

احتمال تعمیم‌پذیر نبودن روش‌های ارزیابی کیفیت مبتنی بر بردار پشتیبان

آرمین احمدزاده* و پوریا چراغی†

مرکز پردازش سریع، پژوهشگاه علوم کامپیوتر، پژوهشگاه دانش‌های بنیادی، تهران، ایران

چکیده

مطلوب،) دچار تخریب شده و یک نسخه‌ی تخریب‌شده از آن بوجود می‌آید. وقتی هردوی این دو تصویر در دسترس باشند، می‌توان برای ارزیابی کیفیت تصویر تخریب‌شده، به اطلاعات موجود در تصویر سالم رجوع کرد. به این کار اکت «مرجع کامل» گفته می‌شود [۵]. یک مثال از کاربرد اکت مرجع کامل، فشرده‌سازی تصاویر است. الگوریتم فشرده‌کننده با کاهش اطلاعات نشانک^۶، حجم آن را کاهش داده و از طرفی کیفیت آن را نیز خدشه‌دار می‌نماید. الگوریتم اکت مرجع کامل می‌تواند در هر لحظه تصویر فشرده‌شده را با تصویر اصلی مقایسه کرده و در صورت تخریب بیش از حد کیفیت، این مورد را به الگوریتم فشرده‌کننده بازخورد دهد [۶]. اگر $x(i, j)$ نشانک مرجع و $y(i, j)$ نشانک تخریب‌شده باشد، میانگین مربعات خطا، «MSE»^۷، اختلاف مقادیر آن دو را خواهد سنجید:

$$MSE(x, y) = \frac{1}{M \times N} \sum_{i=1}^M \sum_{j=1}^N (x(i, j) - y(i, j))^2 \quad (1)$$

M و N در (۱)، ابعاد تصاویر هستند. نشان داده می‌شود که «MSE» یا بیشینه‌ی نسبت نشانک به نوفه^۸، «PSNR»^۹، نمی‌تواند درک انسان از کیفیت تصویر را شبیه‌سازی نمایند [۷]. وانگ^{۱۰} و همکارانش نشان دادند که تخریب ساختارهای تصویر در کاهش کیفیت آن مؤثر است [۷]. در سال ۲۰۰۴، روشی به نام «SSIM»^{۱۱} برای کمی‌سازی شباهت ساختاری دو تصویر ارائه کردند [۸] که با اختلاف از روش‌های مبتنی بر «MSE» بهتر بود [۹].

با موفقیت SSIM، محققین سعی کردند ساختارهای تصویر را با ویژگی‌های دیگری، مثل لبه‌ها، مدل کنند. آماره‌های تصاویر طبیعی^{۱۲} [۱۰] و پایداری اطلاعاتی^{۱۳} [۱۱]، دیگر معیارهای پیشنهاد شده برای کیفیت تصویر هستند. روش‌های یادگیری ماشین سر-تا-سری^{۱۴} نیز، برای اکت محاسباتی استفاده شده‌اند [۱۲، ۱۳، ۱۴].

یک ساز و کار رایج برای اکت، استفاده از ماشین بردار پشتیبان، «SVR»^{۱۵} [۱۵] است [۱۶]. کلیات بکارگیری «SVR» در شکل ۱ دیده می‌شود. در روش‌های مبتنی بر بردار پشتیبان، نوآوری اصلی، در طراحی ویژگی‌ها صورت می‌گیرد. چنین روشی، ابتدا از تصویر ورودی یک بردار ویژگی استخراج می‌کند. اگر روش مرجع کامل باشد، می‌تواند برای استخراج ویژگی به نسخه‌ی سالم نیز رجوع کند. بردار ویژگی، f ،

یک مدل محاسباتی که بتواند نظر انسان در مورد کیفیت تصاویر را پیش‌بینی نماید، کاربردهای فراوانی در مسائل مربوط به پردازش تصویر دارد. به دلیل پیچیدگی دستگاه بینایی بشر، مدل‌های یادگیری ماشین برای شبیه‌سازی نحوه‌ی قضاوت انسان‌ها استفاده می‌شوند. یکی از این مدل‌ها، وایزش بردار پشتیبان است. در این مقاله، نشان می‌دهیم که روش رایج برای بکارگیری بردار پشتیبان در ارزیابی کیفیت تصویر، لزوماً به مدل‌های قابل تعمیم نمی‌انجامد.

واژه‌های کلیدی: هوش مصنوعی، پردازش تصویر، دستگاه بینایی انسان، یادگیری ماشین، ارزیابی کیفیت تصویر، نامفاد، ماشین بردار پشتیبان، وایزش

۱ مقدمه

تصاویر دیجیتال بخش قابل توجهی از مصرف رسانه‌ای بشر را تشکیل می‌دهند [۱]. بدیهی است که کیفیت کم این رسانه، منجر به نارضایتی کاربران و ناکامی ما در دریافت اطلاعات مدنظرمان خواهد شد. لذا، مطلوب است که وضعیت تصاویر ارسالی، از منظر کیفیت‌شان، تحت نظارت باشد.

مطمئن‌ترین راه برای ارزیابی کیفیت تصویر (که به اختصار «اکت» می‌خوانیم)، پرسش از انسان‌هاست [۲]. به این ترتیب که تصویری را به جمعی از ناظرین^۱ نشان داده و نظر آن‌ها در مورد کیفیت تصویر را در قالب نمره‌ای از یک بازه‌ی مشخص (مثلاً [0, 100])، جویا می‌شویم. میانگین نظرات^۲، یک کمیت قابل اطمینان از کیفیت تصویر است [۳]. چنین آزمایشی نیز اکت «انسانی»^۳ نام دارد.

بدیهی است که اکت انسانی برای کاربردهای بر-خط^۴ و حجیم امکان‌پذیر نیست. در این موارد، نرم‌افزاری مطلوب است که بتواند قضاوت دستگاه بینایی بشر در مورد کیفیت یک تصویر را پیش‌بینی نماید. دقت و سرعت این پیش‌بینی، معیارهای عملکرد یک مدل محاسباتی اکت^۵ هستند. ساخت چنین مدلی، موضوع مورد مطالعه‌ی بخش قابل توجهی از تحقیقات پردازش تصویر است [۴]. برای ارزیابی و آموزش روش‌های محاسباتی، از نتایج اکت انسانی استفاده می‌شود.

ساده‌ترین راه برای اکت محاسباتی، محاسبه‌ی اختلاف دو تصویر است. موارد زیادی وجود دارند که یک تصویر سالم (با کیفیت احتمالاً

^۶معادل فارسی «signal»

^۷Mean Squared Error

^۸معادل فارسی «noise»

^۹Peak Signal-to-Noise Ratio

^{۱۰}Zhou Wang

^{۱۱}Structural Similarity

^{۱۲}Natural Scene Statistics- NSS

^{۱۳}Information Fidelity

^{۱۴}end-to-end

^{۱۵}Support Vector Regression

a.ahmadzadeh@ipm.ir*

cheraaqee@ipm.ir†

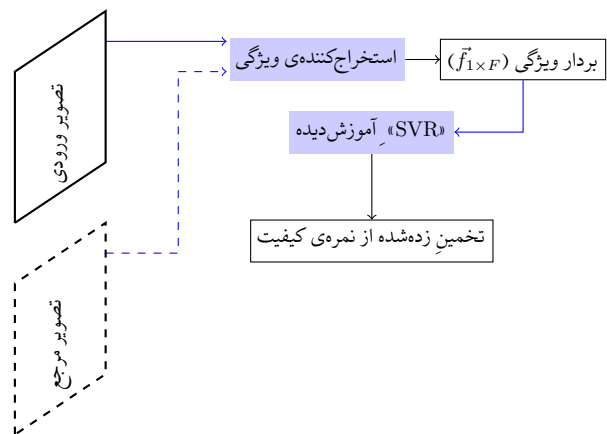
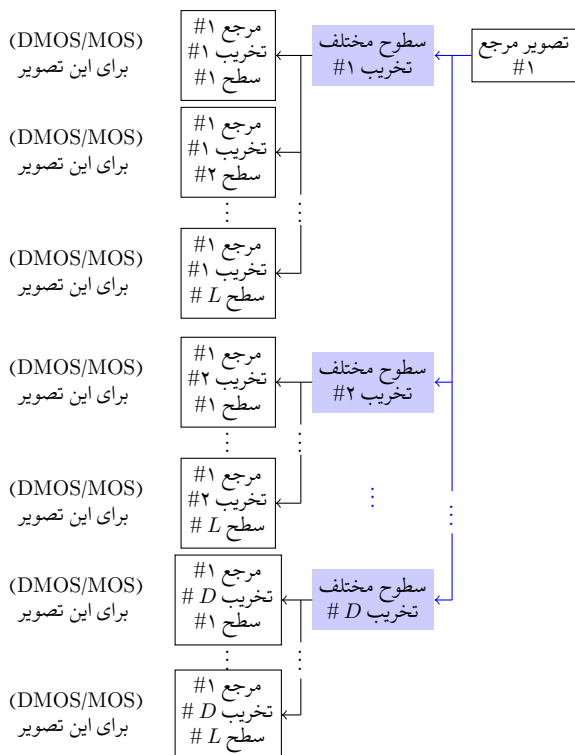
^۱همان «subject»ها در آزمایش‌های علمی

^۲که به اختصار «MOS» نامیده می‌شود و مخفف عبارت «Mean Opinion Score» است. لازم به توضیح است که «MOS» نسبت مستقیمی با کیفیت ادراکی و «Difference MOS» یا «DMOS» نسبتی عکس با کیفیت تصویر دارد.

^۳subjective image quality assessment

^۴معادل فارسی «on-line»

^۵objective image quality assessment



شکل ۱: بکارگیری «SVR» برای اکت

آرایه‌ای از اعداد است. اگر F ویژگی استخراج شوند، ابعاد این بردار، $1 \times F$ خواهد بود. یک مدل مبتنی بر وایزش^{۱۶} بردار پشتیبان، آموخته است که این بردار ویژگی را به نمره‌ی کیفیت نگاشت کند. برای آموزش چنین مدلی، از تصاویر و نمرات ارزیابی‌های انسانی استفاده می‌شود. در این مقاله نشان می‌دهیم که مدل حاصل از روش رایج برای آموزش بردار پشتیبان که در اکت‌های محاسباتی استفاده می‌شود، لزوماً قابل تعمیم^{۱۷} به تصاویر مختلف نیست.

مبانی مورد نیاز و برخی از کارهای مرتبط در قسمت ۲ مرور می‌شوند. قسمت ۳، آزمایش‌های انجام شده را تشریح خواهد کرد و مقاله در قسمت ۴ جمع‌بندی می‌شود.

۲ مبانی و مرور ادبیات

در این قسمت برخی مفاهیم و قراردادهای رایج در تحقیقات اکت محاسباتی تشریح می‌شوند. همچنین، یک دسته‌بندی از حوزه‌های ارزیابی کیفیت ارائه می‌گردد.

۱.۲ مجموعه داده‌ها

در اکثر کاربردها، انسان ناظر نهایی تصاویر است. لذا، نظر انسان بهترین معیار اکت برای این کاربردها خواهد بود. روش‌های محاسباتی سعی می‌کنند که نمراتی مثل نتایج اکت انسانی تولید کنند. مجموعه داده‌هایی وجود دارد که نظرات انسان به همراه تصاویر مربوطه را در اختیار محققین قرار می‌دهند، تا روش‌های اکت محاسباتی خود را محک بزنند [۳]. به این مجموعه داده‌ها، مجموعه داده‌های اکت گفته می‌شود.

یک مجموعه داده‌ی اکت، شامل یک سری تصویر مرجع، نسخ تخریب‌شده و نمرات انسانی است (شکل ۲). برخی اتفاقات برای تصاویر رایج هستند؛ مثلاً تار شدن به علت فاصله‌ی کانونی نامناسب، یا تصنعی^{۱۸} که به دلیل فشرده‌سازی ایجاد می‌شود. به دلیل رایج بودن این تخریب‌ها، مهم است که توانایی الگوریتم‌های اکت در سنجش آن‌ها بررسی شود. به همین دلیل، یک سری تصویر سالم، با کیفیت قابل قبول، انتخاب کرده و این تخریب‌ها را به صورت مصنوعی روی آن‌ها اعمال می‌کنند. هر تخریب هم می‌تواند با شدت متفاوتی اعمال شود. یک تصویر می‌تواند مقداری تار، یا خیلی تار باشد. بدیهی است که نمره‌ی

^{۱۶}معادل فارسی «regression»

^{۱۷}generalization
^{۱۸}artifacts

شکل ۲: قسمتی از آنچه که در یک مجموعه داده‌ی اکت وجود دارد. فرض می‌کنیم که R تصویر مرجع و D نوع تخریب وجود دارد که هر تخریب با L سطح از شدت اعمال می‌شود. شکل، موارد را برای مرجع #۱ نشان می‌دهد. همین موارد برای مراجع #۲ تا # R نیز تکرار می‌شوند.

انسانی برای شدت‌های مختلف، متفاوت خواهد بود. انتظار می‌رود که نمره‌ی حاصل‌شده از روش محاسباتی نیز با نمره‌ی انسانی مطابق باشد [۱۷].

بنابراین، هر مجموعه داده‌ی اکت، شامل مجموعه‌ای از «نمونه»‌ها است. هر نمونه، یک زوج مرتب به شکل (نمره‌ی انسانی، تصویر) است. نمره‌ی انسانی می‌تواند به صورت MOS یا DMOS ذخیره شده باشد. به این ترتیب که، «تصویر»، ورودی یک روش محاسباتی، و «نمره‌ی انسانی»، پاسخ صحیح مورد انتظار^{۱۹} است. حال اگر روش مرجع کامل باشد، تصویر مرجع متناظر نیز به عنوان ورودی به الگوریتم داده می‌شود. برای ارزیابی یک الگوریتم، نمرات آن برای تصاویر موجود در یک مجموعه داده محاسبه شده و سپس همبستگی نمرات الگوریتم با نمرات انسانی اندازه‌گیری می‌شود. معیار همبستگی، همان شاخص‌های آماری (پیرسن^{۲۰}، اسپیرمن^{۲۱} و غیره) هستند. هرچه این همبستگی بیشتر باشد، الگوریتم در پیش‌بینی نظر انسان، دقیق‌تر عمل کرده است [۹].

۲.۲ آموزش و آزمون روش‌های مبتنی بر SVR

طبق روش رایج برای بکارگیری SVR در اکت [۱۸]، ابتدا باید صحنه‌های یک مجموعه داده را افزایش^{۲۲} کنیم. منظور از «صحنه»^{۲۳} در یک

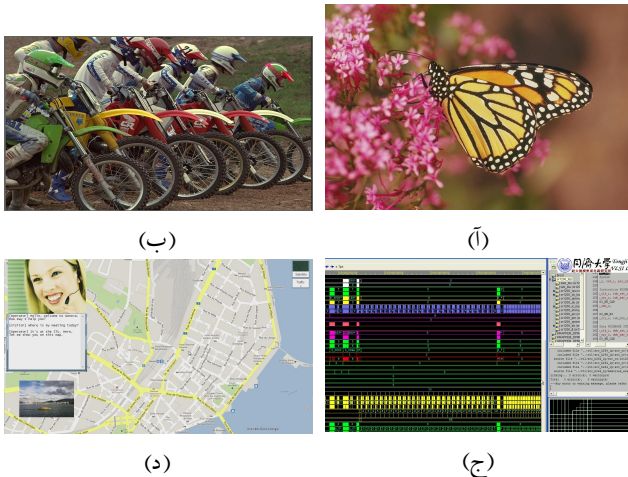
^{۱۹}همان «ground truth» در ادبیات یادگیری ماشین

^{۲۰}Pearson Correlation Coefficient- PLCC

^{۲۱}Spearman Rank Order Correlation Coefficient- SROCC

^{۲۲}افزاد یک مجموعه به معنی تقسیم آن به زیرمجموعه‌هایی است که با هم اشتراک ندارند. یعنی اگر یک مجموعه داده را به زیرمجموعه‌هایی برای آموزش و آزمون افزایش کنیم، آن زیرمجموعه‌ها عضو مشترکی ندارند.

^{۲۳}scene



شکل ۳: دو نمونه تصویر طبیعی (در (آ) و (ب)) و دو نمونه نامفاد (در (ج) و (د))

یک الگوریتم در ارزیابی این تخریب‌ها، همان اعتبار سنجی ۱۰۰۰-لایه را انجام می‌دهیم. تفاوت این است که هنگام سنجش الگوریتم روی مجموعه‌ی آزمون، تنها نمونه‌هایی را در نظر می‌گیریم که دارای تخریب مورد نظر هستند [۱۸].

۳.۲ حوزه‌های مختلف اکت

مشخص است که دقت الگوریتم مبتنی بر بردار پشتیبان، به طور اصلی، به ویژگی‌هایی که استخراج می‌شوند بستگی دارد. هر چه این ویژگی‌ها گویاتر^{۳۱} باشند، SVR هم کار راحت‌تری برای یادگیری نحوه‌ی نگاشت ویژگی‌ها به نمره‌ی کیفیت خواهد داشت [۱۷]. ویژگی‌های گویا برای هر تصویری متفاوت خواهند بود و طراحی روشی که بتواند برای هر تصویری ویژگی‌های مربوط به کیفیت را محاسبه کند، کار دشواری است [۱۹]. غالب تحقیقات اکت محاسباتی، برای تصاویر طبیعی انجام گرفته است [۲۰]. تصویر طبیعی، تصویری است که با دوربین‌های حساس به طیف مرئی موج الکترومغناطیس و از جهان فیزیکی پیرامون گرفته شده باشد [۲۱].

با شیوع تصویربرداری دیجیتال، تصاویر زیادی از اسناد و نوشته‌ها بوجود آمده است. ارزیابی کیفیت این تصاویر به روش‌هایی نیاز دارد که با روش‌های مربوط به تصاویر طبیعی متفاوت هستند [۲۲]. دلیل این تفاوت، وجود مقادیر زیاد نویسه‌ها و نواحی سیاه-سفید است. ترکیب تصاویر طبیعی و تصاویر اسناد، «نمادها»^{۳۲} را می‌سازد (شکل ۳). محاسبات از راه دور^{۳۳} و آموزش مجازی، شرایطی را بوجود می‌آورند که نیاز است مفاصل صفحه‌نمایش یک رایانه به عنوان تصویر مخابره شود. در [۲۰] نشان داده می‌شود که نمادها هم به روش‌های مختص به خود نیازمند هستند. مجموعه‌داده‌های کیفیت متخصص نمادها نیز برای تحقیقات اکت ارائه شده است [۲۳، ۲۰]. در این قسمت، نحوه‌ی آموزش و آزمون روش‌های اکت مبتنی بر بردار پشتیبان را مرور کردیم. همچنین، دیدیم که اکت محاسباتی، کاربردهای مختلفی برای انواع تصاویر دیجیتال دارد. در قسمت بعد، برخی جنبه‌های مربوط به یادگیری ماشین در آموزش SVRها را تحلیل می‌نماییم.

مجموعه‌داده، مجموعه‌ی تمامی تصاویر تخریب‌شده‌ای است، که متعلق به یک مرجع هستند. (یعنی یک منظره‌ی یکسان را نشان می‌دهند، منتها با تخریب‌های متفاوت و شدت‌های متفاوت). اگر ۸۰٪ صحنه‌ها را، به طور تصادفی، برای آموزش، و ۲۰٪ را برای آزمون کنار گذاشته باشیم، یک «افراز ۸۰-۲۰» انجام داده‌ایم.

روش مبتنی بر SVR، ابتدا از تصویر یک بردار ویژگی استخراج می‌کند. اگر در مجموعه‌داده، S نمونه، به شکل DS_{image} در (۲) داشته باشیم، استخراج‌کننده‌ی ویژگی، S بردار ویژگی محاسبه می‌کند که مجموعه‌ی $DS_{feature\ vector}$ در (۳) را تشکیل می‌دهند.

$$DS_{image} = \{(S_1 \text{ نمره‌ی انسانی}, \text{تصویر}_1), \dots, (S_S \text{ نمره‌ی انسانی}, \text{تصویر}_S)\} \quad (2)$$

$$DS_{vector\ feature} = \{(\vec{f}_1, S_1 \text{ نمره‌ی انسانی}), \dots, (\vec{f}_S, S_S \text{ نمره‌ی انسانی})\} \quad (3)$$

وقتی افراز تصادفی ۸۰-۲۰ را انجام دهیم، برخی اعضای $DS_{feature\ vector}$ مربوط به صحنه‌هایی هستند که در مجموعه‌ی آموزش^{۲۴} قرار گرفته‌اند و سایر اعضای آن مربوط به صحنه‌هایی هستند که در مجموعه‌ی آزمون^{۲۵} قرار گرفته‌اند. می‌توانیم عامل‌های^{۲۶} یک بردار پشتیبان را با استفاده از بردارهای ویژگی و نمرات انسانی مجموعه‌ی آموزش بهینه کرده و عملکرد آن را روی مجموعه‌ی آزمون ارزیابی کنیم. از آنجایی که انتخاب صحنه‌ها برای افراز به مجموعه‌های آموزش و آزمون به صورت تصادفی صورت می‌گیرد، می‌توانیم افراز تصادفی ۸۰-۲۰ را چند بار انجام دهیم و میانگی^{۲۷} عملکرد وایازنده^{۲۸} در این چند بار را به عنوان دقت الگوریتم در نظر بگیریم. اگر این کار را n دفعه انجام دهیم، اصطلاحاً می‌گویند که «اعتبار سنجی متقابل n -لایه^{۲۹}» انجام داده‌ایم. مقدار پیشنهاد شده برای n ، در اکت، ۱۰۰۰ است [۱۸]. یعنی هزار بار، افراز ۸۰-۲۰ را به صورت تصادفی انجام دهیم، آموزش و آزمون را اجرا نمائیم و میانگی عملکرد مدل در این هزار دفعه را به عنوان دقت الگوریتم گزارش کنیم.

البته باید این نکته را هم در نظر بگیریم که وایازنده‌ی مبتنی بر بردار پشتیبان، دو فراعامل^{۳۰}، به نام‌های $c(OST)$ و γ ، دارد. می‌توانیم نحوه‌ی کار وایازنده را به صورت رابطه‌ی (۴) بیان نماییم.

$$SVR_{c,\gamma}(\vec{W} \cdot \vec{f}^T + \vec{b}) = \text{پیش‌بینی} \quad (4)$$

که \vec{W} و \vec{b} وزن‌های بهینه شده با استفاده از داده‌های آموزشی و \vec{f} بردار ویژگی ورودی است. مقادیر c و γ قابل یادگیری نبوده و به صورت دستی تعیین می‌شوند. می‌توان بازه‌ای از مقادیر برای این دو فراعامل در نظر گرفت و اعتبار سنجی متقابل ۱۰۰-لایه را روی یک مجموعه‌داده برای تمامی این مقادیر انجام داد. زوج (c, γ) ای که بهترین عملکرد را داشته باشند، به عنوان فراعامل‌های انتخابی در نظر گرفته می‌شوند [۱۸].

۱.۲.۲ ارزیابی جداگانه‌ی عملکرد الگوریتم روی هر یک از تخریب‌ها

همانطور که در قسمت ۱.۲ گفتیم، برخی تخریب‌ها رایج هستند و مجموعه‌داده‌ها نمونه‌هایی با این تخریب‌ها دارند. برای اندازه‌گیری دقت

^{۲۴}training set

^{۲۵}test set

^{۲۶}معادل فارسی «parameter»

^{۲۷}median

^{۲۸}regressor

^{۲۹} n -fold cross validation

^{۳۰}hyper-parameter

^{۳۱}expressive

^{۳۲}معادل فارسی «screen content image». دقت شود که «ناگرفت»ها، «screen-shots»، یک حالت خاص نمادها هستند.

^{۳۳}remote computing

۳ آزمایش‌های پیشنهادی و نتایج

در این قسمت می‌بینیم که SVRهای آموزش دیده برای آکت، ممکن است مشکلاتی از نظر تعمیم‌پذیری داشته باشند. ابتدا یک آزمایش تعریف کرده و سپس به تحلیل عملکرد SVR می‌پردازیم.

۱.۳ آزمایش پیشنهادی

گفتیم که نامفادها از نوشته‌ها، تصاویر طبیعی و نگاشته‌ها^{۳۴} تشکیل شده‌اند (قسمت ۳.۲). یک راه برای ارزیابی کیفیت نامفادها، تحلیل جداگانه‌ی این نواحی است [۲۴]. اگر یک روش، مختص ارزیابی کیفیت تصاویر اسناد داشته باشیم، به نظر می‌رسد که ترکیب آن با روشی مختص تصاویر طبیعی، برای ارزیابی کیفیت نامفاد، شانس داشته باشد. $LBPSI(.,.)$ [۲۵]، یک روش مرجع کامل، برای آکت اسناد و نوشته‌هاست. به این ترتیب که اگر x و y مثل (۱) تعریف شده باشند، خواهیم داشت:

$$LBPSI = LBPSI(x, y) \quad (۵)$$

(نمره کیفیت تخمینی به وسیله LBPSI)

همینطور، $SQMS(.,.)$ [۲۶]، یک روش مرجع کامل برای نامفاد است. برای ارزیابی نواحی طبیعی، از $HaarPSI(.,.)$ [۲۷] و $GMSD(.,.)$ [۲۸] استفاده می‌کنیم که جنبه‌های مختلف ساختاری را می‌سنجند. از آنجایی که ارزیابی تباین^{۳۵} کار دشواری است [۲۹]، می‌توانیم از $VIF(.,.)$ [۱۱] هم استفاده کنیم که در این زمینه دقیق است.

با استفاده از این روش‌ها، می‌توانیم یک روش ترکیبی مرجع کامل [۳۰]، برای ارزیابی نامفادها بسازیم. $HaarPSI$ ، $LBPSI$ ، $SQMS$ ، $GMSD$ و VIF را برای x و y محاسبه می‌کنیم. در این صورت، ۵ عدد خواهیم داشت که می‌توانیم با آن‌ها یک بردار ویژگی تشکیل دهیم:

$$\vec{f} = [SQMS, LBPSI, HaarPSI, GMSD, VIF] \quad (۶)$$

با داشتن این بردار ویژگی، می‌توانیم یک SVR را برای نگاشت آن به نمره کیفیت آموزش دهیم. برای داده‌ی آموزشی، از مجموعه داده‌های کیفیت نامفاد، به نام‌های SIQAD [۲۰] و SCID [۲۳] استفاده می‌کنیم.

SIQAD شامل ۲۰ تصویر مرجع است، که هر کدام با ۷ نوع تخریب، تغییر یافته‌اند. هر تخریب در ۷ سطح اعمال شده است. بنابراین، $7 \times 7 = 49$ تصویر تخریب شده در SIQAD وجود دارد. علاوه بر تخریب‌های SIQAD، SCID دو تخریب دیگر را نیز به ۴۰ تصویر مرجع خود اعمال کرده است. هر تخریب در ۵ سطح شدت شبیه‌سازی شده است، لذا در SCID، ۱۸۰۰ تصویر تخریب شده موجود است. برای تصاویر تخریب شده در این دو مجموعه داده، نمرات کیفیت انسانی تهیه شده است.

علاوه بر (۶)، می‌توانیم ترکیب‌های دیگری، مثل $\vec{f}_1 = [SQMS, LBPSI, GMSD]$ را در نظر بگیریم. از آنجایی که $SQMS(.,.)$ ، خود یک روش مختص نامفاد است، تأثیر ترکیب سایر روش‌ها با آن را بررسی می‌کنیم. به این ترتیب، ترکیباتی که می‌توان امتحان کرد را در جدول ۱ خلاصه نموده‌ایم. برای مثال، \vec{f}_1 ، حاصل ترکیب $SQMS(.,.)$ و $HaarPSI(.,.)$ است. یعنی: $\vec{f}_1 = [SQMS, HaarPSI]$. یا $\vec{f}_6 = [SQMS, HaarPSI, VIF]$. بدیهی است که برای هر یک از بردارهای ویژگی \vec{f}_1 تا \vec{f}_{14} ، باید SVR مجزایی آموزش داده شود. در ادامه، عملکرد هر یک از این مدل‌ها بررسی می‌شوند.

³⁴graphics
³⁵contrast

جدول ۱: ترکیب‌های در نظر گرفته شده برای تشکیل بردار ویژگی. وجود \checkmark به معنی استفاده شدنِ نمره‌ی روشِ آن ستون است.

بردار	SQMS	HaarPSI	GMSD	VIF	LBPSI
\vec{f}_1	\checkmark	\checkmark			
\vec{f}_2	\checkmark		\checkmark		
\vec{f}_3	\checkmark			\checkmark	
\vec{f}_4	\checkmark				\checkmark
\vec{f}_5	\checkmark	\checkmark	\checkmark		
\vec{f}_6	\checkmark	\checkmark		\checkmark	
\vec{f}_7	\checkmark	\checkmark			\checkmark
\vec{f}_8	\checkmark		\checkmark	\checkmark	
\vec{f}_9	\checkmark		\checkmark		\checkmark
\vec{f}_{10}	\checkmark			\checkmark	\checkmark
\vec{f}_{11}	\checkmark	\checkmark	\checkmark	\checkmark	
\vec{f}_{12}	\checkmark	\checkmark	\checkmark		\checkmark
\vec{f}_{13}	\checkmark	\checkmark		\checkmark	\checkmark
\vec{f}_{14}	\checkmark		\checkmark	\checkmark	\checkmark

۲.۳ موفقیت ترکیب روش‌ها، در هر یک از مجموعه داده‌ها

بردارهای \vec{f}_1 تا \vec{f}_{14} را برای نمونه‌های SIQAD و SCID محاسبه کرده، و اعتبارسنجی متقابل ۱۰۰۰-لایه را برای یک یک آن‌ها انجام می‌دهیم. شکل ۴، بهره‌وری عملکرد مدل‌های ترکیبی، نسبت به استفاده‌ی مستقیم از SQMS را نشان می‌دهد. منظور از بهره‌وری، محاسبه‌ی مقدار زیر است:

$$f_i = |SROCC_{f_i}| - |SROCC_{SQMS}| \quad (۷)$$

به‌رووری f_i

که $i \in \{1, \dots, 14\}$ و $SROCC_{f_i}$ ، ضریب همبستگی [سپیرمن] محاسبه شده برای SVR است که روی f_i آموزش دیده باشد. همانطور که در قسمت ۱.۲ گفته شد، این ضریب همبستگی، بین نمرات الگوریتم و نمرات انسانی محاسبه می‌شود. $SROCC_{SQMS}$ هم دقت $SQMS(.,.)$ نشان می‌دهد. می‌بینیم که تمام مدل‌ها عملکرد بهتری نسبت به $SQMS(.,.)$ داشته‌اند.

۳.۳ عدم تعمیم‌پذیری به تخریب‌های مختلف

با افزایش دقت $SQMS(.,.)$ ، روی کل مجموعه داده، انتظار می‌رود که ترکیب‌ها بهبودی مشابه را روی تک تک تخریب‌ها نیز داشته باشند. میزان دقت مدل‌ها روی تخریب‌ها طبق ۱.۲.۲ محاسبه گردیده و بهره‌وری آن نسبت به $SQMS(.,.)$ در شکل ۵ نشان داده می‌شود. می‌بینیم که نه تنها همه‌ی مدل‌ها عملکرد مشابهی نداشتند، بلکه در بسیاری از موارد، بهره‌وری منفی بوده است. کاملاً منطقی است که یکی از علت‌ها را ضعف بردارهای ویژگی بدانیم. فارغ از این نظریه، مسئله‌ی دیگری وجود دارد که باید در نظر گرفته شود و آن فراعامل‌ها هستند. همانطور که در انتهای ۲.۲ گفته شد، بهینه‌سازی c و γ با اعتبارسنجی متقابل ۱۰۰-لایه صورت می‌گیرد:

$$(c^*, \gamma^*) = \text{opt}(r_c, r_\gamma, \text{dataset}_{80, \text{all}}, \text{dataset}_{20, \text{all}}, 100) \quad (۸)$$

رابطه‌ی ۸ مواردی که باید در نظر گرفته شوند را خلاصه می‌کند. c^* و γ^* فراعامل‌های بهینه هستند. تابع opt نتایج بهینه‌سازی را برمی‌گرداند.

شکل ۴ نتایج را برای آزمایش زیر نشان می‌دهد:

$$cv(dataset_{80,all}, \text{opt}(r_c, r_\gamma, dataset_{80,all}, dataset_{20,all}, 100), dataset_{20,all}, 1000) \quad (10)$$

که dataset در هر بار، یا SIQAD یا SCID است. برای ارزیابی روی هر یک از تخریب‌ها (شکل ۵)، رابطه (۱۰) به شکل زیر تغییر می‌کند:

$$cv(dataset_{80,all}, \text{opt}(r_c, r_\gamma, dataset_{80,all}, dataset_{20,all}, 100), dataset_{20,dst}, 1000) \quad (11)$$

که dataset مطابق (۱۰) و {dataset های موجود در dst}. اگر در آزمایشی که رابطه (۱۱) نشان می‌دهد، فراعامل‌ها را روی تخریب‌های متناظر بهینه کنیم، نتایج شکل ۶ حاصل می‌شوند. به طور رسمی، نتایج شکل ۶، حاصل آزمایش زیر هستند:

$$cv(dataset_{80,all}, \text{opt}(r_c, r_\gamma, dataset_{80,all}, dataset_{20,dst}, 100), dataset_{20,dst}, 1000) \quad (12)$$

پس می‌بینیم که علاوه بر تأثیر احتمالی بردارهای ویژگی، فراعامل‌ها نیز نقش مهمی در عملکرد SVR ها دارند.

۴.۳ عدم تعمیم‌پذیری به سایر مجموعه داده‌ها

برای بررسی عدم وابستگی یک مدل به صحنه‌های یک مجموعه داده، آن را روی یک مجموعه داده آموزش داده و روی مجموعه داده‌ی دیگر می‌آزمایند:

$$cv(SIQAD_{all,all}, \text{opt}(r_c, r_\gamma, SIQAD_{80,all}, SIQAD_{20,all}, 100), SCID_{all,all}, 1) \quad (13)$$

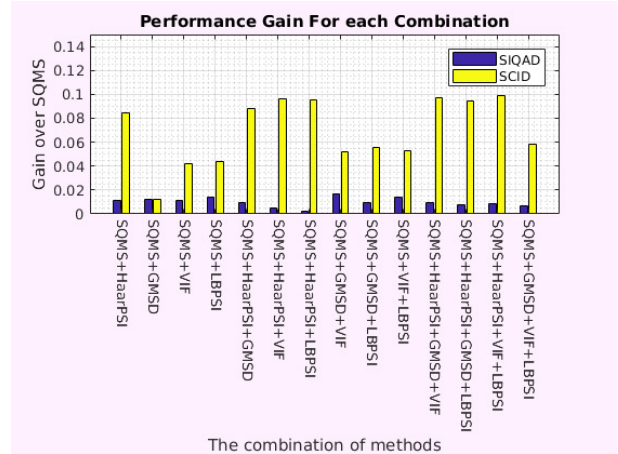
$$cv(SCID_{all,all}, \text{opt}(r_c, r_\gamma, SCID_{80,all}, SCID_{20,all}, 100), SIQAD_{all,all}, 1) \quad (14)$$

شکل ۷ نتایج را برای آزمایش‌های (۱۴) و (۱۳) نشان می‌دهد. می‌بینیم که بهبود مشاهده شده در آزمایش (۱۰) (با نتایج قابل مشاهده در شکل ۴) برای SIQAD اتفاق نمی‌افتد. به غیر از این مسئله، اگر مدل آموزش دیده را برای تک تک تخریب‌ها بیازمائیم، باز هم با افت عملکردی مشابه قسمت ۳.۳ مشابه می‌شویم. بیان رسمی این آزمایش‌ها در روابط (۱۵) و (۱۶) و نتایج آن‌ها در شکل ۸ ارائه شده‌اند.

$$cv(SCID_{all,dst}, \text{opt}(r_c, r_\gamma, SCID_{80,all}, SCID_{20,all}, 100), SIQAD_{all,dst}, 1) \quad (15)$$

$$cv(SIQAD_{all,dst}, \text{opt}(r_c, r_\gamma, SIQAD_{80,all}, SIQAD_{20,all}, 100), SCID_{all,dst}, 1) \quad (16)$$

dst در روابط بالا، عضو مجموعه‌ی تخریب‌هایی است که در SCID و SIQAD مشترک هستند. می‌بینیم که نتایج بهتر آزمون تعمیم‌پذیری برای مجموعه‌ی SCID هم به تک تک تخریب‌های آن قابل تعمیم نیست.



شکل ۴: بهره‌وری هر یک از روش‌های ترکیبی نسبت به $SQMS(.,.)$. برای هر ترکیب، بهره‌وری عملکرد بر روی دو مجموعه داده، با دو رنگ متفاوت نشان داده شده است.

r_c و r_γ ، به ترتیب، مجموعه‌ی مقادیری هستند که برای c و γ امتحان می‌شوند. مثلاً اگر $\{1, 10, 10^2, 10^3, 10^4, 10^5, 10^6\}$ باشد، r_c بهترین مقدار برای c از این مجموعه پیدا می‌شود. مقادیر مورد بررسی به صورت دستی انتخاب شده و تعدادشان به توان پردازشی در دسترس بستگی دارد. در هر بار آموزش و ارزیابی، نیاز به یک مجموعه‌ی آموزش و یک مجموعه‌ی آزمون داریم.

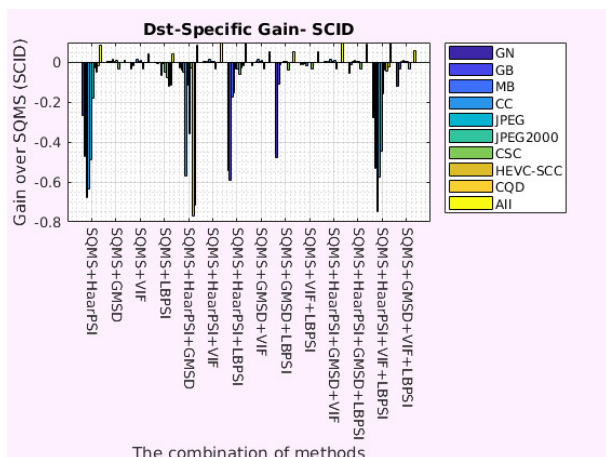
طبق رابطه‌ی ۸، مجموعه‌ی آموزش، $dataset_{80,all}$ است. این عبارت یک زیرمجموعه از نمونه‌های مجموعه داده‌ی $dataset$ را مشخص می‌نماید. این زیرمجموعه، عبارت است از همه‌ی (all) تخریب‌های موجود در زیرمجموعه‌ی آموزش افزاز ۸۰-۲۰. برخی نمونه‌ها از این عبارت، می‌توانند $SIQAD_{80,JPEG}$ و $SCID_{20,GB}$ باشند. که اولی، یعنی تمامی نمونه‌های SIQAD؛ که در افزاز ۸۰-۲۰، در مجموعه‌ی آموزش قرار گرفته‌اند؛ و تخریب آن‌ها از نوع JPEG است. دومی، نمونه‌هایی از SCID را مشخص می‌کند که در مجموعه‌ی آزمون قرار گرفته‌اند و تخریب آن‌ها از نوع GB^{۳۶} است. آخرین ورودی تابع opt هم، تعداد لایه‌های (دفعات) اعتبارسنجی متقابل است.

اگر به جای $\text{opt}(r_c, r_\gamma, SIQAD_{80,all}, SIQAD_{20,all}, 100)$ از $\text{opt}(r_c, r_\gamma, SIQAD_{80,all}, SIQAD_{20,JPEG}, 100)$ استفاده کنیم، نتایج متفاوت خواهند بود (شکل ۶). قبل از بررسی نتایج، یک رابطه مانند ۸ برای اعتبارسنجی متقابل قرارداد می‌کنیم، که نقش فراعامل‌ها را نیز در نظر بگیرد:

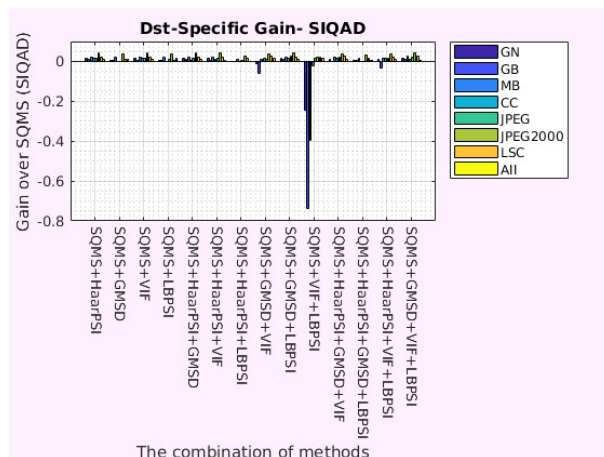
$$\begin{aligned} \text{دقت الگوریتم} = \\ cv(dataset_{80,all}, \text{opt}(r_c, r_\gamma, dataset_{80,all}, dataset_{20,all}, 100), dataset_{20,all}, 1000) \end{aligned} \quad (9)$$

تابع $cv(.,.,.,.)$ ، چهار ورودی دارد که عبارت‌اند از: مجموعه‌ی آموزش، فراعامل‌های بهینه، مجموعه‌ی آزمون و تعداد لایه‌های اعتبارسنجی متقابل. این تابع میانه‌ی SROCC های محاسبه شده در ۱۰۰۰ مرتبه را بر می‌گرداند.

^{۳۶} نوعی تار شدگی تصویر. برای شرح سایر تخریب‌ها به مقاله‌ی خود مجموعه داده‌ها رجوع شود.

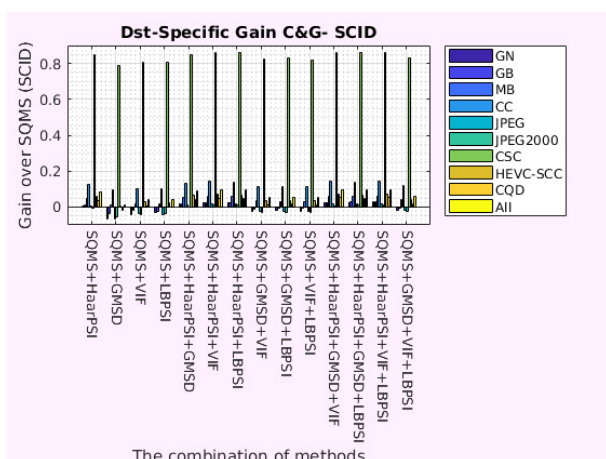


(ب)

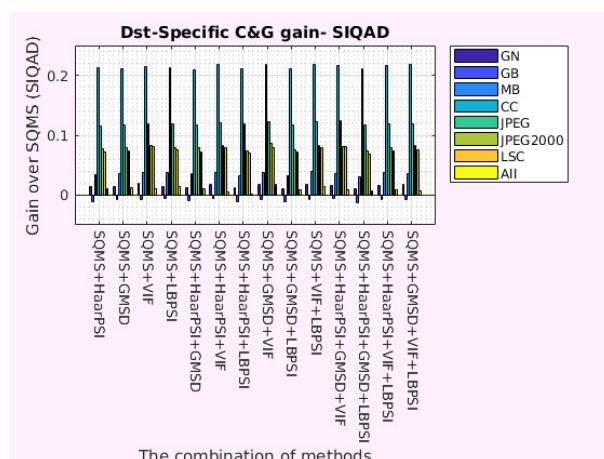


(آ)

شکل ۵: نتیجه‌ی ارزیابی به ازای هر تخریب. نتایج SIQAD در (آ) و SCID در (ب) قابل ملاحظه هستند. بهره‌وری به ازای هر تخریب، با رنگ متفاوتی مشخص شده است.



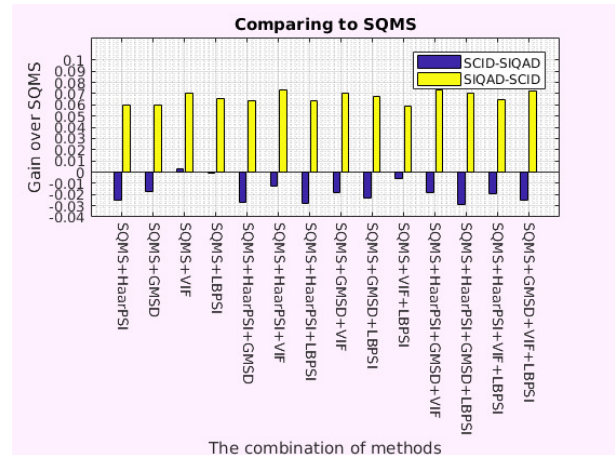
(ب)



(آ)

شکل ۶: نتیجه‌ی ارزیابی به ازای هر تخریب بعد از بهینه‌سازی فراعامل‌ها روی هر یک از تخریب‌ها. نتایج SIQAD در (آ) و SCID در (ب) قابل ملاحظه هستند.

- [5] E. C. Larson and D. M. Chandler, "Most apparent distortion: full-reference image quality assessment and the role of strategy," *Journal of electronic imaging*, vol.19, no.1, p.011006, 2010.
- [6] M. P. Eckert and A. P. Bradley, "Perceptual quality metrics applied to still image compression," *Signal processing*, vol.70, no.3, pp.177–200, 1998.
- [7] Z. Wang and A. C. Bovik, "Mean squared error: Love it or leave it? a new look at signal fidelity measures," *IEEE signal processing magazine*, vol.26, no.1, pp.98–117, 2009.
- [8] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol.13, no.4, pp.600–612, 2004.
- [9] H. R. Sheikh, M. F. Sabir, and A. C. Bovik, "A statistical evaluation of recent full reference image quality assessment algorithms," *IEEE Transactions on image processing*, vol.15, no.11, pp.3440–3451, 2006.
- [10] D. L. Ruderman and W. Bialek, "Statistics of natural images: Scaling in the woods," in *Advances in neural information processing systems*, pp.551–558, 1994.
- [11] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on image processing*, vol.15, no.2, pp.430–444, 2006.
- [12] J. Kim, H. Zeng, D. Ghadiyaram, S. Lee, L. Zhang, and A. C. Bovik, "Deep convolutional neural models for picture-quality prediction: Challenges and solutions to data-driven image quality assessment," *IEEE Signal processing magazine*, vol.34, no.6, pp.130–141, 2017.
- [13] X. Yang, F. Li, and H. Liu, "A survey of dnn methods for blind image quality assessment," *IEEE Access*, vol.7, pp.123788–123806, 2019.
- [14] P. Ye, J. Kumar, L. Kang, and D. Doermann, "Unsupervised feature learning framework for no-reference image quality assessment," in *2012 IEEE conference on computer vision and pattern recognition*, pp.1098–1105, IEEE, 2012.
- [15] V. Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- [16] پوریا چراغی، "ارائه‌ی یک روش بدون مرجع برای ارزیابی کیفیت تصاویر با تخریب چندگانه،" پایان‌نامه کارشناسی ارشد، گروه مهندسی برق و کامپیوتر دانشگاه خوارزمی، شهریور ۱۳۹۸.
- [17] D. M. Chandler, "Seven challenges in image quality assessment: past, present, and future research," *International Scholarly Research Notices*, vol.2013, 2013.



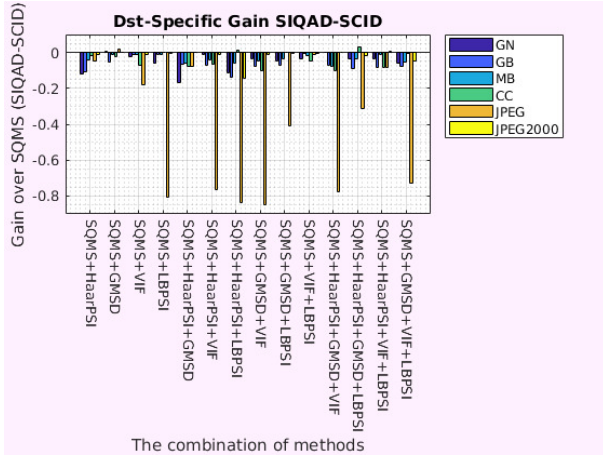
شکل ۷: نتایج آزمون تعمیم‌پذیری. ستون‌های زرد مربوط به آموزش روی SIQAD و آزمون روی SCID هستند. ستون‌های سُرْمه‌ای مربوط به آموزش و آزمون عکس هستند.

۴ جمع‌بندی

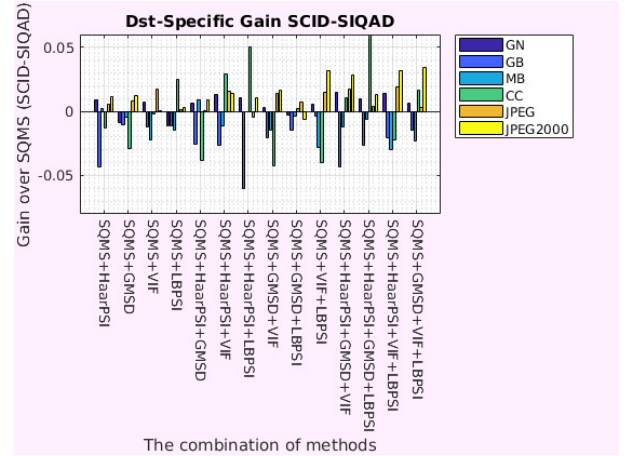
در این مقاله بردارهای پشتیبان را، طبق روش رایج، برای اکت آموزش دادیم. دیدیم که وقتی گزارش عملکرد به روش رایج انجام می‌شود، نقش فراعامل‌ها هم باید بازتاب گردد. یک فرمول‌بندی ارائه شد که این موارد را لحاظ نماید. تأثیر بهینه‌سازی فراعامل‌ها را مشاهده کردیم و دیدیم که به این مسئله در آزمایش‌های ارزیابی کیفیت پرداخته نمی‌شود. محدود بودن تصاویر مجموعه داده‌ها، نگرانی در مورد تعمیم‌پذیری مدل‌ها به تصاویر دنیای واقعی را بیشتر می‌کند. مخصوصاً وقتی که تخریب‌های مجموعه داده‌ها مصنوعی بوده و ممکن است با تخریب‌های طبیعی متفاوت باشند. شاید بهتر باشد که نحوه بهینه‌سازی فراعامل‌های مدل‌های یادگیری نیز در گزارش‌های عملکرد مقالات تشریح گردد. آزمایش‌های انجام شده با کدهای موجود در https://github.com/cheraaqee/fusion_iqa قابل شبیه‌سازی هستند.

مراجع

- [1] CISCO, "Cisco annual internet report," <https://www.cisco.com/c/en/us/solutions/executive-perspectives/annual-internet-report/index.html>, 2022. Accessed: 2022-July-15.
- [2] E. Allen and S. Triantaphillidou. *The manual of photography*. CRC Press, 2012.
- [3] P. Mohammadi, A. Ebrahimi-Moghadam, and S. Shirani, "Subjective and objective quality assessment of image: A survey," *arXiv preprint arXiv:1406.7799*, 2014.
- [4] G. Zhai and X. Min, "Perceptual image quality assessment: a survey," *Science China Information Sciences*, vol.63, no.11, pp.1–52, 2020.



(ب)



(آ)

شکل ۸: نتایج ارزیابی هر تخریب، وقتی مجموعه داده‌های آموزش و آزمون متفاوت‌اند. (آ): برای SCID → SIQAD و (ب) برای SIQAD → SCID

- [25] A. Alaei, D. Conte, M. Blumenstein, and R. Raveaux, "Document image quality assessment based on texture similarity index," in *2016 12th IAPR Workshop on Document Analysis Systems (DAS)*, pp.132–137, IEEE, 2016.
- [26] K. Gu, S. Wang, H. Yang, W. Lin, G. Zhai, X. Yang, and W. Zhang, "Saliency-guided quality assessment of screen content images," *IEEE Transactions on Multimedia*, vol.18, no.6, pp.1098–1110, 2016.
- [27] R. Reisenhofer, S. Bosse, G. Kutyniok, and T. Wiegand, "A haar wavelet-based perceptual similarity index for image quality assessment," *Signal Processing: Image Communication*, vol.61, pp.33–43, 2018.
- [28] W. Xue, L. Zhang, X. Mou, and A. C. Bovik, "Gradient magnitude similarity deviation: A highly efficient perceptual image quality index," *IEEE transactions on image processing*, vol.23, no.2, pp.684–695, 2013.
- [29] A. Shokrollahi, B. M.-N. Maybodi, and A. Mahmoudi-Aznavah, "Histogram modification based enhancement along with contrast-changed image quality assessment," *Multimedia Tools and Applications*, pp.1–22, 2020.
- [30] K. Okarma, P. Lech, and V. V. Lukin, "Combined full-reference image quality metrics for objective assessment of multiply distorted images," *Electronics*, vol.10, no.18, p.2256, 2021.
- [18] W. Xue, X. Mou, L. Zhang, A. C. Bovik, and X. Feng, "Blind image quality assessment using joint statistics of gradient magnitude and laplacian features," *IEEE Transactions on Image Processing*, vol.23, no.11, pp.4850–4862, 2014.
- [19] L. Kang, P. Ye, Y. Li, and D. Doermann, "Convolutional neural networks for no-reference image quality assessment," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.1733–1740, 2014.
- [20] H. Yang, Y. Fang, and W. Lin, "Perceptual Quality Assessment of Screen Content Images," *IEEE Transactions on Image Processing*, vol.24, no.11, pp.4408–4421, 2015.
- [21] A. Mittal, *Natural scene statistics-based blind visual quality assessment in the spatial domain*. Phd thesis, The University of Texas at Austin, 2013.
- [22] P. Ye and D. Doermann, "Document image quality assessment: A brief survey," in *2013 12th International Conference on Document Analysis and Recognition*, pp.723–727, IEEE, 2013.
- [23] Z. Ni, L. Ma, H. Zeng, J. Chen, C. Cai, and K. K. Ma, "ESIM: Edge Similarity for Screen Content Image Quality Assessment," *IEEE Transactions on Image Processing*, vol.26, no.10, pp.4818–4831, 2017.
- [24] Y. Zhang, D. M. Chandler, and X. Mou, "Quality Assessment of Screen Content Images via Convolutional-Neural-Network-Based Synthetic/Natural Segmentation," *IEEE Transactions on Image Processing*, vol.27, no.10, pp.5113–5128, 2018.