

Least Squares Method from the View Point of Deep Learning II: Generalization

Kazuyuki Fujii^{1,2}

¹International College of Arts and Sciences, Yokohama City University, Yokohama, Japan

²Department of Mathematical Sciences, Shibaura Institute of Technology, Saitama, Japan

Email: fujii@yokohama-cu.ac.jp

How to cite this paper: Fujii, K. (2018) Least Squares Method from the View Point of Deep Learning II: Generalization. *Advances in Pure Mathematics*, 8, 782-791. <https://doi.org/10.4236/apm.2018.89048>

Received: August 30, 2018

Accepted: September 22, 2018

Published: September 25, 2018

Copyright © 2018 by author and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

The least squares method is one of the most fundamental methods in Statistics to estimate correlations among various data. On the other hand, Deep Learning is the heart of Artificial Intelligence and it is a learning method based on the least squares method, in which a parameter called learning rate plays an important role. It is in general very hard to determine its value. In this paper we generalize the preceding paper [K. Fujii: Least squares method from the view point of Deep Learning: *Advances in Pure Mathematics*, 8, 485-493, 2018] and give an admissible value of the learning rate, which is easily obtained.

Keywords

Least Squares Method, Statistics, Deep Learning, Learning Rate, Gerschgorin's Theorem

1. Introduction

This paper is a sequel to the preceding paper [1].

The least squares method in Statistics plays an important role in almost all disciplines, from Natural Science to Social Science. When we want to find properties, tendencies or correlations hidden in huge and complicated data we usually employ the method. See for example [2].

On the other hand, Deep Learning is the heart of Artificial Intelligence and will become a most important field in Data Science in the near future. As to Deep Learning see for example [3]-[10].

Deep Learning may be stated as a successive learning method based on the least squares method. Therefore, to reconsider it from the view point of Deep Learning is natural and instructive. We carry out the calculation thoroughly of

the successive approximation called gradient descent sequence, in which a parameter ϵ called learning rate plays an important role.

One of main points is to determine the range of the learning rate, which is a very hard problem [8]. We showed in [1] that a difference in methods between Statistics and Deep Learning leads to different results when the learning rate changes.

We generalize the preceding results to the case of the least squares method by polynomial approximation. Our results may give a new insight to both Statistics and Data Science.

2. Least Squares Method

Let us explain the least squares method by polynomial approximation [9]. The model function $f(x)$ is a polynomial in x of degree M given by

$$f(x) = w_0 + w_1x + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j. \quad (1)$$

For N pieces of two dimensional real data

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

we assume that their scatter plot is given like **Figure 1**.

The coefficients of (1)

$$\mathbf{w} = (w_0, w_1, \dots, w_M)^T \quad (2)$$

must be determined by the data set (T denotes the transposition of a vector or a matrix).

For this set of data the error function is given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{y_i - f(x_i)\}^2 = \frac{1}{2} \sum_{i=1}^N \left(y_i - \sum_{j=0}^M w_j x_i^j \right)^2. \quad (3)$$

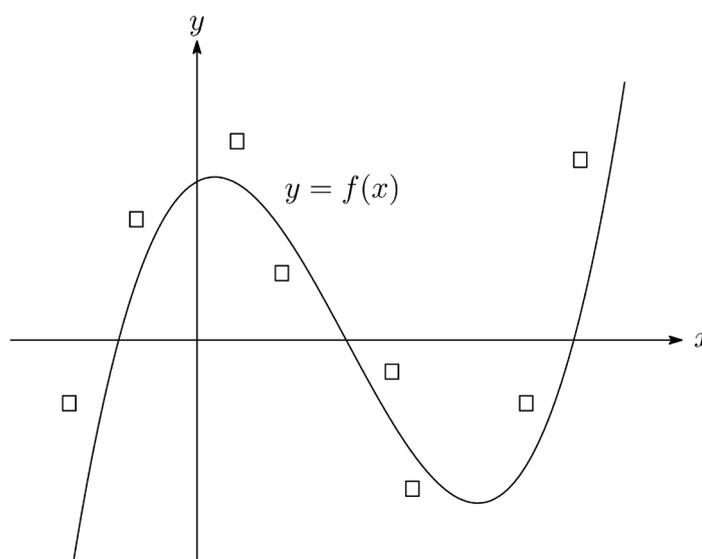


Figure 1. Scatter plot.

The aim of least squares method is to minimize the error function (3) with respect to \mathbf{w} in (2). Usually it is obtained by solving the simultaneous differentiable equations

$$\begin{cases} \frac{\partial E(\mathbf{w})}{\partial w_0} = 0, \\ \frac{\partial E(\mathbf{w})}{\partial w_1} = 0, \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w_M} = 0. \end{cases}$$

However, in this paper another approach based on quadratic form is given, which is instructive.

Let us calculate the error function (3). By using the definition of inner product

$$\langle \mathbf{a} | \mathbf{b} \rangle = a_1 b_1 + a_2 b_2 + \dots + a_n b_n = \mathbf{a}^T \mathbf{b}$$

it is not difficult to see

$$2E(\mathbf{w}) = \langle \mathbf{y} - \Phi \mathbf{w} | \mathbf{y} - \Phi \mathbf{w} \rangle \tag{4}$$

where

$$\mathbf{y} = (y_1, y_2, \dots, y_N)^T, \quad \mathbf{w} = (w_0, w_1, \dots, w_M)^T$$

and

$$\Phi = \begin{pmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^M \\ 1 & x_2^1 & x_2^2 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^1 & x_N^2 & \dots & x_N^M \end{pmatrix}.$$

Here we make an important

Assumption $N > M + 1$ and $\text{rank}(\Phi) = M + 1$ (full rank).

Let us deform (4). From

$$\begin{aligned} \langle \mathbf{y} - \Phi \mathbf{w} | \mathbf{y} - \Phi \mathbf{w} \rangle &= (\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) \\ &= (\mathbf{y}^T - \mathbf{w}^T \Phi^T) (\mathbf{y} - \Phi \mathbf{w}) \\ &= (\mathbf{w}^T \Phi^T - \mathbf{y}^T) (\Phi \mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{y} - \mathbf{y}^T \Phi \mathbf{w} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

we set for simplicity

$$\mathbf{x} = \mathbf{w}, \quad A = \Phi^T \Phi, \quad \mathbf{b} = \Phi^T \mathbf{y}, \quad c = \mathbf{y}^T \mathbf{y}.$$

Namely, we have a general quadratic form

$$2E(\mathbf{w}) = \langle \mathbf{y} - \Phi \mathbf{w} | \mathbf{y} - \Phi \mathbf{w} \rangle = \mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b} - \mathbf{b}^T \mathbf{x} + c. \tag{5}$$

On the other hand, the deformation of (5) is well-known.

Formula For a symmetric and invertible matrix A ($\exists A^{-1}$) we have

$$\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b} - \mathbf{b}^T \mathbf{x} + c = (\mathbf{x} - A^{-1} \mathbf{b})^T A (\mathbf{x} - A^{-1} \mathbf{b}) - \mathbf{b}^T A^{-1} \mathbf{b} + c. \tag{6}$$

The proof is easy. Since $A^T = A \Rightarrow (A^{-1})^T = A^{-1}$ we obtain

$$\begin{aligned} (\mathbf{x} - A^{-1}\mathbf{b})^T A(\mathbf{x} - A^{-1}\mathbf{b}) &= (\mathbf{x}^T - \mathbf{b}^T A^{-1}) A(\mathbf{x} - A^{-1}\mathbf{b}) \\ &= \mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b} - \mathbf{b}^T \mathbf{x} + \mathbf{b}^T A^{-1} \mathbf{b} \end{aligned}$$

and this gives (6).

Therefore, our case becomes

$$\begin{aligned} 2E(\mathbf{w}) &= \left\{ \mathbf{w} - (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \right\}^T \Phi^T \Phi \left\{ \mathbf{w} - (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \right\} \\ &\quad - \mathbf{y}^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \end{aligned} \quad (7)$$

because $\Phi^T \Phi$ is symmetric and invertible by the assumption.

If we choose

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad (8)$$

then the minimum is given by

$$2E(\mathbf{w})_{\min} = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} = \mathbf{y}^T \left\{ E_N - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right\} \mathbf{y} \quad (9)$$

where E_N is the N -dimensional identity matrix.

Our method is simple and clear (“smart” in our terminology).

3. Least Squares Method from Deep Learning

In this section we reconsider the least squares method in Section 2 from the view point of Deep Learning.

First we arrange the data in Section 2 like

$$\text{Input data : } \left\{ (1, x_j, x_j^2, \dots, x_j^M) \mid 1 \leq j \leq N \right\}$$

$$\text{Teacher signal : } \{y_1, y_2, \dots, y_N\}$$

and consider a simple neuron model in [11] (see **Figure 2**).

Here we use the polynomial (1) instead of the sigmoid function $z = \sigma(x)$.

In this case the square error function becomes

$$L(\mathbf{w}) = \frac{1}{2} \langle \mathbf{y} - \Phi \mathbf{w} \mid \mathbf{y} - \Phi \mathbf{w} \rangle = \frac{1}{2} \left(\mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{y} - \mathbf{y}^T \Phi \mathbf{w} + \mathbf{y}^T \mathbf{y} \right). \quad (10)$$

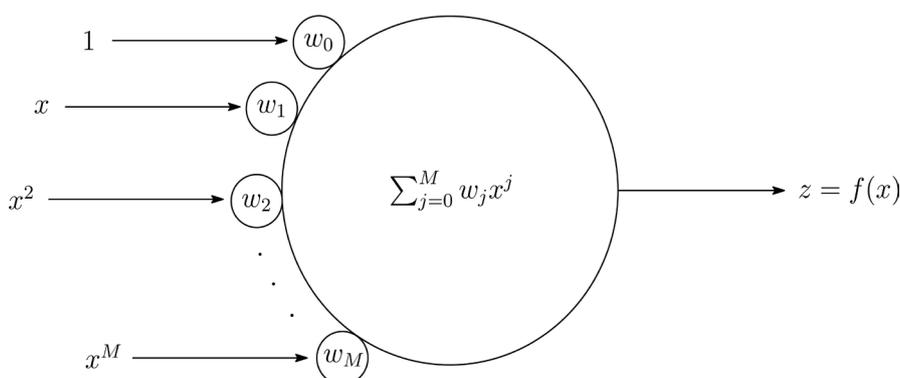


Figure 2. Simple neuron model.

We in general use $L(\mathbf{w})$ instead of $E(\mathbf{w})$ in (3).

Our aim is also to determine the parameters \mathbf{w} in order to minimize $L(\mathbf{w})$. However, the procedure is different from the least squares method in Section 2. This is an important and interesting point.

The parameters \mathbf{w} are determined successively by the gradient descent method (see for example [12]): For $t = 0, 1, 2, \dots$

$$\mathbf{w}(0) \rightarrow \mathbf{w}(1) \rightarrow \dots \rightarrow \mathbf{w}(t) \rightarrow \mathbf{w}(t+1) \rightarrow \dots$$

and

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \epsilon \frac{\partial L}{\partial \mathbf{w}(t)} \tag{11}$$

where

$$L = L(\mathbf{w}(t)) = \frac{1}{2} \left\{ \mathbf{w}(t)^T \Phi^T \Phi \mathbf{w}(t) - \mathbf{w}(t)^T \Phi^T \mathbf{y} - \mathbf{y}^T \Phi \mathbf{w}(t) + \mathbf{y}^T \mathbf{y} \right\} \tag{12}$$

and $\epsilon (0 < \epsilon < 1)$ is a small parameter called the **learning rate**.

The initial value $\mathbf{w}(0)$ is given appropriately. Pay attention that t is discrete time and T is the transposition.

Let us calculate (11) explicitly. Since

$$\frac{\partial L}{\partial \mathbf{w}(t)} = \Phi^T \Phi \mathbf{w}(t) - \Phi^T \mathbf{y}$$

from (12) we have

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \epsilon \left\{ \Phi^T \Phi \mathbf{w}(t) - \Phi^T \mathbf{y} \right\} = (E_{M+1} - \epsilon \Phi^T \Phi) \mathbf{w}(t) + \epsilon \Phi^T \mathbf{y}. \tag{13}$$

This equation is easily solved to be

$$\mathbf{w}(t) = (E_{M+1} - \epsilon \Phi^T \Phi)^t \mathbf{w}(0) + \left\{ E_{M+1} - (E_{M+1} - \epsilon \Phi^T \Phi)^t \right\} (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \tag{14}$$

for $t = 0, 1, 2, \dots$.

The proof is left to readers.

Since this is not a final form let us continue the calculation. From (14) we have

$$\lim_{t \rightarrow \infty} \mathbf{w}(t) = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \tag{15}$$

if

$$\lim_{t \rightarrow \infty} (E_{M+1} - \epsilon \Phi^T \Phi)^t = O_{M+1} \tag{16}$$

where O_{M+1} is the N -dimensional zero matrix. (15) is just the equation (8) and it is independent of ϵ .

Let us evaluate (14) further. The matrix $\Phi^T \Phi$ is positive definite, so all eigenvalues are positive. This can be shown as follows. Let us consider the eigenvalue equation

$$\Phi^T \Phi \mathbf{v} = \lambda \mathbf{v} \quad (\mathbf{v} \neq \mathbf{0}).$$

Then we have

$$\lambda \langle \mathbf{v} | \mathbf{v} \rangle = \langle \lambda \mathbf{v} | \mathbf{v} \rangle = \langle \Phi^T \Phi \mathbf{v} | \mathbf{v} \rangle = \langle \Phi \mathbf{v} | \Phi \mathbf{v} \rangle > 0 \Rightarrow \lambda > 0.$$

Therefore we can arrange all eigenvalues like

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{M+1} > 0.$$

Since $\Phi^T \Phi$ is symmetric, it is diagonalized as

$$\Phi^T \Phi = Q D Q^T \quad (17)$$

where Q is an element in $O(M+1)$ ($Q^T = Q^{-1}$) and D is a diagonal matrix

$$D = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_{M+1} \end{pmatrix}.$$

See for example [13].

By substituting (17) into (14) and using the equation

$$E_{M+1} - \epsilon \Phi^T \Phi = Q (E_{M+1} - \epsilon D) Q^T = Q \begin{pmatrix} 1 - \epsilon \lambda_1 & & & \\ & 1 - \epsilon \lambda_2 & & \\ & & \ddots & \\ & & & 1 - \epsilon \lambda_{M+1} \end{pmatrix} Q^T$$

we finally obtain

Theorem I A general solution to (14) is

$$\begin{aligned} \mathbf{w}(t) = & Q \begin{pmatrix} (1 - \epsilon \lambda_1)^t & & & \\ & (1 - \epsilon \lambda_2)^t & & \\ & & \ddots & \\ & & & (1 - \epsilon \lambda_{M+1})^t \end{pmatrix} Q^T \mathbf{w}(0) \\ & + Q \begin{pmatrix} \frac{1 - (1 - \epsilon \lambda_1)^t}{\lambda_1} & & & \\ & \frac{1 - (1 - \epsilon \lambda_2)^t}{\lambda_2} & & \\ & & \ddots & \\ & & & \frac{1 - (1 - \epsilon \lambda_{M+1})^t}{\lambda_{M+1}} \end{pmatrix} Q^T \Phi^T \mathbf{y}. \end{aligned} \quad (18)$$

This is our main result.

Next, let us show how to choose the learning rate ϵ ($0 < \epsilon < 1$), which is a very important problem in Deep Learning [7] [8].

Let us remember

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{M+1} > 0.$$

From (16) and (18) the equations

$$\lim_{t \rightarrow \infty} (E_{M+1} - \epsilon \Phi^T \Phi)^t = O_{M+1} \Leftrightarrow \lim_{t \rightarrow \infty} (1 - \epsilon \lambda_j)^t = 0 \text{ for } 1 \leq j \leq M+1 \quad (19)$$

determine the range of ϵ . Noting

$$|1 - x| < 1 (\Leftrightarrow 0 < x < 2) \Rightarrow \lim_{n \rightarrow \infty} (1 - x)^n = 0$$

and

$$\epsilon \lambda_1 \geq \epsilon \lambda_2 \geq \dots \geq \epsilon \lambda_{M+1} > 0$$

we obtain

Theorem II The learning rate ϵ must satisfy an inequality

$$0 < \epsilon \lambda_1 < 2 \Leftrightarrow 0 < \epsilon < \frac{2}{\lambda_1}. \tag{20}$$

The greater the value of ϵ , the sooner goes the gradient descent (11) so long as the convergence (19) is guaranteed. Let us note that the choice of the initial values $w(0)$ is irrelevant when the convergence condition (20) is satisfied.

Comment For example, if we choose ϵ like

$$\frac{2}{\lambda_1} < \epsilon < 1$$

then we cannot recover (15), which shows a difference in methods between Statistics and Deep Learning.

4. How to Estimate the Learning Rate

How do we calculate λ_1 ? Since $\{\lambda_j\}$ are the eigenvalues of the matrix $\Phi^T \Phi$, they satisfy the equation

$$F(\lambda_j) = 0, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{M+1} > 0$$

where $F(\lambda)$ is the characteristic polynomial of $\Phi^T \Phi$ given by

$$F(\lambda) = |\lambda E_{M+1} - \Phi^T \Phi|. \tag{21}$$

This is abstract, so let us deform (21). For simplicity we write Φ as

$$\Phi = \begin{pmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^M \\ 1 & x_2^1 & x_2^2 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^1 & x_N^2 & \dots & x_N^M \end{pmatrix} \equiv (\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}). \tag{22}$$

Then it is easy to see

$$\Phi^T \Phi = \left(\langle \mathbf{x}^{(i)} | \mathbf{x}^{(j)} \rangle \right)_{0 \leq i, j \leq M}, \quad \langle \mathbf{x}^{(i)} | \mathbf{x}^{(j)} \rangle = \sum_{k=1}^N x_k^{i+j}$$

where the notation $\langle \mathbf{a} | \mathbf{b} \rangle$ is the (real) inner product of vectors.

For clarity let us write down (21) explicitly.

$$F(\lambda) = \begin{vmatrix} \lambda - \langle \mathbf{x}^{(0)} | \mathbf{x}^{(0)} \rangle & -\langle \mathbf{x}^{(0)} | \mathbf{x}^{(1)} \rangle & \dots & -\langle \mathbf{x}^{(0)} | \mathbf{x}^{(M)} \rangle \\ -\langle \mathbf{x}^{(1)} | \mathbf{x}^{(0)} \rangle & \lambda - \langle \mathbf{x}^{(1)} | \mathbf{x}^{(1)} \rangle & \dots & -\langle \mathbf{x}^{(1)} | \mathbf{x}^{(M)} \rangle \\ \vdots & \vdots & \ddots & \vdots \\ -\langle \mathbf{x}^{(M)} | \mathbf{x}^{(0)} \rangle & -\langle \mathbf{x}^{(M)} | \mathbf{x}^{(1)} \rangle & \dots & \lambda - \langle \mathbf{x}^{(M)} | \mathbf{x}^{(M)} \rangle \end{vmatrix}.$$

As far as we know there is no viable method to determine the greatest root of

$F(\lambda) = 0$ if M is very large¹. Therefore, let us get satisfied by obtaining an approximate value which is both greater than λ_1 and easy to calculate.

For the purpose the Gerschgorin's theorem is very useful². Let $A = (a_{ij})$ be an $n \times n$ complex (real in our case) matrix, and we set

$$R_i = \sum_{j=1, j \neq i}^n |a_{ij}| \quad (23)$$

and

$$D(a_{ii}; R_i) = \{z \in \mathbb{C} \mid |z - a_{ii}| \leq R_i\} \quad (24)$$

for each i . This is a closed disc centered at a_{ii} with radius R_i called the Gerschgorin's disc.

Theorem (Gerschgorin [14]) For any eigenvalue λ of A we have

$$\lambda \in \bigcup_{i=1}^n D(a_{ii}; R_i). \quad (25)$$

The proof is simple. See for example [7].

Our case is real and $n = M + 1$ and

$$A = \Phi^T \Phi = \left(\left\langle \mathbf{x}^{(i)} \mid \mathbf{x}^{(j)} \right\rangle \right)_{0 \leq i, j \leq M}.$$

Therefore, all eigenvalues $\{\lambda\}$ satisfy

$$\lambda \in \bigcup_{i=1}^n D(a_{ii}; R_i) = \bigcup_{i=1}^{M+1} [a_{ii} - R_i, a_{ii} + R_i] \quad (26)$$

where $[A, B]$ is a closed interval and

$$a_{ii} = \left\langle \mathbf{x}^{(i-1)} \mid \mathbf{x}^{(i-1)} \right\rangle \text{ and } R_i = \sum_{k=1, k \neq i}^{M+1} \left| \left\langle \mathbf{x}^{(i-1)} \mid \mathbf{x}^{(k-1)} \right\rangle \right|.$$

If we define

$$\mathcal{F}_{M+1} = \max_{1 \leq i \leq M+1} \{a_{ii} + R_i\} \quad (27)$$

then it is easy to see

$$\lambda_1 \leq \mathcal{F}_{M+1} \Rightarrow \frac{2}{\mathcal{F}_{M+1}} \leq \frac{2}{\lambda_1}.$$

from (26).

Thus we arrive at an admissible value of the learning rate ϵ which is easily obtained.

Theorem III An admissible value of ϵ is

$$\epsilon = \frac{2}{\mathcal{F}_{M+1}}. \quad (28)$$

Let us show an example in the case of $M = 1$ ([1]), which is very instructive for non-experts.

Example In this case it is easy to see and we set

¹ $\Phi^T \Phi$ is not a sparse matrix.

² In my opinion this theorem is not so popular. Why?

$$\Phi^T \Phi = \begin{pmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{pmatrix} \equiv \begin{pmatrix} a & x \\ x & X \end{pmatrix}$$

for simplicity. Moreover, we may assume $x > 0$. Then from (21) we have

$$f(\lambda) = |\lambda E_2 - \Phi^T \Phi| = \lambda^2 - (a + X)\lambda + (aX - x^2)$$

and

$$\begin{aligned} \lambda_1 &= \frac{a + X + \sqrt{(a + X)^2 - 4(aX - x^2)}}{2} \\ &= \frac{a + X + \sqrt{(a - X)^2 + 4x^2}}{2}. \end{aligned}$$

On the other hand, from (27) we have

$$\mathcal{F}_2 = \max\{a + x, X + x\}$$

because $x > 0$.

Then it is easy to show

$$\max\{a + x, X + x\} \geq \frac{a + X + \sqrt{(a - X)^2 + 4x^2}}{2}.$$

To check this inequality is left to readers. Therefore, from (28) the admissible value becomes

$$\epsilon = \frac{2}{\max\{a + x, X + x\}}.$$

We emphasize once more that \mathcal{F}_{M+1} is easy to evaluate, while to calculate λ_1 is very hard if M is large.

5. Concluding Remarks

In this paper we have discussed the least squares method by polynomial approximation from the view point of Deep Learning and carried out calculation of the gradient descent thoroughly. A difference in methods between Statistics and Deep Learning delivers different results when the learning rate ϵ is changed. Theorem III is the first result to provide an admissible value of ϵ as far as we know.

Deep Learning plays an essential role in Data Science and maybe in almost all fields of Science. Therefore it is desirable for undergraduates to master it in the early stages. To master it they must study Calculus, Linear Algebra and Statistics from Mathematics. My textbook [7] is recommended.

Acknowledgements

We wish to thank Ryu Sasaki for useful suggestions and comments.

Conflicts of Interest

The author declares no conflicts of interest regarding the publication of this paper.

References

- [1] Fujii, K. (2018) Least Squares Method from the View Point of Deep Learning. *Advances in Pure Mathematics*, **8**, 485-493.
- [2] Wikipedia: Least Squares. https://en.m.wikipedia.org/wiki/Least_Squares
- [3] Wikipedia: Deep Learning. https://en.m.wikipedia.org/wiki/Deep_Learning
- [4] Goodfellow, I., Bengio, Y. and Courville, A. (2016) Deep Learning. The MIT Press, Cambridge.
- [5] Patterson, J. and Gibson, A. (2017) Deep Learning: A Practitioner's Approach, O'Reilly Media, Inc., Sebastopol.
- [6] Alpaydin, J. (2014) Introduction to Machine Learning. 3rd Edition, The MIT Press, Cambridge.
- [7] Fujii, K. (2018) Introduction to Mathematics for Understanding Deep Learning. Scientific Research Publishing Inc., Wuhan.
- [8] Okaya, T. (2015) Deep Learning (In Japanese). Kodansha Ltd., Tokyo.
- [9] Nakai, E. (2015) Introduction to Theory of Machine Learning (In Japanese). Gijutsu-Hyouronn Co., Ltd., Tokyo.
- [10] Amari, S. (2016) Brain Heart Artificial Intelligence (In Japanese). Kodansha Ltd., Tokyo.
- [11] Fujii, K. (2018) Mathematical Reinforcement to the Minibatch of Deep Learning. *Advances in Pure Mathematics*, **8**, 307-320. <https://doi.org/10.4236/apm.2018.83016>
- [12] Wikipedia: Gradient Descent. https://en.m.wikipedia.org/wiki/Gradient_descent
- [13] Kasahara, K. (2000) Linear Algebra (In Japanese). Saiensu Ltd., Tokyo.
- [14] Gerschgorin, S. (1931) Über die Abgrenzung der Eigenwerte einer Matrix. *Izv. Akad. Nauk. USSR Otd. Fiz.-Mat. Nauk*, **6**, 749-754.