# Transfer Faster, Price Smarter: Minimax Dynamic Pricing under Cross-Market Preference Shift

Yi Zhang<sup>♯</sup>, Elynn Chen<sup>†</sup>, Yujun Yan<sup>♭</sup>

<sup>‡</sup>Columbia University, <sup>†</sup>New York University, <sup>b</sup>Dartmouth College

October 21, 2025







#### Introduction

- Challenge: In Dynamic pricing, new markets suffer from data scarcity, while mature markets generate abundant logs.
- Gap: Existing transfer methods assume identical utilities across markets, failing under preference shift.
- Model: A general random utility model for the market value of the product is given by

$$v_t^{(0)} = \mathring{g}^{(0)}(x_t^{(0)}) + \varepsilon_t,$$
 (1)

In cross-market transfer learning, we observe additional samples from K sources markets indexed by superscript  $^{(k)}$  for  $k \in [K]$ :

$$\mathbf{v}_t^{(k)} = \mathring{\mathbf{g}}^{(k)}(\mathbf{x}_t^{(k)}) + \varepsilon_t, \tag{2}$$

 Goal: Design a framework to leverage auxiliary markets while provably handling structured model shift.



## Online-to-Online (O2O<sub>on</sub>)

The source and target markets operate concurrently; streaming data from large markets must be incorporated into the pricing decisions of small markets in real time.

```
Algorithm 1: CM-TDP-O2Oon
   Input: Streaming source data \{(p_t^{(k)}, x_t^{(k)}, y_t^{(k)})\}_{t\geq 1} for k \in [K]; streaming target contexts
              \{x_t^{(0)}\}_{t>1}
1 Initialisation: \ell_1 \leftarrow 1, \mathcal{T}_1 = \{1\}, \hat{q}_0^{(0)} = 0.
2 for m = 1, 2, ... do
                                                                                                                                // episodes
        Compute \ell_m = 2^{m-1}, \mathcal{T}_m = \{\ell_m, \dots, \ell_{m+1} - 1\}.
        // (i) aggregate previous episode's source data
        \widehat{\widehat{g}}_m^{(\mathrm{ag})} \leftarrow \mathtt{MLE\_or\_KRR}\big(\{(p_t^{(k)}, x_t^{(k)}, y_t^{(k)})\}_{t \in \mathcal{T}_{m-1}, k \in [K]}\big).
        // (ii) debias with previous episode's target data
       \hat{\delta}_m \leftarrow \text{Debias}(\hat{\hat{g}}_m^{(\text{ag})}, \{(p_t^{(0)}, x_t^{(0)}, y_t^{(0)})\}_{t \in \mathcal{T}_{m-1}}).
5
        // For functions MLE or KRR and Debias, call Algorithm 2 for linear utility (or
             Algorithm 3 for non-parametric utility)
        Set \hat{\hat{g}}_m^{(0)} \leftarrow \hat{\hat{g}}_m^{(ag)} + \hat{\delta}_m.
        for t \in \mathcal{T}_m do
7
                                                                                                                    // episode pricing
            Post price \widehat{p}_t^{(0)} = h(\widehat{q}_m^{(0)}(x_t^{(0)})); observe y_t^{(0)} and store data.
8
```

## Offline-to-Online (O2O<sub>off</sub>)

The firm holds a fixed log of source-market data gathered before the target market opens, and this static information is used once the target goes live.

```
Algorithm 4: CM-TDP-O2O<sub>off</sub>
   Input: Offline source market data \{(p_t^{(k)}, x_t^{(k)}, y_t^{(k)})\}_{t \in \mathcal{H}(k)} for k \in [K]; feature matrix \{x_t^{(0)}\}_{t \in \mathbb{N}}
             for the target market
   /* ****** Phase 1: Update with transfer learning ******
1 Call Algorithm 2 or 3 to calculate the initial aggregated estimate \widehat{\hat{g}}^{(ag)} using entire source market
    data \{(p_t^{(k)}, X_t^{(k)}, y_t^{(k)})\}_{t \in \mathcal{H}^{(k)}} for k \in [K]
2 Apply the price \widehat{p}_1^{(0)} := h(\widehat{\mathring{g}}^{(ag)}(x_1^{(0)})) and collect data (\widehat{p}_1^{(0)}, x_1^{(0)}, y_1^{(0)}).
3 for each episode m=2,\ldots,m_0 do
        Set the length of the m-th episode: \ell_m := 2^{m-1}
        Call Algorithm 2 or 3 to calculate the debiasing estimate \widehat{\delta}_m using target market data
         \{(p_t^{(0)}, x_t^{(0)}, y_t^{(0)})\}_{t \in [2^{m-2}, 2^{m-1}-1]} and aggregated estimate \hat{g}^{(ag)}.
                                                           \hat{a}_{m}^{(0)} := \hat{a}^{(ag)} + \hat{\delta}_{m}
        For each time t, apply price \widehat{p}_t^{(0)} := h(\widehat{g}_m^{(0)}(x_t^{(0)})) and collect data (\widehat{p}_t^{(0)}, x_t^{(0)}, y_t^{(0)}).
  /* ****** Phase 2: Update without transfer learning ******
                                                                                                                                       */
s for each m \geq m_0 + 1 do
        Set the length of the m-th episode: \ell_m := 2^{m-1}
        Call Algorithm 2 or 3 to calculate \hat{\hat{g}}_m^{(0)} using target market data
         \{(p_t^{(0)}, x_t^{(0)}, y_t^{(0)})\}_{t \in [2^{m-2}, 2^{m-1}-1]}.
     For each time t, apply price and collect data.
  Output: Offered price \widehat{p}_t^{(0)}, t \ge 1
```

## Two Utility Models

#### Parametric Utility Model

• Model: Consider a linear model for the mean utility:

$$v_t^{(k)} = \mathbf{x}_t^{(k)} \cdot \boldsymbol{\beta}^{(k)} + \varepsilon_t, \quad k \in \{0\} \cup [K]$$
 (3)

• Similarity Characterization: The maximum h-norm of the difference between target and source coefficients is bounded:

$$\max_{k \in [K]} \|\beta^{(0)} - \beta^{(k)}\|_0 \le s_0.$$

• Estimation: Maximum Likelihood Estimation

### Nonparametric Utility Model

Model: Utilities in Reproducing kernel Hilbert Space:

$$v_t^{(k)} = g^{(k)}(x_t^{(k)}) + \varepsilon_t, \quad g^{(k)} \in \mathcal{H}_k, \ k \in 0 \cup [K],$$
 (4)

• Similarity Characterization: The discrepancy between target task and source task in the RKHS norm is uniformly bounded as

$$\max_{k \in [K]} \|g^{(0)} - g^{(k)}\|_K \le H.$$

• Estimation: Kernel Logistic Regression



#### THM. Regret Bound (O2O<sub>on</sub>, linear)

Upper Bound:

$$\operatorname{Regret}(T;\pi) = \mathcal{O}(\frac{d}{K}\log d\log T + s_0\log d\log T).$$

Lower Bound:

$$\inf_{\pi} \text{ sup } \operatorname{Regret}(T; \pi) \geq c_1 \frac{d}{K} \log T + c_2 s_0 \log \frac{d}{s_0} \log T,$$

#### THM. Regret Bound (O2O<sub>on</sub>, RKHS)

Upper Bound:

$$\operatorname{Regret}(T;\pi) = \mathcal{O}\left(K^{-\frac{2\alpha\beta}{2\alpha\beta+1}}T^{\frac{1}{2\alpha\beta+1}} + H^{\frac{2}{2\alpha+1}}T^{\frac{1}{2\alpha+1}}\right)$$

Lower Bound:

$$\inf_{\pi} \ \mathrm{sup} \ \mathrm{Regret}(T;\pi) \ \geq \ c \left\{ \mathcal{K}^{-\frac{2\alpha\beta}{2\alpha\beta+1}} \ T^{\frac{1}{2\alpha\beta+1}} \ + \ H^{\frac{2}{2\alpha+1}} \ T^{\frac{1}{2\alpha+1}} \right\}.$$

### Experiments

Configurations: (i) Identical Markets, (ii) Sparse difference Markets, (iii) Dense difference Markets

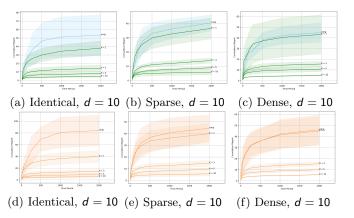


Figure 1: Cumulative regret across experimental conditions in O2O<sub>on</sub> with linear (up) and RKHS (down) utility models.

- Unified transfer pricing framework under utility shifts.

  CM-TDP is the first dynamic pricing framework that allows multiple source markets whose utilities differ from the target by a structured shift, working in both Online-to-Online and Offline-to-Online regimes.
- Minimax-optimal guarantees. We establish minimax regret rate under linear mean utilities and RKHS-smooth utilities.
- Bias-corrected aggregation architecture. Our two-step aggregate

   → debias pipeline cleanly connects meta-learning, robust statistics, and
   exploration-driven bandits, and can plug in MLE, Lasso, or kernel
   ridge as well as black boxes.
- Large empirical gains. Simulations show up to 50% lower cumulative regret, 28% lower standard error and  $5\times$  faster learning relative to single-market learning, with the largest gains in data-scarce targets under Online-to-Online transfer.