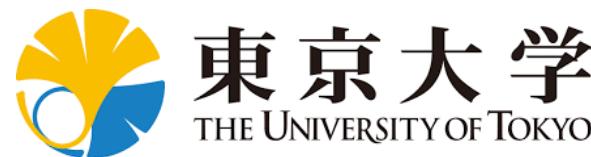


Semi-Supervised Ordinal Regression Based on Empirical Risk Minimization

Taira Tsuchiya^{1,2} Nontawat Charoenphakdee^{1,2}

Issei Sato^{1,2} Masashi Sugiyama^{2,1}

1. The University of Tokyo 2. RIKEN AIP



Ordinal Regression Problem

- Medical diagnosis [Bender+ 1997]
 - E.g., diabetic retinopathy: eye disease associated with long-standing diabetes
 - State of patient (retina) is expressed with five **discrete but ordered label** by severity.

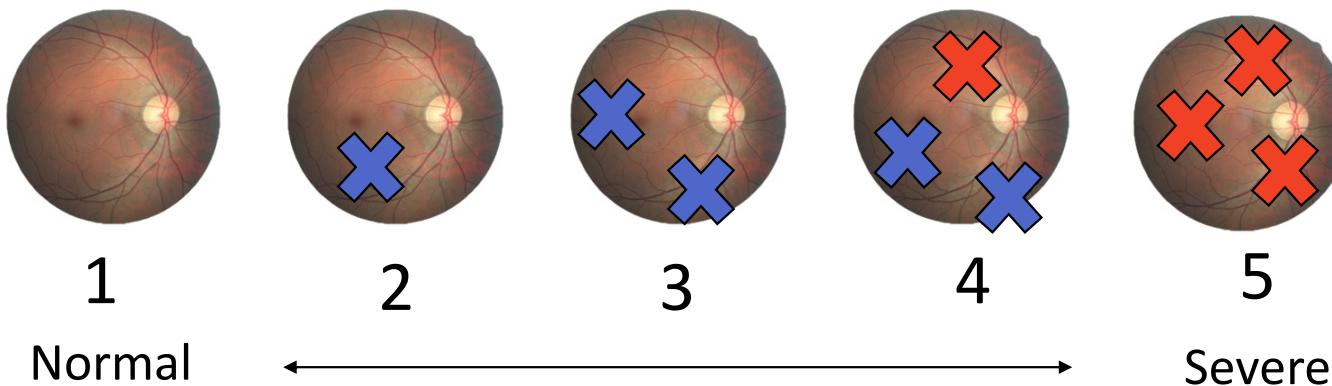


Figure 1. Example of ordinal labels for retina images¹

Misclassifying 5 to 1 can be more harmful than 5 to 4

1. <https://www.kaggle.com/c/diabetic-retinopathy-detection>

Ordinal Regression Formulation

3

[Pedregosa+ JMLR2017]

- Estimate discrete but ordered label y from data point \mathbf{x} by function g

$$y = g(\mathbf{x}),$$

where $g : \mathcal{X} \rightarrow \mathcal{Y} = \{1, \dots, K\}$.

- Q. How to predict **discrete but ordered labels?**
- A. Use ***prediction function***: g to predict label

$$g(\mathbf{x}; f, \boldsymbol{\theta}) := 1 + \sum_{i=1}^{K-1} \mathbb{I}\{f(\mathbf{x}) > \theta_i\} \quad \text{(? Discrete nature)}$$

$$f : \mathcal{X} \rightarrow \mathbb{R}$$

$$\boldsymbol{\theta} = [\theta_1, \dots, \theta_{K-1}]^\top \in \mathbb{R}^{K-1} \quad (\theta_1 \leq \theta_2 \leq \dots \leq \theta_{K-1})$$

- Goal: find g minimizing following ***task risk***

$$\mathcal{R}(g) := \mathbb{E}_{X,Y} [\mathcal{L}(g(X), Y)]$$

$\mathcal{L}(g(\mathbf{x}), y)$: ***task loss*** (e.g., absolute $|y - g(\mathbf{x})|$, squared $(y - g(\mathbf{x}))^2$)

Task Losses and their Task Surrogate Losses

4

[Pedregosa+ JMLR2017]

- Three task losses and corresponding task surrogate losses

Task loss: $\mathcal{L}(g(\mathbf{x}), y)$

Absolute

$$|y - g(\mathbf{x})|$$

Zero-one

$$\mathbb{I}\{f(\mathbf{x}) \leq \theta_{y-1}\} + \mathbb{I}\{f(\mathbf{x}) > \theta_y\}$$

Squared

$$(y - g(\mathbf{x}))^2$$

Task surrogate loss: $\psi(\alpha(\mathbf{x}), y)$

$$\sum_{i=1}^{y-1} \ell(-\alpha_i) + \sum_{i=y}^{K-1} \ell(\alpha_i) \quad \text{all threshold (AT)}$$

$$\ell(-\alpha_{y-1}) + \ell(\alpha_y)$$

Immediate threshold (IT)

$$\left(y + \alpha_1 - \frac{3}{2}\right)^2$$

Least squares (LS)

⌚ Minimizing task loss can be computationally intractable

- Surrogate goal: minimize **task surrogate risk**

⌚ Computationally efficient

$\ell(\cdot)$: binary surrogate loss

$\alpha_i(\mathbf{x}) := \theta_i - f(\mathbf{x})$

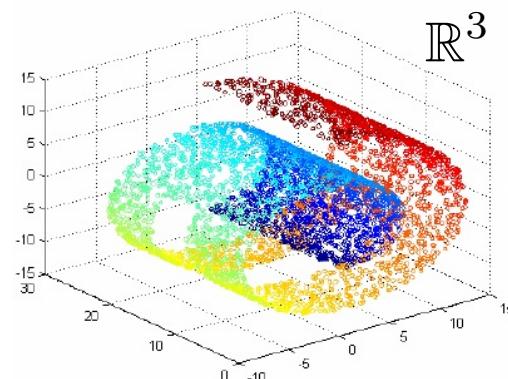
$$\mathcal{R}_\psi(g) := \mathbb{E}_{X,Y} [\psi(\alpha(X), Y)]$$

Semi-Supervised Learning

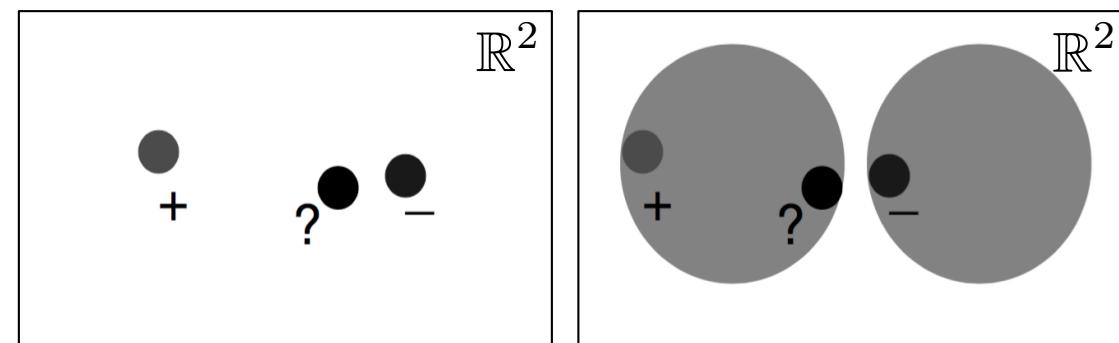
5

- Labeling process can be costly and time-consuming (especially in medical area.)
-> Make use of unlabeled data (semi-supervised ordinal regression)
- Usually ***geometric assumption*** on unlabeled data is required.

[Chapelle+ 2010]



Manifold assumption



Cluster assumption
 $\hat{=}$ Low-density separation

- 😞 Performance degrades if the geometric assumption does not hold [Sakai+ ICML2017].

Existing Work

We are only aware of three studies for semi-supervised OR.

- Liu et al. (2011) relies on a **manifold assumption** for image classification.
- Seah et al. (2012) relies on a **cluster assumption** and cannot predict unseen unlabeled data (transductive).
- Srijith et al. (2013) relies on **low-density separation principle** using Gaussian process.

	Geometric assumption	Loss agnostic	Theoretical guarantee	Inductive	Computational cost
Liu et al. [3]	Yes	No	No	😊 Yes	High
Seah et al. [4]	Yes	😊 Yes	No	No	Medium
Srijith et al. [5]	Yes	No	No	😊 Yes	High
Ours	😊 No	😊 Yes	😊 Yes	😊 Yes	😊 Low

Problem Setup

Dataset

- Suppose that we are given the following data:

$$\mathcal{X}_L := \{(x_j^L, y_j)\}_{j=1}^{n_L} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}(X, Y) \quad n_L : \# \text{ of labeled data}$$

- Also given unlabeled data drawn from marginal distribution:

$$\mathcal{X}_U := \{x_j^U\}_{j=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}(X) = \sum_{y=1}^K \pi_y \mathbb{P}(X|Y=y), \quad \pi_y := \mathbb{P}(Y=y)$$

$n_U : \# \text{ of unlabeled data}$

Empirical risk minimization (ERM) [Vapnik 1998]

(Expected) task surrogate risk:

$$\mathcal{R}_\psi(g) := \mathbb{E}_{X,Y} [\psi(\alpha(X), Y)]$$

[Estimate risk by finite samples]

Empirical task surrogate risk:

$$\widehat{\mathcal{R}}_\psi(g) := \frac{1}{n_L} \sum_{j=1}^{n_L} \psi(\alpha(x_j), y_j)$$



unbiased

↑ Minimize this

Proposed Risk Estimator

(Our goal)

Idea: use a large portion of U data to estimate risk (inspired by Sakai+ 2017)

Task surrogate risk : $\mathcal{R}_\psi(g) := \mathbb{E}_{X,Y} [\psi(\alpha(X), Y)]$

unbiased

Risk estimator : $\widehat{\mathcal{R}}_{\psi,LU}^{\setminus k}(g) := \sum_{y \in \mathcal{Y} \setminus k} \frac{\pi_y}{n_y} \sum_{j=1}^{n_y} \psi(\alpha(x_j^y), y)$

- * Holds for any $k \in \{1, \dots, K\}$
- * k : removed class

$$+ \frac{1}{n_U} \sum_{j=1}^{n_U} \psi(\alpha(x_j^U), k)$$

😊 Use unlabeled data

$$- \sum_{y \in \mathcal{Y} \setminus k} \frac{\pi_y}{n_y} \sum_{j=1}^{n_y} \psi(\alpha(x_j^y), k)$$

😢 Data labeled as class k is not used

-> Combine with the supervised risk via convex sum

$$\mathcal{R}_{\psi,SEMI-\gamma}^{\setminus k}(g) := \gamma \mathcal{R}_{\psi,LU}^{\setminus k}(g) + (1 - \gamma) \mathcal{R}_\psi(g)$$

Theoretical Analysis (informal)

Prop. (Fisher-consistency)

Task surrogate loss: AT, IT, or LS (Bayes optimal risk)

⇒ Every minimizer g of $\mathcal{R}_{\psi, \text{SEMI}-\gamma}^{\setminus k}(g)$ reaches $\mathcal{R}^* := \inf_g \mathcal{R}(g)$

😊 Minimizing $\mathcal{R}_{\psi, \text{SEMI}-\gamma}^{\setminus k}(g) \iff$ minimizing $\mathcal{R}(g)$

Thm. (Estimation error bound)

$$\begin{aligned}\hat{g}^{\setminus k} &:= \arg \min_{g \in \mathcal{G}} \hat{\mathcal{R}}_{\psi, \text{LU}}^{\setminus k}(g) \\ g^* &:= \arg \min_{g \in \mathcal{G}} \mathcal{R}_{\psi, \text{LU}}^{\setminus k}(g)\end{aligned}$$

$$\frac{\mathcal{R}_\psi(\hat{g}^{\setminus k}) - \mathcal{R}_\psi(g^*)}{\text{Estimation error of empirical risk minimizer}} = \mathcal{O}_p \left(\sum_{y \in \mathcal{Y}^{\setminus k}} \frac{\pi_y}{\sqrt{n_y}} + \frac{1}{\sqrt{n_U}} \right)$$

L data U data
not labeled as class k

😊 Estimation error $\rightarrow 0$ w/ optimal rate w/o strong assumption

Thm. (Variance reduction)

For fixed $g \in \mathcal{G}$, if we select γ properly

$$\mathcal{R}_{\psi, \text{SEMI}-\gamma}^{\setminus k}(g) := \gamma \mathcal{R}_{\psi, \text{LU}}^{\setminus k}(g) + (1 - \gamma) \mathcal{R}_\psi(g)$$

$$\text{Var}[\hat{\mathcal{R}}_{\psi, \text{SEMI}-\gamma}^{\setminus k}(g)] < \text{Var}[\hat{\mathcal{R}}_\psi(g)]$$

😊 Unlabeled data reduce variance

What Class to Pick Out?

10

We can arbitrarily select k to estimate risk by our risk $\widehat{\mathcal{R}}_{\psi, \text{LU}}^{(k)}(g)$.

Q. How can we choose “best” k ?

Strategy 1: pick k with that has smallest number of data to reduce variance:

$$k = \arg \min_{y \in \mathcal{Y}} n_y$$

Strategy 2: pick k that minimizes upper bound of estimation error:

$$\mathcal{R}_\psi(\widehat{g}^{(k)}) - \mathcal{R}_\psi(g^*) = \mathcal{O}_p \left(\sum_{y \in \mathcal{Y}^{(k)}} \frac{\pi_y}{\sqrt{n_y}} + \frac{1}{\sqrt{n_U}} \right)$$

Minimize here

$$\implies k = \arg \max_{y \in \mathcal{Y}} n_y$$

Experiments 1: Variance Reduction

11

- Setup
 - Model: linear-in-parameter model
 - **Task surrogate loss: AT, IT, LS**
 - # of trials: 20
- Experiment 1: Variance reduction

See the ratio of variance of empirical risk for random regressor g_{rand}

$$\text{Var}[\widehat{\mathcal{R}}_{\psi, \text{SEMI}-\gamma}^k(g_{\text{rand}})] / \text{Var}[\widehat{\mathcal{R}}_\psi(g_{\text{rand}})] \quad <- \text{ good if smaller than 1}$$

Surr.	car	era	lev	swd	winequality	bank	calhousing	census
AT	0.108	0.115	0.044	0.041	0.058	0.313	0.122	0.234
IT	0.109	0.116	0.042	0.048	0.058	0.360	0.140	0.196
LS	0.157	0.081	0.070	0.065	0.071	1.303	0.275	0.736

Surr.	computer	autompq	abalone	boston	machinecpu	stocks	triazines	wisconsin
AT	0.122	0.165	0.100	0.110	0.097	0.099	0.227	0.181
IT	0.124	0.165	0.104	0.126	0.130	0.078	0.284	0.222
LS	0.387	0.223	0.195	0.065	1.021	0.159	0.176	0.198

😊 Unlabeled data indeed help reduce the variance

Experiments 2: Task Loss-agnostic Property

12

- Experiment 2: comparison with supervised learning

Task surrogate loss: AT, evaluation metrics: mean squared error

Dataset	SV-Linear	SEMI1-Linear	SEMI2-Linear	SV-Kernel	SEMI1-Kernel	SEMI2-Kernel
car	0.529 (0.09)	0.566 (0.06)	0.449 (0.08)	0.401 (0.12)	0.614 (0.14)	0.345 (0.02)
era	0.541 (0.11)	0.464 (0.11)	0.410 (0.12)	0.411 (0.15)	0.320 (0.06)	0.303 (0.04)
lev	0.493 (0.08)	0.309 (0.16)	0.176 (0.10)	0.168 (0.08)	0.192 (0.14)	0.127 (0.02)
swd	0.625 (0.06)	0.670 (0.05)	0.379 (0.09)	0.253 (0.01)	0.509 (0.13)	0.252 (0.01)
winequality	0.424 (0.08)	0.457 (0.11)	0.034 (0.02)	0.020 (0.01)	0.291 (0.19)	0.018 (0.00)
bank	0.651 (0.07)	0.640 (0.06)	0.639 (0.05)	0.438 (0.11)	0.464 (0.10)	0.406 (0.02)
calhousing	0.655 (0.06)	0.619 (0.10)	0.519 (0.08)	0.463 (0.09)	0.466 (0.11)	0.418 (0.04)
census	0.635 (0.09)	0.624 (0.08)	0.515 (0.10)	0.447 (0.11)	0.451 (0.12)	0.405 (0.02)
computer	0.609 (0.16)	0.602 (0.16)	0.499 (0.16)	0.471 (0.18)	0.522 (0.19)	0.412 (0.04)
autompq	0.439 (0.22)	0.467 (0.21)	0.359 (0.17)	0.346 (0.24)	0.267 (0.02)	0.263 (0.03)
abalone	0.612 (0.13)	0.470 (0.15)	0.370 (0.11)	0.347 (0.10)	0.307 (0.00)	0.307 (0.00)
boston	0.536 (0.15)	0.431 (0.16)	0.387 (0.15)	0.450 (0.25)	0.275 (0.16)	0.229 (0.03)
machinecpu	0.494 (0.26)	0.536 (0.29)	0.375 (0.16)	0.431 (0.19)	0.587 (0.19)	0.346 (0.07)
stocks	0.520 (0.12)	0.488 (0.15)	0.329 (0.10)	0.346 (0.10)	0.320 (0.12)	0.274 (0.03)
triazines	0.783 (0.06)	0.774 (0.06)	0.733 (0.09)	0.349 (0.00)	0.635 (0.23)	0.349 (0.00)
wisconsin	0.578 (0.07)	0.557 (0.07)	0.568 (0.08)	0.512 (0.09)	0.454 (0.01)	0.483 (0.06)

- ☺ Estimation error bound based strategy works well
- ☺ Works well for most of datasets, all evaluation metrics

Same tendency for (IT, mean zero-one error), (LS, mean squared error)

Conclusion

- Proposed a novel semi-supervised ordinal regression framework based on unbiased risk estimator

Evaluation metrics agnostic

	Geometric assumption	Loss agnostic	Theoretical guarantee	Inductive	Computational cost
Liu et al. [3]	Yes	No	No	😊 Yes	High
Seah et al. [4]	Yes	😊 Yes	No	No	Medium
Srijith et al. [5]	Yes	No	No	😊 Yes	High
Ours	😊 No	😊 Yes	😊 Yes	😊 Yes	😊 Low

Use unlabeled data to estimate risk

- Consistency
- Estimation error bound
- Variance reduction

References

- [1] Fabian Pedregosa, Francis Bach, and Alexandre Gramfort. On the consistency of ordinal regression methods. *Journal of Machine Learning Research*, 18:1 – 35, 2017.
- [2] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. MIT Press, 2006.
- [3] Yang Liu, Yan Liu, Shenghua Zhong, and Keith CC Chan. Semi-supervised manifold ordinal regression for image ranking. In ACMMM, pages 1393–1396, 2011.
- [4] Chun-Wei Seah, Ivor W Tsang, and Yew-Soon Ong. Transductive ordinal regression. *IEEE Transactions on Neural Networks and Learning Systems*, 23(7):1074–1086, 2012.
- [5] PK Sripathi, Shirish Shevade, and S Sundararajan. Semi-supervised Gaussian process ordinal regression. In ECML-PKDD, pages 144–159, 2013.
- [6] Tomoya Sakai, Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Semi-supervised classification based on classification from positive and unlabeled data. In ICML, volume 70, pages 2998–3006, 2017.