

# Computational Optimization

## with Applications to Machine Learning

A Series of Lectures at Missouri S&T

Wenqing Hu <sup>1</sup>

<sup>1</sup>Department 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](mailto:huwen@mst.edu)

# Contents

<b>1</b>	<b>Motivating Examples</b>	<b>3</b>
1.1	Supervised Learning . . . . .	3
1.2	Matrix Optimizations . . . . .	7
1.3	Experiment: Activation Functions . . . . .	8
<b>2</b>	<b>Convex Functions</b>	<b>9</b>
2.1	Optimization problem: set-up and its solutions . . . . .	9
2.2	Convexity . . . . .	9
2.3	Taylor's theorem and convexity in Taylor's expansion . . . . .	11
2.4	Optimality Conditions . . . . .	17
2.5	Experiment: Nonconvexity of the Loss Landscape of Neural Networks .	19
<b>3</b>	<b>Gradient Descent and Line Search Methods</b>	<b>21</b>
3.1	Gradient Descent . . . . .	21
3.2	Back Propagation and GD on Neural Networks . . . . .	23
3.3	Online Principle Component Analysis (PCA) . . . . .	25
3.4	Line Search Methods . . . . .	26
3.5	Convergence to Approximate Second Order Necessary Points . . . . .	27

# 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  $\omega$  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  $\omega$ ).

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  $R(\omega)$ :

$L^1$ -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 [5].

(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  $\omega$  and  $\beta$ , we can then formulate the problem as looking for optimal  $\omega$  and  $\beta$  such that  $\omega^T x_j + \beta \geq 1$  when  $y_j = 1$  and  $\omega^T x_j + \beta \leq -1$  when  $y_i = -1$ . 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 [7].

(3) Neural Network: “activation function”  $\sigma$ .

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 [3].

## 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 [4] 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  $\geq 0$ ). 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  $\lambda_1$  (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.

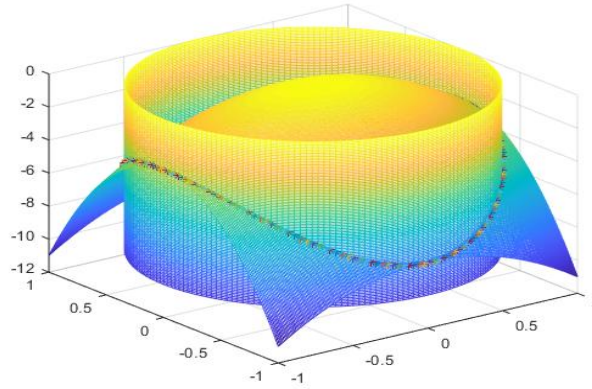


Figure 1.1: Loss Landscape of 1-dimensional PCA.

Online PCA: see [6].

(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  $\Omega$  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  $\Omega$  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  $\phi$  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  $\phi$ .

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

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

- $\phi$  is a “proper convex function” if  $\phi(x) < \infty$  for some  $x \in \Omega$  and  $\phi(x) > -\infty$  for all  $x \in \Omega$ ;
- $\phi$  is a “closed proper convex function” if  $\phi$  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  $\nabla f$  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  $\nabla f$  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).  $\square$

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)  *$\nabla f$  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  $(1 - \alpha)f(x) + \alpha f(y)$  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  $\nabla f$  is  $L$ -Lipschitz continuous. Assume  $\nabla f$  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  $\nabla f$  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 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  $\alpha$  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  $\alpha$ , 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

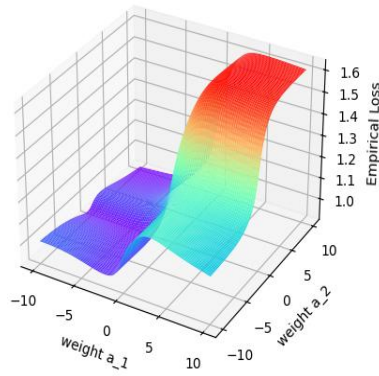


Figure 2.1: Loss Landscape of a neural network.

## 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

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

**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:*

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

or

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

or

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

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  $\nabla f$  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.* Use fixed point iteration. Let  $\Phi(x) = x - \alpha \nabla f(x)$ ,  $\alpha > 0$ . Then  $x^{k+1} = \Phi(x^k)$ . Calculate  $\|\Phi(x) - \Phi(z)\| \leq \beta \|x - z\|$ .

$$\text{In order } \|x^T - x^*\| \leq \varepsilon \text{ we want } T \geq \frac{\ln \left( \frac{\|x^0 - x^*\|}{\varepsilon} \right)}{|\ln \beta|}.$$

Want  $-\beta \leq 1 - \alpha L \leq 1 - \alpha m \leq \beta$ , and when  $1 - \alpha L = 1 - \alpha m$ , we have  $\beta = \frac{L - m}{L + m}$  and  $\alpha = \frac{2}{L + m}$ . □

Standard procedure in optimization theory: make use of Taylor's expansion.

Consider

$$\begin{aligned}
 f(x + \alpha d) &= f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 [\nabla f(x + \gamma \alpha d) - \nabla f(x)]^T d \, d\gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 \|\nabla f(x + \gamma \alpha d) - \nabla f(x)\| \|d\| \, d\gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha^2 \frac{L}{2} \|d\|^2.
 \end{aligned} \tag{3.1}$$

From (3.1) and the line search  $d = -\nabla f(x^k)$ ,  $x = x^k$  on RHS of (3.1), then the  $\alpha = \frac{1}{L}$  minimizes RHS of (3.1), and thus for GD in particular

$$f(x^{k+1}) \leq f(x^k) + \left(-\alpha + \alpha^2 \frac{L}{2}\right) \|\nabla f(x^k)\|^2 = f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2. \tag{3.2}$$

Estimates of type (3.2) 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

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

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

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

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

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

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

## 3.2 Back Propagation and GD on Neural Networks

Classical algorithm of training multi-layer neural networks: Back-propagation.

Reference: Catherine F. Higham, Desmond J. Higham, Deep Learning: An Introduction for Applied Mathematicians, arXiv:1801.05894[math.HO]

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

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

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  $\omega$  and  $b$ , that is  $\frac{\partial C}{\partial w_{jk}^{[l]}}$  and  $\frac{\partial C}{\partial b_j^{[l]}}$ . Set the *error* in the  $j$ -th neuron at layer  $l$  to be  $\delta^{[l]} = (\delta_j^{[l]})_j$  (viewed as a column vector) with  $\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}}$ . Then one can calculate directly to obtain the following lemma

**Lemma 3.2.** *If  $x, y \in \mathbb{R}^n$ , let  $x \circ y \in \mathbb{R}^n$  to be the Hadamard product defined by  $(x \circ y)_i = x_i y_i$ . We can directly calculate*

$$\delta^{[L]} = \sigma'(z^{[L]}) \circ (a^{[L]} - y) , \quad (3.3)$$

$$\delta^{[l]} = \sigma'(z^{[l]}) \circ (W^{[l+1]})^T \delta^{[l+1]} , \text{ for } 2 \leq l \leq L-1 , \quad (3.4)$$

$$\frac{\partial C}{\partial b_j^{[l]}} = \delta_j^{[l]} , \text{ for } 2 \leq l \leq L , \quad (3.5)$$

$$\frac{\partial C}{\partial \omega_{jk}^{[l]}} = \delta_j^{[l]} a_k^{[l-1]} , \text{ for } 2 \leq l \leq L . \quad (3.6)$$

---

**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 (3.3)
  - 4: **Backward pass:** Form (3.4) 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 (3.5) and (3.6)
- 

Training GD on Over-parametrized Neural Networks has been understood only recently. See the following references:



Zhang, X., Yu, Y., Wang, L., Gu, Q., Learning One-hidden-layer ReLU networks via Gradient Descent, arXiv:1806.07808v1[stat.ML]

Du, S.S., Zhai, X., Póczos, B., Singh, A., Gradient Descent Provably Optimizes Over-parameterized Neural Networks, arXiv:1810.02054v1[cs.LG]

Allen-Zhu, Z., Li, Y., Song, Z., A Convergence Theory for Deep Learning via Over-Parametrization, arXiv:1811.03962v3[cs.LG]

Du, S.S., Lee, J.D., Li, H., Wang, L., Zhai, X., Gradient Descent Finds Global Minima of Deep Neural Networks, arXiv:1811.03804v1[cs.LG]

Neural Tangent Kernel. See the following paper:

Jacot, A., Gabriel, F., Hongler, C., Neural tangent kernel: Convergence and generalization in neural networks, NIPS, 2018.

Computational Graph and General Backpropagation for taking derivatives: basic idea in tensorflow. <https://www.tensorflow.org/tutorials/customization/autodiff>

### 3.3 Online Principle Component Analysis (PCA)

Reference: Li, C.J., Wang, M., Liu, H. and Zhang, T., Near-optimal stochastic approximation for online principal component estimation, *Mathematical Programming, Ser. B* (2018) 167:75-97.

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  $\Sigma$  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  $\Sigma$  that corresponds to the largest eigenvalue  $\lambda_1$ , 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:

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

Online PCA:  $\Pi \mathbf{u} = \frac{\mathbf{u}}{\|\mathbf{u}\|}$  retraction.

---

**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:**  $\mathbf{u}^{(N)}$

---

Oja's algorithm. See

Oja, E. (1982). Simplified neuron model as a principal component analyzer. *Journal of mathematical biology*, **15**(3), 267-273.

### 3.4 Line Search Methods

Line search methods in general

$$x^{k+1} = x^k + \alpha_k d^k, \quad k = 0, 1, 2, \dots \quad (3.7)$$

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

Preliminary conditions:

$$\begin{aligned} (1) \quad & -(d^k)^T \nabla f(x^k) \geq \bar{\varepsilon} \|\nabla f(x^k)\| \cdot \|d^k\|; \\ (2) \quad & \gamma_1 \|\nabla f(x^k)\| \leq \|d^k\| \leq \gamma_2 \|\nabla f(x^k)\| \text{ where } \bar{\varepsilon}, \gamma_1, \gamma_2 > 0. \end{aligned} \quad (3.8)$$

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

In fact one can expand

$$\begin{aligned} f(x^{k+1}) &= f(x^k + \alpha d^k) \\ &= f(x^k) + \alpha \nabla f(x^k)^T d^k + \alpha \int_0^1 [\nabla f(x^k + \gamma \alpha d^k) - \nabla f(x^k)] d^k d\gamma \\ &\leq f(x^k) - \alpha \left( \bar{\varepsilon} - \alpha \frac{L}{2} \gamma_2 \right) \|\nabla f(x^k)\| \cdot \|d^k\|. \end{aligned} \quad (3.9)$$

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

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

$$f(x^{k+1}) \leq f(x^k) - C\|\nabla f(x^k)\|^2 \text{ for some } C > 0. \quad (3.10)$$

Various Line Search Methods: (1)  $d^k = -S^k \nabla f(x^k)$ ,  $S^k$  symmetric positively definite,  $\lambda(S^k) \in [\gamma_1, \gamma_2]$ ; (2) Gauss-Southwell; (3) Stochastic coordinate descent; (4) Stochastic Gradient Methods; (5) Exact line search  $\min_{\alpha > 0} f(x^k + \alpha d^k)$ ; (6) Approximate Line Search, “weak Wolfe conditions”; (7) Backtracking Line Search.

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

$$\|\nabla f(x)\| \leq \varepsilon_g \text{ and } \lambda_{\min}(\nabla^2 f(x)) \geq -\varepsilon_H \quad (3.11)$$

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

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq M\|x - y\| \text{ for all } x, y \in \text{dom}(f). \quad (3.12)$$

Cubic upper bound

$$f(x + p) \leq f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + \frac{1}{6} M \|p\|^3. \quad (3.13)$$

For steepest descent step we have from (3.2) 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 (3.12), that

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

---

**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  $x^k$ 

```

---

Number of iterations

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

# 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](https://github.com/huwenqing0606/Computational_Optimization).
- [3] 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.
- [4] 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.
- [5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [6] 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.
- [7] Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 2000.