



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

دانشکده ریاضی و علوم کامپیوتر

گزارش پروژه نهایی درس جبرخطی عددی

موضوع پروژه  
گرادیان کاهشی طبیعی

نگارش  
رقیه مهدوی لاین

استاد  
دکتر مهدی دهقاز  
بهار ۱۴۰۴

## چکیده

در این تحقیق سه روش مهم بهینه‌سازی شامل گرادیان با تندترین شیب، گرادیان مزدوج و گرادیان طبیعی مورد بررسی قرار گرفته‌اند. ابتدا روش گرادیان با تندترین شیب به عنوان ساده‌ترین الگوریتم کاهش تابع هدف معرفی شده و نحوه‌ی جهت‌گیری آن در راستای منفی گرادیان تابع هدف توضیح داده شده است. سپس روش گرادیان مزدوج، که با در نظر گرفتن حافظه‌دار بودن جهت‌ها و استفاده از اطلاعات پیشین گرادیان به سرعت همگرایی بالاتری دست می‌یابد، تحلیل شده است. در ادامه، روش گرادیان طبیعی بررسی شده که با تکیه بر ساختار هندسی فضای پارامتر و استفاده از متریک Fisher، جهت واقعی‌ترین کاهش را در فضاها پیچیده مانند شبکه‌های عصبی ارائه می‌دهد. در این تحقیق مزایا، پایداری عددی، سرعت همگرایی و کاربردهای هر روش مقایسه شده و نقاط قوت و ضعف آنها تحلیل شده‌اند.

2.....	چکیده
4.....	مقدمه
4.....	مروری بر بهینه سازی محدب
4.....	مروری بر روش گرادیان کاهشی
5.....	مشکل اصلی
5.....	تعریف بد حالتی
5.....	تعریف عدد حالت ماتریس
5.....	ماتریس هسین
7.....	نمونه تأثیر ماتریس هسین بد حالت
9.....	راهکار های پیشنهادی
9.....	توضیح فرم درجه دوم (Quadratic Form)
12.....	روش گرادیان با بیشترین شیب (The method of Steepest Descent)
15.....	پایداری روش گرادیان کاهشی با بیشترین شیب
16.....	روش گرادیان مزدوج
20.....	آنالیز پایداری گرادیان مزدوج
20.....	مقایسه پیچیدگی
21.....	روش گرادیان طبیعی
23.....	پایداری گرادیان طبیعی
23.....	نتیجه
23.....	مقایسه
24.....	تحلیل و آنالیز پیاده سازی
25.....	مراجع

## مقدمه

### مروری بر بهینه سازی محدب

مسئله بهینه سازی محدب به صورت زیر تعریف می شود:

$$\begin{aligned} &\text{minimize } f_0(x) \\ &\text{subject to } f_i(x) \ll b_i \quad i = 1, \dots, m \end{aligned}$$

که در آن توابع  $R \rightarrow R^n: f_0, f_1, \dots, f_m$  همگی محدب هستند به این معنا که شرط زیر را ارضا میکنند.

$$f_i(\alpha x + \beta y) \ll \alpha f_i(x) + \beta f_i(y)$$

برای همه بردارهای  $x, y \in R^n$  و ضرایب  $\alpha, \beta \in R$  که جمع آنها برابر با یک است، یعنی  $\alpha + \beta = 1$ ، و نیز  $\alpha, \beta \geq 0$ . مسائل کمترین مربعات و برنامه ریزی خطی هر دو نمونه هایی خاص از این قالب کلی بهینه سازی محدب هستند.

### مروری بر روش گرادیان کاهشی

در طول اجرای الگوریتم گرادیان کاهشی (Gradient Descent)، دستگاه به صورت تکراری نقطه ی بعدی به روزرسانی شده را با استفاده از گرادیان تابع هدف در موقعیت فعلی محاسبه می کند. این به روزرسانی از طریق کم کردن حاصل ضرب نرخ یادگیری در گرادیان نسبت به پارامتر فعلی صورت می گیرد. رابطه ی این به روزرسانی به صورت زیر بیان می شود:

$$p_{n+1} = p_n - \eta \nabla f(p_n)$$

$$\nabla f = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]^T$$

که در آن:

- $\eta$  نرخ یادگیری
- $p_n$  پارامتر مودر بهینه سازی در مرحله  $n$ -ام است
- $\nabla f(p_n)$  گرادیان تابع هزینه مورد انتظار در نقطه ی  $p_n$  را نشان می دهد.

## مشکل اصلی

### تعریف بد حالتی

یک ماتریس زمانی **بد حالت** گفته می‌شود که تغییرات جزئی در داده‌های ورودی منجر به تغییرات بسیار بزرگ در خروجی شود.

### تعریف عدد حالت ماتریس

برای ماتریس شرطی  $A$  به مقدار عددی  $\|A\| \cdot \|A^{-1}\|$  عدد شرطی یا عدد وضعیت یا عدد حالت (Condition number)

ماتریس  $A$  گفته می‌شود و با نماد  $\kappa(A)$  نمایش داده می‌شود.  $\kappa(A) = \|A\| \cdot \|A^{-1}\|$ .

فرض کنید  $\lambda_{(max)}$  و  $\lambda_{(min)}$  به ترتیب بزرگترین و کوچکترین مقدار ویژه ماتریس متقارن  $A$  باشند. آنگاه:

$$\kappa_2(A) = \frac{\lambda_{(max)}}{\lambda_{(min)}}$$

### ماتریس هسین

ماتریس هسین یک ماتریس مربعی از مشتق‌های مرتبه دوم یک تابع اسکالر است. اگر تابع  $f: R \rightarrow R^n$  باشد آنگاه ماتریس

هسین  $H$  به صورت زیر تعریف می‌شود:

$$H = (\nabla f)' = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

این ماتریس که در اینجا با  $H$  نشان داده شده است، به ماتریس هسین تابع  $f$  معروف است و درایه‌های آن به صورت زیر تعریف

می‌شوند:

$$H_{i,j} = \frac{\partial^2 f}{\partial x_j \partial x_i} = \frac{\partial^2 f}{\partial x_i \partial x_j} = H_{j,i}$$

این که می‌توان مشتق‌های جزئی را به هر ترتیبی محاسبه کرد، یک نکته از «تقارن مشتق‌های جزئی مرتبه دوم» یا «برابری مشتق‌های

جزئی مختلط» شناخته می‌شود و به این معناست که ماتریس هسین یک ماتریس متقارن است، یعنی:  $H = H^T$

بسیاری از کاربردهای مهم مشتق‌های مرتبه دوم، ماتریس‌های هسین، و تقریب‌های درجه دوم در زمینه بهینه‌سازی ظاهر می‌شوند؛

یعنی در مسائلی مانند کمینه‌سازی (یا بیشینه‌سازی) توابع  $f(x)$ .

در مسئله یافتن مینیمم (یا ماکزیمم) محلی یک تابع پیچیده  $f(x)$ ، یک رویکرد متداول ریاضی آن است که مقدار  $f(x + \delta x)$

را برای بردار جابه‌جایی کوچک  $\delta x$  با یک تابع ساده‌تر موسوم به "مدل" تقریب زده و سپس آن مدل را بهینه کنیم تا یک گام بهینه‌سازی

پیشنهاد شود.

به طور مشخص استفاده از تقریب:

$$f(x + \delta x) \approx f(x) + f'(x)[\delta x]$$

که یک مدل خطی (affine) به شمار می‌رود، منجر به روش‌های گرادیان کاهشی و سایر الگوریتم‌های مشابه می‌شود.

با این حال تقریب دقیق‌تر  $f(x + \delta x)$  می‌تواند به الگوریتم‌هایی با نرخ همگرایی سریع‌تر بینجامد. از این رو، یک ایده‌ی طبیعی آن است که از مشتق دوم  $f''$  برای ساخت یک مدل درجه دوم استفاده کنیم، یعنی:

$$f'(x + \delta x) \approx f'(x) + f''(x)[\delta x]$$

که در فضای  $\mathbb{R}^n$  به صورت  $\delta(\nabla f) \approx H\delta x$  بیان می‌شود؛ که در آن  $H$  ماتریس هسین (Hessian) است.

از این رو، کمینه‌سازی مدل درجه دوم متناظر با گام نیوتن است:

$$\delta x \approx -H^{(-1)}\nabla f$$

که تلاشی است برای یافتن ریشه‌ی  $\nabla f$  بر اساس تقریب مرتبه اول..

با این حال دقت و پایداری این تقریب به شرایط عددی ماتریس هسین  $H$  بستگی دارد. در شرایطی که ماتریس هسین بد حالت (*ill-conditioned*) باشد، عدد شرطی آن (*condition number*) بسیار بالا باشد، محاسبه‌ی

$$\delta x \approx -H^{(-1)}\nabla f$$

با خطای عددی قابل توجهی مواجه می‌شود. در این حالت، جهت و اندازه‌ی گام پیشنهادی  $\delta x$  می‌تواند به صورت ناپایدار و نادقیق تخمین زده شود.

به عبارت دیگر، هنگامی که  $H$  بد شرط است، بردار گرادیان  $\nabla f$  ممکن است در جهاتی با خمیدگی زیاد (بردار ویژه‌های متناظر با مقادیر ویژه بزرگتر) مؤثرتر از جهاتی با خمیدگی کم باشد، و همین موضوع باعث می‌شود مشتق دوم (و در نتیجه مدل درجه دوم) رفتار نامناسبی در مقیاس‌های مختلف فضا داشته باشد. این موضوع در عمل به الگوریتم‌هایی منجر می‌شود که در مسیر بهینه‌سازی دچار زیگ زاگ و نوسانات شدید می‌شوند و نرخ همگرایی بسیار کاهش می‌یابد.

بنابراین، بدحالتی هسین نه تنها محاسبه‌ی گام نیوتنی را بی‌ثبات می‌کند، بلکه کیفیت تقریب مشتق دوم را نیز به شدت تحت تأثیر قرار می‌دهد، چرا که بردار جابه‌جایی  $\delta x$  به شدت نسبت به خطای عددی حساس خواهد بود. به همین دلیل، در این شرایط استفاده از روش‌هایی چون گرادیان مزدوج (Conjugate Gradient) یا گرادیان طبیعی (*Natural Gradient*) که این بدشرطی را اصلاح یا به نحوی پیش شرط‌گذاری (Preconditioning) می‌کنند، می‌تواند به بهبود پایداری و سرعت همگرایی بینجامد.

نمونه تأثیر ماتریس هسین بد حالت

تابع اول:  $f(x, y) = \frac{1}{2}(x^2 + y^2)$

گرادیان:  $\nabla f(x, y) = \begin{bmatrix} x \\ y \end{bmatrix}$

ماتریس هسین تابع:  $H = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

عدد حالت ماتریس:  $\kappa_2(A) = \frac{\lambda_{(max)}}{\lambda_{(min)}} = \frac{1}{1} = 1$  ماتریس بسیار خوش حالت هست.

نقطه شروع:  $(x_0, y_0) = (2, 2)$

نرخ یادگیری:  $\alpha = 0.1$

تکرار ۱

$$\nabla f = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \quad (x_1, y_1) = (2, 2) - 0.1 \times \begin{bmatrix} 2 \\ 2 \end{bmatrix} = (1.8, 1.8)$$

تکرار ۲

$$\nabla f = \begin{bmatrix} 1.8 \\ 1.8 \end{bmatrix}, \quad (x_2, y_2) = (1.8, 1.8) - 0.1 \times \begin{bmatrix} 1.8 \\ 1.8 \end{bmatrix} = (1.62, 1.62)$$

تکرار ۳

$$\nabla f = \begin{bmatrix} 1.62 \\ 1.62 \end{bmatrix}, \quad (x_3, y_3) = (1.62, 1.62) - 0.1 \times \begin{bmatrix} 1.62 \\ 1.62 \end{bmatrix} = (1.458, 1.458)$$

...

هر گام به طور یکنواخت و مستقیم به سوی مبدأ، که نقطه‌ی مینیمم است، حرکت می‌کند. همگرایی سریع و پایدار است.

تابع دوم:  $f(x, y) = \frac{1}{2}(x^2 + 1000y^2)$

گرادیان:  $\nabla f(x, y) = \begin{bmatrix} x \\ 1000y \end{bmatrix}$

ماتریس هسین تابع:  $H = \begin{bmatrix} 1 & 0 \\ 0 & 1000 \end{bmatrix}$

عدد حالت ماتریس:  $\kappa_2(A) = \frac{\lambda_{(max)}}{\lambda_{(min)}} = \frac{1000}{1} = 1000$  ماتریس بسیار بد حالت هست.

نقطه شروع:  $(x_0, y_0) = (2, 2)$

نرخ یادگیری:  $\alpha = 0.001$  به دلیل انحنای زیاد، نرخ یادگیری باید بسیار کوچک باشد.

تکرار ۲

### تکرار ۳

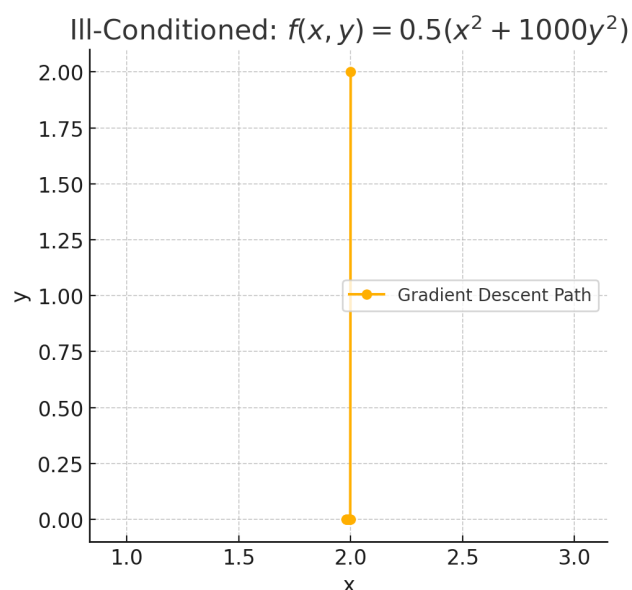
$$\nabla f = \begin{bmatrix} 1.996002 \\ 0 \end{bmatrix}, \quad (x_3, y_3) = (1.996002, 0) - 0.001 \times \begin{bmatrix} 1.996002 \\ 0 \end{bmatrix} = (1.994006, 0)$$

■■■■■■

Well-Conditioned:  $f(x, y) = 0.5(x^2 + y^2)$

Gradient Descent Path

x	y
0.75	0.70
0.85	0.77
0.95	0.86
1.05	0.96
1.15	1.07
1.25	1.18
1.35	1.30
1.45	1.46
1.60	1.62
1.80	1.80
2.00	2.00



8



## راهکار های پیشنهادی

### توضیح فرم درجه دوم (Quadratic Form)

بیایید با چند تعریف و توضیح درباره‌ی نمادگذاری‌ها شروع کنیم. الگوریتم‌هایی که در ادامه به آن می‌پردازیم در واقع برای حل دستگاه‌های معادلات به صورت زیر به کار می‌روند:

$$Ax = b$$

در این معادله،  $x$  بردار مجهول است که باید آن را بیابیم،  $b$  بردار معلوم و  $A$  یک ماتریس مربعی، متقارن و معین مثبت (یا نیمه معین مثبت) است. چنین دستگاه‌هایی در بسیاری از مسائل مهم مهندسی و علمی پدیدار می‌شوند؛ از جمله در روش‌های اختلاف محدود و اجزای محدود برای حل معادلات دیفرانسیل جزئی، تحلیل سازه‌ها، تحلیل مدارها و حتی تمرین‌های ریاضی.

ماتریس  $A$  معین مثبت است اگر به ازای هر بردار ناصفر  $x$  داشته باشیم:

$$x^T A x > 0$$

**فرم درجه دوم (Quadratic Form)** یک تابع اسکالر درجه دوم از یک بردار است که به صورت زیر نوشته می‌شود:

$$f(x) = \frac{1}{2} x^T A x - b^T x + c \quad (1)$$

که در آن  $A$  یک ماتریس، و  $x$  و  $b$  بردارهایی هستند. اگر ماتریس  $A$  متقارن و معین مثبت (positive-definite) باشد، آنگاه این تابع در نقطه‌ای کمینه می‌شود که همان پاسخ دستگاه معادلات خطی زیر است:

$$Ax = b$$

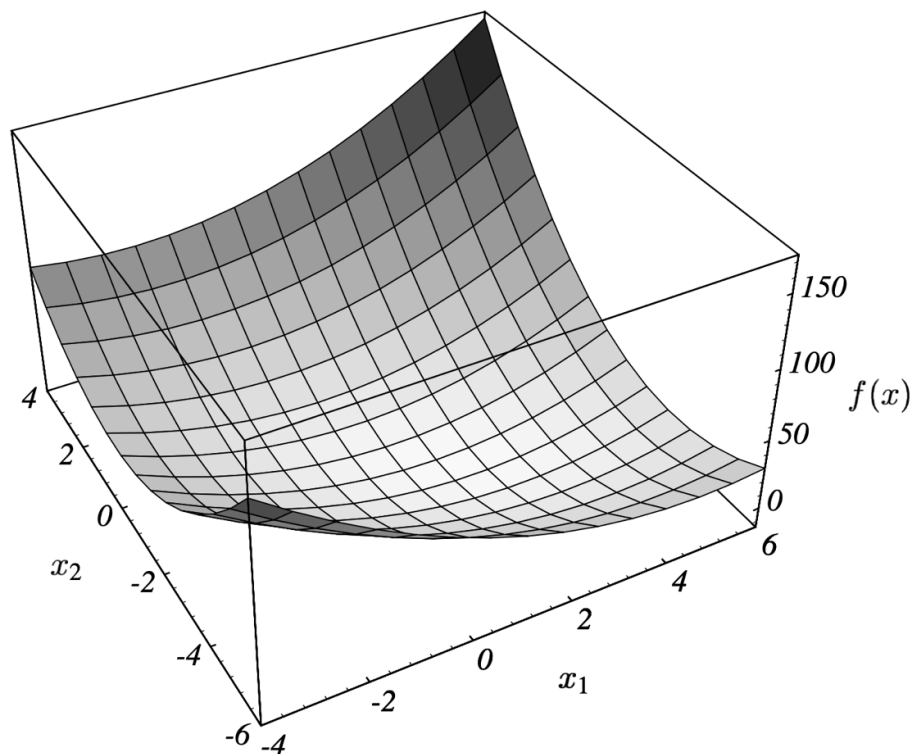
این ایده را با یک مسئله‌ی نمونه‌ی ساده نمایش خواهیم داد.

$$A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}, \quad b = \begin{bmatrix} 2 \\ -8 \end{bmatrix}, \quad c = 0$$

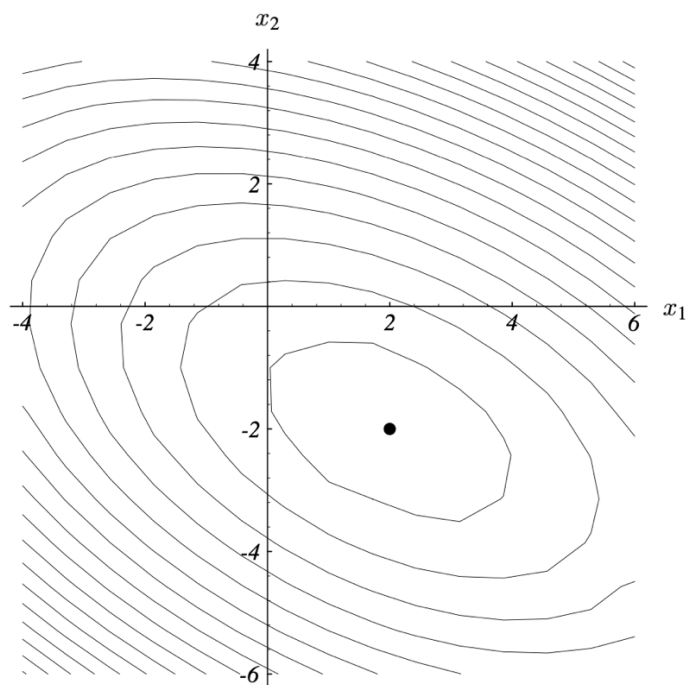
دستگاه  $Ax = b$  در شکل ۱ نمایش داده شده است. به‌طور کلی، جواب  $x$  در نقطه‌ی تقاطع  $n$  ابرصفحه (hyperplane) قرار

دارد، که هر یک بعد  $n - 1$  دارند. در این مسئله، جواب برابر است با  $x = [2, -2]^T$ . فرم درجه دوم متناظر  $f(x)$  در شکل ۲

نمایش داده شده است. نمودار کانتور (خطوط تراز)  $f(x)$  تابع نیز در شکل ۳ ارائه شده است.



شکل (۱) نمودار شکل درجه دوم تابع  $f(x)$  می‌باشد، نقطه مینیمم این سطح همان پاسخ  $Ax = b$  است.

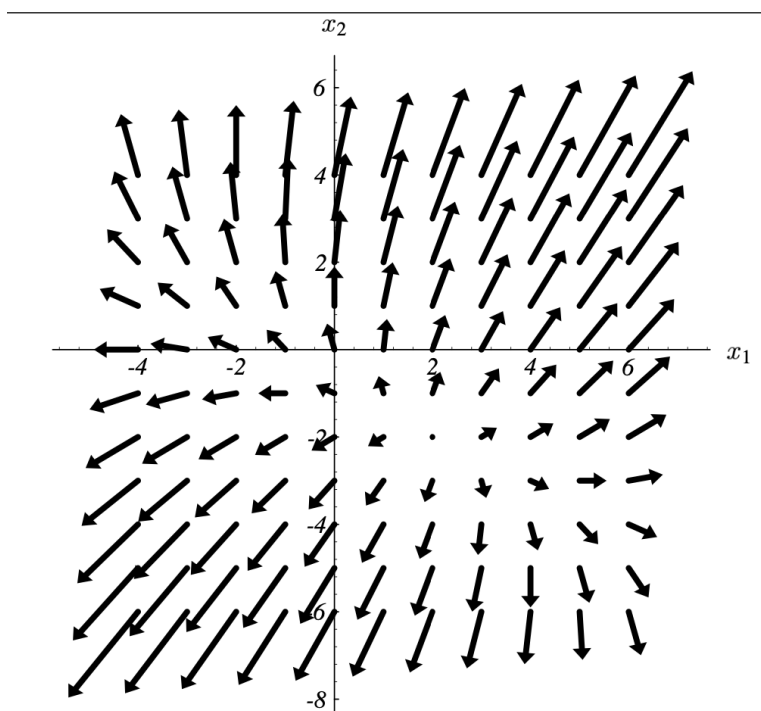


شکل (۲) کانتورهای فرم درجه دوم. هر منحنی بیضی شکل دارای مقدار ثابتی از  $f(x)$  است.

نکته: نمودار کانتور، منحنی‌هایی در صفحه دوبعدی (برای دو متغیر) هستند که در آن مقدار تابع  $f(x)$  برابر با یک مقدار ثابت

$$f(x) = c \quad \text{است: } c$$

برای تابع درجه دوم بالا، این معادله معمولاً منجر به بیضی‌هایی در صفحه می‌شود. مرکز این بیضی‌ها نقطه‌ای است که  $f(x)$  در آن مینیمم می‌شود (یعنی جواب دستگاه)



شکل (۳) گرادیان  $f'(x)$  برای تابع درجه دوم، جهت بیشترین افزایش مقدار تابع  $f(x)$  را در هر نقطه  $x$  نشان می‌دهد و همواره بر خطوط کانتور عمود است.

نکته: نقطه مرکزی که همه فلش‌ها از آن دور می‌شوند (یا به آن نزدیک می‌شوند)، همان نقطه بهینه (کمینه) تابع و در نتیجه جواب بهینه دستگاه  $Ax = b$  است. این نقطه جایی است که گرادیان صفر است، یعنی مشتق تابع در آن صفر شده و هیچ جهتی برای افزایش تابع وجود ندارد.

$$f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \vdots \\ \frac{\partial}{\partial x_n} f(x) \end{bmatrix} \quad (2) \quad \text{گرادیان فرم درجه دوم به شکل زیر است:}$$

با کمی محاسبات می‌توان با اعمال کردن معادلات (۱) و (۲) به رابطه زیر رسید:

$$f'(x) = \frac{1}{2}A^T x + \frac{1}{2}Ax - b \quad (۳)$$

و اگر ماتریس  $A$  متقارن باشد (یعنی  $A^T = A$ ) آنگاه میتوان این رابطه را به صورت زیر نوشت:

$$f'(x) = Ax - b \quad (۴)$$

و اگر این رابطه را مساوی با صفر قرار دهیم آنگاه جواب بهینه دستگاه را بدست می‌آوریم. بنابراین، جواب دستگاه یک نقطه بحرانی از تابع  $f(x)$  است. اگر ماتریس  $A$  هم متقارن و هم معین مثبت باشد، آنگاه این جواب یک نقطه مینیمم برای تابع  $f(x)$  خواهد بود.

### روش گرادیان با بیشترین شیب (The method of Steepest Descent)

در روش گرادیان کاهشی تندترین شیب، از یک نقطه دلخواه  $x(0)$  شروع می‌کنیم و به سمت پایین یک سطح سهمی‌وار (paraboloid) حرکت می‌کنیم. ما مجموعه‌ای از قدم‌ها به صورت  $x(1), x(2), \dots$  برمی‌داریم تا زمانی که اطمینان یابیم به اندازه کافی به جواب  $x$  نزدیک شده‌ایم.

در هر گام، جهتی را انتخاب می‌کنیم که تابع  $f$  با بیشترین سرعت کاهش می‌یابد، یعنی در جهت منفی گرادیان  $f'(x(i))$ . معادله (۴)، این جهت برابر است با:

$$-f'(x(i)) = b - Ax(i)$$

چند تعریف اولیه:

- خطا:  $e(i) = x(i) - x$  برداری است که نشان می‌دهد چقدر از جواب واقعی فاصله داریم.
- باقیمانده:  $r(i) = b - Ax(i)$  نشان می‌دهد که چقدر با مقدار صحیح بردار  $b$  فاصله داریم. (به راحتی می‌توان دید که  $r(i) = Ae(i)$  و باید به باقیمانده به عنوان خطا که با  $A$  به فضای  $b$  نگاشته شده فکر کنید.)
- $r(i) = -f'(x(i))$  پس باید باقیمانده را برابر با جهت تندترین کاهش تابع  $f$  در نظر بگیرید.

فرض کنید از نقطه  $x = [2, -2]^T$  اولین گام ما، در جهت تندترین کاهش، روی خط صافی قرار می‌گیرد که در شکل ۴ (a) نشان داده شده است. نقطه بعدی به صورت زیر به دست می‌آید:

$$x(1) = x(0) + ar(0)$$

سؤال این است که چه اندازه گام (step size) باید برداریم؟

**جستجوی خطی (line search)** روشی است که مقدار  $\alpha$  را طوری انتخاب می‌کند که مقدار  $f$  را در طول یک خط کمینه کند. شکل ۴ (b) این فرآیند را نشان می‌دهد: فقط اجازه داریم نقطه‌ای را روی تقاطع صفحه عمودی و سطح سهمی‌وار انتخاب کنیم. شکل ۴ (c) سهمی‌ای (parabola) را نشان می‌دهد که از تقاطع آن دو سطح به دست آمده است. اکنون این پرسش مطرح می‌شود که مقدار  $\alpha$  در پایه‌ی این سهمی چقدر است؟

می‌دانیم برای کمینه کردن  $f$  باید مشتق جهتی  $\frac{d}{d\alpha} f(x_{(1)})$  برابر صفر باشد، با استفاده از مشتق زنجیره‌ای داریم:

$$\frac{d}{d\alpha} f(x_{(1)}) = f'(x_{(1)})^T \frac{d}{d\alpha} x_{(1)} = f'(x_{(1)})^T r_{(0)}$$

با برابر قرار دادن این عبارت با صفر، نتیجه می‌گیریم که  $\alpha$  باید طوری انتخاب شود که  $r_{(0)}$  بر  $f'(x_{(0)})$  عمود باشد.

برای تعیین  $\alpha$ ، توجه کنید که  $f'(x_{(1)}) = -r_{(1)}$  و داریم:

$$r_{(1)}^T r_{(0)} = 0$$

$$(b - Ax_{(1)})^T r_{(0)} = 0$$

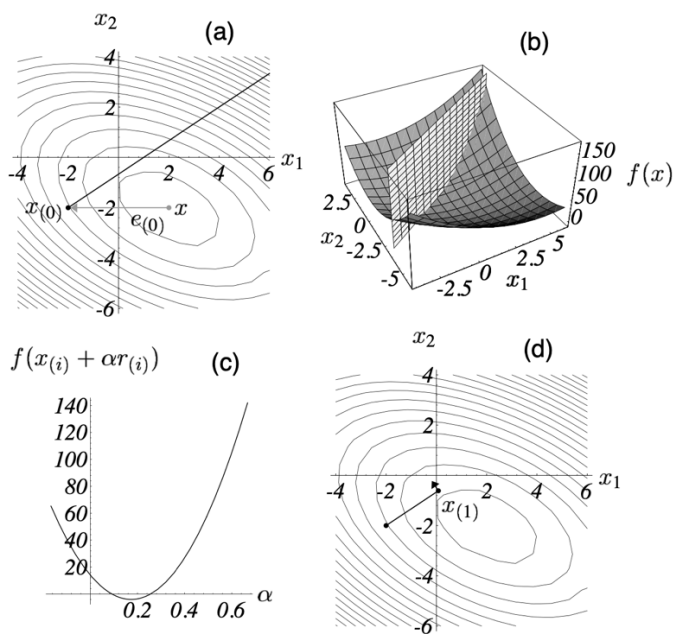
$$(b - A(x_{(0)} + \alpha r_{(0)}))^T r_{(0)} = 0$$

$$(b - Ax_{(0)})^T r_{(0)} - \alpha (Ar_{(0)})^T r_{(0)} = 0$$

$$(b - Ax_{(0)})^T r_{(0)} = \alpha (Ar_{(0)})^T r_{(0)}$$

$$r_{(0)}^T r_{(0)} = \alpha r_{(0)}^T (Ar_{(0)})$$

$$\alpha = \frac{r_{(0)}^T r_{(0)}}{r_{(0)}^T Ar_{(0)}}$$



شکل (۴) روش گرادیان نزولی (Steepest Descent):

- (a) با شروع از نقطه  $[-2, -2]$ ، یک گام در جهت بیشترین کاهش تابع  $f$  برداشته می‌شود.
- (b) نقطه‌ای روی تقاطع این دو سطح پیدا می‌شود که مقدار  $f$  را کمینه می‌کند.
- (c) این سهمی، حاصل تقاطع سطوح است؛ پایین‌ترین نقطه‌ی آن هدف ماست.
- (d) گرادیان در پایین‌ترین نقطه عمود بر گرادیان گام قبلی است.

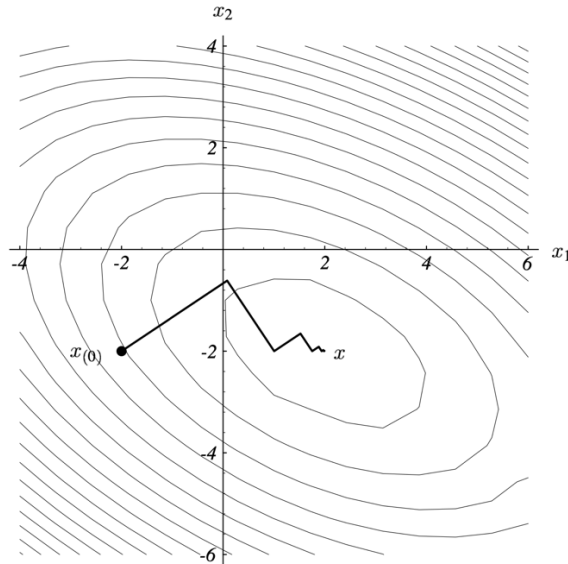
در مجموع، روش گرادیان نزولی (Steepest Descent) به صورت زیر است:

$$r(i) = b - Ax(i)$$

$$\alpha(i) = \frac{r(i)^T r(i)}{r(i)^T A r(i)}$$

$$x(i+1) = x(i) + \alpha(i) r(i)$$

مثال زیر را در نظر بگیرید.



شکل (۵) در این جا، الگوریتم از نقطه‌ی حدودی  $[-2, 2]$  شروع می‌شود و در نقطه‌ی حدودی  $[2, 2]$  همگرا می‌شود.

این مثال تا زمانی که همگرا شود، در شکل ۸ اجرا می‌شود. به مسیر زیگزاگی توجه کنید؛ این مسیر به این دلیل ظاهر می‌شود که هر گرادیان بر گرادیان قبلی عمود است.

پایداری روش گرادیان کاهش با بیشترین شیب

پایداری این روش به عوامل متعددی وابسته است از جمله:

### ۱. پایداری عددی و شرطی بود ماتریس هسین

در مسائل مربعی (مانند کمینه‌سازی فرم درجه دوم)، عملکرد روش به شدت به ماتریس  $A$  بستگی دارد. اگر ماتریس بد حالت باشد، یعنی نسبت بزرگ‌ترین به کوچک‌ترین مقدار ویژه آن زیاد باشد، مسیر حرکت الگوریتم به صورت زیگزاگ و ناپایدار درمی‌آید و نرخ همگرایی کند می‌شود. این رفتار در نمودارهای کانتور به صورت حرکت‌های عمودشونده پیاپی روی ایزوکانتورهای بیضوی مشاهده می‌شود.

### ۲. حساسیت به انتخاب گام

روش نزول تندترین نیاز به یک جستجوی خطی دقیق (line search) دارد تا اندازه بهینه گام در هر تکرار تعیین شود. اگر این گام به درستی انتخاب نشود (مثلاً خیلی کوچک یا خیلی بزرگ باشد)، روش ممکن است همگرا نشود یا نوسان کند. در نتیجه، پایداری عددی این روش نسبت به روش‌هایی مانند گرادیان مزدوج یا روش نیوتن ضعیف‌تر است.

### ۳. همگرایی آهسته برای مسائل بد شرط

برای سیستم‌هایی با ایزوکانتورهای کشیده (یعنی مقادیر ویژه‌ی بسیار متفاوت)، جهت‌های گرادیان در هر مرحله تقریباً عمود بر یکدیگر هستند. این موضوع باعث می‌شود الگوریتم با هر بار حرکت، بخش کمی از خطا را کاهش دهد و نیاز به تکرارهای زیاد داشته باشد. بنابراین، پایداری و سرعت همگرایی در این حالت کاهش می‌یابد.

## روش گرادیان مزدوج

روش گرادیان کاهشی با بیشترین شیب (Steepest Descent) اغلب در گام‌های متوالی خود، در همان جهتی حرکت می‌کند که در گام‌های قبلی نیز طی شده است (نگاه کنید به شکل ۸). اما آیا بهتر نبود اگر هر بار که گامی برمی‌داریم، آن گام از همان ابتدا در جهت بهینه باشد؟

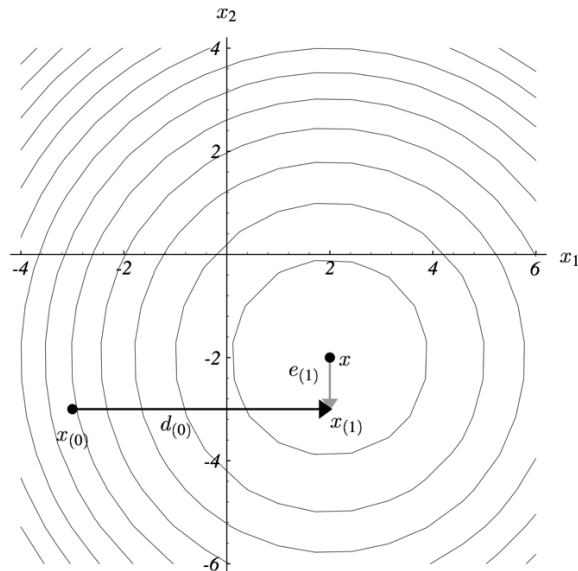
در این جا ایده‌ای مطرح می‌شود: فرض کنیم مجموعه‌ای از بردارهای جستجوی متعامد  $d_{(0)}, d_{(1)}, \dots, d_{(n-1)}$  را در اختیار داریم. در هر یک از این جهت‌های جستجو فقط یک گام برمی‌داریم، و این گام با دقت طوری انتخاب می‌شود که نتیجه‌ی آن دقیقاً در راستای نقطه‌ی هدف  $x$  قرار گیرد. پس از  $n$  گام، به پاسخ نهایی خواهیم رسید.

شکل ۶ این ایده را به تصویر می‌کشد؛ در آن محورهای مختصات به‌عنوان جهت‌های جستجو در نظر گرفته شده‌اند. گام اول (در راستای محور افقی) ما را به مقدار صحیح مولفه‌ی اول  $x_1$  می‌رساند؛ گام دوم (در راستای محور عمودی) نیز دقیقاً به مقصد می‌رسد. توجه کنید که خطای حاصل از گام دوم، یعنی  $e_{(1)}$ ، بر جهت گام اول یعنی  $d_{(0)}$  عمود است.

در حالت کلی، برای هر گام، نقطه‌ای را به گونه‌ای انتخاب می‌کنیم که:

$$x_{(i+1)} = x_{(i)} + \alpha_{(i)} d_{(i)}$$

برای تعیین مقدار بهینه‌ی  $\alpha_{(i)}$  از این نکته استفاده می‌کنیم که خطای حاصل از گام بعدی یعنی  $e_{(i+1)}$  باید بر جهت گام فعلی  $d_{(i)}$  عمود باشد؛ در نتیجه، هیچ‌گاه نیازی نخواهد بود در همان جهت دوباره حرکت کنیم.



شکل (۶) روش جهت‌های متعامد متأسفانه، این روش تنها زمانی کارآمد است که پاسخ مسئله را از قبل بدانید.



CG صرفاً همان روش جهت‌های مزدوج (Conjugate Directions) است، با این تفاوت که جهت‌های جستجو از باقی‌مانده‌ها (residuals) ساخته می‌شوند (یعنی با قرار دادن  $u(i) = r(i)$ ).

این انتخاب از چند جهت منطقی است. اول آنکه باقی‌مانده‌ها در روش گرادینز کاهشی تند (Steepest Descent) عملکرد خوبی داشتند، پس چرا در روش جهت‌های مزدوج نیز از آنها استفاده نکنیم؟ دوم آنکه باقی‌مانده دارای ویژگی مطلوبی است: با جهت‌های جستجوی قبلی متعامد است. بنابراین، استفاده از آن تضمین می‌کند که در هر مرحله، جهت جستجوی جدیدی تولید می‌شود که مستقل خطی از جهت‌های قبلی است — مگر آنکه باقی‌مانده صفر باشد، که در این صورت مسئله قبلاً حل شده است.

حالا به بررسی پیامدهای این انتخاب بپردازیم. از آنجایی که بردارهای جستجو از باقی‌مانده‌ها ساخته می‌شوند، زیرفضای تولید شده توسط (مجموعه  $\{r(0), r(1), \dots, r(i-1)\}$  span برابر است با  $D(i)$  یعنی زیرفضای جستجو در مرحله  $i$  ام). همچنین، از آنجایی که هر باقی‌مانده بر جهت‌های جستجوی قبلی متعامد است، بر باقی‌مانده‌های قبلی نیز متعامد خواهد بود. بنابراین معادله به صورت زیر ساده می‌شود:

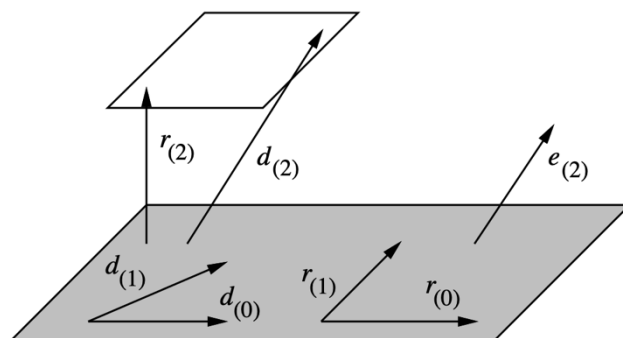
$$r(i)^T r(j) = 0 \quad i \neq j$$

نکته‌ی جالب این‌جاست که معادله‌ی زیر:

$$\begin{aligned} r_{(i+1)} &= -Ae_{(i+1)} \\ &= -A(e_{(i+1)} + \alpha_{(i)}d_{(i)}) \\ &= r_{(i)} - \alpha_{(i)}Ad_{(i)} \end{aligned}$$

نشان می‌دهد باقی‌مانده‌ی جدید  $r_{(i)}$  صرفاً ترکیبی خطی از باقی‌مانده‌ی قبلی و بردار  $Ad_{(i)}$  است. با توجه به اینکه  $d_{(i)}$  این حقیقت نشان می‌دهد که هر زیرفضای جدید  $D_{(i+1)}$  از اجتماع زیرفضای قبلی  $D_{(i)}$  و زیرفضای  $AD_{(i)}$  تشکیل می‌شود. بنابراین:

$$\begin{aligned} D_{(i)} &= \text{span} \{d_{(0)}, Ad_{(0)}, A^2d_{(0)} \dots, A^{i-1}d_{(0)}\} \\ &= \text{span} \{r_{(0)}, Ar_{(0)}, A^2r_{(0)} \dots, A^{i-1}r_{(0)}\} \end{aligned}$$



شکل (۷) در روش گرادینز مزدوج، هر باقی مانده‌ی جدید با تمام باقی مانده‌ها و جهت‌های جستجوی قبلی عمود است؛ و هر جهت جستجوی جدید نیز از باقی مانده ساخته می‌شود، به گونه‌ای که با تمام باقی مانده‌ها و جهت‌های جستجوی قبلی، نسبت به ماتریس  $A$  عمود باشد. نقاط پایانی بردارهای  $r(2)$  و  $d(2)$  روی صفحه‌ای قرار می‌گیرند که با زیرفضای  $D(2)$  (که در شکل به صورت ناحیه‌ی سایه‌خورده نشان داده شده) موازی است. در روش CG، بردار  $d(2)$  ترکیبی خطی از  $r(2)$  و  $d(2)$  است.

این زیرفضا، به نام زیرفضای کریلوف (Krylov Subspace) شناخته می‌شود. زیرفضایی که با اعمال مکرر یک ماتریس بر یک بردار تولید می‌شود. این زیرفضا ویژگی جذابی دارد: از آنجا که  $AD(i)$ ، و با توجه به اینکه باقی مانده‌ی بعدی  $r(i)$  بر  $D(i)$  عمود است، نتیجه می‌گیریم که  $r(i+1)$  نسبت به زیرفضای  $D(i)$  نیز  $A$ -متعامد خواهد بود. در نتیجه، فرآیند متعامدسازی گرم-اشمیت (Gram-Schmidt Conjugation) در اینجا ساده می‌شود، چرا که  $r(i+1)$  از پیش نسبت به تمام جهت‌های جستجوی قبلی (به جز  $d(i)$  نسبت به  $A$  متعامد است). ثابت‌های گرام اشمیت به صورت زیر است بیایید آن را ساده تر کنیم.

$$\beta_{ij} = -\frac{r_{(i)}^T A d_{(j)}}{d_{(j)}^T A d_{(j)}};$$

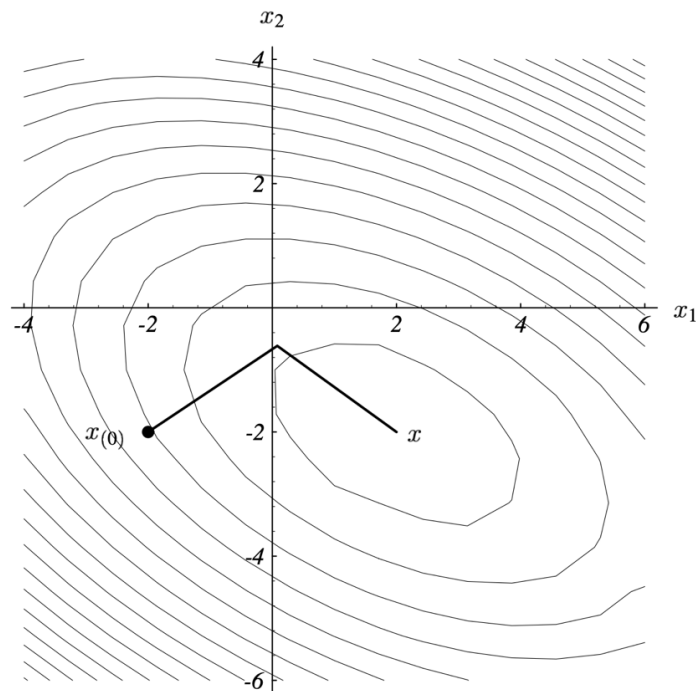
$$r_{(i)}^T r_{(j+1)} = r_{(i)}^T r_{(j)} - \alpha_{(j)} r_{(i)}^T A d_{(j)}$$

$$\Rightarrow \alpha_{(j)} r_{(i)}^T A d_{(j)} = r_{(i)}^T r_{(j)} - r_{(i)}^T r_{(j+1)}$$

$$r_{(i)}^T A d_{(i)} = \begin{cases} \frac{1}{\alpha_{(i)}} r_{(i)}^T r_{(i)}, & i = j \\ -\frac{1}{\alpha_{(i-1)}} r_{(i)}^T r_{(i)}, & i = j + 1 \\ 0 & otherwise \end{cases}$$

دیگر نیازی به ذخیره سازی بردارهای جستجوی قبلی برای تضمین  $A$ -متعامد بودن بردارهای جستجوی جدید وجود ندارد. این پیشرفت عمده همان چیزی است که الگوریتم  $CG$  (گرادیان مزدوج) را تا این حد مهم کرده است، زیرا پیچیدگی فضایی و زمانی هر تکرار از  $O(n^2)$  به  $O(m)$  کاهش یافته است، که در آن  $m$  تعداد مؤلفه های ناصفر ماتریس  $A$  است. از این پس، از نماد اختصاری  $\beta_{(i)} = \beta_{i,i-1}$  استفاده خواهیم کرد. بیایید بیشتر ساده سازی کنیم:

$$\begin{aligned} \beta_i &= \frac{r_{(i)}^T r_{(i)}}{d_{(i-1)}^T r_{(i-1)}} \\ &= \frac{r_{(i)}^T r_{(i)}}{r_{(i-1)}^T r_{(i-1)}} \end{aligned}$$



شکل (۸) روش گرادیان مزدوج

پس به طور کل داریم:

$$\begin{aligned}
 d_{(0)} &= r_{(0)} = b - Ax_{(0)}, \\
 \alpha_{(i)} &= \frac{r_{(i)}^T r_{(i)}}{d_{(i)}^T A d_{(i)}} \quad (\text{by Equations 32 and 42}), \\
 x_{(i+1)} &= x_{(i)} + \alpha_{(i)} d_{(i)}, \\
 r_{(i+1)} &= r_{(i)} - \alpha_{(i)} A d_{(i)}, \\
 \beta_{(i+1)} &= \frac{r_{(i+1)}^T r_{(i+1)}}{r_{(i)}^T r_{(i)}}, \\
 d_{(i+1)} &= r_{(i+1)} + \beta_{(i+1)} d_{(i)}.
 \end{aligned}$$

### آنالیز پایداری گرادیان مزدوج

الگوریتم گرادیان مزدوج (CG) به صورت نظری پس از  $n$  تکرار کامل می شود، اما در عمل، خطاهای عددی مانند گرد کردن و حذف موجب کاهش دقت پسماند و از دست رفتن خاصیت  $A$ -متعامد بودن بردارهای جستجو می شوند. این مشکلات، به ویژه از دست رفتن مزدوج بودن در دهه ی ۱۹۶۰ باعث شد جامعه ی ریاضی CG را کنار بگذارد، اما در دهه ی ۱۹۷۰ با اثبات کارایی آن در مسائل بزرگ، دوباره مورد توجه قرار گرفت. امروزه تحلیل همگرایی اهمیت زیادی دارد، زیرا CG برای مسائلی استفاده می شود که اجرای کامل  $n$  تکرار در آنها ممکن نیست. بنابراین، تحلیل همگرایی بیشتر به عنوان اثباتی بر مفید بودن CG در مسائل بسیار بزرگ تلقی می شود تا ابزاری برای مقابله با خطای عددی. همچنین، چون اولین تکرار CG مانند روش *steepest descent* است، شرایط همگرایی آن مشابه بررسی می شود.

### مقایسه پیچیدگی

نتیجه می گیریم که روش Steepest Descent دارای پیچیدگی زمانی  $O(mk)$  است، در حالی که روش گرادیان مزدوج (CG) دارای پیچیدگی زمانی  $O(m\sqrt{k})$  می باشد. هر دو الگوریتم دارای پیچیدگی فضایی  $O(m)$  هستند. تقریبات تفاضل محدود و المان محدود برای مسائل مقدار مرزی بیضوی مرتبه دوم که روی حوزه هایی با ابعاد  $d$  تعریف شده اند، معمولاً شرط وضعیتی  $\kappa$  از مرتبه  $O(n^{2/d})$  دارند. بنابراین، روش نزول شیب دار برای مسائل دوبعدی دارای پیچیدگی زمانی  $O(n^2)$  است، در حالی که CG دارای پیچیدگی زمانی  $O(n^{3/2})$  می باشد. همچنین، برای مسائل سه بعدی، پیچیدگی زمانی نزول شیب دار  $O(n^{5/3})$  و برای CG برابر با  $O(n^{4/3})$  است.

## روش گرادیان طبیعی

گرادیان طبیعی (Natural Gradient) نسخه‌ای اصلاح‌شده از گرادیان معمولی است که با در نظر گرفتن ساختار هندسی فضای پارامترها، جهت واقعی بیشترین کاهش تابع هزینه را تعیین می‌کند. برخلاف گرادیان کلاسیک که در فضاهای پیچیده ممکن است باعث نوسان یا همگرایی کند شود، گرادیان طبیعی با استفاده از هندسه‌ی اطلاعات، مسیر یادگیری را بهینه‌تر و پایدارتر می‌سازد. در بسیاری از کاربردهای یادگیری ماشین مانند پرسپترونها، جداسازی کور منابع، و سیستم‌های دینامیکی، ثابت شده است که این روش دارای کارایی فشر است؛ یعنی در بلندمدت به اندازه‌ی بهترین روش‌های تخمین دسته‌ای عملکرد دارد. همچنین، استفاده از گرادیان طبیعی می‌تواند پدیده‌ی توقف ناگهانی یادگیری (پلاتو) در الگوریتم‌هایی مانند backpropagation را کاهش داده یا حذف کند.

اجازه دهید  $S = w \in \mathbb{R}^n$  یک فضای پارامتر باشد که تابعی به نام  $L(w)$  روی آن تعریف شده است.

هنگامی که  $S$  یک فضای اقلیدسی با دستگاه مختصات متعامد و یک‌نواخت باشد، طول مربع یک بردار افزایشی کوچک  $dw$  که نقاط  $w$  و  $w + dw$  را به هم متصل می‌کند، به صورت زیر بیان می‌شود:

$$|dw|^2 = \sum_{i=1}^n (dw_i)^2,$$

که در آن  $dw_i$  مؤلفه‌های بردار  $dw$  هستند. با این حال زمانی که دستگاه مختصات غیرمتعامد باشد، طول مربع این بردار با استفاده از یک فرم درجه دوم (quadratic form) داده می‌شود.

$$|dw|^2 = \sum_{i,j} g_{ij}(w) dw_i dw_j.$$

وقتی  $S$  یک منیفلد خمیده باشد، دیگر دستگاه مختصات خطی متعامد وجود ندارد و طول بردار  $dw$  همیشه به صورت معادله بالا نوشته می‌شود. چنین فضایی یک فضای ریمانی (Riemannian space) نامیده می‌شود. در فضای پارامترهای شبکه‌های عصبی دارای ویژگی ریمانی است. ماتریس  $n \times n$  به صورت  $G = (g_{ij})$  که به طور کلی به  $w$  وابسته است، تانسور متریک ریمانی (Riemannian metric tensor) نامیده می‌شود. این ماتریس در شرایط خاص به حالت ساده‌تری کاهش می‌یابد.

$$g_{ij}(w) = \delta_{ij} = \begin{cases} 1, & i = j, \\ 0, & i \neq j \end{cases}$$

در حالت متعامد اقلیدسی، ماتریس  $G$  برابر با ماتریس یک  $I$  است.

جهت بیشترین کاهش (شیب تندترین کاهش) تابع  $L(w)$  در نقطه  $w$  توسط بردار  $dw$  تعریف می‌شود که مقدار  $L(w + dw)$  را کمینه می‌کند، به شرطی که طول  $|dw|$  ثابت باشد؛ یعنی تحت این قید که:

$$|dw|^2 = \varepsilon^2$$

قضیه ۱. جهت بیشترین کاهش تابع  $L(w)$  در یک فضای ریمانی توسط رابطه زیر داده می‌شود:

$$-\tilde{\nabla} L(w) = -G^{-1}(w) \nabla L(w)$$

که در آن  $G^{-1} = (g^{ij})$  (معکوس متریک  $G = (g_{ij})$  است و  $\nabla L$  گرادیان معمولی تابع  $L$  می باشد).

$$\nabla L(\mathbf{w}) = \left( \frac{\partial}{\partial w_1} L(\mathbf{w}), \dots, \frac{\partial}{\partial w_n} L(\mathbf{w}) \right)^T,$$

Proof. We put

$$d\mathbf{w} = \varepsilon \mathbf{a},$$

and search for the  $\mathbf{a}$  that minimizes

$$L(\mathbf{w} + d\mathbf{w}) = L(\mathbf{w}) + \varepsilon \nabla L(\mathbf{w})^T \mathbf{a}$$

under the constraint

$$|\mathbf{a}|^2 = \sum g_{ij} a_i a_j = 1.$$

By the Lagrangean method, we have

$$\frac{\partial}{\partial a_i} \{ \nabla L(\mathbf{w})^T \mathbf{a} - \lambda \mathbf{a}^T G \mathbf{a} \} = 0.$$

This gives

$$\nabla L(\mathbf{w}) = 2\lambda G \mathbf{a}$$

or

$$\mathbf{a} = \frac{1}{2\lambda} G^{-1} \nabla L(\mathbf{w}),$$

where  $\lambda$  is determined from the constraint.

We call

$$\tilde{\nabla} L(\mathbf{w}) = G^{-1} \nabla L(\mathbf{w})$$

گرادیان طبیعی تابع  $L$  در فضای ریمانی، جهت  $-\tilde{\nabla} L$  است. بنابراین، جهت بیشترین کاهش تابع  $L$  را نشان می دهد. (اگر از نمادگذاری تانسوری استفاده کنیم، این دقیقاً شکل متضاد گرادیان  $-\nabla L$  است.) زمانی که فضا اقلیدسی و دستگاه مختصات، متعامد و نرمال باشد، در این حالت ماتریس متریک  $G$  برابر با ماتریس همانی است، بنابراین گرادیان طبیعی با گرادیان معمولی برابر خواهد بود.

$$\tilde{\nabla} L = \nabla L$$

این بیانگر الگوریتم گرادیان نزولی طبیعی با فرم زیر است:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \tilde{\nabla} L(\mathbf{w}_t),$$

که در آن  $\eta_t$  نرخ یادگیری است که اندازه ی گام را تعیین می کند.

## پایداری گرادیان طبیعی

یکی از مهم‌ترین مزایای گرادیان طبیعی نسبت به گرادیان معمولی، پایداری عددی بالاتر آن در هنگام یادگیری است. در فضاهای پارامتری که شرطساز (conditioning) ضعیف دارند یا متغیرها با مقیاس‌های مختلف ظاهر می‌شوند، گرادیان معمولی ممکن است باعث نوسان یا کندی همگرایی شود. اما گرادیان طبیعی با در نظر گرفتن ساختار ریمانی فضای پارامتر، جهت بهینه را به گونه‌ای تنظیم می‌کند که حرکت در آن فضا پایدارتر و کاراتر باشد. همچنین، مطالعات نظری نشان داده‌اند که گرادیان طبیعی Efficient Fisher است؛ یعنی به صورت مجانبی (asymptotically) عملکردی معادل با برآورد بهینه‌ی دسته‌ای (Batch Estimation) دارد. به همین دلیل، پدیده‌هایی مانند "پلاتو" که در یادگیری شبکه‌های عصبی با گرادیان معمولی رخ می‌دهد، در گرادیان طبیعی یا دیده نمی‌شوند یا بسیار کاهش می‌یابند.

## نتیجه

می‌توان نتیجه گرفت که هر یک از این روش‌ها برای موقعیت‌ها و ساختارهای خاصی از مسائل مناسب هستند. روش گرادیان با تندترین شیب ساده و قابل فهم است، اما در مسائل بدشرط (ill-conditioned) و با توابع پیچیده، همگرایی کندی دارد و مسیر زیگ‌زاگی طی می‌کند. روش گرادیان مزدوج با بهینه‌سازی جهت حرکت و استفاده از اطلاعات قبلی، سرعت همگرایی بهتری در مسائل مربعی یا خطی دارد و اغلب در مسائل بزرگ‌مقیاس مؤثرتر است. در مقابل، گرادیان طبیعی که بر پایه‌ی هندسه‌ی اطلاعاتی و متریک Fisher عمل می‌کند، در فضاهای پیچیده مانند آموزش شبکه‌های عصبی عملکرد بهتری دارد.

## مقایسه

ویژگی	گرادیان با تندترین شیب	گرادیان مزدوج	گرادیان طبیعی
پیچیدگی محاسباتی	پایین	متوسط	بالا (به دلیل محاسبه‌ی متریک Fisher)
سرعت همگرایی	کند در مسائل بدحالت	سریع‌تر در مسائل مربعی	پایدار و مناسب در فضاهای پیچیده
پایداری عددی	پایین	متوسط	بالا
کاربردها	مسائل ساده و آموزش اولیه	مسائل بزرگ‌مقیاس و خطی	شبکه‌های عصبی، سیستم‌های دینامیکی، جداسازی کور منبع

## تحلیل و آنالیز پیاده سازی

در ابتدا از نظر خروجی نهایی، هر سه الگوریتم توانستند جواب دقیقی برای مسئله کمینه سازی تابع درجه دوم  $f(x) = \frac{1}{2}x^T Ax - b^T x + c$  پیدا کنند. الگوریتم Steepest Descent معمولاً برای چنین توابعی به خوبی عمل می کند، اما در عمل به دلیل استفاده از گرادیان مستقیم ممکن است نیاز به تعداد بیشتری تکرار داشته باشد، خصوصاً اگر ماتریس  $A$  دارای عدد حالت بالایی باشد. الگوریتم Conjugate Gradient معمولاً بسیار سریع تر همگرا می شود، چون از ساختار خاص ماتریس متقارن مثبت معین استفاده می کند و جهت های جستجوی خود را بهینه تر تنظیم می کند. در مورد Natural Gradient Descent، این الگوریتم از اطلاعات هندسی (یعنی ساختار ماتریس  $A$ ) برای اصلاح گرادیان استفاده می کند و در واقع با "پیش شرطی سازی" بهبود یافته، به سمت جواب حرکت می کند. در این مورد خاص که  $A$  شناخته شده و کوچک است، Natural Gradient نیز جواب دقیقی را با نرخ همگرایی خوب ارائه می دهد.

از نظر پایداری عددی، Conjugate Gradient برتری دارد چون نیاز به محاسبه معکوس ماتریس ندارد و تنها ضرب ماتریسی و برداری انجام می دهد. در مقابل، Natural Gradient نیاز به محاسبه معکوس  $A$  دارد، که در مسائل بزرگ تر یا با شرط عددی ضعیف می تواند منجر به ناپایداری عددی شود. Steepest Descent در شرایطی که نرخ یادگیری بد انتخاب شود یا هندسه تابع بد باشد (مثل وقتی که بیضی های تراز بسیار کشیده اند)، دچار نوسان یا همگرایی کند می شود و از این جهت نسبت به بقیه الگوریتم ها حساس تر است.

از نظر پیچیدگی زمانی، Steepest Descent و Natural Gradient هر دو دارای پیچیدگی  $O(n^2)$  در هر گام هستند (با فرض ضرب ماتریسی کامل)، اما Natural Gradient نیاز به معکوس گیری اولیه با پیچیدگی  $O(n^3)$  دارد. این هزینه بالا در مسائل بزرگ بسیار محسوس است. در مقابل، Conjugate Gradient نه تنها از معکوس گیری اجتناب می کند بلکه اگر الگوریتم در کمتر از  $n$  گام همگرا شود، بسیار کارا تر خواهد بود؛ در واقع پیچیدگی آن تقریباً  $O(nk)$  و برای مسائل بزرگ مقیاس بهترین گزینه است، به ویژه وقتی  $A$  اسپارس باشد.



1. Natural Gradient Works Efficiently in Learning Shun-ichi Amari  
RIKEN Frontier Research Program, Saitama 351-01, Japan
2. An Introduction to the Conjugate Gradient Method Without the Agonizing Pain  
Jonathan Richard Shewchuk
3. S096 Matrix Calculus for Machine Learning and Beyond  
Independent Activities Period (IAP) 2023
4. Numerical Optimization, Jorge Nocedal