

ارائه روشی جدید و کم هزینه مبتنی بر شبکه عصبی پیچشی برای مقابله با حمله خصمانه بر مدل‌های یادگیری عمیق

علی سرآبادانی

دانشجوی دکتری، دانشکده فنی و مهندسی - دانشگاه قم - قم - ایران
پست الکترونیکی: alisarabadani14@gmail.com

امیر جلالی بیدگلی*

استادیار، دانشکده فنی و مهندسی - دانشگاه قم - قم - ایران
پست الکترونیکی: jalaly@qom.ac.ir

چکیده

و همچنین برای مقابله با حملات عمدی در حین حفظ کیفیت تصویر اشاره کرد. از نتایج روش پیشنهادی می‌توان به دستیابی به درصد موفقیت استحکام در برابر حملات خصمانه، ۸۲/۱۱ در مجموعه داده هدا، ۷۵/۱۹ در مجموعه داده CIFAR-10 و ۷۹/۳۴ در مجموعه داده MNIST اشاره کرد که دارای بهبود ۰٫۲ در مقایسه با سایر مقالات مروری است. همچنین این مقاله در مقایسه با روشهای پیشین فاقد پیچیدگی محاسباتی بالا و دارای زمان محاسباتی کمتر است.

واژه‌های کلیدی: حملات خصمانه، شبکه عصبی

پیچشی، نطفه نمک‌وفلفل، مدل‌های یادگیری عمیق

۱- مقدمه

حملات خصمانه اقدامات خرابکارانه‌ای هستند که هدفشان تضعیف عملکرد یادگیری عمیق، ایجاد رفتار نادرست مدل یا کسب اطلاعات محافظت شده است. یادگیری عمیق خصمانه در اوایل سال ۲۰۰۴ مورد مطالعه قرار گرفت [۱]. رشد یادگیری عمیق و ادغام آن در بسیاری

روش‌های حمله به مدل‌های یادگیری عمیق می‌توانند برچسب رده‌ها را تغییر دهند یا این‌که مخاطره‌ای به وجود آورند و خطرهای جدی امنیتی را ایجاد کنند. در حملات خصمانه مهاجمان با ایجاد تغییراتی اندک و البته حساب‌شده در داده‌های تصویری، بدون این‌که توجه کاربر جلب شود، الگوریتم را به اشتباه می‌اندازند. از طرفی به واسطه حمله، نطفه‌ها کیفیت تصویر را کاهش می‌دهند و باعث از بین رفتن اطلاعات می‌شوند. نطفه نمک‌وفلفل یکی از محبوب‌ترین نطفه‌هایی است که کیفیت تصویر را تحت تاثیر قرار می‌دهد. روش‌های زیادی برای حذف نطفه نمک‌وفلفل از تصویر با حداقل از دست دادن اطلاعات پیشنهاد شده‌است ولی در این مقاله روشی پیشنهاد شده‌است که نطفه نمک‌وفلفل با روشی مبتنی بر شبکه عصبی پیچشی در مجموعه داده‌های هدا، CIFAR-10 و MNIST برای مقابله با حملات خصمانه به تصویر اضافه شده‌است. از مزایای این روش پیشنهادی می‌توان به استفاده از کمترین میزان نطفه نمک‌وفلفل، جلوگیری از تعداد بیشتری از حملات خصمانه

از برنامه‌های کاربردی در سال‌های اخیر باعث تجدید علاقه به یادگیری عمیق متخاصم شده‌است.

نگرانی‌های زیادی در حوزه امنیت وجود دارد. می‌توان سیستم‌های مورد حمله را به دستاوردهای هوش مصنوعی مجهز کرد [۲]. برای رده‌بندی تصاویر مدل از کاربر تصاویری دریافت شده است که این تصاویر ممکن است مورد حمله متخاصم قرار بگیرد به همین جهت پیش از این بعضی از مدل‌ها برای مقابله با حملات متخاصم اقدام به اضافه کردن نوبه نمک‌وفلفل به تصاویر قبل از پردازش کردن آن‌ها می‌پرداختند که این روش از دو جهت ایراد داشت. در صورت کم یا ناکافی بودن میزان نوبه اضافه شده، مقاومت تصویر در مقابل حملات متخاصم پایین می‌آید. در صورت بیش از اندازه بالا بودن میزان نوبه نمک‌وفلفل کیفیت تصویر و در نتیجه دقت رده‌بندی کاهش می‌یابد. این روش بر مبنای خروجی شبکه عصبی پیچشی است. در روش پیشنهادی میزان مناسبی نوبه نمک‌وفلفل در پیکسل‌های مناسب اضافه شده است. با این روش ضمن حفظ کیفیت تصویر، امنیت مدل یادگیری عمیق در مقابل حملات خصمانه بالاتر می‌رود. اضافه کردن تعداد زیاد نوبه نمک‌وفلفل باعث کاهش قابل توجه کیفیت تصویر می‌شود. از این رو یافتن میزان مناسب نوبه نمک‌وفلفل و جایگاه پیکسل‌های مناسب قرار گرفتن آن‌ها اهمیت ویژه‌ای دارد. این چالش در روش پیشنهادی بر طرف شده است. میزان و جایگاه پیکسل نوبه نمک‌وفلفل اضافه شده به تصویرها را شبکه عصبی پیچشی تعیین می‌کند که در این مقاله به‌طور میانگین به تعداد ۸۰ پیکسل است. ساختار مقاله نیز به این صورت است که، در قسمت ۲ ادبیات موضوع مقاله مطرح شده است. در قسمت ۳ پیشینه پژوهش ارائه گردیده است. در قسمت ۴ روش پیشنهادی قرار گرفته است. قسمت ۵ ارزیابی و قسمت ۶ جمع‌بندی می‌باشد. از نتایج روش پیشنهادی می‌توان به دستیابی به درصد موفقیت استحکام در برابر حملات خصمانه، ۸۲/۱۱ در مجموعه

داده‌ها، ۷۵/۱۹ در مجموعه داده CIFAR-10 و ۷۹/۳۴ در مجموعه داده MNIST اشاره کرد.

۲- ادبیات موضوع

در بخش ۱-۲ توضیحاتی در مورد حملات خصمانه علیه مدل‌های یادگیری عمیق، نحوه عمل‌کردشان و همچنین محافظت از سیستم‌های یادگیری عمیق در مقابل این حملات داده شده است. در بخش ۲-۲ به معرفی نوبه نمک‌وفلفل که در روش پیشنهادی از آن استفاده می‌شود، شده است. در بخش ۲-۳ نیز تعاریفی از شبکه عصبی پیچشی مطرح شده است.

۲-۱- حملات خصمانه علیه مدل‌های یادگیری عمیق

یادگیری عمیق بخشی جدایی‌ناپذیر از بسیاری از برنامه‌هایی است که ما هر روز استفاده می‌کنیم، از قفل تشخیص چهره در آیفون گرفته تا عملکرد تشخیص صدای الکسا و صافی‌های اسپم در ایمیل‌های ما که انجام می‌پذیرد. اما فراگیر بودن یادگیری عمیق همچنین منجر به حملات خصمانه شده است. در این حملات، مهاجمان داده‌ها را دچار آشفتگی^۱ می‌کنند و آن‌ها را به‌عنوان ورودی به مدل می‌دهند، سپس از طریق این آشفتگی، ورودی‌های تخصصی ساخته می‌شود که عملکرد مدل را مختل می‌کنند. بنابراین تغییرات و آشفتگی‌هایی که در داده‌ها ایجاد می‌شوند تا آن‌ها را به نمونه‌های متخاصم تبدیل کند، غیر قابل تشخیص نیستند ولی بسیار به نمونه اولیه شبیه هستند که می‌تواند تشخیص آن‌ها را سخت کند. در سال‌های اخیر، مسئله‌ای تحت عنوان آسیب‌پذیری مدل‌های مبتنی بر یادگیری ماشین مطرح گردیده است که نشان می‌دهد مدل‌های یادگیری در مواجهه با آسیب‌پذیری‌ها از مقاومت بالایی برخوردار نیستند. یکی از معروف‌ترین آسیب‌ها یا به بیان دیگر حملات، تزریق مثال‌های تخصصی به مدل می‌باشد که در این مورد، شبکه‌های عصبی و به ویژه شبکه‌های عصبی عمیق بیشترین میزان آسیب‌پذیری

1- filter

2- perturbation

را دارند. مثال‌های تخصصی، از طریق افزودن اندکی نوفه هدفمند به مثال‌های اصلی تولید می‌شوند، به طوری که از منظر کاربر انسانی تغییر محسوس در داده‌ها مشاهده نمی‌شود اما مدل‌های یادگیری ماشینی در دسته‌بندی داده‌ها به اشتباه می‌افتند [۳].

برخلاف سیستم‌های نرم‌افزاری کلاسیک، که در آن تولیدکنندگان به صورت دستی دستورالعمل‌ها و قوانین را می‌نویسند، الگوریتم‌های یادگیری عمیق رفتار خود را از طریق تجربه توسعه می‌دهند.

به عنوان مثال، برای ایجاد یک سیستم تشخیص خط، تولیدکنندگان یک الگوریتم‌های یادگیری عمیق ایجاد می‌کند و با ارائه تصاویر برچسب‌گذاری شده از خطوط خیابان از زوایای مختلف و در شرایط نوری مختلف، آن را آموزش می‌دهد. سپس مدل یادگیری عمیق پارامترهای خود را تنظیم می‌کند تا الگوهای رایجی را که در تصاویر حاوی خطوط خیابان رخ می‌دهد، ثبت کند. با ساختار الگوریتم مناسب و نمونه‌های آموزشی کافی، این مدل قادر خواهد بود خطوط را در تصاویر و فیلم‌های جدید با دقت قابل توجهی تشخیص دهد [۴]. اما علیرغم موفقیت آن‌ها در زمینه‌های پیچیده مانند بینایی رایانه و تشخیص صدا، الگوریتم‌های یادگیری عمیق موتورهای استنتاج آماری هستند، به عبارت دیگر توابع پیچیده ریاضی هستند که ورودی‌ها را به خروجی تبدیل می‌کنند. اگر یک یادگیری عمیق تصویری را به عنوان حاوی یک شیء خاص برچسب‌گذاری کند، متوجه می‌شود که مقادیر پیکسل در آن تصویر از نظر آماری مشابه سایر تصاویر شیئی است که در طول آموزش پردازش کرده است [۵]. بنابراین با افزودن تکه‌های کوچک و نامحسوس پیکسل‌ها به یک تصویر، یک عامل مخرب می‌تواند باعث شود الگوریتم یادگیری ماشین آن را به عنوان یک تصویر نادرست در گروه نادرست رده‌بندی کند. انواع اختلالات اعمال شده در حملات خصمانه به نوع داده هدف و اثر مورد نظر بستگی دارد [۶]. مدل تهدید باید برای روش‌های مختلف داده سفارشی‌سازی شود تا

به طور منطقی خصمانه باشد. به عنوان مثال، برای تصاویر و فایل‌های صوتی، منطقی است که اغتشاش داده‌های کوچک را به عنوان یک مدل تهدید در نظر بگیریم، زیرا به راحتی توسط انسان قابل درک نخواهد بود، اما ممکن است مدل هدف را بدرفتار کند و باعث ناهماهنگی بین انسان و ماشین شود. با این حال، برای برخی از انواع داده‌ها مانند متن، به سادگی با تغییر یک کلمه یا یک نویسه، ممکن است معنایی را مختل کند و به راحتی توسط انسان شناسایی شود. بنابراین، مدل تهدید برای متن باید به طور طبیعی با تصویر یا صدا متفاوت باشد. همچنین این میزان برای تصاویر نیز متفاوت می‌باشد [۷].

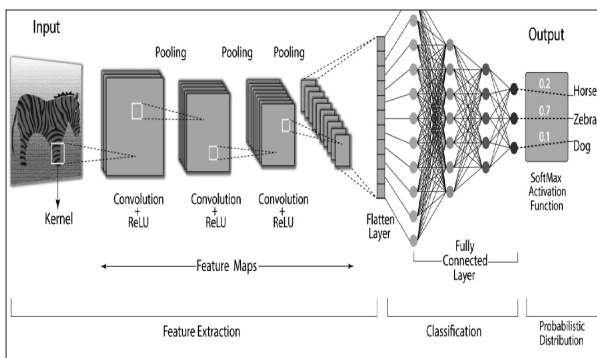
۲-۲- نوفه نمک و فلفل

هر حسگر دنیای واقعی تحت تاثیر درجه خاصی از نوفه قرار می‌گیرد، خواه حرارتی، الکتریکی یا غیره. به طور کلی، نوفه می‌تواند توسط حسگر تصویر و مدار یک اسکنر یا دوربین دیجیتال تولید شود. نوفه تصویر همچنین می‌تواند از دانه‌های فیلم و نوفه غیرقابل اجتناب یک آشکارساز فوتون ایده‌آل ایجاد کند. نوفه تصویر یک محصول جانبی نامطلوب ثبت تصویر است که اطلاعات مورد نظر را مبهم می‌کند. معنی اصلی نوفه، سیگنال ناخواسته^۳ است. نوفه نمک و فلفل نوعی نوفه ضربه‌ای در تصاویر است. ما نوفه نمک و فلفل را در نظر می‌گیریم که برای آن مقدار مشخصی از پیکسل‌های تصویر سیاه یا سفید (نقاط سیاه یا سفید) هستند. به طور معمول اگر نقاط سیاهی در تصویر وجود داشته باشد آن را نوفه فلفل و اگر نقاط سفید در تصویر وجود داشته باشد آن را نوفه نمک می‌نامیم که عمدتاً در فرآیند انتقال اطلاعات رخ می‌دهند. احتمال رخ دادن این نوفه فقط در دو مقدار بوده، یا صفر یا ۲۵۵ (تصویر ۸ بیتی)، یا یک سیگنال را نابود می‌کند (صفر می‌کند)، یا یک سیگنال را کاملاً یک می‌کند و چیزی بین آن وجود ندارد. (مقدار نوفه را با مقدار سیگنال جایگزین می‌کند). این نوفه عموماً به دلیل خطا در انتقال داده، خرابی سلول

3- Unwanted signal



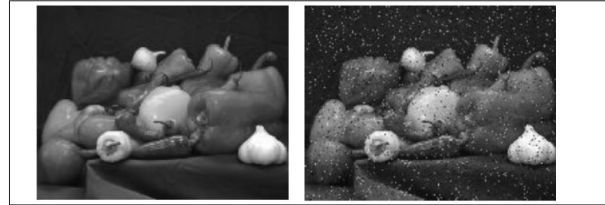
شکل ۲: تاثیر نوفه نمک و فلفل بر تصویر مبتنی بر میزان نوفه اضافه شده [۸]



شکل ۳: معماری شبکه عصبی پیچشی

خاصی از لبه‌ها وجود داشته‌باشد. برخی از نورون‌ها وقتی لبه‌های عمودی نشان داده می‌شوند، فعال می‌شوند، در حالی که برخی دیگر هنگامی که لبه‌های افقی یا مورب نشان داده می‌شوند، واکنش نشان می‌دهند [۱۰]. شبکه عصبی پیچشی نوعی شبکه عصبی مصنوعی است که در یادگیری عمیق برای ارزیابی اطلاعات بصری استفاده می‌شود. این شبکه‌ها می‌توانند طیف وسیعی از وظایف شامل تصاویر، صداها، متون، ویدئوها و سایر رسانه‌ها را انجام دهند. پروفیسور Yann LeCun از آزمایشگاه Bell اولین شبکه‌های پیچشی موفق را در اواخر دهه ۱۹۹۰ ایجاد کرد [۱۱]. در چند دهه گذشته، یادگیری عمیق به دلیل توانایی آن در مدیریت حجم زیادی از داده‌ها، ابزار بسیار قدرتمندی بوده است. علاقه به استفاده از لایه‌های پنهان، به ویژه در تشخیص الگو، از روش‌های سنتی پیشی گرفته است. یکی از محبوب‌ترین شبکه‌های عصبی عمیق، شبکه‌های عصبی پیچشی است.

5- hidden layers



شکل ۱: تاثیر نوفه نمک و فلفل بر تصویر [۹]

حافظه یا خطاهای مبدل آنالوگ به دیجیتال ایجاد می‌شود. اگر تصویر ۸ بیتی را در نظر گرفته‌شود، نوفه نمک و فلفل به‌طور تصادفی مقدار معینی از پیکسل‌ها را در دو حد ۰ یا ۲۵۵ ایجاد می‌کند. نوفه به‌طور قابل توجهی به اطلاعات تصویر آسیب می‌زند که منجر به مشکلاتی در موفقیت در انجام وظایف پردازش تصویر مانند تشخیص لبه یا تقسیم‌بندی تصویر می‌شود. وظایف تشخیص تصویر از آنجایی که پیکسل نوفه با اکثر همسایگان محلی خود متفاوت است، مقدار گرادینان بزرگی مانند پیکسل لبه دارد [۸]. در شکل ۱ اعمال نوفه نمک و فلفل را بر روی تصویر میوه مشاهده می‌کنید:

در شکل ۲ تاثیر نوفه نمک و فلفل بر تصویر مبتنی بر میزان نوفه اضافه شده را مشاهده می‌کنیم. اگر دقت کنید از تصویر k تا تصویر a به مرور میزان نوفه نمک و فلفل بیشتر است. در حدی که در تصویر a کیفیت تصویر به شدت کاهش پیدا کرده‌است و پیکسل‌های اصلی تصویر به وسیله نوفه‌های نمک و فلفل از بین رفته‌اند. به احتمال فراوان اگر تصویر را به مدلی مبتنی بر یادگیری عمیق بدهیم در رده‌بندی تصویر دچار مشکل می‌شود.

۲-۳- شبکه عصبی پیچشی

شبکه عصبی پیچشی که با نام‌های CNN یا ConvNets نیز شناخته می‌شوند، نوعی شبکه عصبی مصنوعی پیش‌خورنده هستند که ساختار اتصال آن از سازمان‌دهی قشر بینایی حیوانات الهام گرفته شده‌است. خوشه‌های کوچکی از سلول‌ها در قشر بینایی به مناطق خاصی از میدان بینایی حساس هستند. سلول‌های عصبی منفرد در مغز تنها زمانی واکنش نشان می‌دهند که جهت‌گیری‌های

4- feed-forward

شبکه‌های عصبی پیچشی دارای یک لایه ورودی، یک لایه خروجی، لایه‌های پنهان متعدد و میلیون‌ها پارامتر هستند که به آن‌ها اجازه می‌دهد اشیاء و الگوهای پیچیده را یاد بگیرند. از فرآیندهای پیچشی و ادغام برای نمونه برداری فرعی ورودی داده شده قبل از اعمال تابع فعال‌سازی استفاده می‌کند، جایی که همه آن‌ها لایه‌های پنهانی هستند که تا حدی به هم متصل هستند و لایه کاملاً متصل در انتها منجر به لایه خروجی می‌شود. شکل خروجی شبیه به اندازه تصویر ورودی است. پیچش^۱ فرآیند ترکیب دو تابع برای تولید خروجی تابع دیگر است. تصویر ورودی با استفاده از صافی‌ها در شبکه‌های عصبی پیچشی پیچیده می‌شود که منجر به یک نقشه ویژگی می‌شود. صافی‌ها وزن‌ها و سوگیری‌هایی هستند که به صورت تصادفی بردارهایی در شبکه تولید می‌شوند. شبکه‌های عصبی پیچشی به جای داشتن وزن‌ها و سوگیری‌های جداگانه برای هر نورون، از وزن‌ها و سوگیری‌های یکسان برای همه نورون‌ها استفاده می‌کند. صافی‌های زیادی را می‌توان ایجاد کرد که هر کدام جنبه متفاوتی از ورودی را می‌گیرد. هسته‌ها^۲ نام دیگری برای صافی‌ها هستند [۱۲].

۳- پیشینه پژوهش

محافظت از سیستم‌های یادگیری عمیق در مقابل حملات خصمانه

در برابر حملات متخاصم در چند سال گذشته، محققان هوش مصنوعی روش‌های مختلفی را توسعه داده‌اند تا مدل‌های یادگیری عمیق را در برابر حملات خصمانه قوی‌تر کنند. شناخته‌شده‌ترین روش دفاعی آموزش خصمانه^۳ است که در آن یک توسعه‌دهنده آسیب‌پذیری‌ها را با آموزش مدل یادگیری عمیق بر روی نمونه‌های متخاصم اصلاح می‌کند. سایر روش‌های دفاعی شامل تغییر یا بهینه‌سازی ساختار مدل است، مانند افزودن لایه‌های تصادفی و برون‌یابی بین

6- Convolution
7- Kernel
8- adversarial training

چندین مدل یادگیری عمیق برای جلوگیری از سوء استفاده از آسیب‌پذیری‌های متخاصم هر مدل حملات خصمانه را می‌شود راهی هوشمندانه برای انجام آزمایش فشار^۹ و اشکال‌زدایی^{۱۰} روی مدل‌های یادگیری عمیق دانست که بالغ^{۱۱} در نظر گرفته می‌شوند. اگر شما معتقدید که یک فناوری قبل از تبدیل شدن به یک محصول باید به‌طور کامل آزمایش و اشکال‌زدایی شود، پس حمله خصمانه - به منظور آزمایش قوی و بهبود - یک گام اساسی در خط لوله توسعه فناوری یادگیری عمیق خواهد بود. محققان تعدادی راهکار دفاع از حمله متخاصم را پیشنهاد کرده‌اند که می‌توان آن‌ها را به سه دسته اصلی، یعنی اصلاح داده‌ها^{۱۲}، اصلاح مدل‌ها^{۱۳} و استفاده از ابزارهای کمکی^{۱۴} تقسیم کرد. ما به ترتیب آن‌ها را با جزئیات شرح می‌دهیم:

۳-۱- اصلاح داده‌ها

این راهکارها به اصلاح مجموعه داده‌های آموزشی در مرحله آموزش یا تغییر داده‌های ورودی در مرحله آزمایش، از جمله آموزش خصمانه، پنهان‌سازی گرادیان، مسدود کردن قابلیت انتقال، فشرده‌سازی داده‌ها و تصادفی‌سازی داده‌ها اشاره دارد.

۱- آموزش خصمانه

در آموزش خصمانه، نمونه‌های متخاصم به مجموعه داده‌های آموزشی معرفی می‌شوند تا با استفاده از مدل آموزشی با نمونه‌های متخاصم قانونی، استحکام مدل هدف را بهبود بخشند. سگدی و همکاران [۱۳] ابتدا نمونه‌های متخاصم را تزریق کرده و برچسب‌های آن را اصلاح کرد تا مدل در برابر دشمنان قوی‌تر شود.

۲- پنهان‌سازی گرادیان

پنهان‌سازی گرادیان^{۱۵}، یک دفاع طبیعی در برابر حملات مبتنی بر گرادیان و حملات با استفاده از روش ساخت

9- Pressure
10- Debugging
11- mature
12- modifying data
13- modifying models
14- using auxiliary tools
15- Gradient Hiding

خصمانه مانند FGSM در [۱۴] ارائه شده است. این روش اطلاعات مربوط به گرادیان مدل را از دشمنان پنهان می‌کند، یعنی اگر یک مدل غیرقابل تمایز باشد (به عنوان مثال، یک درخت تصمیم^{۱۶}، یک رده‌بندی‌کننده نزدیک‌ترین همسایه^{۱۷}، یا یک جنگل تصادفی^{۱۸}، حمله مبتنی بر گرادیان نامعتبر است. با این حال، با یادگیری مدل جعبه‌سیاه پیشکار^{۱۹} با گرادیان و استفاده از نمونه‌های متخاصم تولید شده توسط این مدل [۱۵]، می‌توان روش را به راحتی در این مورد فریب داد.

۳- مسدود کردن قابلیت انتقال

در مسدود کردن قابلیت انتقال^{۲۰} از آنجایی که ویژگی انتقال پذیری حتی اگر رده‌بندی‌کننده‌ها معماری متفاوتی داشته باشند یا بر روی مجموعه داده‌های مجزا آموزش دیده باشند، برقرار است، کلید جلوگیری از حمله جعبه سیاه، جلوگیری از قابلیت انتقال نمونه‌های متخاصم است. حسینی و همکاران [۱۶] یک روش سه مرحله‌ای NULL Labeling را پیشنهاد کردند تا از نمونه‌های متخاصم از یک شبکه به شبکه دیگر جلوگیری کند. ایده اصلی آن افزودن یک برچسب NULL جدید به مجموعه داده و رده‌بندی آن‌ها به برچسب NULL با آموزش رده‌بندی‌کننده برای مقاومت در برابر حملات متخاصم است. این روش به طور کلی شامل سه مرحله اصلی، یعنی آموزش اولیه رده‌بندی‌کننده هدف، محاسبه احتمالات NULL و آموزش خصمانه می‌باشد.

۴- متراکم‌سازی داده‌ها

برای متراکم‌سازی داده‌ها^{۲۱}، Dziugaite و همکاران [۱۷] دریافتند که روش فشرده‌سازی JPG می‌تواند دقت تشخیص مدل شبکه را که ناشی از اختلال حمله FGSM

16- decision tree
17- nearest neighbor classifier
18- random forest
19- proxy
20- Blocking the Transferability
21- Data Compression

کاهش یافته است بهبود بخشید. داس و همکاران [۱۸] از یک روش فشرده‌سازی JPEG مشابه برای مطالعه یک روش دفاعی در برابر حملات FGSM و DeepFool استفاده کرد. با این حال، این فناوری‌های فشرده‌سازی تصویر هنوز نمی‌توانند به عنوان یک دفاع موثر در برابر حملات قدرتمندتر مانند حملات Carlini و Wagner عمل کنند [۱۹] به‌طور مشابه، روش فشرده‌سازی فناوری فشرده‌سازی نمایشگر که در مبارزه با حملات اختلال جهانی استفاده می‌شود [۲۰] نیز ناکافی است. بزرگترین محدودیت این روش‌های دفاعی مبتنی بر فشرده‌سازی داده‌ها این است که مقدار زیاد فشرده‌سازی منجر به کاهش دقت رده‌بندی تصویر اصلی می‌شود، در حالی که مقدار کمی فشرده‌سازی اغلب برای حذف تأثیر اختلال کافی نیست.

۵- تصادفی سازی داده‌ها

در تصادفی‌سازی داده‌ها^{۲۲} زی و همکاران [۲۱] نشان داد که عملیات تغییر اندازه تصادفی نمونه‌های متخاصم می‌تواند اثربخشی نمونه‌های متخاصم را کاهش دهد. به‌طور مشابه، افزودن برخی بافت‌های تصادفی^{۲۳} به نمونه‌های متخاصم نیز می‌تواند فریب‌دهی آن‌ها را به مدل شبکه کاهش دهد. وانگ و همکاران [۲۲] از یک ماژول تبدیل داده جدا شده از مدل شبکه برای از بین بردن اختلال احتمالی متخاصم در تصویر استفاده کرد و عملیات گسترده‌ای را در فرآیند آموزش انجام داد، مانند افزودن برخی پردازش‌های تصادفی گاوسی، که می‌تواند اندکی استحکام مدل شبکه را بهبود بخشد.

۳-۲- اصلاح مدل‌ها

در روش اصلاح مدل شبکه‌عصبی را اصلاح می‌شود، مانند با قاعده‌سازی، تقطیر دفاعی، فشردن ویژگی، شبکه انقباضی عمیق، پوشش دفاعی و شبکه‌های تجزیه.

22- Data Randomization
23- random textures

۱- با قاعده‌سازی

هدف روش با قاعده‌سازی^{۲۴}، بهبود توانایی تعمیم مدل هدف با افزودن عبارات با قاعده است که به عنوان عبارات مجازات^{۲۵} شناخته می‌شوند به تابع هزینه^{۲۶} و باعث می‌شود که مدل دارای سازگاری خوبی برای مقاومت در برابر حملات به مجموعه داده ناشناخته در پیش‌بینی باشد. بیگیو و همکاران [۲۳] از یک روش با قاعده‌سازی برای محدود کردن آسیب‌پذیری داده‌ها هنگام آموزش مدل SVM استفاده کرد. مقالات [۲۴-۲۶] از روش با قاعده‌سازی برای بهبود استحکام الگوریتم استفاده کردند و به نتایج خوبی در مقاومت در برابر حملات متخاصم دست یافتند.

۲- تقطیر دفاعی

برای تقطیر دفاعی^{۲۷}، پاپرنوت و همکاران [۲۷] یک روش تقطیر دفاعی برای مقاومت در برابر حملات بر اساس فناوری تقطیر پیشنهاد کرد. در مقاله [۲۸] هدف فناوری تقطیر اصلی فشرده‌سازی مدل در مقیاس بزرگ به مقیاس کوچک و حفظ دقت اولیه است، در حالی که تقطیر دفاعی مقیاس مدل را تغییر نمی‌دهد و مدلی با سطح خروجی صاف‌تر و حساسیت کمتر به اختلال ایجاد می‌کند تا استحکام مدل را بهبود ببخشد.

۳- فشردن ویژگی

فشردن ویژگی^{۲۸} یک روش بهبود مدل است [۲۹] که ایده اصلی آن کاهش پیچیدگی نمایش داده‌ها است و در نتیجه کاهش تداخل متخاصم به دلیل حساسیت کم است. دو روش اکتشافی وجود دارد، یکی کاهش عمق رنگ در سطح پیکسل، یعنی کدگذاری رنگ با مقادیر کمتر و دیگری استفاده از یک صافی مسطح بر روی تصویر است، به عنوان مثال، ورودی‌های متعدد به یک مقدار نگاهت می‌شوند، بنابراین مدل را تحت نوفه و حملات دفاعی^{۲۹} ایمن‌تر می‌کند. اگرچه این روش می‌تواند به‌طور موثری

24- Regularization
25- penalty terms
26- penalty terms
27- Defensive Distillation
28- Feature Squeezing
29- confrontational attack

از حملات خصمانه جلوگیری کند، اما دقت رده بندی نمونه‌های واقعی را نیز کاهش می‌دهد.

۴- شبکه انقباضی عمیق

در شبکه انقباضی عمیق^{۳۰} گو و همکاران [۳۰] نوعی شبکه فشرده‌سازی عمیق^{۳۱} را معرفی کرد که از رمزگذار خودکار کاهش نوفه برای کاهش نوفه دشمن استفاده می‌کند [۳۱]. بر اساس این پدیده، DCN ثابت کرد که اثر دفاعی خاصی در برابر حملاتی مانند L-BGFS دارد.

۵- پوشش دفاعی

در پوشش دفاعی^{۳۲}، گائو و همکاران [۳۲] پیشنهاد کرد که یک لایه پوشش قبل از پردازش مدل شبکه رده‌بندی شده درج شود. این لایه پوشش تصاویر اصلی و نمونه‌های متخاصم مربوطه را آموزش داده و تفاوت‌های بین این تصاویر و ویژگی‌های خروجی لایه مدل شبکه قبلی را کدگذاری می‌کند. به‌طور کلی اعتقاد بر این است که مهم‌ترین وزن در لایه اضافی مربوط به حساس‌ترین ویژگی در شبکه است. بنابراین، در رده‌بندی نهایی، این ویژگی‌ها با فشار دادن لایه‌های اضافی با وزن اولیه صفر پوشانده می‌شوند. به این ترتیب می‌توان از انحراف نتایج رده‌بندی ناشی از نمونه‌های متخاصم محافظت کرد.

۶- شبکه‌های تجزیه

در شبکه‌های تجزیه^{۳۳}، سیسه و همکاران [۳۳] شبکه‌ای به نام شبکه‌های تجزیه را به عنوان روشی دفاعی در برابر حملات متخاصم پیشنهاد کرد. این شبکه با کنترل ثابت Lipschitz شبکه، قاعده‌سازی سلسله مراتبی را اتخاذ می‌کند.

۳-۳- استفاده از ابزارهای کمکی

این رویکرد به استفاده از ابزارهای اضافی به عنوان ابزار کمکی برای مدل شبکه‌عصبی، از جمله Defense-MagNet، GAN و HGD اشاره دارد.

30- Deep Contractive Network (DCN)
31- deep compression network
32- Mask Defense
33- Parseval Networks

۱- Defense-GAN

سمنگوی و همکاران [۳۴-۳۵] سازوکاری را پیشنهاد کرد که هم برای حملات جعبه‌سفید و هم برای حملات جعبه‌سیاه قابل اجراست تا کارایی اغتشاش خصمانه را کاهش دهد. این روش از قدرت شبکه متخاصم مولد استفاده می‌کند.

۲- MagNet

منگ و همکاران [۳۶] چارچوبی به نام MagNet پیشنهاد کرد که خروجی آخرین لایه رده‌بندی‌کننده را به‌عنوان یک جعبه‌سیاه بدون خواندن داده‌های لایه داخلی یا تغییر رده‌بندی‌کننده می‌خواند. MagNet از یک آشکارساز برای شناسایی نمونه‌های قانونی و متخاصم استفاده می‌کند. آشکارساز^{۳۴} فاصله بین یک نمونه معین تحت آزمایش و نمونه‌های متنوع را اندازه‌گیری می‌کند و اگر فاصله از آستانه بیشتر شود، نمونه را رد می‌کند. همچنین از یک اصلاح‌کننده برای تبدیل نمونه متخاصم از طریق یک رمزگذار خودکار به نمونه قانونی مشابه استفاده می‌کند. در حملات جعبه‌سفید، عملکرد MagNet به‌طور قابل توجهی کاهش می‌یابد زیرا دشمنان پارامترهای MagNet را می‌دانند. بنابراین، نویسنده ایده استفاده از چندین رمزگذار خودکار و انتخاب تصادفی یکی در یک زمان را پیشنهاد کرد تا پیش‌بینی رمزگذار خودکار مورد استفاده دشمن را دشوار کند.

۳- High-Level Representation Guided Denoiser

برخلاف دستگاه استاندارد حذف نوفه مانند تابع از دست دادن^{۳۵} بازسازی در سطح پیکسل، که دارای مشکل تقویت خطا است، HGD^{۳۶} می‌تواند به‌طور موثری بر این مشکل با استفاده از یک تابع از دست دادن غلبه کند تا خروجی‌های مدل هدف را با تصویر تمیز مقایسه کند. لیاؤ و همکاران [۳۷]، HGD را برای طراحی یک مدل هدف قوی در برابر حملات خصمانه جعبه‌سفید و جعبه‌سیاه معرفی کرد. دیگر استفاده از HGD این است که می‌توان آن را

34- detector
35- loss function
36- High-Level Representation Guided Denoiser

بر روی یک مجموعه داده نسبتاً کوچک آموزش داد و می‌تواند برای محافظت از مدل‌هایی غیر از مدلی که آن را هدایت می‌کند، استفاده شود. در این مقاله پژوهشگران ۲۵ مقاله را در دامنه موضوعی پژوهش بررسی کرده و پس از بررسی روش‌های قبلی به یک جمع‌بندی برای مدل پیشنهادی خود دست یافته‌اند. جدول (۱) کارهای پیشین بررسی شده را معرفی می‌کند.

۴- روش پیشنهادی

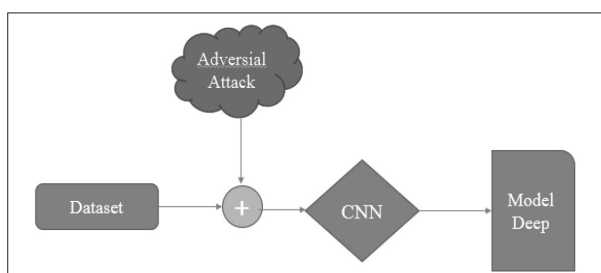
قابل ذکر است که روش پیشنهادی، دارای نوآوری در مقابل سایر مقالات و پژوهش‌ها در موضوع استحکام مدل یادگیری عمیق در برابر حملات خصمانه می‌باشد. همان‌طور که می‌دانید در برابر حملات متخاصم در چند سال گذشته، محققان هوش مصنوعی روش‌های مختلفی را توسعه داده‌اند تا مدل‌های یادگیری عمیق را در برابر حملات خصمانه قوی‌تر کنند. شناخته‌شده‌ترین روش دفاعی آموزش خصمانه است که در آن یک توسعه‌دهنده آسیب‌پذیری‌ها را با آموزش مدل یادگیری عمیق بر روی نمونه‌های متخاصم اصلاح می‌کند. سایر روش‌های دفاعی شامل تغییر یا بهینه‌سازی ساختار مدل است. هیچ‌یک از این روش‌ها مبتنی بر اضافه کردن نوفه^{۳۷} نمک و فلفل به تصاویر به وسیله شبکه عصبی پیچشی نبوده است. در مقاله در روش پیشنهادی، قصد بر آن است که در وهله اول یک شبکه عصبی پیچشی را طبق مراحل زیر آموزش داده شود:

مرحله اول - مجموعه داده اولیه را در معرض حمله خصمانه قرار داده و یک مجموعه داده ثانویه خواهیم داشت.

مرحله دوم - مجموعه داده اولیه و مجموعه داده ثانویه را با هم مقایسه کرده و حاصل اختلاف آن‌ها را به صورت یک مجموعه داده نهایی برای آموزش دادن شبکه عصبی پیچشی به کار می‌بریم. می‌دانیم، هر پیکسل عددی معادل ۰ تا ۲۵۵ می‌باشد ولی عدد دو پیکسل (عدد

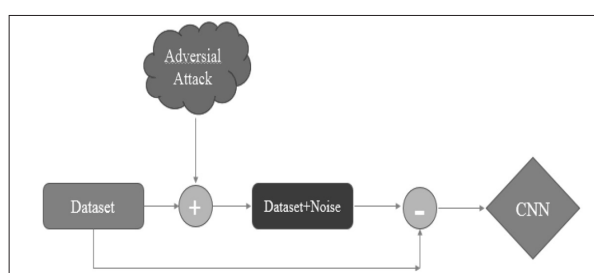
جدول ۱: خلاصه مقالات مروری

منابع	نقاط ضعف	نقاط قوت	انواع راهکارها	دسته اصلی
[۱۳]	تزیق نمونه‌های متخاصم	قوی تر شدن در برابر دشمنان، استحکام مدل هدف	آموزش خصمانه	اصلاح داده‌ها
[۱۴-۱۵]	فریب سریع	استحکام بیشتر در برابر حملات	پنهان‌سازی گرادبان	
[۱۶]	زمانبر بودن	عدم ورود نمونه‌های متخاصم از یک شبکه به شبکه دیگر	مسدود کردن قابلیت انتقال	
[۱۷-۲۰]	مناسب نبودن در برابر حملات قدرتمند، کاهش دقت رده بندی	روش فشرده‌سازی JPG، افزایش دقت	متراکم‌سازی داده‌ها	
[۲۱-۲۲]	کیفیت بالایی در برابر حملات سخت ندارد.	کاهش اثر بخشی نمونه‌های متخاصم، کاهش فریب دهی	تصادفی سازی داده‌ها	
[۲۳-۲۶]	نیاز به اعمال قاعده و قوانین دارد.	سازگاری خوب در برابر حملات به مجموعه داده ناشناخته، مقاومت در برابر حملات متخاصم	قاعده‌سازی	اصلاح مدل
[۲۸-۲۷]	مناسب نبودن در برابر حملات قدرتمند	عدم تغییر در مقیاس مدل، حساسیت کمتر به اختلال، استحکام مدل	تقطیر دفاعی	
[۲۹]	کاهش دقت رده بندی نمونه‌های واقعی	کاهش پیچیدگی نمایش داده‌ها، کاهش تداخل متخاصم، ایمن بودن در برابر نوفه	فشردن ویژگی	
[۳۰-۳۱]	عملیات رمزگذاری خود پیچیدگی دارد.	رمزگذار خودکار کاهش نوفه برای کاهش نوفه دشمن	شبکه انقباضی عمیق	
[۳۲]	مناسب نبودن در برابر حملات قدرتمند	محافظت از انحراف نتایج رده‌بندی نمونه‌های متخاصم	پوشش دفاعی	
[۳۳]	زمان بر بودن پیاده سازی	استفاده از قاعده سازی سلسله مراتبی	شبکه‌های تجزیه	
[۳۵-۳۴]	به نسبت بقیه کارایی خوبی دارد	کاهش کارایی اغتشاش خصمانه	Defense-GAN	ابزارهای کمکی
[۳۶]	کاهش عملکرد در حملات جعبه سفید	استفاده از آشکار ساز برای شناسایی نمونه‌های قانونی و متخاصم	Magnet	
[۳۷]	پیچیدگی پیاده سازی	مدل هدف قوی در برابر حملات خصمانه جعبه سفید و جعبه سیاه، محافظت از مدل‌هایی غیر از مدل مورد استفاده	High-Level Representation Guided Denoiser	



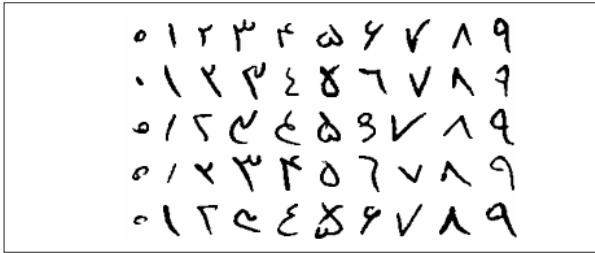
شکل ۵: اضافه کردن نوفه نمک و فلفل به وسیله شبکه عصبی پیچشی و مقابله با حملات خصمانه

مورد حمله خصمانه قرار گرفته است را برای آزمایش به آن می‌دهیم. شبکه عصبی پیچشی در مرحله آموزش یاد گرفته است که کدام پیکسل‌های تصاویر مورد حمله خصمانه قرار گرفته است. حال با گرفتن ورودی آزمایش، پیکسل‌های مورد حمله را تشخیص داده و در آن پیکسل‌ها



شکل ۴: آموزش شبکه عصبی پیچشی

منتاسب به پیکسل تصویر بدون نوفه، عدد منتاسب به پیکسل تصویر با نوفه را از هم کم می‌کنیم عددی بین ۲۵۵ تا ۲۵۵- می‌شود. شکل ۵ این مراحل را نشان می‌دهد. در وهله دوم، حال که شبکه عصبی پیچشی ما آموزش داده شده است به عنوان ورودی آن مجموعه داده‌ای که



شکل ۶: نمونه‌هایی از دستخط‌های مختلف موجود در مجموعه ارقام دست‌نویس مجموعه داده هدا

۵-۱-۲- مجموعه داده CIFAR-10

مجموعه داده CIFAR-10 (انستیتوی تحقیقات پیشرفته کانادا) مجموعه‌ای از تصاویر است که معمولاً برای آموزش الگوریتم‌های یادگیری ماشین و بینایی رایانه استفاده می‌شود. مجموعه داده CIFAR-10 مجموعه تصاویری است که معمولاً برای آموزش الگوریتم‌های یادگیری ماشین و بینایی رایانه استفاده می‌شود. این یکی از پرکاربردترین مجموعه‌های داده برای تحقیقات یادگیری ماشین است. مجموعه داده CIFAR-10 شامل ۶۰,۰۰۰ تصویر رنگی ۳۲×۳۲ در ۱۰ رده مختلف است. ۱۰ رده مختلف هواپیما، ماشین، پرنده، گربه، آهو، سگ، قورباغه، اسب، کشتی و کامیون را نشان می‌دهد. از هر رده ۶۰۰۰ تصویر وجود دارد. شبکه‌های عصبی پیچشی بهترین روش تشخیص تصاویر در CIFAR-10 هستند. CIFAR-10 زیرمجموعه‌ای از ۸۰ میلیون عکس که تمامی آن‌ها با استفاده از برچسب نام‌گذاری شده‌است. در شکل ۸، نمونه‌هایی از کیفیت دست‌خط‌های موجود در مجموعه ارقام دست‌نویس مجموعه داده هدا را نشان داده شده است.

۵-۱-۳- مجموعه داده MNIST

پایگاه داده MNIST^{۳۷}، یک پایگاه داده بزرگ از ارقام دست‌نویس است که معمولاً برای آموزش سیستم‌های مختلف پردازش تصویر استفاده می‌شود. این پایگاه داده همچنین به‌طور گسترده‌ای برای آموزش و آزمایش در زمینه یادگیری ماشین استفاده می‌شود. پایگاه داده MNIST شامل ۶۰,۰۰۰ تصویر آموزشی و ۱۰,۰۰۰ تصویر آزمایشی است. سازندگان اصلی پایگاه داده فهرستی از برخی از

نوفه نمک و فلفل اضافه کرده و خروجی را به مدل یادگیری عمیق می‌دهد. شکل ۶ این مرحله را نشان می‌دهد. قابل ذکر است که شبکه عصبی پیچشی ما دارای ۵ لایه می‌باشد و تابع فعال ساز ۴ لایه اولیه ReLU می‌باشد اما تابع فعال ساز لایه آخر از نوع سیگموئید می‌باشد و همان‌طور که می‌دانیم این تابع یک منحنی S شکل است. زمانی که می‌خواهیم خروجی مدل احتمال باشد، از تابع سیگموئید استفاده می‌کنیم؛ چون تابع سیگموئید مقادیر را به بازه صفر تا ۱ می‌برد و احتمالات هم میان همین بازه قرار دارند.

۵-ارزیابی

در این بخش، ابتدا در بخش ۵-۱ به معرفی مجموعه داده‌های استفاده شده در مقاله می‌پردازیم سپس در بخش ۵-۲، ملاحظات و مقایسه در قالب جدول ۱ و نمودارهای مختلف به بررسی خروجی آزمایش‌ها بر روی مجموعه داده‌های معرفی شده در بخش ۵-۱ با استفاده از روش‌های نامبرده در بخش روش‌های پیشنهادی می‌پردازیم.

۵-۱- مجموعه داده‌های مورد استفاده

در این بخش به معرفی مجموعه داده‌های هدا، CIFAR-10 و MNIST می‌پردازیم.

۵-۱-۱- مجموعه داده ارقام دست‌نویس هدا

مجموعه ارقام دست‌نویس هدا که اولین مجموعه بزرگ ارقام دست‌نویس فارسی است، مشتمل بر ۱۰۲۳۵۳ نمونه دست‌نوشته سیاه سفید است. خصوصیات این مجموعه داده به شرح زیر است [۲۸]:

درجه تفکیک نمونه‌ها: ۲۰۰ نقطه بر اینچ

تعداد کل نمونه‌ها: ۱۰۲۳۵۲ نمونه

تعداد نمونه‌های آموزش: ۶۰۰۰ نمونه از هر رده

تعداد نمونه‌های آزمایش: ۲۰۰۰ نمونه از هر رده

سایر نمونه‌ها: ۲۲۳۵۲ نمونه

در شکل ۶، نمونه‌هایی از دستخط‌های موجود در مجموعه ارقام دست‌نویس مجموعه داده هدا را نشان داده شده است.



شکل ۸: نمونه‌هایی کیفیتی از دستخط‌های مختلف موجود در مجموعه ارقام دست‌نویس مجموعه داده MNIST

میزان نوفه درصد موفقیت مقابله با حملات خصمانه، افزایش یافته (مقدار ۰/۳۱) و همزمان مقدار Drop Accuracy نیز افزایش می‌یابد و به مقدار ۱۲/۱۱ می‌رسد که تغییر اولی مطلوب و تغییر دومی نامطلوب است. در روش ۳، مشاهده می‌شود با استفاده از روش هوشمند اضافه کردن نوفه نمک و فلفل توسط شبکه عصبی پیچشی با مقدار مناسب و کافی که می‌تواند متغیر باشد ولی به‌طور میانگین به میزان ۸۰ می‌باشد، به درصد موفقیت بالایی برای مقابله با حملات خصمانه در حین نگه داشتن مقدار عددی Drop Accuracy در کم‌ترین حالت خود می‌رسد. در سایر مجموعه داده‌ها نیز می‌توان به طریق مشابه داده‌های به‌دست آمده را بررسی کرد و با توجه به این‌که مجموعه داده HDA و MNIST تقریباً مشابه یکدیگر هستند انتظار داریم داده‌های جدولشان نیز نزدیک به هم باشند و همچنین به دلیل رنگی بودن تصاویر در مجموعه داده CIFAR-10 توقع داریم که درصد موفقیت برای مقابله با حملات خصمانه، Drop Accuracy در تمامی روش‌ها و در مقایسه با سایر مجموعه داده‌ها کمتر باشد. در روش اصلاح داده، این راهکارها به اصلاح مجموعه داده‌های آموزشی در مرحله آموزش یا تغییر داده‌های ورودی در مرحله آزمایش می‌پردازد ولی در روش پیشنهادی عملیاتی بر مجموعه داده‌های آموزشی انجام نشده است. همچنین در روش اصلاح مدل شبکه عصبی اصلاح می‌شود که در روش پیشنهادی هیچگونه تغییر یا جابه‌جایی در لایه‌های مدل یا پارامترهای دیگر مدل صورت نگرفته



شکل ۷: نمونه‌هایی از تصاویر مختلف موجود در مجموعه داده CIFAR-10

روش‌های آزمایش شده بر روی آن را نگه می‌دارند.

۲-۵- ملاحظات و مقایسه

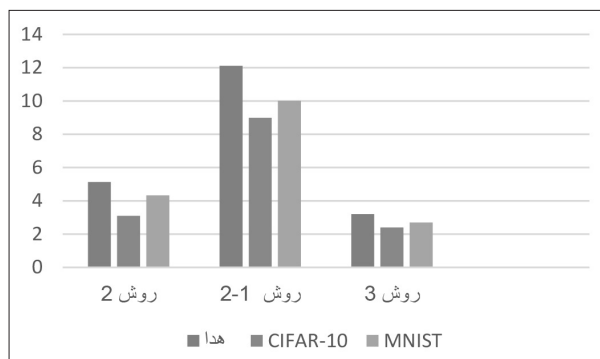
در این بخش نتایج آزمایش‌ها در قالب جدول ۱ و چهار نمودار مختلف نشان داده و مورد تجزیه و تحلیل قرار گرفته است.

۵-۲-۱- جدول نتایج

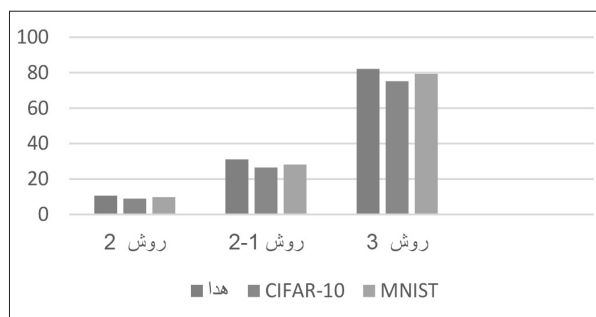
در مجموعه داده HDA، روش ۱ یعنی عدم انجام هیچگونه تدبیر برای مقابله با حمله خصمانه مفاهیم درصد موفقیت، Drop Accuracy و نوفه اضافه شده را تعریف نمی‌کند. پس ما مقدار عددی صفر را در جدول به جهت مقایسه با سایر مقادیر عددی قرار می‌دهیم. در روش ۲، قبل از ارسال تصاویر به مدل یادگیری عمیق برای مقابله با حملات خصمانه به صورت تصادفی مجموعه داده نوفه نمک و فلفل را اضافه می‌کنیم در این روش مقدار نوفه اضافه شده در حدود مقدار عددی ۲۰۰ پیکسل و در نتیجه یک مقدار ناکافی نوفه می‌باشد که باعث می‌شود درصد موفقیت در حدود مقدار عددی ۱۰/۶۳ باشد. با این حال در روش ۲ مقدار Drop Accuracy کم و قابل قبولی می‌باشد. در روش ۱-۲ مجدداً نوفه نمک و فلفل به جهت مقابله به جهت مقابله با حملات خصمانه به صورت تصادفی ولی این بار با مقدار نوفه بیشتر یعنی حدود مقدار عددی ۸۰۰ پیکسل به مجموعه داده اضافه می‌شود. در نتیجه، مشاهده می‌کنیم به دلیل افزایش

جدول ۱: مقایسه نتایج روش های مختلف در مجموعه داده های CIFAR-10 و MNIST

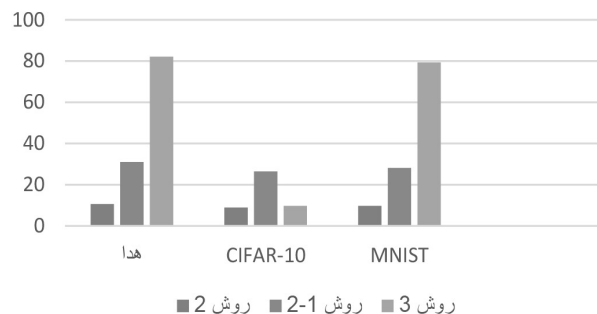
	مجموعه داده هدا			مجموعه داده CIFAR-10			مجموعه داده MNIST		
	درصد موفقیت	Accuracy Drop	نوفه اضافه شده	درصد موفقیت	Accuracy Drop	نوفه اضافه شده	درصد موفقیت	Accuracy Drop	نوفه اضافه شده
روش ۱	۰	۰	۰	۰	۰	۰	۰	۰	۰
روش ۲	۱۰/۶۳	۵/۱۳	۲۰۰	۸/۹۳	۳/۱۰	۲۰۰	۹/۷۴	۴/۳۳	۲۰۰
روش ۲-۱	۳۱/۰۵	۱۲/۱۱	۸۰۰	۲۶/۴۴	۸/۹۹	۸۰۰	۲۸/۱۵	۱۰/۰۱	۸۰۰
روش ۳	۸۲/۱۱	۳/۲	هوشمندانه (۸۰)	۷۵/۱۹	۲/۴	هوشمندانه (۸۰)	۷۹/۳۴	۲/۷	هوشمندانه (۸۰)



شکل ۱۰: Drop Accuracy روش های مختلف برای هر گروه از مجموعه داده ها



شکل ۹: درصد موفقیت روش های مختلف برای هر گروه از مجموعه داده ها



شکل ۱۱: درصد موفقیت مجموعه داده های مختلف برای هر گروه از روش ها

۲- نمودار دوم

مشاهده می شود که به صورت میانگین در روش ۱-۲ برای تمامی گروه های مجموعه داده بیشترین Drop Accuracy و در روش ۳ کمترین Drop Accuracy را مشاهده می شود.

۳- نمودار سوم

مشاهده می شود که به صورت میانگین در مجموعه داده هدا برای تمامی روش ها بیشترین درصد موفقیت را داریم و در مجموعه داده CIFAR-10 کمترین درصد موفقیت را مشاهده می شود.

است. در راهکار استفاده از ابزارهای کمکی این رویکرد به استفاده از ابزارهای اضافی به عنوان کمکی برای مدل شبکه عصبی، به مدل های یادگیری عمیق در مقابل حملات خصمانه کمک می کند ولی در روش پیشنهادی از هیچ ابزار کمکی استفاده نشده است.

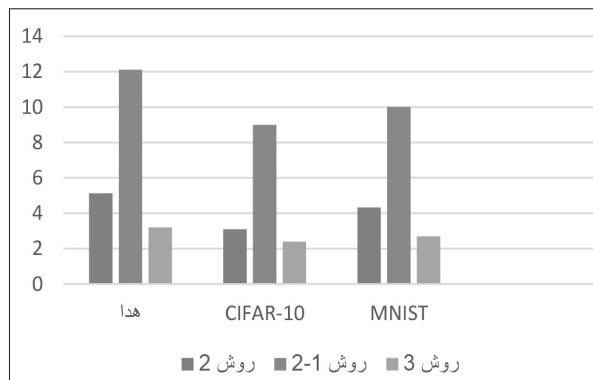
۵-۲-۲- نمودارها

حال سراغ بررسی چهار نمودار صفحه بعدی می رویم، که در آن دو نمودار اول به ترتیب درصد موفقیت و Drop Accuracy روش های مختلف برای هر گروه از مجموعه داده ها می باشد. همچنین در نمودارهای سوم و چهارم به ترتیب درصد موفقیت و Drop Accuracy مجموعه داده های مختلف برای هر گروه از روش های موجود است.

۱- نمودار اول

مشاهده می شود که در روش ۳ برای تمامی گروه های مجموعه داده بیشترین درصد موفقیت و در روش ۲ کمترین درصد موفقیت را مشاهده می شود.

۸۰ مورد می‌باشد و به وسیله شبکه عصبی پیچشی‌ای به دست آمده است، سبب می‌شود که ضمن استفاده از میزان کم و مناسب نوفه کیفیت تصویر حفظ شده و درصد موفقیت بالایی را برای مقابله با حملات متخاصم به همراه داشته باشد. این روش پیشنهادی را بر روی ۳ مجموعه داده معروف هدا، CIFAR-10 و MNIST آزمایش شده است و روش خود را با ۳ روش مرسوم دیگر مقایسه شده است و با این مقایسه قابل درک می‌باشد که روش پیشنهادی، روشی کارآمد و کم هزینه برای مقابله با حملات خصمانه است.



شکل ۱۲: Accuracy Drop مجموعه داده‌های مختلف برای هر گروه از روش‌ها

۴- نمودار چهارم

مشاهده می‌شود که به صورت میانگین در مجموعه داده هدا برای تمامی گروه‌های روش‌های موجود بیشترین Drop Accuracy و در مجموعه داده CIFAR-10 کم‌ترین Drop Accuracy را مشاهده می‌شود.

۶- جمع بندی

همان‌طور که گفته شد نوفه به اختلالات ناخواسته به وجود آمده روی تصویر گفته می‌شود، به گونه‌ای که بر روی کیفیت تصویر تاثیر منفی می‌گذارد. یکی از دلایل حذف نوفه از تصویر افزایش کیفیت تصویر جهت نمایش بهتر تصویر می‌باشد. از دلایل دیگر حذف نوفه این است که این نوفه مشکلی برای یکی از اساسی‌ترین مراحل پردازش تصویر یعنی رده بندی می‌باشد. در واقع اگر تصویر نوفه داشته باشد احتمال موفقیت مرحله رده‌بندی کمتر است. در حملات خصمانه مهاجمان با ایجاد تغییرات اندک و البته حساب شده در داده‌های تصویری، بدون این‌که توجه کاربر جلب شود، الگوریتم را به اشتباه می‌اندازند و منجر به ایجاد خطرهای جدی امنیتی می‌شود، برای مقابله با این حملات از روش اضافه کردن نوفه نمک و فلفل استفاده شده است. نوفه نمک و فلفل یکی از محبوب‌ترین نوفه‌هایی است که کیفیت تصویر را تحت تاثیر قرار می‌دهد. بنابراین در این مقاله روشی پیشنهاد شده که مدعی است که توزیع هوشمند نوفه نمک و فلفل بر روی تصویر که به طور میانگین

۶- منابع

- [1] Cezara Benegui, Radu Tudor Ionescu. Adversarial Attacks on Deep Learning Systems for User Identification based on Motion Sensors. arXiv:2009.01109.
- [2] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, Debdeep Mukhopadhyay. Adversarial Attacks and Defences: A Survey. <https://doi.org/10.48550/arXiv.1810.00069>. 28 Sep 2018.
- [3] Jianhe Yuan; Zhihai He, "Adversarial Dual Network Learning With Randomized Image Transform for Restoring Attacked Images," IEEE Access · 24 January 2020.
- [4] Anirban Chakraborty, Manaar Alam, "Adversarial Attacks and Defences: A Survey," in <https://arxiv.org/abs/1810.00069>. 28 Sep 2018.
- [5] Gu, S.; Rigazio, L. "Towards deep neural network architectures robust to adversarial examples", 2014. arXiv:1412.5068
- [6] Akhtar, N.; Mian, A. "Threat of adversarial attacks on deep learning in computervision: A survey", IEEE Access, vol. 6, pp. 14410-14430, 2018.
- [7] Tramèr, F.; Kurakin, A.; Papernot, N.; Goodfellow, I.; Boneh, D.; McDaniel, P. "Ensemble adversarial training: Attacks and defenses, 2017. arXiv:1705.07204.
- [8] Ziad Alqadi, "Salt and Pepper Noise: Effects and Removal," in International Journal on Electrical Engineering and Informatics · July 2018.
- [9] Raymond Hon Fu Chan, Mila Nikolova, "Salt-and-Pepper Noise Removal by Median-Type Noise Detectors and Detail-Preserving Regularization," in IEEE Transactions on Image Processing · November 2005.
- [10] "Convolutional Neural Networks (LeNet) – DeepLearning 0.1 documentation". DeepLearning 0.1. LISA Lab. Retrieved 31 August 2013.
- [11] Matusugu, Masakazu; Katsuhiko Mori; Yusuke Mitari; Yuji Kaneda (2003). "Subject independent facial expression recognition with robust face detection using a convolutional neural network" (PDF). Neural Networks. 16 (5): 555–559.

- deep neural networks. In Proceedings of the 2016 IEEE Symposium on Security and Privacy (SP), San Jose, CA, USA, 22–26 May 2016; pp. 582–597.
- [28] Hinton, G.; Vinyals, O.; Dean, J. Distilling the knowledge in a neural network. arXiv 2015, arXiv:1503.02531.
- [29] Xu, W.; Evans, D.; Qi, Y. Feature squeezing: Detecting adversarial examples in deep neural networks. arXiv 2017, arXiv:1704.0115.
- [30] Gu, S.; Rigazio, L. Towards deep neural network architectures robust to adversarial examples. arXiv 2014, arXiv:1412.5068.
- [31] Rifai, S.; Vincent, P.; Muller, X.; Glorot, X.; Bengio, Y. Contractive auto-encoders: Explicit invariance during feature extraction. In Proceedings of the 28th International Conference on International Conference on Machine Learning, Bellevue, WA, USA, 28 June–2 July 2011; pp. 833–840.
- [32] Gao, J.; Wang, B.; Lin, Z.; Xu, W.; Qi, Y. Deepcloak: Masking deep neural network models for robustness against adversarial samples. arXiv 2017, arXiv:1702.06763.
- [33] Cisse, M.; Adi, Y.; Neverova, N.; Keshet, J. Houdini: Fooling deep structured prediction models. arXiv 2017, arXiv:1707.05373.
- [34] Samangouei, P.; Kabkab, M.; Chellappa, R. DefenseGAN: Protecting classifiers against adversarial attacks using generative models. arXiv 2018, arXiv:1805.06605.
- [35] Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Generative adversarial nets. In Advances in Neural Information Processing Systems, Proceedings of the Annual Conference on Neural Information Processing Systems, Lake Tahoe, NV, USA, 5–10 December 2013; Curran Associates, Inc.: Red Hook, NY, USA, 2014; pp. 2672–2680.
- [36] Meng, D.; Chen, H. Magnet: A two-pronged defense against adversarial examples. In Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, Dallas, TX, USA, 30 October–3 November 2017; pp. 135–147.
- [37] Liao, F.; Liang, M.; Dong, Y.; Pang, T.; Zhu, J.; Hu, X. Defense against adversarial attacks using high-level representation guided denoiser. arXiv 2017, arXiv:1712.02976.
- [38] Hossein Khosravi, Ehsanollah Kabir, “Introducing a very large dataset of handwritten Farsi digits and a study on their varieties,” in Pattern Recognition Letters 28 1133–1141, (2007). doi:10.1016/S0893-6080(03)00115-1. Retrieved 17 November 2013.
- [12] Alzubaidi, L., Zhang, J., Humaidi, A.J. et al. Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. J Big Data 8, 53 (2021). <https://doi.org/10.1186/s40537-021-00444-8>.
- [13] Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; Fergus, R. Intriguing properties of neural networks. arXiv 2013, arXiv:1312.6199.
- [14] Tramèr, F.; Kurakin, A.; Papernot, N.; Goodfellow, I.; Boneh, D.; McDaniel, P. Ensemble adversarial training: Attacks and defenses. arXiv 2017, arXiv:1705.07204.
- [15] Papernot, N.; McDaniel, P.; Goodfellow, I.; Jha, S.; Celik, Z.B.; Swami, A. Practical black-box attacks against machine learning. In Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, Abu Dhabi, UAE, 2–6 April 2017; pp. 506–519.
- [16] Hosseini, H.; Chen, Y.; Kannan, S.; Zhang, B.; Poovendran, R. Blocking transferability of adversarial examples in black-box learning systems. arXiv 2017, arXiv:1703.04318.
- [17] Dziugaite, G.K.; Ghahramani, Z.; Roy, D.M. A study of the effect of jpg compression on adversarial images. arXiv 2016, arXiv:1608.00853.
- [18] Das, N.; Shanbhogue, M.; Chen, S.T.; Hohman, F.; Chen, L.; Kounavis, M.E.; Chau, D.H. Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression. arXiv 2017, arXiv:1705.02900.
- [19] Akhtar, N.; Liu, J.; Mian, A. Defense against Universal Adversarial Perturbations. arXiv 2017, arXiv:1711.05929.
- [20] Carlini, N.; Wagner, D. Towards evaluating the robustness of neural networks. In Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP), San Jose, CA, USA, 22–26 May 2017; pp. 39–57.
- [21] Xie, C.; Wang, J.; Zhang, Z.; Zhou, Y.; Xie, L.; Yuille, A. Adversarial examples for semantic segmentation and object detection. arXiv 2017, arXiv:1703.08603.
- [22] Wang, Q.; Guo, W.; Zhang, K.; Ororbia, I.; Alexander, G.; Xing, X.; Liu, X.; Giles, C.L. Learning adversary-resistant deep neural networks. arXiv 2016, arXiv:1612.01401.
- [23] Biggio, B.; Nelson, B.; Laskov, P. Support vector machines under adversarial label noise. In Proceedings of the Asian Conference on Machine Learning, Taoyuan, Taiwan, 13–15 November 2011; pp. 97–112.
- [24] Lyu, C.; Huang, K.; Liang, H.N. A unified gradient regularization family for adversarial examples. In Proceedings of the 2015 IEEE International Conference on Data Mining (ICDM), Atlantic City, NJ, USA, 14–17 November 2015; pp. 301–309.
- [25] Zhao, Q.; Griffin, L.D. Suppressing the unusual: Towards robust cnns using symmetric activation functions. arXiv 2016, arXiv:1603.05145.
- [26] Rozsa, A.; Gunther, M.; Boulton, T.E. Towards robust deep neural networks with BANG. arXiv 2016, arXiv:1612.00138.
- [27] Papernot, N.; McDaniel, P.; Wu, X.; Jha, S.; Swami, A. Distillation as a defense to adversarial perturbations against