

# Transfer Learning Optimization

## بررسی تکنیک‌های نرمال‌سازی و دینامیک گرادیان در MobileNetV2 بر روی CIFAR-10

هیوا ابوالهادی زاده

۴۰۰۴۰۵۰۰۴

مباحث ویژه در هوش مصنوعی  
استاد: جناب آقای دکتر افتخاری

۲۲ فروردین ۱۴۰۴

چکیده

در این پژوهش، تأثیر سه روش نرمال‌سازی شامل BatchNorm، LayerNorm و Filter Response Normalization (FRN) در چارچوب یادگیری انتقالی با مدل MobileNetV2 بر مجموعه داده CIFAR-10 بررسی شده است. همچنین نقش Gradient Clipping در پایداری گرادیان‌ها تحلیل شده است. نتایج نشان می‌دهد که ساختار FRN بدون Clipping بهترین عملکرد را ارائه می‌دهد و Clipping در ترکیب با BatchNorm موجب بهبود همگرایی می‌شود. LayerNorm نیز بدون نیاز به Clipping عملکردی پایدار از خود نشان داده است. این یافته‌ها انتخاب نرمال‌سازی مناسب را به عنوان عاملی کلیدی در موفقیت fine-tuning برجسته می‌سازد.

## فهرست مطالب

۳	۱	مقدمه
۴	۲	مروری بر کارهای پیشین
۵	۳	روش‌شناسی
۵	۱.۳	معماری مدل و تغییرات
۸	۲.۳	تنظیمات آزمایش‌ها
۱۱	۴	تحلیل تجربی
۱۱	۱.۴	دقت و خطا
۱۱	۲.۴	بررسی چشم‌انداز تابع هزینه
۱۲	۵	بحث
۱۲	۶	نتیجه‌گیری و پیشنهادات
۱۳	۷	خلاصه مدیریتی
۱۳	۸	ضمائم
۱۳	۹	منابع

## ۱ مقدمه

یادگیری انتقالی<sup>۱</sup> یکی از رویکردهای کلیدی در یادگیری ماشین مدرن است که به مدل‌ها امکان می‌دهد دانش حاصل از حل یک مسئله را برای حل مسئله‌ای دیگر به کار ببرند. این روش به‌ویژه در شرایطی که داده‌های آموزشی محدود هستند، بسیار مؤثر عمل می‌کند. به جای آموزش مدل از ابتدا، از یک مدل از پیش آموزش دیده استفاده شده و آن را برای وظیفه‌ی جدید بهینه‌سازی می‌کنیم. این فرآیند که اصطلاحاً fine-tuning نام دارد، در سال‌های اخیر به‌طور گسترده در کاربردهای بینایی ماشین، پردازش زبان طبیعی، و سایر حوزه‌ها مورد استفاده قرار گرفته است.

در مسیر آموزش مدل‌های عمیق، یکی از عوامل تعیین کننده در پایداری و سرعت همگرایی، نرمال‌سازی داده‌ها و لایه‌ها است. تکنیک‌هایی نظیر Batch Normalization، Layer Normalization و Filter Response Normalization (FRN) با تنظیم آماری ویژگی‌ها در طول آموزش، باعث می‌شوند مدل بتواند سریع‌تر و پایدارتر یاد بگیرد. با این حال، رفتار این روش‌های نرمال‌سازی در زمینه‌ی یادگیری انتقالی، به‌ویژه هنگام تطبیق مدل‌های از پیش آموزش دیده با داده‌های جدید، هنوز نیازمند بررسی دقیق‌تر است.

هدف این پروژه، تحلیل تجربی و سیستماتیک عملکرد روش‌های مختلف نرمال‌سازی در فرآیند انتقال یادگیری است. در این راستا، از مدل MobileNetV2 به عنوان مدل پایه استفاده شده است؛ مدلی سبک و بهینه برای دستگاه‌های با توان پردازشی پایین که بر روی مجموعه‌های داده‌ی بزرگ مانند ImageNet آموزش دیده است. داده‌های مورد استفاده در این پروژه، مربوط به مجموعه‌ی تصویری CIFAR-10 هستند که شامل ۶۰۰۰۰ تصویر رنگی در ۱۰ کلاس مختلف با ابعاد  $32 \times 32$  پیکسل می‌باشند. برای سازگاری با ورودی MobileNetV2، این تصاویر به ابعاد  $224 \times 224$  باز نمونه‌گیری شده‌اند.

در این گزارش، سه نوع head نرمال‌سازی مختلف طراحی و با مدل پایه تلفیق شده‌اند. همچنین، تأثیر تکنیک Gradient Clipping نیز بررسی شده تا تأثیر آن بر پایداری و رفتار گرادینت‌ها در طی فرآیند fine-tuning مشخص گردد. برای تحلیل جامع، از ابزارهایی مانند تحلیل چشم‌انداز تابع هزینه<sup>۲</sup> و پایش گرادینت استفاده شده است.

سوالات اصلی این پروژه عبارت‌اند از:

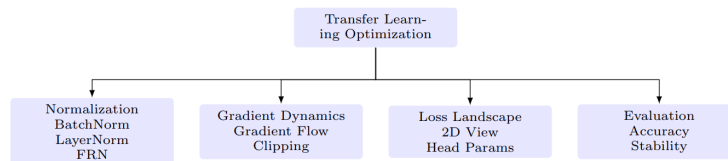
کدام نوع نرمال‌سازی در انتقال یادگیری بر روی CIFAR-10 عملکرد بهتری دارد؟

آیا استفاده از Gradient Clipping باعث بهبود پایداری آموزش می‌شود؟

رفتار گرادینت‌ها در لایه‌های مختلف مدل چه تفاوت‌هایی در شرایط نرمال‌سازی مختلف دارد؟

---

<sup>۱</sup>Transfer Learning  
<sup>۲</sup>Loss Landscape



شکل ۱: نمای کلی اجزای مورد بررسی در این گزارش

## ۲ مروری بر کارهای پیشین

یادگیری انتقالی، به‌ویژه در حوزه بینایی ماشین، در سال‌های اخیر مورد توجه گسترده‌ای قرار گرفته است. استفاده از مدل‌های از پیش آموزش دیده مانند ResNet، EfficientNet و MobileNet در مسائل جدید، موجب صرفه‌جویی در منابع محاسباتی و بهبود عملکرد مدل شده است. با این حال، چگونگی تطبیق این مدل‌ها با داده‌های جدید و نقش تکنیک‌های نرمال‌سازی در فرآیند fine-tuning، همچنان از موضوعات باز و پرچالش در ادبیات هوش مصنوعی هستند. در مطالعه‌ی [Singh and Krishnan, 2020]، نویسندگان تکنیک جدیدی به نام Filter Response Normalization (FRN) را معرفی می‌کنند که برخلاف Batch Normalization، مستقل از اندازه‌ی batch و توزیع آماری داده‌ها عمل می‌کند. این ویژگی، FRN را به گزینه‌ای مناسب برای یادگیری انتقالی، به‌ویژه در شرایط با batch size کوچک یا داده‌های نامتوازن، تبدیل کرده است.

از سوی دیگر، تحلیل چشم‌انداز تابع هزینه<sup>۳</sup> به‌عنوان ابزاری برای درک بهتر رفتار مدل‌ها در طول آموزش پیشنهاد شده است. مقاله‌ی مشهور [Li et al., 2018] نشان می‌دهد که شکل و همواری نواحی کمینه‌ی تابع هزینه می‌تواند با قابلیت تعمیم مدل ارتباط مستقیم داشته باشد. در زمینه‌ی پایداری آموزش شبکه‌های عمیق، تکنیک Gradient Clipping نیز به‌عنوان روشی برای کنترل gradient explosion مطرح شده است. یافته‌های مطالعاتی مانند [Zhang et al., 2019] نشان می‌دهد که در شرایطی خاص، به‌ویژه هنگام fine-tuning لایه‌های عمیق، clipping می‌تواند منجر به همگرایی پایدارتر و کاهش نوسانات در مقدار گرادیان‌ها شود.

مطالعه‌ی حاضر، در ادامه‌ی این پژوهش‌ها، تلاش دارد تا با تمرکز بر ترکیب MobileNetV2، داده CIFAR-10، و تکنیک‌های نرمال‌سازی متنوع، به ارزیابی تجربی نقش نرمال‌سازی و clipping در یادگیری انتقالی بپردازد. در این مسیر، از تحلیل loss landscape و پراکندگی گرادیان نیز برای درک عمیق‌تر استفاده خواهد شد.

<sup>۳</sup> Loss Landscape

## ۳ روش شناسی

این بخش به تشریح فرآیند آماده سازی داده، طراحی مدل انتقال یادگیری، پیاده سازی Head های نرمال سازی، و تنظیمات دقیق آزمایش ها می پردازد. کلیه مراحل در محیط Google Colab با استفاده از کتابخانه های PyTorch و Torchvision پیاده سازی شده اند.

### ۱.۳ معماری مدل و تغییرات

مدل پایه استفاده شده در این پروژه MobileNetV2 است که از پیش بر روی مجموعه داده ImageNet آموزش دیده و به عنوان Feature Extractor عمل می کند.<sup>۴</sup> برای تطبیق این مدل با وظیفه ی طبقه بندی تصاویر CIFAR-10، بخش classifier آن حذف شده و خروجی آن با Head جدیدی ترکیب شده است. این Head مسئول تولید خروجی ۱۰ کلاسه نهایی است. سه نوع Head مختلف با هدف بررسی تأثیر روش های نرمال سازی طراحی و پیاده سازی شده اند:

□ **BatchNormHead**: این Head شامل یک لایه Fully Connected با خروجی ۲۵۶، یک لایه Batch Normalization<sup>۵</sup>، تابع فعال ساز ReLU، Dropout و در نهایت یک لایه FC خروجی است. میانگین و واریانس را در طول batch محاسبه می کند. نرمال سازی به صورت  $\hat{x} = \frac{x - \mu_{batch}}{\sqrt{\sigma_{batch}^2 + \epsilon}}$  انجام شده و سپس توسط مقیاس و انتقال  $(\gamma, \beta)$  تنظیم می شود.

□ **LayerNormHead**: ساختاری مشابه BatchNormHead دارد، با این تفاوت که از Layer Normalization استفاده می کند. آماره ها به صورت مستقل از batch و در طول ویژگی ها (در هر نمونه) محاسبه می شوند، و نرمال سازی روی feature dimension داخلی انجام می گیرد.

□ **FRNHead**: در این ساختار از تکنیک Filter Response Normalization (FRN) استفاده شده است که به طور خاص طراحی شده تا از وابستگی به آماره های دسته جلوگیری کند. لایه FRN در این Head به صورت سفارشی پیاده سازی شده و با یک پارامتر  $\tau$  ترکیب می شود:  $\max(\gamma x + \beta, \tau)$ . برخلاف دو مورد قبلی، هیچ وابستگی به batch یا channel خاص ندارد.

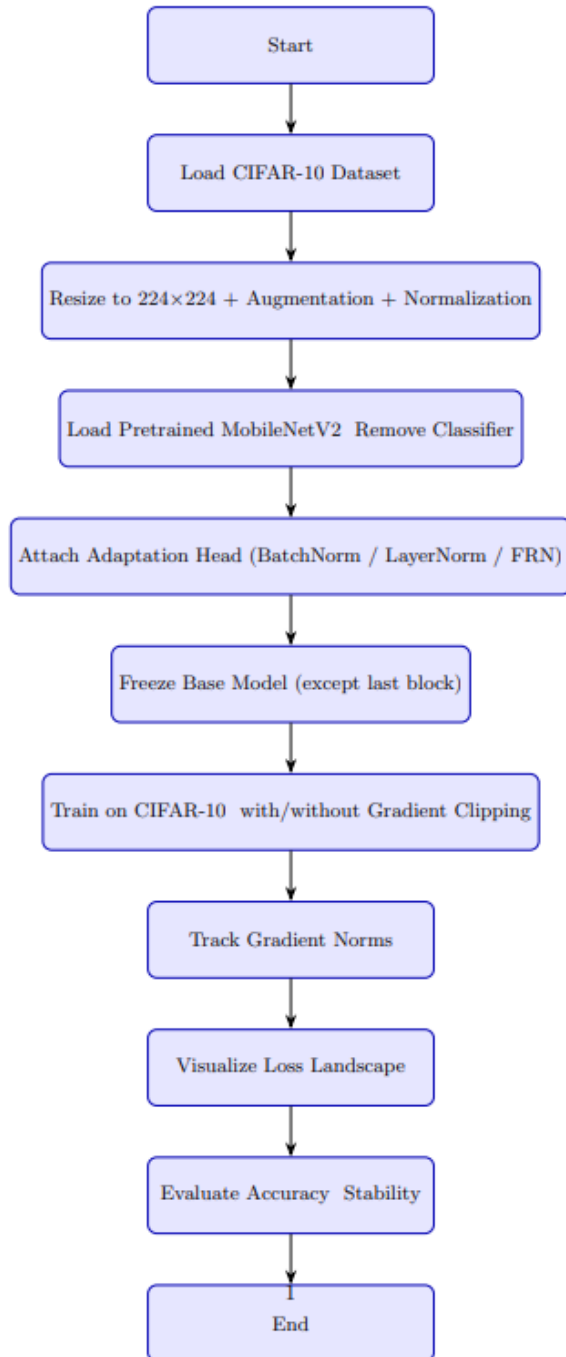
از منظر یادگیری انتقالی، BatchNorm در شرایطی با batch size کوچک<sup>۶</sup> با نوسان آماره ها دچار افت عملکرد می شود.<sup>۷</sup> FRN به دلیل مستقل بودن از می تواند در این شرایط عملکرد

<sup>۴</sup> استفاده از مدل های از پیش آموزش دیده یکی از رایج ترین روش ها در یادگیری انتقالی است [Pan and Yang, 2010].

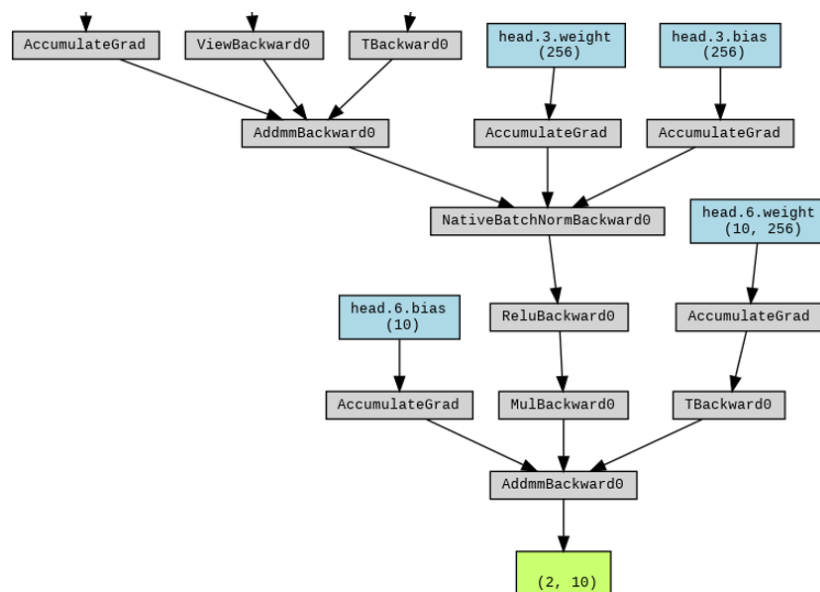
<sup>۵</sup> نرمال سازی دسته ای نقش مهمی در پایداری گرادین ها و تسریع همگرایی دارد [Singh and Krishnan, 2020].

<sup>۶</sup> مانند fine-tuning با داده کم

<sup>۷</sup> رجوع شود به [Singh and Krishnan, 2020]

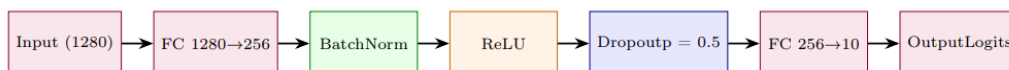


شکل ۲: جریان کار



شکل ۳: بخشی از معماری مدل

بهتری داشته باشد. در مقابل، LayerNorm گرچه بیشتر در حوزه‌های مبتنی بر توالی مانند NLP رایج است، اما در برخی ساختارهای Vision Transformer نیز موفق بوده و در این پروژه نیز برای مقایسه استفاده شده است. در تمامی Headها، ابتدا روی خروجی Feature Extractor عملیات Global Average Pooling انجام می‌شود، سپس ویژگی‌ها با Flatten به برداری با اندازه ثابت تبدیل شده و به لایه FC ورودی Head داده می‌شوند. برای جلوگیری از Overfitting و حفظ ویژگی‌های عمومی آموخته‌شده توسط MobileNetV2، تمامی لایه‌های مدل پایه فریز شده‌اند به جز بلوک انتهایی آن (features.18)، که برای فرآیند fine-tuning باز شده و با نرخ یادگیری پایین‌تر آموزش دیده است.



شکل ۴: معماری BatchNorm Head



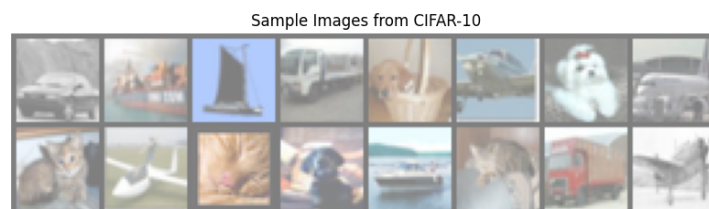
شکل ۵: معماری LayerNorm Head



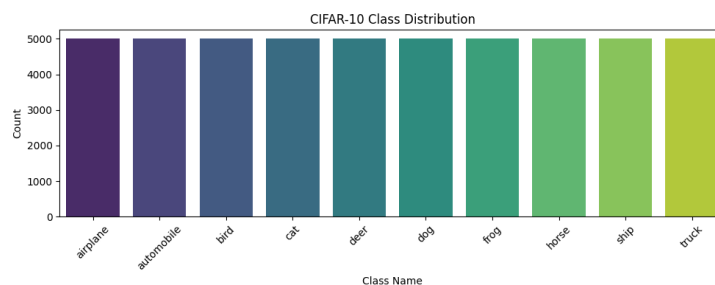
شکل ۶: معماری FRN Head

## ۲.۳ تنظیمات آزمایش‌ها

مجموعه داده مورد استفاده در این پروژه، CIFAR-10 است که شامل ۶۰۰۰۰ تصویر رنگی  $32 \times 32$  در ده کلاس است. برای سازگاری با ورودی مدل MobileNetV2، تمام تصاویر به اندازه  $224 \times 224$  بازنمونه‌گیری شده‌اند.<sup>۸</sup>



شکل ۷: نمونه‌هایی از داده‌ها



شکل ۸: توزیع داده‌ها به ازای هر کلاس

برای داده‌های آموزش، عملیات Data Augmentation شامل موارد زیر انجام شده است:

RandomHorizontalFlip □

RandomCrop با padding برابر ۴ □

نرمال‌سازی با میانگین و انحراف معیار CIFAR-10: □

$$\mu = (0.4914, 0.4822, 0.4465)$$

$$\sigma = (0.2023, 0.1994, 0.2010)$$

<sup>۸</sup>ورودی اصلی MobileNetV2 برای تصویر با اندازه  $224 \times 224$  طراحی شده است.



برای داده‌های اعتبارسنجی و تست، تنها Resize و نرمال‌سازی اعمال شده است. تقسیم‌بندی داده‌ها به صورت ۸۰٪ آموزش و ۲۰٪ اعتبارسنجی انجام شده است. یک کلاس اختصاصی به نام CIFAR10DataLoader مسئول بارگذاری، تقسیم، و بصری‌سازی داده‌هاست (از جمله توزیع کلاس‌ها و پیش‌نمایش نمونه‌ها). آزمایش‌های طراحی شده به صورت ترکیبی از سه نوع Head<sup>۹</sup> و دو حالت استفاده یا عدم استفاده از Gradient Clipping (با مقدار آستانه ۱/۰) می‌باشند. در نتیجه ۶ حالت مختلف بررسی شده‌اند. برای آموزش از تابع هزینه CrossEntropyLoss و بهینه‌ساز Adam استفاده شده است. نرخ یادگیری برای لایه‌های Head برابر  $10^{-3}$  و برای لایه‌ی باز MobileNetV2 برابر  $10^{-5}$  در نظر گرفته شده است. پارامترهای کلیدی آزمایش‌ها:

□ تعداد epoch: ۱۵

□ batch size: ۶۴ (و ۳۲ برای تحلیل گرادیان)

□ بهینه‌ساز: Adam با دو نرخ یادگیری مجزا

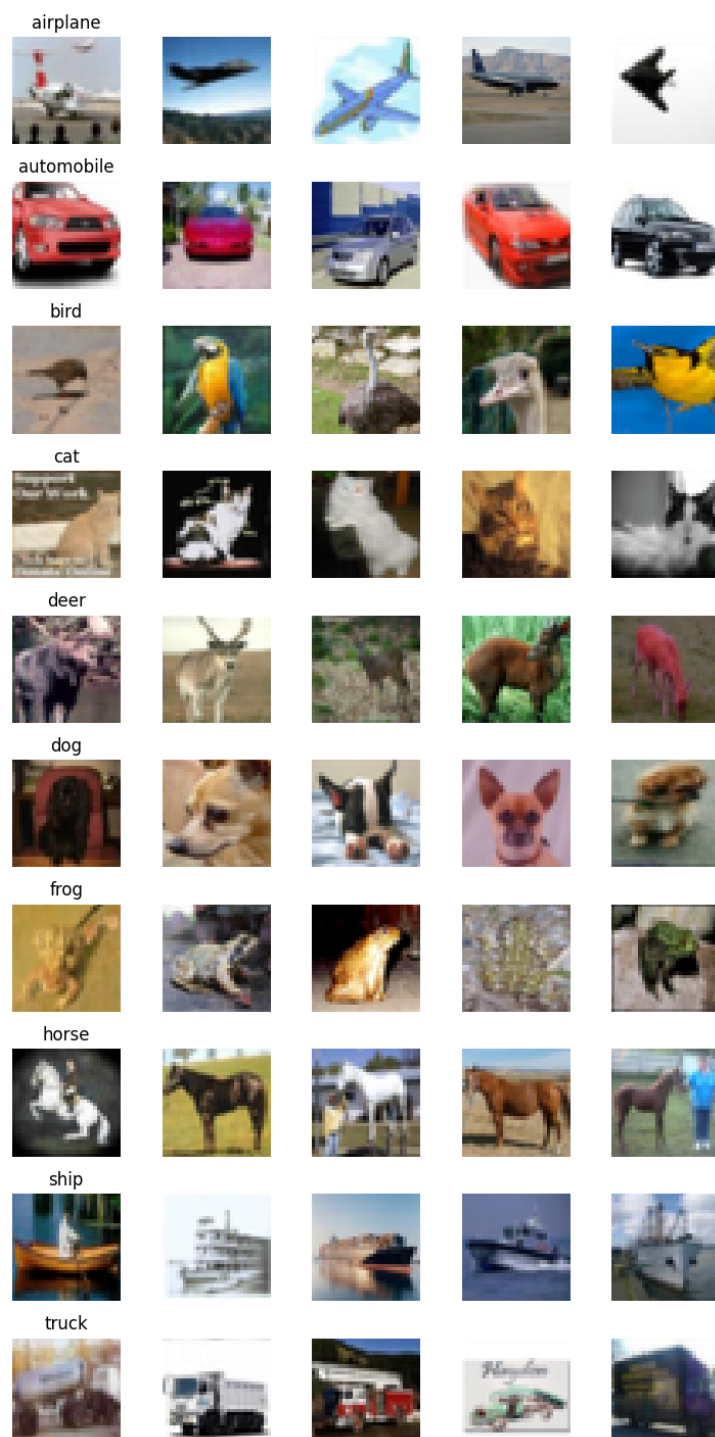
□ Gradient Clipping: فقط در نیمی از آزمایش‌ها با مقدار ۱/۰

فرآیند آموزش شامل پایش دقت و مقدار Loss در هر epoch، ذخیره‌ی نتایج، و در صورت نیاز، تحلیل گرادیان‌ها و Loss Landscape برای مدل‌های مختلف است. همچنین، ابزار اختصاصی برای ثبت گرادیان لایه‌ها و مقایسه‌ی توزیع آن‌ها در دو حالت clipping و بدون clipping توسعه داده شده است [Zhang et al., 2019, Li et al., 2018].

---

BatchNorm, LayerNorm, FRN<sup>۹</sup>

5 Samples per CIFAR-10 Class



شکل ۹: نمونه هایی از هر کلاس

## ۴ تحلیل تجربی

### ۱.۴ دقت و خطا

برای ارزیابی اثر ساختارهای مختلف نرمال‌سازی، شش آزمایش مجزا شامل ترکیب سه روش نرمال‌سازی<sup>۱۰</sup> با یا بدون تکنیک Gradient Clipping انجام شد. تمامی مدل‌ها روند صعودی نسبتاً همواری در دقت اعتبارسنجی طی ۲۵ دوره آموزش نشان دادند. به‌طور خاص، بالاترین دقت نهایی متعلق به FRN-NoClip با مقدار 90.57% بود. پس از آن، مدل BatchNorm-Clip با 90.71% و مدل LayerNorm-NoClip با 90.22% قرار گرفتند. پایین‌ترین دقت نهایی مربوط به BatchNorm-NoClip با 89.64% بود. از نظر همگرایی، مدل‌های دارای Clipping معمولاً در دوره‌های ابتدایی همگرایی سریع‌تری داشتند. به‌طور مثال، دقت BatchNorm-Clip در اولین اپک 75.16% بود، درحالی‌که نسخه بدون کلیپینگ 73.73% را ثبت کرد. با این حال، در دوره‌های پایانی، تفاوت عملکرد کاهش یافته و در برخی موارد مانند LayerNorm اثر Clipping حتی ناچیز یا منفی بود.

### ۲.۴ بررسی چشم‌انداز تابع هزینه

با وجود خطای فنی در رسم نمودار چشم‌انداز تابع هزینه، روند همگرایی پایدار مدل‌های FRN و LayerNorm بدون نیاز به Clipping می‌تواند نشانه‌ای از ساختار نرم‌تر loss surface در این معماری‌ها باشد. همچنین ثبات نرم‌گرادیان در طی زمان (کاهش تدریجی از حدود ۴۴ به زیر ۲۹) در این ساختارها مؤید همین موضوع است.



شکل ۱۰: آنالیز آموزش مدل‌ها

BatchNorm, LayerNorm, FRN<sup>۱۰</sup>

## ۵ بحث

نتایج تجربی نشان دادند که نوع نرمال سازی تأثیر مستقیمی بر پایداری و کیفیت آموزش دارد. ساختار FRN بدون کلیپینگ بالاترین دقت را به دست آورد و نشان داد که به طور ذاتی پایدار است.

مدل های LayerNorm نیز به ویژه در حالت بدون Clipping رفتار همگرایی روان و پایداری را از خود نشان دادند، که می توان آن را به ماهیت مستقل بودن آن از آماره های Batch نسبت داد. روند تغییر نرُم گرادیان در این ساختار نیز بسیار منظم بود.

در مقابل، BatchNorm بدون کلیپینگ دارای نوسانات بیشتری بود. استفاده از Clipping در این ساختار باعث بهبود پایداری گرادیان و افزایش دقت نهایی از 89.64% به 90.71% شد. به طور مشخص، نرُم گرادیان در BatchNorm-Clip کنترل شده تر بود (از حدود ۱۵.۴۶ تا ۰.۲۹).

نقش گرادیان کلیپینگ در تمامی ساختارها یا مثبت بود (مانند BatchNorm) یا تأثیر خنثی داشت (مانند LayerNorm) و در هیچ کدام کاهش عملکرد ایجاد نکرد.

## ۶ نتیجه گیری و پیشنهادات

یافته های اصلی به شرح زیر است:

- ساختار FRN-NoClip بهترین عملکرد نهایی را با دقت 90.57% نشان داد.
- Clipping در ساختار BatchNorm منجر به بهبود عملکرد شد (افزایش 1.07%).
- ساختارهای FRN و LayerNorm بدون نیاز به Clipping به عملکرد پایدار و دقیقی رسیدند.
- محدودیت ها:
- آموزش فقط روی بخش Head مدل، با فریز بودن لایه های پایه
- استفاده از منابع سخت افزاری محدود (مانند Google Colab)
- عدم امکان ترسیم چشم انداز تابع هزینه برای برخی مدل ها به دلیل خطای list index out of range
- پیشنهادهای توسعه آتی:
- اعمال نرمال سازی ترکیبی یا بررسی تکنیک های جدید مانند GroupNorm
- استفاده از مدل های سبک تر مانند EfficientNet-Lite
- انجام Fine-Tuning کامل و نه فقط Head
- ارزیابی روی داده های پیچیده تر مانند Tiny-ImageNet

## ۷ خلاصه مدیریتی

در این پروژه، سه ساختار نرمال‌سازی متفاوت به همراه یا بدون گرادیان کلیپینگ در فرآیند انتقال یادگیری ارزیابی شدند. یافته‌های کلیدی عبارت‌اند از:

- ساختار FRN-NoClip بهترین دقت را ارائه داد (90.57%).
  - گرادیان کلیپینگ در ساختار BatchNorm باعث بهبود عملکرد شد (+1.07%).
  - LayerNorm و FRN در اکثر موارد بدون نیاز به Clipping عملکرد پایداری داشتند.
- پاسخ به سه سؤال کلیدی:

۱. بهترین تکنیک نرمال‌سازی: FRN-NoClip با دقت نهایی 90.57%
۲. تأثیر Clipping: در BatchNorm مؤثر و مثبت، در سایر روش‌ها تأثیر کم یا خنثی
۳. راهکارهای عملی: استفاده از LayerNorm/FRN در پروژه‌های واقعی با batch کوچک؛ فعال‌سازی Clipping در BatchNorm؛ توجه به محدودیت منابع در طراحی معماری Head

## ۸ ضمائم

لینک گوگل کولب تست‌ها

## ۹ منابع

## References

- Hao Li, Zheng Xu, Gavin Taylor, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *NeurIPS*, 2018.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- Aarush Singh and Dilip Krishnan. Filter response normalization layer: Eliminating batch dependence in the training of deep neural networks. In *CVPR*, 2020.

Zihan Zhang, Hao He, and Ruslan Salakhutdinov. Understanding and improving gradient clipping. *arXiv preprint arXiv:1905.11881*, 2019.