# A Method to Reduce Data Dimensions in Machine Learning Programming

Yong Suk Suh[*], Seung Ki Shin, Dane Baang, Sang Mun Seo, Jong Bok Lee
*KJRR Project Div., Korea Atomic Energy Research Institute (KAERI), 111, Daedeok-Daero 989Beon-Gil,*
*Yuseong-Gu, Daejeon, 34057, Korea*
[*]*Corresponding author: yssuh@kaeri.re.kr*

## 1. Introduction

When programming a machine learning application, a large amount of input data can collected. It might be needed to segregate the data into several dimensions of data. While the data is being collected and segregated, some of the data might be unnecessarily collected or the data dimensions be unnecessarily segregated. In this case, we need to eliminate the unnecessary data or reduce the dimensions of the data. This paper deals with the reduction of data dimensions.

The machine learning programming for simple and multiple linear regressions were briefly reviewed in previous papers [1] and [2]. The machine learning method and the algebra method for the linear regressions were also compared in previous paper [3].

When we want to apply machine learning to nuclear facilities, we collect much data from field instruments, such as temperature, level, pressure, flow, etc. While the domain of a nuclear facility is analyzed, some of the field data are not needed for the machine learning. In this case, we want to determine how much data impacts the results of machine learning. If some of the data has a low impact, we can eliminate the data from the original data. This elimination will increase the computing time of the machine learning. In programming the linear regressions using the machine learning and/or algebra method, the reduction of data dimensions produces many benefits in computing time.

This paper introduces the principal component analysis (PCA) as a method to reduce the data dimensions and presents the test results of the PCA with three dimensional simple and arbitral sample data.

## 2. Principal Component Analysis

Data projection is one of ways to reduce data dimensions by orthogonally moving data in higher dimensional data space to lower dimensional data space [4]. The lower or higher dimensional data space then its original data space is called a hyper plane [5].

When we suppose that the data vector $x$ is projected onto a unit vector $w$, the projected data is represented in the form of $x^T w$, which is a scalar value. The projected data $x^T w$ lies on the $w$ so it can be represented in the vector form of $(x^T w)w$. When n numbers of $x$ data are projected onto the $w$, the mean of all projected data is represented in the form of $\frac{1}{n}\sum_{i=1}^{n} x_i^T w$. In this paper, we assume that all the $x$ data are normalized so that the mean of them becomes zero. Thus, $\frac{1}{n}\sum_{i=1}^{n} x_i^T w$ is also zero because $\frac{1}{n}\sum_{i=1}^{n} x_i^T$ is zero. This makes many of our

equations simple. The normalization of data enables the machining learning program to run fast and produce accurate outcomes.

When the data $x$ is projected onto a hyperplane $w$, we need to find the optimized vector to minimize the distance between the $x$ and $w$ and maximize the variance of the projected data $x^T w$.

The distance between the $x$ and $(x^T w)w$ can be represented in the square of the L2-norm form of $\|[x - (x^T w)w]\|^2$. An optimal $w$ can be obtained using the mean square error (MSE) called a least square fit for all $x$ vectors projected onto the $w$. The MSE with respect to $w$ is represented in the form of

$$\begin{aligned}
\text{MSE}(w) &= \frac{1}{n}\sum_{i=1}^{n}[x_i - (x_i^T w)w]^2 \\
&= \frac{1}{n}\sum_{i=1}^{n}[x_i^T x_i - 2x_i^T(x_i^T w)w + (x_i^T w)w^T(x_i^T w)w] \\
&= \frac{1}{n}\sum_{i=1}^{n}[x_i^T x_i - (x_i^T w)^2] \\
&= \frac{1}{n}\sum_{i=1}^{n} x_i^T x_i - \frac{1}{n}\sum_{i=1}^{n}(x_i^T w)^2.
\end{aligned} \tag{1}$$

From the Eq. (1), we know that the MSE$(w)$ is minimized when the $\frac{1}{n}\sum_{i=1}^{n}(x_i^T w)^2$ is maximized because the $\frac{1}{n}\sum_{i=1}^{n} x_i^T x_i$ is not related to $w$.

We represent the variance of projected data $x^T w$ in the form of

$$\text{Var}(x^T w) = \frac{1}{n}\sum_{i=1}^{n}(x_i^T w)^2 - \left(\frac{1}{n}\sum_{i=1}^{n} x_i^T w\right)^2. \tag{2}$$

In Eq. (2), the square of mean of $x^T w$, $(\frac{1}{n}\sum_{i=1}^{n} x_i^T w)^2$, is zero because we already assumed all the $x$ data were normalized in this paper. Thus, the variance is represented in the simple form of

$$\text{Var}(x^T w) = \frac{1}{n}\sum_{i=1}^{n}(x_i^T w)^2. \tag{3}$$

In Eq. (3), when the $\frac{1}{n}\sum_{i=1}^{n}(x_i^T w)^2$ is maximized, the Var$(x^T w)$ is also maximized. Thus, we know that the MSE$(w)$ is minimized when the Var$(x^T w)$ is maximized. We need to find a way to maximize the Var$(x^T w)$. In order to eliminate the $\Sigma$ in Eq. (3) and make the equation a simple form, we expand the data

vector $x$ to an n-by-m or n-by-n matrix $X$. In this case, the $w$ will be n number of vectors. When the $X$ is projected onto the $w$, the variance of $Xw$ is represented in the form of

$$\text{Var}(Xw) = \frac{1}{n}(Xw)^T(Xw) = \frac{1}{n}w^T X^T Xw$$
$$= w^T \frac{X^T X}{n} w = w^T Vw, \tag{4}$$
where $V$ is $\frac{X^T X}{n}$.

In Eq. (4), the $V$ is a matrix that represents a variance of $X$. When the matrix $X$ is given, we want to find an optimal (i.e., stationary) $w$ to maximize the $\text{Var}(Xw)$. Thus, the function of maximizing the $\text{Var}(Xw)$ is defined as an objective function. Along with the objective function, when $w$ is a unit vector, we define an constraint function

$$g(w) = w^T w - 1 \tag{5}$$

With Eqs. (4) and (5), we can adopt the Lagrange multiplier $\lambda$ to find an optimal $w$ so as to define a Lagrangian function with respect to $w$ and $\lambda$

$$\mathcal{L}(w, \lambda) = \text{Var}(Xw) - \lambda g(w)$$
$$= w^T Vw - \lambda(w^T w - 1) \tag{6}$$

In order to solve the Lagrangian function, we take partial derivatives

$$\frac{\partial \mathcal{L}}{\partial w} = 2Vw - 2\lambda w \tag{7}$$
$$\frac{\partial \mathcal{L}}{\partial \lambda} = w^T w - 1 \tag{8}$$

Thus, we obtain an optimal $w$ value at zero from Eq. (7)

$$Vw = \lambda w \tag{9}$$

In Eq. (9), the $w$ is an eigenvector of matrix $V$ and the $\lambda$ is an eigenvalue of eigenvector $w$. The linear transformation of $w$ depends on the $V$. The maximal variance of $w$ is in $\lambda$. Thus, the $\lambda$ is also the diagonal value of the $V$. The $\lambda$ is obtained by the eigenvalue decomposition for n-by-n matrix or singular value decomposition for n-by-m matrix of $V$. The largest $\lambda$ represents the largest variance of projected data.

### 3. Sample Test Cases

Sample test cases to reduce the data dimensions are arbitrarily chosen in this paper as shown in Table 1.

**Table 1: Sample data and normalized data**

| Given data | | | | | |
|---|---|---|---|---|---|
| $x_1$ | 1 | 3 | 5 | 8 | 10 |
| $x_2$ | 1 | 12 | 33 | 54 | 75 |
| $x_3$ | 1 | 1 | 33 | 33 | 99 |
| y | 1 | 2 | 3 | 4 | 5 |
| Normalized data of the given $x_1, x_2$ and $x_3$ | | | | | |
| $x_1$ | -0.49 | -0.27 | -0.04 | 0.29 | 0.51 |
| $x_2$ | -0.33 | -0.33 | -0.01 | -0.01 | 0.67 |
| $x_3$ | 0.46 | -0.31 | -0.03 | 0.26 | 0.54 |

Using the Sckit-learn PCA library, we can obtain projected sample data of $x_1, x_2$ and $x_3$ and eigenvalues of $\lambda$ as shown in Table 2. In Table 2, the $x_1$ is the most dominant characteristic sample data because its eignevalue is the largest which represents the largest variance.

**Table 2: Results of PCA**

| Projected sample data after PCA | | | | | |
|---|---|---|---|---|---|
| $x_1$ | -0.74 | -0.52 | -0.04 | 0.32 | 0.99 |
| $x_2$ | -0.11 | 0.04 | -0.02 | 0.22 | -0.13 |
| $x_3$ | -0.01 | 0.02 | -0.01 | -0.01 | 0.01 |
| $\lambda$ | 1.38 | 0.29 | 0.03 | | |

When we perform the linear regression with the normalized data in Table 1, we obtain the values of $w$ (i.e., coefficient of $x_1, x_2$ and $x_3$) and the bias (i.e., y-intercept) for the regression linear line, and predicted $y$ values and R-squared value as shown in Table 3.

**Table 3: Linear regression with the normalized data in Table 1**

| $w$ | 4.89 | 0.43 | -1.38 | | |
|---|---|---|---|---|---|
| bias | 3 | | | | |
| Predicted y | 1.10 | 1.98 | 2.82 | 4.06 | 5.04 |
| $R^2$ | 0.995 | | | | |

When we perform the linear regression with the projected data in Table 2, we obtain the $w$ and the bias for each data dimensions, and predicted $y$ values and R-squared value as shown in Table 4.

**Table 4: Linear regression with the projected data in Table 2**

| Case 1: Linear regression with $x_1, x_2$ and $x_3$ (i.e., Data dimensions are 3.) | | | | | |
|---|---|---|---|---|---|
| $w$ | 2.26 | 1.68 | 4.25 | | |
| bias | 3 | | | | |
| Predicted y | 1.10 | 1.98 | 2.82 | 4.06 | 5.04 |
| $R^2$ | 0.995 | | | | |
| Case 2: Linear regression with $x_1$ and $x_2$ (i.e., Data dimensions are 2.) | | | | | |
| $w$ | 2.26 | 1.68 | | | |
| bias | 3 | | | | |
| Predicted y | 1.14 | 1.89 | 2.86 | 4.09 | 5.02 |
| $R^2$ | 0.994 | | | | |
| Case 3: Linear regression with $x_1$ (i.e., Data dimension is 1.) | | | | | |
| $w$ | 2.26 | | | | |
| bias | 3 | | | | |
| Predicted y | 1.33 | 1.82 | 2.90 | 3.71 | 5.24 |
| $R^2$ | 0.971 | | | | |

In Table 4, Case 1 does not reduce the dimensions but Case 2 and 3 do. Case 1 generates the same results as the linear regression with the normalized sample data because the data dimensions are not reduced. Case 2 dose not significantly degrade the characteristic of the original data but Case 3 degrades it a little bit. If we allow the threshold of R-squared value is 9.5, Case 3 is also effective. We need to study or establish a formal method to measure how much the data dimensions can be reduced while still maintaining the characteristic of the original data. The programming is done using Spyder 3.3.6 in Anaconda 3 1.9.12, Python 3.7.4, and Sckit-learn library.

### 4. Conclusions

This paper introduced the principal component analysis (PCA), which we used to reduce data dimensionality, and presented test results of the reduction of data dimensions using the PCA method.

The PCA method showed an effective method for the reduction. There is a further need to establish a formal method to determine the level of reduction.

A further study will explore the application of singular value decomposition.

### REFERENCES

[1] Yong Suk Suh, et al., A Brief Review of a Machine Learning Programming of Simple Linear Regression, Transactions of the Korean Nuclear Society Spring Meeting, Jeju, Korea, May 17-18, 2018.
[2] Yong Suk Suh, et al., Considerations on Machine Learning Programming of Multiple Linear Regression, Transactions of the Korean Nuclear Society Spring Meeting, Jeju, Korea, May 23-24, 2019
 [3] Yong Suk Suh, et al., Comparisons of Linear Regression Methods: Machine Learning Method and Linear Algebra Method, Transactions of the Korean Nuclear Society Virtual Spring Meeting, July 9-10, 2020.
[4] Davie C. Lay, et al., LINEAR ALGEBRA AND ITS APPLICATIONS, 5th Edition, PEARSON, 2014.
[5] http://en.wikipedia.org/wiki/Hyperplane