

الگوریتم‌های بهینه‌سازی روش گرادیان نزولی*

چکیده: برای تحلیل الگوریتم‌های بهینه‌سازی کاهش گرادیان، با وجود محبوبیت فزاینده، از آنالیز جعبه سیاه (BlackBox) استفاده می‌شود زیرا ارائه‌ی توضیح عملی از نقاط ضعف و قوت آن‌ها دشوار است. هدف این مقاله ایجاد درکی از رفتار الگوریتم‌های متفاوت برای خواننده است تا بتواند از آن‌ها بهره بگیرد. در طول این مقاله، به فرم‌های مختلف روش کاهش گرادیان نگاه می‌کنیم، چالش‌ها را به طور مختصر بیان می‌کنیم، پرکاربردترین الگوریتم‌های بهینه‌سازی را معرفی می‌کنیم، معماری سیستم‌های موازی و توزیع شده را مرور می‌کنیم و به دنبال روش‌های جدیدی برای بهینه‌سازی روش کاهش گرادیان خواهیم گشت.

۱ مقدمه

روش گرادیان نزولی یکی از محبوب‌ترین الگوریتم‌ها برای بهینه‌سازی و با فاصله‌ی زیاد، پرکاربردترین روش برای بهینه‌سازی شبکه‌های عصبی است. در عین حال، پیشرفته‌ترین کتابخانه‌ی یادگیری عمیق ماشین، شامل نمونه‌هایی از پیاده‌سازی الگوریتم‌های متفاوت برای بهینه‌سازی روش گرادیان نزولی است (مثل متن‌های راهنمای لازانیا^۱، کافه^۲ و کراس^۳). با اینحال از این الگوریتم‌ها برای بهینه‌سازی با روش جعبه سیاه (BlackBox) استفاده می‌شود زیرا ارائه‌ی توضیح عملی از نقاط ضعف و قوت آنها مشکل است.

هدف این مقاله ایجاد درکی از رفتار الگوریتم‌های مختلف بهینه‌سازی روش گرادیان نزولی برای خواننده است تا بتواند از آنها استفاده کند. در بخش ۲، ابتدا به بررسی فرم‌های مختلف روش گرادیان نزولی می‌پردازیم. سپس به طور خلاصه چالش‌های این الگوریتم را در خلال آموزش‌های بخش ۳ مرور می‌کنیم. پس از آن در بخش ۴، معمول‌ترین الگوریتم‌های بهینه‌سازی را معرفی می‌کنیم و نشان می‌دهیم که چگونه تلاش برای حل چالش‌های بهینه‌سازی، سبب ایجاد تغییراتی در رویکردهای استفاده شده در این الگوریتم‌ها شده است. سپس، در بخش ۵، نگاهی کوتاه به الگوریتم‌ها و معماری‌هایی برای بهینه‌سازی گرادیان نزولی در سیستم‌های موازی و توزیع شده می‌اندازیم. در نهایت، در بخش ۶ رویکردهای تازه‌ای که برای بهینه‌سازی گرادیان نزولی مفید هستند را در نظر می‌گیریم.

* این مقاله در اصل به صورت یک مطلب در یک وبلاگ به آدرس <http://sebastianruder.com/optimizing-gradient-descent/index.html> در ۱۹ ژانویه ۲۰۱۶ منتشر شده است.

^۱ Lasagne's Documentation: <http://lasagne.readthedocs.org/en/latest/modules/updates.html>

^۲ Caffe's Documentation: <http://lasagne.readthedocs.org/en/latest/modules/updates.html>

^۳ Keras' Documentation: <http://caffe.berkeleyvision.org/tutorial/solver.html>

روش گرادیان نزولی راهی برای یافتن مینیمم یک تابع هدف مثل $J(\theta)$ است که بر اساس بردار پارامترهای مدل یعنی $\theta \in \mathbb{R}^d$ نوشته شده است؛ به این صورت که پارامترها در خلاف جهت گرادیان تابع نسبت به پارامترها یعنی $\nabla_{\theta} J(\theta)$ تغییر می‌کنند. نرخ یادگیری شبکه عصبی یا η اندازه‌ی گام تغییر پارامترها برای یافتن مینیمم (یا مینیمم نسبی) را تعیین می‌کند. به عبارت دیگر، در جهت شیب سطح ایجاد شده توسط تابع به سمت پایین حرکت می‌کنیم تا به یک دره برسیم.^۴

۲ فرم‌های مختلف روش گرادیان نزولی

سه فرم از روش گرادیان نزولی وجود دارند که در حجم داده‌ای که برای محاسبه‌ی گرادیان تابع هدف استفاده می‌شود با هم تفاوت دارند. بسته به مقدار داده، تعادلی بین دقت تغییر پارامتر و مدت زمانی که برای تغییر پارامتر نیاز است ایجاد می‌کنیم.

۲.۱ گرادیان نزولی انبوه

روش گرادیان نزول ساده، که به نام گرادیان نزولی انبوه هم شناخته می‌شود، گرادیان تابع هزینه نسبت به بردار پارامتر θ برای کل مجموعه داده‌های ورودی سیستم محاسبه می‌شود:

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

برای اینکه برای ایجاد هر تغییر باید گرادیان کل داده‌ها محاسبه شود، روش گرادیان نزولی انبوه می‌تواند بسیار کند باشد و برای مجموعه‌ای از داده‌ها که به قدری بزرگند که در حافظه جا نمی‌شوند قابل قبول نیست. درضمن این روش به ما اجازه نمی‌دهد که مدل را /حظه‌ای تغییر دهیم؛ به این معنا که نمی‌توان داده‌های جدید را حین کار به آن اضافه کرد.

در کدنویسی، روش گرادیان نزولی انبوه چیزی شبیه به این است:

for i in range (nb_epochs):

params_grad = evaluate_gradient (loss_function, data, params)

params = params – learning_rate * params_grad

برای تعداد تکرار مشخص، ابتدا بردار گرادیان یعنی params_grad را نسبت به بردار پارامترها یعنی params برای تابع ضرر محاسبه می‌کنیم. دقت کنید که جدیدترین کتابخانه‌های یادگیری ژرف شامل دیفرانسیل گیری خودکار هستند که گرادیان نسبت به پارامترهای مشخص را به صورت بهینه محاسبه می‌کنند. اگر گرادیان را خودتان حساب کرده اید، بهتر است که صحت گرادیان‌ها را بررسی کنید.^۵

^۴ اگر با روش گرادیان نزولی آشنایی ندارید، مطالب خوبی در مورد بهینه سازی سیستم های عصبی را می توانید در اینجا بیابید:

<http://cs231n.github.io/optimization-1>

^۵ برای چند توصیه عالی برای بررسی درست صحت گرادیان گیری، به آدرس زیر مراجعه کنید:

<http://cs231n.github.io/neural-networks-3/>

سپس پارامترها را در جهت گرادیان به اندازه‌ی نرخ یادگیری که تعیین می‌کند تغییرات ایجاد شده چقدر بزرگ باشند، تغییر می‌دهیم. روش گرادیان نزولی در سطوح محدب حتماً به مینیمم مطلق و در سطوح غیر محدب حتماً به مینیمم نسبی همگرا می‌شود.

۲.۲ گرادیان نزولی تصادفی

روش گرادیان نزولی تصادفی (Stochastic Gradient Descent: SGD) برخلاف روش قبلی، نسبت به هرکدام از داده‌ها مثل $X^{(i)}$ و مقدار متناظر آن $y^{(i)}$ یک گام تغییر ایجاد می‌کند:

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)})$$

روش گرادیان نزولی انبوه محاسبات اضافی و نالازم برای مجموعه داده‌های بزرگ انجام می‌دهد زیرا گرادیان داده‌های یکسانی را به طور مکرر قبل از ایجاد هر تغییر محاسبه می‌کند. روش تصادفی این محاسبات اضافی را با ایجاد یک تغییر با هر گرادیان‌گیری حذف می‌کند. به همین دلیل هم این روش بسیار سریع‌تر است و می‌تواند یادگیری لحظه‌ای هم داشته باشد. روش تصادفی تغییرات را با واریانس بالایی ایجاد می‌کند که باعث می‌شود تابع هدف مثل شکل ۱ به شدت نوسان کند.

در حالیکه روش گرادیان انبوه همواره به مینیمم همان قوسی که پارامترها را در بر می‌گیرد همگرا می‌شود، نوسانات روش تصادفی از یک طرف این روش را قادر می‌سازد که مینیمم‌های نسبی بهتری را بیابد و از طرف دیگر باعث می‌شود که درنهایت همگرایی به مقدار دقیق مینیمم سخت‌تر شود زیرا این روش در هر تکرار انحراف زیادی نسبت به مقدار اصلی دارد. با اینحال، ثابت شده که وقتی نرخ یادگیری را به آرامی زیاد می‌کنیم، این روش همان همگرایی روش انبوه را از خود نشان می‌دهد و با قابلیت اطمینان بالایی، در سطوح محدب به مینیمم مطلق و در سطوح غیر محدب به مینیمم نسبی همگرا می‌شود. تکه کد مربوط به این الگوریتم خیلی ساده یک حلقه‌ی تکرار (loop) در داده‌ها ایجاد کرده و گرادیان نسبت به هر نمونه را حساب می‌کند. دقت کنید همانطور که در بخش ۱.۶ توضیح داده شده، در هر تکرار، ترتیب داده‌ها به صورت تصادفی تغییر داده می‌شود.

```
for i in range ( nb_epochs):
```

```
    np . random . shuffle ( data)
```

```
    for example in data:
```

```
        params_grad = evaluate_gradient ( loss_function, example, params )
```

```
        params = params - learning_rate*params_grad
```

۲.۳ گرادیان نزولی نیمه انبوه

در نهایت روش گرادیان نزولی نیمه انبوه بهترین ویژگی‌های هر دو روش را باهم ترکیب کرده و به ازای گرادیان هر زیرمجموعه n عضوی از نمونه‌ها، یک گام تغییر ایجاد می‌کند:

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)})$$

با اینکار، الف) واریانس تغییرات پارامترها کمتر می‌شود، که باعث می‌شود نرخ همگرایی سریع‌تر شود (ب) می‌تواند از روش‌های بهینه‌سازی ماتریس که در جدیدترین کتابخانه‌های یادگیری ژرف به سادگی یافت می‌شوند و به شدت بهینه هستند برای محاسبه‌ی بسیار بهینه‌ی گرادیان نسبت زیرمجموعه‌ها بهره بگیرد. اندازه‌ی معمول زیرمجموعه‌ها بین ۵۰ و ۲۵۶ نمونه است، ولی می‌تواند برای کاربردهای مختلف اندازه‌های متفاوتی داشته باشد. در آموزش دادن شبکه‌های عصبی به طور معمول از روش گرادیان نزولی نیمه انبوه استفاده می‌شود و عبارت SGD معمولاً برای این روش هم به کار برده می‌شود. نکته‌ی مهم: در ادامه‌ی متن در جاهایی که فرمول‌های مربوط به SGD آمده‌اند، برای سادگی پارامترهای $x^{(i:i+n)}$ و $y^{(i:i+n)}$ نوشته نشده‌اند.

در کدنویسی، به جای تکرار کردن حلقه برای کل داده‌ها، حلقه در یک زیرمجموعه با اندازه‌ی ۵۰ تکرار می‌شود:

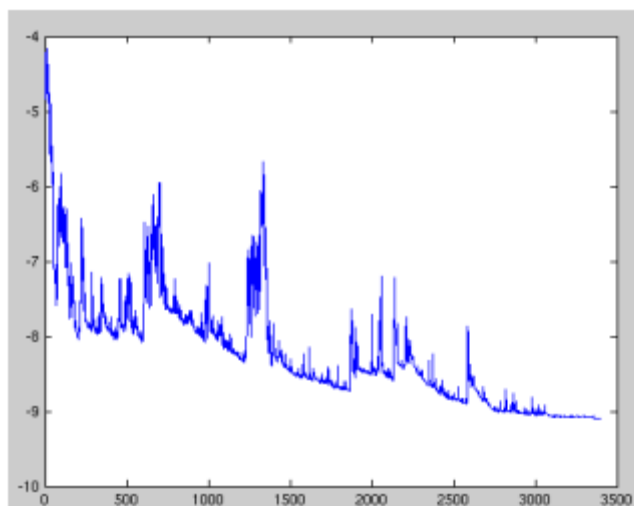
```
for i in range ( nb_epochs ) :
```

```
    np . random . shuffle (data)
```

```
    for batch in get_batches ( data, batch_size=50 )
```

```
        params_grad = evaluate_gradient ( loss_function, batch, params )
```

```
        params = params – learning_rate * params_grad
```



شکل 1 تناوب شدید روش تصادفی (منبع: ویکی‌پدیا)

با وجود موارد گفته شده، روش گرادیان نزولی نیمه انبوه ساده، همگرایی سریع را تضمین نمی‌کند، و در این زمینه چالش‌هایی دارد که باید مورد توجه قرار گیرند:

- انتخاب نرخ یادگیری مناسب می‌تواند سخت باشد. اگر نرخ یادگیری زیادی کوچک باشد باعث می‌شود همگرایی به شدت کند باشد، و از طرف دیگر اگر زیادی بزرگ باشد برای همگرایی مزاحمت ایجاد کرده و باعث شود تابع ضرر در اطراف مقدار مینیمم نوسان کند یا حتی واگرا شود.
- روندهای نرخ یادگیری [۱۸]، در طول آموزش سیستم عصبی، سعی می‌کنند با روش‌هایی مثل الگوریتم تبرید نرخ یادگیری را اصلاح کنند. الگوریتم تبرید، نرخ یادگیری را بر اساس روند تعریف شده، یا زمانی که تغییرات تابع هدف بین دو تکرار کمتر از آستانه‌ی مشخصی می‌شود، کاهش می‌دهد. با اینحال روند استفاده شده و آستانه‌ی تغییرات، باید از قبل مشخص شده باشند و بنابراین نمی‌توان آن‌ها را متناسب با ویژگی‌های مجموعه داده‌ها تغییر داد. [۴]
- علاوه بر اینها، یک نرخ یادگیری ثابت برای همه‌ی تغییرات ایجاد شده در پارامترها استفاده می‌شود. اگر داده‌ها خیلی پراکنده باشند و ویژگی‌های مورد بررسی تعداد رخدادهای متفاوتی داشته باشند، نباید برای همه‌ی آنها را به یک اندازه تغییر دهیم، بلکه باید ویژگی‌هایی که به ندرت رخ داده‌اند را با گام بزرگتری تغییر دهیم.
- یکی دیگر از چالش‌های کلیدی مینیمم کردن توابع ضرر پرکاربرد در سیستم‌های عصبی که به شدت غیرمحدب هستند، پیشگیری از گیر کردن در مینیمم‌های نسبی است که نسبت به مینیمم‌های اصلی بهینه نیستند. داوفین و همکاران [۵] عنوان کرده‌اند که چالش اصلی مینیمم‌های نسبی نیستند بلکه نقاط زینی‌اند؛ یعنی نقاطی که شیب در یکی از جهت‌ها مثبت و در جهت دیگر منفی است. این نقاط زینی معمولاً با سطحی صاف پر از نقاط مشابه احاطه شده‌اند که خروج از آن‌ها را برای روش SGD دشوار می‌کند زیرا گرادیان در همه‌ی جهات نزدیک صفر است.

۴ الگوریتم‌های بهینه‌سازی با روش گرادیان نزولی

در ادامه، مبانی اصلی برخی الگوریتم‌ها را که به طور گسترده‌ای توسط افراد فعال در حوزه‌ی یادگیری شرف برای غلبه بر چالش‌های عنوان شده استفاده می‌شوند را بیان می‌کنیم. الگوریتم‌هایی مثل روش‌های مرتبه دو از جمله روش نیوتون که محاسبات آن‌ها در عمل برای مجموعه داده‌های چند بعدی عملی نیست را در اینجا مطرح نمی‌کنیم.^۶

۴.۱ روش شتاب (گشتاور)

روش SGD برای حرکت در دره‌ها، یعنی نواحی که در آن فرو رفتگی سطح در یک جهت بسیار بیشتر از جهت‌های دیگر است [۲۰] خوب عمل نمی‌کند در حالیکه که چنین نواحی در اطراف نقاط اکسترمم بسیار معمول هستند. در چنین مواردی، روش

^۶ https://en.wikipedia.org/wiki/Newton%27s_method_in_optimization

SGD بین دو شیب دره تناوب می‌کند و همانطور که در شکل ۲ الف نشان داده شده بسیار کند به سمت اکسترمم موضعی حرکت می‌کند.



(الف) SGD بدون استفاده از روش شتاب



(ب) SGD با استفاده از روش شتاب

شکل ۲ منبع: جنویو ب. اور

روش شتاب [۱۷] روشی است که به الگوریتم SGD در جهت مناسب سرعت می‌بخشد و تناوب ها را میرا می‌کند. این عملکرد در شکل ۲ ب نشان داده شده است. روش شتاب این کار را با اضافه کردن کسری از بردار تغییرات مرحله‌ی قبل یعنی γ به بردار تغییرات فعلی انجام می‌دهد.^۷

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta = \theta - v_t$$

ضریب شتاب یعنی γ معمولا ۰/۹ یا همین حدود در نظر گرفته می‌شود.

اساسا وقتی از روش شتاب استفاده می‌کنیم، مثل این است که یک توپ را از یک تپه به پایین هل می‌دهیم. توپ حین قل خوردن به پایین تپه شتاب می‌گیرد و در راه سریع‌تر و سریع‌تر می‌شود (تا جایی که به سرعت نهایی برسد، البته در صورتی که مقاومت هوا وجود داشته باشد که به معنی γ کمتر از ۱ است). همین اتفاق برای تغییرات ایجاد شده در پارامترها هم می‌افتد: ضریب شتاب در جهت هایی که گرادیان نسبت به آنها تغییر علامت نمی‌دهد افزایش می‌یابد و در جهت هایی که گرادیان تغییر علامت می‌دهد کاهش می‌یابد. در نتیجه، همگرایی سریع‌تر شده و تناوب‌ها کمتر می‌شود.

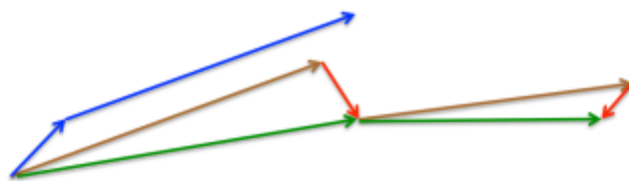
۴.۲ گرادیان شتاب یافته‌ی نستروف

^۷ در بعضی جاها علامت ضرایب در معادله متفاوت است.

با اینحال، توپی که روی یک تپه قل می‌خورد و تنها شیب روی حرکت آن اثر می‌گذارد، چندان رضایت بخش نیست. ما می‌خواهیم توپی هوشمندتر داشته باشیم، توپی که بداند قرار است به کجا برسد داشته باشد تا قبل از اینکه دوباره به سربالایی تپه برسد سرعتش را کم کند.

روش گرادیان شتاب یافته‌ی نستروف (Nesterov Accelerated Gradient: NAG) [۱۴] روشی است که با آن می‌توان چنی خاصیتی به ضریب شتاب داد. می‌دانیم که از ضریب شتاب به صورت γv_{t-1} برای حرکت بردار پارامتر θ استفاده می‌کنیم. بنابراین محاسبه‌ی $\theta - \gamma v_{t-1}$ تقریبی از مقدار بعدی بردار پارامترها را به دست می‌دهد (گرادیان که در میزان تغییرات موثر است در اینجا در نظر گرفته نشده)؛ با اینکار درکی مناسب از مقادیر بعدی پارامترها به دست می‌آوریم. اکنون می‌توانیم با محاسبه گرادیان نسبت به مقدار تقریبی بعدی پارامترها، نه مقدار فعلی آن‌ها، به طور موثری مقدار بعدی داده‌ها را در محاسبات دخیل کنیم:

$$\begin{aligned} v_t &= \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1}) \\ \theta &= \theta - v_t \end{aligned}$$



شکل ۳ تغییرات نستروف (منبع: اسلاید 6c ج. هینتون)

باز هم به ضریب شتاب γ مقداری حدود ۰/۹ می‌دهیم. درحالی‌که روش شتاب ابتدا گرادیان فعلی را حساب کرده (بردار کوچک آبی در شکل ۳) و سپس جهشی بزرگ در جهت گرادیان جدید تجمع یافته (بردار قهوه‌ای) می‌کند، روش نستروف ابتدا جهشی بزرگ در جهت گرادیان تجمع یافته‌ی قبلی می‌کند (بردار قهوه‌ای)، گرادیان را اندازه می‌گیرد و سپس یک اصلاح کوچک انجام می‌دهد (بردار سبز). این تغییرات ایجاد شده با استفاده از پیش بینی مقادیر بعدی، از زیادی سرعت گرفتن تغییرات جلوگیری می‌کند. این امر باعث پاسخگویی بیشتر الگوریتم می‌شود که کارایی شبکه‌هایی عصبی بازگشتی در انجام بسیاری از فعالیت‌ها را به شدت بهبود داده است. [۲]^۸

اکنون که می‌توانیم تغییرات را با شیب تابع ضرر منطبق کنیم، و لذا سرعت روش SGD را بهبود ببخشیم، بد نیست که تغییرات را با تک تک پارامترها هم انطباق دهیم تا با توجه به میزان اهمیت آنها، تغییرات بزرگتر یا کوچتری بدهیم.

^۸ برای توضیحات بیشتر پیرامون مفاهیم پشت روش نستروف به لینک زیر مراجعه کنید. ایلیا استاتسکور هم مروری با جزئیات بیشتر در تز دکترای خود نسبت به این مفهوم ارائه کرده است. [۱۹]

AdaGrad [۸] الگوریتمی برای بهینه سازی بر پایه‌ی گرادیان است که از این روش استفاده می‌کند: نرخ یادگیری را با پارامترها انطباق می‌دهد، به این معنی که به ازای پارامترهایی که میزان رخداد کمتری دارند تغییرات بزرگتری داده و به ازای پارامترهایی که بیشتر رخ داده اند تغییرات کمتری می‌دهد. به همین دلیل، برای محاسبات روی داده‌های پراکنده بسیار مناسب است. دین و همکاران [۶] دریافته‌اند که روش AdaGrad اطمینان پذیری روش SGD را به شدت افزایش داده است و از این روش برای آموزش دادن شبکه‌های عصبی گسترده‌ای در گوگل استفاده کردند که —علاوه بر کارهای دیگر— یاد گرفت که گربه‌ها را در ویدئوهای یوتیوب تشخیص دهد.^۹ بعلاوه، پنینگتون و همکاران [۱۶] از این روش برای یاد دادن تکنیک "کلمه تعبیه شده" (الگوریتمی برای تشخیص جملات توهین آمیز حتی اگر شامل کلمات توهین آمیز نباشند) به سیستم GloVe شدند زیرا در این روش کلمات کم استفاده تر باید تاثیر بیشتری نسبت به کلمات پر استفاده داشته باشند.

در روش های قبلی، ما روی همه‌ی پارامترها یک تغییر یکسان را اعمال می‌کردیم زیرا همه‌ی پارامترها از یک نرخ یادگیری η استفاده می‌کردند. از آنجایی که روش AdaGrad نرخ یادگیری متفاوتی برای هر پارامتر θ_i با هر گام زمانی t استفاده می‌کند، ابتدا تغییرات هر پارامتر مستقل در AdaGrad را نشان می‌دهیم، سپس پارامترها را برداری می‌کنیم. برای اختصار، $g_{t,i}$ را گرادیان تابع هدف نسبت به θ_i و در لحظه‌ی t در نظر می‌گیریم:

$$g_{t,i} = \nabla_{\theta_i} J(\theta_{t,i})$$

پس تغییرات ایجاد شده توسط تابع SGD بر روی هر پارامتر θ_i در هر لحظه‌ی t برابر خواهد بود با:

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

در اینجا $G_t \in \mathbb{R}^{d \times d}$ یک ماتریس قطری است که در آن هر درایه‌ی قطری $\hat{1}, \hat{1}$ جمع مربعات گرادیان‌ها نسبت به θ_i در گام زمانی t است^{۱۰}، و ϵ یک ضریب هموارسازی است که از ایجاد صفر در مخرج جلوگیری می‌کند (معمولا در مقیاس 10^{-8} است). جالب است که بدون جذر گرفتن، الگوریتم عملکرد بسیار ضعیف تری دارد.

از آنجایی که قطر اصلی ماتریس G_t شامل جمع مربعات گرادیان‌های قبلی نسبت به همه‌ی پارامترها یعنی بردار θ است، می‌توانیم فرمول را با یک ضرب درایه به درایه‌ی ماتریس و بردار \odot بین G_t و g_t تبدیل به یک فرمول برداری کنیم:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

^۹ <http://www.wired.com/2012/06/google-x-neural-network/>

^{۱۰} دوچی و همکاران [۸] این ماتریس را به عنوان جایگزینی برای ماتریس کامل شامل ضرب خارجی تمام گرادیان‌های قبلی ارائه می‌کند زیرا در غیر این صورت محاسبه‌ی ماتریس جذر حتی با تعداد کم پارامترها یعنی d هم عملی نیست.

یکی از مزایای اصلی روش AdaGrad این است که نیاز به تعیین نرخ یادگیری را از بین می برد. در بیشتر کاربردها مقدار 0.01 انتخاب شده و تغییر داده نمی شود.

ضعف اصلی روش AdaGrad وجود ماتریس تجمع مربعات گرادیانها در مخرج است: از آنجایی که هر عبارت اضافه شده مثبت است، حاصل جمع همواره در طول آموزش سیستم عصبی افزایش می یابد. این امر به نوبه ی خود باعث می شود که نرخ یادگیری کاهش یابد و درنهایت، بی نهایت کوچک شود که پس از آن الگوریتم نمی تواند چیز جدیدی یاد بگیرد. الگوریتم های زیر به هدف از بین بردن این نقص به وجود آمده اند.

۴.۴ AdaDelta

روش AdaDelta [۲۲] بسطی از روش AdaGrad است که هدف آن کمتر کردن کاهش فزاینده و یکنواخت نرخ یادگیری در این روش است. به جای تجمع مربعات همه ی گرادیان های قبلی، روش AdaDelta تعداد مربعات جمع شده را به تعداد محدودی مثل W محدود می کند.

به جای ذخیره ی غیربهرینه ی W تا از مربعات گرادیانها، جمع گرادیان به صورت بازگشتی، میانگین میراشونده ی همه ی مربعات گرادیان های قبلی تعریف می شود. میانگین در لحظه ی t یعنی $E[g^2]_t$ فقط (با یک کسر γ که شبیه به ثابت شتاب است) به میانگین قبلی و گرادیان فعلی بستگی دارد:

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2$$

به γ مقداری شبیه به ثابت سرعت یعنی حدود 0.9 می دهیم. برای شفاف تر شدن موضوع، میزان تغییرات پارامترها روش SGD پایه را بر اساس بردار تغییرات یعنی $\Delta\theta_t$ می نویسیم:

$$\begin{aligned}\Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$

پس بردار تغییرات روش AdaGrad که قبلا به دست آوردیم برابر است با:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

اکنون به سادگی در این فرمول ماتریس قطری G_t را با میانگین میراشونده ی مربعات گرادیان های قبلی یعنی $E[g^2]_t$ جایگزین می کنیم:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

از آنجایی که مخرج همان معیار خطای جذر میانگین مربعات گرادیان است، می توانیم آن را با علامت اختصاری این معیار جایگزین کنیم:

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t$$

سازندگان این روش متوجه شدند که واحد بردار تغییرات (همانطور که در روش‌های SGD، روش شتاب و AdaGrad هم وضع به همین صورت است) با پارامترها همخوانی ندارد؛ چرا که تغییرات باید واحد فرضی یکسانی با پارامترها داشته باشند. برای حل این مشکل، ابتدا میانگین دیگری تعریف کردند که آن هم به صورت نمایی میراشونده بود. اما این بار از میانگین مربعات استفاده نکردند بلکه از میانگین تغییرات پارامترها استفاده کردند:

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta_t^2$$

بنابراین خطای جذر میانگین مربعات تغییرات پارامترها به صورت زیر است:

$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$$

از آنجایی که $RMS[\Delta\theta]_t$ معلوم نیست، آن را با مقدار RMS تغییرات پارامترها تا گام زمانی قبلی، تقریب می‌زنیم. با جایگزین کردن نرخ یادگیری η در فرمول تغییرات پارامتر قبلی با $RMS[\Delta\theta]_{t-1}$ در نهایت فرمول تغییرات برای روش AdaDelta به دست می‌آید:

$$\begin{aligned}\Delta\theta_t &= -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$

با روی AdaDelta، حتی نیاز نیست نرخ یادگیری پیشفرض را تعیین کنیم، زیرا از فرمول تغییرات حذف شده است.

۴.۵ RMSprop

روش RMSprop یک روش بر مبنای نرخ یادگیری انطباق یافته است که در جایی منتشر نشده و جاف هینتون در درس‌هایی که در وبسایت Coursera ارائه کرده آن را معرفی کرده است.^{۱۱}

^{۱۱} http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

هم RMSprop و هم AdaDelta به طور مستقل و تقریباً در یک زمان در پی نیاز به راه حلی برای حل مشکل کاهش شدید نرخ یادگیری در روش AdaGrad به وجود آمدند. روش RMSprop در واقع کاملاً مشابه اولین بردار تغییرات روش AdaDelta که در بالا به دست آوردیم است:

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

روش کسر RMS (RMSprop) نیز نرخ یادگیری را بر میانگین نمایی میراشوندهی مربعات گرادیانها تقسیم می‌کند. هینتون پیشنهاد کرده است که مقدار γ برابر 0.9 باشد و یک مقدار پیشفرض مناسب برای نرخ یادگیری یا η برابر 0.001 است.

۴.۶ Adam

روش تخمین تطبیق پذیر گشتاور (Adam) [۱۰] یکی دیگر از روش‌هایی است که نرخ یادگیری را منطبق بر داده‌ها محاسبه می‌کند. علاوه بر ذخیره کردن میانگین میراشوندهی مربعات گرادیانهای قبلی، یعنی v_t مثل روش‌های AdaDelta و RMSprop، روش Adam میانگین میراشوندهی نمایی گرادیانها m_t را هم مثل روش شتاب حفظ می‌کند:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

m_t و v_t به ترتیب تخمین گشتاور اول (میانگین) و گشتاور دوم (واریانس غیرمرکزی) گرادیانها هستند که نام‌گذاری روش هم به همین خاطر است. به خاطر این که حالت اولیه ی بردارهای m_t و v_t صفر است، سازندگان روش مشاهده کردند که نتایج به صفر متمایل می‌شوند، به خصوص در گام‌های اولیه و مخصوصاً وقتی نرخ میراشوندگی کوچک است (یا به عبارت دیگر β_1 و β_2 نزدیک ۱ هستند).

آنها این مشکل را با تخمین‌های اصلاح شدهی گشتاور اول و دوم حل کردند:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

سپس از این دو فرمول همانطور که در روش‌های AdaDelta و RMSprop هم دیدیم، برای محاسبه‌ی تغییرات پارامترها استفاده کردند که فرمول تغییرات را برای این روش به دست می‌دهد:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

سازندگان روش مقادیر پیشفرض β_1 برای 0.9 ، β_2 برای 0.999 و ϵ را برای 10^{-8} پیشنهاد کرده اند. آنها به صورت تجربی نشان داده اند که روش Adam در عمل به خوبی کار می کند و نسبت به دیگر روش های یادگیری انطباق پذیر ارجحیت دارد.

۴.۷ AdaMax

ضریب v_t در فرمول تغییرات روش Adam گرادیان را با نسبت معکوس نسبت به میانگین گرادیان های قبلی یعنی l_2 (با ضریب v_{t-1}) و گرادیان فعلی مربوط می کند.

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) |g_t|^2$$

می توانیم این تغییر را به میانگین l_p هم تعمیم دهیم. دقت کنید که کینگما و با، ضریب β_2 را هم به صورت پارامتری یعنی β_2^p نوشته اند:

$$v_t = \beta_2^p v_{t-1} + (1 - \beta_2^p) |g_t|^p$$

میانگین برای مقادیر بزرگ p معمولاً از لحاظ عددی ناپایدار می شود. به همین خاطر است که در عمل میانگین های l_1 و l_2 پرکاربردتر هستند. با این وجود، l_∞ نیز معمولاً از خود پایداری نشان می دهد. به همین دلیل، سازندگان روش AdaMax را پیشنهاد می کنند [۱۰] و نشان می دهند که v_t نیز مانند l_∞ به مقدار پایدارتر زیر میل می کند. برای جلوگیری از اشتباه گرفتن متغیرها با روش Adam، ما برای نشان دادن حد در بینهایت v_t از عبارت u_t استفاده کرده ایم:

$$\begin{aligned} u_t &= \beta_2^\infty v_{t-1} + (1 - \beta_2^\infty) |g_t|^\infty \\ &= \max(\beta_2 \cdot v_{t-1}, |g_t|) \end{aligned}$$

اکنون می توانیم با جایگزین کردن u_t با $\sqrt{\hat{v}_t} + \epsilon$ این معادله را با معادله ی تغییرات در روش آدام ادغام کنیم تا فرمول تغییرات پارامترها برای روش AdaMax به دست آید:

$$\theta_{t+1} = \theta_t - \frac{\eta}{u_t} \hat{m}_t$$

دقت کنید از آنجایی که u_t حاصل از یک عملیات ماکسیمم گیری است، مثل m_t و v_t در روش Adam به صفر همگرا نمی شود و نیازی به محاسبه ی ضریب تصحیح برای جلوگیری از این امر نیست. مقادیر پیشفرض مناسب برای این روش هم عبارتند از: 0.002 برای η ، مقدار 0.9 برای β_1 و مقدار 0.999 برای β_2 .

۴.۸ Nadam

همانطور که قبلا هم دیدیم، روش Adam را می‌توان ترکیبی از روش RMSprop و روش شتاب دانست: میانگین میراث‌شونده‌ی مربعات گرایان‌های قبلی یعنی ∇_t از روش RMSprop و میانگین میراث‌شونده‌ی نمایی میانگین گرایان‌های قبلی یعنی m_t . هم چنین دیدیم که روش گرایان شتاب یافته‌ی نستروف نیز به روش شتاب معمولی ارجحیت دارد.

روش Nadam (Nesterov-accelerated Adaptive Moment Estimation) [۷] روش Adam و نستروف را با هم ترکیب می‌کند. برای استفاده از روش نستروف در الگوریتم Adam، باید عبارت m_t را اصلاح کنیم.

ابتدا، بایید فرمول تغییرات در روش شتاب را با ضرایب جدید بنویسیم:

$$\begin{aligned}g_t &= \nabla_{\theta_t} J(\theta_t) \\m_t &= \gamma m_{t-1} + \eta g_t \\\theta_{t+1} &= \theta_t - m_t\end{aligned}$$

که در آن J همان تابع هدف، γ ثابت میرایی شتاب و η اندازه‌ی گام است. با ادغام معادله‌ی دم در معادله‌ی سوم داریم:

$$\theta_{t+1} = \theta_t - (\gamma m_{t-1} + \eta g_t)$$

این فرمول نیز نشان می‌دهد که در روش شتاب یک گام در جهت بردار شتاب قبلی و سپس یک یک گام در جهت گرایان فعلی برداشته می‌شود.

روش نستروف به ما اجازه می‌دهد که با تغییر دادن پارامترها و گشتاورها پیش‌از محاسبه‌ی گرایان گامی دقیق‌تر در جهت گرایان برداریم. برای به کارگیری روش نستروف فقط باید گرایان ∇_t را اصلاح کنیم:

$$\begin{aligned}g_t &= \nabla_{\theta_t} J(\theta_t - \gamma m_{t-1}) \\m_t &= \gamma m_{t-1} + \eta g_t \\\theta_{t+1} &= \theta_t - m_t\end{aligned}$$

دو‌بارت پیشنهاد کرده است که روش نستروف به این صورت اصلاح شود: به جای دوبار دخیل کردن گشتاور در محاسبات – یک بار برای به تغییر دادن گرایان ∇_t و یک بار برای تغییر دادن پارامترها یعنی بردار θ_{t+1} – مستقیماً از تخمین بردار گشتاور برای محاسبه‌ی تغییرات پارامترها استفاده می‌کنیم:

$$\begin{aligned}g_t &= \nabla_{\theta_t} J(\theta_t) \\m_t &= \gamma m_{t-1} + \eta g_t \\\theta_{t+1} &= \theta_t - (\gamma m_t + \eta g_t)\end{aligned}$$

دقت کنید که به جای استفاده از بردار گشتاور قبلی یعنی m_{t-1} در فرمول تغییرات قبلی، اکنون از بردار گشتاور m_t یعنی بردار گشتاور فعلی برای رسیدن به مقادیر بعدی استفاده می‌کنیم. به طور مشابه برای اضافه کردن روش شتاب نستروف به الگوریتم Adam، می‌توانیم بردار گشتاور قبلی را با بردار گشتاور فعلی جایگزین کنیم. ابتدا، بیایید فرمول تغییرات روش Adam را به یاد بیاوریم (دقت کنید که نیاز داریم \hat{v}_t را اصلاح کنیم):

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \end{aligned}$$

با جایگذاری کردن فرمول \hat{m}_t و m_t در فرمول آخر داریم:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left(\frac{\beta_1 m_{t-1}}{1 - \beta_1^t} + \frac{(1 - \beta_1) g_t}{1 - \beta_1^t} \right)$$

توجه داشته باشید که $\frac{\beta_1 m_{t-1}}{1 - \beta_1^t}$ تنها یک تخمین اصلاح‌شده از بردار گشتاور در گام زمانی قبلی است. لذا می‌توانیم آن را با \hat{m}_{t-1} جایگزین کنیم:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left(\beta_1 \hat{m}_{t-1} + \frac{(1 - \beta_1) g_t}{1 - \beta_1^t} \right)$$

۴.۹ تجسم تصویری الگوریتم‌ها

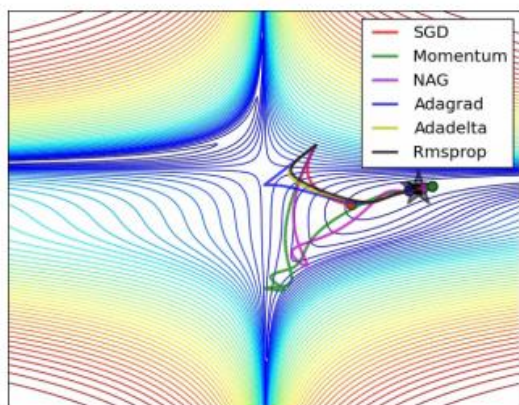
دو شکل زیر درکی از عملکرد بهینه‌سازی الگوریتم‌های معرفی شده ارائه می‌کند.^{۱۲}

در شکل ۴ الف، مسیری الگوریتم‌ها بر روی نمای دو بعدی سطح ضرر (تابع Beale) نشان داده شده است. همه‌ی الگوریتم‌ها از یک نقطه شروع کرده‌اند و مسیرهای متفاوتی را برای رسیدن به مینیمم در پیش گرفته‌اند. توجه داشته باشید که روش AdaGrad، AdaDelta، RMSprop و سریعا به جهت درست رفته‌اند و نسبتاً سریع همگرا شده‌اند در حالیکه روش‌های شتاب و نستروف از مسیر اصلی منحرف شده‌اند که حرکت یک توپ در دامنه‌ی یک تپه را تداعی می‌کند. با اینحال روش نستروف

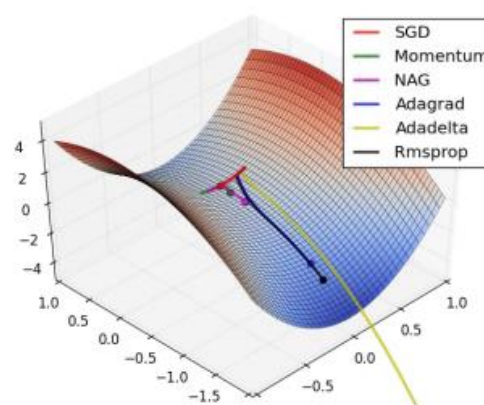
^{۱۲} به لینک زیر هم نگاهی بیندازید. در این لینک کارپائی توضیحات دیگری برای همین تصاویر و مرور مختصر دیگری بر الگوریتم‌های ذکر شده ارائه کرده است.

قادر بوده است که مسیرش را به خاطر حساسیت بیشتر نسبت به مسیر و استفاده از اطلاعات نقطه‌ی بعدی، قادر بوده است که سریع‌تر مسیرش را اصلاح کند.

شکل ۴ ب رفتار الگوریتم‌ها را در نقطه‌ی زینی نشان می‌دهد، یعنی نقطه‌ای که شیب در یک جهت مثبت و در جهت دیگر منفی است. این نقطه برای روش SGD همانطور که قبلاً اشاره کردیم، دشواری‌هایی ایجاد کرده است. توجه کنید که روش‌های SGD، شتاب و نستروف برای خارج شدن از سطح متقارن اطراف این نقطه با دشواری هستند، اگرچه دو الگوریتم آخر توانسته‌اند از نقطه‌ی زینی خارج شوند، در حالیکه AdaGrad، RMSprop، AdaDelta و سریعا به سمت شیب منفی سرازیر شده‌اند، و روش AdaDelta در بین این روش‌ها پیشگام بوده است.



(الف) نمای دو بعدی بهینه‌سازی SGD بر روی سطح ضرر



(ب) بهینه‌سازی SGD در نقطه‌ی زینی

شکل ۴ منبع و انیمیشن‌های کامل: الک ردفورد

همانطور که دیده می‌شود، روش‌های با نرخ یادگیری انطباق یافته، یعنی AdaGrad، AdaDelta، RMSprop و Adam از همه مناسب‌تر هستند و در این مثال‌ها بهترین همگرایی را از خود نشان می‌دهند.

۴.۱۰ کدام روش بهتر است؟

با این اوصاف، باید از کدام روش بهینه‌سازی استفاده کرد؟ اگر داده‌ی ورودی پراکنده است، پس احتمالا بهترین نتایج با استفاده از یکی از روش‌های با نرخ یادگیری انطباق یافته به دست می‌آورد. یکی از مزیت‌های اضافی این روش‌ها این است که دیگر نیاز نیست نرخ یادگیری را تنظیم کنید ولی احتمالا بهترین نتایج را با مقادیر پیش‌فرض به دست خواهید آورد.

به طور خلاصه، روش RMSprop تعمیمی از روش AdaGrad است که کاهش شدید نرخ یادگیری در این روش را حل کرده است. این روش کاملا شبیه به روش AdaDelta است، با این تفاوت که روش AdaDelta از RMS تغییرات پارامترها در صورت فرمول تغییرات استفاده کرده است. در نهایت، روش Adam روشی برای تصحیح میل به صفر و گشتاور را به روش RMSprop اضافه کرده است. تا اینجا کار، RMSprop، AdaDelta و Adam الگوریتم‌های بسیار مشابهی هستند

که در شرایط مشابه همگی به خوبی کار می‌کنند. کینگما و همکاران [۱۰] نشان داده‌اند که تصحیح میل به صفر در الگوریتم Adam باعث شده که کارایی آن در پایان بهینه سازی که گرادیان ها پراکنده تر می‌شوند کمی از RMSprop بیشتر باشد. پس تا به حال، روش Adam بهترین انتخاب در بین همه‌ی الگوریتم‌هاست.

نکته‌ی جالب اینجاست که بیشتر مقالات اخیر از روش SGD معمولی بدون روش شتاب و یک روند تیرید نرخ یادگیری ساده استفاده می‌کنند. همانطور که نشان داده شد، روش SGD معمولاً در یافتن یک مینیمم موفق است، ولی ممکن است به طور قابل ملاحظه‌ای بیش از سایر الگوریتم‌ها طول بکشد، به مقادیر اولیه‌ای که برای داده‌های مختلف کارایی را تضمین کنند وابسته‌تر است و ممکن است به جای رسیدن به مینیمم نسبی در نقطه‌ی زینی گیر بیفتد. در نتیجه، اگر همگرایی سریع برایتان مهم است یا شبکه عصبی عمیق یا پیچیده‌ای را آموزش می‌دهید، باید از یکی روش‌های با نرخ یادگیری انطباق یافته استفاده کنید.

اجرای SGD به صورت موازی و توزیع شده

با فراگیری راه حل های بر پایه‌ی داده های بزرگ مقیاس برای چالش های صنعتی و ظهور خوشه‌های پردازش ابری با تعداد کاربر کم، توزیع الگوریتم SGD به منظور افزایش سرعت آن رویکرد مسلمی است. روش SGD به خودی خود ذاتاً ترتیبی است: قدم به قدم، به سمت مینیمم حرکت می‌کنیم. استفاده از این الگوریتم همگرایی خوبی فراهم می‌کند ولی می‌تواند به خصوص برای مجموعه داده‌های بزرگ کند باشد. در عوض، اجرای غیرهمزمان (آسنکرون) SGD سریع‌تر است، ولی ارتباط غیر بهینه بین اجزای اجرای کننده‌ی الگوریتم می‌تواند باعث همگرایی ضعیف شود. به علاوه، می‌توانیم SGD را تنها بر روی یک دستگاه بدون نیاز به خوشه‌های پردازشی بزرگ هم به صورت موازی اجرا کنیم. در زیر الگوریتم‌ها و معماری‌هایی که برای اجرای موازی و توزیع‌شده‌ی SGD پیشنهاد شده‌اند آمده‌اند.

۵.۱ HOGWILD!

نیو و همکاران [۱۵] روندی برای ایجاد تغییرات به اسم HOGWILD! را معرفی کرده‌اند که اجرای تغییرات SGD به صورت موازی در پردازنده را امکان‌پذیر می‌کند. هسته‌های پردازشی مجازند بدون ثابت نگه داشتن پارامترها به حافظه‌ی اشتراکی دسترسی داشته باشند. این روش تنها در صورتیکه کارآمد است که داده‌های ورودی پراکنده باشند، زیرا هر تغییر تنها بر روی کسری از داده‌ها پیاده می‌شود. آن‌ها نشان دادند که در این صورت، روند تغییرات تقریباً به نرخ همگرایی بهینه‌ای می‌رسد زیرا احتمال پاک کردن و بازنویسی داده‌های مفید توسط پردازنده‌ها بسیار پایین است.

۵.۲ روش SGD بارشی

روش SGD بارشی یک فرم غیرهمزمان (آسنکرون) از روش SGD است که توسط دین و همکاران [۶] در فریم‌ورک DistBelief (که TensorFlow از روی آن ساخته شده است) در گوگل از آن استفاده شده است. این الگوریتم چندین کپی از مدل را بر روی زیرمجموعه‌های داده‌های آموزش سیستم به صورت موازی اجرا می‌کند. این مدل‌ها تغییرات خود را به سرور پارامترها ارسال می‌کنند، که بین دستگاه‌های زیادی پخش شده است. هر دستگاه مسئول ذخیره و تغییر دادن کسری از پارامترهای

مدل است. با اینحال، چون کپی های مدل با یکدیگر ارتباط برقرار نمی کنند؛ یعنی مثلاً وزن داده ها یا تغییرات را به اشتراک نمی گذارند، پارامترهای آنها همواره در معرض خطر واگرایی یا همگرایی کند هستند.

۵.۳ الگوریتم های مقاوم در برابر تاخیر برای SGD

مک ماهان و استریتر [۱۲] الگوریتم AdaGrad را با توسعه ی الگوریتم های مقاوم در برابر تاخیر، به سیستم های موازی تعمیم دادند. این الگوریتم ها نه تنها بر گرادین های قبلی، بلکه با تاخیر در تغییرات هم انطباق می یابند. این روش عملکرد خوبی را در عمل نشان داده است.

۵.۴ TensorFlow

TensorFlow^{۱۳}، [۲۳] فریم ورک گوگل برای طراحی و پیاده سازی مدل های بزرگ مقیاس یادگیری ماشین است که اخیراً برای استفاده ی عموم آزاد شده است. این فریم ورک بر اساس تجربه ی گوگل با DistBelief است هم اکنون به صورت داخلی در بسیاری از دستگاه های همراه برای انجام پردازش های مختلف، هم چنین در سیستم های توزیع شده ی بزرگ مقیاس استفاده می شود. نسخه ی توزیع یافته، که در آوریل ۲۰۱۶ منتشر شده است^{۱۴}، به یک گراف محاسباتی وابسته است که به چند زیرگراف در دستگاه ها تقسیم شده است، و ارتباط بین دستگاه ها با جفت گره های ارسال/دریافت برقرار می شود.

۵.۵ SGD با میانگین گیری فنی

ژانگ و همکاران [۲۳] روش SGD با میانگین گیری فنی (Elastic Averaging SGD:EASGD) را پیشنهاد کرده اند، که پارامترهای دستگاه های دخیل در SGD غیرهمزمان (آسنکرون) را به یک نیروی فنی مرتبط می کند، به این معنا که همان متغیر مرکزی ذخیره شده در سرور پارامترهاست. این به متغیرهای محلی (Local) اجازه می دهد که فراتر از متغیر مرکزی تناوب کنند، که روی کاغذ جستجوی بیشتر در فضای برداری پارامترها را ممکن می کند. آنها به صورت تجربی نشان دادند که این افزایش امکان جستجو به یافتن اسکترم های نسبی جدید و در نتیجه عملکرد بهتر منتهی می شود.

۶ رویکردهای بیشتر برای بهینه سازی SGD

در نهایت، ما رویکردهای بیشتری معرفی می کنیم که می توانند در کنار همه ی الگوریتم هایی که قبلاً نام برده شدند استفاده شوند تا عملکرد SGD را بهبود دهند. برای مروری عالی بر روش های معمول دیگر، به مرجع [۱۱] مراجعه کنید.

۶.۱ یادگیری تصادفی و ترتیبی

^{۱۳} <https://www.tensorflow.org/>

^{۱۴} <http://googlesearch.blogspot.ie/2016/04/announcing-tensorflow-08-now-with.htm>

ما عموماً ترجیح می‌دهیم که نمونه‌های آموزشی را با یک ترتیب معنادار به مدل ندهیم زیرا این امر ممکن است در الگوریتم بهینه‌سازی اختلال ایجاد کند. لذا معمولاً توصیه می‌شود که پس از هر تکرار، چینش داده‌ها به صورت تصادفی تغییر کند.

از طرف دیگر، برای بعضی موارد که می‌خواهیم مسائلی که به طور فزاینده دشوارند را حل کنیم، دادن نمونه‌های آموزشی با یک ترتیب معنادار به سیستم، ممکن است واقعا به بهبود عملکرد و همگرایی منجر شود. روش ایجاد این ترتیب معنادار یادگیری ترتیبی (Curriculum Learning) نامیده می‌شود.

زارمبا و ساتسکور [۲۱] برای ارزیابی یادگیری ترتیبی در برنامه‌های ساده تنها توانستند حافظه‌های طولانی کوتاه مدت (LSTM) را با این روش آموزش دهند. آنها نشان دادند که رویکردی ترکیبی که نمونه‌ها را با افزایش دشواری مساله مرتب می‌کند، نسبت به رویکرد خام بهتر است.

۶.۲ نرمالیزاسیون بچ

برای تسهیل یادگیری، معمولاً مقادیر اولیه‌ی پارامترها را با میانگین صفر و واریانس واحد نرمال می‌کنیم. همینطور که آموزش شبکه عصبی جلوتر می‌رود و پارامترها را به مقادیر جدیدی تغییر می‌دهیم، نرمال‌سازی انجام شده از دست می‌رود که باعث کند شدن یادگیری و شدت یافتن تغییرات با عمیق تر شدن شبکه می‌شود.

نرمالیزاسیون بچ [۹] هر دسته از داده‌ها دوباره نرمال‌سازی کرده و تغییرات را نیز به عقب بازمی‌گرداند. با تبدیل نرمال‌سازی به بخشی از معماری مدل، می‌توانیم از نرخ‌های یادگیری بزرگتری استفاده کنیم و کمتر درگیر مقادیر اولیه‌ی پارامترها شویم. نرمالیزاسیون بچ هم چنین مثل یک رگیولایزر عمل می‌کند و نیاز به دراپ اوت (Dropout) را کم می‌کند (یا به کلی از بین می‌برد).

۶.۳ توقف پیش از موعد

بر اساس گفته‌های جاف هینتون، "توقف پیش از موعد، یک رایگان و زیبا است"^{۱۵} بنابراین باید همیشه باید خطا را با یک مجموعه داده‌ی صحت سنجی در طول آموزش بسنجید و (با به خرج دادن کمی صبر) اگر خطا به اندازه‌ی کافی زیاد نشد، آموزش را متوقف کنید.

۶.۴ نويز گراديان

^{۱۵} اسلایدهای آموزشی NIPS سال ۲۰۱۵، اسلاید ۶۳، قابل دسترسی در لینک زیر:

<http://www.iro.umontreal.ca/~bengioy/talks/DL-Tutorial-NIPS2015.pdf>

نیلاکاتان و همکاران [۱۳] به هر تغییرات گرادیان، نویز با توزیع گاوسی $N(0, \sigma_t^2)$ اضافه کردند:

$$g_{t,i} = g_{t,i} + N(0, \sigma_t^2)$$

آنها از روند تبرید زیر برای واریانس استفاده کردند:

$$\sigma_t^2 = \frac{\eta}{(1+t)^\nu}$$

آنها نشان داده اند که اضافه کردن این نویز شبکه‌ها را نسبت به مقادیر اولیه نامناسب مقاوم‌تر کرده و به خصوص به آموزش شبکه‌های عمیق و پیچیده کمک می‌کند. آن‌ها حدس زده‌اند که نویز به مدل برای فرار و یافتن مینیمم‌های نسبی جدید، که در مدل‌های عمیق‌تر تعدادشان بسیار بیشتر است، شانس بیشتری می‌دهد.

۷ جمع‌بندی

در این مقاله، ابتدا سه فرم از روش گرادیان نزولی را مرور کردیم، که از بین آنها روش گرادیان نزولی نیمه انبوه از همه محبوب‌تر است. الگوریتم‌هایی را بررسی کردیم که برای بهینه‌سازی SGD بیشترین استفاده را دارند: روش شتاب، گرادیان شتاب‌یافته‌ی نستروف، AdaGrad، AdaDelta، RMSprop، Adam، AdaMax، Nadam، و همچنین الگوریتم‌های متفاوتی برای بهینه‌سازی SGD غیرهمزمان (آسنکرون). درنهایت، رویکردهای دیگری برای بهبود روش SGD مثل یادگیری تصادفی و ترتیبی، نرمالیزاسیون بچ تو توقف پیش از موعد را نیز بررسی کردیم.