

Inferring Heterogeneous Private Valuations from Offline Market Data via Entropic Risk-Sensitive Utility Maximization (Supplementary Material)

Xingyu Qian, Haoran Yu*

School of Computer Science & Technology, Beijing Institute of Technology
3120241000@bit.edu.cn, yhrhawk@gmail.com

1 Comparison Baselines

We introduce the details and pseudocode of five baseline methods: **CE**, **BLUE**, **SL**, **DL**, **PGM**, **RvS**.

1.1 CE

This method addresses valuation inference using supervised learning to train a decision function $\pi_\theta(\cdot|\mathbf{v}, \mathbf{o}^t)$ that maps valuation \mathbf{v} and market observation \mathbf{o}^t to buyer behavior probabilities. Since true valuations are unavailable during training, we substitute \mathbf{v} with samples from feasible ranges $\mathcal{V}_{\text{feas}}$. The inference phase solves an optimization problem to find valuations \mathbf{v}_n^* that maximize the likelihood of observed behaviors.

During the Model Training Phase, we minimize the cross-entropy between predicted and actual behaviors:

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \pi_\theta(b_n^t | \tilde{\mathbf{v}}_n, \mathbf{o}_n^t). \quad (1)$$

where $\tilde{\mathbf{v}}_n$ are feasible valuations sampled from $\mathcal{V}_{n,\text{feas}}$. This trains π_θ to predict behaviors given sampled valuations and market observations.

In the Valuation Inference Phase, for fixed θ , we find \mathbf{v}_n^* that minimizes the negative log-likelihood of observed behaviors:

$$\mathbf{v}_n^* = \arg \min_{\mathbf{v}_n \in \mathcal{V}_{n,\text{feas}}} - \sum_{t=0}^{T-1} \log \pi_\theta(b_n^t | \mathbf{v}_n, \mathbf{o}_n^t). \quad (2)$$

This optimization recovers valuations that best explain observed behaviors under the trained model. The pseudocode is shown in Algorithm 1.

1.2 BLUE

BLUE extends the **CE** method by introducing a rationality constraint during training to ensure economically plausible behavior predictions. This constraint penalizes the decision function when it assigns probability to strictly dominated actions. Specifically, actions should only be plausible for valuations within the buyer's feasible valuation range $\mathcal{V}_{\text{feas}}$.

*Corresponding author.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Algorithm 1: Cross-Entropy Based Valuation Inference

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, learning rates $\eta_{\text{train}}, \eta_{\text{inf}}$, the number of training epochs T_{train} , the number of inference iterations T_{inf} .

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

- 1: Compute the feasible valuation ranges $\mathcal{V}_{n,\text{feas}}$ for each buyer.
- 2: // **Model Training Phase:** Train behavior predictor π_θ .
- 3: Initialize θ randomly.
- 4: **for** epoch = 1 **to** T_{train} **do**
- 5: Shuffle \mathcal{D} and partition into batches \mathcal{B} .
- 6: **for** each batch $\mathcal{B} \subset \mathcal{D}$ **do**
- 7: Sample $\tilde{\mathbf{v}}_n \sim \mathcal{V}_{n,\text{feas}}$ for each buyer n .
- 8: Compute \mathcal{L}_{CE} .
- 9: Update $\theta \leftarrow \theta - \eta_{\text{train}} \nabla_\theta \mathcal{L}_{\text{CE}}$.
- 10: **end for**
- 11: **end for**
- 12: // **Valuation Inference Phase:** Recover valuations \mathbf{v}_n^* .
- 13: **for** each buyer $n \in \mathcal{N}$ **do**
- 14: Initialize $\mathbf{v}_n^{(0)}$.
- 15: **for** iter = 1 **to** T_{inf} **do**
- 16: Compute $\mathcal{L} = - \sum_{t=0}^{T-1} \log \pi_\theta(b_n^t | \mathbf{v}_n^{(\text{iter})}, \mathbf{o}_n^t)$.
- 17: Update $\mathbf{v}_n^{(\text{iter}+1)} \leftarrow \mathbf{v}_n^{(\text{iter})} - \eta_{\text{inf}} \nabla_{\mathbf{v}} \mathcal{L}$.
- 18: **end for**
- 19: $\mathbf{v}_n^* \leftarrow \mathbf{v}_n^{(T_{\text{inf}})}$.
- 20: **end for**
- 21: **return** $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

The loss function combines standard cross-entropy with a rationality penalty:

$$\mathcal{L}_{\text{BLUE}} = (1 - \alpha) \mathcal{L}_{\text{CE}} + \alpha \mathcal{L}_{\text{ratio}}, \quad (3)$$

where $\alpha \in [0, 1]$ balances the two objectives. The theoretical formulation represents the negative log probability that valuations belong to the feasible set given observed behaviors:

$$\mathcal{L}_{\text{ratio}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \Pr(\mathbf{v}_n \in \mathcal{V}_{n,\text{feas}} | \mathbf{o}_n^t, b_n^t). \quad (4)$$

In practical implementation, we approximate this using discrete valuation sampling and probability ratios. The com-

Algorithm 2: Bayesian Learning-based Valuation Inference (BLUE)

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, learning rates $\eta_{\text{train}}, \eta_{\text{inf}}$, the number of training epochs T_{train} , the number of inference iterations T_{inf} .

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

- 1: Compute the feasible valuation ranges $\mathcal{V}_{n,\text{feas}}$ for each buyer.
- 2: Initialize θ randomly.
- 3: **for** epoch = 1 **to** T_{train} **do**
- 4: Shuffle \mathcal{D} and create batches $\{\mathcal{B}\}$.
- 5: **for** each batch \mathcal{B} **do**
- 6: Sample $\tilde{\mathbf{v}}_n \sim \mathcal{V}_{n,\text{feas}}$ for each buyer.
- 7: Compute \mathcal{L}_{CE} using $\tilde{\mathbf{v}}_n$.
- 8: Compute $\mathcal{L}_{\text{ratio}}$ using (5).
- 9: $\mathcal{L}_{\text{total}} = (1 - \alpha)\mathcal{L}_{\text{CE}} + \alpha\mathcal{L}_{\text{ratio}}$.
- 10: Update $\theta \leftarrow \theta - \eta_{\text{train}} \nabla_{\theta} \mathcal{L}_{\text{total}}$.
- 11: **end for**
- 12: **end for**
- 13: Infer each buyer's value according to (2).
- 14: **return** $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

putationally efficient version used in code is:

$$\mathcal{L}_{\text{ratio}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \frac{\sum_{v_c \in \mathcal{V}_{n,\text{feas}}} \pi_{\theta}(b_n^t | \mathbf{v}_c, \mathbf{o}_n^t)}{\sum_{v_c \in \mathcal{V}_{\text{all}}} \pi_{\theta}(b_n^t | \mathbf{v}_c, \mathbf{o}_n^t)}. \quad (5)$$

Here, \mathcal{V}_{all} is the complete valuation range considered in the analysis, defined as the discrete set ranging from 20 to 250, containing all possible valuation points.

Algorithm 2 shows the pseudocode. The Valuation Inference Phase in **BLUE** remains identical to the **CE** method.

1.3 SL

SL learns a direct mapping from buyer observations and behaviors to their private valuations. Unlike **CE** and **BLUE** which model the decision function and invert it, **SL** uses a supervised regression approach to predict valuations directly from market observations. The method incorporates critical rationality constraints to ensure predicted valuations are economically plausible.

The loss function aims to maximize the probability that inferred valuations belong to the feasible set $\mathcal{V}_{n,\text{feas}}$ given observed behaviors:

$$\mathcal{L}_{\text{SL}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \Pr(\mathbf{v}_n \in \mathcal{V}_{n,\text{feas}} | \mathbf{o}_n^t, b_n^t). \quad (6)$$

In code, we implement a computationally efficient approximation using constraint penalties. Let $\hat{\mathbf{v}}_n$ denote the valuation predicted by the neural network. Let $\mathbf{v}_{n,\text{min}}$ and $\mathbf{v}_{n,\text{max}}$ denote the values in the feasible set $\mathcal{V}_{n,\text{feasible}}$ with the lowest L_2 -norm and the highest L_2 -norm, respectively. The loss function is defined as follows (the operation max is taken

Algorithm 3: Supervised Valuation Inference (SL)

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, learning rate η , the number of epochs T_{train} ,

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

- 1: Compute the feasible valuation ranges $\mathcal{V}_{n,\text{feas}}$ for each buyer.
- 2: Initialize the neural network f_{θ} .
- 3: **for** epoch = 1 **to** T_{train} **do**
- 4: Shuffle \mathcal{D} and create batches $\{\mathcal{B}\}$.
- 5: **for** each batch \mathcal{B} **do**
- 6: Predict valuations $\hat{\mathbf{v}}_n$ according to $f_{\theta}(\mathbf{o}_n^t, b_n^t)$.
- 7: Compute lower penalty $\mathcal{L}_{\text{lower}}$.
- 8: Compute upper penalty $\mathcal{L}_{\text{upper}}$.
- 9: $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{lower}} + \mathcal{L}_{\text{upper}}$.
- 10: Update $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{total}}$.
- 11: **end for**
- 12: **end for**
- 13: Infer each buyer's value with f_{θ} .
- 14: **return** $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

elementwise).

$$\mathcal{L}_{\text{lower}} = \sum_n \|\max(\mathbf{0}, \mathbf{v}_{n,\text{min}} - \hat{\mathbf{v}}_n)\|_2^2. \quad (7)$$

$$\mathcal{L}_{\text{upper}} = \sum_n \|\max(\mathbf{0}, \hat{\mathbf{v}}_n - \mathbf{v}_{n,\text{max}})\|_2^2. \quad (8)$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{lower}} + \mathcal{L}_{\text{upper}}. \quad (9)$$

We show the pseudocode in Algorithm 3.

1.4 DL

DL employs a dual learning framework that alternately tackles two interconnected tasks: (i) valuation inference, estimating valuations from observations and behaviors, and (ii) behavior prediction, forecasting behaviors from valuations and observations. This approach creates a mutually constraining relationship where each task provides supervisory signals to the other, enhancing both through iterative refinement.

DL employs two interconnected models with distinct loss functions derived from our baseline methods:

- Behavior Prediction Model (π_{θ_b}) uses the cross-entropy loss \mathcal{L}_{CE} from Equation (1). It maps observations and valuations to behavior probabilities.
- Valuation Inference Model (f_{θ_c}) uses the feasibility constraint loss \mathcal{L}_{SL} from Equation (6). It maps observations and behaviors to valuations.

In the code implementation, we utilize four specialized loss functions that build upon the theoretical foundations of \mathcal{L}_{CE} and \mathcal{L}_{SL} , but with specific inputs tailored to each optimization step.

The first loss function is defined as

$$\begin{aligned} \mathcal{L}_{c1} = & \sum_n \|\max(\mathbf{0}, \mathbf{v}_{n,\text{min}} - \hat{\mathbf{v}}_n^{(1)})\|_2^2 \\ & + \sum_n \|\max(\mathbf{0}, \hat{\mathbf{v}}_n^{(1)} - \mathbf{v}_{n,\text{max}})\|_2^2, \end{aligned} \quad (10)$$

where $\hat{v}_n^{(1)} = f_{\theta_c}(\mathbf{o}_n^t, b_n^t)$ is the predicted valuation from the valuation inference model given the true observation and true behavior.

The second loss function is given by

$$\mathcal{L}_{b1} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \pi_{\theta_b}(b_n^t | \hat{v}_n^{(1)}, \mathbf{o}_n^t). \quad (11)$$

This loss uses the valuation prediction $\hat{v}_n^{(1)}$ from Step 1 as input to the behavior prediction model.

The third loss function is given by

$$\mathcal{L}_{b2} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \pi_{\theta_b}(b_n^t | \tilde{v}_n, \mathbf{o}_n^t), \quad (12)$$

where \tilde{v}_n is a valuation sampled from the feasible range $\mathcal{V}_{n,\text{feas}}$. The output \hat{b}_n^t is the predicted behavior generated by the behavior model using this sampled valuation.

The fourth loss function is given by

$$\begin{aligned} \mathcal{L}_{c2} = & \sum_n \|\max(\mathbf{0}, \mathbf{v}_{n,\min} - \hat{v}_n^{(2)})\|_2^2 \\ & + \sum_n \|\max(\mathbf{0}, \hat{v}_n^{(2)} - \mathbf{v}_{n,\max})\|_2^2, \end{aligned} \quad (13)$$

where $\hat{v}_n^{(2)} = f_{\theta_c}(\mathbf{o}_n^t, \hat{b}_n^t)$ is the valuation predicted using the behavior output \hat{b}_n^t from Step 3 as input to the valuation inference model. Algorithm 4 shows the pseudocode.

1.5 PGM

PGM searches for valuations that minimize the norm of the bidding policy's gradient, motivated by the principle that the optimal bidding policy has a zero gradient when it takes the true valuation as input. This method combines inverse reinforcement learning with gradient-based optimization to recover valuations by analyzing policy stability.

The method exploits the equilibrium condition where the optimal bidding policy π_{θ}^E satisfies:

$$\nabla_{\theta} \pi_{\theta}^E(\mathbf{v}^*, \mathbf{o}) = \mathbf{0} \quad (14)$$

for true valuation \mathbf{v}^* , meaning that the policy gradient vanishes at optimality. **PGM** inverts this relationship to find valuations that minimize the policy gradient norm.

The objective function for valuation recovery is

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathcal{V}_{n,\text{feas}}} \|\nabla_{\theta} \mathcal{L}_{\theta}(\mathbf{v})\|_2^2, \quad (15)$$

where the policy loss \mathcal{L} is computed as

$$\mathcal{L}_{\theta}(\mathbf{v}_n) = -\frac{1}{T} \sum_{t=1}^T \log \pi_{\theta}^E(b_n^t | \mathbf{v}_n, \mathbf{o}_n^t) \cdot G_t, \quad (16)$$

with $G_t = \sum_{k=t}^T \gamma^{k-t} r_k$ being the discounted accumulated return. Algorithm 5 shows the pseudocode.

Algorithm 4: Dual Learning Framework

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, learning rates η_c and η_b , the number of epochs T_{train} .

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

```

1: Compute the feasible valuation ranges  $\mathcal{V}_{n,\text{feas}}$  for each buyer.
2: Initialize valuation inference model  $f_{\theta_c}$  and behavior prediction model  $\pi_{\theta_b}$ .
3: for epoch = 1 to  $T_{\text{train}}$  do
4:   Shuffle  $\mathcal{D}$  and create batches  $\{\mathcal{B}\}$ .
5:   for each batch  $\mathcal{B}$  do
6:     // Step 1: Update valuation model with true behavior.
7:      $\hat{v}_n^{(1)} \leftarrow f_{\theta_c}(\mathbf{o}_n^t, b_n^t)$ .
8:     Compute  $\mathcal{L}_{c1}$ .
9:      $\theta_c \leftarrow \theta_c - \eta_c \nabla_{\theta_c} \mathcal{L}_{c1}$ .
10:    // Step 2: Update behavior model with predicted valuations.
11:    Compute  $\mathcal{L}_{b1}$  using  $\hat{v}_n^{(1)}$ .
12:     $\theta_b \leftarrow \theta_b - \eta_b \nabla_{\theta_b} \mathcal{L}_{b1}$ .
13:    // Step 3: Update behavior model with sampled valuations.
14:    Sample  $\tilde{v}_n \in \mathcal{V}_{n,\text{feas}}$  for each buyer.
15:    Compute  $\mathcal{L}_{b2}$  using  $\tilde{v}_n$ .
16:     $\theta_b \leftarrow \theta_b - \eta_b \nabla_{\theta_b} \mathcal{L}_{b2}$ .
17:    // Step 4: Update valuation model with predicted behavior.
18:     $\hat{b}_n^t \leftarrow f_{\theta_b}(\tilde{v}_n, \mathbf{o}_n^t)$ .
19:     $\hat{v}_n^{(2)} \leftarrow f_{\theta_c}(\mathbf{o}_n^t, \hat{b}_n^t)$ .
20:    Compute  $\mathcal{L}_{c2}$ .
21:     $\theta_c \leftarrow \theta_c - \eta_c \nabla_{\theta_c} \mathcal{L}_{c2}$ .
22:  end for
23: end for
24: Infer each buyer's value with  $f_{\theta_c}$ .
25: return  $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$ .

```

1.6 RvS

This method addresses valuation inference using supervised learning to train a decision function $\pi_{\theta}(\cdot | r, \mathbf{v}, \mathbf{o}^t)$ that maps reward r , valuation \mathbf{v} , and market observation \mathbf{o}^t to buyer behavior probabilities. During training, since true valuations are unavailable, we sample \mathbf{v} . And compute rewards r based on the current observation \mathbf{o}^t , the sampled valuation \tilde{v} , and the next observation \mathbf{o}^{t+1} . The inference phase involves sampling rewards from a positive range R_{pos} and then solving an optimization problem to find valuations \mathbf{v}_n^* that maximize the likelihood of observed behaviors.

During the Model Training Phase, we minimize the cross-entropy between predicted and actual behaviors:

$$\mathcal{L}_{\text{RvS}} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{T-1} \log \pi_{\theta}(b_n^t | \tilde{r}_n, \tilde{v}_n, \mathbf{o}_n^t). \quad (17)$$

where \tilde{v}_n are sampled and \tilde{r}_n are computed as $\tilde{r}_n = f_r(\mathbf{o}_n^t, \tilde{v}_n, \mathbf{o}_n^{t+1})$. This trains π_{θ} to predict behaviors given computed rewards, sampled valuations, and market observations.

Algorithm 5: Policy Gradient Minimization (PGM)

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, the number of inference iterations T_{inf} , a generative model as the market environment simulator.

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

- 1: Use \mathcal{D} to train expert policy π_θ^E .
- 2: **for** each buyer **do**
- 3: Initialize Bayesian optimizer.
- 4: **for** $i = 1$ to T_{iter} **do**
- 5: Sample valuation candidate $\mathbf{v}^{(i)}$ according to the Bayesian optimization method.
- 6: Simulate trajectories τ using $\pi_\theta^E(\cdot | \mathbf{v}^{(i)}, \cdot)$ and the environment simulator.
- 7: Compute policy loss \mathcal{L} using (16).
- 8: Calculate gradient norm $g^{(i)} = \|\nabla_\theta \mathcal{L}\|_2^2$.
- 9: Update Bayesian optimizer using the collected values of $(\mathbf{v}^{(1)}, g^{(1)}), \dots, (\mathbf{v}^{(i)}, g^{(i)})$.
- 10: **end for**
- 11: Select \mathbf{v}_n^* as the $\mathbf{v}^{(i)}$ with the lowest $g^{(i)}$.
- 12: **end for**
- 13: **return** $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

In the Valuation Inference Phase, for fixed θ , we first sample reward r_n from a positive range \mathcal{R}_{pos} , then find \mathbf{v}_n^* that minimizes the negative log-likelihood of observed behaviors:

$$\mathbf{v}_n^* = \arg \min_{\mathbf{v}_n \in \mathcal{V}_{n, \text{feas}}} - \sum_{t=0}^{T-1} \log \pi_\theta(b_n^t | r_n, \mathbf{v}_n, \mathbf{o}_n^t). \quad (18)$$

This optimization recovers valuations that best explain observed behaviors under the trained model and sampled rewards. The pseudocode is shown in Algorithm 6.

Algorithm 6: RL via Supervised Learning Inference

Input: Offline datasets of buyer $\mathcal{D} = \{(\mathbf{o}_n^t, b_n^t)\}_{n \in \mathcal{N}, t \in \mathcal{T}}$, learning rates $\eta_{\text{train}}, \eta_{\text{inf}}$, the number of training epochs T_{train} , the number of inference iterations T_{inf} .

Output: A set of inferred valuations, i.e., $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.

- 1: // **Model Training Phase:** Train behavior predictor π_θ .
- 2: Initialize θ randomly.
- 3: **for** epoch = 1 to T_{train} **do**
- 4: Shuffle \mathcal{D} and partition into batches \mathcal{B} .
- 5: **for** each batch $\mathcal{B} \subset \mathcal{D}$ **do**
- 6: Sample $\tilde{\mathbf{v}}_n$ for each buyer n .
- 7: Compute $\tilde{r}_n = f_r(\mathbf{o}_n^t, \tilde{\mathbf{v}}_n, \mathbf{o}_n^{t+1})$ using reward function.
- 8: Compute \mathcal{L}_{RvS} .
- 9: Update $\theta \leftarrow \theta - \eta_{\text{train}} \nabla_\theta \mathcal{L}_{\text{RvS}}$.
- 10: **end for**
- 11: **end for**
- 12: // **Valuation Inference Phase:** Recover valuations \mathbf{v}_n^* .
- 13: **for** each buyer $n \in \mathcal{N}$ **do**
- 14: Sample reward r_n from a positive range \mathcal{R}_{pos} and initialize $\mathbf{v}_n^{(0)}$.
- 15: **for** iter = 1 to T_{inf} **do**
- 16: Compute $\mathcal{L} = - \sum_{t=0}^{T-1} \log \pi_\theta(b_n^t | \mathbf{v}_n^{(\text{iter})}, \mathbf{o}_n^t, r_t)$.
- 17: Update $\mathbf{v}^{(\text{iter}+1)} \leftarrow \mathbf{v}^{(\text{iter})} - \eta_{\text{inf}} \nabla_{\mathbf{v}} \mathcal{L}$.
- 18: **end for**
- 19: $\mathbf{v}_n^* \leftarrow \mathbf{v}_n^{(T_{\text{inf}})}$.
- 20: **end for**
- 21: **return** $\{\mathbf{v}_n^*\}_{n \in \mathcal{N}}$.
