# Semi-Supervised AUC Optimization without Guessing Labels of Unlabeled Data

Zheng Xie and Ming Li AAAI, 2018

#### Introduction

#### □ Semi-supervised Learning

- > Obtain the labels for the collected data is expensive
- Semi-supervised learning are based on certain distributional assumption, e.g. cluster assumption and manifold assumption
- Estimate the labels of unlabeled instances
- **AUC** Optimization for SSL
  - Rely on the distributional assumptions
  - Require prior probability
    - Sakai, Niu, and Sugiyama (2017) proposed an unbiased semi-supervised AUC optimization method based on positive-unlabeled learning

#### □ Proposed Method

- > Do not depend on any distributional assumptions
- Do not need to know prior probability

#### AUC



The AUC is the probability the model will score a randomly chosen positive class higher than a randomly chosen negative class.

#### **AUC Optimization**

• Supervised learning

• Since AUC is equivalent to the probability of a randomly drawn positive instance being ranked before a randomly drawn negative instance, it can be formulated as

. . .

AUC = 1 - 
$$\mathbb{E}_{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}} [\mathbb{E}_{\boldsymbol{x}' \in \mathcal{X}_{\mathrm{N}}} [\ell_{01}(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}'))]]$$
$$R_{\mathrm{PN}}$$

#### **AUC Optimization in SSL**

 $\succ$  Rely on the distributional assumptions

• These assumptions may be violated in many real-world applications

➢ Require an accurate estimation of the prior probability

• It is usually difficult especially when the number of labeled data is extremely small.

We can achieve unbiased semi-supervised AUC optimization without distributional assumptions or prior knowledge about the distribution or class prior probabilities.





$$R_{\mathrm{PU}} = \mathop{\mathbb{E}}_{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}} \left[ \mathop{\mathbb{E}}_{\boldsymbol{x}'' \in \mathcal{X}_{\mathrm{U}}} \left[ \ell_{01} (\boldsymbol{w}^{\top} (\boldsymbol{x} - \boldsymbol{x}'')) \right] \right],$$
$$R_{\mathrm{NU}} = \mathop{\mathbb{E}}_{\boldsymbol{x}'' \in \mathcal{X}_{\mathrm{U}}} \left[ \mathop{\mathbb{E}}_{\boldsymbol{x}' \in \mathcal{X}_{\mathrm{N}}} \left[ \ell_{01} (\boldsymbol{w}^{\top} (\boldsymbol{x}'' - \boldsymbol{x}')) \right] \right].$$

Proof:  

$$R_{\mathrm{PU}} = \underset{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}}{\mathbb{E}} [\underset{\boldsymbol{x}'' \in \mathcal{X}_{\mathrm{U}}}{\mathbb{E}} [\ell_{01}(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}''))]] = \underbrace{\mathbb{E}}_{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}} [\underset{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}}{\mathbb{E}} [\ell_{01}(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}))]] + \theta_{N} \underset{\boldsymbol{x}' \in \mathcal{X}_{\mathrm{N}}}{\mathbb{E}} [\ell_{01}(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}'))]] = \frac{1}{2} \theta_{P} + \theta_{N} \underset{\boldsymbol{x} \in \mathcal{X}_{\mathrm{P}}}{\mathbb{E}} [\underset{\boldsymbol{x}' \in \mathcal{X}_{\mathrm{N}}}{\mathbb{E}} [\ell_{01}(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}'))]] = \frac{1}{2} \begin{bmatrix} l_{01}(\boldsymbol{w}(\boldsymbol{x}_{1} - \boldsymbol{x}_{2})) + l_{01}(\boldsymbol{w}(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}))] = \frac{1}{2} \end{bmatrix}$$

#### **Unbiased Estimation**

$$\theta_P + \theta_N = 1$$
$$R_{PU} = \theta_N R_{PN} + \frac{1}{2} \theta_P$$
$$R_{NU} = \theta_P R_{PN} + \frac{1}{2} \theta_N$$
$$R_{PU} + R_{NU} - \frac{1}{2} = R_{PN}$$

Conduct unbiased AUC risk estimation without knowing the class prior probabilities  $\theta_P$  and  $\theta_N$ .

# Alg-1. SAMULT

Semi-supervised AUC Maximization by treating the UnLabeled data in Two ways

Objective: 
$$\min_{\boldsymbol{w}} \gamma \widehat{R}_{PN} + (1 - \gamma) \left( \widehat{R}_{PU} + \widehat{R}_{NU} - \frac{1}{2} \right) + \lambda ||\boldsymbol{w}||^{2}$$
supervised AUC  
risk
Where,
$$\widehat{R}_{PN} = \frac{1}{n_{P}n_{N}} \sum_{\boldsymbol{x} \in \mathcal{X}_{P}} \sum_{\boldsymbol{x}' \in \mathcal{X}_{N}} \ell(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}'))$$

$$\widehat{R}_{PU} = \frac{1}{n_{P}n_{U}} \sum_{\boldsymbol{x} \in \mathcal{X}_{P}} \sum_{\boldsymbol{x}'' \in \mathcal{X}_{U}} \ell(\boldsymbol{w}^{\top}(\boldsymbol{x} - \boldsymbol{x}'))$$

$$\widehat{R}_{NU} = \frac{1}{n_{U}n_{N}} \sum_{\boldsymbol{x}'' \in \mathcal{X}_{U}} \sum_{\boldsymbol{x}' \in \mathcal{X}_{N}} \ell(\boldsymbol{w}^{\top}(\boldsymbol{x}'' - \boldsymbol{x}'))$$

$$\gamma \in [0, 1]$$

$$\lambda \geq 0$$

#### **Analytical Solution**

$$\widehat{\boldsymbol{w}} = (\gamma \boldsymbol{H}_{\mathrm{PN}} + (1 - \gamma)(\boldsymbol{H}_{\mathrm{PU}} + \boldsymbol{H}_{\mathrm{NU}}) + \lambda \boldsymbol{I}_d)^{-1} (\gamma \boldsymbol{h}_{\mathrm{PN}} + (1 - \gamma)(\boldsymbol{h}_{\mathrm{PU}} + \boldsymbol{h}_{\mathrm{NU}})) ,$$

where

$$\begin{split} \boldsymbol{h}_{\mathrm{PN}} &= \frac{1}{n_{\mathrm{P}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{P}}} - \frac{1}{n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{N}}}, \\ \boldsymbol{h}_{\mathrm{PU}} &= \frac{1}{n_{\mathrm{P}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{P}}} - \frac{1}{n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{U}}}, \\ \boldsymbol{h}_{\mathrm{NU}} &= \frac{1}{n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{U}}} - \frac{1}{n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{N}}}, \\ \boldsymbol{H}_{\mathrm{PN}} &= \frac{1}{n_{\mathrm{P}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{P}} - \frac{1}{n_{\mathrm{P}} n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{P}}} \boldsymbol{1}_{n_{\mathrm{P}}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{N}}} \boldsymbol{1}_{n_{\mathrm{N}}}, \\ \boldsymbol{H}_{\mathrm{PN}} &= \frac{1}{n_{\mathrm{P}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{P}} - \frac{1}{n_{\mathrm{P}} n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{P}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{P}}} \boldsymbol{1}_{n_{\mathrm{P}}} \boldsymbol{1}_{n_{\mathrm{P}}} \boldsymbol{1}_{n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}, \\ &- \frac{1}{n_{\mathrm{P}} n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{N}}} \boldsymbol{1}_{n_{\mathrm{P}}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{P}} + \frac{1}{n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{N}}, \\ &- \frac{1}{n_{\mathrm{P}} n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{U}}} \boldsymbol{1}_{n_{\mathrm{P}}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{P}} + \frac{1}{n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{U}}, \\ &- \frac{1}{n_{\mathrm{P}} n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{U}}} \boldsymbol{1}_{n_{\mathrm{V}}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{U}} + \frac{1}{n_{\mathrm{U}}} \boldsymbol{X}_{\mathrm{U}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{N}}, \\ &- \frac{1}{n_{\mathrm{U}} n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{1}_{n_{\mathrm{N}}} \boldsymbol{1}_{n_{\mathrm{U}}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{U}} + \frac{1}{n_{\mathrm{N}}} \boldsymbol{X}_{\mathrm{N}}^{\mathsf{T}} \boldsymbol{X}_{\mathrm{N}}, \end{split}$$

#### SAMULT<sup>P+U</sup>

When only the positive and unlabeled data are available

## Alg-2. SAMPURA

Semi-supervised AUC Maximization by Partitioning Unlabeled data at RAndom

$$\mathcal{X}_{U} \implies \stackrel{\mathcal{X}_{U^{+}}}{\longrightarrow} \stackrel{\mathcal{X}_{P'}}{\longrightarrow} = \mathcal{X}_{P} \cup \mathcal{X}_{U^{+}} \qquad \qquad \mathcal{X}_{P'} \text{ should be ranked before } \mathcal{X}_{N}$$
$$\mathcal{X}_{U^{-}} \implies \mathcal{X}_{N'} = \mathcal{X}_{N} \cup \mathcal{X}_{U^{-}} \qquad \qquad \mathcal{X}_{P} \text{ should be ranked before } \mathcal{X}_{N'}$$

$$\min_{\boldsymbol{w}} \quad \widehat{R}_{\mathrm{P'N}} + \widehat{R}_{\mathrm{PN'}} + \lambda ||\boldsymbol{w}||^2$$

take the average of those **w**s to construct the final ensemble

The model that minimizes the PU-AUC risk or the PNU-AUC risk converges to the supervised case



 $\widehat{\boldsymbol{w}}^*$  is learned from all available data with label

Compare our methods to the state-of-the-art semi-supervised AUC optimization methods.

- SAMULT
- SAMPURA
- **SSRankBoost** (Amini, Truong, and Goutte 2008), a boosting based algorithm for learning a bipartite ranking function with partially labeled data
- **PNU-AUC** (Sakai, Niu, and Sugiyama 2017), which is a semi-supervised AUC optimization method based on positive-unlabeled learning.
- Supervised AUC optimization
- Logistic regression

Dataset	Supervised	Log. Reg.	SSRankBoost	PNU-AUC	SAMULT	SAMPURA
australian	.879±.029	.860±.027	.886±.013	.903±.009	.903±.009	.903±.009
breast	$.655 {\pm} .097$	$.625 {\pm} .095$	$.647 {\pm} .065$	.701±.029	.701±.029	.704±.026
breastw	$.987 {\pm} .009$	$.980 {\pm} .006$	$.984 {\pm} .013$	$.992 \pm .001$	<b>.996±.001</b>	<b>.996±.001</b>
clean1	$.760 \pm .062$	$.725 {\pm} .060$	$.737 {\pm} .050$	$.767 {\pm} .042$	$.777 \pm .039$	.782±.038
colic	$.829 \pm .112$	$.818 {\pm} .074$	$.721 \pm .062$	$.858 {\pm} .013$	.869±.013	.870±.013
colic.orig	$.647 \pm .093$	$.645 {\pm} .076$	$.612 \pm .081$	$.644 {\pm} .048$	.658±.049	.663±.044
credit-a	$.893 \pm .024$	$.886 \pm .023$	.885±.013	<b>.906±.008</b>	<b>.906±.007</b>	.906±.008
credit-g	$.719 \pm .043$	$.709 {\pm} .030$	$.665 \pm .027$	$.748 {\pm} .018$	$.748 {\pm} .018$	.761±.017
fourclass	$.825 \pm .023$	$.826 \pm .026$	$.692 \pm .029$	.827±.008	.827±.008	.828±.006
german	$.683 \pm .057$	$.672 \pm .048$	$.709 {\pm} .025$	.727±.019	.727±.019	.729±.017
haberman	$.551 \pm .086$	$.530 {\pm} .075$	<b>.582±.067</b>	$.547 {\pm} .051$	$.551 \pm .045$	.556±.043
heart	$.857 {\pm} .065$	$.842 \pm .060$	$.823 \pm .042$	.876±.025	.876±.025	.878±.024
house	<b>.975±.038</b>	$.961 \pm .015$	.972±.034	<b>.979±.012</b>	.979±.012	.979±.011
ijcnn1	.912±.003	$.901 \pm .004$	$.902 \pm .002$	$.904 \pm .009$	.913±.005	.915±.004
madelon	$.510 \pm .037$	$.541 \pm .020$	.571±.023	$.528 {\pm} .029$	$.517 {\pm} .027$	.530±.022
parkinsons	$.848 \pm .129$	$.826 {\pm} .082$	$.799 {\pm} .051$	.860±.023	.860±.023	.863±.011
phishing	$.975 {\pm} .097$	$.972 \pm .001$	$.983 {\pm} .003$	$.974 \pm .002$	$.976 \pm .002$	.985±.002
vehicle	$.932 {\pm} .038$	$.922 \pm .022$	$.912 \pm .039$	.965±.020	.965±.020	.970±.014
vote	$.965 {\pm} .038$	$.951 {\pm} .015$	<b>.972±.034</b>	<b>.979±.012</b>	<b>.979±.012</b>	<b>.979±.011</b>
wdbc	$.971 \pm .014$	$.963 {\pm} .006$	$.964 \pm .016$	<b>.983±.006</b>	<b>.983±.006</b>	.983±.005
# Best/Comp.	2	0	4	11	15	18

Dataset	PU-RSVM	PU-AUC	SAMULT <sup>P+U</sup>
australian	$.844 \pm .034$	$.898 {\pm} .021$	.900±.019
breast	$.615 \pm .104$	<b>.701±.077</b>	<b>.701±.077</b>
breastw	$.987 {\pm} .009$	$.993 {\pm} .002$	<b>.996±.002</b>
clean1	$.709 {\pm} .072$	$.786 {\pm} .050$	.796±.050
colic	$.807 \pm .103$	<b>.877±.060</b>	<b>.877±.060</b>
colic.orig	$.621 \pm .078$	$.650 {\pm} .068$	.670±.061
credit-a	$.876 {\pm} .028$	<b>.912±.015</b>	.912±.015
credit-g	$.688 {\pm} .047$	$.755 {\pm} .028$	.757±.027
fourclass	$.823 \pm .031$	.832±.026	.832±.025
german	$.642 \pm .045$	$.734 \pm .034$	.736±.034
haberman	.572±.083	$.561 \pm .081$	$.555 {\pm} .080$
heart	$.835 {\pm} .073$	<b>.883±.039</b>	.883±.040
house	$.945 \pm .029$	$.980 {\pm} .012$	<b>.983±.011</b>
ijcnn1	<b>.927±.005</b>	$.900 {\pm} .011$	$.905 {\pm} .012$
madelon	$.470 \pm .015$	<b>.533±.031</b>	$.514 \pm .032$
parkinsons	$.797 {\pm} .089$	.870±.033	.870±.032
phishing	$.960 {\pm} .008$	$.966 {\pm} .005$	<b>.970±.005</b>
vehicle	$.942 \pm .034$	$.959 {\pm} .030$	.966±.025
vote	$.945 {\pm} .029$	$.980 {\pm} .012$	<b>.983±.011</b>
wdbc	$.967 {\pm} .021$	$.984 {\pm} .007$	<b>.985±.007</b>
#Best/Com.	2	7	17

#### Degeneration to AUC Optimization for Positive and Unlabeled Data

**PU-RSVM** (Sellamanickam, Garg, and Selvaraj 2011): A ranking SVM based method for positive and unlabeled data

**PU-AUC** (Sakai, Niu, and Sugiyama 2017): A positive unlabeled AUC optimization method by optimizing an unbiased AUC risk estimator relies only on positive and unlabeled data

Do not need to know the prior.

#### Conclusion

- In semi-supervised AUC optimization, it is unnecessary to design sophisticated strategies to estimate the possible labels of the unlabeled data or the class prior probabilities.
- Proposed two semi-supervised AUC optimization methods: SAMULT and SAMPURA
- The positive-unlabeled AUC optimization problem can be addressed by a degenerated version of proposed method that simply treats the unlabeled data as negative.

#### Reference

- Zheng Xie, Ming Li. Semi-Supervised AUC Optimization without Guessing Labels of Unlabeled Data. AAAI, 2018.
- Sundararajan Sellamanickam, Priyanka Garg, Sathiya Keerthi Selvaraj. A pairwise ranking based approach to learning with positive and unlabeled examples. CIKM, 2011.
- Tomoya Sakai, Gang Niu, and Masashi Sugiyama. Semi-Supervised AUC Optimization Based on Positive-Unlabeled Learning. *Machine Learning*, 2017.
- Massih-Reza Amini, Tuong-Vinh Truong, Cyril Goutte. A Boosting Algorithm for Learning Bipartite Ranking Functions with Partially Labeled Data. SIGIR, 2008.