

بهینه‌سازی برون‌سپاری محاسباتی شبکه‌های عصبی عمیق برای تشخیص فعالیت انسانی

پریسا سعادت

دانشجوی کارشناسی ارشد گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه بجنورد، بجنورد، ایران
پست الکترونیکی: parisa.saadati@stu.ub.ac.ir

حمید فدیشه‌ای*

استادیار گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه بجنورد، بجنورد، ایران
پست الکترونیکی: fadishei@ub.ac.ir

چکیده

سمت ابر با نگاشت معکوس به صورت تخمینی بازیابی می‌شوند. در روش فشرده‌سازی، نمونه‌های داده قبل از ارسال، با یک الگوریتم فشرده‌سازی سبک دلتا به یکی از دو روش با اتلاف یا بدون اتلاف فشرده می‌شوند. نتایج عملی نشان می‌دهد اگر چه روش‌های کاهش تعداد نمونه و کاهش دقت اعشار سبب کاهش حجم داده ارسالی بدون تأثیر قابل توجه بر روی دقت تشخیص می‌شوند، روش کاهش دقت اعشار به دلیل میزان کاهش بیشتر حجم داده نسبت به روش کاهش تعداد نمونه‌ها برتری دارد. ضمناً روش فشرده‌سازی دلتا به اندازه دو روش دیگر موثر نمی‌باشد.

واژه‌های کلیدی: برون‌سپاری محاسباتی، شبکه‌های عصبی عمیق، تشخیص فعالیت انسانی

مقدمه

در سال‌های اخیر شاهد افزایش چشمگیر تحقیقات

با توجه به محدودیت‌های منابع محاسباتی موجود در دستگاه‌های تلفن همراه، این دستگاه‌ها در اجرای برخی وظایف روزمره با چالش مواجه هستند. یک راه‌حل برای این مشکل، برون‌سپاری است که در آن، دستگاه پردازش خود را برای اجرا روی ابر محاسباتی ارسال می‌کند. مقاله حاضر با تمرکز بر روی کاربرد تشخیص فعالیت انسانی، روش‌هایی برای کاهش حجم داده‌های ارسالی وظایف به ابر، با تکیه بر نقاط مطلوب در معامله بین دقت استنتاج یادگیری و هزینه ارتباطات برون‌سپاری ارائه می‌کند. سه روش پیشنهادی کاهش تعداد نمونه‌های داده، کاهش دقت اعشار نمونه‌های داده و فشرده‌سازی نمونه‌های داده ارائه شده است که در روش اول نمونه‌های داده قبل از ارسال به صورت یکی در میان یا بیشتر حذف شده و در سمت ابر با تخمین درونیابی مجدداً بازیابی می‌شوند. در روش کاهش دقت اعشار، نمونه‌های داده قبل از ارسال با یک نگاشت به عدد صحیح با تعداد بیت محدود تبدیل و در

* نویسنده مسؤل

و توسعه شبکه‌های تلفن همراه بوده‌ایم. با پیشرفت در پایانه‌های تلفن همراه و از آنجا که تلفن‌های هوشمند محبوبیت زیادی پیدا کردند، برنامه‌های جدید تلفن همراه مانند تشخیص چهره، پردازش تصویر، بازی‌های تعاملی و واقعیت افزوده توجه زیادی را به خود جلب کرده و انتظار از دستگاه‌های تلفن همراه برای اجرای برنامه‌های سنگین‌تر در حال افزایش است [۱]. اکنون کاربران تلفن همراه استفاده‌های روزمره فراوانی مانند جستجو در میان آهنگ‌ها، انجام بازی‌های ویدیویی، ضبط، ویرایش و بارگذاری فیلم، تجزیه و تحلیل مجموعه عکس‌های خود، فهرست‌بندی محتوا و مدیریت امور مالی دارند [۲]. با وجود نقش پررنگ این دستگاه‌ها در زندگی افراد، اجرای برنامه‌های پیچیده در دستگاه‌های تلفن همراه به دلیل محدودیت منابع آن‌ها از قبیل ظرفیت حافظه، سرعت پردازش گرافیکی و قدرت باتری چالش برانگیز است [۳]. دستگاه‌های تلفن همراه دارای قدرت محاسباتی نسبتاً ضعیف، محدودیت باتری و منابع سخت افزاری می‌باشند. همچنین برنامه‌های کاربردی دستگاه‌های تلفن همراه معمولاً نیاز به محاسبات فشرده و مصرف انرژی بالایی دارند. با توجه به محدودیت‌های منابع محاسباتی موجود در دستگاه‌های تلفن همراه، ممکن است این دستگاه‌ها نتوانند برنامه کاربردی را به‌طور موثر اجرا کنند [۴]. امروزه به‌خاطر افزونی کاربردهایی با حجم بالای پردازش، نیاز به محیط و منابع قدرتمندی که بتواند این محاسبات سنگین را به عهده بگیرد، می‌باشد. یک راه حل برای این مشکل به نام برون‌سپاری وجود دارد که در آن دستگاه‌های کم قدرت مانند تلفن همراه، اشیاء اینترنتی و... پردازش خود را بر روی ابر محاسباتی ارسال کرده و به کارسازهای ابری می‌سپارند [۵]. برون‌سپاری می‌تواند موجب صرفه‌جویی در مصرف انرژی و بهبود عملکرد شود و همین‌طور می‌تواند قابلیت محاسبات سیستم‌های تلفن همراه را تقویت کند [۲]. بسیاری از کاربردهای محاسباتی ذکر شده از یادگیری

ماشین سود می‌برند که بنا به تعریف، قابلیت بهبود کارکرد یک برنامه بدون دخالت برنامه‌نویس، بلکه از طریق تجربه می‌باشد [۶]. یادگیری عمیق نوع خاص و پرکاربرد یادگیری ماشین می‌باشد که اخیراً مورد توجه محققان قرار گرفته است و قابلیت خود را در تشخیص با دقت بالا در کاربردهایی مثل پردازش تصویر، تشخیص فعالیت انسانی، پردازش گفتار و... نشان داده است [۷]. برای یادگیری عمیق معمولاً از شبکه‌های عصبی عمیق استفاده می‌شود که یک نمودار جریان داده هستند و از تعدادی لایه تشکیل شده‌اند، هر کدام از لایه‌ها عملکرد خود را بر روی داده‌های ورودی انجام می‌دهند و داده‌های خروجی را به لایه بعدی منتقل می‌کنند [۸]. شبکه‌های با حافظه دور و نزدیک (Long Short Term Memory) نوع خاصی از شبکه‌های عصبی عمیق هستند که توانایی یادگیری وابستگی‌های کوتاه مدت و بلند مدت را دارند. این شبکه‌ها برای اولین بار در سال ۱۹۹۷ توسط Hochreiter و Schmidhuber معرفی شدند [۹]. این شبکه‌ها به دلیل ساختار خاص و دارا بودن حافظه، در مورد داده‌هایی مثل مسئله تحقیق حاضر که وابستگی زمانی دارند معمولاً دقت بهتری نسبت به سایر انواع شبکه‌های عصبی نشان می‌دهند [۱۰].

عملیات محاسبات شبکه‌های عصبی عمیق شامل دو مرحله آموزش و استنتاج می‌باشد که مسئله برون‌سپاری می‌تواند در هر دو مرحله مذکور مطرح شود. در مرحله آموزش با استفاده از داده‌های ورودی از قبل برچسب خورده، پارامترهای شبکه‌های عصبی عمیق (نظیر وزن یال‌ها) تعیین می‌شود، تا شبکه عصبی عمیق بتواند در مرحله کاربرد بر روی داده‌هایی که تا به حال مشاهده نشده است، استنتاج انجام دهد. پردازش هر لایه می‌تواند به‌عنوان یک عملیات بردار در نظر گرفته شود که پارامترهای آن به‌طور تکراری در حالی که شبکه‌های عصبی عمیق با داده‌های برچسب خورده آموزش می‌یابند، به‌روز می‌شوند. با توجه به این‌که کاربردهای عملی گوشی‌های هوشمند در مرحله استنتاج نقش پررنگ‌تری دارد، اکثر

محققین بر روی برون‌سپاری محاسبات مرحله استنتاج تمرکز دارند. از طرف دیگر، مرحله آموزش به دلیل نیاز گسترده به منابع محاسباتی به‌طور معمول در کارسازهای قدرتمند انجام می‌شود [۳]. دلیل دیگر کم‌رنگ بودن برون‌سپاری مرحله آموزش این است که پس از انجام این مرحله، پارامترهای هر لایه ثابت هستند، بنابراین شبکه‌های عصبی عمیق مادامی که داده‌های آموزش تغییر نیافته‌اند، همان پارامترها را برای استنتاج روی داده‌های ورودی به کار می‌برند. برخی از محققان بر روی برون‌سپاری همه محاسبات شبکه‌های عصبی عمیق یا بخش‌هایی از آن تمرکز کرده و راه‌حلی برای این کار ارائه داده‌اند، تا بر موانعی از قبیل محدودیت‌های استفاده از باتری در سمت تلفن همراه و منابع محدود محاسباتی آن فایز آیند. برون‌سپاری شبکه‌های عصبی عمیق معمولاً با معامله همراه است و در ازای صرفه‌جویی در زمان اجرا و انرژی مصرفی، هزینه‌ای در کاهش دقت استنتاج پرداخت می‌گردد [۱۱]. این معامله تمرکز مقاله حاضر است و در این تحقیق بر روی اولین نقطه ورودی شبکه‌های عصبی عمیق یعنی داده‌های ارسال شده به ابر در مرحله استنتاج تمرکز شده است تا حجم آن‌ها به طرق مختلفی کاهش یابد.

کاربردی که در این تحقیق به آن پرداخته شده است، تشخیص فعالیت‌های انسانی می‌باشد. تشخیص فعالیت‌های انسانی برای ارائه خدمات در دنیای اینترنت اشیا ضروری است. با توجه به همه‌گیر بودن، قابلیت سنجش و قدرت پردازش، تلفن‌های هوشمند مدرن به دستگاه‌های مورد توجه برای تشخیص فعالیت انسانی تبدیل شده‌اند. با این حال، ظرفیت محدود باتری و منابع تلفن هوشمند مانع بهره‌برداری از چنین قدرت سنجش و پردازش می‌شود [۱۲]. تشخیص فعالیت‌های انسانی امکان می‌دهد بسیاری از فعالیت‌های بدنی را که یک کاربر تلفن هوشمند انجام می‌دهد (نظیر پیاده‌روی، دویدن و...) براساس حرکات کاربر تشخیص داده شود. در این تحقیق از داده‌های جمع‌آوری شده از حسگرهای شتاب‌سنج، مجموعه داده

WISDM استفاده شده است [۱۳]. این مجموعه داده حاوی بیش از دو میلیون نمونه شامل پنج فعالیت پیاده‌روی، دویدن، فعالیت پله‌ها، نشستن و ایستادن می‌باشد. این تحقیق، تلاشی برای بهینه‌سازی ارتباطات در برون‌سپاری مسئله تشخیص فعالیت انسانی با گوشی هوشمند از طریق بررسی معامله بین حجم داده ارسالی و دقت تشخیص است. نوآوری‌های این مقاله به شرح زیر می‌باشد:

- مطالعه روش‌های مختلف کاهش حجم داده‌های ارسالی به ابر برای برون‌سپاری استنتاج شبکه‌های عصبی عمیق.
- پیشنهاد سه روش کاهش تعداد نمونه‌ها، کاهش دقت اعشار و فشرده‌سازی نمونه‌های داده برای بهینه‌سازی برون‌سپاری مذکور.
- بررسی و تحلیل روش‌های مختلف کاهش حجم داده‌های ارسالی برای یافتن نقاط مطلوب در معامله بین دقت و هزینه ارتباطات.

۲- کارهای گذشته

از آنجا که در سال‌های اخیر برنامه‌های تلفن همراه به‌طور فزاینده‌ای مورد توجه همه جانبه قرار گرفته‌اند و عملکرد قدرتمندی را در دستگاه‌های تلفن همراه ارائه کرده‌اند و با توجه به این‌که دستگاه‌های تلفن همراه با چالش‌های محاسبات محدود و کم‌قدرت روبرو هستند [۱۴]، توسعه‌دهندگان معمولاً تمایل دارند فرآیند محاسبه را به ابر منتقل کنند. تحقیقات زیادی برای بهینه‌سازی برون‌سپاری شبکه‌های عصبی عمیق انجام شده است. یک رویکرد کلی برای برون‌سپاری شبکه‌های عصبی عمیق، بخش‌بندی می‌باشد که در آن گراف شبکه‌های عصبی عمیق به چند بخش تقسیم می‌شود و در هنگام اجرا و با در نظر گرفتن وابستگی‌ها، فرآیند محاسبه بر روی ابر و یا به صورت محلی انجام می‌شود. در پژوهش Huang و همکارانش [۴]، برای استفاده بهینه از منابع موجود، بار محاسباتی بین دستگاه، لبه و ابر توزیع شده است که این امر می‌تواند تأخیر ارتباطی را به‌طور قابل توجهی کاهش

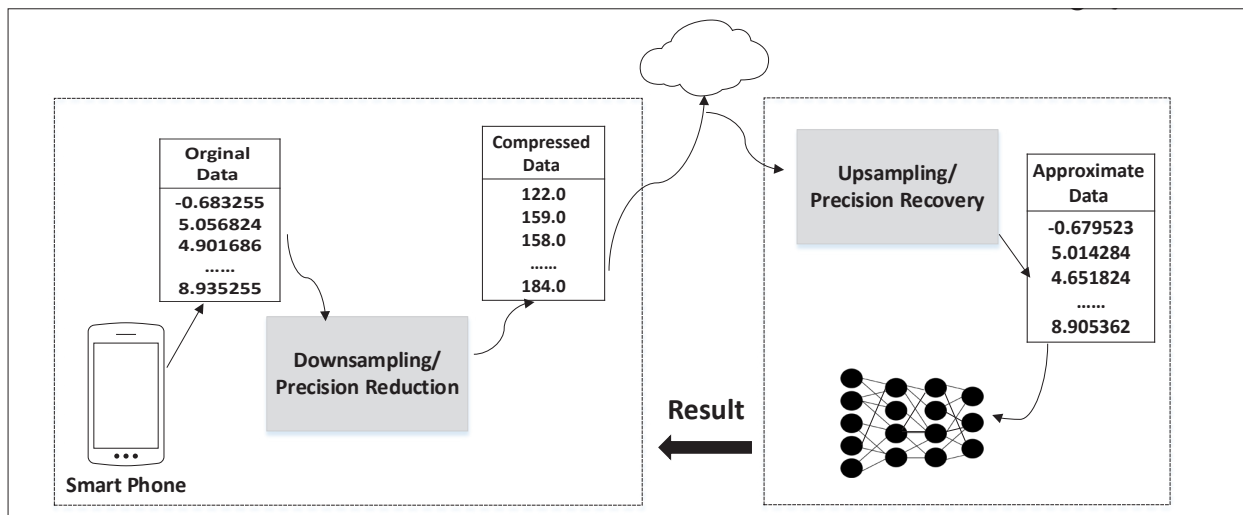
دهد. بخش‌بندی می‌تواند با در نظر گرفتن خصوصیات پویای بارکاری و شرایط باری لبه در هنگام برون‌سپاری و زمان اجرا همراه باشد که علاوه بر بهبود زمان تأخیر، بهره‌وری انرژی در سیستم‌هایی با منابع محدود را افزایش دهد [۳]. برای غلبه بر محدودیت‌های منابع در پژوهش [۱۵] روشی برای بخش‌بندی شبکه‌های عصبی عمیق با هدف کاهش زمان محاسبه و افزایش استفاده از منابع در شبکه ارائه شده است.

در پژوهش [۱۶] که بر یادگیری عمیق برای اینترنت اشیا در محیط محاسباتی لبه تمرکز دارد، به دلیل توان پردازش محدود گره‌های لبه، یک روش برون‌سپاری محاسباتی جدید برای بهینه‌سازی عملکرد برنامه یادگیری عمیق ارائه شده که محافظت از حریم خصوصی کاربر در برون‌سپاری داده‌ها را نیز به همراه دارد. روش IONN [۱۷] مدل شبکه عصبی عمیق کارخواه را به چند قسمت تقسیم و سپس یک به یک در کارساز لبه بارگذاری می‌کند. کارساز با رسیدن هر قسمت از شبکه عصبی عمیق، مدل را به صورت تدریجی ایجاد می‌کند و به کارخواه اجازه می‌دهد تا حتی قبل از بارگذاری کل مدل، برون‌سپاری جزئی از اجرای شبکه عصبی عمیق را آغاز کند، این امر مصرف انرژی و بهره‌وری کوئری در برون‌سپاری محاسبات شبکه‌های عصبی عمیق از دستگاه تلفن همراه به کارساز را بهبود می‌بخشد. همچنین بخش‌بندی بین ابر، لبه و دستگاه‌های انتهایی با استفاده از مقیاس‌پذیری شبکه‌های عصبی عمیق به لحاظ اندازه و مقطع جغرافیایی آن‌ها، هزینه‌های ارتباطی را کاهش می‌دهد و سبب بهبود دقت می‌شود [۱۸].

روش کلی دیگر که در مقابل رویکرد بخش‌بندی قرار می‌گیرد، انتقال کل ساختار شبکه عصبی به ابر است، به‌عنوان مثال پژوهش [۱۹] از رویکرد انتقال کل ساختار شبکه عصبی از کارخواه به یک کارساز لبه استفاده کرده است. تحقیق مذکور در پیاده‌سازی روش پیشنهادی خود از بستر اجرای برنامه‌های وب بهره‌برداری می‌کند تا شبکه عصبی عمیق را در قالب وضعیت اجرای یک برنامه وب از

کارخواه به کارساز لبه انتقال دهد. بعد از انجام محاسبات در لبه، حالت اجرای جدید از کارساز لبه به کارخواه مهاجرت می‌کند تا کارخواه بتواند اجرای برنامه را ادامه دهد. پژوهش [۱۱] بین تفکیک‌پذیری و نرخ قاب معامله می‌کند تا بتواند بین تجزیه و تحلیل وضوح بالا با تأخیر بالا و وضوح پایین با تأخیر کم تصمیم بگیرد و کارکرد برون‌سپاری در این است که حداقل نرخ قاب و دقت را با توجه به شرایط شبکه بهبود بخشد. رویکرد فشرده‌سازی یکی دیگر از روش‌های برون‌سپاری شبکه‌های عصبی عمیق است. در پژوهش [۲۰] پدram و همکارانش با ارائه یک معماری یادگیری عمیق برای خدمات محاسباتی ابر و دستگاه تلفن هوشمند توانسته‌اند اندازه ویژگی‌های مورد نیاز برای ارسال به ابر را کاهش دهند و تأخیر و مصرف انرژی شبکه‌های عصبی عمیق در برنامه‌های تلفن همراه را بهبود ببخشند. همینطور در پژوهش [۲۱] با طراحی یک چارچوب که نسبت به تغییرات شبکه مقاوم است، اندازه داده‌های منتقل شده در شبکه را کاهش داده و سبب بهبود ترافیک شبکه، توان عملیاتی و دقت شده‌اند. در تحقیق دیگری مدلی تحت عنوان JointDNN ارائه شده است [۲۲] که در آن علاوه بر این‌که با پردازش برخی لایه‌ها بر روی دستگاه تلفن همراه و برخی بر روی کارساز ابر برون‌سپاری شبکه‌های عصبی عمیق را بهینه کرده‌اند، توانسته‌اند با استفاده از روش فشرده‌سازی با کاهش ابعاد ویژگی‌ها هزینه‌های ارتباطی و زمان لازم برای ارسال داده به ابر را کاهش دهند.

با تأمل در کارهای گذشته می‌توان دریافت که بخش عظیمی از تحقیقات بر روی بخش‌بندی معطوف هستند که این رویکرد با چالش‌هایی مواجه است. معمولاً مدل شبکه‌های عصبی عمیق آن‌قدر بزرگ است که حتی در نظر گرفتن یک قسمت از مدل روی تلفن همراه سربار زیادی را به وجود می‌آورد. همچنین با توجه به این‌که معمولاً مرحله آموزش روی ابر انجام می‌شود، بهتر است مدل شبکه‌های عصبی عمیق نیز روی ابر باشد. مدل‌های شبکه



شکل ۱: چارچوب روش پیشنهادی در این مقاله

نظیر مورد مطالعاتی تشخیص فعالیت انسانی که در ادامه توضیح داده شده است)، می‌تواند منجر به کاهش انرژی نمونه‌برداری داده‌های حسگر شود. دلیل این امر این است که انرژی مصرفی لایه حسگر در گوشی‌های هوشمند متناسب با نرخ نمونه‌برداری افزایش می‌یابد. در نتیجه هر چه تعداد نمونه‌هایی که در واحد زمان دریافت می‌شود کمتر باشد انرژی کمتری در سمت تلفن همراه مصرف می‌شود. البته کاهش نمونه‌ها، دقت استنتاج شبکه‌های عصبی عمیق را پایین خواهد آورد. هدف از این تحقیق مطالعه معامله بین دقت و هزینه‌های ارتباطی و پیدا کردن نقطه بهینه در این معامله است و باید دید آیا با بازگرداندن تقریبی نمونه‌های حذف شده امکان ترمیم دقت از دست رفته تا حدودی وجود دارد یا خیر. برای این منظور، نگارندگان داده‌ها را با حالت‌های مختلف حذف نمونه‌های میانی یعنی حذف نمونه داده به صورت حفظ یکی در میان (نگهداری یک نمونه داده از هر دو نمونه)، حفظ دو در میان (نگهداری یکی از هر سه نمونه داده) و سه در میان کاهش داده و سپس آن‌ها را با سیاست درونیابی به دو روش خطی و مکعبی بازیابی نموده‌اند.

درونیابی فرآیند شناخته شده‌ای برای تخمین مقادیر موجود بین نقاط مشخص داده است. برای پیدا کردن مقادیر در نقاط بین نقاط داده شده می‌توان با فرض خطی

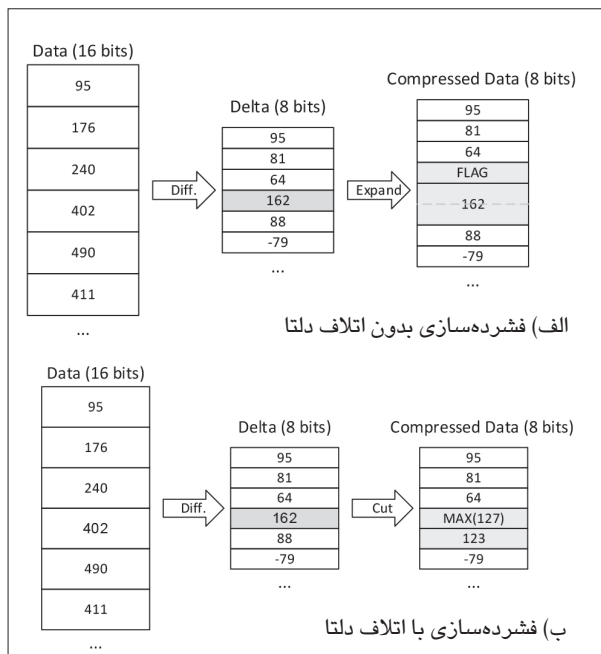
عصبی عمیق در عمل معمولاً پویا هستند و اگر پارامترها و وزن‌های آن‌ها بخواهند مرتباً به روز شوند، فرستادن بخشی یا کل مدل بر روی تلفن همراه هزینه‌بر خواهد بود. ضمناً شبکه‌های عصبی عمیق آن‌قدر بزرگ هستند که ارسال کل آن به صورت برنامه وب به ابر با هزینه زیادی مواجه خواهد بود.

۳- روش پیشنهادی

هدف از این پژوهش، بهینه‌سازی برون‌سپاری محاسبات شبکه‌های عصبی عمیق از طریق کاهش حجم داده‌های ارسالی به ابر می‌باشد. چارچوب کلی روش پیشنهادی برای رسیدن به این هدف در شکل (۱) نشان داده شده است و شامل سه رویکرد مختلف است که در ادامه توضیح داده شده است.

کاهش تعداد نمونه‌های داده

اولین رویکرد پیشنهادی این مقاله برای بهینه‌سازی برون‌سپاری محاسبات شبکه‌های عصبی عمیق، کاهش تعداد نمونه‌های ارسالی به ابر می‌باشد. کاهش تعداد نمونه داده‌های ورودی همیشه حجم داده‌های ارسالی به ابر را کاهش می‌دهد و منجر به کاهش هزینه ارتباطات (زمان، قیمت و انرژی مصرفی) می‌شود. ضمناً در بعضی از موارد که داده‌های ورودی از حسگرها دریافت می‌شوند



شکل ۲: روش پیشنهادی در فشرده‌سازی نمونه داده‌ها

فشرده‌سازی نمونه‌های داده

رویکرد دیگر مورد مطالعه در این مقاله برای بهینه‌سازی برون‌سپاری محاسبات شبکه‌های عصبی عمیق، فشرده‌سازی نمونه‌های داده است. باید توجه داشت که اگر چه الگوریتم‌های فشرده‌سازی مختلفی با ضریب فشرده‌سازی بالا وجود دارند، باید روشی با سربار محاسباتی خیلی پایین انتخاب شود. یکی از روش‌های سبک برای فشرده‌سازی جریان داده‌ها، روش دلتا می‌باشد. این روش با فرض این‌که دامنه تغییر هر نمونه داده نسبت به مقدار قبلی به حد کافی کوچک است، به جای ارسال نمونه‌ها، تفاضل آن‌ها را با تعداد بیت کمتر ارسال می‌کند. بدیهی است که در صورت برقرار نبودن این فرض، الگوریتم مذکور قادر به فشرده‌سازی موثر داده‌ها نمی‌باشد. در این تحقیق از دو روش فشرده‌سازی با اتلاف و بدون اتلاف دلتا همانند شکل (۲) استفاده شده است.

در روش فشرده‌سازی بدون اتلاف، فشرده‌سازی و غیر فشرده‌سازی داده‌ها سبب هیچ گونه از دست رفتگی در آن‌ها نمی‌شود. در این روش همان گونه که در شکل (۲) الف مشاهده می‌شود، با فرض این‌که نمونه داده‌ها n

بودن تغییرات، طبق رابطه (۱) میان‌گیری حسابی انجام داد:

$$y = y_a + (y_b - y_a) \frac{x - x_a}{x_b - x_a} \quad (1)$$

که در این رابطه x_a و x_b دو نقطه موجود و x نقطه مورد نظر برای یافتن مقدار درونیابی شده است. همچنین می‌توان با در نظر گرفتن تعداد نقاط بیشتر، درونیابی با چند جمله‌ای از درجات بالاتر انجام داد. در این پژوهش از دو نوع سیاست درونیابی خطی و مکعبی استفاده شده است [۲۳].

کاهش دقت اعشار نمونه‌های داده

رویکرد پیشنهادی دوم در این مقاله برای بهینه‌سازی برون‌سپاری محاسبات شبکه‌های عصبی عمیق، کاهش دقت اعشار (Precision) نمونه داده‌های ارسالی به ابر می‌باشد. برای بهینه‌سازی برون‌سپاری می‌توان دقت اعشار داده‌ها را کاهش داد تا با کم شدن تعداد بیت‌های لازم برای هر نمونه داده، حجم داده‌های ارسالی به ابر کاهش یافته و در مصرف انرژی صرفه‌جویی شود. بدیهی است که در این روش معامله‌ای مانند روش قبل بین دقت و حجم داده ارسالی وجود دارد. در این روش با استفاده از رابطه نگاشت خطی (۲)، که در آن x متغیر اولیه و n تعداد بیت می‌باشد، داده‌های حسگرها را به جای نمایش ممیز شناور می‌توان با یک عدد صحیح ۸ بیتی، ۱۶ بیتی و... کد کرد. به‌عنوان مثال برای کدگذاری ۸ بیتی، کوچک‌ترین عدد خوانده شده از حسگرها بر روی صفر و بزرگ‌ترین عدد بر روی ۲۵۵ نگاشت می‌شود و بقیه به تناسب در این فاصله قرار می‌گیرند.

$$f(x, n) = \text{round} \left(\frac{(x - \min(x))}{(\max(x) - \min(x))} \times (2^n - 1) \right) \quad (2)$$

پس از دریافت داده‌های کاهش‌یافته در سمت ابر می‌توان با استفاده از رابطه (۳) که معکوس نگاشت قبلی می‌باشد، آن‌ها را به داده‌هایی نزدیک به داده‌های قبلی برگرداند:

$$g(\hat{x}, n) = \min(x) + \frac{(\max(x) - \min(x))}{(2^n - 1)} \times \hat{x} \quad (3)$$

در این رابطه x متغیر اولیه و \hat{x} متغیر نگاشت شده بعد از کاهش دقت و n تعداد بیت مورد نظر برای اختصاص به هر نمونه داده است.

True Label	Jog	3396	1	46	1	3
	Sit	0	570	4	3	0
	Sta	1	2	2237	3	22
	Stan	0	3	1	432	1
	W	4	0	46	0	4205
	Predicted Label	Jog	Sit	Sta	Stan	W

شکل ۳: دقت تشخیص به تفکیک فعالیت بدون حذف داده

مذکور شامل دو لایه کاملاً متصل و دو لایه LSTM با ۶۴ واحد است.

آزمایش‌ها و نتایج

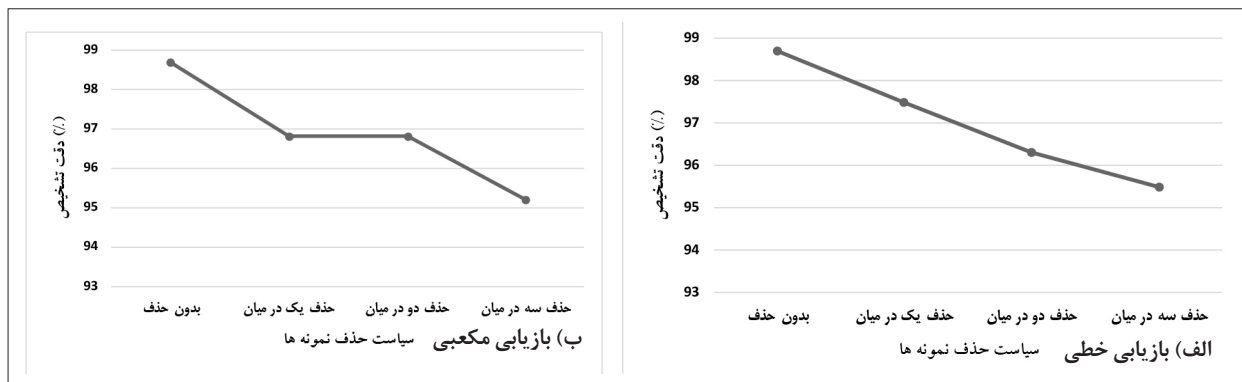
در این بخش تلاش می‌شود با توجه به روش‌های پیشنهادی، نقطه بهینه در معامله بین دقت و هزینه ارتباطی مطالعه شود. برای این منظور روش‌های پیشنهادی این مقاله بر روی مجموعه داده WISDM [۱۳] مورد ارزیابی قرار گرفته است. ابتدا دقت استنتاج شبکه عصبی عمیق برای داده‌های اصلی بدون روش‌های پیشنهادی با اختصاص ۸۰٪ از داده‌ها به آموزش و مابقی آن به آزمایش، آزمایش شد. شکل (۳) جزئیات میزان دقت تشخیص به تفکیک فعالیت را نشان می‌دهد. در این شکل فعالیت‌ها شامل دویدن (Jog)، نشستن (Sit)، بالا و پایین رفتن از پله‌ها (Sta)، ایستادن (Stan) و قدم زدن (W) می‌باشد. دقت کلی تشخیص در این حالت (بدون صرفه‌جویی در حجم داده ارسالی) برابر با ۹۸٪ می‌باشد.

در اولین آزمایش، روش پیشنهادی کاهش حجم داده‌های ارسالی به ابر از طریق کاهش تعداد نمونه‌ها مورد ارزیابی قرار گرفت. بدین منظور، کسری از داده‌ها با حالت‌های مختلف حذف نمونه‌های میانی حذف شده و با سیاست درونیابی به دو روش درونیابی خطی و مکعبی بازیابی شد. هر دو گام یادگیری و استنتاج با داده‌های

بیت (4, 8, 16, 32, n)، دلتا m بیت (2, 4, 8, 16, m) و $m < n$ است (به‌عنوان مثال: $n = 32, m = 16$)، به‌جای ارسال اعداد n بیتی به ابر، تفاضل آن‌ها در قالب m بیت ارسال می‌شود. اگر تفاضل نمونه داده‌ها از m بیت (ظرفیت دلتا) بیشتر باشد، اجباراً باید نمونه داده اصلی ارسال شود و برای تفکیک این حالت، قبل از نمونه داده اصلی یک علامت (داده ذخیره‌شده) قرار داده شود. کوچک‌ترین عدد منفی بازه دلتا به‌عنوان علامت ذخیره استفاده شده است. در نتیجه سررابط این علامت‌ها، با وجود فشرده‌سازی ممکن است حجم نهایی بعد از فشرده‌سازی از حجم اولیه داده بزرگ‌تر شود. این حالت زمانی اتفاق خواهد افتاد که دامنه تغییرات بین نمونه داده‌های پشت سر هم زیاد باشد و به‌طور مکرر نیاز به استفاده از علامت باشد. در مقابل، روش فشرده‌سازی با اتلاف دلتا برآوردی تقریبی از داده‌های اصلی را نگه می‌دارد و بخشی از داده اصلی را به ازای درجه فشرده‌سازی از دست می‌دهد. این روش همان‌طور که شکل (۲) نشان می‌دهد، هنگامی که تفاضل نمونه داده‌ها از m بیت بیشتر باشد (برای حالت خاص $m = 8$ در این شکل) حداکثر مقداری که در ۸ بیت می‌تواند نمایش یابد را ارسال نموده و از نمونه‌های بعدی برای اصلاح آن استفاده می‌کند. در نتیجه در فشرده‌سازی با اتلاف، افزایش حجم مشاهده نمی‌شود ولی به دلیل وجود اتلاف، بهای کاهش حجم با از دست رفتن دقت پرداخت می‌شود.

۴- پیاده‌سازی

تقریباً هر گوشی هوشمند مدرن دارای یک شتاب‌سنج سه محوری است که شتاب را در هر سه بعد فضایی اندازه‌گیری می‌کند. در این مقاله نگارندگان از داده‌های جمع‌آوری شده از حسگرهای شتاب‌سنج استفاده کردند و یک شبکه عصبی LSTM، که با زبان پایتون در کتابخانه Tensorflow [۲۴] اجرا شده است را به‌عنوان مدل مرجع آزمایش‌ها برای شناسایی فعالیت‌های انسانی (HAR) از داده‌های شتاب‌سنج آموزش دادند. مدل شبکه عصبی



شکل ۴: اثر حذف نمونه‌های داده قبل از ارسال به ابر بر روی دقت

همان‌گونه که انتظار می‌رود، کاهش دقت اعشار با کاهش دقت تشخیص همراه است. اما این معامله دقت در ازای کاهش حجم سودمند است و به‌عنوان مثال در حالت ۸ بیتی که به میزان ۸ برابر صرفه‌جویی در حجم ارسالی (در مقایسه با دقت اعشار ۶۴ بیتی ممیز شناور) حاصل می‌شود، تنها در حدود یک درصد کاهش دقت تشخیص اتفاق افتاده است. نکته جالب توجه این است که حتی با کاهش بیشتر دقت اعشار و استفاده از دقت ۴ بیتی که منجر به کاهش ۱۶ برابری حجم داده می‌شود، دقت تشخیص در همین محدوده باقی‌مانده است. البته کاهش دقت اعشار به میزان کمتر از ۴ بیت دقت تشخیص را به حد غیر قابل قبولی پایین می‌آورد.

شکل (۷) جزئیات اثر کاهش دقت اعشار نمونه‌های داده را در دقت تشخیص به تفکیک فعالیت در حالت‌های مختلف نشان می‌دهد. همان‌طور که در شکل مشاهده می‌شود، مشابه روش اول، بیشترین اثر منفی کاهش دقت اعشار بر روی تشخیص فعالیت بالا و پایین رفتن از پله‌ها اتفاق افتاده است. همچنین کاهش بیش از اندازه دقت اعشار سبب کاهش دقت تشخیص اغلب فعالیت‌ها شده است که در آزمایش‌های ۱ بیتی به وضوح قابل مشاهده است.

در رویکرد فشرده‌سازی نمونه‌های داده، به جای ارسال اعداد n بیتی، تفاضل آن‌ها در قالب دلتای m بیتی ارسال می‌شود. البته قبل از این کار، داده‌ها از حالت نمایش اعشاری به حالت صحیح n بیتی تبدیل می‌شوند. اثر فشرده‌سازی به هر دو صورت بدون اتلاف و با

کاهش یافته انجام شد. شکل (۴) اثر حذف نمونه‌های داده قبل از ارسال به ابر را بر روی دقت نشان می‌دهد که قسمت (الف) حذف یک در میان، دو در میان و سه در میان با باز یابی به روش خطی و قسمت (ب) برای باز یابی به روش مکعبی را نشان می‌دهد. همان‌طور که در شکل‌ها مشاهده می‌شود به‌طور کلی رویکرد حذف نمونه‌ها می‌تواند سودمند قلمداد شود. به‌عنوان مثال حذف یک در میان با این‌که باعث کاهش چشمگیر ۵۰ درصدی حجم داده ارسالی می‌شود تنها به میزان یک درصد در دقت پیش‌بینی‌ها اثر منفی دارد. ضمناً تفاوت معناداری بین دو روش باز یابی خطی و مکعبی مشاهده نمی‌شود.

شکل (۵) اثر حذف نمونه‌های داده در دقت تشخیص به تفکیک فعالیت برای باز یابی خطی و مکعبی در حالت‌های مختلف حذف را نشان می‌دهد. همان‌طور که در شکل مشاهده می‌شود عمده اثر منفی دقت تشخیص مربوط به فعالیت بالا و پایین رفتن از پله می‌باشد. بنابراین دقت تشخیص به نوع فعالیت وابسته می‌باشد.

یکی از روش‌های دیگر برای کاهش حجم داده‌های ارسالی به ابر، کاهش دقت اعشار می‌باشد. قبل از ارسال به ابر، داده‌ها با تابع نگاشت خطی به عدد صحیح n بیتی تبدیل می‌شود ($n = 1, 2, 4, 8, 16, 32$) و پس از دریافت داده‌ها در ابر محاسباتی، با رابطه معکوس به داده‌هایی نزدیک به داده‌های اولیه برگردانده می‌شود. شکل (۶) اثر کاهش دقت نمونه‌های داده قبل از ارسال به ابر به روی دقت را برای حالت‌های مختلف نگاشت خطی نشان می‌دهد.



شکل ۶: اثر کاهش دقت اعشار نمونه‌های داده قبل از ارسال به ابر بر روی دقت تشخیص برای حالت‌های مختلف نگاشت خطی

وجود فشرده‌سازی ممکن است حجم نهایی از حجم اولیه بیشتر شود. این مراحل را برای داده‌های ۱۶، ۸، ۴ و ۲ بیت با دلتای ۸، ۴ و ۲ بیت آزمایش شده و حجم داده‌ها مورد بررسی قرار گرفته است. همان‌طور که در شکل (۸) مشاهده می‌شود در اغلب حالت‌ها حجم نهایی بعد از فشرده‌سازی از حجم اولیه داده‌های اصلی بیشتر شده است. در نتیجه، روش فشرده‌سازی بدون اتلاف دلتا به جز در حالت داده ۴ بیتی با ظرفیت دلتا ۲ بیتی، روش مناسبی برای کاهش حجم داده‌های ارسالی به ابر نمی‌باشد.

روش فشرده‌سازی با اتلاف دلتا که برآوردی تقریبی از داده اصلی را نگه می‌دارد، همانند فشرده‌سازی با اتلاف به جای ارسال داده‌های صحیح ۳۲ بیتی، تفاضل آن‌ها در قالب دلتای ۱۶، ۸، ۴ و ۲ بیتی ارسال می‌شود با این تفاوت که در حالت بیشتر بودن ظرفیت دلتا حداکثر مقداری که در ظرفیت دلتا استفاده می‌شود را قرار می‌دهد و از نمونه‌های بعدی برای اصلاح آن استفاده می‌کند. در نتیجه نسبت فشرده‌سازی همیشه برابر با نسبت ایده‌آل اندازه داده به اندازه دلتا است، اما به دلیل از دست رفتن بخشی از داده اصلی، دقت تشخیص کمتر خواهد شد. در مواردی که اختلاف اندازه داده و اندازه دلتا قابل توجه باشد، دقت به میزان فاحشی از دست خواهد رفت. از این رو آزمایش‌های این بخش با فشرده‌سازی داده‌های ۱۶، ۸، ۴ و ۲ بیت به ترتیب با دلتای ۸، ۴ و ۲ بیت انجام شد که نتایج دقت حاصل در جدول (۱) قابل مشاهده است. با توجه به این نتایج، فشرده‌سازی دلتا با اتلاف بهتر از حالت بدون اتلاف است.

True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6829	0	46	0	19
Sit	1	1175	8	8	2
Sta	77	2	4142	1	226
Stan	0	6	8	950	1
W	15	0	115	0	8332

True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6824	0	21	0	49
Sit	1	1172	9	8	4
Sta	85	3	3929	6	425
Stan	0	3	6	953	3
W	6	1	56	0	8399

True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6783	0	62	0	49
Sit	0	1172	14	7	1
Sta	110	0	3882	3	453
Stan	0	3	5	954	3
W	7	0	94	1	8361

True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6740	0	112	0	42
Sit	0	1173	9	9	3
Sta	71	0	4220	3	154
Stan	0	3	11	951	0
W	5	0	277	0	8180

True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6781	0	54	0	59
Sit	0	1177	8	8	1
Sta	154	0	3932	3	359
Stan	0	4	4	954	3
W	31	0	281	1	8149

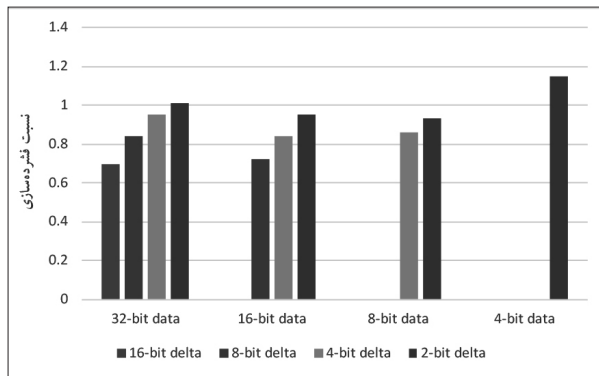
True Label \ Predicted Label	Jog	Sit	Sta	Stan	W
Jog	6769	0	98	0	27
Sit	0	1170	16	7	1
Sta	127	1	4006	1	313
Stan	0	3	14	947	1
W	24	0	406	0	8032

شکل ۵: اثر حذف نمونه‌های داده بر دقت تشخیص به تفکیک فعالیت برای بازیابی خطی (سمت چپ) و مکعبی (سمت راست) با حذف یک در میان (بالا)، دو در میان (وسط) سه در میان (پایین).

اتلاف آزمایش شده است. معیار اندازه‌گیری مؤثر بودن فشرده‌سازی، نسبت فشرده‌سازی است که طبق رابطه (۴) تعریف می‌شود.

$$(۴) \quad \text{نسبت فشرده سازی} = \frac{\text{اندازه داده غیر فشرده}}{\text{اندازه داده فشرده}}$$

اگرچه فشرده‌سازی بدون اتلاف اثری بر روی دقت ندارد، باید مؤثر بودن آن از طریق اندازه‌گیری نسبت فشرده‌سازی بررسی شود. از طرف دیگر، فشرده‌سازی با اتلاف پیشنهادی، اگرچه نسبت ثابتی از فشرده‌سازی ارائه می‌کند، باید اثر آن بر دقت تشخیص بررسی شود. در حالت بدون اتلاف ابتدا به جای ارسال داده‌های صحیح ۳۲ بیتی، تفاضل آن‌ها در قالب دلتای ۱۶، ۸، ۴ و ۲ بیتی ارسال می‌شود. در حالتی که تفاضل داده‌ها از ظرفیت دلتا بیشتر باشد نمونه داده اصلی ارسال می‌شود و برای تفکیک این حالت قبل از داده اصلی یک علامت قرار داده می‌شود و با



شکل ۸: اثر فشردسازی بدون اتلاف دلتا بر حجم نمونه داده

جدول ۱: تاثیر فشردسازی با اتلاف نمونه داده ها در دقت

نسبت فشردسازی	دقت با فشردسازی	دقت بدون فشردسازی	اندازه دلتا (بیت)	اندازه نمونه داده (بیت)
۰.۲۷	۹۴.۷	۹۶.۷	۶	۸
۰.۰۹	۹۵.۴	۹۷.۰	۳	۴
۰.۲۹	۹۵.۰	۹۷.۰	۲	۴

بدون اتلاف فشرد شده و در سمت ابر بازیابی می‌شوند. در دو روش پیشنهادی کاهش حجم تعداد نمونه‌ها و کاهش دقت اعشار نمونه‌های داده‌ها، کاهش حجم داده ارسالی تنها باعث کاهش اندکی در دقت تشخیص فعالیت می‌گردد. روش کاهش دقت اعشار به دلیل کاهش چشمگیرتر در حجم داده نسبت به روش اول برتری نشان می‌دهد. با استفاده از تنها ۴ بیت برای ارسال هر نمونه داده که باعث کاهش ۱۶ برابری حجم داده‌ها می‌شود، دقت از دست رفته محدود به یک درصد می‌باشد. اگر چه روش فشردسازی با اتلاف نتایج بهتری نسبت به روش بدون اتلاف نشان می‌دهد، هیچ یک از این دو روش به کارآمدی روش کاهش دقت و کاهش تعداد نمونه‌های داده نیستند. به‌عنوان کار آینده این تحقیق، نیاز به بررسی اثر ترکیبی روش‌های کاهش حجم در معامله بین دقت تشخیص و هزینه ارتباطات خواهد بود. ضمناً شایستگی روش پیشنهادی در سایر کاربردهای مختلف نیازمند بررسی می‌باشد. بررسی روش کاهش هزینه ارتباطی نظیر استفاده از روش کاهش دقت اعشار نمونه‌ها با نگاشت‌های مختلف غیر خطی و یا بررسی سایر روش‌های فشردسازی از دیگر کارهای پیش روی این تحقیق می‌باشد.

True Label	Predicted Label				
	Jog	Sit	Sta	Stan	W
Jog	3430	0	10	1	6
Sit	0	531	5	41	0
Sta	41	2	2004	6	212
Stan	0	19	17	401	0
W	0	1	78	1	4167

True Label	Predicted Label				
	Jog	Sit	Sta	Stan	W
Jog	3417	1	7	1	21
Sit	1	533	4	39	0
Sta	45	3	1875	0	342
Stan	0	1	2	432	2
W	5	0	19	1	4230

True Label	Predicted Label				
	Jog	Sit	Sta	Stan	W
Jog	3414	0	21	0	12
Sit	0	567	2	8	0
Sta	110	1	2036	2	216
Stan	0	33	1	403	0
W	5	0	47	2	4203

True Label	Predicted Label				
	Jog	Sit	Sta	Stan	W
Jog	3323	0	95	0	29
Sit	0	435	19	123	0
Sta	37	26	1810	19	373
Stan	0	89	60	283	5
W	2	20	236	13	3984

شکل ۷: اثر کاهش دقت اعشار نمونه‌ها بر دقت تشخیص به تفکیک فعالیت در نگاشت خطی ۳۲ بیتی (بالا سمت چپ)، ۸ بیتی (بالا سمت راست)، ۱۶ بیتی (پایین سمت چپ) و ۱ بیتی (پایین سمت راست)

به‌عنوان مثال در حالت داده ۵ بیتی در ازای دو درصد کاهش دقت، اندازه داده ارسالی به نصف کاهش یافته است. اما در هر حال روش فشردسازی دلتا به اندازه دو روش قبل مؤثر نمی‌باشد.

۵- نتیجه‌گیری

تحقیق حاضر به مطالعه امکان بهینه‌سازی برون‌سپاری محاسبات از طریق کاهش حجم داده ارسالی به ابر محاسباتی با تمرکز بر کاربرد تشخیص فعالیت انسانی با یادگیری عمیق اختصاص دارد. برای این منظور سه روش کاهش تعداد نمونه‌های داده، کاهش دقت اعشار نمونه‌های داده و فشردسازی نمونه‌های داده پیشنهاد شده است. در روش اول، داده‌ها قبل از ارسال به صورت یکی در میان یا بیشتر حذف شده و در سمت ابر به روش تخمینی بازیابی می‌شوند. در روش دوم، نمونه داده‌ها قبل از ارسال با تابع نگاشت خطی به عدد صحیح با تعداد بیت کمتر تبدیل شده و در سمت ابر با رابطه معکوس نگاشت تخمین زده می‌شوند. در روش سوم، داده‌ها با یک الگوریتم فشردسازی با سربار کم و به یکی از دو صورت فشردسازی با اتلاف و

- eration with Adaptive DNN Partitioning and Offloading”, IEEE INFOCOM 2020-IEEE Conference on Computer Communications. IEEE, 2020.
16. Li, He, Kaoru Ota, and Mianxiong Dong., “Learning IoT in edge: Deep learning for the Internet of Things with edge computing”, IEEE network, 2018.
 17. Jeong, Hyuk-Jin, et al., “IONN: Incremental offloading of neural network computations from mobile devices to edge servers”, Proceedings of the ACM Symposium on Cloud Computing, 2018.
 18. Teerapittayanon, Surat, Bradley McDanel, and Hsiang-Tsung Kung., “Distributed deep neural networks over the cloud, the edge and end devices”, IEEE 37th International Conference on Distributed Computing Systems (ICDCS). IEEE, 2017.
 19. Jeong, Hyuk-Jin, et al., “Computation offloading for machine learning web apps in the edge server environment”, 2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS). IEEE, 2018.
 20. Eshratifar, Amir Erfan, Amirhossein Esmaili, and Massoud Pedram., “Bottlenet: A deep learning architecture for intelligent mobile cloud computing services”, IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED). IEEE, 2019.
 21. Chung, Jae-Won, Jae-Yun Kim, and Soo-Mook Moon., “ShadowTutor: Distributed Partial Distillation for Mobile Video DNN Inference”, arXiv preprint arXiv:2003.10735, 2020.
 22. Eshratifar, Amir Erfan, Mohammad Saeed Abrishami, and Massoud Pedram., “JointDNN: an efficient training and inference engine for intelligent mobile cloud computing services”, IEEE Transactions on Mobile Computing, 2019.
 23. Rosloniec, S., “Fundamental numerical methods for electrical engineering”, Vol. 18. Springer Science & Business Media, 2008.
 24. <https://www.tensorflow.org/>.
 1. Graves, Alex, Abdel-rahman Mohamed, and Geoffrey Hinton., “Speech recognition with deep recurrent neural networks”, 2013 IEEE international conference on acoustics, speech and signal processing. IEEE, 2013.
 2. Eshratifar, Amir Erfan, and Massoud Pedram., “Energy and performance efficient computation offloading for deep neural networks in a mobile cloud computing environment”, Proceedings of the 2018 on Great Lakes Symposium on VLSI. 2018.
 3. Dey, Swarnava, Jayeeta Mondal, and Arijit Mukherjee., “Offloaded execution of deep learning inference at edge: Challenges and insights”, 2019 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops). IEEE, 2019.
 4. Huang, Yutao, et al., “DeePar: A hybrid device-edge-cloud execution framework for mobile deep learning applications”, IEEE INFOCOM 2019-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS). IEEE, 2019.
 5. Kemp, Roelof, et al., “Cuckoo: a computation offloading framework for smartphones”, International Conference on Mobile Computing, Applications, and Services. Springer, Berlin, Heidelberg, 2010.
 6. Mitchell, Tom M. “Machine learning”, Burr Ridge, IL: McGraw Hill, 1997.
 7. Karki, Aajna, et al., “Tango: A deep neural network benchmark suite for various accelerators”, IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS). IEEE, 2019.
 8. Guo, Kaiyuan, et al., “[DL] A survey of FPGA-based neural network inference accelerators”, ACM Transactions on Reconfigurable Technology and Systems (TRETS), 2019.
 9. Qin, Dongming, et al., “A novel combined prediction scheme based on CNN and LSTM for urban PM 2.5 concentration”, IEEE Access 7: 20050-20059, 2019.
 10. Hochreiter, Sepp, and Jürgen Schmidhuber., “Long short-term memory”, Neural computation 9.8 : 1735-1780, 1997.
 11. Ran, Xukan, et al., “Delivering deep learning to mobile devices via offloading”, Proceedings of the Workshop on Virtual Reality and Augmented Reality Network, 2017.
 12. Fadishei, H., “Energy-Efficient Human Activity Recognition on Smartphones: A Test-Cost Sensitive Approach”, International Journal of Information and Communication Technology Research, 2018.
 13. Weiss, Gary M., “WISDM Smartphone and Smartwatch Activity and Biometrics Dataset”, UCI Machine Learning Repository: WISDM Smartphone and Smartwatch Activity and Biometrics Dataset Data Set, 2019.
 14. Wang, Jianyu, et al., “Edge cloud offloading algorithms: Issues, methods, and perspectives”, ACM Computing Surveys (CSUR), 2019.
 15. Mohammed, Thaha, et al., “Distributed Inference Accel-