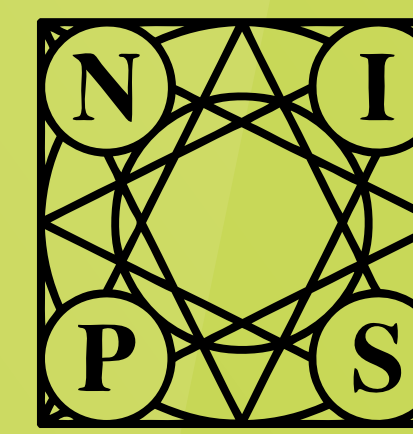


RANDOM PERMUTATION ONLINE ISOTONIC REGRESSION

WOJCIECH KOTŁOWSKI WOUTER M. KOOLEN ALAN MALEK



MOTIVATION



? distant goal: online isotonic regression on **partial orders**

... Current solution for **linear orders** does **not scale**

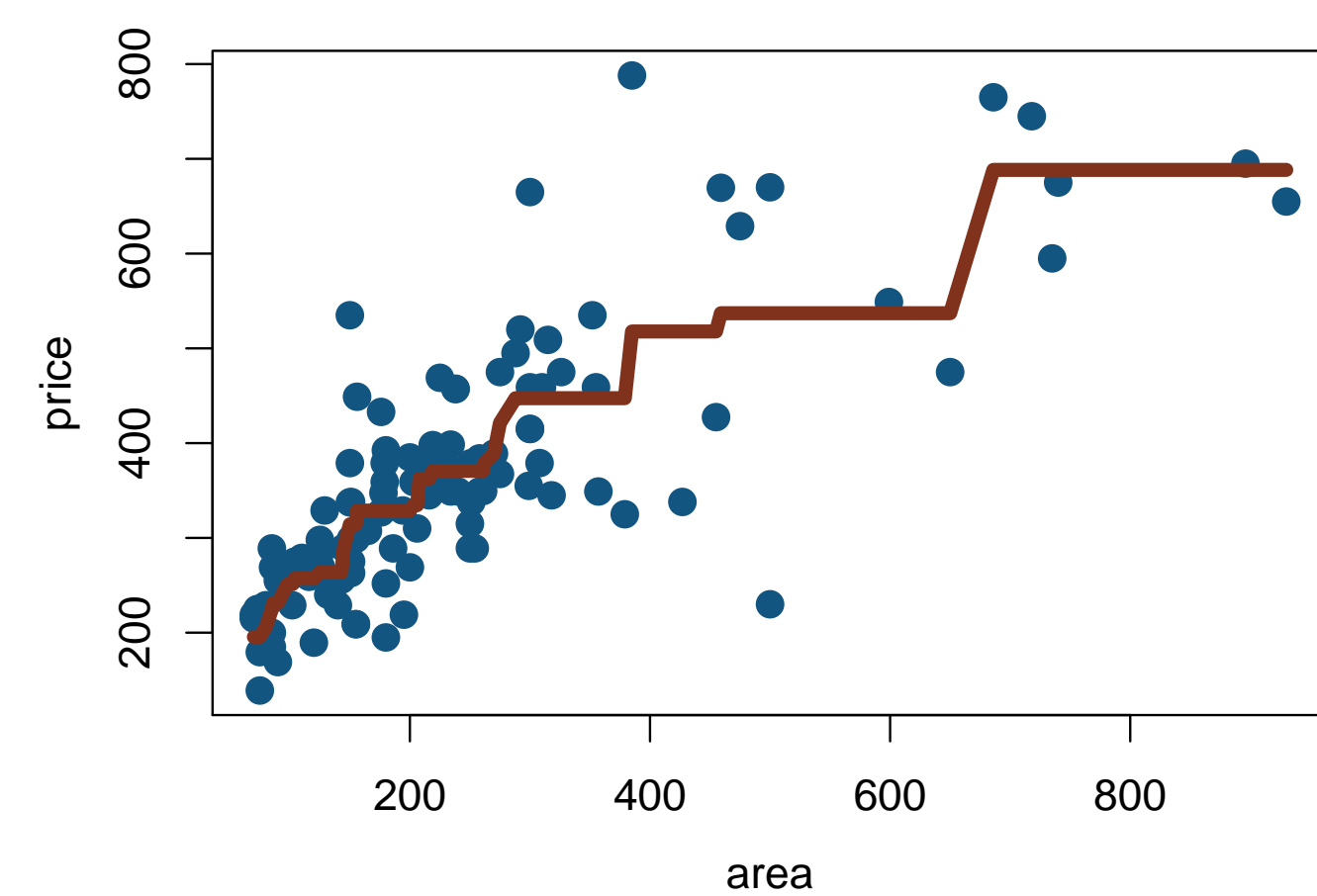
! New **model** and **algorithms** for **linear case**

(OFFLINE) ISOTONIC REGRESSION

Fit an *isotonic* (non-decreasing) function to the data:

$$f^* = \operatorname{argmin}_{\text{isotonic } f} \sum_{t=1}^T (y_t - f(x_t))^2$$

isotonic regression function



Pool Adjacent Violators Algorithm (PAVA) [Ayer et al., 1955]:

- Iteratively merge data into blocks until no violator of isotonic constraints exists
- Assign to data in each block the average of their labels y_t
- Blocks correspond to *level sets* of f^*

ONLINE ISOTONIC REGRESSION

At trial $t = 1 \dots T$:

Adversary chooses covariate x_t
 Learner predicts $\hat{y}_t \in [0, 1]$
 Adversary reveals label $y_t \in [0, 1]$
 Learner suffers squared loss $(y_t - \hat{y}_t)^2$

- Regret: $\sum_{t=1}^T (y_t - \hat{y}_t)^2 - \underbrace{\min_{\text{isotonic } f} \sum_{t=1}^T (y_t - f(x_t))^2}_{\text{total loss of offline IR function}}$
- Linear regret without restriction on x_t

RANDOM PERMUTATION MODEL

Random Permutation Model

- Adversary chooses data instances $x_1 < \dots < x_T, y_1, \dots, y_T$
- Sample UAR a permutation $\sigma = (\sigma_1, \dots, \sigma_T)$ of $\{1, \dots, T\}$
- Round t : covariate x_{σ_t} , true label y_{σ_t} , and loss $(\hat{y}_{\sigma_t} - y_{\sigma_t})^2$

Learner minimizes *expected regret*,

$$R_T := \mathbb{E}_\sigma \left[\sum_{t=1}^T (y_{\sigma_t} - \hat{y}_{\sigma_t})^2 \right] - L_T^* = \sum_{t=1}^T r_t,$$

where $r_t := \mathbb{E}_\sigma [(y_{\sigma_t} - \hat{y}_{\sigma_t})^2 - L_t^* + L_{t-1}^*]$ is the *per-round regret* and $L_t^* = L^*(\{(x_{\sigma_1}, y_{\sigma_1}), \dots, (x_{\sigma_t}, y_{\sigma_t})\})$ is the optimal loss of the first t labeled instances.

LEAVE-ONE-OUT LOSS

With Data $D = \{(x_1, y_1), \dots, (x_t, y_t)\}$, the *loo* of a t round game is

$$\text{loo}_t(D) := \frac{1}{t} \left(\left(\sum_{i=1}^t (y_i - \hat{y}_i(D \setminus (x_i, y_i)))^2 \right) - L^*(D) \right).$$

Lemma 1. $r_t(D) \leq \text{loo}_t(D)$ for any t and any data set $D = \{(x_1, y_1), \dots, (x_t, y_t)\}$.

LOWER BOUND

Adversarial lower bound [Kotłowski, Koolen, and Malek, 2016] applies to random permutation model: $\text{loo}_t = \Omega(t^{-2/3})$.

MATCHING BOUNDS

Theorem 2. There is an algorithm for the random-permutation model with excess leave-one-out loss $\text{loo}_t = \tilde{O}(t^{-2/3})$ and hence expected regret $R_T \leq \sum_t \tilde{O}(t^{-2/3}) = \tilde{O}(T^{1/3})$, which matches the lower bound of $\text{loo}_t = \Omega(t^{-2/3})$.

Caveat: algorithm is not efficient (on partial orders)!

FORWARD ALGORITHMS

Two observations:

- PAVA is efficient and generalizes to partial orders
- Follow The Leader algorithms are common in practice

Forward Algorithm: To predict at x_t , imagine $y'_t \in [0, 1]$, compute f^* on $\{(x_1, y_1) \dots (x_{t-1}, y_{t-1})\} \cup \{(x_t, y'_t)\}$, and play $\hat{y}_t = f^*(x_t)$.

FORWARD ALGORITHM EXAMPLES

- IR-Int: Compute f^* on past data. Predict with average of f^* at nearest x_i .
- Interpolation $\hat{y}_i = \lambda_i \hat{y}_i^0 + (1 - \lambda_i) \hat{y}_i^1$ (where \hat{y}_i^0 and \hat{y}_i^1 are „plug-in $y'_t = 0$ ” and „plug-in $y'_t = 1$ ”)
- Last step minimax:

$$\hat{y}_i = \operatorname{argmin}_{\hat{y} \in [0, 1]} \max_{y_i \in [0, 1]} \{(\hat{y} - y_i)^2 - L^*(\mathbf{y})\}$$

- IVAP predictors [Vovk et al., 2015]:

$$\hat{y}_i^{\log} = \frac{\hat{y}_i^1}{\hat{y}_i^1 + 1 - \hat{y}_i^0}, \quad \hat{y}_i^{\text{Brier}} = \frac{1 + (\hat{y}_i^0)^2 - (1 - \hat{y}_i^1)^2}{2}$$

REGRET BOUNDS

Theorem 3. Any forward algorithm has $\text{loo}_t = O(\sqrt{\frac{\log t}{t}})$.

Lemma 4. The IR-Int algorithm has $\text{loo}_t = \Omega(t^{-1/2})$ when run on

$$\underbrace{10000}_{1/5} \quad \underbrace{11000}_{2/5} \quad \underbrace{11100}_{3/5} \quad \underbrace{11110}_{4/5} \quad \underbrace{11111}_{5/5}.$$

HEAVY- γ

Parameters: Weight $c > 0$ and label $\gamma \in [0, 1]$.

Algorithm: To predict at x_t

- Compute isotonic regression f' on **weighted** dataset

$$D' := \{(x_s, y_s, 1) \mid 1 \leq s < t\} \cup \{(x_t, \gamma, c)\}$$

- Predict $y_t = f'(x_t)$

Efficient weighted algorithms available [Kyng et al., 2015].

TUNING HEAVY- γ

Any fixed label γ works. We like $\gamma = 1$. (Not all adaptive labels work. Fixed point + lower bound.)

Theorem 5. Heavy- γ has sub-optimal loo_t loss unless $c = \Theta(t^{1/3})$.

Conjecture 6. Heavy- γ with weight $c = \Theta(t^{1/3})$ has $\text{loo}_t = \tilde{O}(t^{-2/3})$.