

Integrating Inference and Experimental Design for Contextual Behavioral Model Learning (Supplemental Material)

Gongtao Zhou, Haoran Yu

A Derivation of EIG

We derive Expected Information Gain (EIG) and its approximation. Recall that Information Gain (IG) is defined as follows:

$$\begin{aligned}\text{IG}(\mathbf{d}, \mathbf{y}) &\triangleq H[p(\boldsymbol{\theta})] - H[p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d})] \\ &= \int p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) \log p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} - \int p(\boldsymbol{\theta}) \log p(\boldsymbol{\theta}) d\boldsymbol{\theta}.\end{aligned}$$

Information Gain captures the change in information entropy under the conditions where the platform selects \mathbf{d} and the investor chooses \mathbf{y} . The investor's choice \mathbf{y} is determined by the context \mathbf{X} and the design \mathbf{d} . Given \mathbf{X} and \mathbf{d} , its distribution is given by $p(\mathbf{y}|\mathbf{X}, \mathbf{d}) = \mathbb{E}_{p(\boldsymbol{\theta})}[p(\mathbf{y}|\mathbf{X}, \mathbf{d}, \boldsymbol{\theta})]$. We can further take the expectation of $\text{IG}(\mathbf{d}, \mathbf{y})$ with respect to the uncertainty of \mathbf{y} to obtain EIG:

$$\begin{aligned}\text{EIG}(\mathbf{d}) &\triangleq \mathbb{E}_{p(\mathbf{y}|\mathbf{X}, \mathbf{d})}[\text{IG}(\mathbf{d}, \mathbf{y})] \\ &= \sum_{\mathbf{y}} \int p(\mathbf{y}|\mathbf{X}, \mathbf{d}) p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) \log p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} \\ &\quad - \sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{X}, \mathbf{d}) \int p(\boldsymbol{\theta}) \log p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \sum_{\mathbf{y}} \int p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}, \mathbf{d}) \log p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} - \int p(\boldsymbol{\theta}) \log p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \sum_{\mathbf{y}} \int p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}, \mathbf{d}) \log p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d}) d\boldsymbol{\theta} - \sum_{\mathbf{y}} \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d}) p(\boldsymbol{\theta}) \log p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &\stackrel{(a)}{=} \sum_{\mathbf{y}} \int p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}, \mathbf{d}) \log \frac{p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \mathbf{d})}{p(\boldsymbol{\theta})} d\boldsymbol{\theta} \\ &\stackrel{(b)}{=} \sum_{\mathbf{y}} \int p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}, \mathbf{d}) \log \frac{p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{y}|\mathbf{X}, \mathbf{d})} d\boldsymbol{\theta} \\ &= \sum_{\mathbf{y}} \int p(\boldsymbol{\theta}) p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d}) \log \frac{p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d})}{p(\mathbf{y}|\mathbf{X}, \mathbf{d})} d\boldsymbol{\theta}.\end{aligned}$$

We use the fact that $\boldsymbol{\theta}$ is independent of \mathbf{X} and \mathbf{d} (i.e., $p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{d})$) to derive equation (a), and use Bayes' theorem to derive equation (b).

Since there is no closed-form solution for $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d})$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{d})$, obtaining the exact computation of the EIG is intractable. Therefore, we approximate the value using the nested Monte Carlo method:

$$\text{EIG}(\mathbf{d}) \approx \widehat{\text{EIG}}(\mathbf{d}) = \frac{1}{I} \sum_{i=1}^I \log \frac{p(\mathbf{y}^i|\mathbf{X}, \mathbf{d}, \boldsymbol{\theta}^{i,0})}{\frac{1}{J} \sum_{j=1}^J p(\mathbf{y}^j|\mathbf{X}, \mathbf{d}, \boldsymbol{\theta}^{i,j})},$$

where $\boldsymbol{\theta}^{i,0} \sim p(\boldsymbol{\theta})$, $\mathbf{y}^i \sim p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta} = \boldsymbol{\theta}^{i,0}, \mathbf{d})$, $\boldsymbol{\theta}^{i,j} \sim p(\boldsymbol{\theta})$. The outer sampling estimates $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \mathbf{d})$ and the inner sampling estimates $p(\mathbf{y}|\mathbf{X}, \mathbf{d})$. Parameters I and J denote the number of samples of the outer and inner loops.

B Updating ϕ^τ in Iterative Optimization

We introduce the optimization of ϕ^τ in the iterative optimization of our **I-ID-LP** method.

$$\begin{aligned} \phi^\tau &= \arg \max_{\phi} \int q(\boldsymbol{\theta}|\phi) \log \frac{p(\boldsymbol{\theta}|\mathcal{H}_{t-1})g(\boldsymbol{\theta}, \mathbf{d}_t^{\tau-1})}{q(\boldsymbol{\theta}|\phi)} d\boldsymbol{\theta} \\ &= \arg \min_{\phi} \int q(\boldsymbol{\theta}|\phi) \log \frac{q(\boldsymbol{\theta}|\phi)}{p(\boldsymbol{\theta}|\mathcal{H}_{t-1})g(\boldsymbol{\theta}, \mathbf{d}_t^{\tau-1})} d\boldsymbol{\theta} \\ &\stackrel{(a)}{=} \arg \min_{\phi} \int q(\boldsymbol{\theta}|\phi) \log \frac{q(\boldsymbol{\theta}|\phi)p(\mathcal{H}_{t-1})}{p(\mathcal{H}_{t-1}|\boldsymbol{\theta})p(\boldsymbol{\theta})} d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}|\phi) \log g(\boldsymbol{\theta}, \mathbf{d}_t^{\tau-1}) d\boldsymbol{\theta} \\ &\stackrel{(b)}{=} \arg \min_{\phi} \int q(\boldsymbol{\theta}|\phi) \log \frac{q(\boldsymbol{\theta}|\phi)}{p(\boldsymbol{\theta})} d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}|\phi) \log p(\mathcal{H}_{t-1}|\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &\quad - \int q(\boldsymbol{\theta}|\phi) \log g(\boldsymbol{\theta}, \mathbf{d}_t^{\tau-1}) d\boldsymbol{\theta} \\ &= \arg \min_{\phi} \text{KL}[q(\boldsymbol{\theta}|\phi)||p(\boldsymbol{\theta})] - \mathbb{E}_{q(\boldsymbol{\theta}|\phi)}[\log p(\mathcal{H}_{t-1}|\boldsymbol{\theta})] - \mathbb{E}_{q(\boldsymbol{\theta}|\phi)}[\log g(\boldsymbol{\theta}, \mathbf{d}_t^{\tau-1})]. \end{aligned}$$

We use the fact that $p(\boldsymbol{\theta}|\mathcal{H}_{t-1}) = \frac{p(\mathcal{H}_{t-1}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{H}_{t-1})}$ to derive equation (a), and the fact that $p(\mathcal{H}_{t-1})$ is independent of $\boldsymbol{\theta}$ to derive equation (b).

C Data Settings

We introduce the remaining four data settings.

Setting B: We use a three-layer neural network to generate data. Specifically, we compute

$$(\bar{\lambda}(\mathbf{x}); \bar{r}(\mathbf{x})) = \text{Sigmoid}(\mathbf{w}_3^T \text{Sigmoid}(\mathbf{w}_2^T \text{Sigmoid}(\mathbf{w}_1^T \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_3),$$

where $\mathbf{w}_1 \in \mathbb{R}^{16 \times 6}$, $\mathbf{w}_2 \in \mathbb{R}^{6 \times 4}$, $\mathbf{w}_3 \in \mathbb{R}^{4 \times 2}$ and $\mathbf{b}_1 \in \mathbb{R}^6$, $\mathbf{b}_2 \in \mathbb{R}^4$, $\mathbf{b}_3 \in \mathbb{R}^2$ are randomly chosen. Then, we normalize $\bar{r}(\mathbf{x}), \bar{\lambda}(\mathbf{x})$ across all data to get $r(\mathbf{x}), \lambda(\mathbf{x})$.

Setting C: We use a two-layer network to generate data. We compute

$$(\bar{\lambda}(\mathbf{x}); \bar{r}(\mathbf{x})) = \text{Sigmoid}(\mathbf{w}_2^T (\text{Sigmoid}(\mathbf{w}_1^T \mathbf{x}) + \mathbf{b}_1) + \mathbf{b}_2 + \boldsymbol{\epsilon}),$$

where $\mathbf{w}_1 \in \mathbb{R}^{16 \times 16}$, $\mathbf{w}_2 \in \mathbb{R}^{16 \times 2}$ and $\mathbf{b}_1 \in \mathbb{R}^{16}$, $\mathbf{b}_2 \in \mathbb{R}^2$ are randomly chosen. $\boldsymbol{\epsilon} \in \mathbb{R}^2$ represents the noise vector. We normalize $\bar{r}(\mathbf{x}), \bar{\lambda}(\mathbf{x})$ across all data to get $r(\mathbf{x}), \lambda(\mathbf{x})$.

Setting D: Setting D employs the same network structure as Setting C to generate data. The difference is that setting D normalizes the values of $\bar{\lambda}(\mathbf{x})$ to a smaller range to get $\lambda(\mathbf{x})$. Table 1 in the paper shows that our methods can more accurately predict investor behavior, compared with Setting C.

Setting E: We generate data according to

$$(\bar{\lambda}(\mathbf{x}); \bar{r}(\mathbf{x})) = \text{Sigmoid}(\mathbf{w}_2^T (\text{Sigmoid}(\mathbf{w}_1^T \mathbf{x}) + \mathbf{b}_1) + \mathbf{b}_2),$$

where $\mathbf{w}_1 \in \mathbb{R}^{16 \times 16}$, $\mathbf{w}_2 \in \mathbb{R}^{16 \times 2}$ and $\mathbf{b}_1 \in \mathbb{R}^{16}$, $\mathbf{b}_2 \in \mathbb{R}^2$ are randomly chosen. Then, we normalize $\bar{r}(\mathbf{x}), \bar{\lambda}(\mathbf{x})$ across all data to get $r(\mathbf{x}), \lambda(\mathbf{x})$. Compared with Setting D, Setting E does not include any noise in the data generation.

D Comparison of ID Methods

We show the comparison of ID methods under data settings C~E as follows.

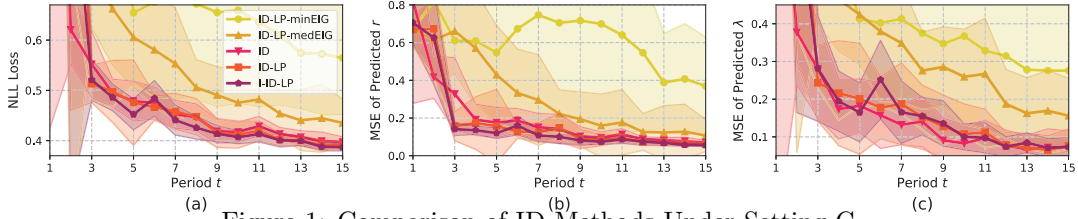


Figure 1: Comparison of ID Methods Under Setting C.

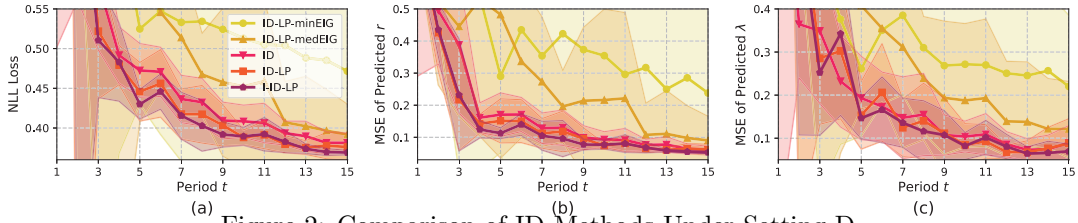


Figure 2: Comparison of ID Methods Under Setting D.

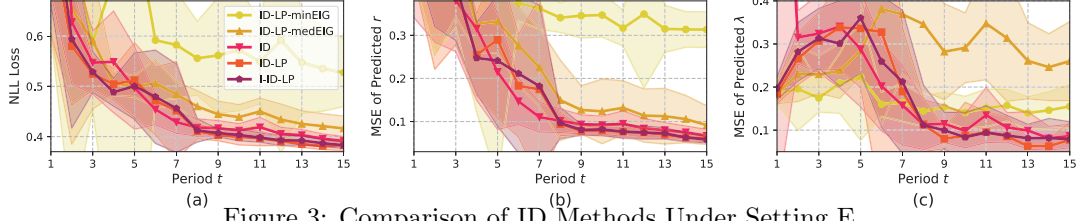


Figure 3: Comparison of ID Methods Under Setting E.

E Mean Squared Errors of Predicted r and λ

We evaluate the performance of predicting r and λ based on the mean squared errors (MSEs). We show the MSEs achieved by different methods under settings A~E as follows.

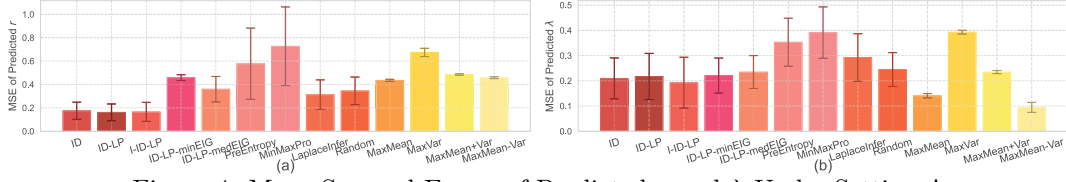


Figure 4: Mean Squared Errors of Predicted r and λ Under Setting A.

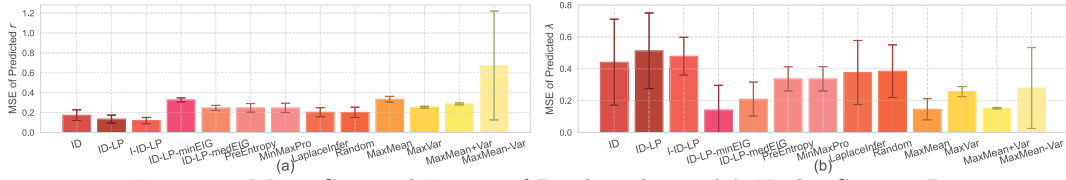


Figure 5: Mean Squared Errors of Predicted r and λ Under Setting B.

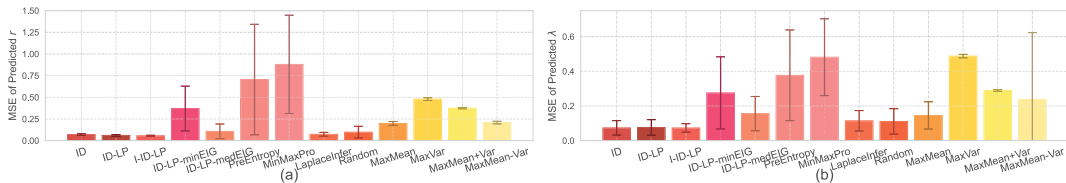


Figure 6: Mean Squared Errors of Predicted r and λ Under Setting C.

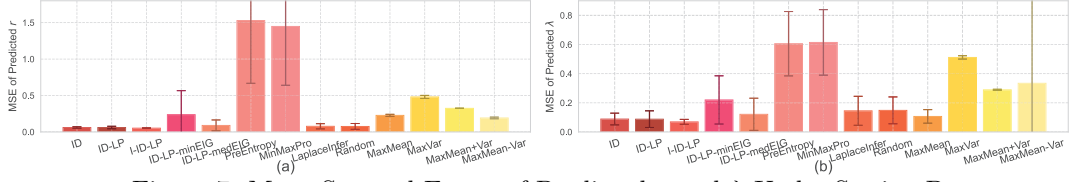


Figure 7: Mean Squared Errors of Predicted r and λ Under Setting D.

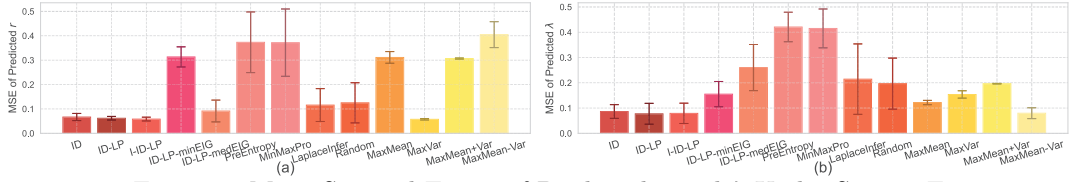


Figure 8: Mean Squared Errors of Predicted r and λ Under Setting E.