



دانشگاه صنعتی امیرکبیر  
(پلی‌تکنیک تهران)  
دانشکده علوم ریاضی

گزارش پروژه درس داده کاوی محاسباتی  
پروژه 3  
رشته علوم کامپیوتر گرایش هوش مصنوعی و محاسبات نرم  
هندسه یادگیری (انحنا، تعامد و هزینه محاسباتی در شبکه های عصبی)

نگارش  
هومن ذوالفقاری

استاد راهنما  
دکتر مهدی قطعی

تدریس یار  
مهندس بهنام یوسفی مهر

آذر 1404

اینجانب هومن ذوالفقاری متعهد می‌شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است.

نقل مطالب با ذکر مآخذ بلامانع است.

در صفحه تعهدنامه اصالت اثر، در قسمت بالا سمت چپ، تاریخ دفاع خود را جایگزین تاریخ نوشته شده کنید.

همچنین در صفحه تعهدنامه اصالت اثر، در خط اول، نام و نام خانوادگی خود را به صورت کامل با نام و نام خانوادگی نمونه، جایگزین کنید. در انتهای متن تعهد، در قسمت امضا نیز باید نام و نام خانوادگی کامل خود را وارد نمایید.

هومن ذوالفقاری

امضا

## چکیده

این پروژه به بررسی رفتار الگوریتم‌های بهینه‌سازی از دیدگاه هندسی می‌پردازد. ابتدا نشان داده می‌شود که در توابع بدحالت، گرادیان نزولی به‌دلیل عدد وضعیت بالا دچار زیگ‌زاگ و همگرایی کند می‌شود، در حالی‌که روش‌های مرتبه‌دوم مانند نیوتون و گرادیان مزدوج عملکرد بسیار بهتری دارند. سپس با آزمایش روی یک شبکه کوچک مشاهده می‌شود که روش‌های شبه‌نیوتونی در ابعاد پایین کارآمدتر از روش‌های مرتبه‌اول‌اند، اما در شبکه‌های عمیق هزینه محاسباتی هسین استفاده از آنها را غیرعملی می‌کند و روش‌هایی مانند SGD و Adam تنها گزینه‌های قابل استفاده هستند. در پایان، با اعمال تجزیه QR و ایجاد داده‌های متعامد نشان داده شد که بهبود هندسه ورودی موجب کاهش عدد وضعیت و تسریع همگرایی الگوریتم‌های مرتبه‌اول می‌شود. این نتایج نقش اساسی هندسه داده و مدل را در کارایی فرایند یادگیری برجسته می‌کند.

## واژه‌های کلیدی:

هندسه یادگیری، گرادیان کاهشی، روش نیوتون، گرادیان مزدوج، شبکه‌های عصبی عمیق، هسین، تجزیه QR، تعامد، پیش‌شرطی‌سازی، همگرایی.

چکیده .....	ا
فصل اول مقدمه .....	1
فصل دوم تحلیل های ریاضی (سطوح بدحالت) .....	3
تحلیل های ریاضی (سطوح بدحالت) .....	4
2-1- تعریف تابع آزمایشی .....	4
2-2- روش های پیاده سازی شده .....	5
2-3- مشاهدات تجربی .....	5
2-4- وقتی عدد وضعیت (Condition Number) زیاد می شود: .....	6
فصل سوم شبکه عصبی کلاسیک (فضای نیوتونی) شبکه عصبی کلاسیک (فضای نیوتونی) ....	7
2-5- مدل و داده .....	8
2-6- بهینه سازی مقایسه شده .....	8
2-7- نتایج .....	8
فصل چهارم شبکه عمیق و تله ی مقیاس پذیری شبکه عمیق و تله ی مقیاس پذیری .....	11
2-8- معماری: .....	12
2-9- محاسبه ابعاد مدل و هسین: .....	12
2-10- مقایسه SGD و Adam .....	14
فصل پنجم تعامد و QR (رویکرد داده کاوی) تعامد و QR (رویکرد داده کاوی) .....	15
2-11- ساخت دیتاست همبسته .....	16
فصل ششم جمع بندی و نتیجه گیری .....	18
منابع و مراجع .....	20

## فصل اول

### مقدمه

## مقدمه

مباحث هندسه یادگیری در سال‌های اخیر تبدیل به یکی از پایه‌های درک رفتار الگوریتم‌های بهینه‌سازی در شبکه‌های عصبی شده‌اند. با وجود اینکه روش‌های مرتبه‌دوم مانند نیوتون از نظر تئوری دارای همگرایی بسیار سریع‌تری هستند، در عمل مشاهده می‌شود که مدل‌های عمیق همچنان عمدتاً با روش‌های مرتبه‌اول (گرادیان ساده یا انواع مومنتوم مثل Adam) آموزش داده می‌شوند.

هدف این پروژه بررسی این پرسش است که چرا روش‌های مرتبه‌اول در یادگیری عمیق بهتر هستند، و چگونه هندسه داده‌ها (نه فقط الگوریتم) می‌تواند سرعت همگرایی را تغییر دهد.

برای این کار چهار محور اصلی بررسی شده است:

1. رفتار روش‌ها روی توابع بدحالت (Ill-conditioned)
2. کارایی روش‌های مرتبه‌اول و شبه‌نیوتونی روی یک شبکه کوچک
3. محل شکست روش نیوتون روی شبکه‌های عمیق
4. نقش تعامد (Orthogonality) و QR در بهبود هندسه مسئله

## فصل دوم

### تحلیل های ریاضی (سطوح بدحالت)

## تحلیل های ریاضی (سطوح بدحالت)

### 2-1- تعریف تابع آزمایشی

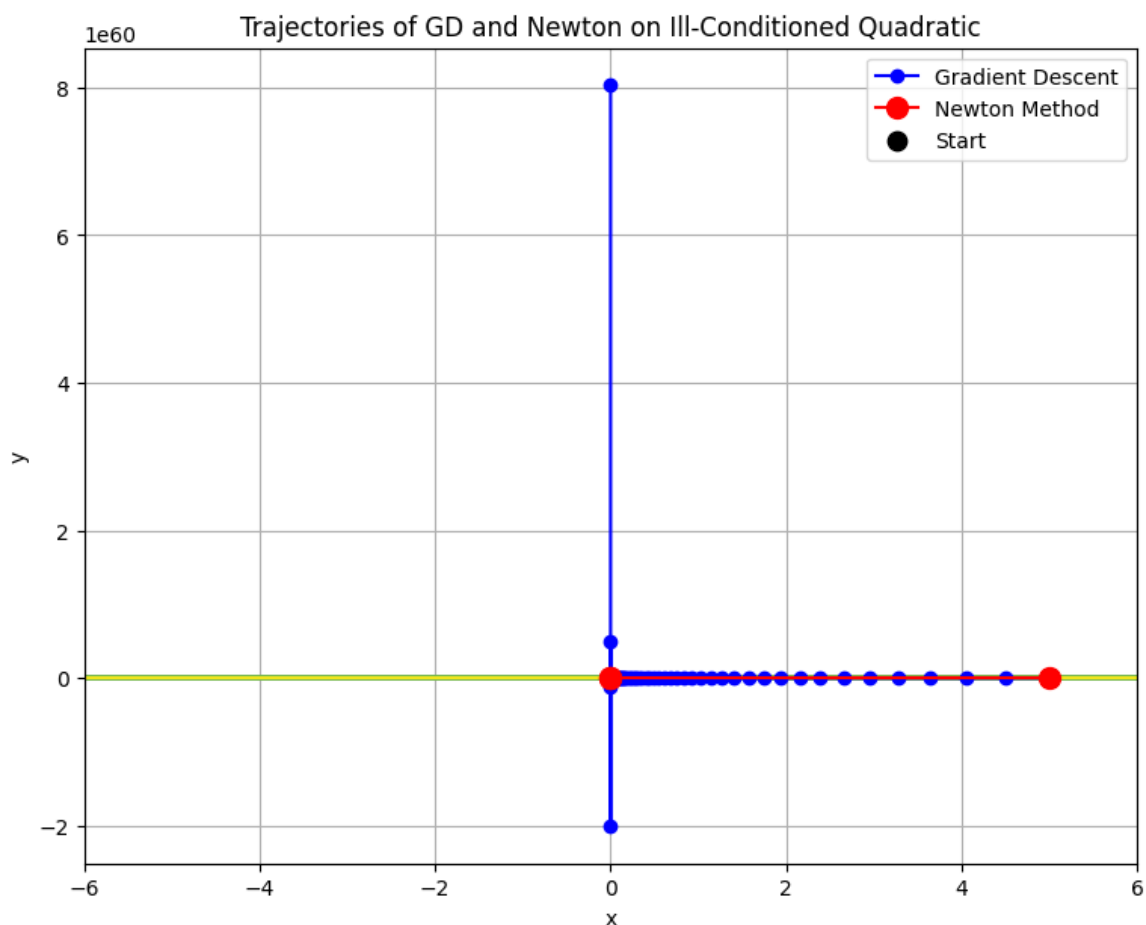
برای مشاهده رفتار بهینه‌سازها یک تابع درجه‌دو مصنوعی با ماتریس هسین زیر ساخته شد:

$$H = \begin{pmatrix} 1 & 0 \\ 0 & 50 \end{pmatrix}$$

این ماتریس دارای عدد وضعیت بزرگ است:

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} = 50$$

بنابراین کانتورهای تابع بیضی‌های بسیار کشیده هستند و گرادیان ساده در چنین دره‌هایی دچار زیگ زاگ می‌شود.





## 2-2- روش های پیاده سازی شده

سه روش از پایه پیاده سازی شدند:

### 1. گرادیان کاهشی (GD)

$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta_t)$$

در این مسئله، گرادیان برابر  $H\theta$  است GD فقط جهت شیب را می بیند و نسبت به عدد وضعیت بسیار حساس است.

### 2. روش نیوتون

$$\theta_{t+1} = \theta_t - H^{-1} \nabla J(\theta_t)$$

با وارد کردن  $H^{-1}$ ، گام ها از «بیضی» به «دایره» نگاشت می شوند؛ بنابراین از زیگزاگ جلوگیری می شود.

### 3. گرادیان مزدوج (CG)

روش میانی که گام های جستجو را نسبت به ماتریس  $H$  متعامد مزدوج انتخاب می کند؛ بدون نیاز به ذخیره هسین.

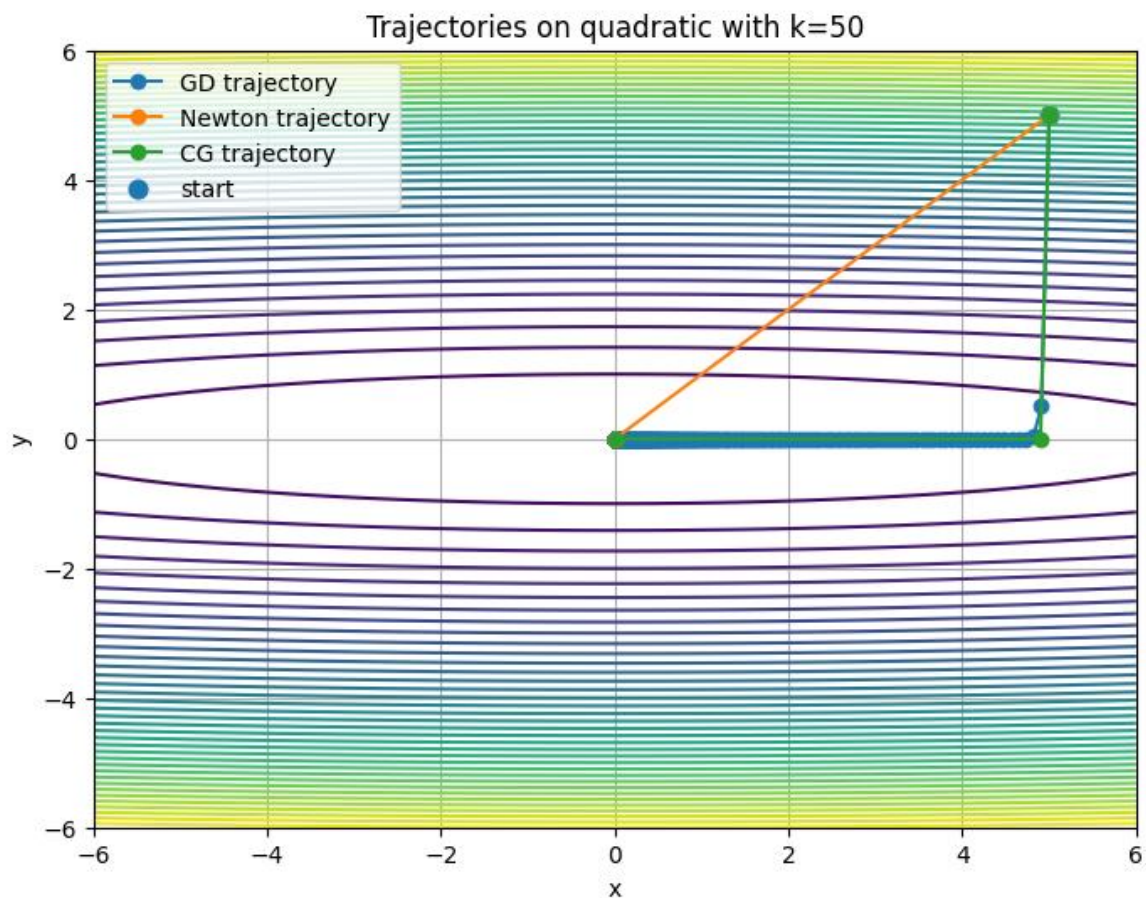
## 2-3- مشاهدات تجربی

- GD مسیر به وضوح زیگزاگی است؛ هر بار در جهت عمود بر مسیر واقعی منحرف می شود. دلیل: محور  $y$  انحنای بسیار بزرگ دارد و گرادیان در جهت عمودی بسیار قوی تر است.
- Newton در یک گام مستقیماً به مرکز می رود؛ همگرایی درجه دو. دلیل: نیوتون اطلاعات انحنا ( $H^{-1}$ ) را دارد و مستقیماً «جهت بهینه» را می داند.
- CG عملکرد بینابینی دارد؛ سریع تر از GD ولی کندتر از نیوتون. در حالت درجه دوم باید در 2 گام (تعداد ابعاد) همگرا شود. این روش مثل نیوتون رفتار می کند اما  $H$  را ذخیره نمی کند.

## 4-2- وقتی عدد وضعیت (Condition Number) زیاد می‌شود:

GD به شدت کند می‌شود (تا صدها یا هزاران تکرار ممکن است).

نیوتون عملاً مستقل از  $k$  در چند تکرار همگراست.  
CG سرعتش کمتر از نیوتون اما تقریباً پایدار است.



## فصل سوم شبکه عصبی کلاسیک (فضای نیوتونی)

## شبکه عصبی کلاسیک (فضای نیوتونی)

### 5-2- مدل و داده

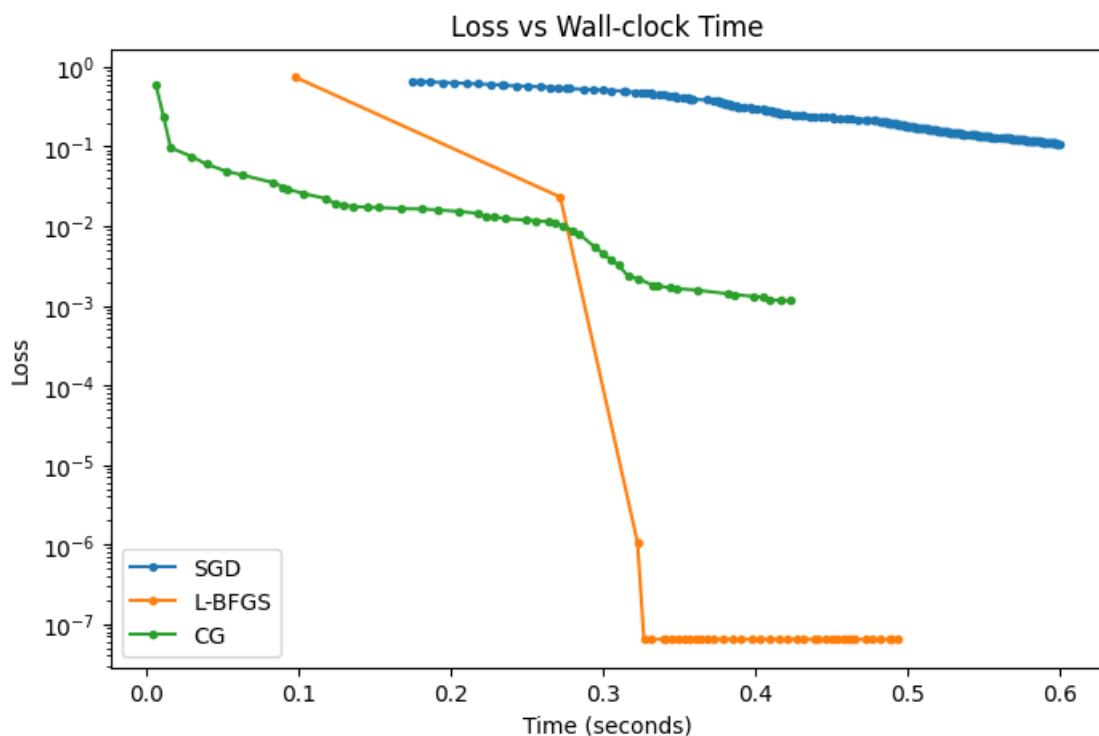
- دیتاست (Breast Cancer Wisconsin) دودویی
- شبکه کوچک: (Shallow MLP) یک Hidden layer با ۵ نورون
- تعداد پارامترها: ۱۱۵ → مناسب برای روش‌های شبه‌نیوتونی

### 6-2- بهینه‌سازهای مقایسه شده

1. SGD
2. L-BFGS تقریب نیوتون
3. Conjugate Gradient با استفاده از Scipy

### 7-2- نتایج

در مدل‌های کوچک، استفاده از اطلاعات انحنای کاملاً سودمند است و حتی L-BFGS از نظر زمان واقعی از SGD جلو می‌زند. این حالت مشابه تنوری است و نشان می‌دهد مشکل اصلی روش‌های مرتبه دوم هزینه محاسباتی آنهاست، نه کیفیت گام‌ها.



### چرا در مدل کوچک، L-BFGS بهتر عمل می‌کند؟

در این ابعاد (۱۵۰ پارامتر)، حافظه لازم برای ذخیره ماتریس تقریب هسین L-BFGS بسیار کم است. هر iteration آن هزینه زیادی دارد اما تعداد iterations آن کم است. به همین دلیل اگر loss سطحی و بدون نویز باشد (مثل دیتاست سرطان پستان)، L-BFGS معمولاً سریع‌تر از SGD به صفر می‌رسد.

### رفتار CG

CG نسبت به L-BFGS کمی کندتر است، اما از SGD بسیار پایدارتر و بدون نوسان همگرا می‌شود. این روش نیازی به ذخیره هسین ندارد (فقط گرادیان). در روش CG مبتنی بر SciPy، به دلیل اینکه شبکه عصبی غیرخطی با تابع زیان غیرکوادرانتیک است و ماتریس هسین در طول آموزش تغییر می‌کند، الگوریتم CG کلاسیک همگرایی تضمین‌شده ندارد و معمولاً با گام‌های بسیار بزرگ و افزایش شدید گرادیان باعث overflow می‌شود. اما نسخه Trust-Region CG (trust-ncg) یا استفاده از مدل پایدارسازی‌شده با Softplus و BCEWithLogitsLoss باعث رفع این ناپایداری می‌شود.

### چرا SGD کندتر است؟

هر آپدیت ارزان است اما تعداد آپدیت‌ها زیاد. در دیتاست واقعی که کمی بدحالت (ill-conditioned) است، مسیر SGD زیگزاگی و کند می‌شود. با این حال در مدل‌های deep به دلیل ارزان بودن هر iteration از L-BFGS بهتر است.

سریع‌ترین؟	بازه زمانی کل (ثانیه)	الگوریتم
کندترین	0.21~ثانیه	SGD
سریع	0.064~ثانیه	L-BFGS
متوسط	0.162~ثانیه	CG



## فصل چهارم

### شبکه عمیق و تله ی مقیاس پذیری

## شبکه عمیق و تله ی مقیاس پذیری

### 8-2- معماری:

شبکه عمیق با:

- ۳ لایه مخفی
- هر کدام ۱۰۰ نورون
- دیتاست Fashion-MNIST

### 9-2- محاسبه ابعاد مدل و هسین:

هر لایه شامل وزن ها و بایاس هاست. محاسبه پارامترها گام به گام:

- لایه ۱ : وزن ها  $784 \times 100 = 78\,400$  بایاس ها ۱۰۰ .

$$\text{مجموع لایه ۱} = 78\,400 + 100 = 78\,500$$

- لایه ۲ وزن ها  $100 \times 100 = 10\,000$  . بایاس ها ۱۰۰ .

$$\text{مجموع لایه ۲} = 10\,000 + 100 = 10\,100$$

- لایه ۳ : مانند لایه ۲ است

- لایه خروجی: وزن ها  $100 \times 10 = 1\,000$  بایاس ها ۱۰ .

$$\text{مجموع خروجی} = 1\,000 + 10 = 1\,010$$

- حالا جمع کل پارامترها  $N$

$$N = 78\,500 + 10\,100 + 10\,100 + 1\,010 = 99\,710$$

حافظه لازم برای ذخیره ماتریس هسین  $N \times N$ :

$$N^2 = 99,710^2$$

تعداد مؤلفه ها در هسین:



محاسبه با روش جبری (کوتاه):  $99,710 = 100,000 - 290$  :

$$\begin{aligned} 99,710^2 &= 100,000^2 - 2 \cdot 100,000 \cdot 290 + 290^2 \text{ پس} \\ &= 10,000,000,000 - 58,000,000 + 84,100 = 9,942,084,100 \end{aligned}$$

هر عدد float چهار بایت است. پس کل بایت ها:

$$\text{bytes} = 9,942,084,100 \times 4 = 39,768,336,400 \text{ bytes}$$

تبدیل به گیگابایت:

بر حسب «گیبی بایت»:  $(\text{GiB} = 2^{30} = 1,073,741,824)$  »

$$\frac{39,768,336,400}{1,073,741,824} \approx 37.04 \text{ GiB}$$

بر حسب «گیگابایت ده دهی»:  $(\text{GB} = 10^9)$  »

$$\frac{39,768,336,400}{10^9} \approx 39.77 \text{ GB}$$

بنابراین ذخیره هسین کامل نیازمند حدود ۳۷ GiB ( $\approx 40 \text{ GB}$ ) حافظه RAM است.

اگر فقط بخش مثلثی متقارن هسین را ذخیره کنیم (چون هسین متقارن است) حداقل نصف این مقدار لازم است  $\approx 18.5 \text{ GiB}$  که باز هم بزرگ و غیر عملی است برای اکثر کارت ها/سرورها در عمل (و این فقط حافظه نگهداری، بدون احتساب نیاز محاسباتی برای معکوس کردن است)

**هزینه محاسباتی معکوس/حل هسین:** عملگر معکوس/حل هسین (مثلاً چولسکی یا معکوس مستقیم)

در پیچیدگی تقریباً  $O(N^3)$  عملیات عددی است. برای  $N \approx 10^5$  داریم که  $N^3 \approx 9.913 \times 10^{14}$  عمل

که عددی عظیم (صدها تریلیون فلاپ)، به علاوه هزینه حافظه بالا و تبادل با حافظه. نتیجه عملی:

**نیوتون خالص** برای این مدل عمیق ( $N \approx 100k$ ) غیر قابل اجرا است. همین طور، هسین در نقطه های

مختلف آموزش تغییر می کند — پس حتی اگر یکبار هسین را محاسبه کنی، باید باز هم آن را

به روز رسانی کنی تا روش مؤثر بماند.

**نتیجه:** روش نیوتون در شبکه های عمیق غیر قابل اجراست؛ نه به دلیل ضعف تئوریک، بلکه به

دلیل هزینه حافظه و زمان.

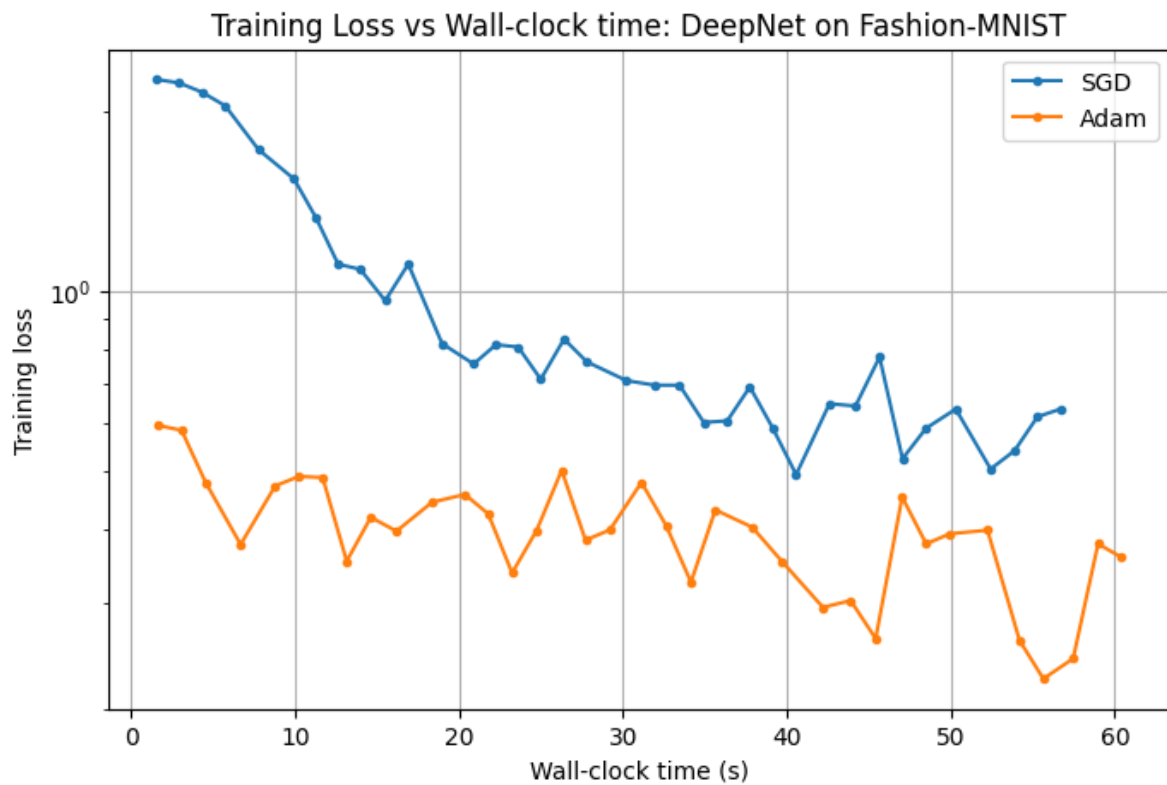
## 10-2- مقایسه Adam و SGD

### مشاهدات:

- Adam در ابتدا بسیار سریع‌تر پایین می‌آید.
- SGD رفتار لرزانی دارد.
- پس از چند Epoch ممکن است SGD برسد، اما Adam در زمان کمتر عملکرد بهتر دارد.

### نتیجه:

استفاده از اطلاعات آماری لحظه‌ای (moment estimates) در Adam باعث می‌شود الگوریتم بدون هسین، «اثر شبه‌نیوتونی» داشته باشد.



## فصل پنجم

### تعامل و QR (رویکرد داده کاوی)

## تعامل و QR (رویکرد داده کاوی)

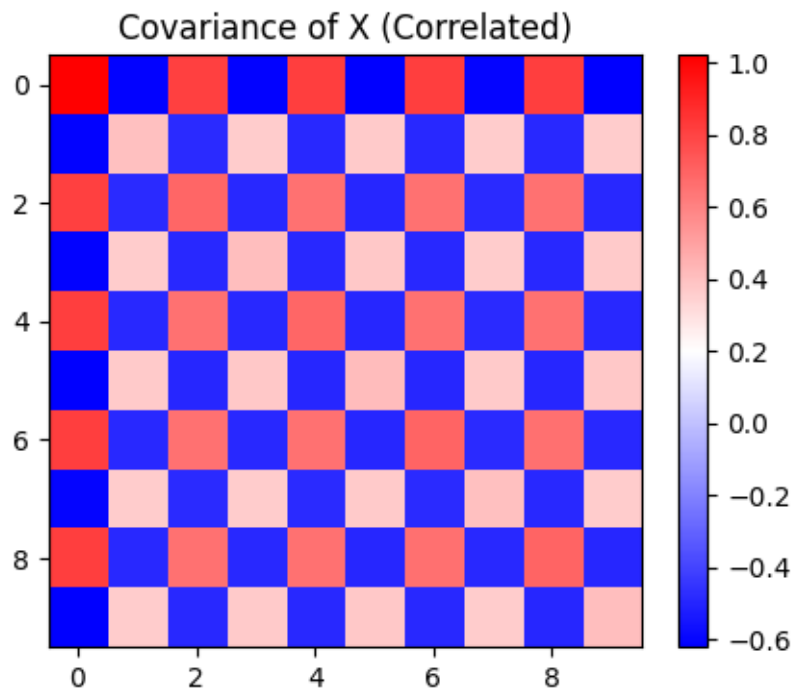
### 11-2- ساخت دیتاست همبسته

ویژگی‌ها به صورت خطی به یکدیگر وابسته طراحی شدند. (Multicollinearity)  
برای دیتاست اولیه:

$$\kappa(\text{Cov}(X)) \gg 1$$

این وضعیت باعث می‌شود سطح خطا در رگرسیون دره‌ای و کشیده شود و GD عملکرد ضعیف داشته باشد. در دیتاست داخل کد داریم:

Condition number of  $\text{Cov}(X) = 678.0130182609772$

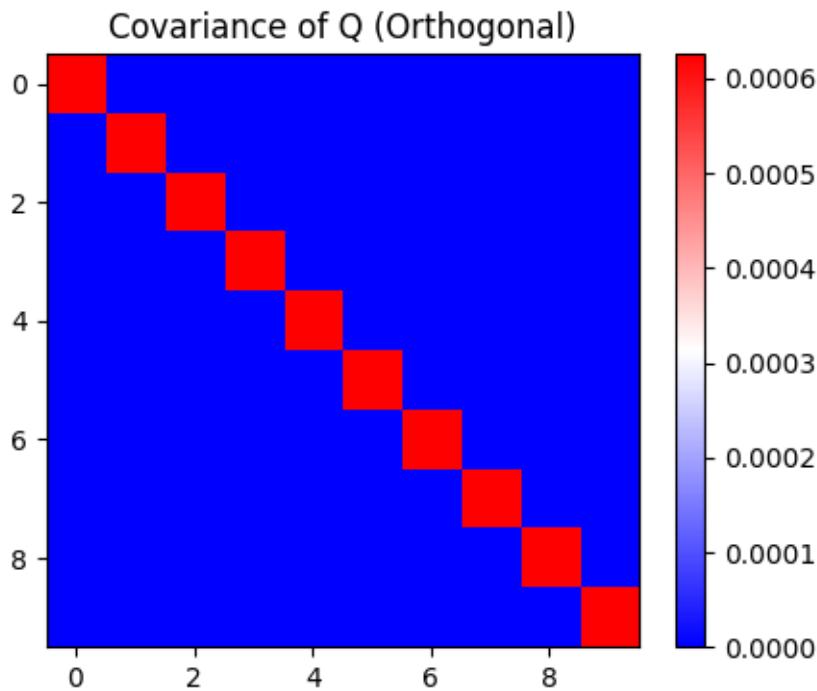


بعد از اجرای QR :

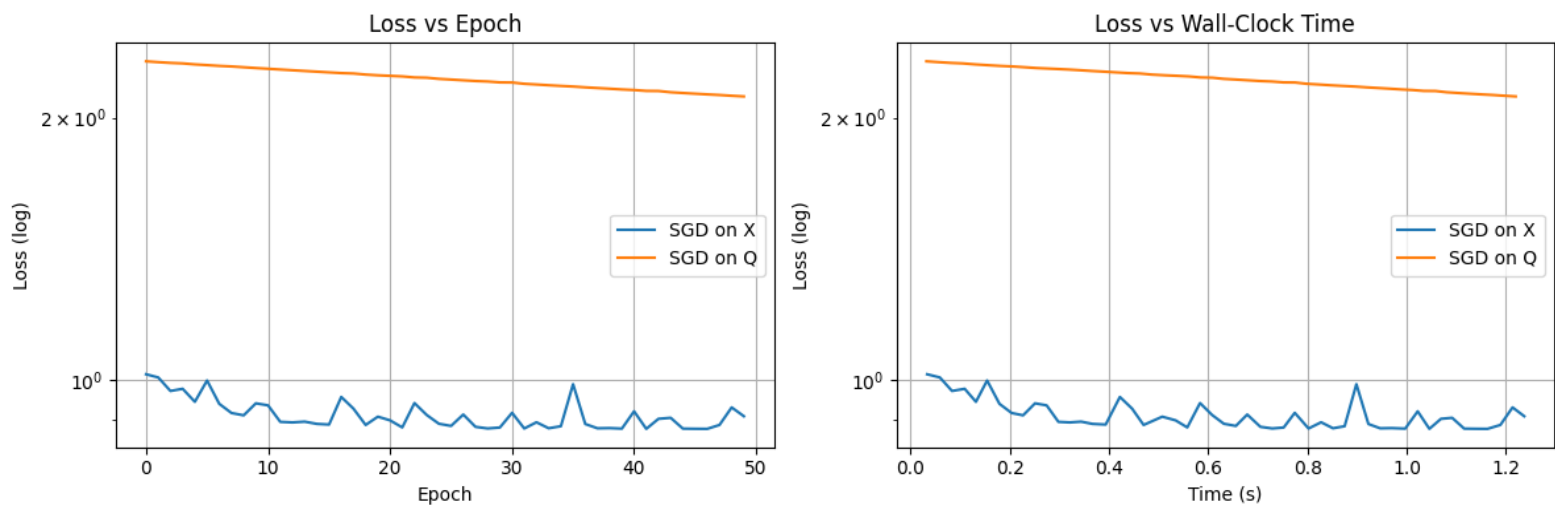
- ماتریس Q ویژگی‌های متعامد دارد
- کوواریانس Q تقریباً قطری می‌شود
- عدد **Condition Number** از مقدار بسیار بزرگ به نزدیک 1 کاهش می‌یابد
- در نتیجه سطح هزینه (Loss landscape) از کشیده به گرد (well-conditioned) می‌رود
- گرادیان کاهشی دیگر زیگزاگ نمی‌کند

- همگرایی بسیار سریع تر و نرم تر می‌شود
- رفتار SGD روی Q شبیه یک روش مرتبه دوم کارآمد می‌شود
- سطح خطا تقریباً ایزوتروپیک شده است:

Condition number of  $\text{Cov}(Q) = 1.0018074320188883$



اگرچه Loss اولیه روی داده Q به دلیل تغییر مقیاس ویژگی‌ها و کوچک بودن نرَم ستون‌ها کمی بزرگ‌تر است، اما مسیر کاهش Loss بسیار یکنواخت تر و سریع تر است. QR هندسه سطح خطا را اصلاح می‌کند، نه مقدار Loss اولیه را. بنابراین معیار صحیح مقایسه، سرعت همگرایی و پایداری آن است نه مقدار اولیه Loss



## فصل ششم جمع‌بندی و نتیجه‌گیری

## جمع‌بندی و نتیجه‌گیری

این پروژه به صورت تجربی چند واقعیت اساسی را در یادگیری عمیق نشان داد:

- روش نیوتون از نظر تئوری عالی است
- ولی هزینه محاسباتی  $O(N^2)$  آن باعث می‌شود در شبکه‌های بزرگ عملاً غیرقابل استفاده باشد.
- GD به شدت به عدد وضعیت وابسته است
- در توابع دره‌ای، GD کند و زیگ‌زاگی است؛ اما Newton و CG این مشکل را ندارند.
- در شبکه‌های کوچک، روش‌های شبه‌نیوتونی واقعاً بهتر عمل می‌کنند
- L-BFGS در مدل Breast Cancer بسیار سریع‌تر از SGD بود.
- در شبکه‌های بزرگ، Adam و SGD باقی می‌مانند
- Adam با استفاده از مومنت‌ها اثر «پیش‌شرطی‌سازی» ایجاد می‌کند.
- هندسه داده‌ها همان‌قدر مهم است که انتخاب بهینه‌ساز
- تعامد (QR) عدد وضعیت را بهبود می‌دهد و بهینه‌سازی را به شدت تسریع می‌کند.

## منابع و مراجع

Nocedal & Wright — Numerical Optimization (ویرایش دوم، 2006) [1]

(2004) Boyd & Vandenberghe — Convex Optimization [2]

(2016) Goodfellow, Bengio, Courville — Deep Learning [3]

(2022) Murphy — Probabilistic Machine Learning [4]