

**Learning with slight forgetting optimizes sensorimotor transformation in
redundant motor systems**

Masaya Hirashima, Daichi Nozaki

Graduate School of Education, The University of Tokyo,
Bunkyo-ku, Tokyo 113-0033, Japan

1. Optimal synaptic weights

To find the optimal synaptic weights, we must define the criteria for optimality (i.e., the cost function). It is reasonable to assume that the optimal synaptic weights would minimize the movement error and the motor effort, wherever a target appears in the task space. Mathematically, the aim is to find the synaptic weights that minimize the expected value of the weighted sum of the error cost and the motor effort cost:

$$\begin{aligned} E[J] &= E[\alpha J_e + \beta J_m] \\ &= E\left[\alpha\left(\frac{1}{2}\mathbf{e}^T\mathbf{e}\right) + \beta\left(\frac{1}{2}\mathbf{r}^T\mathbf{r}\right)\right] \end{aligned} \quad (\text{A1})$$

with the condition that the mean (\mathbf{m}) and covariance matrix (\mathbf{C}) of target $\boldsymbol{\tau}$ are:

$$\begin{cases} \mathbf{m} = (0 \ 0)^T \\ \mathbf{C} = c \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \end{cases} \quad (\text{A2})$$

Because the error cost (J_e) can be represented as $J_e = \frac{1}{2}(\mathbf{B}\boldsymbol{\tau})^T\mathbf{B}\boldsymbol{\tau}$ using $\mathbf{B} = \mathbf{M}\mathbf{W} - \mathbf{I} \in \mathfrak{R}^{2 \times 2}$, the expected value of the error cost can be written as:

$$\begin{aligned} E[J_e] &= E\left[\frac{1}{2}\mathbf{e}^T\mathbf{e}\right] = E\left[\frac{1}{2}(\mathbf{B}\boldsymbol{\tau})^T\mathbf{B}\boldsymbol{\tau}\right] \\ &= \frac{1}{2}\left(\text{Tr}(\mathbf{B}\mathbf{C}\mathbf{B}^T) + (\mathbf{B}\mathbf{m})^T\mathbf{B}\mathbf{m}\right) \\ &= \frac{1}{2}\text{Tr}(\mathbf{B}\mathbf{C}\mathbf{B}^T) \\ &= \frac{c}{2}(B_{11}^2 + B_{12}^2 + B_{21}^2 + B_{22}^2) \end{aligned} \quad (\text{A3})$$

The expected value of the motor effort cost (J_m) can be written as:

$$\begin{aligned}
E[J_m] &= E\left[\frac{1}{2}\mathbf{r}^T\mathbf{r}\right] = E\left[\frac{1}{2}(\mathbf{W}\boldsymbol{\tau})^T\mathbf{W}\boldsymbol{\tau}\right] \\
&= \frac{1}{2}\left(\text{Tr}(\mathbf{W}\mathbf{C}\mathbf{W}^T) + (\mathbf{W}\mathbf{m})^T\mathbf{W}\mathbf{m}\right) \\
&= \frac{1}{2}\text{Tr}(\mathbf{W}\mathbf{C}\mathbf{W}^T) \tag{A4} \\
&= \frac{c}{2}\left(\sum W_{i1}^2 + \sum W_{i2}^2\right) \\
&= \frac{c}{2}\left(\sum_{ij} W_{ij}^2\right)
\end{aligned}$$

It is now apparent that the Moore-Penrose pseudoinverse ($\mathbf{M}^+ \in \mathfrak{R}^{n \times 2}$) of \mathbf{M} is the optimal solution that minimizes the cost function $E[J]$ among the many solutions satisfying zero error, for the following two reasons. First, since $\mathbf{B} = \mathbf{M}\mathbf{W} - \mathbf{I} = \mathbf{M}\mathbf{M}^+ - \mathbf{I} = \mathbf{0}$, $E[J_e]$ becomes 0. Second, the 1st column vector ($\mathbf{M}_{(1)}^+$) of \mathbf{M}^+ is the solution with the smallest Euclidean norm among the many solutions of $\mathbf{W}_{(1)} \in \mathfrak{R}^{n \times 1}$ that satisfy the equation $\mathbf{M}\mathbf{W}_{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, and the 2nd column vector ($\mathbf{M}_{(2)}^+$) of \mathbf{M}^+ is the solution with the smallest Euclidean norm among the many solutions of $\mathbf{W}_{(2)} \in \mathfrak{R}^{n \times 1}$ that satisfy the equation $\mathbf{M}\mathbf{W}_{(2)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$. Therefore, $E[J_m]$ is also minimized, because $E[J_m]$ is represented as the sum of squares of the elements of the two vectors $\mathbf{W}_{(1)}$ and $\mathbf{W}_{(2)}$ in Eq. (A4).

2. Mathematical proof of convergence

In this section, we prove that the synaptic weight matrix \mathbf{W} converges to the pseudoinverse of the matrix \mathbf{M} if the synaptic weights are modified by the feedback-with-decay rule (Eq. (3) in the main text) in the linear neural network model (Figure 1B). \mathbf{W} at the $(t+1)$ th trial can be written as:

$$\begin{aligned}
\mathbf{W}(t+1) &= \mathbf{W}(t) + \Delta\mathbf{W}(t) \\
&= \mathbf{W}(t) + \left(-\alpha \frac{\partial J_e}{\partial \mathbf{W}(t)} - \beta \mathbf{W}(t) \right) \\
&= \mathbf{W}(t) - \alpha \mathbf{M}^T (\mathbf{M}\mathbf{W}(t) - \mathbf{I}) \boldsymbol{\tau}(t) \boldsymbol{\tau}(t)^T - \beta \mathbf{W}(t)
\end{aligned} \tag{A5}$$

This equation can be averaged to give:

$$E[\mathbf{W}(t+1)] = \mathbf{W}(t) - \alpha \mathbf{M}^T (\mathbf{M}\mathbf{W}(t) - \mathbf{I}) \mathbf{C} - \beta \mathbf{W}(t) \tag{A6}$$

where $\mathbf{C} = E[\boldsymbol{\tau}(t)\boldsymbol{\tau}(t)^T]$.

2.1. Uniform presentation of targets

If we assume that the targets are presented randomly and uniformly in space, \mathbf{C} becomes $c\mathbf{I} \in \mathfrak{R}^{2 \times 2}$ ($c = 0.5$, if the eight targets in Figure 1A are used) and Eq. (A6) becomes:

$$E[\mathbf{W}(t+1)] = ((1-\beta)\mathbf{I} - c\alpha\mathbf{M}^T\mathbf{M})\mathbf{W}(t) + c\alpha\mathbf{M}^T \tag{A7}$$

Using following substitutions:

$$\begin{cases} (1-\beta)\mathbf{I} - c\alpha\mathbf{M}^T\mathbf{M} = \mathbf{Q} \\ c\alpha\mathbf{M}^T = \mathbf{R} \end{cases} \tag{A8}$$

we can express Eq. (A7) as follows:

$$E[\mathbf{W}(t+1)] = \mathbf{Q}\mathbf{W}(t) + \mathbf{R} \tag{A9}$$

Thus, $\mathbf{W}(t)$ at the t^{th} trial can be generally written as:

$$\mathbf{W}(t) = \mathbf{Q}^t \mathbf{W}(0) + (\mathbf{Q}^{t-1} + \mathbf{Q}^{t-2} + \dots + \mathbf{Q} + \mathbf{I}) \mathbf{R} \tag{A10}$$

Because \mathbf{Q}^t can be expressed in the following form:

$$\mathbf{Q}^t = \mathbf{V} \begin{bmatrix} ((1-\beta) - c\alpha\lambda_1)^t & & \dots & 0 \\ & ((1-\beta) - c\alpha\lambda_2)^t & & \vdots \\ & & (1-\beta)^t & \\ \vdots & & & \ddots \\ 0 & \dots & & (1-\beta)^t \end{bmatrix} \mathbf{V}^{-1} \tag{A11}$$

where λ_1 and λ_2 are the eigenvalues of the symmetric matrix $\mathbf{M}^T\mathbf{M}$, as the number of trials approaches infinity ($t \rightarrow \infty$),

$$\mathbf{Q}^t \rightarrow \mathbf{0} \quad (\text{A12})$$

if the following conditions are satisfied:

$$\begin{cases} -1 < (1-\beta) - c\alpha\lambda_1 < 1 \\ -1 < (1-\beta) - c\alpha\lambda_2 < 1 \\ -1 < (1-\beta) < 1 \end{cases} \quad (\text{A13})$$

Because α , λ_1 , and λ_2 are positive values, the conditions in (A13) can be more simply written as follows:

$$\begin{cases} \beta < 2 - c\alpha\lambda_1 \\ \beta < 2 - c\alpha\lambda_2 \\ \beta > 0 \end{cases} \quad (\text{A14})$$

Furthermore, under these conditions, as the number of trials approaches infinity ($t \rightarrow \infty$):

$$\mathbf{Q}^{t-1} + \mathbf{Q}^{t-2} + \dots + \mathbf{Q} + \mathbf{I} \rightarrow (\mathbf{I} - \mathbf{Q})^{-1} \quad (\text{A15})$$

From (A10), (A12), and (A15), as $t \rightarrow \infty$:

$$\begin{aligned} \mathbf{W}(t) &\rightarrow (\mathbf{I} - \mathbf{Q})^{-1} \mathbf{R} \\ &= (\mathbf{I} - \{(1-\beta)\mathbf{I} - c\alpha\mathbf{M}^T\mathbf{M}\})^{-1} c\alpha\mathbf{M}^T \\ &= \left(\frac{\beta}{c\alpha}\mathbf{I} + \mathbf{M}^T\mathbf{M}\right)^{-1} \mathbf{M}^T \\ &= \mathbf{M}^T \left(\frac{\beta}{c\alpha}\mathbf{I} + \mathbf{M}\mathbf{M}^T\right)^{-1} \end{aligned} \quad (\text{A16})$$

When $c\alpha \gg \beta$, $\mathbf{W}(t)$ converges to $\mathbf{M}^T(\mathbf{M}\mathbf{M}^T)^{-1}$, which is the Moore-Penrose pseudo-inverse of \mathbf{M} . Taken together, the conditions for convergence to the Moore-Penrose pseudo-inverse are:

$$\begin{cases} \beta < 2 - c\alpha\lambda_1 \\ \beta < 2 - c\alpha\lambda_2 \\ 0 < \beta \ll c\alpha \end{cases} \quad (\text{A17})$$

The conditions in (A17) are satisfied, because $\alpha = 20$, $\beta = 1.0 \times 10^{-4}$, $c = 0.5$, $\lambda_1 = 0.0036$, and $\lambda_2 = 0.00072$.

Although we have proven the convergence of \mathbf{W} for the 2-dimensional task, the proof also holds true for 3-dimensional or much higher-dimensional tasks. For a D -dimensional task, $\mathbf{W} \in \mathfrak{R}^{n \times D}$ converges to the Moore-Penrose pseudo-inverse of $\mathbf{M} \in \mathfrak{R}^{D \times n}$ under the following conditions:

$$\begin{cases} \beta < 2 - c\alpha\lambda_1 \\ \beta < 2 - c\alpha\lambda_2 \\ \vdots \\ \beta < 2 - c\alpha\lambda_D \\ 0 < \beta \ll c\alpha \end{cases} \quad (\text{A18})$$

2.2. Non-uniform presentation of targets

If the distribution of the targets is not uniform in space, \mathbf{C} does not take the form $c\mathbf{I} \in \mathfrak{R}^{2 \times 2}$. A transformation is conducted as follows:

$$\hat{\boldsymbol{\tau}}(t) = \boldsymbol{\sigma}\mathbf{V}\boldsymbol{\tau}(t) \quad (\text{A19})$$

where $\mathbf{V} = (\mathbf{v}_1 \mathbf{v}_2)$ is a matrix that consists of the eigenvectors of the matrix \mathbf{C} , and

$$\boldsymbol{\sigma} = \begin{pmatrix} 1/\sqrt{\lambda_1} & 0 \\ 0 & 1/\sqrt{\lambda_2} \end{pmatrix} \quad (\text{A20})$$

where λ_1 and λ_2 are the eigenvalues. Then, $\hat{\mathbf{C}} = E[\hat{\boldsymbol{\tau}}(t)\hat{\boldsymbol{\tau}}(t)^T]$ becomes $\mathbf{I} \in \mathfrak{R}^{2 \times 2}$, if there are at least two target vectors that are linearly independent.

In addition to the targets, the MDVs must be transformed in the same manner, $\hat{\mathbf{M}} = \boldsymbol{\sigma}\mathbf{V}\mathbf{M}$. When the synaptic weights from the new input vectors $\hat{\boldsymbol{\tau}}(t)$ to

the actuators are defined as $\hat{\mathbf{W}}$ and the output vectors as $\hat{\mathbf{T}}$, the transformed network has the same form as the original network. Therefore $\hat{\mathbf{W}}(t)$ converges to the pseudo-inverse of $\hat{\mathbf{M}}$:

$$\hat{\mathbf{W}}(t) \rightarrow \hat{\mathbf{M}}^T (\hat{\mathbf{M}} \hat{\mathbf{M}}^T)^{-1} \quad (\text{A21})$$

By substituting $\hat{\mathbf{M}} = \sigma \mathbf{V} \mathbf{M}$ into Eq. (A21), we obtain:

$$\begin{aligned} \hat{\mathbf{W}}(t) &\rightarrow (\sigma \mathbf{V} \mathbf{M})^T (\sigma \mathbf{V} \mathbf{M} (\sigma \mathbf{V} \mathbf{M})^T)^{-1} \\ &= \mathbf{M}^T (\mathbf{M} \mathbf{M}^T)^{-1} (\sigma \mathbf{V})^{-1} \end{aligned} \quad (\text{A22})$$

Because the original \mathbf{W} is expressed as:

$$\mathbf{W} = \hat{\mathbf{W}} \sigma \mathbf{V} \quad (\text{A23})$$

$\mathbf{W}(t)$ converges to the Moore-Penrose pseudo-inverse of \mathbf{M} as follows:

$$\begin{aligned} \mathbf{W}(t) &\rightarrow \hat{\mathbf{W}} \sigma \mathbf{V} \\ &= \mathbf{M}^T (\mathbf{M} \mathbf{M}^T)^{-1} (\sigma \mathbf{V})^{-1} \sigma \mathbf{V} \\ &= \mathbf{M}^T (\mathbf{M} \mathbf{M}^T)^{-1} \end{aligned} \quad (\text{A24})$$

2.3. “Feedback-only” rule

Here, we mathematically assess where the synaptic weight matrix $\mathbf{W}(t)$ converges if the decay is not incorporated. As has been already shown, $\mathbf{W}(t)$ at the t^{th} trial can be written as:

$$\mathbf{W}(t) = \mathbf{Q}^t \mathbf{W}(0) + (\mathbf{Q}^{t-1} + \mathbf{Q}^{t-2} + \dots + \mathbf{Q} + \mathbf{I}) \mathbf{R} \quad (\text{A10})$$

If the decay is not incorporated (i.e., $\beta = 0$), \mathbf{Q}^t is expressed as follows:

$$\mathbf{Q}^t = \mathbf{V} \begin{bmatrix} (1 - c\alpha\lambda_1)^t & & \dots & 0 \\ & (1 - c\alpha\lambda_2)^t & & \vdots \\ & & 1^t & \\ \vdots & & & \ddots \\ 0 & \dots & & 1^t \end{bmatrix} \mathbf{V}^{-1} \quad (\text{A25})$$

where λ_1 and λ_2 are the eigenvalues of the symmetric matrix $\mathbf{M}^T\mathbf{M}$ and the matrix \mathbf{V} ($= (\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_n)$) consists of the eigenvectors of the symmetric matrix $\mathbf{M}^T\mathbf{M}$. When (A14) is satisfied, as the number of trials approaches infinity ($t \rightarrow \infty$):

$$\mathbf{Q}^t \rightarrow \mathbf{V} \begin{bmatrix} 0 & & \dots & 0 \\ & 0 & & \vdots \\ & & 1 & \\ \vdots & & & \ddots \\ 0 & \dots & & 1 \end{bmatrix} \mathbf{V}^{-1} = \mathbf{A} (\neq \mathbf{0}) \quad (\text{A26})$$

Thus, $|\mathbf{A}\mathbf{W}(0)| \leq |\mathbf{W}(0)|$, which means that \mathbf{A} never increases $|\mathbf{W}(0)|$. As for the second term of Eq.(A10), as the number of trials approaches infinity ($t \rightarrow \infty$):

$$\begin{aligned} (\mathbf{Q}^{t-1} + \mathbf{Q}^{t-2} + \dots + \mathbf{Q} + \mathbf{I})\mathbf{R} &\rightarrow (\mathbf{I} - \mathbf{Q})^{-1}\mathbf{R} \\ &= \mathbf{M}^T \left(\frac{\beta}{c\alpha} \mathbf{I} + \mathbf{M}\mathbf{M}^T \right)^{-1} \end{aligned} \quad (\text{A27})$$

When $\beta = 0$, it converges to $\mathbf{M}^T(\mathbf{M}\mathbf{M}^T)^{-1}$. Taken together, the synaptic weight matrix converges as:

$$\mathbf{W}(t) \rightarrow \mathbf{A}\mathbf{W}(0) + \mathbf{M}^T(\mathbf{M}\mathbf{M}^T)^{-1} \quad (t \rightarrow \infty) \quad (\text{A28})$$

Finally, $\mathbf{M}\mathbf{W}(t)$ converges on:

$$\begin{aligned} \mathbf{M}\mathbf{W}(t) &\rightarrow \mathbf{M}[\mathbf{A}\mathbf{W}(0) + \mathbf{M}^T(\mathbf{M}\mathbf{M}^T)^{-1}] \quad (t \rightarrow \infty) \\ &= \mathbf{M}\mathbf{A}\mathbf{W}(0) + \mathbf{I} \end{aligned} \quad (\text{A29})$$

Because \mathbf{M} can be decomposed using singular value decomposition as follows:

$$\mathbf{M} = \mathbf{U} \begin{bmatrix} \sqrt{\lambda_1} & 0 & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & 0 & \dots & 0 \end{bmatrix} \mathbf{V}^T \quad (\text{A30})$$

where the matrix \mathbf{U} ($= (\mathbf{u}_1 \mathbf{u}_2)$) includes the eigenvectors of the symmetric matrix $\mathbf{M}\mathbf{M}^T$. Therefore, $\mathbf{M}\mathbf{A}\mathbf{W}(0)$ is calculated as:

$$\mathbf{MAW}(0) = \left(\mathbf{U} \begin{bmatrix} \sqrt{\lambda_1} & 0 & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & 0 & \dots & 0 \end{bmatrix} \mathbf{V}^T \right) \left(\mathbf{V} \begin{bmatrix} 0 & \dots & 0 \\ & 0 & \vdots \\ & 1 & \\ \vdots & & \ddots \\ 0 & \dots & & 1 \end{bmatrix} \mathbf{V}^{-1} \mathbf{W}(0) \right)$$

Because \mathbf{V} is a matrix containing the eigenvectors of the symmetric matrix $\mathbf{M}^T\mathbf{M}$, \mathbf{V} is an orthogonal matrix satisfying $\mathbf{V}^T\mathbf{V} = \mathbf{I}$. Therefore, it is further calculated as:

$$\begin{aligned} \mathbf{MAW}(0) &= \mathbf{U} \begin{bmatrix} \sqrt{\lambda_1} & 0 & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} 0 & \dots & 0 \\ & 0 & \vdots \\ & 1 & \\ \vdots & & \ddots \\ 0 & \dots & & 1 \end{bmatrix} \mathbf{V}^{-1} \mathbf{W}(0) \\ &= \mathbf{U} \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \mathbf{V}^{-1} \mathbf{W}(0) = \mathbf{0} \end{aligned} \quad (\text{A31})$$

Note that $\mathbf{AW}(0)$ itself is not $\mathbf{0}$ and depends on the initial synaptic weight $\mathbf{W}(0)$, but it always satisfies $\mathbf{MAW}(0) = \mathbf{0}$, irrespective of the initial synaptic weight matrix. From (A29) and (A31), as $t \rightarrow \infty$:

$$\mathbf{MW}(t) \rightarrow \mathbf{I} \quad (\text{A32})$$

Thus, the converged weight matrix produces zero error for the arbitrary target $\boldsymbol{\tau}$ because it always satisfies $\mathbf{T} = \mathbf{MW}\boldsymbol{\tau} = \boldsymbol{\tau}$.

In summary, by the “feedback-only” rule, the synaptic weight matrix $\mathbf{W}(t)$ converges with a different matrix depending on the initial synaptic weight matrix (i.e., $\mathbf{W}(t) \rightarrow \mathbf{AW}(0) + \mathbf{M}^T(\mathbf{MM}^T)^{-1}$), while the converged weight matrices always satisfy $\mathbf{MW} = \mathbf{I}$ and hence produce zero error for the arbitrary target.