



به نام خدا

دانشگاه تهران

دانشکده علوم و فناوری های میان رشته ای

# Machine Learning

تمرین سوم

نام و نام خانوادگی	امید مرادی
شماره دانشجویی	830402071
تاریخ ارسال گزارش	13 دی

## سوال اول

### Vision Transformers (ViTs)

Vision Transformers (ViTs) معماری‌های شبکه عصبی هستند که برای پردازش تصاویر با استفاده از اصول مدل‌های ترانسفورماتور طراحی شده‌اند که در اصل برای پردازش زبان طبیعی (NLP) توسعه یافته‌اند. ترانسفورماتورهای بینایی کاربردهای گسترده‌ای در کارهای رایج تشخیص تصویر مانند تشخیص اشیاء، تقسیم بندی تصویر، طبقه بندی تصویر و تشخیص عمل دارند. علاوه بر این، ViT ها در مدل سازی تولیدی و وظایف چند مدلی، از جمله زمینه سازی بصری، پاسخ گویی به سوال بصری، و استدلال بصری استفاده می شوند.

در ViTs، تصاویر به صورت توالی نمایش داده می شوند و برچسب های کلاس برای تصویر پیش بینی می شوند که به مدل ها اجازه می دهد تا ساختار تصویر را به طور مستقل یاد بگیرند. تصاویر ورودی به عنوان دنباله ای از وصله ها در نظر گرفته می شوند که در آن هر وصله با به هم پیوستن کانال های همه پیکسل ها در یک وصله و سپس نمایش خطی آن به بعد ورودی مورد نظر، به یک بردار واحد تبدیل می شود.

ساختار ViT:

1. Image Tokenization
2. Linear Embedding
3. Positional Embeddings
4. Transformer Encoder
5. Classification Token
6. Output Layer

ویژگی های اصلی:

ورودی مبتنی بر پیچ:

ViT ها وصله های تصویر را به جای کل تصویر پردازش می کنند و به مدل این امکان را می دهند که اندازه های متغیر تصویر را مدیریت کند و به طور مستقل روی مناطق فضایی تمرکز کند.

مکانیسم توجه:

ماژول MHSA به مدل اجازه می دهد تا وابستگی های دوربرد را بیاموزد و بر روی قسمت های خاصی از تصویر که مربوط به یک کار است تمرکز کند.

## مقیاس پذیری:

ViT ها با مجموعه داده ها و مدل های بزرگ (به عنوان مثال، ترانسفورماتورهای عمیق تر و گسترده تر) به خوبی مقیاس می شوند، زیرا مکانیسم توجه به طور مؤثر اطلاعات جهانی را مدیریت می کند.

## پیش آموزش و تنظیم دقیق:

ViT ها زمانی بهترین عملکرد را دارند که روی مجموعه داده های مقیاس بزرگ (مانند ImageNet-21k، JFT-300M) از قبل آموزش داده شوند و سپس روی مجموعه داده های کوچکتر و مختص کار تنظیم شوند.

## تعصبات استقرایی کمتر:

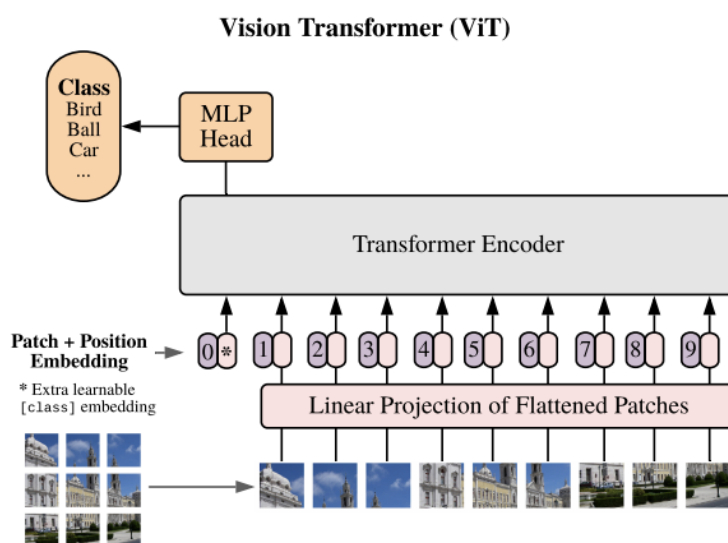
برخلاف شبکه های عصبی کانولوشنال (CNN)، ViT ها به میدان های دریافتی محلی یا اشتراک وزن متکی نیستند. این باعث می شود آنها انعطاف پذیرتر باشند، اما در عین حال تشنه اطلاعات هستند.

## عملکرد بالا:

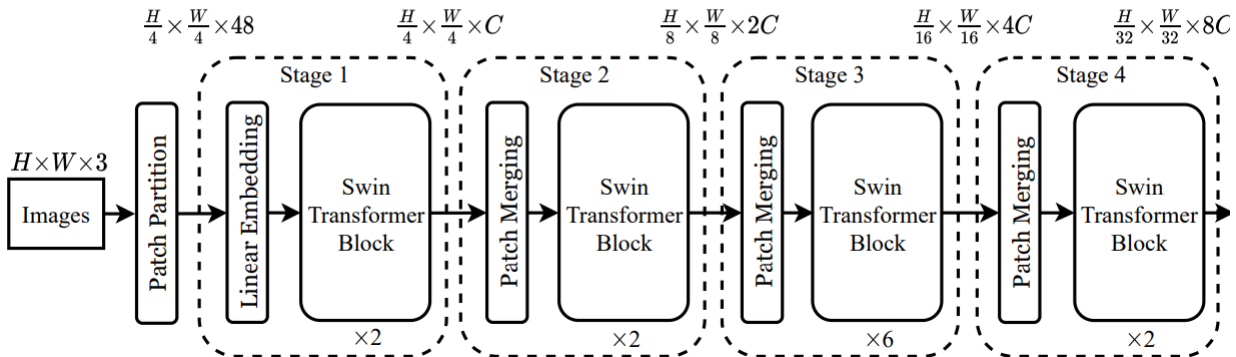
هنگامی که ViT ها به اندازه کافی از قبل آموزش دیده باشند، در بسیاری از وظایف بینایی کامپیوتری، از جمله طبقه بندی، تشخیص اشیاء و تقسیم بندی، به عملکردی پیشرفته دست می یابند.

## انعطاف پذیری در بین روش ها:

معماری ترانسفورماتور را می توان با سایر مدالیت ها مانند متن تطبیق داد و ViTs را به یک مدل متحد کننده برای یادگیری چندوجهی تبدیل کرد.



## Swin Transformers



این نمودار معماری ترانسفورماتور Swin را نشان می دهد. با تصویری به ابعاد  $H \times W \times 3$  شروع می شود که به تکه های  $4 \times 4$  غیر همپوشانی (Patch Partition) تقسیم می شود. سپس این تکه ها به صورت خطی در بردارهای مشخصه بعد  $C$  جاسازی می شوند و در نتیجه نمایشی میانی با اندازه  $H/4 \times W/4 \times C$  ایجاد می شود.

معماری به چهار مرحله تقسیم می شود:

۱. مرحله ۱: بلوک های ترانسفورماتور Swin Patch های تعبیه شده را بدون تغییر اندازه فضایی آنها ( $H/4 \times W/4$ ) پردازش می کنند، اما نمایش ویژگی را افزایش می دهند.
۲. مرحله ۲: ادغام وصله ابعاد فضایی را به  $H/8 \times W/8$  کاهش می دهد، ابعاد ویژگی را دو برابر می کند و به  $2C$  می رساند و به دنبال آن بلوک های ترانسفورماتور Swin بیشتر می شود.
۳. مرحله ۳: یکی دیگر از مراحل ادغام وصله، ابعاد فضایی را به  $H/16 \times W/16$  کاهش می دهد، ابعاد ویژگی را به  $4C$  افزایش می دهد و بلوک ها این نمایش را پردازش می کنند.
۴. مرحله ۴: ادغام نهایی ابعاد فضایی را به  $H/32 \times W/32$  با ابعاد ویژگی  $8C$  کاهش می دهد و به دنبال آن بلوک های ترانسفورماتور Swin.

معماری به تدریج وضوح فضایی را کاهش می دهد در حالی که ابعاد ویژگی را برای استخراج کارآمد ویژگی افزایش می دهد.

### ویژگی های شبکه:

۱. نمایش سلسله مراتبی: شبکه تصاویر را در چهار مرحله پردازش می کند و به تدریج ابعاد فضایی را کاهش می دهد (از  $H/4 \times W/4$  به  $H/32 \times W/32$ ) در حالی که ابعاد ویژگی را

- افزایش می دهد) از C به ۸. (این رویکرد سلسله مراتبی به طور موثر ویژگی های سطح پایین و بالا را به تصویر می کشد.
۲. ورودی مبتنی بر وصله: تصویر ورودی به وصله های غیر همپوشانی ( $4 \times 4$ ) تقسیم می شود که به صورت خطی در بردارهای ویژگی جاسازی شده اند. این به ترانسفورماتور اجازه می دهد تا تصاویر را به صورت توالی، مشابه مدل های پردازش زبان طبیعی، مدیریت کند.
۳. توجه مبتنی بر پنجره: هر بلوک ترانسفورماتور Swin توجه خود را در پنجره های محلی محاسبه می کند و پیچیدگی محاسباتی را در مقایسه با توجه جهانی کاهش می دهد. این باعث می شود شبکه برای تصاویر با وضوح بالا کارآمد باشد.
۴. مکانیسم پنجره جابجا شده: با جابجایی پنجره ها در لایه های متناوب، شبکه تعاملات بین پنجره ای را فعال می کند و توانایی آن را برای مدل سازی وابستگی های جهانی افزایش می دهد.
۵. ادغام وصله: در هر انتقال مرحله، وصله های مجاور ادغام می شوند و اندازه فضایی را کاهش می دهند و بعد ویژگی را افزایش می دهند. این پردازش چند مقیاسی شبکه های عصبی کانولوشنال (CNN) را تقلید می کند.
۶. مقیاس پذیری: شبکه را می توان با تنظیم تعداد بلوک ها در هر مرحله، بعد ویژگی C و اندازه وصله های ورودی، مقیاس پذیر کرد و آن را با وظایف مختلف و بودجه های محاسباتی سازگار کرد.
۷. تطبیق پذیری: ساختار سلسله مراتبی و مکانیسم های توجه آن را برای کارهای بینایی متنوع، از جمله طبقه بندی تصویر، تشخیص اشیا و تقسیم بندی معنایی مناسب می کند.
۸. کارایی: با محدود کردن توجه به پنجره های کوچکتر و کاهش تدریجی اندازه فضایی، ترانسفورماتور Swin به پیچیدگی خطی نسبت به اندازه تصویر دست می یابد و از نظر محاسباتی کارآمد می شود.

## EfficientNet

- این تصویر معماری EfficientNet را نشان می دهد که یک تصویر ورودی با اندازه  $3 \times 224 \times 224$  را در یک نقشه ویژگی پردازش می کند. از اجزای زیر تشکیل شده است:
۱. پیچیدگی اولیه ( $\text{Conv } 3 \times 3$ ): یک لایه کانولوشن استاندارد با هسته  $3 \times 3$  تصویر ورودی را پردازش می کند و ویژگی های فضایی اساسی را استخراج می کند.

۲. بلوک های MBConv: اینها لایه های پیچیدگی گلوگاه معکوس موبایل (MBConv) با پیچیدگی های قابل تفکیک عمیق هستند که برای کاهش پیچیدگی محاسباتی طراحی شده اند.

۱. MBConv1 با هسته  $3 \times 3$  عملیات سبک وزن را در اوایل شبکه معرفی می کند.

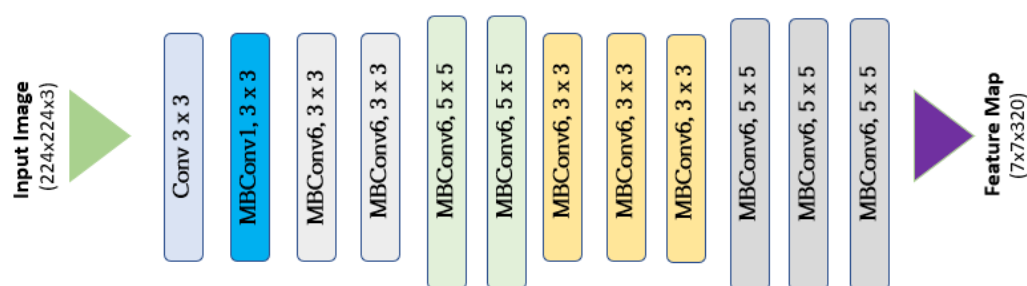
۲. MBConv6 با هسته های  $3 \times 3$  یا  $5 \times 5$  عمق می افزاید و فضای ویژگی را با استفاده از ضریب گسترش ۶ برابری برای ابعاد ویژگی میانی افزایش می دهد.

۳. استخراج ویژگی های پیشرونده: بلوک های MBConv به طور متناوب بین هسته های  $3 \times 3$  و  $5 \times 5$  برای گرفتن اطلاعات متنی ریز و بزرگتر قرار می گیرند. وضوح فضایی به تدریج کاهش می یابد و کانال های ویژگی با عمیق شدن شبکه افزایش می یابد.

۴. فیچر مپ نهایی: خروجی یک نقشه ویژگی به ابعاد  $7 \times 7 \times 320$  است که ویژگی های سطح بالای تصویر ورودی را خلاصه می کند.

این طراحی با مقیاس بندی سیستماتیک عمق، عرض و وضوح، کارایی و عملکرد مدل را متعادل می کند و آن را برای کارهای طبقه بندی بسیار موثر می سازد.

## EfficientNet Architecture



## ویژگی ها

EfficientNet برای عملکرد بالا با حداقل هزینه محاسباتی طراحی شده است. از مقیاس بندی ترکیبی برای تعادل سیستماتیک عمق (تعداد لایه ها)، عرض (تعداد کانال ها) و

وضوح (اندازه تصویر ورودی) برای مقیاس بندی کارآمد مدل استفاده می کند. این معماری دارای بلوک های MBConv است که از پیچیدگی های قابل تفکیک عمیق برای کاهش محاسبات و حفظ دقت استفاده می کنند. این هسته های  $3 \times 3$  و  $5 \times 5$  را برای استخراج ویژگی های متنوع ترکیب می کند و به تدریج ابعاد فضایی را کاهش می دهد در حالی که کانال های ویژگی را برای نمایش سلسله مراتبی ویژگی افزایش می دهد. EfficientNet با پارامترهای کمتر و FLOP کمتر در مقایسه با سایر شبکه ها به نتایج پیشرفته ای در طبقه بندی تصاویر دست می یابد.

## سوال دوم

روش ارائه شده در مقاله برای تخمین عدم قطعیت تصمیم گیری شبکه های عصبی، از تزریق دراپ آوت در زمان تست به عنوان روشی برای کمی سازی پس پرداز عدم قطعیت استفاده می کند. برخلاف دراپ آوت تعبیه شده سنتی که لایه های دراپ آوت در هر دو مرحله آموزش و تست فعال هستند، این روش لایه های دراپ آوت را به شبکه ای که قبلاً آموزش دیده است اضافه کرده و آن ها را تنها در مرحله تست فعال می کند. این امکان را فراهم می آورد تا بدون نیاز به آموزش شبکه با استفاده از دراپ آوت، عدم قطعیت تخمین زده شود و از بار محاسباتی ناشی از آموزش مجدد جلوگیری شود.

در این روش، از مونت کارلو دراپ آوت در زمان تست استفاده می شود که طی آن چندین پیش گذر تصادفی در شبکه انجام می شود. این فرآیند به تخمین توزیع پیش بینی های شبکه و عدم قطعیت مربوطه کمک می کند. میانگین و واریانس خروجی های تصادفی، به ترتیب پیش بینی و عدم قطعیت را نشان می دهند. برای بهبود کیفیت این تخمین ها، نرخ دراپ آوت بهینه از طریق کمینه سازی خطای لگاریتمی منفی (NLL) روی یک مجموعه اعتبارسنجی تنظیم می شود. علاوه بر این، یک ضریب مقیاس بندی معرفی می شود تا عدم قطعیت ها بر اساس رابطه بین خطای پیش بینی و عدم قطعیت تخمینی تنظیم شوند. این مقیاس بندی اطمینان حاصل می کند که تخمین های عدم قطعیت قوی و کالیبره شده باقی می ماند.

این روش از این جهت حائز اهمیت است که می توان آن را بر روی هر شبکه ای که قبلاً آموزش دیده اعمال کرد و نیاز به دراپ آوت در مرحله آموزش را از بین برد و بدین ترتیب تلاش محاسباتی را به طور قابل توجهی کاهش داد. نتایج تجربی در مسائل رگرسیون و وظایف چالش برانگیزی مانند شمارش جمعیت نشان می دهد که این روش به طور مؤثری عدم قطعیت را کمی سازی کرده و در عین حال دقت پیش بینی رقابتی

را حفظ می‌کند. با مقیاس‌بندی مناسب مقادیر عدم قطعیت، این روش جایگزینی عملی و کارآمد برای دراپ‌آوت تعبیه‌شده در زمینه کمی‌سازی عدم قطعیت ارائه می‌دهد.

## سوال سوم

**الف)** لایه‌های پیچشی (Convolutional Layers) و لایه‌های تماماً متصل (Fully Connected Layers) در شبکه‌های عصبی پیچشی (CNN) تفاوت‌های کلیدی دارند. لایه‌های پیچشی تنها با یک بخش کوچک از ورودی ارتباط دارند و وزن‌ها به صورت اشتراکی بین تمام بخش‌های تصویر اعمال می‌شوند، در حالی که در لایه‌های تماماً متصل، هر نورون به تمام نورون‌های لایه قبلی متصل است که باعث افزایش تعداد پارامترها می‌شود. لایه‌های پیچشی اطلاعات مکانی داده‌ها را حفظ کرده و برای استخراج ویژگی‌های محلی مانند لبه‌ها و بافت‌ها استفاده می‌شوند، در حالی که لایه‌های تماماً متصل ارتباط مکانی را از بین می‌برند و بیشتر برای ترکیب ویژگی‌های استخراج‌شده و انجام تصمیم‌گیری نهایی، مانند طبقه‌بندی، مناسب‌اند. همچنین لایه‌های پیچشی به دلیل استفاده از کرنل‌های کوچک و وزن‌های اشتراکی تعداد پارامترهای کمتری دارند، اما در لایه‌های تماماً متصل به دلیل اتصال کامل، پارامترهای بیشتری وجود دارد. لایه‌های پیچشی معمولاً در مراحل اولیه شبکه برای یادگیری ویژگی‌های سطح پایین به کار می‌روند، در حالی که لایه‌های تماماً متصل در انتهای شبکه برای یادگیری روابط پیچیده‌تر و تصمیم‌گیری نهایی استفاده می‌شوند.

**ب)** فیلترها یا هسته‌ها در لایه‌های پیچشی (Convolutional Layers) با اعمال عملیات کانولوشن بر روی ورودی، به استخراج ویژگی‌های مختلف تصویر کمک می‌کنند. هر فیلتر به یک بخش کوچک از تصویر (receptive field) اعمال می‌شود و با جابجایی در سراسر تصویر، الگوهای خاصی مانند لبه‌ها، گوشه‌ها، بافت‌ها یا ویژگی‌های سطح بالاتر را شناسایی می‌کند. وزن‌های فیلترها در طول فرآیند آموزش به گونه‌ای به‌روزرسانی می‌شوند که فیلترها بتوانند ویژگی‌های خاص و مهم تصویر را یاد بگیرند. در لایه‌های اولیه شبکه، فیلترها معمولاً ویژگی‌های ساده‌تر مانند لبه‌ها یا رنگ‌ها را شناسایی می‌کنند، در حالی که در لایه‌های عمیق‌تر، این فیلترها قادر به شناسایی ویژگی‌های پیچیده‌تر مانند اشکال، اجزا و ساختارهای کلی تصویر می‌شوند. بنابراین، فیلترها نقش کلیدی در تفکیک اطلاعات مفید از تصویر و آماده‌سازی آن برای پردازش و طبقه‌بندی ایفا می‌کنند.

**ج)** یکی از معایب و محدودیت‌های اصلی شبکه‌های عصبی پیچشی (CNN) در وظایف شناسایی اشیاء (Object Detection) این است که این شبکه‌ها به طور ذاتی برای طبقه‌بندی طراحی شده‌اند و مکان اشیاء مختلف را در تصویر مشخص نمی‌کنند. در مدل‌های اولیه مانند R-CNN، مشکل اصلی سرعت پایین در



پردازش تصاویر بود، زیرا این مدل برای هر ناحیه پیشنهادی (Region Proposal) یک بار شبکه عصبی را اجرا می کرد. این باعث افزایش زمان پردازش به ویژه برای تصاویر با تعداد زیادی ناحیه پیشنهادی می شد. علاوه بر این، فرآیند استخراج ویژگی ها و طبقه بندی ناحیه ها به صورت جداگانه انجام می شد که کارایی کلی مدل را کاهش می داد.

Fast R-CNN این مشکل را با اجرای عملیات کانولوشن برای کل تصویر به جای هر ناحیه پیشنهادی حل کرد و ویژگی های مشترک بین ناحیه های پیشنهادی را استخراج کرد، اما همچنان فرآیند تولید ناحیه های پیشنهادی (Region Proposals) خارج از شبکه و با استفاده از روش های سنتی مانند Selective Search انجام می شد که محدودیت سرعت را ایجاد می کرد.

این مشکلات باعث توسعه مدل هایی مانند Faster R-CNN شد که با معرفی شبکه ای به نام RPN (Region Proposal Network) فرآیند تولید ناحیه های پیشنهادی را به صورت یکپارچه در داخل شبکه انجام داد. بنابراین، محدودیت های سرعت و جداسازی مراحل پردازش در CNN ها دلیل اصلی توسعه این مدل های پیشرفته تر بوده است.

د) در شبکه های عصبی باقی مانده (ResNet)، مفهوم اتصالات باقی مانده (Residual Connections) به گونه ای است که ورودی هر لایه یا بلوک مستقیماً به خروجی آن اضافه می شود. این ساختار باعث می شود که به جای یادگیری مستقیم نگاشت اصلی، شبکه تنها تغییرات یا "باقی مانده" بین ورودی و خروجی را مدل سازی کند. این رویکرد با فرمول  $y = F(x) + x$  بیان می شود، که در آن  $x$  ورودی،  $F(x)$  تابع یادگیری شده توسط لایه، و  $y$  خروجی است. اتصالات باقی مانده با فراهم کردن مسیر مستقیم برای انتقال گرادیان، مشکلاتی مانند ناپدید شدن گرادیان (Vanishing Gradient) را کاهش می دهند و یادگیری در شبکه های عمیق را تسهیل می کنند. علاوه بر این، این اتصالات یادگیری را بهینه تر می کنند زیرا شبکه تنها نیاز به مدل سازی تغییرات جزئی دارد. این معماری به شبکه امکان می دهد که عمیق تر شود (حتی تا ۱۵۲ لایه) بدون کاهش دقت یا مشکلات بهینه سازی، و دقت مدل را به طور چشمگیری افزایش می دهد. اتصالات باقی مانده به عنوان یک نوآوری کلیدی، امکان پذیری آموزش شبکه های بسیار عمیق را فراهم کرده و عملکرد مدل های یادگیری عمیق را بهبود بخشیده اند.

ه) برای مقابله با مسئله داده های کم در شبکه های عمیق، می توان از روش های مختلف regularization استفاده کرد. در ورودی داده ها، روش هایی مانند افزایش داده ها (Data Augmentation) با ایجاد تغییرات جزئی روی داده ها (مانند چرخش، تغییر مقیاس، برش، وارونگی یا تغییر روشنایی) و نرمال سازی (Normalization) برای مقیاس بندی ویژگی ها، تنوع داده ها را افزایش داده و عملکرد شبکه را بهبود

می‌دهند. در ساختار شبکه، استفاده از **Dropout** برای غیرفعال کردن تصادفی نورون‌ها در هر مرحله، نرمال‌سازی دسته‌ای (**Batch Normalization**) برای تثبیت یادگیری، و انتظام‌بخشی وزن‌ها (**Weight Regularization**) با اعمال جریمه L1 یا L2 روی وزن‌ها، از پیچیدگی بیش از حد مدل جلوگیری می‌کند. در خروجی شبکه، روش **توقف زودهنگام (Early Stopping)** با متوقف کردن آموزش هنگام ثابت شدن عملکرد روی داده‌های اعتبارسنجی و استفاده از روش‌های ترکیبی (**Ensemble Methods**) برای ترکیب چندین مدل و افزایش دقت، از بیش‌برازش جلوگیری می‌کنند. این روش‌ها با توجه به نوع داده و مسئله می‌توانند تعمیم‌پذیری مدل را به طور مؤثری افزایش دهند.

## سوال چهارم

(۵)

### تابع هزینه انتخابی

تابع هزینه‌ای که در این مسئله استفاده شده است، **sparse\_categorical\_crossentropy** می‌باشد. این تابع برای مسائل دسته‌بندی چندکلاسه مناسب است و با داده‌هایی که برچسب‌های آنها به صورت عددی (مانند ۰ تا ۹ برای CIFAR-10) ذخیره شده‌اند، به خوبی کار می‌کند. در این حالت نیازی به تبدیل برچسب‌ها به فرمت یک-داغ (one-hot) نیست. این تابع خطای پیش‌بینی مدل را با استفاده از احتمال پیش‌بینی شده کاهش می‌دهد، به طوری که مدل به سمت پیش‌بینی دقیق‌تر حرکت کند. استفاده از این تابع باعث بهینه‌سازی عملکرد در مسائل طبقه‌بندی چندکلاسه می‌شود.

### تعداد نورون‌ها و تابع فعال‌ساز لایه‌ی آخر

در لایه‌ی آخر مدل، **10 نورون** قرار دارد که معادل تعداد کلاس‌های موجود در مجموعه داده CIFAR-10 است. این تعداد نورون‌ها تضمین می‌کند که مدل برای هر کلاس یک خروجی مجزا داشته باشد. همچنین از تابع فعال‌ساز **softmax** استفاده شده است. این تابع مقادیر خروجی هر نورون را به احتمال هر کلاس تبدیل می‌کند و خروجی‌ها را به یک توزیع احتمالی نرمالیزه شده (جمع احتمالات برابر ۱) تبدیل می‌کند. این ویژگی برای دسته‌بندی چندکلاسه ضروری است، زیرا با احتمال بالا، کلاس درست پیش‌بینی می‌شود. انتخاب این ترکیب باعث می‌شود مدل بتواند خروجی‌ای قابل تفسیر ارائه دهد.

### نرخ یادگیری و بهینه‌ساز (optimizer)

بهینه‌سازی در این مدل با استفاده از الگوریتم Adam انجام می‌شود. ترکیبی از روش‌های RMSProp و SGD است که با تطبیق نرخ یادگیری به صورت پویا برای هر پارامتر، سرعت همگرایی مدل را افزایش می‌دهد و به ثبات بیشتری در آموزش دست می‌یابد. نرخ یادگیری پیش‌فرض 0.001 استفاده شده است، که در بسیاری از مسائل یادگیری عمیق عملکرد مناسبی دارد. این مقدار در صورت نیاز می‌تواند با استفاده از یک کاهش‌دهنده نرخ یادگیری به تدریج تنظیم شود تا مدل در انتهای آموزش بتواند به نتایج بهتری دست یابد. Adam به دلیل ویژگی‌های تطبیقی خود، انتخاب مناسبی برای مسائل پیچیده مانند طبقه‌بندی تصاویر در CIFAR-10 است.

(۱۰)

عملکرد شبکه روی داده‌های باکیفیت (32x32) نسبت به داده‌های رزولوشن پایین (16x16) به‌طور قابل توجهی بهتر بوده است. در حالت داده‌های باکیفیت، دقت مدل (High-res on High-res) برابر با 0.1166 است. این مقدار در حالت ارزیابی داده‌های رزولوشن پایین روی مدل آموزش‌دیده با داده‌های باکیفیت (Low-res on High-res) به 0.0946 کاهش می‌یابد.

مدل آموزش‌دیده با داده‌های رزولوشن پایین روی داده‌های رزولوشن پایین (Low-res on Low-res) دقتی برابر با 0.1133 نشان داده است، که مشابه عملکرد داده‌های باکیفیت روی رزولوشن پایین (High-res on Low-res) است. این کاهش دقت ناشی از بین رفتن اطلاعات جزئی و بافت تصویر در داده‌های رزولوشن پایین است، که برای تفکیک دقیق کلاس‌ها بسیار حیاتی است. همچنین، مقادیر Precision و F1-Score نیز در داده‌های باکیفیت بیشتر از رزولوشن پایین بوده‌اند؛ به‌طور مثال، مقدار Precision برای High-res on High-res برابر 0.1154 است، در حالی که برای Low-res on High-res این مقدار به 0.0919 کاهش یافته است.

برای بهبود عملکرد مدل در شرایط رزولوشن پایین، تکنیک‌های داده‌افزایی مانند چرخش، نویزگذاری و تغییر روشنایی می‌توانند تنوع داده‌ها را افزایش داده و مدل را در برابر شرایط مختلف مقاوم کنند. استفاده از معماری‌های پیشرفته‌تری مانند ResNet یا EfficientNet که قابلیت یادگیری بهتر و استخراج ویژگی‌های عمیق‌تر را دارند، توصیه می‌شود. همچنین، بازسازی تصاویر رزولوشن پایین با استفاده از شبکه‌های سوپررزولوشن می‌تواند اطلاعات از دست‌رفته را بازیابی کرده و دقت را بهبود بخشد.

به‌عنوان مثال، استفاده از سوپررزولوشن می‌تواند مقادیر دقت را از 0.0946 در حالت Low-res on High-res به مقادیری نزدیک‌تر به 0.1166 در حالت High-res on High-res ارتقا دهد. علاوه بر این، استفاده از مکانیزم‌های توجه مانند SE-Block می‌تواند اطلاعات مهم را حتی در رزولوشن پایین استخراج کند. تنظیم

دقیق نرخ یادگیری و افزایش تعداد اپوک‌ها نیز می‌تواند به بهبود تدریجی دقت کمک کند. در مجموع، تأثیر کاهش رزولوشن بر دقت پیش‌بینی قابل توجه بوده است، اما با روش‌های پیشنهادی می‌توان این مشکل را تا حد زیادی کاهش داد.

## گزارش نهایی

شبکه طراحی شده برای مجموعه داده CIFAR-10 با هدف طبقه‌بندی تصاویر به ده کلاس مختلف مورد ارزیابی قرار گرفت. این گزارش شامل نحوه آماده‌سازی داده‌ها، معماری شبکه، تنظیمات هایپرپارامترها، نتایج، و تحلیل دقیق تأثیر رزولوشن داده‌ها بر عملکرد مدل است.

## تقسیم‌بندی دیتاست و دلایل انتخاب درصدها

برای آماده‌سازی داده‌ها، مجموعه CIFAR-10 به سه بخش تقسیم شد: ۷۰ درصد برای آموزش، ۱۵ درصد برای اعتبارسنجی و ۱۵ درصد برای تست. این تقسیم‌بندی به مدل اجازه می‌دهد که با داده‌های کافی آموزش ببیند، در حین آموزش بر روی داده‌های اعتبارسنجی ارزیابی شود و در نهایت عملکرد آن روی داده‌های تست بررسی شود. ۷۰ درصد داده‌ها برای آموزش انتخاب شدند تا مدل بتواند الگوهای پیچیده‌تری از داده‌ها یاد بگیرد. ۱۵ درصد داده‌های اعتبارسنجی برای جلوگیری از بیش‌برازش به کار گرفته شد و داده‌های تست به منظور ارزیابی نهایی مدل استفاده شدند.

## معماری شبکه

معماری شبکه طراحی شده به صورت زیر است:

جزئیات اضافی	تابع فعال ساز	ابعاد فیلتر	تعداد فیلتر / نرون	نوع	لایه
(32, 32, 3)	-	-	-	-	ورودی
BatchNorm, Dropout	ReLU	(3, 3)	64	Conv2D	کانولوشن ۱
BatchNorm, Dropout	ReLU	(3, 3)	128	Conv2D	کانولوشن ۲
BatchNorm, Dropout	ReLU	(3, 3)	256	Conv2D	کانولوشن ۳
-	-	-	-	Flatten	تخت سازی
BatchNorm, Dropout	ReLU	-	256	Dense	لایه کاملاً متصل
-	Softmax	-	10	Dense	خروجی

این معماری شامل سه لایه کانولوشن برای استخراج ویژگی، یک لایه تخت سازی برای تبدیل داده های چندبعدی به یک بردار و دو لایه کاملاً متصل برای طبقه بندی است.

## توضیح انتخاب های پیرامترها

برای این مدل، نرخ یادگیری اولیه ۰.۰۰۱ انتخاب شده است که مقدار پیش فرض برای بهینه ساز Adam است. این نرخ یادگیری تطبیقی باعث می شود مدل با سرعت مناسب همگرا شود. بهینه ساز Adam انتخاب شد، زیرا ترکیبی از مزایای روش های SGD و RMSProp را دارد و می تواند نرخ یادگیری را برای هر پارامتر به صورت دینامیک تنظیم کند. تابع هزینه sparse\_categorical\_crossentropy برای محاسبه خطا استفاده شد، زیرا این تابع برای مسائل طبقه بندی چندکلاسه بسیار مناسب است و با برچسب های عددی کار می کند.

## نمودارهای Loss و Accuracy

نمودارهای زیر کاهش Loss و افزایش Accuracy را در طول آموزش نشان می دهند. این نمودارها مشخص می کنند که مدل در داده های باکیفیت عملکرد بهتری داشته است و داده های رزولوشن پایین باعث کاهش دقت و افزایش خطا شده اند. (نمودارها در فایل های همراه قرار دارند.)

## نتایج و مقایسه آنها

نتایج ارزیابی مدل در چهار حالت مختلف به شرح زیر است:

حالت	Accuracy	Precision	F1-Score
High-res on High-res	0.116667	0.115470	0.112602
Low-res on High-res	0.094667	0.091958	0.085875
High-res on Low-res	0.113333	0.122447	0.107741
Low-res on Low-res	0.113333	0.122447	0.107741

این نتایج نشان می‌دهند که مدل در داده‌های باکیفیت عملکرد بهتری داشته و کاهش رزولوشن باعث کاهش دقت و کارایی مدل شده است.

### تحلیل دقیق از تأثیر رزولوشن داده‌ها بر عملکرد مدل

کاهش رزولوشن داده‌ها باعث کاهش قابل توجه دقت مدل شده است. اطلاعات مهمی که در داده‌های باکیفیت موجود بودند، در داده‌های رزولوشن پایین از دست رفته‌اند و مدل نمی‌تواند ویژگی‌های کافی برای تفکیک کلاس‌ها استخراج کند. این امر به وضوح در ماتریس‌های آشفتگی مشاهده می‌شود، جایی که مدل در تشخیص صحیح برخی از کلاس‌ها دچار اشتباه شده است. همچنین، کاهش رزولوشن باعث افزایش نویز در ویژگی‌های استخراج‌شده شده و عملکرد مدل را کاهش داده است.

### پیشنهادهایی برای بهبود عملکرد مدل

برای بهبود عملکرد مدل در شرایط رزولوشن پایین، چند پیشنهاد ارائه می‌شود. استفاده از تکنیک‌های داده‌افزایی مانند چرخش، نویزگذاری و برش می‌تواند به بهبود یادگیری کمک کند. بهره‌گیری از معماری‌های پیشرفته‌تر مانند ResNet یا MobileNet که توانایی یادگیری ویژگی‌های پیچیده‌تر را دارند، می‌تواند عملکرد مدل را افزایش دهد. بازسازی تصاویر رزولوشن پایین با استفاده از شبکه‌های سوپررزولوشن نیز می‌تواند اطلاعات از دست‌رفته را بازیابی کند. در نهایت، استفاده از مکانیزم‌های توجه (Attention Mechanisms) و تنظیم دقیق هایپرپارامترها از جمله نرخ یادگیری و تعداد اپوک‌ها می‌تواند به بهبود نتایج منجر شود.

## منابع