. شکل ۳ مقادیر قابل نمایش در این سیستم اعداد به ازای ۲ و ۳ رقم را نشان می‌دهد.



شکل ۳. سطوح چندی‌سازی غیریکنواخت دو حالتی به ازای ۲ و ۳ رقم در مقاله‌های [۲۴, ۲] نشان داده شده که مقادیر وزن‌ها و فعال‌سازها یک توزیع غیریکنواخت دارند، لذا برای نمایش دادن این پارامترها از سیستم نمایش اعدادی باید استفاده کرد که سطوح نمایش آن به صورت غیریکنواخت باشد، تا بتوانیم حداکثر بهره‌وری را از سیستم نمایش اعداد داشته باشیم. لذا سیستم نمایش اعداد غیریکنواخت دو حالتی می‌تواند گزینه‌ی مناسب‌تری برای نمایش پارامترهای یک شبکه عصبی عمیق با توزیع غیریکنواخت باشد [۵۰].



شکل ۲. نمایش عملیات ضرب کانولوشنی بین دو بردار وزن و فعال‌ساز با استفاده از: (الف) بردارهای چندی شده‌ی وزن‌ها و فعال‌سازها. (ب) بردارهای علامت دو حالتی وزن‌ها و فعال‌سازها. (پ) بردارهای کدگذاری شده وزن‌ها و فعال‌سازها و عملیات Xnor_PopCount

۴-۱-۲- چندی‌سازی با سیستم نمایش اعداد غیریکنواخت دو حالتی

هر چند شبکه‌های عصبی عمیق دو حالتی سنتی از نظر حافظه و حجم محاسبات به حداقل میزان ممکن کاهش پیدا کرده بودند ولی اختلاف صحت شبکه در مقایسه با شبکه‌ی دقت کامل باعث شده بود تا در بسیاری از کاربردها که میزان صحت شبکه از اهمیت بالایی برخوردار هست نتوانند مورد استفاده قرار بگیرد. تحقیقات زیادی انجام شد و جهت بهبود این مشکل راه‌های متفاوتی ارائه شده‌است. استفاده از وزن‌های سه حالتی (سه‌گانه) به جای دو حالتی امکان نمایش وزن‌های صفر را برای یادگیری شبکه عصبی عمیق فراهم می‌کند که خود سبب بهبود صحت شبکه در مقایسه با شبکه‌های عصبی عمیق دو حالتی می‌شود. جهت کدگذاری هر پارامتر سه حالتی که در واقع سه مقدار متفاوت را شامل می‌شود، در سخت‌افزار نیاز به استفاده از دو بیت خواهد بود، لذا علاوه بر سربار محاسباتی بیشتر نسبت به مقادیر دو حالتی، نیاز به حافظه‌ی دو برابری نسبت به وزن‌های دو حالتی خواهد داشت. در تحقیقات دیگری، جهت جبران صحت از دست‌رفته‌ی شبکه‌های عصبی عمیق دو حالتی، استفاده از افزایش ابعاد عرضی شبکه (افزایش تعداد فیلترهای هر لایه) پیشنهاد شده است [۲۳]. هر چند این روش‌ها نتایج قابل‌ملاحظه‌ای در جهت بهبود صحت شبکه داشته‌اند، ولی در این روش‌ها نیاز است تا شبکه با فرا پارامترهای جدید طراحی و یادگیری شود و علاوه بر آن هنوز قابلیت نمایش پارامترهای صفر را ندارد. لذا سربارها و محدودیت‌های مخصوص به خود را دارد. نهایتاً برای اینکه دقت شبکه‌های دو حالتی افزایش یابد و همچنان از مزایای محاسبات سبک بیتی دو حالتی بهره برده شود، روش چندی‌سازی غیریکنواخت دو حالتی در مقاله‌های [۲۴, ۲۴, ۴۹] بیان شده است که در ادامه جزئیات بیشتری در رابطه با این سیستم نمایش اعداد ذکر شده است.

سیستم نمایش اعداد غیریکنواخت دو حالتی: ابتدا نحوه‌ی نمایش

مبنای عملیات باقی مانده گیری عمل می کند استفاده کرده است. در این الگوریتم مقادیر ضریب مقیاس جزء پارامترهای باقابلیت آموزش در حین یادگیری شبکه تعریف شده اند و بردار K رقمی دو حالتی متناظر با هر ورودی مطابق این الگوریتم به دست می آید.

الگوریتم ۱ - نحوه ی چندی سازی غیریکنواخت دو حالتی به روش Residual

Inputs: $\gamma_1, \gamma_2, \dots, \gamma_M, x$

Outputs: b_1, b_2, \dots, b_K

1: $r \leftarrow x$

2: **for** $i = 1 \dots K$ **do**

3: $b_i \leftarrow \text{Binarize}(\text{Sign}(r))$

4: $r \leftarrow r - \text{Sign}(r) * \gamma_i$

5: **end for**

مقاله ی [۲] با استفاده از روش چندی سازی جدیدی با استفاده از نمایش اعداد غیریکنواخت دو حالتی صحت های قابل مقایسه باحالت دقت کامل به دست آورده است. [۲] برای به دست آوردن مقادیر بهینه ی ضریب مقیاس ها و مقادیر بردار دو حالتی متناظر معادله ۸ را به صورت معادله ۹ جایگزین کرده است.

$$\gamma^*, B^* = \arg \min \|B^T \gamma - x\|_2^2 \quad (9)$$

در این فرمول $x = [x_1, \dots, x_N] \in R^N$ (وزن ها یا فعال سازها) ممیز شناور، K نشان دهنده تعداد ارقام چندی شده است، همان طور که گفته شد هدف پیدا کردن بردار $\gamma \in R^K$ بهینه به همراه $B = [b_1, \dots, b_N] \in \{-1, 1\}^{K \times N}$ است. در واقع هدف به دست آوردن بردار ضریب مقیاس و بردار دو حالتی است که مقدار خطای چندی سازی را به حداقل برساند. برای بهینه سازی کوانتایزر دو روش مورد بررسی قرار گرفته است. در روش اول از روش سنتی پس انتشار^۱ استفاده شده است به این صورت که مقادیر γ ها همراه با پارامترهای شبکه آموزش می بینند و مقادیر B به صورتی انتخاب می شود که مقدار چندی شده ی $Q(x)$ حداقل فاصله را از مقدار واقعی x داشته باشد. در روش دوم از الگوریتم حداقل خطای کوانتایزر یا QEM استفاده کرده است. برای حل معادله ۹ در حلقه اصلی آموزش شبکه، ابتدا γ را ثابت در نظر گرفته و B را به صورتی انتخاب می کنیم که $Q(x)$ به سطح چندی ای نگاشت شود که نزدیک ترین فاصله را از خود x داشته باشد. سپس برای به دست آوردن γ بهینه از فرمول بسته ی معادله ۱۰ استفاده می شود. نتایج صحت شبکه برای این دو روش در مقاله [۲] نشان داده است که روش QEM بهتر عمل کرده و شبکه می تواند به صحت نزدیک به

حال بدون از دست دادن کلیت مسئله عملیات ضرب داخلی دو بردار w^q و a^q را در نظر می گیریم که به روش مذکور چندی شده و شامل N پارامتر k رقمی دو حالتی است. ضرب داخلی بردار به عنوان سنگین ترین محاسبات لایه های کانولوشنی و تمام اتصال مورد بررسی قرار می گیرد. در معادلات بیان شده، w نشان دهنده وزن ها و a نیز بیانگر فعال سازها هستند. ماتریس $B^w = [b_1^w, b_2^w, \dots, b_N^w]$ و $B^a = [b_1^a, b_2^a, \dots, b_N^a]$ را به صورتی که $B^w \in \{-1, 1\}^{K \times N}$ و $B^a \in \{0, 1\}^{K \times N}$ در نظر بگیرد (توجه شود که مقادیر ماتریس وزن ها $(-1, 1)$ و فعال سازها $(0, 1)$ هستند)، مقدار ماتریس $\gamma^w \in R^K$ برای تمام مقادیر وزن ها و فعال سازها یکسان و به ترتیب با $\gamma^a \in R^K$ و برای یک مجموعه با N زوج عدد به صورت معادله ۶ و ۷ به دست می آید.

$$\text{DotProduct}(w^q, a^q) = \sum_{n=1}^N \sum_{i=1}^K \sum_{j=1}^K (b_{ni}^w * \gamma_i^w * b_{nj}^a * \gamma_j^a) \quad (6)$$

$$\text{DotProduct}(w^q, a^q) = \sum_{i=1}^K \sum_{j=1}^K (\gamma_i^w * \gamma_j^a \sum_{n=1}^N b_{ni}^w * b_{nj}^a) \quad (7)$$

با فاکتورگیری از معادله ۶ با توجه به یکسان بودن مقادیر (i, j) در N زوج عدد مختلف، میتوان به معادله ۷ رسید که $\sum_{n=1}^N b_{ni}^w * b_{nj}^a$ قسمت عملیات بیتی دو حالتی را نشان میدهد که میتواند مطابق قسمت قبل با استفاده از عملیات Xnor_PopCount انجام شود. در ادامه روشهای مختلف چندی سازی به سیستم نمایش اعداد غیریکنواخت دو حالتی و نحوه ی به دست آوردن مقادیر ماتریس B و γ که هم زمان با عملیات یادگیری شبکه به دست می آید توضیح داده شده است.

چندی سازی با استفاده از سیستم نمایش اعداد غیریکنواخت دو حالتی: معادله ۸ مسئله بهینه سازی به حداقل رساندن خطای چندی سازی جهت به دست آوردن مقدار چندی $Q(x)$ به ازای ورودی x را نشان می دهد.

$$Q^*(x) = \arg \min \int (Q(x) - x)^2 dx \quad (8)$$

الگوریتم های متفاوتی جهت به دست آوردن مقدار چندی $Q(x)$ در سیستم نمایش اعداد غیریکنواخت دو حالتی معرفی شده است. مقاله [۲۴] از یک الگوریتم پیشنهادی residual مطابق الگوریتم ۱ که بر

^۱ Back propagation

است که در ابتدا و به یکباره ۷۰ درصد وزن‌ها را هرس کنیم. نتایج صحت در این دو حالت با یکدیگر تفاوتی ندارد [۵۱].

به علاوه زمان اعمال تنکی به شبکه نیز در دو حالت همزمان با یادگیری و بعد از یادگیری می‌باشد، که در روش همزمان با یادگیری میزان افت صحت نسبت به زمان بعد از یادگیری کمتر است، لذا در این مقاله روش هرس کردن ما در هنگام یادگیری و به یکباره می‌باشد.

انواع فشرده‌سازی حافظه با استفاده از پراکندگی وزن‌ها: پس از تنک کردن شبکه، به دلیل اینکه مقادیر غیرصفر وزن‌ها از ۲۰ تا ۸۰ درصد و برای فعال‌سازها از ۵۰ تا ۷۰ درصد در سطح کلمه کاهش می‌یابد، می‌توان از مزایای تنکی ایجاد شده در وزن‌ها و فعال‌سازها بهره برد و علاوه بر کاهش محاسبات، میزان حافظه لازم برای ذخیره‌سازی داده‌ها را در حافظه‌ی ثانویه کاهش داد که این خود سبب کاهش انرژی دسترسی به حافظه‌ی ثانویه می‌شود.

۴-۲-۲- ارائه مدل چندی‌سازی آگاه از توزیع تراکم صفر

در این بخش از مقاله ایده چندی‌سازی آگاه از تراکم صفر با استفاده از نمایش اعداد غیریکنواخت را مورد بررسی قرار می‌دهیم.

۴-۲-۱- چندی‌سازی غیریکنواخت سه حالتی

در این بخش به توضیح نحوه نمایش اعداد استفاده‌شده در این مقاله، که اعداد غیریکنواخت سه حالتی هستند می‌پردازیم. همان‌طور که در بخش‌های گذشته بیان شد برای پر کردن فضای خالی بین صحت شبکه‌های دو حالتی و شبکه‌های ممیز شناور، شبکه‌های غیریکنواخت دو حالتی معرفی شدند که تا حدودی توانستند مسئله مذکور را برطرف نمایند. اما مشکلی که در شبکه‌های غیریکنواخت دو حالتی وجود دارد این است که این شبکه‌ها قدرت نمایش دادن عدد صفر را ندارند در صورتی که در بسیاری از مقاله‌ها نشان داده شده است که شبکه‌های عصبی عمیق از پراکندگی بالای وزن‌های صفر بهره می‌برند. لذا در شبکه‌های چندی شده با استفاده از نمایش اعداد غیریکنواخت دو حالتی نمی‌توانیم از مزایای تنکی استفاده کنیم. این چالش‌انگیزه‌ای شد تا چندی‌سازی غیریکنواخت سه حالتی را معرفی کنیم، که علاوه بر بهبود صحت شبکه و استفاده از محاسبات سبک بیتی، علاوه بر ۱ و -۱ می‌توانیم عدد صفر را نیز در نمایش هر سطح دخیل کنیم تا کلمه‌های کامل صفر قابل نمایش باشند. در ادامه به بررسی روند چندی‌سازی مطرح‌شده در این مقاله خواهیم پرداخت.

اعداد غیریکنواخت سه‌حالتی علاوه بر اینکه تعداد سطوح بیشتری

شبکه با دقت کامل برسد.

$$\gamma^* = (BB^T)^{-1}Bx \quad (10)$$

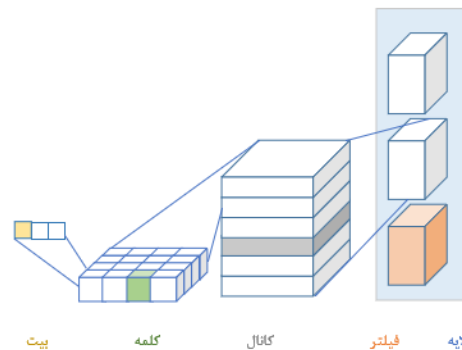
این نوع چندی‌سازی غیریکنواخت دو حالتی برای وزن‌ها به صورت آگاه به کانال و برای فعال‌سازها به صورت آگاه به لایه انجام می‌گیرد. مبنای الگوریتم کوانتایزر استفاده شده در این مقاله، مطابق روش QEM در [۲] هست.

۴-۱-۳- انواع دانه‌بندی پراکندگی و نحوه‌ی تنک کردن

شبکه‌های عصبی عمیق

انواع دانه‌بندی پراکندگی:

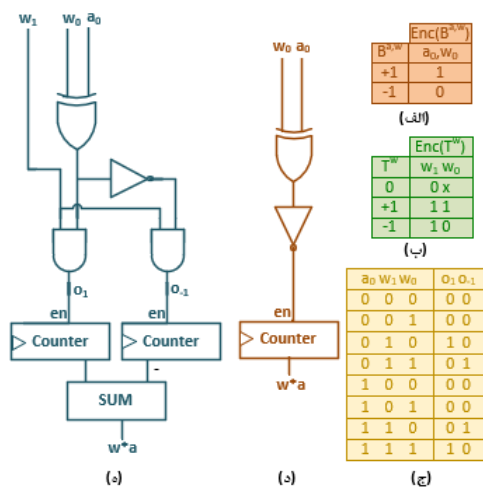
از روش‌های دیگری که باعث بهبود فشرده‌سازی شبکه‌های عصبی می‌شوند، روش هرس کردن است. هرس کردن شبکه‌های عصبی ابعاد گسترده‌ای دارد که در شکل ۴ به‌طور کامل دانه‌بندی تنکی در شبکه‌های عصبی عمیق نشان داده شده است. در این شکل از سمت راست به چپ، در ابتدا تنکی در سطح لایه است که در هر لایه یک سری فیلتر قرار دارد. تنکی می‌تواند در سطح فیلتر اعمال شود و داخل هر فیلتر، می‌توان تنکی در سطح کانال داشته باشیم و داخل هر کانال نیز در سطح کلمه می‌توان تنکی اعمال کرد و برای هر کلمه نیز در سطح بیت می‌توان عملیات هرس کردن را اعمال کرد.



شکل ۴. انواع دانه‌بندی پراکندگی در شبکه‌های عصبی عمیق

نحوه‌ی اعمال تنکی به شبکه: نحوه اعمال تنکی به شبکه به دو صورت افزایشی و به یکباره انجام می‌گیرد. روش افزایشی به اینصورت است که درصد تنکی مورد نظر که می‌خواهیم به شبکه اعمال کنیم را طی چند مرحله با درصدهای کمتر به شبکه تحمیل کنیم تا به درصد مورد نظر دست یابیم. به طور مثال، اگر بخواهیم به شبکه ۷۰ درصد تنکی اعمال کنیم، در روش افزایشی در ابتدا ۱۰ درصد وزن‌ها را تنک کرده سپس از میان وزن‌های باقی‌مانده ۱۰ درصد بعدی را تنک کرده، این کار را تا زمانی که به ۷۰ درصد تنکی در وزن‌ها برسیم ادامه خواهیم داد. اما روش به یکباره به اینصورت

زمانی که n بیت ابتدایی پر شد با توجه به تعداد یک‌هایی که داخل n بیت اول قرار دارد و با توجه به اینکه هرکدام از کلمه‌های وزن‌ها با چند رقم چندی شده‌اند، به تعداد حاصل ضرب یک‌ها در تعداد رقم‌های چندی شده، بیت‌هایی بعد از n بیت اولیه قرار می‌گیرد. به‌طور مثال در شکل ۵ بخش (الف) یک گروه ۸ تایی از مقادیر دو رقمی سه حالتی در نظر گرفته شده که حاوی مقادیر صفر و غیر صفر است و در شکل ۵ (ب) نیز الگوریتم فشرده‌سازی مفروض اعمال شده است. در واقع میزان حافظه لازم برای نشان دادن ۸ عدد دو رقمی سه حالتی، همان طور که در شکل ۵ (ب) نشان داده شده است ۳۲ بیت است، اما میزان حافظه لازم به کمک روش فشرده‌سازی ۲۰ بیت است، به اینصورت که در ابتدا ۸ بیت برای نشان دادن صفر بودن یا نبودن مقادیر تخصیص داده شده و سپس به دلیل اینکه ۳ تا از کلمات این گروه ۸ تایی غیر صفر بود بعد از ۸ بیت اولیه، سه بخش چهار بیتی (۱۲ بیت) برای نشان دادن مقادیر غیر صفر اضافه شده است.



شکل ۶. (الف) کدگذاری ارقام دو حالتی فعال‌سازها و وزن‌ها، (ب) کدگذاری ارقام سه‌حالتی وزن‌ها، (ج) جدول درستی ضرب یک رقم فعال‌ساز دو حالتی کدگذاری شده در یک رقم وزن سه‌حالتی کدگذاری شده متناظر با آن، (د) مدار محاسبات ضرب یک رقم دو حالتی وزن در شده متناظرش، (ه) مدار محاسبات ضرب یک رقم وزن سه‌حالتی در یک رقم فعال‌ساز دو حالتی

۵- نتایج

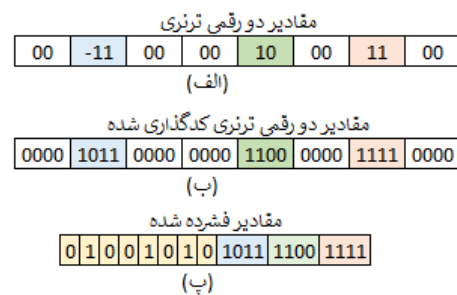
در این بخش از این مقاله، به بررسی نتایج و بهبود شبکه به کمک چندی‌سازی آگاه از تراکم صفر خواهیم پرداخت.

۵-۱- راه‌اندازی

در این مقاله برای به دست آوردن نتایج از مجموعه داده CIFAR10 و شبکه VGG_Small به عنوان یک مجموعه داده و شبکه طبقه‌بند تصویر که محبوبیت زیادی در حوزه ی پردازش تصویر در سیستم‌های نهفته دارد استفاده شده است. در ادامه درباره آن‌ها توضیح مختصری خواهیم داد.

این مجموعه داده شامل نمونه‌هایی از ده کلاس: هواپیما، خودرو،

شکل ۵. الگوریتم فشرده‌سازی برای مقادیر چندی شده دو رقمی تنک



۴- معماری واحدهای محاسباتی دو حالتی و سه حالتی

در این بخش از مقاله به بررسی معماری واحدهای محاسباتی برای عملیات ضرب وزن‌ها در فعال‌سازها در شبکه‌های دو حالتی و سه‌حالتی می‌پردازیم. در ابتدا، برای بررسی معماری واحدهای ضرب، به نحوه کدگذاری مقادیر دو حالتی و سه حالتی خواهیم پرداخت. در شکل ۶ (الف) به کمک یک بیت w_0 یا a_0 یک رقم وزن یا فعال‌ساز دو حالتی (B^a یا B^w) کدگذاری شده‌اند و در شکل ۶ (ب) نیز یک رقم وزن سه حالتی T^w به کمک دو بیت w_0 و w_1 کدگذاری شده است. در ادامه در شکل ۶ (ج) نیز جدول درستی برای ضرب یک رقم دو حالتی در سه حالتی ذکر شده است.

در شکل ۶ (د) معماری ضرب یک رقم وزن دو حالتی در یک رقم فعال‌ساز سه حالتی ترسیم شده است. که در ابتدا عملیات ساده بیتی Xnor بین وزن و فعال‌ساز دو حالتی اعمال می‌شود و سپس

وزن‌ها و فعال‌سازها هستیم به اینصورت که مثلا $A3W2$ به معنای فعال‌ساز سه رقمی و وزن دو رقمی است.

با توجه به نمودار ۱ در حالت $A(f)W(t)$ که فعال‌سازها به صورت ممیزشناور و وزن‌ها به صورت غیریکنواخت سه‌حالتی چندی شده‌اند، صحت شبکه در حالتی که فعال‌سازها 32 بیت و وزن‌ها به صورت یک رقمی سه‌حالتی چندی شده باشند، تقریباً 0.2 درصد از صحت شبکه ممیزشناور کمتر است. اما در صورتی که چندی‌سازی وزن‌ها را از حالت یک رقم سه‌حالتی به چند رقم سه‌حالتی تغییر دهیم، مثل 2 یا 3 رقم سه‌حالتی، صحت شبکه به صحت شبکه ممیز شناور می‌رسد. در واقع در صورتی که فعال‌سازها دقت کامل باشند، با چندی کردن وزن‌ها تا 2 رقم سه‌حالتی به صحت شبکه ممیز شناور خواهیم رسید.

اما در نمودار ۱ در حالتی که هم وزن‌ها و هم فعال‌سازها به صورت دو حالتی $A(b)W(b)$ چندی شده‌اند، زمانی که وزن‌ها به صورت یک رقمی دو حالتی باشند و فعال‌سازها از یک رقم دو حالتی به دو الی سه رقم افزایش پیدا کنند ($A1W1$ و $A2W1$ و $A3W1$) شاهد افزایش صحت در شبکه هستیم. همانطور که در نمودار ۱ نیز مشخص است در حالتی که وزن‌ها را به صورت یک رقمی دو حالتی و فعال‌سازها را به صورت غیریکنواخت دو حالتی چندی کردیم در حالت $A3W1$ به بیشترین صحت خواهیم رسید. به همین منظور در ادامه فعال‌سازها را سه رقمی دو حالتی در نظر می‌گیریم. همانطور که ملاحظه می‌شود در حالت $A3W2$ صحت شبکه 0.05 درصد از صحت شبکه ممیز شناور کمتر است، اما در حالت $A3W3$ صحت شبکه به اندازه صحت شبکه ممیز شناور شده است. بنابراین در حالتی که فعال‌سازها و وزن‌ها به صورت غیریکنواخت دو حالتی چندی می‌شوند، اگر هر کدام از وزن‌ها و فعال‌سازها را با سه رقم دو حالتی چندی کنیم صحت شبکه چندی شده ما به صحت شبکه ممیز شناور خواهد رسید.

اما نوع دیگر چندی‌سازی ارائه شده در نمودار ۱ که ایده این مقاله است، حالت $A(b)W(t)$ است که وزن‌ها به صورت سه‌حالتی و فعال‌سازها به صورت دو حالتی چندی شده‌اند. در این حالت نیز همانند حالت قبل زمانی که وزن‌ها به صورت سه‌حالتی یک رقمی و فعال‌سازها به صورت دو حالتی از یک رقم تا سه رقم افزایش پیدا می‌کند، صحت شبکه همانند شبکه $A(b)W(b)$ افزایش می‌یابد. در این شبکه در صورتی که فعال‌سازها به صورت دو رقمی دو حالتی و وزن‌ها به صورت دو رقمی سه‌حالتی باشند، صحت شبکه با شبکه

پرنده، گربه، آهو، سگ، قورباغه، اسب، کشتی و کامیون است. مجموعه آموزشی شامل 50000 نمونه RGB و مجموعه آزمون شامل 10000 نمونه است. هر نمونه دارای وضوح 32×32 پیکسل است.

شبکه‌ی کانولوشنی VGG-Small که در واقع نسخه‌ی ساده شده‌ی شبکه‌ی VGG-Net استفاده شده در [۴۴، ۵۲] هست، شامل 5 لایه کانولوشن، جمع‌بندی بیشینه^۱، نرمال‌سازی^۲، تابع یکسوساز خطی^۳ و در انتها یک لایه‌ی تمام‌اتصال جهت کلاس‌بندی نتایج است.

نتایج به‌دست‌آمده در بخش بعدی برای 400 دوره^۴ است که در بررسی نتایج از حروف A (بیانگر فعال‌سازها) و W (بیانگر وزن‌ها) و همراه با آن‌ها اعدادی (بیانگر تعداد سطوح) استفاده شده است. همچنین از (b) ، (t) و (f) به ترتیب برای رقم‌های دو حالتی و سه‌حالتی و ممیز شناور استفاده شده است.

۵-۲- نتایج صحت چندی‌سازی غیریکنواخت سه‌حالتی

برای اندازه‌گیری و بررسی صحت چندی‌سازی ارائه شده در این مقاله، چهار حالت در نظر گرفته شد، که نتایج آن‌ها بر روی شبکه VGG_Small با مجموعه داده CIFAR10 در نمودار ۱ نشان داده شده است. اولین حالتی که برای چندی‌سازی در نظر گرفته شد، چندی‌سازی غیریکنواخت دو حالتی است که این نوع چندی‌سازی هم برای وزن‌ها و هم برای فعال‌سازها اعمال شده و در نمودار ۱ آن را با $A(b)W(b)$ نمایش داده‌ایم. در صورتی که تعداد رقم را یک در نظر بگیریم، این حالت نماینده‌ی ای از مقاله‌ی [۴۴] و با افزایش تعداد رقم مشابه مقاله‌ی [۲] خواهد بود. در نوع بعدی، فعال‌سازها به صورت دقت کامل (ممیز شناور) در نظر گرفته شده و فقط وزن‌ها را به صورت غیریکنواخت سه‌حالتی چندی کردیم که با $A(f)W(t)$ در نمودار مشخص شده است. روش بعدی، همان ایده ما در این مقاله است به اینصورت که وزن‌ها به صورت غیریکنواخت سه‌حالتی و فعال‌سازها به صورت غیریکنواخت دو حالتی چندی شدند که در نمودار ۱ به اختصار با $A(b)W(t)$ نشان داده شده است. در این حالت نیز در صورتی که تعداد رقم را یک فرض کنیم نماینده‌ی مقاله‌ی [۲۲] خواهد شد. برای سنجش صحت، مدل‌های مختلف چندی‌سازی‌های ذکر شده را با شبکه ممیزشناور با عنوان $A(f)W(f)$ مقایسه کرده‌ایم. همانطور که در نمودار ۱ مشخص است، در محور عمودی صحت و در محور افقی شاهد تعداد ارقام چندی‌سازی برای

³ Relu

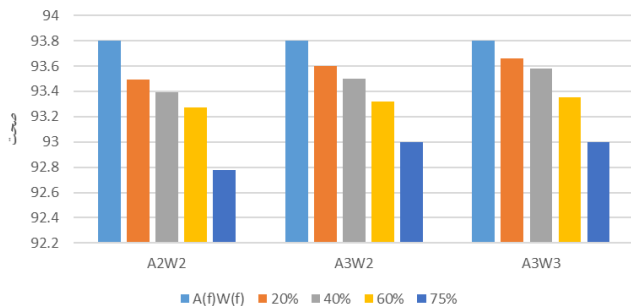
⁴ Epoch

¹ Max-polling

² Batch-normalization

پرداختیم و نشان دادیم که مدل ما در حالت A2W2 صحتی به اندازه شبکه دقت کامل دارد. حال در ادامه می‌خواهیم مدل چندی‌سازی ارائه شده را هرس کنیم. برای این منظور در شبکه A(b)W(t) برای وزن‌ها ۲۰، ۴۰، ۶۰ و ۷۵ درصد، تنگی اعمال کردیم. در اصل کاری که انجام شده به اینصورت است که از میان وزن‌ها درصد مد نظر (۲۰، ۴۰، ۶۰ و ۷۵ درصد) از آنها که نزدیک صفر هستند را همراه با یادگیری شبکه هرس می‌کنیم. در نمودار ۲ صحت شبکه‌هایی که برای درصد‌های مختلف، هرس شده‌اند را گزارش داده‌ایم. همانطور که از نمودار ۲ مشخص است، هر چه درصد هرس کردن را زیاد می‌کنیم، دقت شبکه بیشتر کاهش می‌یابد. به طور مثال در A2W2 در شبکه بدون هرس صحت حدود ۹۳٫۸ است اما وقتی ۲۰ درصد وزن‌ها را هرس می‌کنیم دقت شبکه حدود ۰٫۳ درصد کاهش می‌یابد و برای ۴۰ درصد حدود ۰٫۴ درصد و برای ۶۰ درصد ۰٫۵ درصد و برای ۷۵ درصد حدود ۱ درصد صحت شبکه نسبت به شبکه بدون اعمال هرس کاهش می‌یابد.

صحت شبکه تحت کوانتیزاسیون و تنگی



نمودار ۲. مقایسه صحت شبکه با ترکیب ارقام مختلف وزن و فعال‌ساز تحت درصد‌های تنگی مختلف.

در نمودار ۲ در حالت A3W2 و A3W3 زمانی که ۷۵ درصد وزن‌ها را هرس کنیم، صحت شبکه به ۹۳ درصد خواهد رسید که نسبت به شبکه ممیز شناور تنها ۰٫۸ درصد افت دقت خواهیم داشت. از آن جا که میزان افت صحت در این حالات از A2W2 کمتر است و بعلاوه در حالت A3W2 تعداد بیت کمتری استفاده شده است، بنابراین در ادامه از شبکه‌ی vgg-small با چندی‌سازی A3W2 استفاده خواهیم کرد.

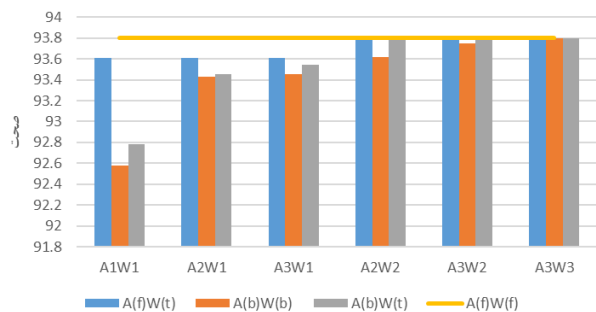
۵-۳-۲- میزان فشردگی‌سازی پارامترها بعد از تنگی شبکه

همانطور که در بخش ۳-۴ نیز بیان شد، به دلیل اینکه نشان دادن ارقام سه حالتی نسبت به ارقام دو حالتی، سربار حافظه ای دو برابر دارند، با استفاده از مزایای تنگی یک نوع فشردگی‌سازی برای ذخیره‌سازی اعداد سه حالتی در حافظه بیان کردیم. نتایج متناظر

ممیز شناور برابر خواهد شد و صحت شبکه‌های A3W2 و A3W3 نیز با شبکه ممیز شناور برابر شد.

در نمودار ۱ صحت شبکه A(b)W(b) در تعداد ارقام مختلف همواره از صحت شبکه A(b)W(t) کمتر است، دلیل آن نیز به این خاطر است که سطوح چندی‌سازی وزن‌ها در مدل ارائه شده در این مقاله، بیشتر از سطوح شبکه دو حالتی است به علاوه صحت شبکه دو حالتی در حالت A3W3 به صحت شبکه ممیز شناور می‌رسد این در حالی است که صحت شبکه سه حالتی در حالت A2W2 به صحت شبکه ممیز شناور خواهد رسید. صحت شبکه A(f)W(t) در اکثر مواقع از مدل ما و شبکه دو حالتی نیز بیشتر است، و دلیل آن هم این است که فعال‌سازها با دقت کامل هستند و فقط وزن‌های آن به صورت غیریکنواخت سه حالتی چندی شده‌اند. اما این شبکه به دلیل آن‌که فعال‌سازهای آن به صورت دقت کامل هستند، مطلوب ما نیست، زیرا این امر به بهبود حافظه و انجام عملیات با پیچیدگی کمتر منجر نخواهد شد.

صحت شبکه تحت کوانتیزاسیون‌های مختلف



نمودار ۱. مقایسه صحت شبکه vgg-small به همراه مجموعه داده cifar10 برای انواع چندی‌سازی‌های مختلف. $A_m W_n$ در محور افقی به معنی m رقم برای چندی‌سازی فعال‌سازها و n رقم برای چندی‌سازی وزن‌ها هست.

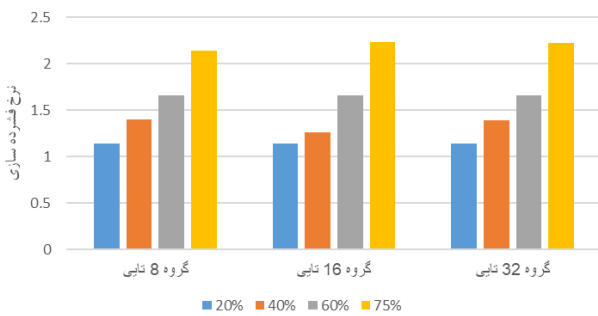
۵-۳-۳- نتایج چندی‌سازی غیریکنواخت سه‌حالتی به همراه هرس کردن

در این بخش ابتدا به بررسی صحت در شبکه‌های چندی شده تنک می‌پردازیم، در ادامه نتایج حاصل از فشردگی‌سازی وزن‌های سه حالتی را نشان می‌دهیم و در نهایت تسریع بالقوه حاصل از شبکه چندی شده‌ی تنک پیشنهاد شده در این مقاله را مورد مطالعه قرار خواهیم داد.

۵-۳-۱- بررسی اثر هرس کردن روی صحت شبکه

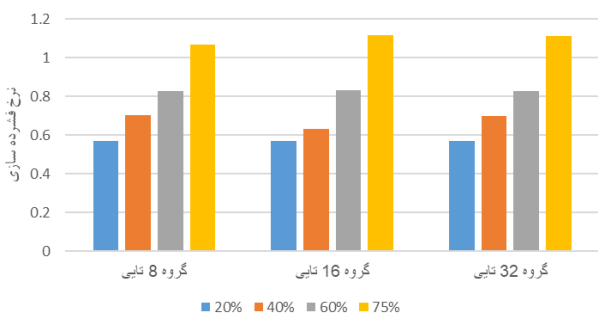
در بخش ۲-۵ به بررسی صحت شبکه در چندی‌سازی‌های مختلف

نرخ فشرده سازي وزن‌ها نسبت به تربي



نمودار ۳. نرخ فشرده‌سازي وزن‌ها در حالت A3W2 سه‌حالي
براي درصدهاي تنكي مختلف به نسبت سه‌حالي

نرخ فشرده سازي وزن‌ها نسبت به باينري



نمودار ۴. نرخ فشرده‌سازي وزن‌ها در حالت A3W2 سه‌حالي
براي درصدهاي تنكي مختلف به نسبت دو‌حالي

در اصل براي حالت A3W2 با هرس كردن وزن‌ها تا ۷۵ درصد، صحت شبكه برابر با ۹۳ شد، كه نسبت به حالت سه‌حالي و دو‌حالي بدون هرس كردن، كه دقت آن ۹۳٫۸ بود، حدود ۰٫۸ درصد صحت آن كمتر شد و بعلاوه همانطور كه در نمودار ۳ مشاهده مي‌كنيد در ۷۵ درصد تنكي، نسبت تعداد بيت فشرده‌سازي از تعداد بيت دو‌حالي كمتر نيز شد. پس در واقع در حالي كه وزن‌هاي سه‌حالي تنك داشته باشيم، تعداد بيت‌هاي لازم براي نشان دادن ارقام سه‌حالي كه دو برابر ارقام دو‌حالي است را به كمك فشرده‌سازي کاهش داده و نشان داديم، در ۷۵ درصد تنكي، تعداد بيت فشرده شده حتي از دو‌حالي نيز كمتر شد. در واقع ميزان حافظه مصرفي شبكه سه‌حالي با ۷۵ درصد تنكي از شبكه دو‌حالي كمتر است. نکته‌ي پاياني هزينه‌ي سخت‌افزاري حاصل از فشرده‌سازي وزن‌ها است. هر الگوريتم فشرده‌سازي، هزينه‌ي اي را به سخت‌افزار تحميل مي‌كند ولي از آن جايي كه روش مذكور تنها از كدگذاري صفرها به صورت يك ساختار مرتب بهره‌مي‌برد هزينه‌ي خروج از فشرده‌سازي خيلي ناچيز هست و مي‌توان با استفاده يك مدار ساده تشكيل شده از يك ثبات تغيير مكان و تسهيم‌كننده، وزن‌ها را از حالت فشرده شده خارج كرد.

با اين عمليات در نمودار ۳ و نمودار ۴ نشان داده شده است. گروه بندي براي ذخيره‌سازي وزن‌ها را ۸، ۱۶ و ۳۲ در نظر گرفتيم، يكبار تعداد بيت سه‌حالي را نسبت به سه‌حالي فشرده شده (نمودار ۳) و بار ديگر تعداد بيت دو‌حالي را نسبت به سه‌حالي فشرده شده (نمودار ۴) اندازه‌گيري كرديم، و اين عمليات را براي درصدهاي مختلف تنكي در حالت A3W2 بدست آورديم. انتظار داريم كه ميزان فشرده‌سازي كه صورت گرفته به اندازه‌اي باشد كه تعداد بيت‌هاي فشرده كه براي نمايش اعداد سه‌حالي نياز داريم نصف تعداد بيت‌هاي لازم براي نمايش اعداد سه‌حالي، يا هم‌اندازه بيت‌هاي لازم براي نشان دادن عدد دو‌حالي متناظرش باشد، تا بتوانيم ادعا كنيم كه فشرده‌سازي كه ارائه داديم مطلوب است و از نظر ما مطلوب است كه نرخ فشرده‌سازي آن نسبت به سه‌حالي ۲ برابر و نسبت به دو‌حالي به اندازه يكسان باشد.

به عنوان مثال در نمودارهاي ۳ و ۴ وزن‌ها را دو رقمي سه‌حالي در نظر گرفتيم و در حالي كه هيچ نوع فشرده‌سازي نداشته باشيم و گروه‌ها ۸ تايي باشند، در كل $2 \times 2 \times 8 = 32$ بيت براي نشان دادن اعداد سه‌حالي و $2 \times 8 = 16$ بيت براي نشان دادن اعداد دو‌حالي نياز خواهيم داشت، اما همانطور كه در نمودار ۳ نشان داده شده است براي گروه‌هاي ۸ تايي در صورتيكه ۲۰ درصد وزن‌ها را هرس كنيم، نرخ فشرده‌سازي نسبت به سه‌حالي تقريبا ۱ است. در همين حالت در نمودار ۴ نرخ فشرده‌سازي نسبت به دو‌حالي تقريبا ۰٫۵ خواهد بود. اين بدان معناست كه تعداد بيت لازم بعد از اعمال عمليات فشرده‌سازي براي وزن‌ها با ۲۰ درصد تنكي تقريبا به اندازه تعداد بيت سه‌حالي تنك نشده است.

همانطور كه در نمودار ۳ و نمودار ۴ مشاهده مي‌شود، هرچه ميزان تنكي افزايش مي‌يابد، ميزان فشرده‌سازي نسبت به سه‌حالي به ۲ نزديك شده و نسبت به دو‌حالي به ۱ نزديك مي‌شود، يعني اگر ميزان تنكي را بيشتر كنيم شاهد برابري تعداد بيت فشرده با دو‌حالي خواهيم بود.

به طور مثال در حالي كه وزن‌ها را تا ۷۵ درصد هرس كنيم و گروه اعداد را ۸ تايي در نظر گرفتيم، ميزان فشرده‌سازي وزن‌ها نسبت به سه‌حالي ۲٫۱۳ و نسبت به دو‌حالي ۱٫۰۶ شده اين بدان معناست كه ميزان فشرده‌سازي در ۷۵ درصد تنكي علاوه بر اينكه تعداد بيت مورد نياز از دو‌حالي بيشتر نيست بلكه توانستيم تعداد بيت مورد نياز را از دو‌حالي نيز كمتر كنيم.

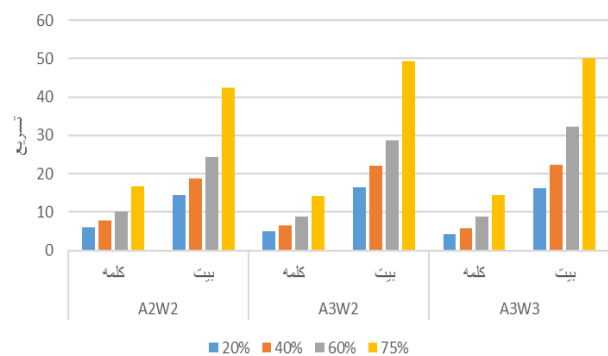
۵-۳-۳- محاسبه‌ی تسریع بالقوه حاصل از تنگی شبکه

برای بدست آوردن تسریع بالقوه ناشی از تنگی وزن‌ها و فعال‌سازها، می‌خواهیم با تشخیص و شناسایی کلمه‌ها و بیت‌های صفر، بررسی کنیم که در صورت انجام ندادن عملیات ضرب با ورودی صفر در شبکه، به چه میزان تسریع بالقوه در محاسبات دست پیدا خواهیم کرد.

برای این منظور زوج‌هایی از وزن‌ها و فعال‌سازها را در نظر گرفتیم و در صورتی که هر کدام از ارقام وزن و فعال‌ساز متناظرش صفر بود، عملیات ضرب بین آن‌ها را نادیده گرفتیم و بدین منظور عملیات ضرب در صفر را در سطح کلمه حذف کردیم. در این حالت تسریع بالقوه را نسبت کل حاصلضرب‌ها به حاصلضرب‌های غیرصفر شده در نظر می‌گیریم (N/Non-zero product). علاوه بر این سعی کردیم عملیات ضرب را در سطح بیتی نیز بررسی کرده و عملیات ضرب بین بیت‌های وزن و فعال‌ساز را در صورتی که یکی از آن‌ها صفر بود در سطح بیت حذف کنیم. نمودار ۵ نشان دهنده میزان تسریع بالقوه عملیات ضرب در شبکه هرس شده، هم در سطح کلمه و هم در سطح بیت می‌باشد. که با توجه به نمودار در تمامی ترکیب‌های ارقام مختلف وزن و فعال‌ساز، هرچقدر که درصد تنگی افزایش می‌یابد، به دلیل حذف بیشتر عملیات ضرب در صفر، میزان تسریع بالقوه هم در سطح کلمه و هم در سطح بیت افزایش خواهد یافت.

به عنوان مثال در حالت A3W3 برای عملیات ضرب بین وزن‌ها و فعال‌سازها بدون در نظر گرفتن تنگی (همانند شبکه دو حالتی) به ۹ سیکل کلاک نیاز است، این در حالی است که اگر عملیات صفر در سطح کلمه را نادیده بگیریم، با توجه به نمودار ۵ در ۶۰ درصد تنگی میزان تسریع بالقوه در سطح کلمه ۸ برابر و در سطح بیت ۳۲ برابر بهتر خواهد شد.

تسریع عملیات



نمودار ۵. بررسی میزان تسریع بالقوه در سطح کلمه و بیت برای درصدهای تنگی مختلف

جدول ۲ نتایج توان مصرفی و مساحت برای واحدهای محاسباتی ارائه شده در شکل ۶ را نشان می‌دهد. این نتایج با استفاده از سنتز مدار مورد نظر به کمک ابزار سنتز تجاری و با کتابخانه تکنولوژی ۴۵ نانومتر FreePDK بدست آمده است. توان مصرفی و مساحت مدار محاسباتی ضرب سه حالتی به نسبت ضرب دو حالتی به ترتیب ۲,۲ و ۳,۱۵ برابر بیشتر است.

جدول ۲. نتایج سخت افزاری واحدهای محاسباتی

نسبت سه حالتی به دو حالتی	سه حالتی	دو حالتی	
2.2	0.40	0.18	توان مصرفی (mW)
3.15	448	142	مساحت (um ²)

با توجه به معماری ارائه شده برای عملیات ضرب بین دو حالتی و سه حالتی، با اینکه مدار ارائه شده در شکل ۶ (ت) تقریباً دو برابر ضرب دو حالتی در دو حالتی سربار محاسباتی دارد، اما با توجه به تسریع‌های بالقوه بیان شده در بالا که همگی از دو برابر خیلی بیشتر هستند، می‌توانیم نشان دهیم که میزان تسریع عملیات ضرب با حذف کردن محاسبات ضرب در صفر می‌تواند بسیار زیاد باشد.

۶- نتیجه‌گیری

چندی‌سازی و هرس کردن دو روش محبوب جهت فشرده‌سازی شبکه‌های عصبی عمیق برای قابل پیاده‌سازی شدن آنها بر روی دستگاه نهفته با محدودیت‌های حافظه و محاسبات هستند. چندی‌سازی غیریکنواخت دو حالتی از محاسبات بیتی بهره می‌برد و افت صحت شبکه‌های دو حالتی پایه را به حداقل می‌رساند. ما در این مقاله چندی‌سازی غیریکنواخت سه حالتی را ارائه دادیم، که صحت آن نزدیک به اندازه صحت شبکه ممیز شناور و بالاتر از صحت شبکه غیریکنواخت دو حالتی است و بر خلاف شبکه غیریکنواخت دو حالتی که وزن‌ها تنگ نیستند، در این نوع از چندی‌سازی وزن‌ها قابلیت تنگ شدن دارند. البته روشن است که هزینه‌ی حافظه و محاسبات غیریکنواخت سه حالتی برای هر محاسبات در حالتی که هیچ بهینه‌سازی انجام نشود، تقریباً دو برابر غیریکنواخت دو حالتی است. اما غیریکنواخت سه حالتی شانس هرس شدن دارد. لذا وزن‌های شبکه چندی شده با نمایش اعداد غیریکنواخت سه حالتی را با درصدهای مختلف هرس کردیم و بررسی کردیم که در حالتی که وزن‌ها تا ۷۵ درصد هرس شوند، صحت شبکه نسبت به شبکه سه حالتی بدون هرس فقط ۰,۸ درصد کاهش می‌یابد. علاوه بر این در شبکه غیریکنواخت سه حالتی با استفاده از فشرده‌سازی وزن‌ها، نه تنها سربار حافظه نسبت به غیریکنواخت دو حالتی نداشتیم، بلکه توانستیم میزان حافظه مصرفی برای ذخیره‌سازی وزن‌های سه

- حالتی را نسبت به شبکه دو حالتی کمتر کنیم. از طرفی برای بهره بردن از تنگی جهت کاهش محاسبات، برای عملیات ضرب وزن سه حالتی در فعال‌ساز دو حالتی یک معماری ارائه دادیم که پیچیدگی محاسبات ضرب، دو برابر شبکه دو حالتی بود، اما به دلیل اینکه شبکه غیریکنواخت سه حالتی ما تنگ است با در نظر گرفتن امکان حذف عملیات ضرب غیر موثر، میزان تسریع بالقوه عملیات ضرب در سطح کلمه و بیت به ترتیب ۱۵ و ۴۵ برابر نسبت به شبکه غیریکنواخت دو حالتی بهتر است. و این میزان تسریع به اندازه‌ای بالا است که به سربار تقریباً دو برابری محاسبات غالب می‌شود.
- مراجع**
- [15] Andri, R., et al. YodaNN: An ultra-low power convolutional neural network accelerator based on binary weights. in 2016 IEEE Computer Society Annual Symposium on VLSI (ISVLSI). 2016. IEEE.
- [16] Jacob, B., et al. Quantization and training of neural networks for efficient integer-arithmetic-only inference. in Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.
- [17] Gysel, P., et al., *Ristretto: A framework for empirical study of resource-efficient inference in convolutional neural networks*. IEEE transactions on neural networks and learning systems, 2018. **29**(11): p. 5784-5789.
- [18] Sharma, H., et al. Bit fusion: Bit-level dynamically composable architecture for accelerating deep neural network. in 2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA). 2018. IEEE.
- [19] Hubara, I., et al., Quantized neural networks: Training neural networks with low precision weights and activations. The Journal of Machine Learning Research, 2017. **18**(1): p. 6869-6898.
- [20] Li, Y. and F. Ren. BNN Pruning: Pruning binary neural network guided by weight flipping frequency. in 2020 21st International Symposium on Quality Electronic Design (ISQED). 2020. IEEE.
- [21] Jin, C., H. Sun, and S. Kimura. Sparse ternary connect: Convolutional neural networks using ternarized weights with enhanced sparsity. in 2018 23rd Asia and South Pacific Design Automation Conference (ASP-DAC). 2018. IEEE.
- [22] Zhu, C., et al., *Trained ternary quantization*. arXiv preprint arXiv:1612.01064, 2016.
- [23] Chin, T.-W., et al. One weight bitwidth to rule them all. in European Conference on Computer Vision. 2020. Springer.
- [24] Ghasemzadeh, M., M. Samragh, and F. Koushanfar. ReBNet: Residual binarized neural network. in 2018 IEEE 26th annual international symposium on field-programmable custom computing machines (FCCM). 2018. IEEE.
- [25] Zhao, Y., et al., *Focused quantization for sparse cnns*. Advances in Neural Information Processing Systems, 2019. **32**.
- [26] Long, X., et al., Learning sparse convolutional neural network via quantization with low rank regularization. IEEE Access, 2019. **7**: p. 51866-51876.
- [27] Long, X., et al. Low Bit Neural Networks with Channel Sparsity and Sharing. in 2022 7th International Conference on Image, Vision and Computing (ICIVC). 2022. IEEE.
- [28] Gadosey, P.K., Y. Li, and P.T. Yamak. On pruned, quantized and compact CNN architectures for vision applications: an empirical study. in Proceedings of the International Conference on Artificial Intelligence, Information Processing and Cloud Computing. 2019.
- [29] Albericio, J., et al., *Cnvlutin: Ineffectual-neuron-free deep neural network computing*. ACM SIGARCH Computer Architecture News, 2016. **44**(3): p. 1-13.
- [30] Han, S., et al., *EIE: Efficient inference engine on compressed deep neural network*. ACM SIGARCH Computer Architecture News, 2016. **44**(3): p. 243-254.
- [31] Parashar, A., et al., *SCNN: An accelerator for compressed-sparse convolutional neural networks*. ACM SIGARCH computer architecture news, 2017. **45**(2): p. 27-40.
- [32] Zhang, S., et al. Cambricon-X: An accelerator for sparse neural networks. in 2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO). 2016. IEEE.
- [33] Delmas, A., et al., DPRed: Making typical activation and weight values matter in deep learning computing. arXiv preprint arXiv:1804.06732, 2018.
- [1] LeCun, Y., Y. Bengio, and G. Hinton, *Deep learning*. nature, 2015. **521**(7553): p. 436-444.
- [2] Zhang, D., et al. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. in Proceedings of the European conference on computer vision (ECCV). 2018.
- [3] Yang, L., Z. He, and D. Fan. Harmonious coexistence of structured weight pruning and ternarization for deep neural networks. in Proceedings of the AAAI Conference on Artificial Intelligence. 2020.
- [4] Burrello, A., et al., *Dory: Automatic end-to-end deployment of real-world dnns on low-cost iot mcus*. IEEE Transactions on Computers, 2021. **70**(8): p. 1253-1268.
- [5] de Prado, M., et al., *Automated Design Space Exploration for Optimized Deployment of DNN on Arm Cortex-A CPUs*. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2020. **40**(11): p. 2293-2305.
- [6] Howard, A.G., et al., Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861, 2017.
- [7] Iandola, F.N., et al., SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size. arXiv preprint arXiv:1602.07360, 2016.
- [8] Courbariaux, M., Y. Bengio, and J.-P. David, *Binaryconnect: Training deep neural networks with binary weights during propagations*. Advances in neural information processing systems, 2015. **28**.
- [9] Cai, Z., et al. Deep learning with low precision by half-wave gaussian quantization. in Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.
- [10] Li, F., B. Zhang, and B. Liu, *Ternary weight networks*. arXiv preprint arXiv:1605.04711, 2016.
- [11] He, Y., X. Zhang, and J. Sun. Channel pruning for accelerating very deep neural networks. in Proceedings of the IEEE international conference on computer vision. 2017.
- [12] Han, S., et al., *Learning both weights and connections for efficient neural network*. Advances in neural information processing systems, 2015. **28**.
- [13] Luo, J.-H., J. Wu, and W. Lin. Thinet: A filter level pruning method for deep neural network compression. in Proceedings of the IEEE international conference on computer vision. 2017.
- [14] Maji, P., et al. Efficient winograd or cook-toom convolution kernel implementation on widely used mobile cpus. in 2019 2nd Workshop on Energy Efficient Machine Learning and Cognitive Computing for Embedded Applications (EMC2). 2019. IEEE.

- [43] Elsken, T., J.H. Metzen, and F. Hutter, *Efficient multi-objective neural architecture search via lamarckian evolution*. arXiv preprint arXiv:1804.09081, 2018.
- [44] Rastegari, M., et al. Xnor-net: Imagenet classification using binary convolutional neural networks. in European conference on computer vision. 2016. Springer.
- [45] Zhang, J., F. Franchetti, and T.M. Low. High performance zero-memory overhead direct convolutions. in International Conference on Machine Learning. 2018. PMLR.
- [46] Lavin, A. and S. Gray. Fast algorithms for convolutional neural networks. in Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [47] De Prado, M., N. Pazos, and L. Benini. Learning to infer: RL-based search for DNN primitive selection on Heterogeneous Embedded Systems. in 2019 Design, Automation & Test in Europe Conference & Exhibition (DATE). 2019. IEEE.
- [48] Rovder, S., J. Cano, and M. O'Boyle, Optimising convolutional neural networks inference on low-powered GPUs. 2019.
- [49] Lin, X., C. Zhao, and W. Pan, *Towards accurate binary convolutional neural network*. Advances in neural information processing systems, 2017. **30**.
- [50] Gholami, A., et al., A survey of quantization methods for efficient neural network inference. arXiv preprint arXiv:2103.13630, 2021.
- [51] Hawks, B., et al., Ps and qs: Quantization-aware pruning for efficient low latency neural network inference. Frontiers in Artificial Intelligence, 2021. **4**: p. 676564.
- [52] Hubara, I., et al., *Binarized neural networks*. Advances in neural information processing systems, 2016. **29**.
- [34] Aimar, A., et al., NullHop: A flexible convolutional neural network accelerator based on sparse representations of feature maps. IEEE transactions on neural networks and learning systems, 2018. **30**(3): p. 644-656.
- [35] Li, J., et al., *SqueezeFlow: A sparse CNN accelerator exploiting concise convolution rules*. IEEE Transactions on Computers, 2019. **68**(11): p. 1663-1677.
- [36] Albericio, J., et al. Bit-pragmatic deep neural network computing. in Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture. 2017.
- [37] Sharify, S., et al. Laconic deep learning inference acceleration. in 2019 ACM/IEEE 46th Annual International Symposium on Computer Architecture (ISCA). 2019. IEEE.
- [38] 38. Chen, Y.-H., et al., Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. IEEE journal of solid-state circuits, 2016. **52**(1): p. 127-138.
- [39] Sandler, M., et al. Mobilenetv2: Inverted residuals and linear bottlenecks. in Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.
- [40] Zhang, X., et al. Shufflenet: An extremely efficient convolutional neural network for mobile devices. in Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.
- [41] Tan, M., et al. Mnasnet: Platform-aware neural architecture search for mobile. in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.
- [42] Wu, B., et al. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.