

Environment Invariant Linear Least Squares, Fan *et al.* (2023)

Xiaohui Yin

University of Connecticut

October 7, 2024

Motivating Example: Risks of Using Spurious Variables

A task to classify cows and camels based on extracted hierarchical features:

- There is a dataset \mathcal{D} containing 10k images of cows and camels from the Internet. When Using 70% to train the classifier, top two features are the back shape x_1 , and the background color x_2 , and it performs well on the remaining 30% testing set. It is contemplated that cows often appear on the grass while most camels appear on the sand from this dataset.



Figure: Camel in brown

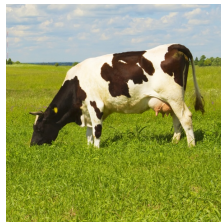


Figure: Cow in green

Motivating Example (Cond)

- However, introducing x_2 is not what we expected: the classifier can not work well in a place farming camels and cows, in which the background color is fixed.
- If we have another dataset $\tilde{\mathcal{D}}$, in which the association between the background color and object label still exists yet slightly perturbs. Intuitively, we may infer that x_2 may be a “spurious” variable for prediction or causation.



Figure: Camel in green

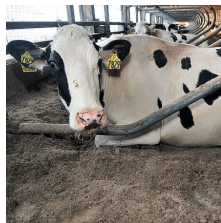


Figure: Cow in brown

Some Definition

This paper propose the multiple-environment version of linear least squares. It utilizes the heterogeneity across datasets. It combines the linear least squares solutions across different datasets and uses their differences to determine the true parameter and important variables in a completely data-driven manner.

Here are some definitions needed:

- response variable $y \in \mathbb{R}$
- explanatory variable $\mathbf{x} \in \mathbb{R}^p$
- set of environments \mathcal{E}
- observations from environment $e \in \mathcal{E}$:
 $(\mathbf{x}_1^{(e)}, y_1^{(e)}), \dots, (\mathbf{x}_n^{(e)}, y_n^{(e)}) \sim \mu^{(e)}$

Model Formulation

The true linear model is assumed as

$$y^{(e)} = (\beta_{S^*}^*)^\top \mathbf{x}_{S^*}^{(e)} + \varepsilon^{(e)} \quad \text{with} \quad \mathbb{E}[\varepsilon^{(e)} | \mathbf{x}_{S^*}^{(e)}] \equiv 0, \quad (1)$$

where the unknown set of important variables $S^* = \{j : \beta_j^* \neq 0\}$ and the model parameters β^* are the same across different environments, while $\mu^{(e)}$, the distribution of $(\mathbf{x}^{(e)}, y^{(e)})$, may vary.

The aim is to estimate β^* and S^* using the $n \cdot |\mathcal{E}|$ data $\{(\mathbf{x}_i^{(e)}, y_i^{(e)})\}_{e \in \mathcal{E}, i \in \{1, \dots, n\}}$.

Challenges

- β^* is not the best linear predictor for each single environment.
- A smaller mean squared error by incorporating the linear spurious variables $\mathbf{x}_{S^*\perp}$ defined outside the set of important variables for all environments. $\mathbf{x}_{S^*\perp}^{(e)}$ can capture information in $\varepsilon^{(e)}$.
- Such spurious variables are not stable when generalized to other environments, because the association between these variables and y are not stable. (background color in the motivating example)

Loss Function

The population-level EILLS objective is

$$\begin{aligned} Q_{\gamma}(\beta) = & \underbrace{\sum_{e \in \mathcal{E}} \mathbb{E} \left[|y^{(e)} - \beta^{\top} \mathbf{x}^{(e)}|^2 \right]}_{R(\beta)} \\ & + \underbrace{\gamma \sum_{j=1}^p \mathbf{1}\{\beta_j \neq 0\} \times \sum_{e \in \mathcal{E}} \left| \mathbb{E}[(y^{(e)} - \beta^{\top} \mathbf{x}^{(e)})x_j^{(e)}] \right|^2}_{J(\beta)}, \end{aligned} \quad (2)$$

where γ is a tuning parameter.

- $R(\beta)$ requires a good overall solution in each single environment.
- $J(\beta)$ discourages selecting variables that have strong correlation with the fitted residuals in some environments.
- From computational perspective, we have

$$J(\beta) = \sum_{j=1}^p \mathbf{1}\{\beta_j \neq 0\} \sum_{e \in \mathcal{E}} \frac{1}{4} |\nabla_j R^{(e)}(\beta)|^2, \quad (3)$$

Empirical Loss Function

$$\hat{J}(\beta) = \sum_{j=1}^p \mathbf{1}\{\beta_j \neq 0\} \sum_{e \in \mathcal{E}} \left| \hat{\mathbb{E}}[x_j^{(e)} (y^{(e)} - \beta^\top \mathbf{x}^{(e)})] \right|^2$$

$$\hat{R}(\beta) = \sum_{e \in \mathcal{E}} \frac{1}{n^{(e)}} \sum_{i=1}^{n^{(e)}} \{y_i^{(e)} - \beta^\top \mathbf{x}_i^{(e)}\}^2$$

$$\hat{Q}(\beta, \gamma) = \hat{R}(\beta) + \lambda \hat{J}(\beta)$$

Some linear exogenous variables which do not contribute to explaining y can be eliminated further by introducing an l_0 penalty as

$$\hat{L}(\beta, \gamma, \lambda) = \hat{R}(\beta) + \lambda \hat{J}(\beta) + \lambda \|\beta\|_0.$$

An Illustration

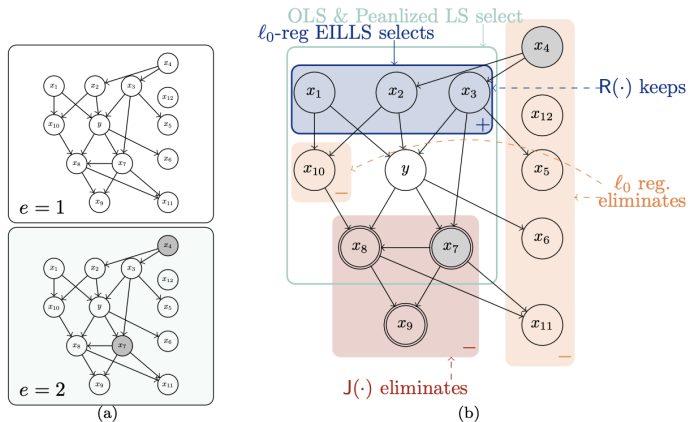


Figure: An Example for Illustration, figure from Fan *et al.* (2023)

References I

Fan, J., Fang, C., Gu, Y. and Zhang, T. (2023) Environment invariant linear least squares. URL <https://arxiv.org/abs/2303.03092>.