

Computational Optimization

with Applications to Machine Learning

A Series of Lectures at Missouri S&T

Wenqing Hu ¹

¹Department of Mathematics and Statistics, Missouri University of Science and Technology (formerly University of Missouri, Rolla), Rolla, Missouri, 65409, USA. Web: <https://huwenqing0606.github.io/>
Email: huwen@mst.edu

Preface

These lecture notes grew out of a series of lectures for the course MATH 5001: Computational Optimization taught at Missouri University of Science and Technology (formerly University of Missouri, Rolla) during Spring 2026.

These notes present a mixture of several topics, including optimization, machine learning, and computation, with a primary focus on optimization.

We thank the participating students: Stanislav Butkovich, Erich Gozebina, Qitong Huang, Nikunj Pradhan, and Mohammed Sleiman.

We also thank our Interim Chair Professor John Singler of the Mathematics & Statistics Department for his support of the course.

Wenqing Hu
Rolla, Missouri
Spring 2026

Contents

1	Motivating Examples	4
1.1	Supervised Learning	4
1.2	Matrix Optimizations	8
1.3	Experiment: Activation Functions	9
2	Convex Functions	10
2.1	Optimization problem: set-up and its solutions	10
2.2	Convexity	10
2.3	Taylor's theorem and convexity in Taylor's expansion	12
2.4	Optimality Conditions	18
2.5	Experiment: Nonconvexity of the Loss Landscape of Neural Networks	19
3	Gradient Descent and Line Search Methods	22
3.1	Gradient Descent	22
3.2	Back Propagation and GD on Neural Networks	27
3.3	Online Principle Component Analysis (PCA)	30
3.4	Line Search Methods	31
3.5	Convergence to Approximate Second Order Necessary Points	32
3.6	Experiment: Learning a Neural Network via Gradient Descent and Backpropagation	34
4	Accelerated Gradient Methods	35
4.1	Polyak's heavy ball method	35
4.2	Nesterov's Accelerated Gradient Descent (AGD)	36
4.3	Conjugate gradient method	40
4.4	Experiment: GD vs Accelerated GD	41

Chapter 1

Motivating Examples

1.1 Supervised Learning

EMPIRICAL RISK MINIMIZATION

Supervised Learning: Given training data points $(x_1, y_1), \dots, (x_n, y_n)$, construct a learning model $y = g(x, \omega)$ that best fits the training data. Here ω stands for the parameters of the learning model, say $\omega = (\omega_1, \dots, \omega_d)$.

Here (x_i, y_i) comes from an independent identically distributed family $(x_i, y_i) \sim p(x, y)$, where $p(x, y)$ is the joint density. The model $x \rightarrow y$ is a black-box $p(y|x)$, which is to be fit by $g(x, \omega)$.

“Loss function” $L(g(x, \omega), y)$, for example, can be $L(g(x, \omega), y) = (g(x, \omega) - y)^2$.

Empirical Risk Minimization (ERM):

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n L(y_i, g(x_i, \omega)) . \quad (1.1)$$

Regularized Empirical Risk Minimization (R-ERM):

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n L(y_i, g(x_i, \omega)) + \lambda R(\omega) . \quad (1.2)$$

For example, we can take $R(\omega) = \|\omega\|^2 = \omega_1^2 + \dots + \omega_d^2$. This regularization helps to control very irregular minimizers (unwanted ω).

In general let $f_i(\omega) = L(y_i, g(x_i, \omega))$ or $f_i(\omega) = L(y_i, g(x_i, \omega)) + \lambda R(\omega)$, then the optimization problem is

$$\omega_n^* = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n f_i(\omega) . \quad (1.3)$$

Key features of nonlinear optimization problem in machine learning: large-scale, nonconvex, ... etc.

Key problems in machine learning: optimization combined with generalization. “Population Loss”: $\mathbf{E}L(g(x, \omega), y)$, minimizer

$$\omega^* = \arg \min_{\omega} \mathbf{E}L(g(x, \omega), y) .$$

Generalization Error: $\mathbf{E}L(y, g(x, \omega_n^*))$. Consistency: Do we have $\omega_n^* \rightarrow \omega^*$? At what speed?

Key problems in optimization: convergence, acceleration, variance reduction.

How can optimization be related to generalization? There are quite abstract notions related to this topic, such as Vapnik–Chervonenkis dimension (VC dimension) and Radmacher complexity, which we might touch later. Also, we want to look at the geometry of the loss landscape, which is closely related to neural network structure.

LOSS FUNCTIONS

Classification Problems: label $y = 1$ or -1 . Choice of Loss function $L(y, g)$, $y = 1, -1$. 0/1 Loss: $\ell_{0/1}(y, g) = 1$ if $yg < 0$ and $\ell_{0/1}(y, g) = 0$ otherwise.

(1) Hinge Loss.

$$L(y, g) = \max(0, 1 - yg) ; \quad (1.4)$$

(2) Exponential Loss.

$$L(y, g) = \exp(-yg) ; \quad (1.5)$$

(3) Cross Entropy Loss.

$$L(y, g) = - \left(I_{\{y=1\}} \ln \frac{e^g}{e^g + e^{-g}} + I_{\{y=-1\}} \ln \frac{e^{-g}}{e^g + e^{-g}} \right) . \quad (1.6)$$

This is to use $p(y) = \frac{e^{yg}}{e^{yg} + e^{-yg}}$, $y = \pm 1$ and the binary cross entropy as

$$-(I_{\{y=1\}} \ln p(1) + I_{\{y=-1\}} \ln p(-1)) .$$

Regression Problems: Choice of Loss function $L(y, g)$.

(1) L^2 -Loss.

$$L(y, g) = |y - g|_2^2 . \quad (1.7)$$

L^2 -norm: $|x|_2^2 = x_1^2 + \dots + x_d^2$

(2) L^1 -Loss.

$$L(y, g) = |y - g|_1 . \quad (1.8)$$

L^1 -norm: $|x|_1 = |x_1| + \dots + |x_d|$.

(3) L^0 -Loss.

$$L(y, g) = |y - g|_0 . \quad (1.9)$$

L^0 -norm: $|x|_0 = \#\{i : x_i \neq 0, 1 \leq i \leq d\}$.

Regularized (penalize) term $R(\omega)$:

$$L^1\text{-regularized } R(\omega) = |\omega|_1;$$

L^0 -regularized $R(\omega) = |\omega|_0$.

LEARNING MODELS

(1) Linear regression: $g(x, \omega) = \omega^T x$. $g(x, \omega) = \frac{1}{1 + \exp(-\omega^T x)}$.

Least squares problem:

$$\min_{\omega \in \mathbb{R}^d} \frac{1}{2m} \sum_{j=1}^m (x_j^T \omega - y_j)^2 = \frac{1}{2m} \|A\omega - y\|_2^2 \quad (1.10)$$

Here training data $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in \mathbb{R}^d$, $y \in \mathbb{R}$, and $A = \begin{pmatrix} x_1^T \\ \dots \\ x_m^T \end{pmatrix}$.

Tikhonov regularization:

$$\min_{\omega} \frac{1}{2m} \|A\omega - y\|_2^2 + \lambda \|\omega\|_2^2. \quad (1.11)$$

LASSO (Least Absolute Shrinkage and Selection Operator):

$$\min_{\omega} \frac{1}{2m} \|A\omega - y\|_2^2 + \lambda \|\omega\|_1. \quad (1.12)$$

See [10].

(2) Support Vector Machines (SVM):

Set-up: $x_j \in \mathbb{R}^n$, $y_j \in \{1, -1\}$. Separating hyperplane $\omega^T x + \beta = 0$ where $\omega \in \mathbb{R}^n$ and $\beta \in \mathbb{R}$.

Classification Problem: Goal is to find a hyperplane $\omega^T x + \beta = 0$ such that it classifies the two kinds of data points most efficiently. The signed distance of any point $x \in \mathbb{R}^n$ to the hyperplane is given by $r = \frac{\omega^T x + \beta}{|\omega|}$. If the classification is good enough, we expect to have $\omega^T x_j + \beta > 0$ when $y = 1$ and $\omega^T x_j + \beta < 0$ when $y = -1$. After rescaling ω and β , we can then formulate the problem as looking for optimal ω and β such that $\omega^T x_j + \beta \geq 1$ when $y_j = 1$ and $\omega^T x_j + \beta \leq -1$ when $y_j = -1$. The closest few data points that match these two inequalities are called “support vectors”. The distance to the separating hyperplane created by two support vectors of opposite type is

$$\left| \frac{1}{|\omega|} \right| + \left| \frac{-1}{|\omega|} \right| = \frac{2}{|\omega|}.$$

So we can formulate the following optimization problem

$$\max_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{2}{|\omega|} \text{ such that } y_j(\omega^T x_j + \beta) \geq 1 \text{ for } j = 1, 2, \dots, m.$$

Or in other words we have the *constrained* optimization problem

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 \text{ such that } y_j(\omega^T x_j + \beta) \geq 1 \text{ for } j = 1, 2, \dots, m. \quad (1.13)$$

“Soft margin”: We allow the SVM to make errors on some training data points but we want to minimize the error. In fact, we allow some training data to violate $y_j(\omega^T x_j + \beta) \geq 1$, so that ideally we minimize

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 + C \sum_{j=1}^m \ell_{0/1}(y_j(\omega^T x_j + \beta) - 1).$$

Here the 0/1 loss is $\ell_{0/1}(z) = 1$ if $z < 0$ and $\ell_{0/1}(z) = 0$ otherwise, and $C > 0$ is a penalization parameter. We can then turn the 0/1 loss to Hinge Loss, that is why Hinge Loss comes in:

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 + C \sum_{j=1}^m \max(0, 1 - y_j(\omega^T x_j + \beta)) . \quad (1.14)$$

We can introduce “slack variables” $\xi_i \geq 0$ to introduce weights to classification errors in the above problem. This leads to “Soft margin SVM with slack variables”:

$$\min_{\omega \in \mathbb{R}^n, \beta \in \mathbb{R}} \frac{1}{2} |\omega|^2 + C \sum_{j=1}^m \xi_j \text{ s.t. } y_j(\omega^T x_j + b) \geq 1 - \xi_j, \xi_j \geq 0, j = 1, 2, \dots, m . \quad (1.15)$$

See [18].

(3) Neural Network: “activation function” σ .

Sigmoid:

$$\sigma(z) = \frac{1}{1 + \exp(-z)} ; \quad (1.16)$$

tanh:

$$\sigma(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} ; \quad (1.17)$$

ReLU (Rectified Linear Unit):

$$\sigma(z) = \max(0, z) . \quad (1.18)$$

Vector-valued activation: if $z \in \mathbb{R}^n$ then $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined by $(\sigma(z))_i = \sigma(z_i)$ where each $\sigma(z_i)$ is the scalar activation function.

Fully connected neural network prediction function

$$g(x, \omega) = a^T \left(\sigma \left(W^{(H)} (\sigma(W^{(H-1)} (\dots (\sigma(W^{(1)} x + b_1)) \dots) + b_{H-1}) + b_H) \right) \right) . \quad (1.19)$$

Optimization

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (g(x_i, \omega) - y_i)^2 .$$

Many other different network structures that we do not expand here: convolutional, recurrent (Gate: GRU, LSTM), ResNet, Transformer, ...

Two layer (one hidden layer) fully connected ReLU neural network has specific loss function

structure: our $W^{(1)} = \begin{pmatrix} \omega_1^T \\ \dots \\ \omega_m^T \end{pmatrix}$ and

$$g(x, \omega) = \sum_{r=1}^m a_r \max(\omega_r^T x, 0) . \quad (1.20)$$

Optimization problem is given by

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n \left(\sum_{r=1}^m a_r \max(\omega_r^T x_i, 0) - y_i \right)^2 .$$

Non-convexity issues: see [5].

1.2 Matrix Optimizations

Many machine learning/statistical learning problems are related to matrix optimizations.

(1) Matrix Completion: Each A_j is $n \times p$ matrix, and we seek for another $n \times p$ matrix \hat{X} such that

$$\hat{X} = \arg \min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 \quad (1.21)$$

where $\langle A, B \rangle = \text{tr}(A^T B)$. We can think of the A_j as “probing” the unknown matrix X . In other words, we want the best X such that $\langle A_j, X \rangle \approx y_j$ for all $1 \leq j \leq m$.

(2) Nonnegative Matrix Factorization: If the full matrix $Y \in \mathbb{R}^{n \times p}$ is observed, then we seek for $L \in \mathbb{R}^{n \times r}$ and $R \in \mathbb{R}^{p \times r}$ such that

$$\min_{L, R} \|LR^T - Y\|_F^2 \text{ subject to } L \geq 0 \text{ and } R \geq 0. \quad (1.22)$$

Here $\|A\|_F = (\sum \sum |a_{ij}|^2)^{1/2}$ is the Frobenius norm of a matrix A . This is used very often in recommendation systems (see [1]).

See [7] for an overview.

(3) Principle Component Analysis (PCA): Let S be a positive-definite (non-negative definite) matrix of size $n \times n$. Then we can diagonalize it as $Se_i = \lambda_i e_i$, $1 \leq i \leq n$, $\lambda_1 \geq \dots \geq \lambda_n > 0$ (or ≥ 0). Let $v \in \mathbb{R}^n$ be written as $v = v_1 e_1 + \dots + v_n e_n$. Then

$$v^T S v = \lambda_1 v_1^2 + \dots + \lambda_n v_n^2$$

is a quadratic form. If we restrict $v_1^2 + \dots + v_n^2 = 1$, then the maximum of above quadratic form will give the direction of e_1 and value λ_1 (principle component).

PCA:

$$\max_{v \in \mathbb{R}^n} v^T S v \text{ such that } \|v\|_2 = 1, \|v\|_0 \leq k. \quad (1.23)$$

Here S is a positive-definite (or non-negative definite) matrix. The objective function is convex, but if you take into account the constraint, then this problem becomes non-convex. A picture for dimension 1 example can be shown below.

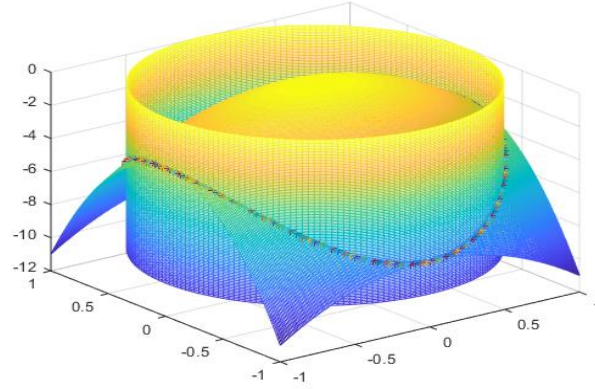


Figure 1.1: Loss Landscape of 1-dimensional PCA.

Online PCA: see [12].

(4) Sparse inverse covariance matrix estimation: Sample covariance matrix $S = \frac{1}{m-1} \sum_{j=1}^m a_j a_j^T$.
 $S^{-1} = X$. “Graphical LASSO”:

$$\min_{X \in \text{Symmetric } \mathbb{R}^{n \times n}, X \succeq 0} \langle S, X \rangle - \ln \det X + \lambda \|X\|_1 \quad (1.24)$$

where $\|X\|_1 = \sum |X_{ij}|$.

1.3 Experiment: Activation Functions

In this experiment we plot the following activation functions as well as their first and second derivatives. (1) Sigmoid; (2) ReLU; (3) Tanh; (4) Exponential. See [2], `src/activations`.

Chapter 2

Convex Functions

2.1 Optimization problem: set-up and its solutions

For the objective function $f : \mathbb{R}^n \supset \mathcal{D} \rightarrow \mathbb{R}$, we have various notions of minimizers:

- Local minimizer: $x^* \in \mathcal{D}$ is a local minimizer if there exist a neighborhood \mathcal{N} containing x^* , such that $f(x) \geq f(x^*)$ for all $x \in \mathcal{N} \cap \mathcal{D}$;
- Global minimizer: $x^* \in \mathcal{D}$ is a global minimizer if $f(x) \geq f(x^*)$ for all $x \in \mathcal{D}$;
- Strict Local Minimizer: $x^* \in \mathcal{D}$ is a strict local minimizer iff x^* is a local minimizer and $f(x) > f(x^*)$ if $x \neq x^*$, $x \in \mathcal{N}$;
- Isolated Local Minimizer: $x^* \in \mathcal{D}$ is an isolated local minimizer iff there exists a neighborhood \mathcal{N} containing x^* such that $f(x) \geq f(x^*)$ for all $x \in \mathcal{N} \cap \mathcal{D}$ and \mathcal{N} does not contain any other local minimizers.

Constrained optimization problem

$$\min_{x \in \Omega} f(x) , \tag{2.1}$$

where $\Omega \subset \mathcal{D} \subset \mathbb{R}^n$ is a closed set.

x^* is a *local solution*: there exist a neighborhood \mathcal{N} containing x^* such that $f(x) \geq f(x^*)$ for any $x \in \mathcal{N} \cap \Omega$.

x^* is a *global solution*: $f(x) \geq f(x^*)$ for any $x \in \Omega$.

2.2 Convexity

Definition 2.1 (Convex Set). A set Ω is called a convex set if for any $x, y \in \Omega$ we have

$$(1 - \alpha)x + \alpha y \in \Omega$$

for all $\alpha \in [0, 1]$.

Given convex set $\Omega \subset \mathbb{R}^n$, the *projection operator* $P : \mathbb{R}^n \rightarrow \Omega$ is given by

$$P(y) = \arg \min_{z \in \Omega} \|z - y\|_2^2 . \quad (2.2)$$

$P(y)$ is the point in Ω that is closest to y in the sense of Euclidean norm. If $x \in \Omega$, then $P(x) = x$.

Definition 2.2 (Convex Function). *A function $\phi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$ is convex if for all $x, y \in \mathbb{R}^n$ we have*

$$\phi((1 - \alpha)x + \alpha y) \leq (1 - \alpha)\phi(x) + \alpha\phi(y) \quad (2.3)$$

for all $\alpha \in [0, 1]$.

Definition 2.3 (Strongly Convex Function). *A convex function ϕ is called strongly convex with modulus of convexity $m > 0$ if*

$$\phi((1 - \alpha)x + \alpha y) \leq (1 - \alpha)\phi(x) + \alpha\phi(y) - \frac{1}{2}m\alpha(1 - \alpha)\|x - y\|_2^2 \quad (2.4)$$

for all x, y in the domain of ϕ .

Let $\Omega \subset \mathbb{R}^n$ be convex. We define the indicator function

$$I_\Omega(x) = \begin{cases} 0 & \text{if } x \in \Omega , \\ +\infty & \text{otherwise .} \end{cases} \quad (2.5)$$

Constrained optimization problem $\min_{x \in \Omega} f(x)$ is the same thing as the unconstrained optimization problem $\min[f(x) + I_\Omega(x)]$. We can set a sequence of functions $F_{\lambda\Omega} \uparrow I_\Omega(x)$ as $\lambda \uparrow \infty$, and then we solve the relaxed problem $\min_{x \in \mathbb{R}^d} [f(x) + F_{\lambda\Omega}(x)]$.

Theorem 2.1 (Minimizers for convex functions). *If the function f is convex and the set Ω is closed and convex, then*

- (a) *Any local solution of (2.1) is also global;*
- (b) *The set of global solutions of (2.1) is a convex set.*

Proof. (a) Suppose x_1^* is a local optimizer that is not global. This means there exists some $x_2^* \neq x_1^*$ such that $f(x_2^*) < f(x_1^*)$. Then by convexity for any $\alpha \in [0, 1]$ we have

$$f((1 - \alpha)x_1^* + \alpha x_2^*) \leq (1 - \alpha)f(x_1^*) + \alpha f(x_2^*) < f(x_1^*) .$$

If α is close to 0 and the above inequality will violate the fact that x_1^* is a local minimizer.

(b) Let S be the set of global minimizers. Then for any $x_1, x_2 \in S$ we have $f(x_1) = f(x_2) = \min_{x \in \Omega} f(x)$. By convexity we have

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2) = \min_{x \in \Omega} f(x) ,$$

which means that S is convex. □

A few more auxiliary notions:

- “Effective domain” of ϕ : $\{x \in \Omega : \phi(x) < \infty\}$;
- “Epigraph” of ϕ :

$$\text{epi}(\phi) = \{(x, t) \in \Omega \times \mathbb{R} : t \geq \phi(x)\} .$$

- ϕ is a “proper convex function” if $\phi(x) < \infty$ for some $x \in \Omega$ and $\phi(x) > -\infty$ for all $x \in \Omega$;
- ϕ is a “closed proper convex function” if ϕ is a proper convex function and $\{x \in \Omega : \phi(x) \leq t\}$ is a closed set for all $t \in \mathbb{R}$.

2.3 Taylor’s theorem and convexity in Taylor’s expansion

Theorem 2.2 (Taylor’s theorem). *Given a continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and given $x, p \in \mathbb{R}^n$ we have*

$$f(x + p) = f(x) + \int_0^1 \nabla f(x + \gamma p)^T p d\gamma , \quad (2.6)$$

$$f(x + p) = f(x) + \nabla f(x + \gamma p)^T p , \quad \text{for some } \gamma \in (0, 1) . \quad (2.7)$$

If f is twice continuously differentiable, then

$$\nabla f(x + p) = \nabla f(x) + \int_0^1 \nabla^2 f(x + \gamma p)^T p d\gamma , \quad (2.8)$$

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + \gamma p) p , \quad \text{for some } \gamma \in (0, 1) . \quad (2.9)$$

Proof. (1) Let $g(\gamma) = f(x + \gamma p)$, then $g'(\gamma) = \nabla f(x + \gamma p) \cdot p = (\nabla f(x + \gamma p))^T p$. Then by Newton-Leibniz we have

$$g(1) - g(0) = \int_0^1 g'(\gamma) d\gamma = \int_0^1 (\nabla f(x + \gamma p))^T p d\gamma ,$$

which gives (2.6).

(2) Using mean-value theorem $\int_a^b h(t) dt = (b - a)h(\xi)$ for some $\xi \in (a, b)$, we get

$$\int_0^1 (\nabla f(x + \gamma p))^T p d\gamma = (\nabla f(x + \gamma_0 p))^T p$$

for some $\gamma_0 \in (0, 1)$, which is (2.7).

Another way is to use the fact that

$$f(y) = f(x) + (\nabla f(x + \gamma_0(y - x)))^T (y - x) .$$

Set $p = y - x$, this gives

$$f(x + p) = f(x) + (\nabla f(x + \gamma_0 p))^T p .$$

(3) We apply (2.8) to each of $\frac{\partial f}{\partial x_i}$, and we get

$$\frac{\partial f}{\partial x_i}(x+p) = \frac{\partial f}{\partial x_i}(x) + \int_0^1 \left[\nabla \left(\frac{\partial f}{\partial x_i} \right) (x + \gamma p) \right]^T p d\gamma .$$

Introduce the Hessian matrix $\nabla^2 f(x) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{1 \leq i, j \leq n}$. Then when we put the above equation for all $1 \leq i \leq n$, we get (2.8).

(4) Define an auxiliary function of a single variable $\phi(t) = f(x+tp)$ for $t \in [0, 1]$. By applying the univariate Taylor's Theorem with the Lagrange form of the remainder to $\phi(t)$ at $t = 0$, we have:

$$\phi(1) = \phi(0) + \phi'(0)(1-0) + \frac{1}{2}\phi''(\gamma)(1-0)^2$$

for some $\gamma \in (0, 1)$.

Now, we compute the derivatives of $\phi(t)$ using the multi-variable chain rule:

$$\phi'(t) = \frac{d}{dt} f(x+tp) = \nabla f(x+tp)^T p$$

At $t = 0$, we obtain $\phi'(0) = \nabla f(x)^T p$.

$$\phi''(t) = \frac{d}{dt} (\nabla f(x+tp)^T p) = p^T \nabla^2 f(x+tp) p$$

where $\nabla^2 f$ is the Hessian matrix of f .

Substituting these expressions back into the Taylor expansion of $\phi(1)$: Since $\phi(1) = f(x+p)$ and $\phi(0) = f(x)$, we get (2.9). \square

Definition 2.4 (L -Lipschitz). *If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is such that*

$$|f(x) - f(y)| \leq L\|x - y\| , \quad (2.10)$$

for $L > 0$ and any $x, y \in \mathbb{R}^n$, then $f(x)$ is L -Lipschitz.

Optimization literatures often assume that ∇f is L -Lipschitz

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| . \quad (2.11)$$

Theorem 2.3. (1) *If f is continuously differentiable and convex then*

$$f(y) \geq f(x) + (\nabla f(x))^T (y - x) \quad (2.12)$$

for any $x, y \in \text{dom}(f)$.

(2) *If f is differentiable and m -strongly convex then*

$$f(y) \geq f(x) + (\nabla f(x))^T (y - x) + \frac{m}{2}\|y - x\|^2 \quad (2.13)$$

for any $x, y \in \text{dom}(f)$.

(3) If ∇f is uniformly Lipschitz continuous with Lipschitz constant $L > 0$ and f is convex then

$$f(y) \leq f(x) + (\nabla f(x))^T(y - x) + \frac{L}{2}\|y - x\|^2 \quad (2.14)$$

for any $x, y \in \text{dom}(f)$.

Proof. (1) Let $z_\alpha = \alpha x + (1 - \alpha)y$. Then by convexity

$$f(z_\alpha) \leq \alpha f(x) + (1 - \alpha)f(y) ,$$

which gives

$$f(z_\alpha) - f(x) \leq (1 - \alpha)(f(y) - f(x)) .$$

We can use (2.7) to get

$$f(z_\alpha) - f(x) = (\nabla f(x + \gamma(z_\alpha - x)))^T(z_\alpha - x)$$

for some $\gamma \in [0, 1]$. Since $z_\alpha - x = (1 - \alpha)(y - x)$, this gives

$$f(y) - f(x) \geq (\nabla f(x + \gamma(z_\alpha - x)))^T(y - x) .$$

Letting $\alpha \rightarrow 1$ we have $z_\alpha \rightarrow x$, that will give (2.12).

(2) We extend the above argument using strong convexity (2.4). This means we have

$$f(z_\alpha) \leq \alpha f(x) + (1 - \alpha)f(y) - \frac{1}{2}m\alpha(1 - \alpha)\|x - y\|_2^2 ,$$

which gives

$$f(z_\alpha) - f(x) + \frac{m}{2}\alpha(1 - \alpha)\|x - y\|^2 \leq (1 - \alpha)(f(y) - f(x)) .$$

By (2.7) again we have

$$(1 - \alpha)(\nabla f(x + \gamma(z_\alpha - x)))^T(y - x) + \frac{m}{2}\alpha(1 - \alpha)\|y - x\|^2 \leq (1 - \alpha)(f(y) - f(x)) .$$

Dividing both sides by $1 - \alpha$ and setting $\alpha \rightarrow 1$ we get (2.13).

(3) We apply (2.6) to get

$$f(y) = f(x) + \int_0^1 (\nabla f(x + \gamma(y - x)))^T(y - x) d\gamma .$$

This gives us

$$\begin{aligned} & f(y) - f(x) - (\nabla f(x))^T(y - x) \\ &= \int_0^1 [(\nabla f(x + \gamma(y - x)))^T - (\nabla f(x))^T](y - x) d\gamma \\ &\leq \int_0^1 \|(\nabla f(x + \gamma(y - x)))^T - (\nabla f(x))^T\| \cdot \|y - x\| d\gamma \\ &\leq \int_0^1 L\gamma\|y - x\|^2 d\gamma = \frac{L}{2}\|y - x\|^2 , \end{aligned}$$

which is (2.14). □

We use the symbol $A \succeq B$ ($A \preceq B$) to mean $\lambda_A \geq \lambda_B$ ($\lambda_A \leq \lambda_B$) for every corresponding pair of eigenvalues of A, B .

Theorem 2.4. *Suppose that the function f is twice continuously differentiable on \mathbb{R}^n . Then*

- (1) *f is strongly convex with modulus of convexity m if and only if $\nabla^2 f(x) \succeq mI$ for all x ;*
- (2) *∇f is Lipschitz continuous with Lipschitz constant L if and only if $\nabla^2 f(x) \preceq LI$ for all x .*

Proof. (1) Suppose that the function f is strongly- m convex. Set $u \in \mathbb{R}^n$ and $\alpha > 0$, then we consider the Taylor's expansion

$$f(x + \alpha u) = f(x) + \alpha \nabla f(x)^T u + \frac{1}{2} \alpha^2 u^T \nabla^2 f(x + t\alpha u) u \quad (2.15)$$

for some $0 \leq t \leq 1$. We apply (2.13) with $y = x + \alpha u$ so that

$$f(x + \alpha u) \geq f(x) + \alpha (\nabla f(x))^T u + \frac{m}{2} \alpha^2 \|u\|^2. \quad (2.16)$$

Comparing (2.15) and (2.16) we see that for arbitrary choice of $u \in \mathbb{R}^n$ we have

$$u^T \nabla^2 f(x + t\alpha u) u \geq m \|u\|^2.$$

This implies that $\nabla^2 f(x) \succeq mI$ as claimed. This shows the “only if” part.

Indeed one can also show the “if” part. To this end assume that $\nabla^2 f(x) \succeq mI$. Then for any $z \in \mathbb{R}^n$ we have that $(x - z)^T \nabla^2 f(z + t(x - z))(x - z) \geq m \|x - z\|^2$. Thus

$$\begin{aligned} f(x) &= f(z) + (\nabla f(z))^T (x - z) + \frac{1}{2} (x - z)^T \nabla^2 f(z + t(x - z))(x - z) \\ &\geq f(z) + (\nabla f(z))^T (x - z) + \frac{m}{2} \|x - z\|^2. \end{aligned} \quad (2.17)$$

Similarly

$$\begin{aligned} f(y) &= f(z) + (\nabla f(z))^T (y - z) + \frac{1}{2} (y - z)^T \nabla^2 f(z + t(y - z))(y - z) \\ &\geq f(z) + (\nabla f(z))^T (y - z) + \frac{m}{2} \|y - z\|^2. \end{aligned} \quad (2.18)$$

We consider $(1 - \alpha)(2.17) + \alpha(2.18)$ and we set $z = (1 - \alpha)x + \alpha y$. This gives

$$\begin{aligned} &(1 - \alpha)f(x) + \alpha f(y) \\ &\geq (\alpha + (1 - \alpha))f(z) + (\nabla f(z))^T ((1 - \alpha)(x - z) + \alpha(y - z)) + \frac{m}{2} ((1 - \alpha)\|x - z\|^2 + \alpha\|y - z\|^2) \\ &= f(z) + (\nabla f(z))^T ((1 - \alpha)(x - z) + \alpha(y - z)) + \frac{m}{2} ((1 - \alpha)\|x - z\|^2 + \alpha\|y - z\|^2). \end{aligned} \quad (2.19)$$

Since $x - z = \alpha(x - y)$ and $y - z = (1 - \alpha)(y - x)$, we see that $((1 - \alpha)(x - z) + \alpha(y - z)) = 0$. Moreover, this means that

$$(1 - \alpha)\|x - z\|^2 + \alpha\|y - z\|^2 = [(1 - \alpha)\alpha^2 + \alpha(1 - \alpha)^2] \|x - y\|^2 = \alpha(1 - \alpha)\|x - y\|^2.$$

From these we see that (2.19) is the same as saying

$$(1 - \alpha)f(x) + \alpha f(y) \geq f((1 - \alpha)x + \alpha y) + \frac{m}{2} \alpha(1 - \alpha)\|x - y\|^2,$$

which is (2.4).

(2) We want to show $\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$ is equivalent to $\nabla^2 f(x) \preceq LI$. Assume $\nabla^2 f(x) \preceq LI$, which means the spectral norm $\|\nabla^2 f(x)\|_2 \leq L$. Using (2.8) we get

$$\begin{aligned} \|\nabla f(y) - \nabla f(x)\| &= \left\| \int_0^1 \nabla^2 f(x + t(y-x))(y-x) dt \right\| \\ &\leq \int_0^1 \|\nabla^2 f(x + t(y-x))\| \cdot \|y-x\| dt \\ &\leq \int_0^1 L\|y-x\| dt = L\|y-x\|, \end{aligned}$$

which proves ∇f is L -Lipschitz continuous. Assume ∇f is L -Lipschitz. For any vector v , consider the limit:

$$\nabla^2 f(x)v = \lim_{t \rightarrow 0} \frac{\nabla f(x+tv) - \nabla f(x)}{t}$$

Taking the norm:

$$\|\nabla^2 f(x)v\| = \lim_{t \rightarrow 0} \frac{\|\nabla f(x+tv) - \nabla f(x)\|}{t} \leq \lim_{t \rightarrow 0} \frac{L\|tv\|}{t} = L\|v\|$$

This implies the operator norm $\|\nabla^2 f(x)\|_2 \leq L$, which for a symmetric Hessian matrix is equivalent to $-LI \preceq \nabla^2 f(x) \preceq LI$. Since we usually deal with convex/semi-convex functions in this context, it implies $\nabla^2 f(x) \preceq LI$. \square

Theorem 2.5 (Strongly convex functions have unique minimizers). *Let f be differentiable and strongly convex with modulus $m > 0$. Then the minimizer x^* of f exists and is unique.*

Proof. Uniqueness. Suppose there exist two distinct minimizers x^* and z^* . Since f is differentiable, we must have $\nabla f(x^*) = 0$ and $\nabla f(z^*) = 0$. Using the strong convexity property (2.13) at x^* :

$$f(z^*) \geq f(x^*) + \nabla f(x^*)^T(z^* - x^*) + \frac{m}{2}\|z^* - x^*\|^2 = f(x^*) + \frac{m}{2}\|z^* - x^*\|^2.$$

Similarly, using strong convexity at z^* :

$$f(x^*) \geq f(z^*) + \frac{m}{2}\|x^* - z^*\|^2.$$

Adding these two inequalities yields $f(x^*) + f(z^*) \geq f(x^*) + f(z^*) + m\|x^* - z^*\|^2$, which simplifies to $0 \geq m\|x^* - z^*\|^2$. Since $m > 0$, it must be that $\|x^* - z^*\| = 0$, hence $x^* = z^*$.

Existence. From the strong convexity definition, $f(y) \rightarrow \infty$ as $\|y\| \rightarrow \infty$ (coercivity). Since f is continuous and coercive, it must attain a minimum on any closed set in \mathbb{R}^n . \square

Below are a few more technical issues regarding convex and strongly convex functions.

Theorem 2.6. *Let f be convex and continuously differentiable and ∇f be L -Lipschitz. Then for any $x, y \in \text{dom}(f)$ the following bounds hold:*

$$f(x) + \nabla f(x)^T(y - x) + \frac{1}{2L}\|\nabla f(x) - \nabla f(y)\|^2 \leq f(y), \quad (2.20)$$

$$\frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x - y) \leq L \|x - y\|^2 . \quad (2.21)$$

In addition, let f be strongly convex with modulus m and unique minimizer x^* . Then for any $x, y \in \text{dom}(f)$ we have that

$$f(y) - f(x) \geq -\frac{1}{2m} \|\nabla f(x)\|^2 . \quad (2.22)$$

Proof. Define an auxiliary function $g(z) = f(z) - \nabla f(x)^T z$. Since f is convex, g is also convex. Furthermore, $\nabla g(z) = \nabla f(z) - \nabla f(x)$, which is also L -Lipschitz continuous. Note that $\nabla g(x) = 0$, implying x is a global minimizer of $g(z)$ via Theorem 2.1. Using (2.14) we get

$$g(z) \leq g(y) + \nabla g(y)^T (z - y) + \frac{L}{2} \|z - y\|^2 .$$

Setting $z = y - \frac{1}{L} \nabla g(y)$, and since x minimizes g , we obtain:

$$g(x) \leq g(y) - \frac{1}{2L} \|\nabla g(y)\|^2 .$$

Substituting $g(z) = f(z) - \nabla f(x)^T z$ back into the inequality:

$$f(x) - \nabla f(x)^T x \leq f(y) - \nabla f(x)^T y - \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2 ,$$

which is (2.20).

We swap x and y in (2.20) to get

$$f(y) + \nabla f(y)^T (x - y) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2 \leq f(x) .$$

Adding this to (2.20) and canceling $f(x) + f(y)$ from both sides we get

$$\nabla f(x)^T (y - x) + \nabla f(y)^T (x - y) + \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq 0 ,$$

which gives the first inequality in (2.21) by rearranging.

The second inequality in (2.21) follows directly from Cauchy-Schwarz and the L -Lipschitz property:

$$(\nabla f(x) - \nabla f(y))^T (x - y) \leq \|\nabla f(x) - \nabla f(y)\| \cdot \|x - y\| \leq L \|x - y\|^2 .$$

Assume f is strongly convex with modulus m . Then we have

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|^2 \\ &= f(x) + \frac{m}{2} \left[\|y - x\|^2 + \frac{2}{m} \nabla f(x)^T (y - x) \right] \\ &= f(x) + \frac{m}{2} \left[\|y - x\|^2 + 2 \left(\frac{1}{m} \nabla f(x) \right)^T (y - x) + \frac{1}{m^2} \|\nabla f(x)\|^2 \right] - \frac{1}{2m} \|\nabla f(x)\|^2 \\ &= f(x) + \frac{m}{2} \left\| y - x + \frac{1}{m} \nabla f(x) \right\|^2 - \frac{1}{2m} \|\nabla f(x)\|^2 , \end{aligned}$$

which gives (2.22). □

Definition 2.5 (Generalized Strong Convexity). *We say the convex function $f(x)$ has Generalized Strong Convexity if*

$$\|\nabla f(x)\|^2 \geq 2m[f(x) - f^*] \text{ for some } m > 0, \quad (2.23)$$

where $f^* = f(x^*)$ is the minimum of f .

The Generalized Strong Convexity holds in situations other than when f is strongly convex.

2.4 Optimality Conditions

We derive optimality conditions for the smooth unconstrained problem

$$\min_{x \in \mathbb{R}^n} f(x), \quad (2.24)$$

where $f(x)$ is a smooth function.

Theorem 2.7 (Necessary conditions for smooth unconstrained optimization). *We have*

- (a) (first-order necessary condition) *Suppose that f is continuously differentiable. Then if x^* is a local minimizer of (2.24), then $\nabla f(x^*) = 0$;*
- (b) (second-order necessary condition) *Suppose that f is twice continuously differentiable. Then if x^* is a local minimizer of (2.24), then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semi-definite.*

Proof. Suppose x^* is a local minimizer of f . Then for any direction $d \in \mathbb{R}^n$ and sufficiently small step size $\alpha > 0$, we have $f(x^* + \alpha d) \geq f(x^*)$.

(a) By (2.9) we have

$$f(x^* + \alpha d) = f(x^*) + \alpha \nabla f(x^*)^T d + O(\alpha^2)$$

Since $f(x^* + \alpha d) \geq f(x^*)$, it follows that $\alpha \nabla f(x^*)^T d + O(\alpha^2) \geq 0$. Dividing by α and letting $\alpha \rightarrow 0$, we obtain $\nabla f(x^*)^T d \geq 0$. Since this holds for any d , we choose $d = -\nabla f(x^*)$, which gives $-\|\nabla f(x^*)\|^2 \geq 0$. Thus, $\nabla f(x^*) = 0$.

(b) Given $\nabla f(x^*) = 0$, by (2.9) we get

$$f(x^* + \alpha d) = f(x^*) + \frac{\alpha^2}{2} d^T \nabla^2 f(x^*) d + o(\alpha^2)$$

For $f(x^* + \alpha d) \geq f(x^*)$ to hold for small α , we must have $d^T \nabla^2 f(x^*) d \geq 0$ for all d . This implies that the Hessian matrix $\nabla^2 f(x^*)$ is positive semi-definite. \square

Theorem 2.8 (Sufficient conditions for convex functions in smooth unconstrained optimization). *If f is continuously differentiable and convex, then if $\nabla f(x^*) = 0$, then x^* is a global minimizer of (2.24). In addition, if f is strongly convex, then x^* is the unique global minimizer.*

Proof. Assume f is convex and $\nabla f(x^*) = 0$. By (2.12), for any $x \in \mathbb{R}^n$ we have

$$f(x) \geq f(x^*) + \nabla f(x^*)^T (x - x^*) .$$

Since $\nabla f(x^*) = 0$, this gives $f(x) \geq f(x^*)$. This proves that x^* is a global minimizer. If f is further assumed to be strongly convex, then x^* is the unique global minimizer due to Theorem 2.5. \square

Theorem 2.9 (Sufficient conditions for nonconvex functions in smooth unconstrained optimization). *Suppose that f is twice continuously differentiable and that for some x^* we have $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite, then x^* is a strict local minimizer of (2.24).*

Proof. Suppose $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite. Let $\lambda_{\min} > 0$ be the smallest eigenvalue of $\nabla^2 f(x^*)$. Then for any $d \neq 0$:

$$d^T \nabla^2 f(x^*) d \geq \lambda_{\min} \|d\|^2 .$$

Using (2.9), the second-order Taylor expansion around x^* , we get

$$f(x^* + d) = f(x^*) + \nabla f(x^*)^T d + \frac{1}{2} d^T \nabla^2 f(x^*) d + o(\|d\|^2) .$$

Since $\nabla f(x^*) = 0$, we have:

$$f(x^* + d) \geq f(x^*) + \frac{\lambda_{\min}}{2} \|d\|^2 + o(\|d\|^2)$$

For sufficiently small $\|d\|$, the term $\frac{\lambda_{\min}}{2} \|d\|^2$ dominates the higher-order terms $o(\|d\|^2)$, ensuring $f(x^* + d) > f(x^*)$. Thus, x^* is a strict local minimizer. \square

2.5 Experiment: Nonconvexity of the Loss Landscape of Neural Networks

Consider a neural network with one hidden layer that consists of p neurons and input $x \in \mathbb{R}^1$, output $y \in \mathbb{R}^1$. The neural network function has the form $y(x) = \sum_{j=1}^p c_j \sigma(a_j x - b_j)$, where $\mathbf{a} = a_j$, $\mathbf{b} = b_j$ and $\mathbf{c} = c_j$ are the neural network weights. Assume $(a_j, b_j, c_j) \sim \mathcal{N}(0, I_3)$, $j = 1, 2, \dots, p$ is a family of i.i.d multivariate normal distributions. For different realizations of (a_j, b_j, c_j) , we plot the function $y(x)$ on x - y graph. We have experimented different hidden layer sizes.

With the above done, assume that the training data (x, y) follows a bivariate normal distribution $\mathcal{N}(0, I_2)$. Let the variables a_1 and a_2 vary in a certain interval. Let all other neural network weights a_j, b_j, c_j follow the same assumption as above. For different (a_1, a_2) , we plot the empirical loss function of this network as a function of (a_1, a_2) . We can see that the loss function may not be convex with respect to (a_1, a_2) . We show experiment results for different hidden layer sizes.

We refer to [2], `src/one_hidden_layer_nn`.

A typical non-convex loss landscape looks as follows:

Sigmoid empirical loss landscape, hidden layer neuron number=20, training size=10

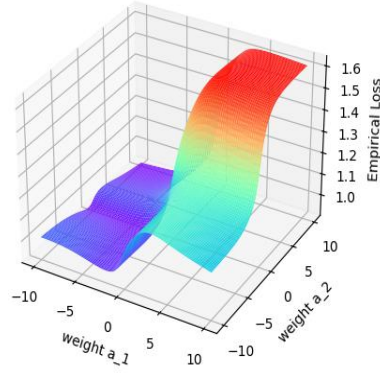


Figure 2.1: Loss Landscape of a neural network.

In the branch `support_gpu` (already merged to `main`), we supported GPU computations, with some optimizations made in the computing flow. This is especially useful when the hidden layer size is huge, where we expect to see the NTK regime (see [4]), and the loss landscape tends to be convex.

Sigmoid empirical loss landscape, hidden layer size=10000, training size=10

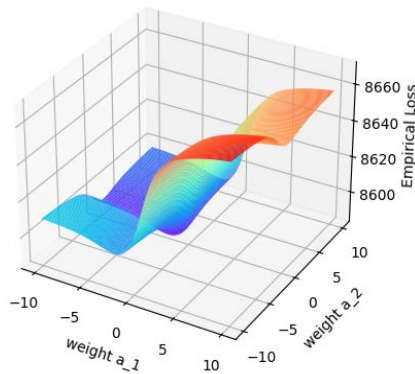


Figure 2.2: Loss Landscape of a neural network, 10000 neurons in the hidden layer.

2.5. EXPERIMENT: NONCONVEXITY OF THE LOSS LANDSCAPE OF NEURAL NETWORKS²¹

Sigmoid empirical loss landscape, hidden layer size=100000, training size=10

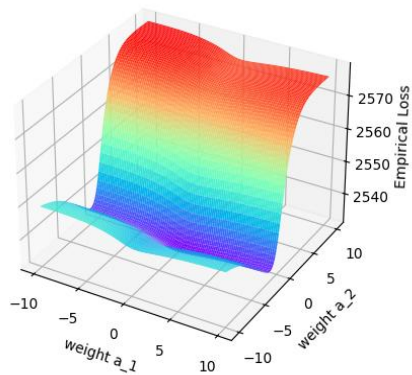


Figure 2.3: Loss Landscape of a neural network, 100000 neurons in the hidden layer.

Sigmoid empirical loss landscape, hidden layer size=1000000, training size=10

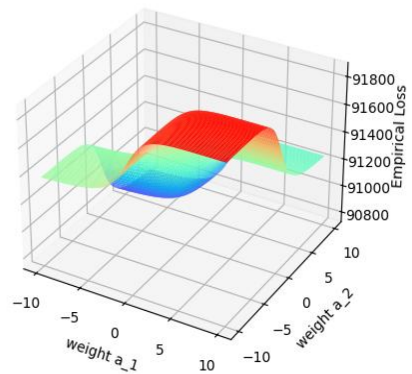


Figure 2.4: Loss Landscape of a neural network, 1000000 neurons in the hidden layer.

Sigmoid empirical loss landscape, hidden layer size=10000000, training size=10

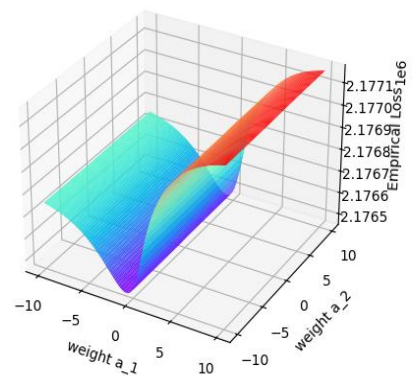


Figure 2.5: Loss Landscape of a neural network, 10000000 neurons in the hidden layer.

Chapter 3

Gradient Descent and Line Search Methods

3.1 Gradient Descent

We want to solve the optimization problem

$$\min_{x \in \mathbb{R}^n} f(x) .$$

The goal is to construct an approximating sequence $\{x^k\}$ such that $f(x^{k+1}) < f(x^k)$, $k = 0, 1, 2, \dots$

“descent direction” $d \in \mathbb{R}^n$ such that

$$f(x + td) < f(x) \text{ for all } t > 0 \text{ sufficiently small.}$$

Say f is continuously differentiable

$$f(x + td) = f(x) + t \nabla f(x)^T d + o(t) \text{ for some } \gamma \in (0, 1) .$$

Definition 3.1. *The vector $d \in \mathbb{R}^n$ is a descent direction if*

$$d^T \nabla f(x) < 0 .$$

It is a steepest descent direction if

$$\inf_{\|d\|=1} d^T \nabla f(x) = -\|\nabla f(x)\| .$$

The inf is achieved when $d = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$. This gives the Gradient Descent Method:

Algorithm 3.1 Gradient Descent

```

1: Input: Input initialization  $x^0$ , stepsize  $\alpha_k > 0$ ,  $k = 0, 1, 2, \dots$ 
2: for  $k = 0, 1, 2, \dots, K - 1$  do
3:   Calculate  $\nabla f(x^k)$ 
4:    $x^{k+1} = x^k - \alpha_k \nabla f(x^k)$ 
5: end for
6: Output: Output  $x^K$ 

```

Definition 3.2 (Convergence of optimization algorithms). *We consider any of the following estimates as characterizing the convergence of optimization algorithms:*

$$\|x_T - x^*\| \leq \varepsilon(T) \quad (3.1)$$

or

$$\mathbf{E}\|x_T - x^*\|^2 \leq \varepsilon(T) \quad (3.2)$$

or

$$\mathbf{E}f(x_T) - f(x^*) \leq \varepsilon(T) . \quad (3.3)$$

Here $\varepsilon(T) \rightarrow 0$ as $T \rightarrow \infty$. The convergence rates are measured by $\ln \varepsilon(T)$. We have the following cases:

- (1) If $\ln \varepsilon(T) = O(-T)$, then we say the algorithm has linear convergence rate;
- (2) If $\ln \varepsilon(T)$ decays slower than $O(-T)$, then we say the algorithm has sub-linear convergence rate;
- (3) If $\ln \varepsilon(T)$ decays faster than $O(-T)$, then we say the algorithm has super-linear convergence rate. In particular, if $\ln \ln \varepsilon(T) = O(-T)$, then we say the algorithm has second-order convergence rate.

Theorem 3.1. *Assume that f is twice continuously differentiable, m -strongly convex and ∇f is L -Lipschitz, so that $mI \preceq \nabla^2 f(x) \preceq LI$ for all $x \in \mathbb{R}^n$. Then the GD Algorithm 3.1 with constant stepsize $\alpha_k = \alpha$ and $\alpha \in \left[\frac{1-\beta}{m}, \frac{1+\beta}{L} \right]$ for some $0 \leq \beta < 1$ has linear convergence rate.*

Proof. Consider the optimization problem $\min_{x \in \mathbb{R}^n} f(x)$. We assume f is twice continuously differentiable and satisfy $mI \preceq \nabla^2 f(x) \preceq LI$ for all $x \in \mathbb{R}^n$.

The Gradient Descent (GD) update rule is given by $x^{k+1} = x^k - \alpha \nabla f(x^k)$. Let $\Phi(x) = x - \alpha \nabla f(x)$ be the iteration operator, so that $x^{k+1} = \Phi(x^k)$. Since x^* is a local minimizer, $\nabla f(x^*) = 0$, which implies x^* is a fixed point:

$$x^* = \Phi(x^*) .$$

By the Mean Value Theorem for vector-valued functions, for any $x, z \in \mathbb{R}^n$, there exists a point c on the line segment $[x, z]$ such that:

$$\begin{aligned}\Phi(x) - \Phi(z) &= \nabla \Phi(c)(x - z) \\ &= (I - \alpha \nabla^2 f(c))(x - z) .\end{aligned}$$

Taking the spectral norm on both sides:

$$\|\Phi(x) - \Phi(z)\| \leq \|I - \alpha \nabla^2 f(c)\|_2 \cdot \|x - z\| .$$

To ensure linear convergence, we require Φ to be a contraction mapping with factor $\beta < 1$, i.e., $\|I - \alpha \nabla^2 f(c)\|_2 \leq \beta$.

Given $mI \preceq \nabla^2 f(x) \preceq LI$, the eigenvalues of the matrix $I - \alpha \nabla^2 f(c)$ lie in the interval $[1 - \alpha L, 1 - \alpha m]$. To find the optimal α that minimizes the maximum absolute eigenvalue β , we set:

$$-\beta = 1 - \alpha L \quad \text{and} \quad \beta = 1 - \alpha m$$

Solving this system yields:

1. **Optimal Stepsize:** $\alpha = \frac{2}{L+m}$
2. **Contraction Factor:** $\beta = \frac{L-m}{L+m}$

Since $L \geq m > 0$, it follows that $0 \leq \beta < 1$.

Applying the contraction property iteratively:

$$\|x^T - x^*\| = \|\Phi(x^{T-1}) - \Phi(x^*)\| \leq \beta \|x^{T-1} - x^*\| \leq \dots \leq \beta^T \|x^0 - x^*\| .$$

To achieve a precision $\|x^T - x^*\| \leq \epsilon$, we require:

$$\beta^T \|x^0 - x^*\| \leq \epsilon \implies T \ln \beta \leq \ln \left(\frac{\epsilon}{\|x^0 - x^*\|} \right) .$$

Since $\ln \beta < 0$, we have:

$$T \geq \frac{\ln(\|x^0 - x^*\|/\epsilon)}{|\ln \beta|} .$$

This confirms that the algorithm converges at a linear rate. □

Taylor's expansion can give us a standard scheme for obtaining the convergence of optimization algorithms. To this end we consider

Lemma 3.2 (Descent Lemma). *We have*

$$f(x^{k+1}) \leq f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2 . \tag{3.4}$$

Proof. We apply Taylor's theorem to get

$$\begin{aligned}
 f(x + \alpha d) &= f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 [\nabla f(x + \gamma \alpha d) - \nabla f(x)]^T d \, d\gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 \|\nabla f(x + \gamma \alpha d) - \nabla f(x)\| \|d\| \, d\gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha^2 \frac{L}{2} \|d\|^2.
 \end{aligned} \tag{3.5}$$

From (3.5) and the line search $d = -\nabla f(x^k)$, $x = x^k$ on RHS of (3.5), then the $\alpha = \frac{1}{L}$ minimizes RHS of (3.5), and thus for GD in particular, we have

$$f(x^{k+1}) \leq f(x^k) + \left(-\alpha + \alpha^2 \frac{L}{2}\right) \|\nabla f(x^k)\|^2 = f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2.$$

□

Estimates of type (3.4) is a standard procedure that leads to convergence analysis of iterative optimization algorithms. Using it, for GD Algorithm 3.1 with $\alpha = \frac{1}{L}$ we have

Theorem 3.3. *Assume ∇f is L -Lipschitz. The convergence rates of GD algorithm has the following cases:*

(a) *If f is general nonconvex, then*

$$\min_{0 \leq k \leq T-1} \|\nabla f(x^k)\| \leq \sqrt{\frac{2L(f(x^0) - f(x^*))}{T}}.$$

Number of iterations $0 \leq k \leq T$ until $\|\nabla f(x^k)\| \leq \varepsilon$ is $T \geq \frac{2L(f(x^0) - f(x^))}{\varepsilon^2}$. This is sublinear convergence;*

(b) *If f is convex, then when the stepsize is $\alpha = \frac{1}{L}$ we have*

$$f(x^k) - f(x^*) \leq \frac{L}{2k} \|x^0 - x^*\|^2.$$

Number of iterations $0 \leq k \leq T$ until $f(x^k) - f(x^) \leq \varepsilon$ is $T \geq \frac{f(x^0) - f(x^*)}{\varepsilon}$. This is sublinear convergence;*

(c) *If f is strongly convex with convexity constant $m > 0$, then*

$$f(x^k) - f(x^*) \leq \left(1 - \frac{m}{L}\right)^k (f(x^0) - f(x^*)).$$

Number of iterations $0 \leq k \leq T$ until $f(x^k) - f(x^) \leq \varepsilon$ is $T \geq \frac{L}{m} \ln \left(\frac{f(x^0) - f(x^*)}{\varepsilon} \right)$. This is linear convergence.*

Proof. (a) Summing the Descent Lemma from $k = 0$ to $T - 1$:

$$\sum_{k=0}^{T-1} (f(x^k) - f(x^{k+1})) \geq \sum_{k=0}^{T-1} \frac{1}{2L} \|\nabla f(x^k)\|^2 .$$

The left side telescopes to $f(x^0) - f(x^T)$. Since $f(x^T) \geq f^*$:

$$f(x^0) - f^* \geq \frac{1}{2L} \sum_{k=0}^{T-1} \|\nabla f(x^k)\|^2 \geq \frac{T}{2L} \min_{0 \leq k \leq T-1} \|\nabla f(x^k)\|^2$$

Rearranging gives the sublinear rate $O(1/\sqrt{T})$ for the gradient norm:

$$\min_{0 \leq k \leq T-1} \|\nabla f(x^k)\| \leq \sqrt{\frac{2L(f(x^0) - f^*)}{T}} .$$

(b) By the first-order condition for convex functions $f(x^*) \geq f(x^k) + \nabla f(x^k)^T(x^* - x^k)$, and using the Descent Lemma:

$$\begin{aligned} f(x^{k+1}) &\leq f(x^*) + \nabla f(x^k)^T(x^k - x^*) - \frac{1}{2L} \|\nabla f(x^k)\|^2 \\ &= f(x^*) + \frac{L}{2} \left[\frac{2}{L} \nabla f(x^k)^T(x^k - x^*) - \frac{1}{L^2} \|\nabla f(x^k)\|^2 \right] . \end{aligned}$$

Using the identity $\|a - b\|^2 = \|a\|^2 - 2a^T b + \|b\|^2$ and the update $x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k)$:

$$f(x^{k+1}) \leq f(x^*) + \frac{L}{2} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2) .$$

Summing from $k = 0$ to $T - 1$ and noting that $f(x^k)$ is non-increasing:

$$T(f(x^T) - f^*) \leq \sum_{k=0}^{T-1} (f(x^{k+1}) - f^*) \leq \frac{L}{2} \|x^0 - x^*\|^2 .$$

This yields the sublinear rate $O(1/T)$ in function value.

(c) For m -strongly convex functions, we have the bound $\|\nabla f(x)\|^2 \geq 2m(f(x) - f^*)$. Substituting this into the Descent Lemma:

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) - \frac{1}{2L} (2m(f(x^k) - f^*)) , \\ f(x^{k+1}) - f^* &\leq \left(1 - \frac{m}{L}\right) (f(x^k) - f^*) . \end{aligned}$$

Applying this recursively T times gives the linear convergence rate:

$$f(x^T) - f^* \leq \left(1 - \frac{m}{L}\right)^T (f(x^0) - f^*) .$$

□

3.2 Back Propagation and GD on Neural Networks

Classical algorithm of training multi-layer neural networks: Back-propagation. This algorithm was proposed in [16]. For a modern approach with the flavor of applied mathematics, we refer to [11].

Back-propagation method goes as follows. Let us set

$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]} \text{ for } l = 2, 3, \dots, L ,$$

in which l stands for the number of layers in the neural network. The neural network iteration can be viewed as

$$a^{[l]} = \sigma(z^{[l]}) \text{ for } l = 2, 3, \dots, L .$$

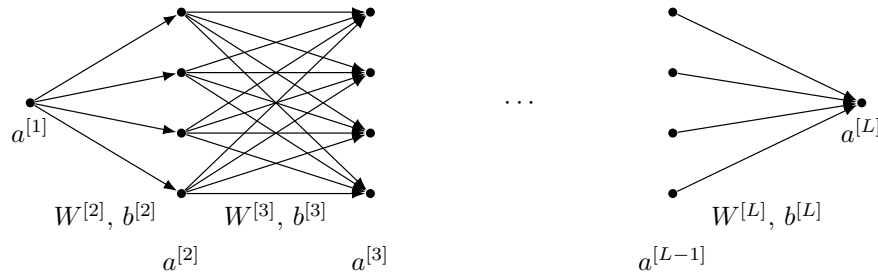


Figure 3.1: Feedforward fully-connected neural network (intermediate layers omitted).

We can view the neural network function as a chain with $a^{[1]} = x$ and $a^{[L]} = g(x; \omega)$. Let the training data be one point (x, y) . Then the loss function is given by

$$C = C(\omega, b) = \frac{1}{2} \|y - a^{[L]}\|^2 .$$

To perform the GD algorithm as in Algorithm 3.1, we shall need the derivatives of the loss function C with respect to the neural-network parameters ω and b , that is $\frac{\partial C}{\partial w_{jk}^{[l]}}$ and $\frac{\partial C}{\partial b_j^{[l]}}$. We first have

Definition 3.3 (Error in Backpropagation). *We say the error in the j -th neuron at layer l to be*

$$\delta^{[l]} = (\delta_j^{[l]})_j$$

(viewed as a column vector) with

$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} .$$

Definition 3.4 (Hadamard Product). *If $x, y \in \mathbb{R}^n$, let $x \circ y \in \mathbb{R}^n$ to be the Hadamard product defined by*

$$(x \circ y)_i = x_i y_i .$$

Then one can calculate directly to obtain the following lemma

Lemma 3.4. *We have*

$$\delta^{[L]} = \sigma'(z^{[L]}) \circ (a^{[L]} - y) , \quad (3.6)$$

$$\delta^{[l]} = \sigma'(z^{[l]}) \circ (W^{[l+1]})^T \delta^{[l+1]} , \text{ for } 2 \leq l \leq L-1 , \quad (3.7)$$

$$\frac{\partial C}{\partial b_j^{[l]}} = \delta_j^{[l]} , \text{ for } 2 \leq l \leq L , \quad (3.8)$$

$$\frac{\partial C}{\partial \omega_{jk}^{[l]}} = \delta_j^{[l]} a_k^{[l-1]} , \text{ for } 2 \leq l \leq L . \quad (3.9)$$

Proof. We prove each identity using the chain rule.

Proof of (3.6). Write the loss function componentwise as

$$C = \frac{1}{2} \sum_{j=1}^{n_L} (a_j^{[L]} - y_j)^2.$$

Then

$$\frac{\partial C}{\partial a_j^{[L]}} = a_j^{[L]} - y_j.$$

Since $a_j^{[L]} = \sigma(z_j^{[L]})$, the chain rule gives

$$\delta_j^{[L]} = \frac{\partial C}{\partial z_j^{[L]}} = \frac{\partial C}{\partial a_j^{[L]}} \frac{\partial a_j^{[L]}}{\partial z_j^{[L]}} = (a_j^{[L]} - y_j) \sigma'(z_j^{[L]}).$$

Writing this in vector form yields

$$\delta^{[L]} = \sigma'(z^{[L]}) \circ (a^{[L]} - y).$$

Proof of (3.7). Fix $l \in \{2, \dots, L-1\}$. Using the chain rule,

$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \sum_{m=1}^{n_{l+1}} \frac{\partial C}{\partial z_m^{[l+1]}} \frac{\partial z_m^{[l+1]}}{\partial z_j^{[l]}}.$$

Since $\frac{\partial C}{\partial z_m^{[l+1]}} = \delta_m^{[l+1]}$, it remains to compute

$$z_m^{[l+1]} = \sum_{k=1}^{n_l} w_{mk}^{[l+1]} a_k^{[l]} + b_m^{[l+1]}, \quad a_k^{[l]} = \sigma(z_k^{[l]}).$$

Therefore,

$$\frac{\partial z_m^{[l+1]}}{\partial z_j^{[l]}} = w_{mj}^{[l+1]} \sigma'(z_j^{[l]}).$$

Substituting back gives

$$\delta_j^{[l]} = \sigma'(z_j^{[l]}) \sum_{m=1}^{n_{l+1}} w_{mj}^{[l+1]} \delta_m^{[l+1]},$$

which is precisely

$$\delta^{[l]} = \sigma'(z^{[l]}) \circ (W^{[l+1]})^T \delta^{[l+1]}.$$

Proof of (3.8). Since

$$z_j^{[l]} = \sum_k w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]},$$

we have $\frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = 1$ and zero otherwise. Hence,

$$\frac{\partial C}{\partial b_j^{[l]}} = \sum_{m=1}^{n_l} \frac{\partial C}{\partial z_m^{[l]}} \frac{\partial z_m^{[l]}}{\partial b_j^{[l]}} = \delta_j^{[l]}.$$

Proof of (3.9). Similarly,

$$\frac{\partial z_j^{[l]}}{\partial w_{jk}^{[l]}} = a_k^{[l-1]},$$

so

$$\frac{\partial C}{\partial w_{jk}^{[l]}} = \sum_{m=1}^{n_l} \frac{\partial C}{\partial z_m^{[l]}} \frac{\partial z_m^{[l]}}{\partial w_{jk}^{[l]}} = \delta_j^{[l]} a_k^{[l-1]}.$$

□

Algorithm 3.2 Gradient Descent on Neural Networks: Backpropagation

- 1: **Input:** Training data (x, y)
 - 2: **Forward pass:** $x = a^{[1]} \rightarrow a^{[2]} \rightarrow \dots \rightarrow a^{[L-1]} \rightarrow a^{[L]}$
 - 3: Obtain $\delta^{[L]}$ from (3.6)
 - 4: **Backward pass:** Form (3.7) obtain $\delta^{[L]} \rightarrow \delta^{[L-1]} \rightarrow \delta^{[L-2]} \rightarrow \dots \rightarrow \delta^{[2]}$
 - 5: **Calculation of the gradient of the loss C :** Use (3.8) and (3.9)
-

Training GD on Over-parametrized Neural Networks has been understood only recently. See [3, 8, 9, 19].

Computational Graph and General Backpropagation for taking derivatives: basic idea in tensorflow (see [17]).

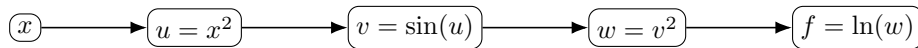
Example. Consider

$$f(x) = \ln(\sin^2(x^2)),$$

defined for $\sin(x^2) \neq 0$. Introduce intermediate variables

$$u = x^2, \quad v = \sin(u), \quad w = v^2, \quad f = \ln(w). \quad (3.10)$$

Computational graph:



$$\frac{\partial u}{\partial x} = 2x$$

$$\frac{\partial v}{\partial u} = \cos(u)$$

$$\frac{\partial w}{\partial v} = 2v$$

$$\frac{\partial f}{\partial w} = \frac{1}{w}$$

Reverse-mode (backprop) computation:

Define adjoints $\bar{u} = \partial f / \partial u$, $\bar{v} = \partial f / \partial v$, $\bar{w} = \partial f / \partial w$, and $\bar{x} = \partial f / \partial x$. Start from the output:

$$\bar{f} = 1, \quad \bar{w} = \bar{f} \frac{\partial f}{\partial w} = \frac{1}{w}.$$

Propagate backward along the edges using the chain rule:

$$\begin{aligned} \bar{v} &= \bar{w} \frac{\partial w}{\partial v} = \frac{1}{w} (2v) = \frac{2}{v}, \\ \bar{u} &= \bar{v} \frac{\partial v}{\partial u} = \frac{2}{v} \cos(u) = 2 \frac{\cos(u)}{\sin(u)} = 2 \cot(u), \\ \bar{x} &= \bar{u} \frac{\partial u}{\partial x} = (2 \cot(u)) (2x) = 4x \cot(u). \end{aligned}$$

Substituting $u = x^2$ gives

$$f'(x) = \frac{df}{dx} = 4x \cot(x^2), \quad \sin(x^2) \neq 0. \quad (3.11)$$

This illustrates the key idea of automatic differentiation: evaluate the forward pass (3.10) once, then compute derivatives by a single backward pass using local Jacobians at each node.

In PyTorch, all computational flows are stored during calculation, so that any gradient can be calculated using back-propagation. This takes memory. So `torch.no_grad()` (inference mode) can be used to accelerate the inference process if no gradient computation is needed.

3.3 Online Principle Component Analysis (PCA)

Let \mathbf{X} be a random vector in \mathbb{R}^d with mean 0 and unknown covariance matrix

$$\Sigma = \mathbf{E} [\mathbf{X} \mathbf{X}^T] \in \mathbb{R}^{d \times d}.$$

Let the eigenvalues of Σ be ordered as $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_d \geq 0$. Principal Component Analysis (PCA) aims to find the principal eigenvector of Σ that corresponds to the largest eigenvalue λ_1 , based on independent and identically distributed sample realizations $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)}$. This can be casted into a *nonconvex stochastic optimization problem* given by

$$\begin{aligned} &\text{maximize } \mathbf{u}^T \mathbf{E} [\mathbf{X} \mathbf{X}^T] \mathbf{u}, \\ &\text{subject to } \|\mathbf{u}\| = 1, \mathbf{u} \in \mathbb{R}^d. \end{aligned}$$

Here $\|\bullet\|$ denotes the Euclidian norm. Assume the principle component \mathbf{u}^* is unique.

Classical PCA:

$$\hat{\mathbf{u}}^{(n)} = \arg \max_{\|\mathbf{u}\|=1} \mathbf{u}^T \hat{\Sigma}^{(n)} \mathbf{u}.$$

Here $\hat{\Sigma}^{(n)}$ is the empirical covariance matrix based on n samples:

$$\hat{\Sigma}^{(n)} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}^{(i)} (\mathbf{X}^{(i)})^T.$$

We can design an online version of the PCA algorithm as in Algorithm 3.3, where we must do at every iteration a normalization step using the retraction operation

$$\Pi \mathbf{u} = \frac{\mathbf{u}}{\|\mathbf{u}\|} .$$

Algorithm 3.3 Online PCA Algorithm

- 1: **Input:** Initialize $\mathbf{u}^{(0)}$ and choose the stepsize $\beta > 0$
- 2: **for** $n = 1, 2, \dots, N$ **do**
- 3: Draw one sample $\mathbf{X}^{(n)}$ from the (streaming) data source
- 4: Update the iterate $\mathbf{u}^{(n)}$ by

$$\mathbf{u}^{(n)} = \Pi \left\{ \mathbf{u}^{(n-1)} + \beta \mathbf{X}^{(n)} (\mathbf{X}^{(n)})^T \mathbf{u}^{(n-1)} \right\} .$$

- 5: **end for**
 - 6: **Output:** $\mathbf{u}^{(N)}$
-

The convergece of online PCA has been proved only recently in [12]. The difficulty here is that this online PCA is non-convex optimization.

Oja's algorithm: Step 4 in Algorithm 3.3 will be modified to

$$\mathbf{u}^{(n)} = \mathbf{u}^{(n-1)} + \beta \left(\mathbf{X}^{(n)} (\mathbf{X}^{(n)})^T \mathbf{u}^{(n-1)} - \left((\mathbf{u}^{(n-1)})^T \mathbf{X}^{(n)} \right)^2 \mathbf{u}^{(n-1)} \right) ,$$

which is expected to control the norm growth. See [14].

3.4 Line Search Methods

Line search methods in general takes the following form

$$x^{k+1} = x^k + \alpha_k d^k , \quad k = 0, 1, 2, \dots \quad (3.12)$$

Here $\alpha_k > 0$ is the stepsize and $d^k \in \mathbb{R}^n$ is the descent direction.

Preliminary conditions: for $\bar{\varepsilon}, \gamma_1, \gamma_2 > 0$ we assume

$$\begin{aligned} (1) \quad & -(d^k)^T \nabla f(x^k) \geq \bar{\varepsilon} \|\nabla f(x^k)\| \cdot \|d^k\|; \\ (2) \quad & \gamma_1 \|\nabla f(x^k)\| \leq \|d^k\| \leq \gamma_2 \|\nabla f(x^k)\| . \end{aligned} \quad (3.13)$$

Notice that if $d^k = -\nabla f(x^k)$, then $\bar{\varepsilon} = \gamma_1 = \gamma_2 = 1$.

In fact one can expand

$$\begin{aligned} f(x^{k+1}) &= f(x^k + \alpha d^k) \\ &= f(x^k) + \alpha \nabla f(x^k)^T d^k + \alpha \int_0^1 [\nabla f(x^k + \gamma \alpha d^k) - \nabla f(x^k)] d^k d\gamma \\ &\leq f(x^k) - \alpha \left(\bar{\varepsilon} - \alpha \frac{L}{2} \gamma_2 \right) \|\nabla f(x^k)\| \cdot \|d^k\| . \end{aligned} \quad (3.14)$$

So if $\alpha \in \left(0, \frac{2\bar{\varepsilon}}{L\gamma_2}\right)$, then $f(x^{k+1}) < f(x^k)$, leading to a descent in objective value.

Schemes of descent type produce iterates that satisfy a bound of the form

$$f(x^{k+1}) \leq f(x^k) - C\|\nabla f(x^k)\|^2 \text{ for some } C > 0. \quad (3.15)$$

Various Line Search Methods:

- (1) $d^k = -S^k \nabla f(x^k)$, S^k symmetric positively definite, $\lambda(S^k) \in [\gamma_1, \gamma_2]$;

In this case we have $-(d^k)^T \nabla f(x_k) = (\nabla f(x_k))^T S^k (\nabla f(x_k)) \geq \gamma_2 \|\nabla f(x^k)\|^2 \geq \gamma_2 \gamma_2^{-1} \|\nabla f(x^k)\| \cdot \|d^k\| = \|\nabla f(x^k)\| \cdot \|d^k\|$.

- (2) Gauss-Southwell: $d^k = -[\nabla f(x^k)]_{i_k}$ where

$$i_k = \arg \min_{i=1,2,\dots,n} |[\nabla f(x^k)]_i|$$

and

$$\nabla f(x) = ([\nabla f(x)]_1, \dots, [\nabla f(x)]_n) ;$$

- (3) Stochastic coordinate descent: $d^k = -[\nabla f(x^k)]_{i_k}$ where i_k is uniformly randomly drawn from $\{1, 2, \dots, n\}$;

- (4) Stochastic Gradient Methods: d^k is random, we will specify it later in Chapter 5;

- (5) Exact line search $\min_{\alpha > 0} f(x^k + \alpha d^k)$;

- (6) Approximate Line Search: choose a certain $\alpha > 0$, not necessarily the exact α in the above. This usually is accompanied with the “weak Wolfe conditions”, which consists of
 (a) Armijo (Sufficient Decrease) Condition: $f(x^k + \alpha p^k) \leq f(x^k) + c_1 \alpha_k (p^k)^T \nabla f(x^k)$, where $c_1 \in (0, 1)$ and is often chosen to be of a small order of magnitude around 10^{-3} . This condition ensures the computed step length can sufficiently decrease the objective function $f(x^k)$. Only using this condition, however, it cannot be guaranteed that x^k will converge in a reasonable number of iterations, since Armijo condition is always satisfied with step length being small enough. Therefore, the second condition below needs to be paired with the sufficient decrease condition to keep α_k from being too short.
 (b) Curvature condition: $(\nabla f(x^k + \alpha p^k))^T p^k \geq c_2 (\nabla f(x^k))^T p^k$ where $c_2 \in (c_1, 1)$ is much greater than c_1 and is typically on the order of 0.1. The condition ensures a sufficient increase of the gradient;

- (7) Backtracking Line Search: Starting from some $\bar{\alpha} > 0$ then look for $\alpha \in \{\bar{\alpha}, \beta \bar{\alpha}, \beta^2 \bar{\alpha}, \dots\}$ where $0 < \beta < 1$.

3.5 Convergence to Approximate Second Order Necessary Points

Second-order necessary point: $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.

Definition 3.5. We call point $x \in \mathbb{R}^n$ an $(\varepsilon_g, \varepsilon_H)$ -approximate second order stationary point if

$$\|\nabla f(x)\| \leq \varepsilon_g \text{ and } \lambda_{\min}(\nabla^2 f(x)) \geq -\varepsilon_H \quad (3.16)$$

for some choices of small constants $\varepsilon_g > 0$ and $\varepsilon_H > 0$.

How to search for such a second-order stationary point?

Set $M > 0$ to be the Lipschitz bound for the Hessian $\nabla^2 f(x)$, so that

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq M\|x - y\| \text{ for all } x, y \in \text{dom}(f) . \quad (3.17)$$

Algorithm 3.4 Search for second-order approximate stationary point based on Gradient Descent

```

1: Input: Input initialization  $x^0$ ,  $k = 0, 1, 2, \dots$ 
2: while  $\|\nabla f(x^k)\| > \varepsilon_g$  or  $\lambda_k < -\varepsilon_H$  do
3:   if  $\|\nabla f(x^k)\| > \varepsilon_g$  then
4:      $x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k)$ 
5:   else
6:      $\lambda^k := \lambda_{\min}(\nabla^2 f(x^k))$ 
7:     Choose  $p^k$  to be the eigenvector corresponding to the most negative eigenvalue of  $\nabla^2 f(x^k)$ 

8:     Choose the size and sign of  $p^k$  such that  $\|p^k\| = 1$  and  $(p^k)^T \nabla f(x^k) \leq 0$ 
9:      $x^{k+1} = x^k + \alpha_k p^k$ ,  $\alpha_k = \frac{2\lambda_k}{M}$ 
10:   end if
11:    $k \leftarrow k + 1$ 
12: end while
13: Output: An estimate of an  $(\varepsilon_g, \varepsilon_H)$ -approximate second order stationary point  $x^k$ 

```

Cubic upper bound

$$f(x + p) \leq f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + \frac{1}{6} M \|p\|^3 . \quad (3.18)$$

For steepest descent step we have from (3.4) that

$$f(x^{k+1}) \leq f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2 \leq f(x^k) - \frac{\varepsilon_g^2}{2L} .$$

Otherwise, the negative direction Hessian eigenvector search will output, by using third-order Taylor expansion and (3.17), that

$$\begin{aligned}
f(x^{k+1}) &\leq f(x^k) + \alpha_k \nabla f(x^k)^T p^k + \frac{1}{2} \alpha_k^2 (p^k)^T \nabla^2 f(x^k) p^k + \frac{1}{6} M \alpha_k^3 \|p^k\|^3 \\
&\leq f(x^k) - \frac{1}{2} \left(\frac{2|\lambda_k|}{M} \right)^2 |\lambda_k| + \frac{1}{6} M \left(\frac{2|\lambda_k|}{M} \right)^3 \\
&= f(x^k) - \frac{2}{3} \frac{|\lambda_k|^3}{M^2} \\
&= f(x^k) - \frac{2}{3} \frac{\varepsilon_H^3}{M^2} .
\end{aligned}$$

Number of iterations

$$K \leq \max \left(2L\varepsilon_g^{-2}, \frac{3}{2}M^2\varepsilon_H^{-3} \right) (f(x^0) - f(x^*)) .$$

3.6 Experiment: Learning a Neural Network via Gradient Descent and Backpropagation

Consider an L -hidden-layer fully connected neural network with hidden layer sizes n_1, \dots, n_L , where the input layer has only one neuron $x \in \mathbb{R}$ and output layer has only one neuron $y \in \mathbb{R}$. Let the training data $(x_1, y_1), \dots, (x_N, y_N)$ be an i.i.d sample of a bivariate unit normal distribution $(x, y) \in \mathcal{N}(0, I_2)$. We use the backpropagation algorithm to calculate the gradient vector of the empirical loss function and the population loss function of the neural network for a quadratic loss.

We experiment

- (1) Different hidden layer sizes: (n_1, \dots, n_L) ;
- (2) Different number of hidden layers: L ;
- (3) Different sizes of the training data set: n ;
- (4) Different activation functions: Sigmoid, ReLU, Tanh.

Following the above, we plot the 3-d graph of the loss function and the trajectory of the gradient descent algorithm by picking two network weight parameters as variables (from a certain layer to its next layer) and randomly choosing other network parameters as i.i.d. samples of normal distributions that follow the LeCun initialization: $w_{ij}^{(l)} \sim \mathcal{N}(0, \frac{1}{n_l})$ and $b_i^{(l)} \sim \mathcal{N}(0, 1)$, where n_l is the network width at level l .

Code can be found at https://github.com/huwenqing0606/Computational_Optimization/tree/main/src/backpropagation

Chapter 4

Accelerated Gradient Methods

4.1 Polyak’s heavy ball method

Why acceleration in GD?

The GD method is “greedy” in the sense that it steps in the direction that is most productive at current iterate – no explicit use of knowledge of the objective function f at earlier iterations.

Adding “momentum” = search direction tends to be similar to that one used in the previous step.

Say the objective function $f(x) = \frac{1}{2}x^T Qx - b^T x + c$ is a quadratic function. It may well happen in deep neural network models, near equilibrium, the loss function is approximated by such a function f with a matrix Q that has a large condition number $\kappa = \frac{L}{m}$, see [6]. Slow convergence for κ large: for the strongly convex case we showed that in GD case the convergence rate is $\beta = 1 - \frac{1}{\kappa}$ which $\rightarrow 1$ when $\kappa \rightarrow \infty$. Fig 4.1 demonstrates oscillations of GD near local min when the Hessian has a large condition number.

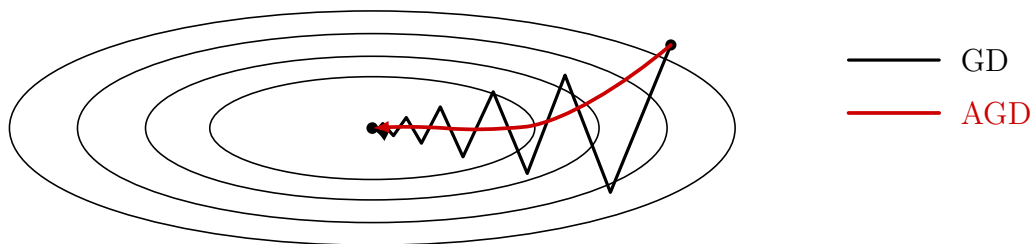


Figure 4.1: Schematic comparison of GD (zig-zag) and Accelerated GD (smoother) on an ill-conditioned valley.

Motivated from classical mechanics, GD is approximated by a first order ODE of the form

$$\text{GD:} \quad \frac{dx}{dt} = -\nabla f(x) .$$

It is natural to modify the above equation to a second order Newton’s equation that models a nonlinear oscillator:

$$\text{AGD:} \quad \mu \frac{d^2 x}{dt^2} = -\nabla f(x) - \mu b \frac{dx}{dt} .$$

“Small mass limit” $\mu \rightarrow 0$. Term $\mu \frac{d^2 x}{dt^2}$ serves as “regularization”.

Finite difference approximation:

$$\mu \frac{x(t + \Delta t) - 2x(t) + x(t - \Delta t)}{(\Delta t)^2} \approx -\nabla f(x(t)) - \frac{\mu b(x(t + \Delta t) - x(t))}{\Delta t} .$$

i.e.

$$x(t + \Delta t) = x(t) - \alpha \nabla f(x(t)) + \beta(x(t) - x(t - \Delta t)) .$$

This gives the “Polyak’s heavy ball method” (see [15]).

$$x^{k+1} = x^k - \alpha \nabla f(x^k) + \beta(x^k - x^{k-1}) , \quad x^{-1} = x^0 .$$

Algorithm 4.1 Polyak’s heavy ball method (GD with momentum)

- 1: **Input:** Input initialization $x^{-1} = x^0$, $p^{-1} = 0$, stepsizes $\alpha, \beta > 0$, $k = 0, 1, 2, \dots$
 - 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 - 3: $p^k = -\alpha \nabla f(x^k) + \beta p^{k-1}$
 - 4: $x^{k+1} = x^k + p^k$
 - 5: **end for**
 - 6: **Output:** Output x^K
-

“Momentum”: $p_k = x^{k+1} - x^k$. The momentum iteration

$$p^k = \beta p^{k-1} - \alpha \nabla f(x^k)$$

can be understood as doing an exponential moving average over the gradients. This is where the “memory” of “inertia” comes in, that inspires the adaptive methods.

4.2 Nesterov’s Accelerated Gradient Descent (AGD)

“first order optimal method”: Instead of evaluating $\nabla f(x^k)$, we evaluate $\nabla f(x^k + \beta(x^k - x^{k-1}))$.

Nesterov’s accelerated gradient descent:

$$x^{k+1} = x^k - \alpha \nabla f(x^k + \beta(x^k - x^{k-1})) + \beta(x^k - x^{k-1}) , \quad x^{-1} = x^0 .$$

See [13].

Algorithm 4.2 Nesterov's accelerated gradient descent

-
- 1: **Input:** Input initialization $x^{-1} = x^0$, stepsizes $\alpha_k, \beta_k > 0$, $k = 0, 1, 2, \dots$
 - 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 - 3: $y^k = x^k + \beta_k(x^k - x^{k-1})$
 - 4: $x^{k+1} = y^k - \alpha_k \nabla f(y^k)$
 - 5: **end for**
 - 6: **Output:** Output x^K
-

By introducing the momentum variable $p^k = x^{k+1} - x^k$ we can also rewrite the Nesterov's accelerated gradient descent in terms of momentum, this is the Nesterov's momentum algorithm:

Algorithm 4.3 Nesterov's momentum (equivalent to Nesterov's accelerated gradient descent)

-
- 1: **Input:** Input initialization $x^{-1} = x^0$, $p^{-1} = 0$, stepsizes $\alpha_k, \beta_k > 0$, $k = 0, 1, 2, \dots$
 - 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 - 3: $p^k = \beta_k p^{k-1} - \alpha_k \nabla f(x^k + \beta_k p^{k-1})$
 - 4: $x^{k+1} = x^k + p^k$
 - 5: **end for**
 - 6: **Output:** Output x^K
-

Convergence proof of Nesterov's method:

(1) Convex quadratic objective function: Spectral method.

$$f(x) = \frac{1}{2} x^T Q x - b^T x + c, \quad Q \succ 0, \quad 0 < m = \lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_2 \leq \lambda_1 = L.$$

Condition number $\kappa = \frac{L}{m}$.

Minimizer of f is $x^* = Q^{-1}b$ and $\nabla f(x) = Qx - b = Q(x - x^*)$.

"Spectral method"

$$x^{k+1} - x^* = (x^k - x^*) - \alpha Q (x^k + \beta(x^k - x^{k-1}) - x^*) + \beta ((x^k - x^*) - (x^{k-1} - x^*)).$$

So we have

$$\begin{bmatrix} x^{k+1} - x^* \\ x^k - x^* \end{bmatrix} = \begin{bmatrix} (1 + \beta)(I - \alpha Q) & -\beta(I - \alpha Q) \\ I & 0 \end{bmatrix} \begin{bmatrix} x^k - x^* \\ x^{k-1} - x^* \end{bmatrix}, \quad k = 1, 2, \dots$$

Set $\omega^k = \begin{pmatrix} x^{k+1} - x^* \\ x^k - x^* \end{pmatrix}$ with $x^{-1} = x^0$, and set $T = \begin{bmatrix} (1 + \beta)(I - \alpha Q) & -\beta(I - \alpha Q) \\ I & 0 \end{bmatrix}$, then $\omega^k = T\omega^{k-1}$. So $\omega^k = T^k\omega^0$.

Theorem 4.1. Pick $\alpha = \frac{1}{L}$, $\beta = \frac{\sqrt{L} - \sqrt{m}}{\sqrt{L} + \sqrt{m}} = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$, then the eigenvalues of T are

$$\nu_{i,1} = \frac{1}{2} \left[(1 + \beta)(1 - \alpha\lambda_i) + i\sqrt{4\beta(1 - \alpha\lambda_i) - (1 + \beta)^2(1 - \alpha\lambda_i)^2} \right],$$

$$\nu_{i,2} = \frac{1}{2} \left[(1 + \beta)(1 - \alpha\lambda_i) - i\sqrt{4\beta(1 - \alpha\lambda_i) - (1 + \beta)^2(1 - \alpha\lambda_i)^2} \right] .$$

In particular, the spectral radius $\rho(T) \equiv \lim_{k \rightarrow \infty} (\|T^k\|)^{1/k} \leq 1 - \frac{1}{\sqrt{\kappa}}$.

By Theorem 4.1, we conclude for any $\varepsilon > 0$ that

$$\|\omega^k\| \leq C\|\omega^0\|(\rho(T) + \varepsilon)^k .$$

Comparison of the rates with GD: when $\alpha = \frac{1}{L}$, GD needs $O\left(\frac{L}{m} |\ln \varepsilon|\right)$ to reach $f(x^k) - f^* \leq \varepsilon$

while Nesterov needs $O\left(\sqrt{\frac{L}{m}} |\ln \varepsilon|\right)$ to reach a factor ε in $\|\omega^k\|$. Key parameter is $\kappa = \frac{L}{m}$.

(2) Strongly convex objective functions: Lyapunov's method.

Lyapunov function $V : \mathbb{R}^d \rightarrow \mathbb{R}$, $V(\omega) > 0$ for all $\omega \neq \omega^* \in \mathbb{R}^d$, $V(\omega^*) = 0$, $V(\omega^{k+1}) < \rho^2 V(\omega^k)$ for $0 < \rho < 1$.

Indeed, $f(x) - f^*$ is a Lyapunov function for GD.

What is the Lyapunov function for Nesterov?

$\tilde{x}^k = x^k - x^*$, $\tilde{y}^k = y^k - y^*$, set

$$V_k = f(x^k) - f(x^*) + \frac{L}{2} \|\tilde{x}^k - \rho^2 \tilde{x}^{k-1}\|^2 .$$

Calculations... Possible to show $\alpha_k = \frac{1}{L}$, $\beta_k = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$, $\rho^2 = 1 - \frac{1}{\sqrt{\kappa}}$ such that $V_{k+1} \leq \rho^2 V_k$.

This is to say we have

$$f(x^k) - f^* \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^k \left\{ f(x_0) - f^* + \frac{m}{2} \|x_0 - x^*\|^2 \right\} ,$$

where one can verify that $V_0 = f(x_0) - f^* + \frac{m}{2} \|x_0 - x^*\|^2$.

(3) Convex objective functions: Lyapunov's method.

Set $\alpha_k = \frac{1}{L}$ but β_k is variable, and thus ρ_k is variable.

$$V_k = f(x^k) - f^* + \frac{L}{2} \|\tilde{x}^k - \rho_{k-1}^2 \tilde{x}^{k-1}\|^2 .$$

(Weakly) Convex objective: $m = 0$.

Obtain

$$V_{k+1} \leq \rho_k^2 V_k + R_k^{(\text{weak})} ,$$

where

$$R_k^{(\text{weak})} = \frac{L}{2} \|\tilde{y}^k - \rho_k^2 \tilde{x}^k\|^2 - \frac{\rho_k^2 L}{2} \|\tilde{x}^k - \rho_{k-1}^2 \tilde{x}^{k-1}\|^2 .$$

Recursive relation: $\rho_k^2 = \frac{(1 - \rho_k^2)^2}{(1 - \rho_{k-1}^2)^2}$, $k = 1, 2, \dots$. We have $1 + \beta_k - \rho_k^2 = \rho_k$, $\beta_k = \rho_k \rho_{k-1}^2$.

This makes $R_k^{(\text{weak})} = 0$.

So $V_k \leq \rho_{k-1}^2 V_{k-1}$ and thus

$$V_k \leq \rho_{k-1}^2 \rho_{k-2}^2 \dots \rho_1^2 V_1 = \frac{(1 - \rho_{k-1}^2)^2}{(1 - \rho_0^2)^2} V_1 .$$

Bound $V_1 \leq (1 - \rho_{k-1}^2)^2 \frac{L}{2} \|x^0 - x^*\|^2$ and inductive shows that $1 - \rho_k^2 \leq \frac{2}{k+2}$.

Finally we obtain the rate

$$f(x^k) - f^* \leq V_k \leq \frac{2L}{(k+1)^2} \|x^0 - x^*\|^2 .$$

See

Algorithm 4.4 Nesterov's optimal algorithm for general convex function f

- 1: **Input:** Input initialization $x^{-1} = x^0$, $\beta_0 = \rho_0 = 0$ and Lipschitz constant L for f
- 2: **for** $k = 0, 1, 2, \dots, K-1$ **do**
- 3: $y^k = x^k + \beta_k(x^k - x^{k-1})$
- 4: $x^{k+1} = y^k - \frac{1}{L} \nabla f(y^k)$
- 5: Define $\rho_{k+1} \in [0, 1]$ to be the root of the quadratic equation

$$1 + \rho_{k+1}(\rho_k^2 - 1) - \rho_{k+1}^2 = 0$$

- 6: Set $\beta_{k+1} = \rho_{k+1} \rho_k^2$
 - 7: **end for**
 - 8: **Output:** Output x^K
-

Connection with differential equations: Su, W., Boyd, S., Candés, E.J., A differential equation for modeling Nesterov's accelerated gradient method: theory and insights. *Journal of Machine Learning Research*, **17**:1–43, 2016. Suggest choosing $\beta_k = \frac{k-1}{k+2}$?

Su-Boyd-Candés differential equation:

$$\ddot{X}_t + \frac{3}{t} \dot{X}_t + \nabla f(X_t) = 0 .$$

Nesterov's method achieves first order optimal convergence rate: Example. See Nesterov, Y., *Introductory Lectures on Convex Optimization*, Springer, 2004, §2.1.2.

Indeed, one can construct a carefully designed function for which *no* method that makes use of all gradient observed up to and including iteration k (namely $\nabla f(x^i)$, $i = 0, 1, 2, \dots, k$) can produce a sequence $\{x^k\}$ that achieves a rate better than the order $O(1/k^2)$.

Set $f(x) = \frac{1}{2} x^T A x - e_1^T x$, where

$$A = \begin{pmatrix} 2 & -1 & 0 & 0 & \dots & \dots & 0 \\ -1 & 2 & -1 & 0 & \dots & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 2 \end{pmatrix} , \quad e_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \dots \\ 0 \end{pmatrix} .$$

Optimization $\min_{x \in \mathbb{R}^n} f(x)$ has a solution x^* such that $Ax^* = e_1$, and its components are $x_i^* = 1 - \frac{i}{n+1}$ for $i = 1, 2, \dots, n$. Let the starting point $x^0 = 0$, and the iterate

$$x^{k+1} = x^k + \sum_{j=0}^k \gamma_j \nabla f(x^j)$$

for some coefficients γ_j , $j = 0, 1, \dots, k$. Then each iterate x^k can have nonzero entries only in its first k components. Thus

$$\|x^k - x^*\|^2 \geq \sum_{j=k+1}^n (x_j^*)^2 = \sum_{j=k+1}^n \left(1 - \frac{j}{n+1}\right)^2.$$

One can show that for $k = 1, 2, \dots, \frac{n}{2} - 1$ we have

$$\|x^k - x^*\|^2 \geq \frac{1}{8} \|x^0 - x^*\|^2.$$

This will lead to the bound that

$$f(x^k) - f(x^*) \geq \frac{3L}{32(k+1)^2} \|x^0 - x^*\|^2, \quad k = 1, 2, \dots, \frac{n}{2} - 1$$

where $L = \|A\|_2$.

4.3 Conjugate gradient method

Conjugate gradient method is a more classical method in computational mathematics. Unlike the Nesterov's accelerated gradient method, the conjugate gradient method does not require knowledge of the convexity parameters L and m to compute the appropriate step-sizes. Unfortunately, it does not admit particularly rigorous analysis when the objective function f is not quadratic.

Originally this method comes from numerical linear algebra: goal is to solve $Qx = p$. But this can be formulated as optimization problem $\min_x f(x)$ with $f(x) = \frac{1}{2}x^T Qx - p^T x + c$.

Algorithm 4.5 Conjugate Gradient Method

- 1: **Input:** Input initialization x^0 , $y^0 = -\nabla f(x^0)$, $f(x) = \frac{1}{2}x^T Qx - p^T x + c$, $k = 0, 1, 2, \dots$
 - 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 - 3: $r^k = Qx^k - p$
 - 4: $\alpha_k = \frac{\langle y^k, r^k \rangle}{\langle y^k, Qy^k \rangle}$
 - 5: $x^{k+1} = x^k - \alpha_k y^k$
 - 6: $y^{k+1} = -\nabla f(x^{k+1}) + \beta_k y^k$
 - 7: $\beta_{k+1} = \frac{\langle r^k, Qy^k \rangle}{\langle y^k, Qy^k \rangle}$
 - 8: **end for**
 - 9: **Output:** Output x^K
-

Exact line search in the direction y^k : $\alpha_k = \arg \min_{\alpha \geq 0} f(x^k + \alpha y^k)$.

Conjugate gradient direction: β_{k+1} is so chosen that $\langle y^{k+1}, Qy^k \rangle = 0$.

4.4 Experiment: GD vs Accelerated GD

Bibliography

- [1] [https://en.wikipedia.org/wiki/Matrix_factorization_\(recommender_systems\)](https://en.wikipedia.org/wiki/Matrix_factorization_(recommender_systems)).
- [2] https://github.com/huwenqing0606/Computational_Optimization.
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.
- [4] Jacot Arthur, Gabriel Franck, and Hongler Clement. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in Neural Information Processing Systems, arXiv:1806.07572*, 31:8571–8580, 2018.
- [5] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. Sgd learns over-parameterized networks that provably generalize on linearly separable data. *ICLR, arXiv:1710.10174*, 12:1–12, 2018.
- [6] P. Chaudhari and S. Soatto. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. *ICLR*, 2018.
- [7] Y. Chi, Y.M. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20), 2019.
- [8] Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.
- [9] Simon S. Du, Xiyu Zhai, Barnabás Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- [10] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [11] Catherine F. Higham and Desmond J. Higham. *Deep Learning: An Introduction for Applied Mathematicians*. arXiv:1801.05894[math.HO], 2018.
- [12] C.J Li, M. Wang, H. Liu, and T. Zhang. Near-optimal stochastic approximation for online principal component estimation. *Mathematical Programming*, 167(1):75–97, 2018.

- [13] Y. Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. *Soviet Math Dokl*, 27(2), 1983.
- [14] E. Oja. Simplified neuron model as a principal component analyzer. *Journal of mathematical biology*, 15(3):267–273, 1982.
- [15] B. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- [16] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986.
- [17] TensorFlow. Automatic differentiation and gradient tape. <https://www.tensorflow.org/tutorials/customization/autodiff>. Accessed: 2026-02-17.
- [18] Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 2000.
- [19] Xiaojin Zhang, Yuchen Yu, Lin Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. *arXiv preprint arXiv:1806.07808*, 2018.