

An Adaptive Alternating-direction-method of Multipliers-Incorporated Approach to Nonnegative Latent Factor Analysis:

Supplementary File

This is the supplementary file for the paper entitled *An Adaptive Alternating-direction-method of Multipliers-incorporated Approach to Non-negative Latent Factor Analysis*. Detailed convergence proof of A²NLF, additional experimental results and several supplementary files are presented here.

I. CONVERGENCE OF A²NLF

A. Proof of Lemma 1

Note that A²NLF's learning objective is non-convex. According to [17], any of its limit points where the gradient becomes zero can be a local/global optimum, or saddle point. Hence, such a limit point can be treated as a solution. Supposing that the optimal solution to $a_{u,k(q)}$ by (4c) is $a_{u,k(q)}^t$. Thus, it fulfills the following condition:

$$\lambda_{(q)}^t |\Lambda(u)| \left(p_{u,k(q)}^t - a_{u,k(q)}^t + \frac{h_{u,k(q)}^{t-1}}{\lambda_{(q)}^t |\Lambda(u)|} \right) = 0. \quad (S1)$$

Following (4e) and (5c), by applying the update rule of $h_{u,k(q)}$ to (S1), we have:

$$h_{u,k(q)}^t = (\eta_{(q)}^t - 1) \lambda_{(q)}^t |\Lambda(u)| (p_{u,k(q)}^t - a_{u,k(q)}^t). \quad (S2)$$

Then (17) stands based on (S2). Following the same principle, we can derive the optimality condition of (5b) related to $x_{i,k(q)}$:

$$\lambda_{(q)}^t |\Lambda(i)| \left(z_{i,k(q)}^t - x_{i,k(q)}^t + \frac{w_{i,k(q)}^{t-1}}{\lambda_{(q)}^t |\Lambda(i)|} \right) = 0, \Rightarrow w_{i,k(q)}^t = (\eta_{(q)}^t - 1) \lambda_{(q)}^t |\Lambda(i)| (z_{i,k(q)}^t - x_{i,k(q)}^t). \quad (S3)$$

Then (18) holds based on (S3). Hence, Lemma 1 holds, and **Step 1** is implemented. \square

B. Proof of Lemma 2

Considering the difference between $g(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1})$ and $g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1})$, we have:

$$\begin{aligned} & g(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1}) - g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1}) \\ &= \left(\sum_{i \in \Lambda(u)} z_{i,k(q)}^{t-1} \left(y_{u,i} - \sum_{l_1=1}^{k-1} p_{u,l_1(q)}^t z_{i,l_1(q)}^t - p_{u,k(q)}^t z_{i,k(q)}^{t-1} - \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} \right) - \lambda_{(q)}^t |\Lambda(u)| \left(p_{u,k(q)}^t - a_{u,k(q)}^{t-1} + \frac{h_{u,k(q)}^{t-1}}{\lambda_{(q)}^t |\Lambda(u)|} \right) \right) (p_{u,k(q)}^{t-1} - p_{u,k(q)}^t) \\ &+ \left(\sum_{u \in \Lambda(i)} p_{u,k(q)}^{t-1} \left(y_{u,i} - \sum_{l_1=1}^{k-1} p_{u,l_1(q)}^t z_{i,l_1(q)}^t - p_{u,k(q)}^{t-1} z_{i,k(q)}^t - \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} \right) - \lambda_{(q)}^t |\Lambda(i)| \left(z_{i,k(q)}^t - x_{i,k(q)}^{t-1} + \frac{w_{i,k(q)}^{t-1}}{\lambda_{(q)}^t |\Lambda(i)|} \right) \right) (z_{i,k(q)}^{t-1} - z_{i,k(q)}^t) \\ &- \frac{1}{2} \left(\sum_{i \in \Lambda(u)} (z_{i,k(q)}^{t-1})^2 + \lambda_{(q)}^t |\Lambda(u)| \right) (p_{u,k(q)}^t - p_{u,k(q)}^{t-1})^2 - \frac{1}{2} \left(\sum_{u \in \Lambda(i)} (p_{u,k(q)}^{t-1})^2 + \lambda_{(q)}^t |\Lambda(i)| \right) (z_{i,k(q)}^t - z_{i,k(q)}^{t-1})^2. \end{aligned} \quad (S4)$$

Considering (5a)'s optimality condition, (S4) is transformed as:

$$\begin{aligned}
& g\left(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1}\right) - g\left(s_{1(q)}^{t-1}, s_{2(q)}^{t-1}\right) \\
&= -\frac{1}{2} \left(\sum_{i \in \Lambda(u)} \left(z_{i,k(q)}^{t-1} \right)^2 + \lambda_{(q)}^t |\Lambda(u)| \right) \left(p_{u,k(q)}^t - p_{u,k(q)}^{t-1} \right)^2 - \frac{1}{2} \left(\sum_{u \in \Lambda(i)} \left(p_{u,k(q)}^{t-1} \right)^2 + \lambda_{(q)}^t |\Lambda(i)| \right) \left(z_{i,k(q)}^t - z_{i,k(q)}^{t-1} \right)^2.
\end{aligned} \tag{S5}$$

Thus, the difference between $g\left(s_{1(q)}^t, s_{2(q)}^{t-1}\right)$ and $g\left(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1}\right)$ is:

$$g\left(s_{1(q)}^t, s_{2(q)}^{t-1}\right) - g\left(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1}\right) = -\left(\lambda_{(q)}^t |\Lambda(u)|/2\right) \left(a_{u,k(q)}^t - a_{u,k(q)}^{t-1}\right)^2 - \left(\lambda_{(q)}^t |\Lambda(i)|/2\right) \left(x_{i,k(q)}^t - x_{i,k(q)}^{t-1}\right)^2. \tag{S6}$$

Moreover, $g\left(s_{1(q)}^t, s_{2(q)}^t\right)$ and $g\left(s_{1(q)}^t, s_{2(q)}^{t-1}\right)$ yields:

$$\begin{aligned}
& g\left(s_{1(q)}^t, s_{2(q)}^t\right) - g\left(s_{1(q)}^t, s_{2(q)}^{t-1}\right) \\
&= \left(p_{u,k(q)}^t - a_{u,k(q)}^t\right) \left(h_{u,k(q)}^t - h_{u,k(q)}^{t-1}\right) + \left(z_{i,k(q)}^t - x_{i,k(q)}^t\right) \left(w_{i,k(q)}^t - w_{i,k(q)}^{t-1}\right) \\
&\stackrel{(I)}{=} \frac{1}{\eta_{(q)}^t \lambda_{(q)}^t |\Lambda(u)|} \left(h_{u,k(q)}^t - h_{u,k(q)}^{t-1}\right)^2 + \frac{1}{\eta_{(q)}^t \lambda_{(q)}^t |\Lambda(i)|} \left(w_{i,k(q)}^t - w_{i,k(q)}^{t-1}\right)^2 \\
&\stackrel{(II)}{\leq} \frac{2|\Lambda(u)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(\left((\eta_{(q)}^t - 1) \lambda_{(q)}^t p_{u,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} p_{u,k(q)}^{t-1} \right)^2 + \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t a_{u,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} a_{u,k(q)}^{t-1} \right)^2 \right) \\
&+ \frac{2|\Lambda(i)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(\left((\eta_{(q)}^t - 1) \lambda_{(q)}^t z_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} z_{i,k(q)}^{t-1} \right)^2 + \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t x_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} x_{i,k(q)}^{t-1} \right)^2 \right),
\end{aligned} \tag{S7}$$

where (I) is based on the update rules of $(h_{u,k(q)}, w_{i,k(q)})$ given in (4e), (4f) and (5c), and (II) is achieved with *Lemma 1*. With (S5)-(S7), we have the following deduction:

$$\begin{aligned}
& g\left(s_{1(q)}^t, s_{2(q)}^t\right) - g\left(s_{1(q)}^{t-1}, s_{2(q)}^{t-1}\right) \\
&\leq -\frac{1}{2} \left(\sum_{i \in \Lambda(u)} \left(z_{i,k(q)}^{t-1} \right)^2 + \lambda_{(q)}^t |\Lambda(u)| \right) \left(p_{u,k(q)}^t - p_{u,k(q)}^{t-1} \right)^2 - \frac{\lambda_{(q)}^t |\Lambda(u)|}{2} \left(a_{u,k(q)}^t - a_{u,k(q)}^{t-1} \right)^2 \\
&- \frac{1}{2} \left(\sum_{u \in \Lambda(i)} \left(p_{u,k(q)}^{t-1} \right)^2 + \lambda_{(q)}^t |\Lambda(i)| \right) \left(z_{i,k(q)}^t - z_{i,k(q)}^{t-1} \right)^2 - \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} \left(x_{i,k(q)}^t - x_{i,k(q)}^{t-1} \right)^2 \\
&+ \frac{2|\Lambda(u)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(\left((\eta_{(q)}^t - 1) \lambda_{(q)}^t p_{u,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} p_{u,k(q)}^{t-1} \right)^2 + \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t a_{u,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} a_{u,k(q)}^{t-1} \right)^2 \right) \\
&+ \frac{2|\Lambda(i)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(\left((\eta_{(q)}^t - 1) \lambda_{(q)}^t z_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} z_{i,k(q)}^{t-1} \right)^2 + \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t x_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} x_{i,k(q)}^{t-1} \right)^2 \right).
\end{aligned} \tag{S8}$$

Owing to (20a), (21a) stands, which indicates that the augmented Lagrangian function (3) related to the q -th particle is non-increasing as $a_{u,k(q)}^t > 0$ and $x_{i,k(q)}^t > 0$. Then after the t -th iteration, the partial objective from (3) related to the q -th particle is formulated as:

$$\begin{aligned}
& g\left(s_{1(q)}^t, s_{2(q)}^t\right) \\
&= \frac{1}{2} \sum_{y_{u,i} \in \Lambda} \left(y_{u,i} - \sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t - \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} \right)^2 \\
&+ \sum_u \left(\left(\sum_{l_1=1}^k h_{u,l_1(q)}^t \left(p_{u,l_1(q)}^t - a_{u,l_1(q)}^t \right) \right) + \left(\sum_{l_2=k+1}^d h_{u,l_2(q)}^{t-1} \left(p_{u,l_2(q)}^{t-1} - a_{u,l_2(q)}^{t-1} \right) \right) \right) + \sum_i \left(\left(\sum_{l_1=1}^k w_{i,l_1(q)}^t \left(z_{i,l_1(q)}^t - x_{i,l_1(q)}^t \right) \right) + \left(\sum_{l_2=k+1}^d w_{i,l_2(q)}^{t-1} \left(z_{i,l_2(q)}^{t-1} - x_{i,l_2(q)}^{t-1} \right) \right) \right) \\
&+ \sum_u \frac{\lambda_{(q)}^t |\Lambda(u)|}{2} \left(\sum_{l_1=1}^k \left(p_{u,l_1(q)}^t - a_{u,l_1(q)}^t \right)^2 + \sum_{l_2=k+1}^d \left(p_{u,l_2(q)}^{t-1} - a_{u,l_2(q)}^{t-1} \right)^2 \right) + \sum_i \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} \left(\sum_{l_1=1}^k \left(z_{i,l_1(q)}^t - x_{i,l_1(q)}^t \right)^2 + \sum_{l_2=k+1}^d \left(z_{i,l_2(q)}^{t-1} - x_{i,l_2(q)}^{t-1} \right)^2 \right).
\end{aligned} \tag{S9}$$

By substituting (S2) and (S3) into (S9), we have:

$$\begin{aligned}
& g(s_{1(q)}^t, s_{2(q)}^t) \\
&= \frac{1}{2} \sum_{y_{u,i} \in \Lambda} \left(y_{u,i} - \sum_{l_1=1}^k p_{u,l_1}^t z_{i,l_1(q)}^t - \sum_{l_2=k+1}^d p_{u,l_2}^{t-1} z_{i,l_2(q)}^{t-1} \right)^2 \\
&+ \sum_u |\Lambda(u)| \left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t \sum_{l_1=1}^k \left(p_{u,l_1}^t - a_{u,l_1(q)}^t \right)^2 + \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} \sum_{l_2=k+1}^d \left(p_{u,l_2}^{t-1} - a_{u,l_2(q)}^{t-1} \right)^2 \right) \\
&+ \sum_i |\Lambda(i)| \left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t \sum_{l_1=1}^k \left(z_{i,l_1}^t - x_{i,l_1(q)}^t \right)^2 + \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} \sum_{l_2=k+1}^d w_{i,l_2}^t \left(z_{i,l_2}^{t-1} - x_{i,l_2(q)}^{t-1} \right)^2 \right) \\
&+ \sum_u \frac{\lambda_{(q)}^t |\Lambda(u)|}{2} \left(\sum_{l_1=1}^k \left(p_{u,l_1}^t - a_{u,l_1(q)}^t \right)^2 + \sum_{l_2=k+1}^d \left(p_{u,l_2}^{t-1} - a_{u,l_2(q)}^{t-1} \right)^2 \right) + \sum_i \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} \left(\sum_{l_1=1}^k \left(z_{i,l_1}^t - x_{i,l_1(q)}^t \right)^2 + \sum_{l_2=k+1}^d \left(z_{i,l_2}^{t-1} - x_{i,l_2(q)}^{t-1} \right)^2 \right).
\end{aligned} \tag{S10}$$

(S10) indicates that if (20b) is fulfilled, (21b) holds, thereby making (3) related to the q -th particle lower-bounded as $a_{u,k(q)}^t > 0$ and $x_{i,k(q)}^t > 0$. Based on the above inferences, *Lemma 2* stands, and **Step 2** is implemented. \square

C. Proof of Theorem 1

Part a. Following *Lemma 2*, $g(s_{1(q)}^t, s_{2(q)}^t)$ converges as $t \rightarrow \infty$, indicating that:

$$\lim_{t \rightarrow \infty} g(s_{1(q)}^t, s_{2(q)}^t) - g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1}) \rightarrow 0. \tag{S11}$$

With (20), when (S1) is fulfilled, the upper-bound of $g(s_{1(q)}^t, s_{2(q)}^t) - g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1})$ is zero as $t \rightarrow \infty$, thereby achieving (22). Following (S8) and (22), we have [24]:

$$\lim_{t \rightarrow \infty} (p_{u,k(q)}^t - p_{u,k(q)}^{t-1}) \rightarrow 0, \tag{S12a}$$

$$\lim_{t \rightarrow \infty} (z_{i,k(q)}^t - z_{i,k(q)}^{t-1}) \rightarrow 0, \tag{S12b}$$

$$\lim_{t \rightarrow \infty} (a_{u,k(q)}^t - a_{u,k(q)}^{t-1}) \rightarrow 0, \tag{S12c}$$

$$\lim_{t \rightarrow \infty} (x_{i,k(q)}^t - x_{i,k(q)}^{t-1}) \rightarrow 0, \tag{S12d}$$

$$\lim_{t \rightarrow \infty} \left(\left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t p_{u,k(q)}^t - \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} p_{u,k(q)}^{t-1} \right)^2 + \left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t a_{u,k(q)}^t - \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} a_{u,k(q)}^{t-1} \right)^2 \right) \rightarrow 0, \tag{S12e}$$

$$\lim_{t \rightarrow \infty} \left(\left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t z_{i,k(q)}^t - \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} z_{i,k(q)}^{t-1} \right)^2 + \left(\left(\eta_{(q)}^t - 1 \right) \lambda_{(q)}^t x_{i,k(q)}^t - \left(\eta_{(q)}^{t-1} - 1 \right) \lambda_{(q)}^{t-1} x_{i,k(q)}^{t-1} \right)^2 \right) \rightarrow 0. \tag{S12f}$$

Based on (17), (18) and (S12), we have the following inferences:

$$\lim_{t \rightarrow \infty} (h_{u,k(q)}^t - h_{u,k(q)}^{t-1}) \rightarrow 0, \tag{S13a}$$

$$\lim_{t \rightarrow \infty} (w_{i,k(q)}^t - w_{i,k(q)}^{t-1}) \rightarrow 0. \tag{S13b}$$

Based on (4e), (4f) and (S13), we conclude that (23) is fulfilled.

Part b. Firstly, following (4a), (4b) and (5a), the update rules of $(p_{u,k(q)}, z_{i,k(q)})$ can be rearranged as:

$$\left(p_{u,k(q)}^{t-1} - p_{u,k(q)}^t\right) \left(\sum_{i \in \Lambda(u)} \left(z_{i,k(q)}^{t-1}\right)^2 + \lambda_{(q)}^t |\Lambda(u)| \right) \quad (\text{S14a})$$

$$= \sum_{i \in \Lambda(u)} z_{i,k(q)}^{t-1} \left(\sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t + \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} - y_{u,i} \right) + \lambda_{(q)}^t |\Lambda(u)| \left(p_{u,k(q)}^{t-1} - a_{u,k(q)}^{t-1} \right) + h_{u,k(q)}^{t-1},$$

$$\left(z_{i,k(q)}^{t-1} - z_{i,k(q)}^t\right) \left(\sum_{u \in \Lambda(i)} \left(p_{i,k(q)}^{t-1}\right)^2 + \lambda_{(q)}^t |\Lambda(i)| \right) \quad (\text{S14b})$$

$$= \sum_{u \in \Lambda(i)} p_{u,k(q)}^{t-1} \left(\sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t + \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} - y_{u,i} \right) + \lambda_{(q)}^t |\Lambda(i)| \left(z_{i,k(q)}^{t-1} - x_{i,k(q)}^{t-1} \right) + w_{i,k(q)}^{t-1}.$$

Then by substituting (23) and (S12) into (S14), we have:

$$\sum_{i \in \Lambda(u)} z_{i,k(q)}^{t-1} \left(\sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t + \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} - y_{u,i} \right) + h_{u,k(q)}^{t-1} \rightarrow 0, \quad (\text{S15a})$$

$$\sum_{u \in \Lambda(i)} p_{u,k(q)}^{t-1} \left(\sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t + \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} - y_{u,i} \right) + w_{i,k(q)}^{t-1} \rightarrow 0. \quad (\text{S15b})$$

Hence, considering a limit point $\{s_{1(q)}^*, s_{2(q)}^*\}$ of a sequence $\{s_{1(q)}^t, s_{2(q)}^t\}$ generated by the update rules of $\{s_{1(q)}, s_{2(q)}\}$ based on (4) and (5), the following KKT conditions are satisfied with (23) and (S15):

$$\sum_{i \in \Lambda(u)} z_{i,k(q)}^* \left(\sum_{k=1}^d p_{u,k(q)}^* z_{i,k(q)}^* - y_{u,i} \right) + h_{u,k(q)}^* \rightarrow 0, \quad (\text{S16a})$$

$$\sum_{u \in \Lambda(i)} p_{u,k(q)}^* \left(\sum_{k=1}^d p_{u,k(q)}^* z_{i,k(q)}^* - y_{u,i} \right) + w_{i,k(q)}^* \rightarrow 0, \quad (\text{S16b})$$

$$p_{u,k}^* - a_{u,k}^* \rightarrow 0, \quad (\text{S16c})$$

$$z_{i,k}^* - x_{i,k}^* \rightarrow 0. \quad (\text{S16d})$$

Afterwards, considering the remaining KKT conditions regarding constraints $a_{u,k(q)} > 0$ and $x_{i,k(q)} > 0$, we extend the original augmented Lagrangian g :

$$g_{(q)}^\# = g_{(q)} - \text{Tr} \left(M_{(q)} \left(A_{(q)} \right)^T \right) - \text{Tr} \left(N_{(q)} \left(X_{(q)} \right)^T \right) = g_{(q)} - \sum_{(u,k)} m_{u,k(q)} a_{u,k(q)} - \sum_{(i,k)} n_{i,k(q)} x_{i,k(q)}, \quad (\text{S17})$$

where the operator $\text{Tr}(\cdot)$ computes the trace of an enclosed matrix, and the definition of $g_{(q)}$ is given by:

$$g_{(q)} = \frac{1}{2} \sum_{y_{u,i} \in \Lambda} \left(y_{u,i} - \sum_{k=1}^d p_{u,k(q)} z_{i,k(q)} \right)^2 + \sum_{(u,k)} h_{u,k(q)} \left(p_{u,k(q)} - a_{u,k(q)} \right) + \sum_{(u,k)} \frac{\lambda_{(q)} |\Lambda(u)|}{2} \left(p_{u,k(q)} - a_{u,k(q)} \right)^2 \quad (\text{S18})$$

$$+ \sum_{(i,k)} w_{i,k(q)} \left(z_{i,k(q)} - x_{i,k(q)} \right) + \sum_{(i,k)} \frac{\lambda_{(q)} |\Lambda(i)|}{2} \left(z_{i,k(q)} - x_{i,k(q)} \right)^2.$$

For the partial derivatives of $g_{(q)}^\#$ with $a_{u,k(q)}$ and $x_{i,k(q)}$, we have:

$$\begin{cases} \frac{\partial g_{(q)}^\#}{\partial a_{u,k}} = -\lambda_{(q)} |\Lambda(u)| \left(p_{u,k(q)} - a_{u,k(q)} + \frac{h_{u,k(q)}}{\lambda_{(q)} |\Lambda(u)|} \right) - m_{u,k(q)} = 0, \\ \frac{\partial g_{(q)}^\#}{\partial x_{i,k}} = -\lambda_{(q)} |\Lambda(i)| \left(z_{i,k(q)} - x_{i,k(q)} + \frac{w_{i,k(q)}}{\lambda_{(q)} |\Lambda(i)|} \right) - n_{i,k(q)} = 0, \end{cases} \Rightarrow \begin{cases} m_{u,k} = -\lambda_{(q)} |\Lambda(u)| \left(p_{u,k(q)} - a_{u,k(q)} + \frac{h_{u,k(q)}}{\lambda_{(q)} |\Lambda(u)|} \right), \\ n_{i,k} = -\lambda_{(q)} |\Lambda(i)| \left(z_{i,k(q)} - x_{i,k(q)} + \frac{w_{i,k(q)}}{\lambda_{(q)} |\Lambda(i)|} \right). \end{cases} \quad (\text{S19})$$

Then, with the KKT conditions of $\forall m_{u,k(q)}, a_{u,k(q)}: m_{u,k(q)} a_{u,k(q)} = 0$ and $\forall n_{i,k(q)}, x_{i,k(q)}: n_{i,k(q)} x_{i,k(q)} = 0$ for (S17), we achieve the following equations based on (S19) [21, 24]:

$$\begin{cases} a_{u,k(q)} \left(-\lambda_{(q)} |\Lambda(u)| \left(p_{u,k(q)} - a_{u,k(q)} + \frac{h_{u,k(q)}}{\lambda_{(q)} |\Lambda(u)|} \right) \right) = 0, \\ x_{i,k(q)} \left(-\lambda_{(q)} |\Lambda(i)| \left(z_{i,k(q)} - x_{i,k(q)} + \frac{w_{i,k(q)}}{\lambda_{(q)} |\Lambda(i)|} \right) \right) = 0, \end{cases} \Rightarrow \begin{cases} a_{u,k(q)} = p_{u,k(q)} + \frac{h_{u,k(q)}}{\lambda_{(q)} |\Lambda(u)|}, \\ x_{i,k(q)} = z_{i,k(q)} + \frac{w_{i,k(q)}}{\lambda_{(q)} |\Lambda(i)|}. \end{cases} \quad (\text{S20})$$

To satisfy the nonnegativity of output LFs $a_{u,k(q)}$ and $x_{i,k(q)}$, (S20) can be rewritten as:

$$\begin{cases} a_{u,k(q)} = \max \left(0, p_{u,k(q)} + \frac{h_{u,k(q)}}{\lambda_{(q)} |\Lambda(u)|} \right), \\ x_{i,k(q)} = \max \left(0, z_{i,k(q)} + \frac{w_{i,k(q)}}{\lambda_{(q)} |\Lambda(i)|} \right). \end{cases} \quad (\text{S21})$$

Note that (S21) is consistent with the update rules of $a_{u,k(q)}$ and $x_{i,k(q)}$ based on (4c) and (4d). Therefore, (S17)-(S21) show that an A²NLF model's learning rules are closely connected with the KKT conditions of its learning objective.

Then considering the KKT conditions related to $a_{u,k(q)}$:

$$\left. \frac{\partial g_{(q)}^\#}{\partial a_{u,k(q)}} \right|_{a_{u,k(q)} = a_{u,k(q)}^*} = -\lambda_{(q)}^* |\Lambda(u)| \left(p_{u,k(q)}^* - a_{u,k(q)}^* + \frac{h_{u,k(q)}^*}{\lambda_{(q)}^* |\Lambda(u)|} \right) - m_{u,k(q)}^* = 0, \quad (\text{S22a})$$

$$m_{u,k(q)}^* a_{u,k(q)}^* = 0, \quad (\text{S22b})$$

$$a_{u,k(q)}^* \geq 0, \quad (\text{S22c})$$

$$m_{u,k(q)}^* \geq 0, \quad (\text{S22d})$$

where $a_{u,k(q)}^*$ is a KKT stationary point of $a_{u,k(q)}$, and $m_{u,k(q)}^*$ is a limit point of the sequence $\{m_{u,k(q)}^t\}$ generated by the update rules of $m_{u,k}$ based on (S19). According to (S17)-(S21) and $a_{u,k(q)}^t = 0$, conditions (S22a)-(S22c) are satisfied. Thus, we have:

$$m_{u,k(q)}^* = -\lambda_{(q)}^* |\Lambda(u)| \left(p_{u,k(q)}^* - a_{u,k(q)}^* + \frac{h_{u,k(q)}^*}{\lambda_{(q)}^* |\Lambda(u)|} \right). \quad (\text{S23})$$

Thus, we focus on condition (S22d). Since $a_{u,k(q)}^t > 0$ in this case, the update rule for $a_{u,k(q)}$ is given as:

$$a_{u,k(q)}^* \leftarrow p_{u,k(q)}^* + \frac{h_{u,k(q)}^*}{\lambda_{(q)}^* |\Lambda(u)|}. \quad (\text{S24})$$

By substituting (S24) into (S23), we have $m_{u,k(q)}^* = 0$. Hence, conditions (S22c) and (S22d) are fulfilled. Note that as $x_{i,k(q)}^t > 0$ in this case, the proof regarding the KKT conditions of $x_{u,k(q)}$ can be achieved similarly. *Theorem 1* stands, and **Step 3** is implemented. \square

D. Proof of Lemma 3

The difference between $g(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1})$ and $g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1})$ in this case is also given by (S5). Considering the fact of $a_{u,k(q)}^t=0$ and, $x_{i,k(q)}^t>0$ the difference between $g(s_{1(q)}^t, s_{2(q)}^{t-1})$ and $g(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1})$ is:

$$g(s_{1(q)}^t, s_{2(q)}^{t-1}) - g(p_{u,k(q)}^t, z_{i,k(q)}^t, a_{u,k(q)}^{t-1}, x_{i,k(q)}^{t-1}, s_{2(q)}^{t-1}) = -\frac{\lambda_{(q)}^t |\Lambda(u)|}{2} (a_{u,k(q)}^{t-1})^2 - \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} (x_{i,k(q)}^t - x_{i,k(q)}^{t-1})^2. \quad (\text{S25})$$

Moreover, $g(s_{1(q)}^t, s_{2(q)}^t)$ and $g(s_{1(q)}^t, s_{2(q)}^{t-1})$ yields:

$$\begin{aligned} & g(s_{1(q)}^t, s_{2(q)}^t) - g(s_{1(q)}^t, s_{2(q)}^{t-1}) \\ &= (p_{u,k(q)}^t - a_{u,k(q)}^t)(h_{u,k(q)}^t - h_{u,k(q)}^{t-1}) + (z_{i,k(q)}^t - x_{i,k(q)}^t)(w_{i,k(q)}^t - w_{i,k(q)}^{t-1}) \\ &\stackrel{(I)}{=} \eta_{(q)}^t \lambda_{(q)}^t |\Lambda(u)| (p_{u,k(q)}^t - a_{u,k(q)}^t)^2 + \frac{1}{\eta_{(q)}^t \lambda_{(q)}^t |\Lambda(i)|} (w_{i,k(q)}^t - w_{i,k(q)}^{t-1})^2 \\ &\stackrel{(II)}{\leq} \eta_{(q)}^t \lambda_{(q)}^t |\Lambda(u)| (p_{u,k(q)}^t)^2 + \frac{2|\Lambda(i)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(((\eta_{(q)}^t - 1) \lambda_{(q)}^t z_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} z_{i,k(q)}^{t-1})^2 + ((\eta_{(q)}^t - 1) \lambda_{(q)}^t x_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} x_{i,k(q)}^{t-1})^2 \right), \end{aligned} \quad (\text{S26})$$

where (I) is based on the update rules of $(h_{u,k(q)}, w_{i,k(q)})$ given in (4e), (4f) and (5c), and (II) is achieved with (18) and $a_{u,k(q)}^t=0$.

With (S5), (S25) and (S26), we have the following deduction:

$$\begin{aligned} & g(s_{1(q)}^t, s_{2(q)}^t) - g(s_{1(q)}^{t-1}, s_{2(q)}^{t-1}) \\ &\leq -\frac{1}{2} \left(\sum_{i \in \Lambda(u)} (z_{i,k(q)}^{t-1})^2 + \lambda_{(q)}^t |\Lambda(u)| \right) (p_{u,k(q)}^t - p_{u,k(q)}^{t-1})^2 - \frac{\lambda_{(q)}^t |\Lambda(u)|}{2} (a_{u,k(q)}^{t-1})^2 \\ &\quad - \frac{1}{2} \left(\sum_{u \in \Lambda(i)} (p_{u,k(q)}^{t-1})^2 + \lambda_{(q)}^t |\Lambda(i)| \right) (z_{i,k(q)}^t - z_{i,k(q)}^{t-1})^2 - \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} (x_{i,k(q)}^t - x_{i,k(q)}^{t-1})^2 \\ &\quad + \eta_{(q)}^t \lambda_{(q)}^t |\Lambda(u)| (p_{u,k(q)}^t)^2 + \frac{2|\Lambda(i)|}{\eta_{(q)}^t \lambda_{(q)}^t} \left(((\eta_{(q)}^t - 1) \lambda_{(q)}^t z_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} z_{i,k(q)}^{t-1})^2 + ((\eta_{(q)}^t - 1) \lambda_{(q)}^t x_{i,k(q)}^t - (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} x_{i,k(q)}^{t-1})^2 \right). \end{aligned} \quad (\text{S27})$$

Owing to (25), (21a) stands, which indicates that the augmented Lagrangian function (3) related to the q -th particle is non-increasing as $a_{u,k(q)}^t=0$ and $x_{i,k(q)}^t>0$ in this case. Then after the t -th iteration, we substitute $a_{u,k(q)}^t=0$ into (S10):

$$\begin{aligned} g(s_{1(q)}^t, s_{2(q)}^t) &= \frac{1}{2} \sum_{y_{u,i} \in \Lambda} \left(y_{u,i} - \sum_{l_1=1}^k p_{u,l_1(q)}^t z_{i,l_1(q)}^t - \sum_{l_2=k+1}^d p_{u,l_2(q)}^{t-1} z_{i,l_2(q)}^{t-1} \right)^2 \\ &\quad + \sum_u |\Lambda(u)| \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t \left(\sum_{l_1=1}^{k-1} (p_{u,l_1(q)}^t - a_{u,l_1(q)}^t)^2 + (p_{u,k(q)}^t)^2 \right) + (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} \sum_{l_2=k+1}^d (p_{u,l_2(q)}^{t-1} - a_{u,l_2(q)}^{t-1})^2 \right) \\ &\quad + \sum_i |\Lambda(i)| \left((\eta_{(q)}^t - 1) \lambda_{(q)}^t \sum_{l_1=1}^k (z_{i,l_1(q)}^t - x_{i,l_1(q)}^t)^2 + (\eta_{(q)}^{t-1} - 1) \lambda_{(q)}^{t-1} \sum_{l_2=k+1}^d w_{i,l_2(q)}^t (z_{i,l_2(q)}^{t-1} - x_{i,l_2(q)}^{t-1})^2 \right) \\ &\quad + \sum_u \frac{\lambda_{(q)}^t |\Lambda(u)|}{2} \left(\sum_{l_1=1}^{k-1} (p_{u,l_1(q)}^t - a_{u,l_1(q)}^t)^2 + (p_{u,k(q)}^t)^2 + \sum_{l_2=k+1}^d (p_{u,l_2(q)}^{t-1} - a_{u,l_2(q)}^{t-1})^2 \right) \\ &\quad + \sum_i \frac{\lambda_{(q)}^t |\Lambda(i)|}{2} \left(\sum_{l_1=1}^k (z_{i,l_1(q)}^t - x_{i,l_1(q)}^t)^2 + \sum_{l_2=k+1}^d (z_{i,l_2(q)}^{t-1} - x_{i,l_2(q)}^{t-1})^2 \right). \end{aligned} \quad (\text{S28})$$

(S28) indicates that if (20b) is fulfilled, (21b) holds, thereby making (3) related to the q -th particle lower-bounded as $a_{u,k(q)}^t=0$, and $x_{i,k(q)}^t>0$ in this case. Based on the above inferences, *Lemma 2* stands, and **Step 4** is implemented. \square

E. Proof of Theorem 2

Part a. Following *Lemma 3*, $g\left(s_{1(q)}^t, s_{2(q)}^t\right)$ converges as $t \rightarrow \infty$, indicating that (S11) is fulfilled. With (22), (25) and (S27), we have (S12a), (S12b), (S12d), (S12f) and the following inferences:

$$\lim_{t \rightarrow \infty} a_{u,k(q)}^{t-1} \rightarrow 0, \quad (\text{S29a})$$

$$\lim_{t \rightarrow \infty} p_{u,k(q)}^t \rightarrow 0. \quad (\text{S29b})$$

Then according to (S12f), (S13b) is fulfilled. Hence, based on (S13b), (S29b) and $a_{u,k(q)}^t=0$, (23) is fulfilled.

Part b. Firstly, considering a limit point $\left\{s_{1(q)}^*, s_{2(q)}^*\right\}$ of a sequence $\left\{s_{1(q)}^t, s_{2(q)}^t\right\}$ generated by the update rules of $\{s_{1(q)}, s_{2(q)}\}$ based on (4) and (5), according to (23) and (S15), (S16) holds when $a_{u,k(q)}^t=0$, and $x_{i,k(q)}^t > 0$ in this case. Then considering the KKT conditions related to $a_{u,k(q)}$, i.e., (S22).

According to (S17)-(S21) and with $a_{u,k(q)}^t=0$, conditions (S22a)-(S22c) are naturally satisfied. Thus, we focus on analyzing condition (S22d). Since we have $a_{u,k(q)}^t=0$ in this case, the following inequality holds:

$$p_{u,k(q)}^* + \frac{h_{u,k(q)}^*}{\lambda_{(q)}^* |\Lambda(u)|} \leq 0. \quad (\text{S30})$$

Note that (S30) indicates that:

$$m_{u,k(q)}^* = -\lambda_{(q)}^* |\Lambda(u)| \left(p_{u,k(q)}^* + \frac{h_{u,k(q)}^*}{\lambda_{(q)}^* |\Lambda(u)|} \right) \geq 0. \quad (\text{S31})$$

Thus, condition (S22) are all fulfilled in this case. Note that as $x_{i,k(q)}^t > 0$ in this case, the proof regarding the KKT conditions of $x_{u,k(q)}$ can be achieved similarly. *Theorem 2* stands, and **Step 5** is implemented. ■

II. SUPPLEMENTARY TABLES

TABLE S1. PSO Parameters Settings in (7).

w	c_1	c_2	r_1, r_2
0.729	2	2	uniform random numbers $\in [0,1]$.

TABLE S2 PSO PARAMETERS SETTINGS IN (13)

$\tilde{\lambda}$	$\hat{\lambda}$	$\tilde{\eta}$	$\hat{\eta}$	\tilde{v}_λ	\hat{v}_λ	\tilde{v}_η	\hat{v}_η
0.2	2	1	2	$-\tilde{v}_\lambda$	$0.2 \times (\hat{\lambda} - \tilde{\lambda})$	$-\tilde{v}_\eta$	$0.2 \times (\hat{\eta} - \tilde{\eta})$

TABLE S3. Grid-search Range of Hyper-parameters in M1-6 on D1-4.

No.	Hyper-parameters	Grid-searching Range
M1	Augmentation coefficient λ	$2^{-6}, 2^{-4}, 2^{-2}, 1, 2^2, 2^4, 2^6, 2^8$
	Learning rate η	$8, 4, 2, 1, 2^{-1}, 2^{-2}, 2^{-3}$
M2	Penalty coefficient λ	$2, 1.8, 1.6, 1.4, 1.2, 1.0, 0.8$
M3	Regularization coefficient μ	$8, 4, 2, 1, 2^{-1}, 2^{-2}, 2^{-3}$
	L_p coefficient p	$1.2, 1.3, 1.4, 1.5, 1.6, 1.7, 1.8, 1.9, 2.0$
M4	Penalty coefficient α	$1, 10, 10^2, 10^3, 10^4, 10^5$
	Penalty coefficient β	$10, 10^2, 10^3, 10^4, 10^5, 10^6$
M5	Penalty coefficient ρ	$2^{-3}, 2^{-5}, 2^{-7}, 2^{-9}, 2^{-11}, 2^{-13}, 2^{-15}, 2^{-17}$
M6	Regularization coefficient λ	$10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1$
	Learning rate η	$10^{-5}, 5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}$
	Batch size bs	$32, 64, 128, 256, 512, 1024$
M7	Regularization coefficient λ	$10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1$
	Learning rate η	$10^{-5}, 5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}$
	Batch size bs	$32, 64, 128, 256, 512, 1024$

D4	MAE	$\lambda=10^{-5}$	$\lambda=10^{-4}$	$\lambda=10^{-5}$	$\lambda=10^{-5}$	$\lambda=10^{-4}$	$\lambda=10^{-4}$	$\lambda=10^{-4}$	$\lambda=10^{-4}$	$\lambda=10^{-4}$
		$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$
		bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128
	RMSE	$\lambda=10^{-3}$	$\lambda=10^{-4}$	$\lambda=10^{-5}$	$\lambda=10^{-4}$	$\lambda=10^{-4}$	$\lambda=10^{-3}$	$\lambda=10^{-5}$	$\lambda=10^{-3}$	$\lambda=10^{-5}$
		$\eta=10^{-4}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$
		bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512
D4	MAE	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-5}$	$\lambda=10^{-2}$	$\lambda=10^{-5}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$
		$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$
		bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512

bs denotes batch size adopted by M6 on an HDI matrix;

TABLE S10. Optimal Hyper-parameters during M7's ten times' training process on D1-4

No.	Type	1	2	3	4	5	6	7	8	9	10
D1	RMSE	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$
		$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-4}$	$\eta=5 \times 10^{-4}$	$\eta=10^{-4}$
		bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512
	MAE	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$	$\lambda=10^{-2}$
		$\eta=10^{-4}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$	$\eta=10^{-4}$	$\eta=10^{-4}$	$\eta=10^{-4}$	$\eta=10^{-4}$
		bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512	bs=512
D2	RMSE	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$
		$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-4}$
		bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64
	MAE	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$	$\lambda=10^{-3}$
		$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$
		bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64	bs=64
D3	RMSE	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$
		$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$
		bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128
	MAE	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$
		$\eta=10^{-3}$	$\eta=10^{-4}$	$\eta=10^{-4}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-4}$	$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=10^{-3}$
		bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128	bs=128
D4	RMSE	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	-	-	-	-	-	-	-
		$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	$\eta=5 \times 10^{-4}$	-	-	-	-	-	-	-
		bs=512	bs=512	bs=512	-	-	-	-	-	-	-
	MAE	$\lambda=0.1$	$\lambda=0.1$	$\lambda=0.1$	-	-	-	-	-	-	-
		$\eta=10^{-3}$	$\eta=10^{-3}$	$\eta=5 \times 10^{-3}$	-	-	-	-	-	-	-
		bs=512	bs=512	bs=512	-	-	-	-	-	-	-

Note that M7 need to consume about 18 days for manually implementing the tuning process of its hyper-parameters on D4. Hence, only third times validation process for M7 on D4 is completed; bs denotes batch size adopted by M7 on an HDI matrix.

TABLE S11. RMSE and Time Cost of M1 and M8 on D1-4.

Dataset	Model	Prediction Accuracy		Tuning Time Cost (Secs)*	Testing Time Cost (Secs)**	Total Time Cost (Secs)
D1	M1	RMSE	0.2373 \pm 2.2E-6	428 \pm 22.7	6 \pm 2.4	434 \pm 25
		MAE	0.1815 \pm 1.1E-6	439 \pm 25.4	6 \pm 2.8	445 \pm 28
	M8	RMSE	0.2339 \pm2.7E-4	-	-	46 \pm4
		MAE	0.1792 \pm3.3E-4	-	-	60 \pm5
D2	M1	RMSE	1.0187 \pm 1.1E-6	271 \pm 45.4	4 \pm 1.7	275 \pm 47
		MAE	0.8079 \pm 2.5E-6	305 \pm 38.3	3 \pm 0.9	308 \pm 39
	M8	RMSE	1.0172 \pm7.4E-4	-	-	25 \pm5
		MAE	0.7856 \pm7.9E-4	-	-	29 \pm8
D3	M1	RMSE	0.8665 \pm7.8E-4	739 \pm 72.3	21 \pm 14.7	759 \pm 87
		MAE	0.6829 \pm 1.7E-6	756 \pm 53.8	3 \pm 0.5	763 \pm 54
	M8	RMSE	0.8675 \pm 8.2E-4	-	-	26 \pm6
		MAE	0.6787 \pm7.3E-4	-	-	30 \pm9
D4	M1	RMSE	0.8096 \pm 2.9E-6	2,934 \pm 353.8	38 \pm 24.6	2,972 \pm 378
		MAE	0.6221 \pm 7.9E-7	9,203 \pm 828.9	33 \pm 20.4	9,236 \pm 849
	M8	RMSE	0.8089 \pm9.4E-4	-	-	334 \pm46
		MAE	0.6193 \pm5.1E-4	-	-	358 \pm31

* Time cost consumed by M1 for manually grid-searching optimal hyper-parameters;

** Time cost consumed by M1 with achieved hyper-parameters.

TABLE S12. Storage Complexity and Memory Requirements of M1-8 on D1-4.

No.	Storage Complexity	Memory Requirements (GB)			
		D1	D2	D3	D4
M1	$\Theta((U + I) \times d + \Lambda)$	0.7	0.6	0.5	1.7
M2	$\Theta(U \times I)$	17.6	1.8	9.3	>256
M3	$\Theta(U \times I)$	15.4	1.9	8.6	>256
M4	$\Theta(U \times I)$	21.9	3.5	11.5	>256
M5	$\Theta(U \times I)$	26.5	3.9	12.7	>256
M6	$\Theta(\text{bs}^* \times U + \Lambda)$	2.9	1.7	1.7	3.1
M7	$\Theta((U + I) \times d + \text{bs} \times 1^{\text{st}}\text{LN}^{**} + \Lambda)$	2.1	1.2	1.8	6.5
M8	$\Theta((U + I) \times d + \Lambda)$	0.7	0.6	0.5	1.7

* bs denotes batch size adopted by M6 and M7 on an HDI matrix;

** 1stLN denotes the number of neuron in the 1st layer and is set at $2 \times d$ according to [10].

TABLE S13. RMSE/MAE of M1-8 on D1-4, including Win/Loss counts and Friedman Rank, where ● indicates that M8 has higher RMSE/MAE than the rival models

No.	CASE	M1	M2	M3	M4	M5	M6	M7	M8
D1	RMSE	0.2373 _{±2.2E-6}	0.3058 _{±5.0E-4}	0.3047 _{±4.3E-5}	0.2384 _{±1.0E-4}	0.4913 _{±1.0E-4}	●0.2302 _{±2.6E-3}	0.2354 _{±2.3E-3}	0.2339 _{±2.7E-4}
	MAE	0.1815 _{±1.1E-6}	0.2439 _{±4.4E-5}	0.2422 _{±3.2E-5}	0.1832 _{±4.8E-5}	0.4111 _{±1.6E-2}	●0.1792 _{±7.9E-5}	0.1842 _{±2.2E-4}	0.1793 _{±3.3E-4}
D2	RMSE	1.0187 _{±1.1E-6}	1.1281 _{±7.2E-3}	1.1257 _{±1.8E-4}	1.0787 _{±8.6E-7}	1.8808 _{±2.8E-2}	1.1494 _{±4.1E-5}	1.0425 _{±3.9E-4}	1.0172 _{±7.4E-4}
	MAE	0.8079 _{±2.5E-6}	0.9256 _{±1.5E-3}	0.9229 _{±1.9E-4}	0.8580 _{±6.2E-6}	1.5403 _{±2.6E-2}	0.9257 _{±2.1E-4}	0.8014 _{±4.9E-2}	0.7856 _{±7.9E-4}
D3	RMSE	●0.8665 _{±7.8E-4}	1.0336 _{±5.4E-4}	1.0848 _{±5.3E-5}	0.8713 _{±3.7E-4}	2.2963 _{±9.8E-3}	0.8845 _{±1.1E-3}	0.8982 _{±4.4E-4}	0.8675 _{±8.2E-4}
	MAE	0.6829 _{±1.7E-6}	0.8832 _{±4.0E-4}	0.9021 _{±4.2E-5}	0.6802 _{±3.3E-4}	1.9190 _{±9.6E-3}	0.7021 _{±5.5E-3}	0.7067 _{±4.1E-4}	0.6787 _{±7.3E-4}
D4	RMSE	0.8096 _{±2.9E-6}	¹ Failure	¹ Failure	¹ Failure	¹ Failure	0.8436 _{±1.2E-3}	^a 0.8657 _{±3.5E-4}	0.8089 _{±9.4E-4}
	MAE	0.6221 _{±7.9E-7}	¹ Failure	¹ Failure	¹ Failure	¹ Failure	0.6530 _{±4.9E-3}	^a 0.6650 _{±2.9E-4}	0.6193 _{±5.1E-4}
Win/Loss		7/1	8/0	8/0	8/0	8/0	6/2	8/0	-
Friedman Rank		2.5	6.38	6.13	4.38	7.63	3.75	3.88	1.38

¹ Note that M2-M5 fails to achieve the final results on D4: their memory requirements are too large to meet on our experimental environment as shown in Table S12.

^a Note that M7 need to consume about 18 days for manually implementing the tuning process of its hyper-parameters on D4. Hence, only third times validation process for M7 on D4 is completed, and the corresponding performance is marked as '^a' on D4.

TABLE S14. Total time cost of M1-8 in RMSE/MAE on D1-4 (Secs), including Win/Loss counts and Friedman Rank

No.	CASE	M1	M2	M3	M4	M5	M6	M7	M8
D1	RMSE	434 _{±25}	38,499 _{±2,332}	43,271 _{±2,962}	209,670 _{±36,625}	3,856 _{±578}	10,804 _{±832}	87,425 _{±9,048}	46 _{±4}
	MAE	445 _{±28}	51,443 _{±2,643}	189,090 _{±29,135}	210,745 _{±27,190}	3,912 _{±103}	13,432 _{±2,219}	78,395 _{±8,698}	60 _{±5}
D2	RMSE	275 _{±47}	236 _{±35}	128 _{±21}	138 _{±39}	55 _{±7}	2,109 _{±295}	1,264 _{±99}	25 _{±5}
	MAE	308 _{±39}	32 _{±3}	128 _{±27}	766 _{±64}	56 _{±6}	2,256 _{±633}	1,156 _{±286}	29 _{±8}
D3	RMSE	759 _{±87}	42,578 _{±3,258}	71,927 _{±2,718}	45,849 _{±2,479}	717 _{±22}	8,359 _{±532}	9,921 _{±2,586}	26 _{±6}
	MAE	763 _{±54}	1,906 _{±291}	12,294 _{±2,478}	48,821 _{±2,363}	716 _{±53}	7,700 _{±219}	8,998 _{±1,446}	30 _{±9}
D4	RMSE	2,972 _{±378}	¹ Failure	¹ Failure	¹ Failure	¹ Failure	266,885 _{±12,775}	^a 1,644,502 _{±98,039}	334 _{±46}
	MAE	9,236 _{±849}	¹ Failure	¹ Failure	¹ Failure	¹ Failure	232,974 _{±13,157}	^a 997,259 _{±84,311}	358 _{±31}
Win/Loss		8/0	8/0	8/0	8/0	8/0	8/0	8/0	8/0
Friedman Rank		3.13	5	6	6.75	3.5	4.88	5.75	1

¹ It means the same with that in Table S13;

^a It means the same with that in Table S13.

TABLE S15. Results of the Wilcoxon signed-ranks test on RMSE/MAE of Table S13.

Comparison	<i>R</i> +	<i>R</i> -	<i>p-value</i> *
M8 vs M1	34	2	0.0117
M8 vs M2	36	0	0.0039
M8 vs M3	36	0	0.0039
M8 vs M4	36	0	0.0039
M8 vs M5	36	0	0.0039
M8 vs M6	33	3	0.0195
M8 vs M7	36	0	0.0039

*With a significance level of 0.1, the accepted hypotheses are highlighted.

TABLE S16. Results of the Wilcoxon signed-ranks test on time cost of Table S14.

Comparison	<i>R</i> +	<i>R</i> -	<i>p-value</i> *
M8 vs M1	36	0	0.0039
M8 vs M2	36	0	0.0039
M8 vs M3	36	0	0.0039
M8 vs M4	36	0	0.0039
M8 vs M5	36	0	0.0039
M8 vs M6	36	0	0.0039
M8 vs M7	36	0	0.0039

*With a significance level of 0.1, the accepted hypotheses are highlighted.