

Uncoupled Regression from Pairwise Comparison Data

Liyuan Xu^{1,2}, Junya Honda^{1,2}, Gang Niu², Masashi Sugiyama^{2,1}
¹The University of Tokyo ²RIKEN

Abstract

- Introduce a novel uncoupled regression problem with pairwise comparison data.
- Propose a new empirical risk minimization methods to solve the problem.
- Propose two estimators for risks for general marginal target distributions.

Problem Setting

Motivating Example

Sensitive Data:

- e.g. Salary, Number of crime committed before,...
- People won't giving an explicit label.
- Containing sensitive data leads the risk of security breach.



Goal: To build a prediction model from marginal distribution of sensitive data

Uncoupled Regression

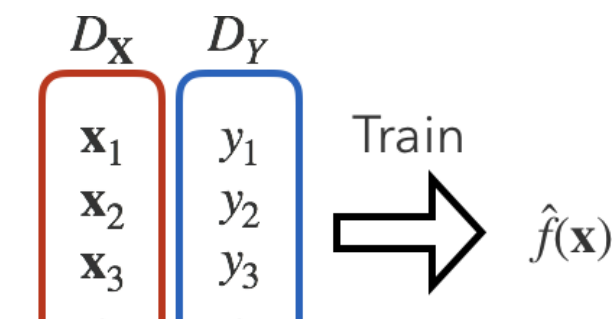
Ordinary Regression:

- Target Y generated from feature \mathbf{X} as $Y = f(\mathbf{X}) + \epsilon$.
- Learn a model \hat{f} from coupled data $\mathcal{D}_{\mathbf{X},Y} = \{(\mathbf{x}_i, y_i)\}$.

Uncoupled Regression:

- Unlabeled data $\mathcal{D}_{\mathbf{X}} = \{\mathbf{x}_i\}$ generated from distribution $P_{\mathbf{X}}$.
- Target $\mathcal{D}_Y = \{y_i\}$ generated from marginal distribution P_Y .
- We try to learn a model \hat{f} to predict Y from \mathbf{X} .

Since no correspondence in $\mathcal{D}_{\mathbf{X}}$ and \mathcal{D}_Y ,
problem is ill-posed without any further assumption.



Pairwise Comparison Data

Difficult to get sensitive data but easier to get their order.

- People might willing to give order information.

Pairwise Comparison Data $\mathcal{D}_R = \{\mathbf{x}_i^+, \mathbf{x}_i^-\}$

- Consists of pairs of features $\{\mathbf{x}_i^+, \mathbf{x}_i^-\}$ such that

$$y(\mathbf{x}_i^+) > y(\mathbf{x}_i^-),$$



where $y(\mathbf{x}_i^+)$, $y(\mathbf{x}_i^-)$ are the target values for \mathbf{x}_i^+ , \mathbf{x}_i^- , respectively.

Generation Process:

- Generate two samples (\mathbf{X}, Y) , (\mathbf{X}', Y') from joint distribution $P_{\mathbf{X},Y}$.
- If $Y \geq Y'$, $\mathbf{X}^+ = \mathbf{X}$ and $\mathbf{X}^- = \mathbf{X}'$. If not, the opposite holds.

- Let $P_{\mathbf{X}^+}$, $P_{\mathbf{X}^-}$ be the distribution of each comparison data.

Question

Can we learn a model \hat{f} from

Unlabeled data $\mathcal{D}_{\mathbf{X}}$, Target values \mathcal{D}_Y , and Comparative Data \mathcal{D}_R ?

Related Work

[Carpentier and Schlueter, 2016]

Uncoupled regression with one-dimensional features and monotonic target function f .

- Requires features to be one-dimensional and target function f to be monotonic.
- Involves complex optimization in learning.

Our Problem

- Applicable to features with multi-dimensions and non-monotonic f .
- Easy to implement.

Algorithm

Empirical Risk Minimization Principle

$$\begin{array}{ccc} \text{Regression Risk} & \xleftarrow{\text{Unbiased}} & \text{Empirical Risk} \\ R(f) = \mathbb{E}_{\mathbf{X},Y}[(Y - f(\mathbf{X}))^2] & & \hat{R}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 \end{array}$$

- Construct an unbiased estimator of risk (e.g. l_2 -risk) from coupled data.
- Minimize the unbiased risk estimator to learn a model \hat{f} .

Distribution of Comparison Data

Lemma Let F_Y be the cumulative distribution function of P_Y . Then,

$$\begin{aligned} \mathbb{E}_{\mathbf{X}^+}[f(\mathbf{X}^+)] &= 2\mathbb{E}_{\mathbf{X},Y}[F_Y(Y)f(\mathbf{X})], \\ \mathbb{E}_{\mathbf{X}^-}[f(\mathbf{X}^-)] &= 2\mathbb{E}_{\mathbf{X},Y}[(1 - F_Y(Y))f(\mathbf{X})], \end{aligned}$$

Therefore, if $F_Y(y) = y$ (i.e. marginal distributions P_Y is uniform on $[0, 1]$),

$$\begin{aligned} \mathbb{E}_{\mathbf{X},Y}[(Y - f(\mathbf{X}))^2] &= \mathbb{E}_Y[Y^2] + \mathbb{E}_{\mathbf{X}}[(f(\mathbf{X}))^2] - 2\mathbb{E}_{\mathbf{X},Y}[Yf(\mathbf{X})] \\ &= \mathbb{E}_Y[Y^2] + \mathbb{E}_{\mathbf{X}}[(f(\mathbf{X}))^2] - \mathbb{E}_{\mathbf{X}^+}[f(\mathbf{X}^+)] \\ &\xleftarrow{\text{Unbiased}} \mathbb{E}_Y[Y^2] + \frac{1}{|\mathcal{D}_{\mathbf{X}}|} \sum_{\mathbf{x}_i \in \mathcal{D}_{\mathbf{X}}} (f(\mathbf{x}_i))^2 - \frac{1}{|\mathcal{D}_R|} \sum_{(\mathbf{x}_i^+, \mathbf{x}_i^-) \in \mathcal{D}_R} f(\mathbf{x}_i^+) \end{aligned}$$

- $\mathbb{E}_Y[Y^2]$ does not depend on f and can be ignored in optimization.
- $\{\mathbf{x}_i^-\}$ can be used for the variance reduction.
- However, we cannot construct unbiased estimators for all marginal distributions.

→ Propose two approaches to construct estimators with small bias.

Risk Approximation Approach

Main Idea

Approximate the expectation $\mathbb{E}_{\mathbf{X},Y}[Yf(\mathbf{X})]$

by the linear combination $w_1\mathbb{E}_{\mathbf{X}^+}[f(\mathbf{X}^+)] + w_2\mathbb{E}_{\mathbf{X}^-}[f(\mathbf{X}^-)]$, $w_1, w_2 \in \mathbb{R}$.

Theorem Let \hat{f}_{RA} be the minimizer of

$$\mathbb{E}_Y[Y^2] + \frac{1}{|\mathcal{D}_{\mathbf{X}}|} \sum_{\mathbf{x}_i \in \mathcal{D}_{\mathbf{X}}} (f(\mathbf{x}_i))^2 - \frac{1}{|\mathcal{D}_R|} \sum_{(\mathbf{x}_i^+, \mathbf{x}_i^-) \in \mathcal{D}_R} (w_1 f(\mathbf{x}_i^+) + w_2 f(\mathbf{x}_i^-)). \quad (1)$$

Then, with an adequate condition, with probability $1 - \delta$,

$$R(\hat{f}_{\text{RA}}) \leq R(f) + O\left(\sqrt{\frac{\log 1/\delta}{n_U}}\right) + O\left(\sqrt{\frac{\log 1/\delta}{n_R}}\right) + \text{Err}(w_1, w_2)$$

holds, where R is l_2 -risk, the CDF F_Y of P_Y and Err is defined as

$$\text{Err}(w_1, w_2) = \mathbb{E}_Y[|Y - w_1 F_Y(Y) - w_2 (1 - F_Y(Y))|].$$

- When marginal distribution P_Y is uniform, estimator (1) is unbiased for R .
 - In this case, $\text{Err} = 0$ with $(w_1, w_2) = (1, 0)$.
- Can be generalized to any risk defined based on Bregman Divergence.

Referece

- A. Carpentier and T. Schlueter. Learning relationships between data obtained independently. In Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, 2016.

Target Transformation Approach

Main Idea

Transform target Y to $F_Y(Y)$, and minimize the risk on transformed variable:

$$\mathbb{E}_{\mathbf{X},Y}[(F_Y(Y) - F_Y(f(\mathbf{X})))^2].$$

Note, marginal distribution of $F_Y(Y)$ is uniform distribution on $[0, 1]$.

Theorem With an appropriate condition, the minimizer \hat{f}_{TT} of

$$\mathbb{E}_Y[Y^2] + \sum_{i=1}^{n_U} (F_Y(f(\mathbf{x}_i)))^2 - \sum_{i=1}^{n_R} F_Y(f(\mathbf{x}_i^+)) \quad (2)$$

satisfies

$$R(\hat{f}_{\text{TT}}) \leq R(f) + O\left(\sqrt{\frac{\log 1/\delta}{n_U}}\right) + O\left(\sqrt{\frac{\log 1/\delta}{n_R}}\right) + \Delta_{\text{TT}}$$

with probability $1 - \delta$, where R is l_2 -risk.

- When marginal is uniform, the estimator (2) is unbiased for R , since $F_Y(y) = y$.
- Δ_{TT} depends on the shape of P_Y and noise level.
- The theorem only holds for l_2 -risks.

Experiments

Settings

- Used benchmark datasets from UCI repository.
- Used original features as unlabeled data and sample 5000 pairs of comparison data.
- Learnt linear models and predicted the target value for unlabelled data.

Methods to be compared

- Linear Regression using fully labeled ordinary coupled data.
- Train SVMRank using pairwise comparison data, predict ranking, and predict value by

$$\hat{f}(\mathbf{x}) = F_Y^{-1}\left(\frac{\hat{n}(\mathbf{x})}{n_U}\right),$$

where $\hat{n}(\mathbf{x})$ is the predicted rank in the data.

Results

Dataset	Supervised Regression		Uncoupled Regression	
	LR	SVMRank	RA	TT
housing	24.5(5.0)	110.3(29.5)	29.5(6.9)	22.5(6.2)
diabetes	3041.9(219.8)	8575.9(883.1)	3087.3(256.3)	3127.3(278.8)
airfoil	23.3(2.2)	62.1(7.6)	23.7(2.0)	22.7(2.2)
concrete	109.5(13.3)	322.9(45.8)	111.7(13.2)	139.1(17.9)
powerplant	20.6(0.9)	372.2(34.8)	21.8(1.1)	22.0(1.0)
mpg	12.1(2.04)	125(15.1)	12.8(2.16)	10.3(2.08)
redwine	0.412(0.0361)	1.28(0.112)	0.442(0.0473)	0.466(0.0412)
whitewine	0.574(0.0325)	1.58(0.0691)	0.597(0.0382)	0.644(0.0414)
abalone	5.05(0.375)	20.9(1.44)	5.26(0.372)	5.54(0.424)

Better than SVMRank, and may be better to ordinal supervised learning