

تکلیف اول درس شبکه‌های عصبی

وحید ملکی

شماره دانشجویی: ۴۰۳۱۳۰۰۴

۱۴۰۴ آبان ۲۶

سوال ۴

الف: محو شدن و انفجار گرادیان چیست؟

در آموزش شبکه‌های عصبی عمیق، گرادیان تابع هزینه نسبت به وزن‌های لایه‌های اولیه با استفاده از قاعده زنجیره‌ای و به صورت ضرب مکرر مشتق توابع فعال‌سازی و وزن‌های لایه‌های بعدی محاسبه می‌شود.

• **محو شدن گرادیان (Vanishing Gradient):** زمانی رخ می‌دهد که مقادیر مطلق مشتق توابع فعال‌سازی (مانند Sigmoid و Tanh) در نواحی اشباع کتر از ۱ باشد. با افزایش تعداد لایه‌ها، ضرب مکرر اعداد کوچکتر از ۱ باعث می‌شود گرادیان در لایه‌های ابتدایی به صورت نمایی به صفر نزدیک شود. نتیجه عملی: لایه‌های اولیه تقریباً هیچ به روزرسانی دریافت نمی‌کنند، ویژگی‌های سطح پایین (مانند لبه‌ها و گوش‌ها در یعنی ماشین) به خوبی یاد گرفته نمی‌شوند و کل شبکه عملاً به یک شبکه کم عمق تبدیل می‌شود.

• **انفجار گرادیان (Exploding Gradient):** زمانی رخ می‌دهد که مقادیر مطلق مشتق‌ها یا وزن‌ها بزرگتر از ۱ باشند. ضرب مکرر این اعداد باعث رشد غایی گرادیان می‌شود و مقدار آن به سرعت به اعداد بسیار بزرگ یا حتی NaN تبدیل می‌شود. نتیجه عملی: وزن‌ها تغییرات بسیار شدیدی می‌کنند، تابع هزینه نوسانات شدید پیدا می‌کند و آموزش کاملاً ناپایدار می‌شود.

ب: روش کلاسیک و بسیار موفق — نرمال‌سازی دسته‌ای (Batch Normalization)

Batch Normalization که در سال ۲۰۱۵ توسط Christian Szegedy و Sergey Ioffe معرفی شد، تا چند سال استاندارد طلایی برای حل این دو مشکل بود. فرمول کامل BatchNorm در زمان آموزش:

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i \quad , \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad , \quad y_i = \gamma \hat{x}_i + \beta$$

که γ و β پارامترهای قابل یادگیری هستند و در زمان استنتاج از میانگین و واریانس متحرک استفاده می‌شود. چرا BatchNorm این مشکلات را حل می‌کند؟

- توزیع ورودی هر لایه را به میانگین صفر و واریانس یک می‌رساند \Rightarrow از ورود به نواحی اشباع توابع فعال‌سازی جلوگیری می‌کند.
- گرادیان‌ها را در مقیاس مناسب نگه می‌دارد \Rightarrow هم از محoshدن و هم از انفجار جلوگیری می‌کند.
- اجازه استفاده از نرخ یادگیری بسیار بالاتر را می‌دهد.
- اثر منظم‌سازی دارد و نیاز به Dropout را کاهش می‌دهد.

معایب مهم :BatchNorm

- وابستگی شدید به سایز پچ (در Batch Size کوچک، تخمین μ و σ نویزدار می‌شود)
- افزایش مصرف حافظه (نیاز به ذخیره آمار پچ)
- مشکل در شبکه‌های بازگشتی و آموزش آنلاین
- تفاوت رفتار بین آموزش و استنتاج

ج: مقاله علمی پیشنهادی – حذف کامل نرمال‌سازی با روشهای هوشمندانه

عنوان مقاله: High-Performance Large-Scale Image Recognition Without Normalization

نویسنده‌گان: Samuel L. Smith، Andrew Brock، Soham De (همه از DeepMind)

کنفرانس/سال: ICML 2021

لینک کامل: <https://arxiv.org/abs/2102.06171>

این مقاله یکی از مهم‌ترین مقالات سال ۲۰۲۱ در حوزه پایداری آموزش شبکه‌های عمیق است و نشان می‌دهد که نه تنها می‌توان BatchNorm را حذف کرد، بلکه با حذف آن می‌توان به عملکرد بهتری هم رسید!

ایده اصلی مقاله: شبکه‌های بدون نرمال‌سازی (NFNet) یا Normalizer-Free Nets

نویسنده‌گان مشاهده کردند که BatchNorm در واقع دو نقش دارد: ۱. پایدار کردن آموزش (جلوگیری از Vanishing/Exploding Gradient) ۲. منظم‌سازی (Normalizing) آنها نشان دادند که نقش اول را می‌توان با روشنی بسیار ساده‌تر و کارآمدتر جایگزین کرد: برش تطبیقی گرادیان (Adaptive Gradient Clipping — AGC)

توضیح کامل روشن (AGC)

روش کلاسیک Gradient Clipping یک آستانه ثابت (مثلاً ۰.۰۱) روی نرم گرادیان کل شبکه اعمال می‌کند. اما AGC هوشمندانه‌تر عمل می‌کند:

برای هر پارامتر w در هر لایه، نسبت نرم گرادیان به نرم وزن محاسبه می‌شود:

$$r = \frac{\|\nabla_w \mathcal{L}\|_2}{\|w\|_2 + \epsilon}$$

اگر این نسبت از یک آستانه λ بیشتر باشد، گرادیان آن لایه به صورت زیر کوچک می‌شود:

$$\nabla_w \mathcal{L} \leftarrow \nabla_w \mathcal{L} \times \min \left(1, \frac{\lambda}{r} \right)$$

به عبارت ساده‌تر: «تغییر وزن در هر قدم باید بیشتر از λ برابر مقدار فعلی وزن باشد.» این کار باعث می‌شود: - وزن‌ها هرگز تغییرات بسیار شدید نداشته باشند \Rightarrow انفجار گرادیان غیرممکن می‌شود. - در عین حال گرادیان‌های مفید و معقول حفظ می‌شوند (برخلاف برش ثابت که ممکن است گرادیان‌های بزرگ اما مفید را هم قطع کند).

نتایج شگفت‌انگیز مقاله

- مدل NFNet-F6 بدون هیچ نوع نرمال‌سازی (نه GroupNorm، نه LayerNorm، نه BatchNorm) به دقت 1% Top-1 در ImageNet 88.5% رسانید که در زمان انتشار رکورد جدیدی بود.
- آموزش پایدار حتی با $Batch\ Size = 1^6$ یا کمتر.
- سرعت آموزش تا ۸۰٪ بالاتر از مدل‌های معادل با BatchNorm.
- مصرف حافظه به شدت کاهش یافت.
- شبکه‌های تا بیش از ۱۰۰۰ لایه بدون مشکل آموزش داده شدند.

مزایای روش AGC و NFNet

- حذف کامل تمام لایه‌های نرمال‌سازی \Rightarrow کد ساده‌تر، سرعت بالاتر، حافظه کمتر
- کاملاً مستقل از سایز بچ \Rightarrow ایده‌آل برای آموزش توزیع شده و دستگاه‌های با حافظه محدود
- امکان استفاده از نرخ یادگیری بسیار بالا
- جلوگیری همزمان و مؤثر از هر دو مشکل Vanishing و Exploding Gradient
- دستیابی به دقت State-of-the-Art بدون هیچ ترفند اضافی
- قابل تعمیم به معماری‌های دیگر (بعداً روی Transformerها و ViT‌ها هم موفق بوده)

محدودیت‌ها و نکات عملی

- نیاز به تنظیم دقیق λ (معمولًاً بین ۰.۱ تا ۰.۱۰)
- محاسبه نرم وزن در هر لایه هزینه محاسباتی جزئی اضافه می‌کند (ولی بسیار کمتر از BatchNorm)
- در مقاله اصلی فقط روی شبکه‌های کانولوشنی تست شده بود (هرچند مقالات بعدی نشان دادند روی Transformerها هم عالی کار می‌کند)

جمع‌بندی نهایی

مشکلات Exploding Gradient و Vanishing Gradient از مهم‌ترین موانع آموزش شبکه‌های عمیق بودند که با Normalization تا حد زیادی حل شدند. اما مقاله ICML 2021 از DeepMind Adaptive Gradient Clipping نشان داد که می‌توان بدون هیچ نوع نرمال‌سازی و فقط با یک خط کد ساده، شبکه‌های عمیق‌تر، سریع‌تر و دقیق‌تری ساخت. این روش امروزه یکی از قدرتمندترین و مدرن‌ترین تکنیک‌های پایدارسازی آموزش شبکه‌های عصبی عمیق محسوب می‌شود و به سرعت در جامعه تحقیقاتی و صنعتی در حال گسترش است.