

پرسش ۲ - پیاده سازی یک سیستم طبقه بندی خودرو با استفاده از VGG16 و SVM

۱-۲ مقدمه

در سال‌های اخیر با رشد سریع حوزه‌ی یادگیری عمیق، مدل‌های زیادی برای تحلیل و طبقه‌بندی تصاویر معرفی شده‌اند. یکی از کاربردهای جالب این مدل‌ها، شناسایی و دسته‌بندی خودروها از روی تصاویر آن‌هاست. این کار می‌تونه در زمینه‌هایی مثل سیستم‌های کنترل ترافیک، پارکینگ‌های هوشمند، تحلیل رفتار راننده‌ها و حتی بازارهای فروش خودرو نقش مهمی داشته باشه.

هدف این پروژه طراحی یک سیستم طبقه‌بندی خودرو بوده که بتونه بین مدل‌های مختلف خودروهای تویوتا تفاوت قائل بشه. برای این کار از مدل VGG16 استفاده شده تا ویژگی‌های تصویری (features) از عکس‌ها استخراج بشن، و بعد این ویژگی‌ها به یک مدل SVM داده شدن تا کار طبقه‌بندی (classification) رو انجام بده.

در کنار این مدل ترکیبی (VGG + SVM)، عملکرد چند مدل دیگه هم بررسی شده؛ از جمله AlexNet و یک مدل ساده‌ی CNN که طراحی شده تا با بقیه مقایسه بشه. همچنین، نسخه‌های مختلفی از مدل SVM با کرنل‌های Linear و RBF آزمایش شدن تا ببینیم کدوم ترکیب عملکرد بهتری داره. در ادامه‌ی این گزارش مراحل آماده‌سازی داده‌ها، استخراج ویژگی، آموزش مدل‌ها و تحلیل نتایج آورده شده. هدف اینه که ببینیم کدوم روش برای طبقه‌بندی خودروهای تصویری بهتر جواب می‌ده و چرا.

استفاده از نسخه GPU-SVM با کتابخانه RAPIDS

در ابتدای کار تصمیم بر این بود که برای پیاده‌سازی مدل Support Vector Machine (SVM) از کتابخانه‌ی معروف scikit-learn استفاده شود. اما به دلیل اینکه استخراج ویژگی‌ها با استفاده از مدل VGG16 منجر به تولید داده‌هایی با ابعاد نسبتاً بالا شده بود، زمان آموزش مدل SVM روی CPU به‌طور محسوسی زیاد بود. علاوه بر این، در برخی مراحل از اجرای مدل، به دلیل حجم بالای ویژگی‌ها، با مشکلات حافظه نیز مواجه شدیم.

در تلاش برای تسریع روند آموزش و بهینه‌سازی مصرف منابع، تصمیم گرفتیم از نسخه‌ی GPU-based الگوریتم SVM استفاده کنیم. در این مرحله، کتابخانه‌ی cuML که بخشی از پلتفرم RAPIDS است، مورد استفاده قرار گرفت. این کتابخانه از GPU برای شتاب‌دهی به محاسبات استفاده می‌کند و نسخه‌ای از SVM را ارائه می‌دهد که بسیار سریع‌تر از نسخه‌ی مرسوم در scikit-learn اجرا می‌شود.

فرآیند نصب RAPIDS و راه‌اندازی آن در محیط Google Colab به راحتی با استفاده از اسکریپت آماده‌ی این پروژه در GitHub انجام شد. بعد از نصب و بارگذاری داده‌ها، مدل SVM با کرنل RBF آموزش داده شد و پیش‌بینی روی داده‌های تست انجام گرفت. خوشبختانه نتایج نهایی از نظر دقت، precision، recall و F1 Score با نتایج مدل CPU کاملاً قابل مقایسه بود، با این تفاوت که زمان آموزش به شکل چشم‌گیری کاهش پیدا کرد.

این تجربه نشان داد که در پروژه‌هایی که با حجم داده‌ی بالا سروکار داریم، استفاده از راهکارهای GPU مانند cuML می‌تواند گزینه‌ی بسیار خوبی برای افزایش کارایی باشد.

مدیریت حافظه در مرحله استخراج ویژگی

در مراحل ابتدایی پروژه، زمانی که تصمیم گرفته شد از مدل‌های VGG16 و AlexNet برای استخراج ویژگی‌های تصویری استفاده شود، یکی از چالش‌های اصلی، محدودیت منابع سخت‌افزاری در محیط Google Colab بود. با توجه به این که هر تصویر باید از طریق مدل عبور می‌کرد تا ویژگی‌های آن استخراج شود، و از طرفی تعداد تصاویر در هر کلاس نیز قابل توجه بود، اجرای همزمان تمام داده‌ها باعث کرش کردن runtime در Colab به دلیل کمبود RAM می‌شد.

در ابتدا برای حل این مشکل، داده‌ها به صورت دستی به چند batch کوچک‌تر تقسیم شدند و هر بار فقط یک بخش از آن‌ها از مدل عبور داده می‌شد. این روش اگرچه موقتی بود و امکان ادامه کار را فراهم می‌کرد، اما هم از نظر زمانی پرهزینه بود و هم احتمال بروز خطاهای انسانی در مدیریت داده‌ها را افزایش می‌داد.

در ادامه‌ی پروژه، با آشنایی بیشتر با کتابخانه‌ی PyTorch، از ساختار استاندارد آن یعنی Dataset و DataLoader استفاده شد. با تعریف یک کلاس Dataset سفارشی برای مجموعه تصاویر پروژه، و سپس استفاده از DataLoader برای ارسال داده‌ها به مدل به صورت batch، هم مشکل کمبود حافظه برطرف شد و هم فرآیند استخراج ویژگی ساختارمند و قابل مدیریت‌تر شد.

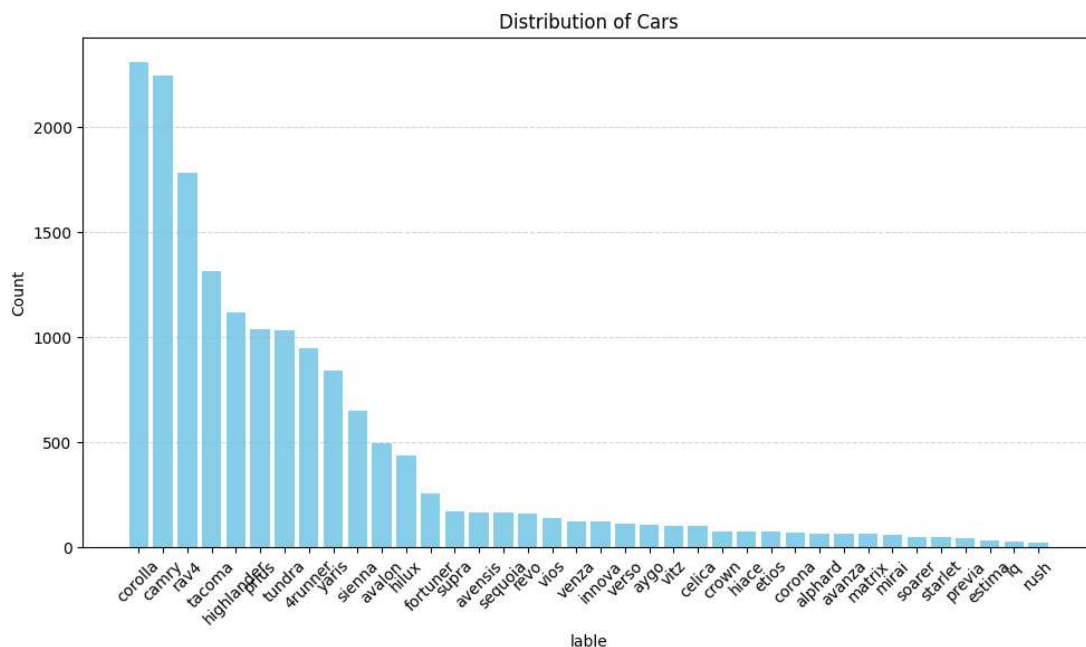
استفاده از DataLoader این امکان را داد تا بدون نگرانی از حجم زیاد داده‌ها، ویژگی‌های تمامی تصاویر به صورت batch و بهینه از مدل CNN استخراج شود. این تجربه‌ی فنی یکی از نقاط قوت اصلی پروژه بود که باعث شد مراحل پردازش و آموزش با اطمینان و پایداری بیشتری پیش بروند.

۲-۲ پیش پردازش داده ها

دیتا های مربوط به خودروهای تویوتا به طور کامل دانلود شدند و چیزی که حائز اهمیت است این است که در سایت ذکر شده است که دیتا ها با دقت ۹۵ درصد لیبل گذاری شده است.

حجم دانلود شده حدود ۲ گیگابایت است که از طریق محیط Google Colab دانلود و در دایرکتوری مربوط به Google Drive ذخیره شدند تا در زمان هایی که ریسورس های colab بدلیل محدودیت گرفته می شود همچنان در دسترس و قابل استفاده باشند.

تعداد عکس های مربوط به هر کلاس خودرو یا توزیع دیتاست از طریق نمودار زیر قابل مشاهده است.



شکل ۲-نمودار توزیع تصاویر خودرو های دیتاست

همانطور که صورت سوال خواسته بود ۱۰ تا از مدل خودرو به صورت دلخواه ، که در اینجا از ۱۰ خودرو با بیشترین توزیع را انتخاب کردیم و لیبل های آن ها را به مقادیر عددی ۰ تا ۹ تغییر داده شد.

برای شروع، مجموعه ای از تصاویر مربوط به ۱۰ مدل مختلف خودروهای تویوتا انتخاب شد. این مدل ها شامل corolla, camry, rav4, tacoma, highlander, prius, tundra, runner4, yaris و sienna هستند که به ترتیب با لیبل های عددی ۰ تا ۹ مشخص شدند.

اولین مرحله در پردازش داده ها، resize کردن تمام تصاویر به اندازه ی ۲۲۴ در ۲۲۴ پیکسل بود. این اندازه همون سایزیه که مدل هایی مثل VGG16 انتظار دارن. بعد از تغییر اندازه، مقدار پیکسلی تصاویر به صورت

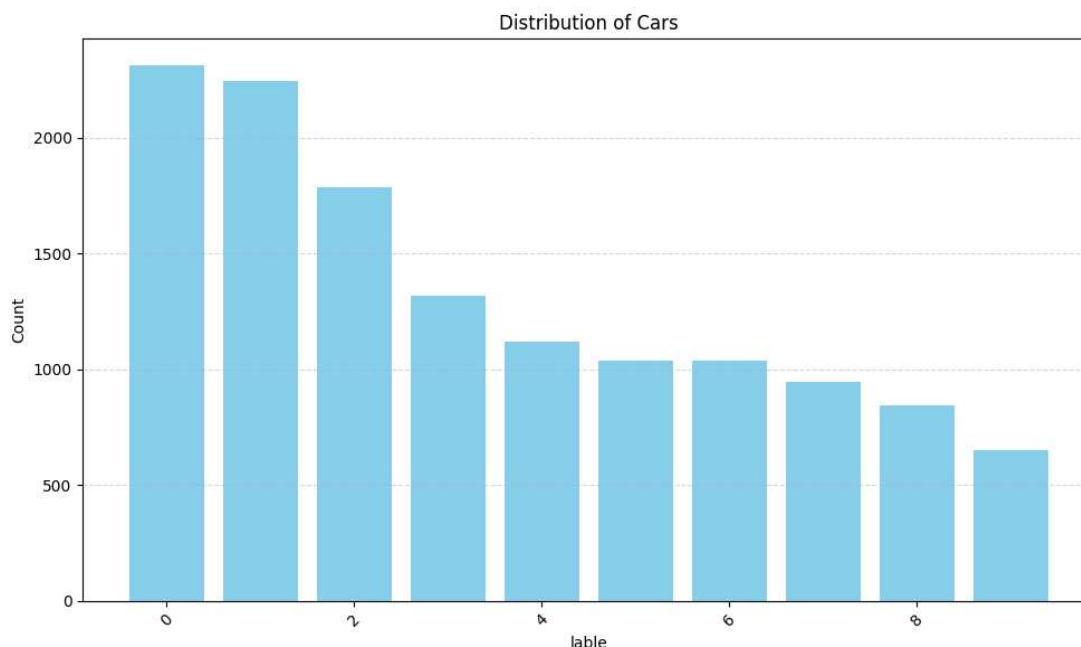
عددی بین صفر و یک نرمال سازی شد. این کار کمک می‌کند تا مدل بتواند بهتر یاد بگیرد و از نوسانات بزرگ عددی در ورودی جلوگیری بکند.

نکته‌ی مهمی که بعد از بررسی دیتاست مشخص شد، این بود که تعداد تصاویر در کلاس‌های مختلف برابر نیست. مثلاً بعضی مدل‌ها چند صد تصویر داشتن ولی بعضی دیگر فقط چند ده تصویر. این عدم تعادل می‌تونه باعث بشه مدل فقط روی کلاس‌های پرتکرار تمرکز کنه و بقیه رو نادیده بگیره.

برای حل این مشکل از روش Oversampling همراه با Data Augmentation استفاده شد. یعنی برای کلاس‌هایی که تصویر کم‌تری داشتن، با ایجاد تغییرات جزئی روی تصاویر موجود، تصاویر جدید تولید شد. این تغییرات شامل چرخش تصویر، فلیپ کردن، تغییرات رنگ و نور و... بودن. این کار هم باعث افزایش تعداد داده‌ها شد و هم باعث شد مدل با تنوع بیشتری از تصاویر روبه‌رو بشه که به یادگیری بهتر کمک می‌کنه.

در نهایت، با اضافه شدن تصاویر جدید، دیتاست به صورت نسبتاً متوازن بین کلاس‌ها تقسیم شد و مدل می‌تونه آموزش عادلانه‌تری ببینه.

طبق مقاله اندازه‌ی تصاویر را به ۲۲۴ در ۲۲۴ تغییر دادیم و پس از تغییر سایز به وسیله نرمال سازی مقادیر پیکسلی را به بین ۰ و ۱ تبدیل کردیم.



شکل ۳ - توزیع کلاس ها پس از ریسایز و تغییر لیبل ها

همانطور که در شکل بالا مشاهده میشود توزیع تعداد تصاویر در کلاس ها ، داری توزیع بالانس نیست و به اصطلاح داده ها نامتوازن توزیع شده اند.

برای بالانس کردن دیتاست میتوان از ۳ تکنیک استفاده کرد:

- UnderSampling : این روش سریع و ساده است ولی استفاده از این روش میتواند باعث از بین رفتن اطلاعات مفید کلاس های پرتکرار شود.
- Oversampling : در این روش برخلاف روش قبل ، اطلاعات یا داده ها حفظ میشوند. ولی احتمال overfitting را افزایش میدهد.
- Oversampling با Data Augmentation : این روش تنوع بالا ، حفظ تعادل را در داده ها به ارمغان می آورد و همچنین به تعمیم مدل کمک میکند. عیب این روش نیازمندی آن به زمان محاسباتی بیشتر نسبت به دو روش دیگر است.

Original



Augmented 1



Augmented 2



Augmented 3

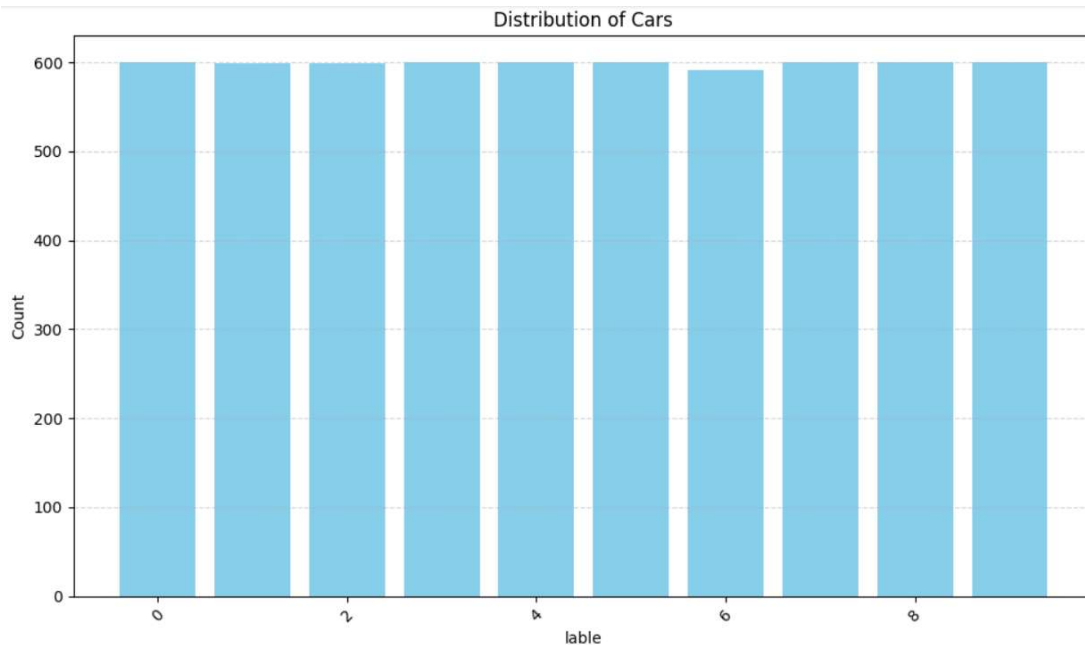


Augmented 4



شکل ۴ - نمایش تصویری نحوه Augmentaion

ما برای برقراری توازن در داده ها از روش DataAugmentation با Oversampling استفاده کردیم. و دلیل این انتخاب ، وجود زمان کافی برای محاسبات بود .



شکل ۵ - توزیع کلاس ها بعد از نرمالسازی

3-2 استخراج ویژگی ها

در این مرحله مدل های VGG16 و AlexNet بارگزاری شدند و با استفاده از لایه های کانولوشنی آن ها ویژگی مربوط به هر تصویر در هر مدل را به صورت جداگانه بدست آوردیم و بعد از یک بعدی سازی (flatten) آن ها را در درایو ذخیره کردیم.

بعد از آماده سازی و نرمال سازی تصاویر، مرحله ی بعدی مربوط به استخراج ویژگی هاست. هدف از این مرحله اینه که قبل از آموزش مدل طبقه بندی کننده (مثل SVM)، ویژگی های معنادار و فشرده ای از تصاویر به دست بیاریم که نماینده ی خوبی از محتویات اون تصویر باشن.

برای این کار از دو مدل معروف در حوزه ی بینایی ماشین استفاده شد VGG16 و AlexNet. هر دو این مدل ها از نوع convolutional neural network (CNN) هستن و روی دیتاست ImageNet آموزش دیدن. به همین دلیل، بدون نیاز به آموزش مجدد، می تونن ویژگی های خیلی خوبی از تصاویر استخراج کنن. استفاده از این مدل ها به صورت از پیش آموزش دیده (pretrained) باعث می شه هم در زمان صرفه جویی بشه، هم نتایج بهتری بگیریم.

در این پروژه، فقط از بخش کانولوشنی (convolutional layers) این مدل ها استفاده شد. یعنی بخش fully connected در انتهای مدل وجود داره و برای classification اصلی مدل طراحی شده، کنار گذاشته شد. دلیل این کار اینه که ما خودمون قصد داشتیم ویژگی های خروجی رو جداگانه به SVM بدیم و classification رو با اون انجام بدیم.

ویژگی‌های خروجی از این لایه‌ها معمولاً به شکل آرایه‌هایی چندبعدی هستن که نماینده‌ی الگوهای مختلف تصویرن (مثل لبه‌ها، بافت‌ها، شکل‌ها و...). این آرایه‌ها بعد از استخراج، به یک بردار یک‌بعدی تبدیل شدن تا بشه اونا رو به مدل SVM داد. این کار رو اصطلاحاً flatten کردن می‌گن.

در نهایت، بردارهای ویژگی استخراج‌شده برای همه‌ی تصاویر ذخیره شدن تا در مراحل بعدی برای آموزش و ارزیابی مدل‌ها استفاده بشن. چون استخراج ویژگی نسبتاً زمان‌بره، این کار فقط یک‌بار انجام شد و نتایج در حافظه‌ی دائمی (مثل Google Drive) ذخیره شدن.

۴-۲ آموزش و ارزیابی مدل

۱. در این مرحله به وسیله ویژگی‌های استخراج شده در مرحله قبل لایه‌های fully connected مربوط به مدل‌های VGG16 و AlexNet آموزش داده شدند. و نتایج به صورت زیر قابل مشاهده اند.

برای ارزیابی عملکرد مدل‌های طبقه‌بندی، معمولاً از چند معیار رایج استفاده می‌شود که مهم‌ترین آن‌ها عبارت‌اند از Precision، Recall و F1 Score. این معیارها به ما کمک می‌کنند تا علاوه بر بررسی دقت کلی مدل (Accuracy)، بتوانیم عملکرد آن را در شناسایی صحیح نمونه‌های هر کلاس به صورت جزئی‌تر بررسی کنیم.

Precision (دقت)

Precision نشان می‌دهد از بین تمام تصاویری که مدل به عنوان «یک کلاس خاص» (مثلاً مدل Corolla) پیش‌بینی کرده، چه درصدی از آن‌ها واقعاً متعلق به همان کلاس بوده‌اند.

$$Precision = \frac{TP}{TP + FP}$$

فرمول ۴-۱

به زبان ساده، اگر مدل در پیش‌بینی کلاس‌ها خیلی سخت‌گیر باشد و فقط در مواقعی که مطمئن است پیش‌بینی کند، معمولاً precision بالاتری خواهد داشت. اما ممکن است برخی نمونه‌ها را اصلاً شناسایی نکنند. در مسئله‌ی ما، اگر مدل فقط خودروهایی را به عنوان «RAV4» تشخیص دهد که کاملاً شبیه به نمونه‌های آموزشی RAV4 هستند، precision بالا می‌رود ولی ممکن است بعضی RAV4های واقعی را نادیده بگیرد.

بازخوانی (Recall)

Recall نشان می‌دهد مدل از بین تمام تصاویری که واقعاً به یک کلاس خاص تعلق داشته‌اند، چه تعداد را به درستی شناسایی کرده است.

$$Recall = \frac{TP}{TP + FN}$$

فرمول ۲-۰

معمولاً با افزایش میزان precision باعث کاهش Recall می‌شود و برعکس.

اگر مدل تمایل داشته باشد تا بیشتر پیش‌بینی کند (حتی با اطمینان کمتر)، معمولاً recall بالا می‌رود. این یعنی احتمال اینکه مدل، تصویر یک خودروی RAV4 را اشتباهی مثلاً «Highlander» تشخیص دهد (False Negative)، کمتر خواهد شد.

رابطه‌ی بین Recall و Precision

بین این دو معیار یک trade-off طبیعی وجود دارد؛ افزایش یکی معمولاً باعث کاهش دیگری می‌شود. به همین دلیل، برای رسیدن به یک ارزیابی متعادل، از معیار ترکیبی F1 Score استفاده می‌شود که میانگین هارمونیک این دو مقدار است:

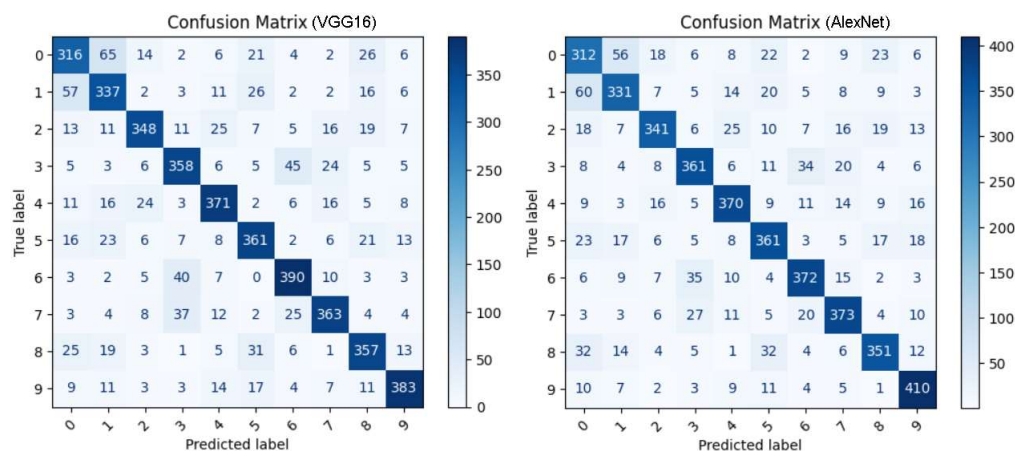
$$f1 - score = \frac{Precision * Recall}{Precision + Recall} * 2$$

فرمول ۳-۰

F1 Score زمانی بیشترین مقدار خود را دارد که precision و recall به یک اندازه بالا باشند. در پروژه‌ی ما، این معیار کمک می‌کند تا عملکرد مدل‌ها را به صورت کلی‌تری ارزیابی کنیم و مدل‌هایی که هم کم‌ترین اشتباه را دارند و هم بیشترین شناسایی صحیح را انجام می‌دهند، بهتر تشخیص داده شوند.

مدل	Accuracy	Precision	Recall	F1 Score
AlexNet	0.7752	0.7759	0.7751	0.7751
VGG16	0.7756	0.7768	0.7756	0.7758

جدول ۴- ارزیابی مدل‌های AlexNet و VGG16



شکل 6 - ماتریس در هم ریختگی مربوط به دو مدل پس آموزش

هر دو مدل در تشخیص برخی کلاس‌ها عملکرد قابل توجهی داشتند. به خصوص کلاس‌های ۳، ۴، ۶، ۷ و ۹ به خوبی از سوی هر دو مدل شناسایی شدند:

- کلاس ۶ در VGG16 با ۳۹۰ پیش‌بینی صحیح و در AlexNet با ۳۷۲ مورد، از بالاترین دقت‌ها برخوردار بود.
 - کلاس ۴ در هر دو مدل عملکرد بالایی داشت (371 در VGG و ۳۷۰ در AlexNet).
 - کلاس ۳ و ۷ نیز با مقادیر پیش‌بینی صحیح بالا، نشان‌دهنده‌ی توانایی مدل‌ها در شناسایی این خودروها هستند.
 - کلاس ۹ بیشترین دقت را در هر دو مدل داشت، به‌ویژه در AlexNet که با ۴۱۰ پیش‌بینی صحیح، بیشترین مقدار در کل ماتریس‌ها بود.
- در مقابل، هر دو مدل در تشخیص کلاس‌های ۰ و ۱ دچار چالش بودند:

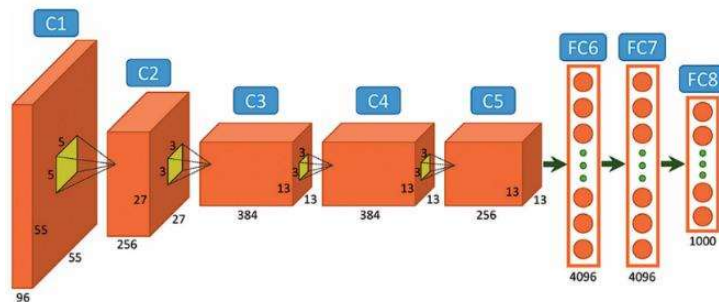
- در VGG16:
 - کلاس ۰ به اشتباه ۶۵ بار به عنوان کلاس ۱ پیش‌بینی شده است.
 - کلاس ۱ نیز ۵۷ بار به اشتباه به عنوان کلاس ۰ طبقه‌بندی شده است.
- در AlexNet:
 - کلاس ۰ با ۵۶ مورد اشتباه در کلاس ۱
 - کلاس ۱ با ۶۰ مورد اشتباه در کلاس ۰

این تقارن در اشتباهات نشان می‌دهد که تمایز بین این دو کلاس برای هر دو مدل دشوار بوده و به احتمال زیاد شباهت‌های ظاهری بین خودروهای این دو کلاس (مثلاً فرم کلی، رنگ یا زاویه عکس) باعث بروز این خطا شده است.

۲. اصلی‌ترین تفاوت بین AlexNet و VGG Net فارغ از نتایج:

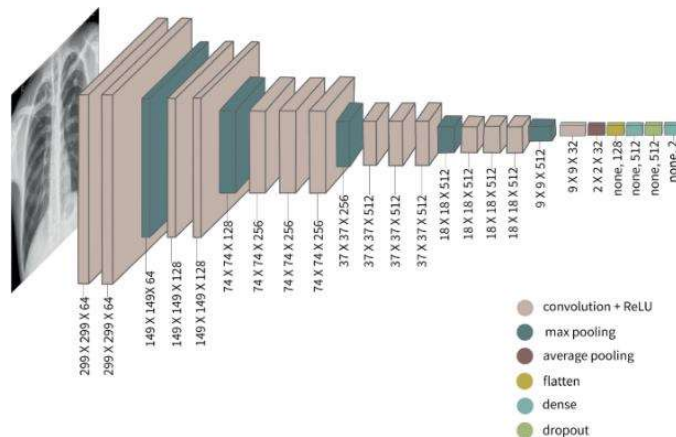
میتوان گفت که اصلی‌ترین تفاوت بین طراحی و عمق لایه‌های کانولوشن است. شبکه AlexNet با ۵ لایه کانولوشن کم عمق تر نسبت به شبکه VGG16 با ۱۳ لایه کانولوشنی است.

در مورد تفاوت این دو مدل میتوان به این نکته توجه کرد که مدل AlexNet از فیلترهای متنوع و بزرگ استفاده میکند ولی VGG16 از یک نوع فیلتر ۳ در ۳ فقط استفاده میکند.



شکل ۷ - تصویر از مدل AlexNet

مدل VGG دارای حجم بیشتر و سرعت پایین اما AlexNet دارای حجم کمتر و سرعت بالا تر ولی از لحاظ دقت مدل VGG عملکرد بهتری دارد.



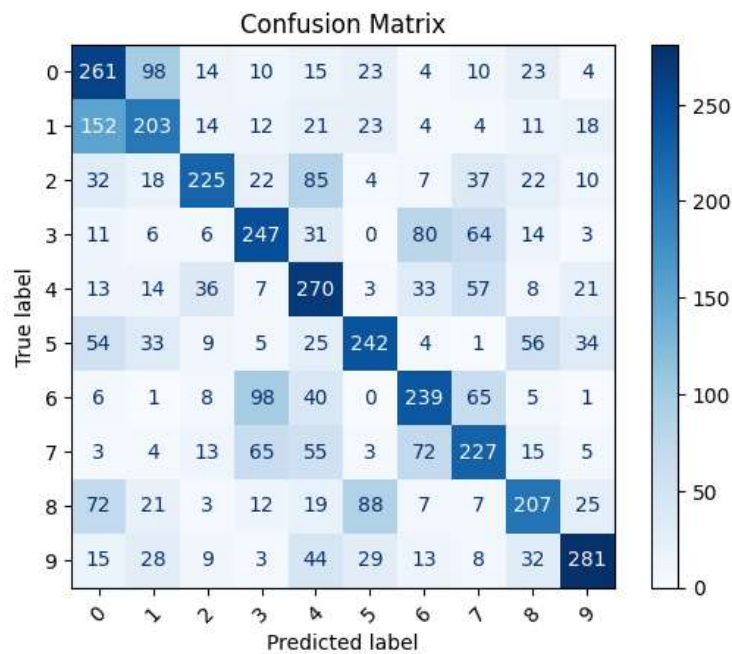
شکل ۸ - تصویر شماتیک مدل VGG16

۳. یک مدل CNN دلخواه در این بخش طراحی شد که دارای ۴ بلوک کانولوشن ، و ۳ لایه Fully Connected است ، لایه های کانولوشنی دارای تابع فعال ساز Relu و تکنیک BatchNorm برای تسریع و پایداری آموزش ، و MaxPooling برای کاهش ابعاد است. و به منظور کاهش Overfitting از Dropout استفاده شده است.

پس از آموزش نتایج به صورت زیر قابل مشاهده است:

مقدار	معیار
0.5198	Accuracy
0.5324	Precision
0.5198	Recall
0.5218	F1 Score

جدول ۵ - ارزیابی مدل CNN



شکل 9- ماتریس در هم ریختگی مربوط به مدل CNN

ماتریس درهم ریختگی مربوط به مدل ساده‌ی CNN طراحی شده در پروژه، ضعف‌هایی در دقت پیش‌بینی را نشان می‌دهد، که عمدتاً به دلیل ساده بودن ساختار آن است.

کلاس‌های ۰، ۱ و ۹ با دقت بهتری نسبت به سایر کلاس‌ها طبقه‌بندی شده‌اند.

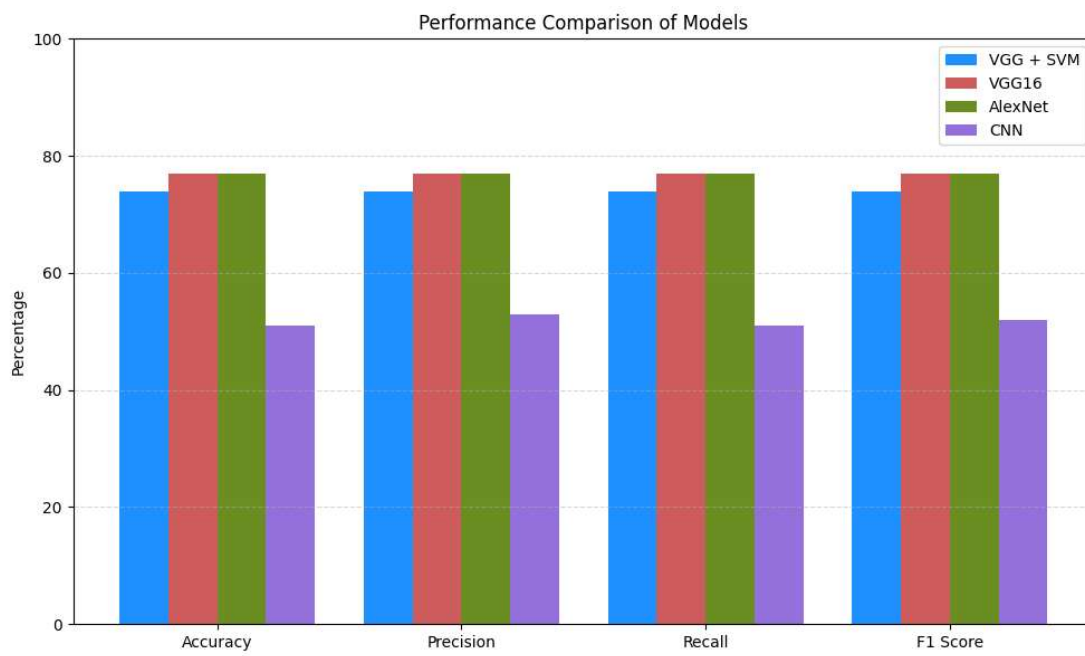
- کلاس ۹: ۲۸۱ پیش‌بینی صحیح
 - کلاس ۱۰: ۲۶۱ پیش‌بینی صحیح
 - کلاس ۱: ۲۰۳ پیش‌بینی درست (هرچند با اشتباه زیاد)
 - کلاس‌های ۲، ۳، ۶، ۷ و ۸ پیش‌بینی‌های بسیار پراکنده‌ای دارند.
 - کلاس ۶، برای مثال فقط ۲۳۹ مورد به درستی تشخیص داده شده و باقی در کلاس‌های دیگر پراکنده شده‌اند.
 - کلاس ۲ و ۳ نیز به‌طور مکرر با سایر کلاس‌ها اشتباه گرفته شده‌اند.
- مدل ساده‌ی CNN در تفکیک مرزهای بین کلاس‌ها ضعیف‌تر عمل کرده است. این موضوع نشان می‌دهد که برای دیتاست‌های واقعی و پیچیده، استفاده از مدل‌های از پیش آموزش‌دیده یا عمیق‌تر، تأثیر چشم‌گیری در بهبود دقت دارند.
۴. مدل پیشنهادی مقاله به این صورت تعریف شده است که با ترکیب VGG و SVM به این‌صورت که لایه‌های کانولوشنی مدل VGG وظیفه خروج feature هارا به عهده داشته باشد و SVM وظیفه Classification این feature هارا به عهده بگیرد.
- برای این منظور یک مدل VGG دانلود و با استفاده از بخش کانولوشنی آن ویژگی هر تصویر را استخراج کردیم و با استفاده از آن مدل SVM را آموزش دادیم.
- نتایج رو داده‌های تست به صورت زیر قابل مشاهده هستند:

مقدار	معیار
0.7453	Accuracy
0.7464	Precision
0.7453	Recall
0.7457	F1 Score

جدول ۶ - ارزیابی مدل VGG+SVM

۵-۲ تحلیل نتایج

معیارهای عملکردی محاسبه شد هاز مرحله‌های قبل به صورت جدول زیر قابل مشاهده است.



شکل ۱۰- نمایش گرافیکی عملکرد به دست آمده توسط مدل های مختلف برای مجموعه داده تویوتا

بعد از آموزش مدل ها و اجرای مرحله ی ارزیابی روی داده های تست، مشخص شد که مدل های CNN کامل، یعنی AlexNet و VGG16، بهترین عملکرد رو داشتن. دقت (accuracy) این دو مدل حدود ۷۷ درصد بود که از سایر روش های آزمایش شده بهتر بود.

مدل VGG16 با دقت ۰.۷۷۵۶ و AlexNet با دقت ۰.۷۷۵۲ تونستن تصاویر خودرو رو با کیفیت مناسبی طبقه بندی کنن. این دو مدل از ابتدا تا انتها شامل لایه های کانولوشنی و fully connected بودن و روی ویژگی هایی که خودشون استخراج کرده بودن آموزش دیدن.

در کنار این دو مدل، ترکیب VGG16 با SVM هم آزمایش شد؛ به این صورت که ویژگی ها از لایه های کانولوشنی VGG استخراج شدن و سپس به یک مدل SVM داده شدن. در این حالت، با استفاده از کرنل Linear، دقت به حدود ۰.۷۴ رسید و با کرنل RBF کمی بهتر شد (حدود ۰.۷۵). هرچند این مقادیر نزدیک به مدل های CNN کامل بودن، ولی در نهایت ترکیب VGG + SVM نتونست عملکردی بهتر از خود مدل VGG16 یا AlexNet ارائه بده.

این موضوع نشون می ده که گاهی اوقات وقتی مدل CNN به صورت end-to-end آموزش ببینه، بهتر می تونه ارتباط بین ویژگی ها و کلاس ها رو یاد بگیره. در مقابل، وقتی feature extraction و classification از هم جدا بشن، به مقدار از هماهنگی بین این دو بخش کم می شه و ممکنه تأثیر منفی بزاره، به خصوص وقتی feature ها خیلی پیچیده ان.

مدل CNN ساده‌ای هم که در این پروژه طراحی شد) با ۴ بلاک کانولوشن و چند لایه‌ی fully connected (عملکرد ضعیف‌تری داشت و به دقتی حدود ۵۲ درصد رسید. این مدل بدون استفاده از pretraining بود، و با توجه به تعداد نسبتاً محدود تصاویر، نتوانست به اندازه‌ی مدل‌های بزرگ‌تر یاد بگیرد. در مجموع:

- VGG16 و AlexNet بهترین عملکرد رو داشتن.
 - VGG + SVM نتایج نسبتاً خوبی داشت ولی نتوانست از مدل‌های اصلی CNN جلو بزنه.
 - مدل ساده‌ی طراحی‌شده برای تمرین و تست خوب بود ولی به بهینه‌سازی بیشتری نیاز داره.
- این نتایج نشون می‌دن که انتخاب بین مدل‌های end-to-end و مدل‌های ترکیبی feature extractor + classifier باید با توجه به نوع داده‌ها، حجم آن‌ها و هدف نهایی پروژه انجام بشه.
- با مقایسه‌ی چهار مدل، مشاهده می‌شود که VGG16 و AlexNet در تشخیص کلاس‌ها عملکرد بهتری دارند، به‌ویژه در کلاس‌هایی که تفاوت‌های بصری مشخص‌تری دارند. اما حتی این مدل‌ها نیز در تمایز بین برخی کلاس‌ها (مثلاً کلاس ۰ و ۱) که ظاهری مشابه دارند، دچار چالش هستند.
- برای بهبود عملکرد مدل‌ها پیشنهاد می‌شود:

۱. Data Augmentation هدفمندتر

- استفاده از augmentation خاص برای کلاس‌هایی که بیشتر اشتباه گرفته می‌شوند. مثلاً برای کلاس‌هایی که از زوایای مختلف سخت تشخیص داده می‌شن، چرخش و تغییر زاویه می‌تونه کمک‌کننده باشه.

۲. Fine-tuning دقیق‌تر روی VGG و AlexNet

- به‌جای استفاده از فقط feature extractor، می‌توان چند لایه‌ی آخر مدل را مجدداً روی دیتاست خود آموزش داد تا دقت مدل بیشتر شخصی‌سازی شود.

۳. افزایش داده برای کلاس‌های ضعیف‌تر

- کلاس‌هایی مثل ۲، ۶ یا ۸ که پراکندگی بیشتری دارند، ممکن است به دلیل کم بودن داده یا تنوع پایین در آموزش باشند.

۴. استفاده از معماری‌های پیشرفته‌تر مثل ResNet یا EfficientNet

- این مدل‌ها به‌طور خاص طراحی شدن تا دقت بیشتری با پارامتر کمتر ارائه بدن، و در پروژه‌های شبیه به این نتایج بهتری داشته‌اند.

۵. تحلیل ویژگی‌های اشتباه‌شده (Misclassification Analysis)

- بررسی بصری نمونه‌هایی که مدل اشتباه تشخیص داده می‌تونه کمک کنه بفهمیم آیا اشکال از شباهت ظاهریه یا کیفیت تصویر یا حتی لیبل‌گذاری.

۲-۶ امتیازی

۱. مدل پیشنهادی مقاله یعنی VGG + SVM را که در مرحله قبل با کرنل Linear پیاده سازی شده بود در این مرحله با استفاده کرنل های Sigmoid, Polynomial , RBF آموزش داده شد و نتایج به صورت زیر قابل مشاهده است.

کرنل SVM	Accuracy	Precision	Recall	F1 Score
Linear	0.7453	0.7464	0.7453	0.7457
RBF	0.7563	0.7598	0.7563	0.7573
Polynomial	0.7475	0.7553	0.7474	0.7480
Sigmoid	0.6624	0.6661	0.6624	0.6635

جدول ۷ - ارزیابی با کرنل های متفاوت

همان‌طور که مشاهده می‌شود، کرنل RBF بهترین عملکرد را در بین تمام گزینه‌ها داشته است، که با توجه به خاصیت انعطاف‌پذیری بالای این کرنل در مدل‌سازی مرزهای غیرخطی، نتیجه‌ی قابل‌انتظاری است. کرنل Polynomial نیز عملکردی نسبتاً نزدیک به RBF داشته و از کرنل Linear کمی بهتر عمل کرده است.

در مقابل، کرنل Sigmoid عملکرد نسبتاً ضعیف‌تری نشان داده است. این کرنل در بسیاری از کاربردهای عملی نسبت به سایر گزینه‌ها دقت پایین‌تری دارد و معمولاً کمتر توصیه می‌شود، مگر در شرایط خاص. این نتایج نشان می‌دهند که برای داده‌هایی که از طریق مدل‌های عمیق مانند VGG16 استخراج می‌شوند و ممکن است دارای مرزهای غیرخطی پیچیده باشند، استفاده از کرنل‌های غیرخطی مانند RBF یا Polynomial می‌تواند منجر به بهبود دقت مدل شود.

تفاوتی که قابل مشاهده است این است که دقت کرنل linear معمولاً کمتر از RBF است دلیل این موضوع هم وجود مرز های غیرخطی در RBF است .

سرعت آموزش در Linear به طور قابل توجه ای بیشتر از RBF بود.

به طور کلی کرنل Linear برای داده هایی با ویژگی های زیاد یا تفکیک پذیری خطی مناسب است و کرنل RBF مناسب برای داده های پیچیده یا با مرز های کلاس نامنظم مناسب است.

اگر داده ها ساده و به خوبی قابل تفکیک باشد ، linear kernel کفایت ولی اگر داده ها پیچیده یا مرز های کلاس ها غیر خطی باشند ، RBF Kernel عملکرد بهتری دارد.

استخراج ویژگی های از طریق مدل هایی مثل VGG و AlexNet میتواند مدل SVM را درگیر پیچیدگی ها تصاویر نشود و داده ها به طور خلاصه و مفید تر به مدل SVM داده میشود. این گونه مدل عملکرد دقیق تری خواهد داشت. این مورد زمانی خود را نشان میدهد که مرز بین کلاس ها پیچیده باشد.

در این پروژه، هدف اصلی طراحی و پیاده سازی یک سیستم طبقه بندی خودرو بر اساس تصویر بود که با استفاده از مدل های از پیش آموزش دیده مانند VGG16 و AlexNet و همچنین ترکیب این مدل ها با الگوریتم SVM انجام شد. برای این کار مراحل مختلفی مثل پیش پردازش داده ها، متعادل سازی کلاس ها با data augmentation ، استخراج ویژگی ها و در نهایت آموزش و ارزیابی مدل ها انجام گرفت.

نتایج نهایی نشان داد که مدل های کامل CNN یعنی VGG16 و AlexNet در این مسئله عملکرد بهتری نسبت به مدل های ترکیبی مثل SVM + VGG داشتند. این تفاوت می تواند به دلیل هماهنگی بهتر بین بخش استخراج ویژگی و طبقه بندی در معماری های end-to-end باشد. با این حال، استفاده از SVM به عنوان طبقه بند نیز عملکرد قابل قبولی ارائه داد، به خصوص در شرایطی که هدف کاهش پیچیدگی مدل یا تفکیک دقیق تر با حجم داده کمتر باشد.

همچنین مشخص شد که استفاده از data augmentation نقش مهمی در بهبود عملکرد مدل ها دارد، مخصوصاً زمانی که داده ها نامتوازن باشند. بالانس کردن کلاس ها باعث شد مدل ها آموزش بهتری ببینند و عملکرد پایدارتری داشته باشند.