September 16, 2017

# Recent Advances on Positive-Unlabeled (PU) Learning

## Gang Niu

Project Assistant Professor, UTokyo Visiting Scientist, RIKEN AIP







#### What is this talk about?

- Fully supervised deep learning from big data is successful
- Nevertheless, massive labeled data is not always available
  - Medicine, manufacturing, disaster, infrastructure ...
- Achieving high accuracy with low labeling costs is always a big challenge (and is our ultimate goal)



#### Overview

- Focus on binary classification the most studied learning problem
- In PU learning, a binary classifier is trained from only P & U data
- Unbiased PU learning
  - State-of-the-art approach to training linear-in-parameter models
  - Unbiased empirical risk estimators are minimized
- Non-negative PU learning
  - State-of-the-art approach to training deep models
  - Non-negative empirical risk estimators are minimized
- PNU learning
  - Convex combinations of the objectives of PU & PN learning

#### Outline

#### Introduction

- Unbiased PU Learning (du Plessis, Niu & Sugiyama, NIPS'14 & ICML'15)
- Non-negative PU Learning (Kiryo, Niu, du Plessis & Sugiyama, NIPS'17)
   (for oral presentation; there are 40 orals/678 acceptance/3240 submissions)
- □ PNU Learning (Sakai, du Plessis, Niu & Sugiyama, *ICML'17*)
- □ Theoretical Analyses (Niu, du Plessis, Sakai, Ma & Sugiyama, NIPS'16)

### Problem settings in a nutshell

#### **PU** learning **PN** learning **PNU** learning (i.e., semi-supervised learning) (i.e., supervised learning) X 0 0 X P, N & U data are P & N data are P & ∪ data are available for training available for training available for training

**O** : positive data

🗙 : negative data

🗆 : unlabeled data

### Motivation of PU learning

- PN learning (whether deep or not) is data demanding
- Then, is PNU learning the most natural choice?
  - Certainly, if the two classes are symmetric
     → Which is P & which is N does not matter
  - Not really, if there is intrinsic difference in them
     Which is P & which is N matters
- PU learning is preferred due to the following reasons
  - I. N data are too expensive
  - II. N data are too diverse
  - III. "N data" are impure

#### Case I: N data are too expensive

- Some data-collecting activity is prohibited by law
  - Suppose I am a market researcher at Apple
  - Let **Samsung** be the imaginary enemy
  - Plan A: Hack the data center of Samsung
  - Plan B: Send corporate spies to Samsung
- Some activity is costly & risky



- Some clinical trials involve healthy subjects with no pre-existing medical conditions → Extreme financial incentives
- While others pertain to patients with specific health conditions who are willing to try an experimental treatment (https://en.wikipedia.org/wiki/Clinical\_trial)

### Case II: N data are too diverse

- In the first running example
  - Android is neither the only competitor
  - Nor is Samsung the only Android vendor

Period	Samsung	Apple	Huawei	OPPO	vivo	Others
	U					
2016Q1	23.8%	15.4%	8.4%	5.9%	4.4%	42.1%
2016Q2	22.7%	11.7%	9.3%	6.6%	4.8%	45.0%
2016Q3	20.9%	12.5%	9.3%	7.1%	5.9%	44.3%
2016Q4	18.0%	18.2%	10.5%	7.3%	5.7%	40.2%
2017Q1	23.3%	14.7%	10.0%	7.5%	5.5%	39.0%
Source: IDC. May 2017						

- In the second running example
  - Too difficult to sample healthy subjects without **selection bias**

promo/smartphone-

market-share/vendor)



#### Case III: "N data" are impure

- Wait, and rethink the definition of N
- For the sake of smartphone advertising
  - A customer is N = S/he hates iPhone and would never buy it in the whole life
  - N = {non-potential user} ∈ {not existing user} = U
- For the purpose of drug testing
  - Many diseases are due to **chromosomal abnormality**
  - P = {specific genetic disorder} ∋ {specific medical condition}
  - N = {not possess this disorder} ∈ {not observe this condition} = U



Ρ

### Previous work in PU learning

- Binary classification (applied to retrieval & novelty/outlier detection)
  - Statistical query model Denis (ALT'98); De Comité+ (ALT'99); Letouzey+ (ALT'00)
  - Linear(-in-parameter) model
     Liu+ (ICML'02); Li & Liu (IJCAI'03); Lee & Liu (ICML'03); Liu+ (ICDM'03);

     Elkan & Noto (KDD'08); ← The first that might be unbiased
     Ward+ (Biometrics 2009); Scott & Blanchard (AISTATS'09); Blanchard+ (JMLR 2010);
     du Plessis+ (NIPS'14) ← The first that must be unbiased

#### Other applications

- Matrix completion Hsieh+ (ICML'15)
- Sequential data Li+ (SDM'09); Nguyen+ (IJCAI'11)

Statistical learning theory

### Outline

#### Introduction

#### Unbiased PU Learning

- Non-negative PU Learning
- PNU Learning
- Theoretical Analyses

#### Empirical risk minimization (Vapnik, Statistical Learning Theory, 1998)

- A "cookbook" procedure of ERM
  - Choose a loss, so that the (expected) learning objective (which is known as the risk) can be defined
  - 2. Choose a **model**, so that the risk can be minimized over this model family (rather than all measurable functions)
  - 3. Approximate the risk by an **empirical risk estimator**
  - 4. Minimize the empirical risk by an **optimization algorithm**
- Nice to have independent learning obj., model & opt. alg.
   Arguably a key reason for the great success of deep learning
- 3. Straightforward for PN; non-trivial for PU

#### Notation

- $X \in \mathbb{R}^d, Y \in \{\pm 1\}$ : Input & output RVs
- p(x, y): Underlying joint density
- $p_p(x) = p(x|Y = +1)$ : P marginal  $p_n(x) = p(x|Y = -1)$ : N marginal

$$p(x): \cup$$
 marginal

- π<sub>p</sub> = p(Y = +1): Class-prior probability Assumed known;
   can be estimated Ramaswamy+ (ICML'16); Jain+ (NIPS'16); du Plessis+ (MLJ 2017)
- $\mathbb{E}_{p}[\cdot] = \mathbb{E}_{X \sim p_{p}}[\cdot], \mathbb{E}_{n}[\cdot] = \mathbb{E}_{X \sim p_{n}}[\cdot]$ : Expectation over P/N marginal

$$\mathcal{X}_{p} = \{x_{i}^{p}\}_{i=1}^{n_{p}} \stackrel{\text{i.i.d.}}{\sim} p_{p}(x): P \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{n} = \{x_{i}^{n}\}_{i=1}^{n_{n}} \stackrel{\text{i.i.d.}}{\sim} p_{n}(x): N \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{\text{i.i.d.}}{\sim} p(x): U \text{ data} \qquad \mathcal{X}_{u} = \{x_{i}^{u$$

#### Empirical risk estimator in PN learning

Let g be a decision function &  $\ell$  be a loss function

■ The **risk** of *g* is

$$R(g) = \mathbb{E}_{(X,Y)\sim p(x,y)} [\ell(Yg(X))]$$
  
=  $\pi_{p} \mathbb{E}_{p} [\ell(g(X))] + \pi_{n} \mathbb{E}_{n} [\ell(-g(X))]$ 

where  $\pi_n = 1 - \pi_p$ 

The risk can be approximated directly by  $\widehat{R}_{pn}(g) = \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} \ell(g(x)) + \frac{\pi_n}{n_n} \sum_{x \in \mathcal{X}_n} \ell(-g(x))$ 

This doesn't work for PU learning!

#### Empirical risk estimator in PU learning (du Plessis+, ICML'15)

#### Key observations

- $\pi_n p_n(x) = p(x) \pi_p p_p(x)$
- $\pi_{n}\mathbb{E}_{n}[\ell(-g(X))] = \mathbb{E}_{X}[\ell(-g(X))] \pi_{p}\mathbb{E}_{p}[\ell(-g(X))]$
- Thus the risk can be expressed as  $R(g) = \pi_{p} \mathbb{E}_{p} \left[ \ell(g(X)) - \ell(-g(X)) \right] + \mathbb{E}_{X} \left[ \ell(-g(X)) \right]$
- This can be approximated indirectly by

$$\widehat{R}_{pu}(g) = \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} \left[ \ell(g(x)) - \ell(-g(x)) \right] + \frac{1}{n_u} \sum_{x \in \mathcal{X}_u} \ell(-g(x))$$

Simple in retrospect!

Non-convex special case (du Plessis+, NIPS'14)

• If 
$$\ell(t) + \ell(-t) = 1$$
  

$$R(g) = 2\pi_p \mathbb{E}_p \left[\ell(g(X))\right] + \mathbb{E}_X \left[\ell(-g(X))\right] - \pi_p$$

• Non-convex in *g* 

Examples



Convex special case (du Plessis+, ICML'15)

• If 
$$\ell(t) - \ell(-t) = -t$$
 (Natarajan+, *NIPS'13*; Patrini+, *ICML'16*)  
 $R(g) = \pi_p \mathbb{E}_p[-g(X)] + \mathbb{E}_X[\ell(-g(X))]$ 

• Convex in g, and convex in  $\theta$  if  $g(x; \theta)$  is linear in  $\theta$ 

Examples



#### Remarks

- $\hat{R}_{pu}(g)$  is unbiased & consistent, similarly to  $\hat{R}_{pn}(g)$ 
  - Why is unbiasedness important? Easy to be consistent
- **Biased SVM** (Liu+, *ICDM'03*) best method prior to Elkan & Noto (*KDD'08*)  $\frac{1}{2} ||w||^2 + C_p \sum_{x \in \mathcal{X}_p} \ell_H(w \cdot x + b) + C_u \sum_{x \in \mathcal{X}_u} \ell_H(-w \cdot x - b)$ 
  - Guess/search C<sub>p</sub> & C<sub>u</sub>
  - No learning guarantee
- **SMO** solver speeding up double hinge loss (Sansone+, arXiv 2016)

Note that theoretical & experimental results omitted for simplicity

### Outline

- Introduction
- Unbiased PU Learning
- Non-negative PU Learning
- PNU Learning
- Theoretical Analyses

#### Revisit ERM

- 2. Choose a model  $\mathcal{G}$ 
  - Approximation error  $\inf_{g \in \mathcal{G}} R(g) \inf_g R(g)$ Smaller for more flexible model
- 3. Approximate R(g) by  $\hat{R}(g)$ 
  - Estimation error R(ĝ) − inf<sub>g∈G</sub>R(g) ← ĝ = arg inf<sub>g∈G</sub>R̂(g)
     Smaller for less flexible model, or bigger training data
     Converge to zero, if learning is consistent
- 4. Minimize  $\hat{R}(g)$  by an optimization algorithm which returns  $\hat{g}'$ 
  - Optimization error  $\hat{R}(\hat{g}') \hat{R}(\hat{g})$

#### Motivation

- $\hat{R}_{pu}(g)$  is nice for training linear-in-parameter models
- However, it cannot be used for training deep networks



#### Non-negative risk estimator (Kiryo+, NIPS'17)

- Rescue with neither changing model nor labeling more data
- Recall  $\pi_{n}\mathbb{E}_{n}[\ell(-g(X))] = \mathbb{E}_{X}[\ell(-g(X))] \pi_{p}\mathbb{E}_{p}[\ell(-g(X))]$ 
  - Approximate left-hand-side  $\rightarrow \hat{R}_{pn}(g) \ge 0$
  - Approximate right-hand-side  $\rightarrow \hat{R}_{pu}(g) \ge 0$
- Force it to be non-negative!

$$\widetilde{R}_{pu}(g) = \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} \left[ \ell(g(x)) \right] + \max\left\{ 0, \frac{1}{n_u} \sum_{x \in \mathcal{X}_u} \ell(-g(x)) - \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} \left[ \ell(-g(x)) \right] \right\}$$

• Minimizing  $\tilde{R}_{pu}(g)$  is no longer embarrassingly parallel

#### Large-scale learning algorithm (Kiryo+, NIPS'17)

• Let 
$$(\mathcal{X}_{p}^{1}, \mathcal{X}_{u}^{1}), \dots, (\mathcal{X}_{p}^{N}, \mathcal{X}_{u}^{N})$$
 be mini-batches  

$$\max\left\{0, \frac{1}{n_{u}}\sum_{x \in \mathcal{X}_{u}}\ell(-g(x)) - \frac{\pi_{p}}{n_{p}}\sum_{x \in \mathcal{X}_{p}}[\ell(-g(x))]\right\}$$

$$= \max\left\{0, \sum_{i=1}^{N} \frac{1}{n_{u}}\sum_{x \in \mathcal{X}_{u}^{i}}\ell(-g(x)) - \sum_{i=1}^{N} \frac{\pi_{p}}{n_{p}}\sum_{x \in \mathcal{X}_{p}^{i}}[\ell(-g(x))]\right\}$$

$$\leq \frac{1}{N}\sum_{i=1}^{N} \max\left\{0, \frac{1}{n_{u}/N}\sum_{x \in \mathcal{X}_{u}^{i}}\ell(-g(x)) - \frac{\pi_{p}}{n_{p}/N}\sum_{x \in \mathcal{X}_{p}^{i}}[\ell(-g(x))]\right\}$$

- Safe to minimize  $\tilde{R}_{pu}(g)$  averaged over mini-batches
- You can choose whatever stochastic opt. alg. you like
  - Adam (Kingma+, ICLR'15) Or AdaGrad (Duchi+, JMLR 2011)

#### Experiments on MNIST



- P = {even digits, i.e., 0, 2, 4, 6 & 8} N = {odd digits, i.e., 1, 3, 5, 7 & 9}  $\pi_p = 0.49$   $n_p = 1,000$   $n_n = (\pi_n/2\pi_p)^2 n_p$   $n_u = 60,000$
- Model: 6-layer MLP with
   ReLU (Nair & Hinton, ICML'10)
   & BN (loffe & Szegedy, ICML'15)

#### Experiments on CIFAR10



- P = {airplane, automobile, ship & truck}
  - N = {bird, cat, deer, dog, frog & horse}

• 
$$\pi_{\rm p} = 0.40$$
  
 $n_{\rm p} = 1,000$   
 $n_{\rm n} = (\pi_{\rm n}/2\pi_{\rm p})^2 n_{\rm p}$   
 $n_{\rm u} = 50,000$ 

 Model: 13-layer CNN –
 All Convolutional Net (Springenberg+, ICLR'15)

#### Remarks

#### • $\tilde{R}_{pu}(g)$ is not a **shrinkage estimator**

- $\tilde{R}_{pu}(g) = \hat{R}_{pu}(g)$  if  $\hat{R}_{pu}(g) \ge 0$  Identical for simple models
- $\tilde{R}_{pu}(g) > \hat{R}_{pu}(g)$  if  $\hat{R}_{pu}(g) < 0 \leftarrow$  Different for complex models
- $\tilde{R}_{pu}(g)$  is biased but still consistent; bias is in  $\mathcal{O}\left(\exp\left(-\frac{1}{1/n_{p}+1/n_{u}}\right)\right)$
- Mean squared error is reduced for certain losses

• Learning is consistent; estimation error bound is in  $\mathcal{O}_p\left(-\frac{1}{\sqrt{2}}\right)$ 

$$\frac{1}{\sqrt{n_p}} + \frac{1}{\sqrt{n_u}}$$

for linear-in-parameter models

### Outline

- Introduction
- Unbiased PU Learning
- Non-negative PU Learning
- PNU Learning
- Theoretical Analyses

#### Motivation

SSL usually needs additional **distributional assumptions** 

- Cluster (Chapelle+, NIPS'02), manifold (Belkin+, JMLR 2006) ...
- U data are used in **regularization** 
  - → They will bias the classifier (learned following ERM)
- Violated assumption → U data become harmful & hurt (Cozman+, ICML'03; Sokolovska+, ICML'08; Li & Zhou, TPAMI 2015; Krijthe & Loog, PR 2017)
- PU learning has no additional distributional assumption
  - U data are used in risk evaluation
    - ➔ They will not bias the classifier

#### Illustration



#### Unbiased risk estimators (Sakai+, ICML'17)

Given 
$$\gamma \in [0,1]$$
, PN+PU learning is to minimize  
 $\hat{R}_{pn,pu}^{\gamma}(g) = (1 - \gamma)\hat{R}_{pn}(g) + \gamma \hat{R}_{pu}(g)$ 

Analogously, PN+NU learning is to minimize

$$\widehat{R}_{\text{pn,nu}}^{\gamma}(g) = (1 - \gamma)\widehat{R}_{\text{pn}}(g) + \gamma\widehat{R}_{\text{nu}}(g)$$

where 
$$\widehat{R}_{nu}(g) = \frac{\pi_n}{n_n} \sum_{x \in \mathcal{X}_n} \left[ \ell \left( -g(x) \right) - \ell \left( g(x) \right) \right] + \frac{1}{n_u} \sum_{x \in \mathcal{X}_u} \ell \left( g(x) \right)$$

Given 
$$\eta \in [-1,1]$$
, we define PNU learning as to minimize  
 $\widehat{R}_{pnu}^{\eta}(g) = \begin{cases} \widehat{R}_{pn,pu}^{\eta}(g), & \text{if } \eta \ge 0 \\ \widehat{R}_{pn,nu}^{-\eta}(g), & \text{if } \eta < 0 \end{cases}$ 
NU learning, if  $\eta = -1$   
 $\widehat{P}$ N learning, if  $\eta = 0$   
PU learning, if  $\eta = +1$ 

### Why not PU+NU?

- Motivated by theoretical comparisons (Niu+, NIPS'16)
- PN learning can never be the worst among PN, PU & NU learning
  - If n<sub>u</sub> is sufficiently large (compared with n<sub>p</sub> & n<sub>n</sub>)
     → PN learning is the second best
  - Otherwise → PN learning is the best
- PU+NU learning is not the best possible combination PN+PU & PN+NU learning are the best combinations

Note that theoretical & experimental results omitted for simplicity

### Outline

- Introduction
- Unbiased PU Learning
- Non-negative PU Learning
- PNU Learning
- Theoretical Analyses

#### To be short and intuitive

- Based on upper bounds on estimation errors
- Find simple conditions, such that
  - PU learning is likely to outperform PN learning if  $\pi_{\rm p}/\sqrt{n_{\rm p}}+1/\sqrt{n_{\rm u}}<\pi_{\rm n}/\sqrt{n_{\rm n}}$
  - NU learning is likely to outperform PN learning if  $\pi_{\rm n}/\sqrt{n_{\rm n}} + 1/\sqrt{n_{\rm u}} < \pi_{\rm p}/\sqrt{n_{\rm p}}$
- Either PU or NU learning (depending on  $\pi_p$ ,  $\pi_n$ ,  $n_p \& n_n$ ) given infinite U data will improve on PN learning

#### Simplified estimation error bounds (Niu+, NIPS'16)

- Let  $g^*/\hat{g}_{pn}$ ,  $\hat{g}_{pu} \& \hat{g}_{nu}$  be true/empirical risk minimizers
- Assume Rademacher complexities (Mohri+, Foundations of ML, 2012) of the function class G (a.k.a. model family) are in  $O(1/\sqrt{n})$
- Bounds below hold separately with probability at least 1 –

$$\begin{split} R(\hat{g}_{\mathrm{pn}}) - R(g^*) &\leq C(\delta) \cdot \left(\frac{\pi_{\mathrm{p}}}{\sqrt{n_{\mathrm{p}}}} + \frac{\pi_{\mathrm{n}}}{\sqrt{n_{\mathrm{n}}}}\right) \\ R(\hat{g}_{\mathrm{pu}}) - R(g^*) &\leq C(\delta) \cdot \left(\frac{2\pi_{\mathrm{p}}}{\sqrt{n_{\mathrm{p}}}} + \frac{1}{\sqrt{n_{\mathrm{u}}}}\right) \\ R(\hat{g}_{\mathrm{nu}}) - R(g^*) &\leq C(\delta) \cdot \left(\frac{1}{\sqrt{n_{\mathrm{u}}}} + \frac{2\pi_{\mathrm{n}}}{\sqrt{n_{\mathrm{n}}}}\right) \\ \end{split}$$
where  $C(\delta)$  is a function of  $\delta$ 

Originally for du Plessis+ (NIPS'14) but can be extended to du Plessis+ (ICML'15)

#### Finite-sample comparisons (Niu+, NIPS'16)

- Tighter PU bound than PN bound, if  $\alpha_{pu,pn} = \frac{\pi_p / \sqrt{n_p + 1 / \sqrt{n_u}}}{\pi_p / \sqrt{n_p}} < 1$
- Tighter NU bound than PN bound, if  $\alpha_{nu,pn} = \frac{\pi_n / \sqrt{n_n} + 1 / \sqrt{n_u}}{\pi_p / \sqrt{n_p}} < 1$
- $\alpha_{pu,pn}$  is monotonically decreasing in  $n_p \& n_u$ 
  - More U data → PU improves
  - More P data → PN improves too, but PU improves more!
- Many other implications discussed as well
  - When  $n_{\rm p}$ ,  $n_{\rm n}$  &  $n_{\rm u}$  are proportional
  - When  $n_{\rm p}/n_{\rm n} \approx \pi_{\rm p}/\pi_{\rm n}$  is further imposed

#### Asymptotic comparisons (Niu+, NIPS'16)

- In practice, we may find PU worse than PN and  $\alpha_{pu,pn} > 1$ Give up PU? Collect more U data (in order to improve PU)?
- Assume  $n_{\rm p}$ ,  $n_{\rm n} < \infty \& n_{\rm u} \to \infty$ , or  $\mathcal{O}(n_{\rm p}) = \mathcal{O}(n_{\rm n}) < \mathcal{O}(n_{\rm u})$ 
  - $\alpha_{pu,pn}^* \& \alpha_{nu,pn}^*$  exist as the limits of  $\alpha_{pu,pn} \& \alpha_{nu,pn}$
- Then,  $\alpha^*_{pu,pn} \cdot \alpha^*_{nu,pn} = 1$ 
  - Either  $\alpha^*_{pu,pn} < 1 \leftarrow$  Limit of PU will improve on that of PN
  - Or  $\alpha^*_{nu,pn} < 1 \leftarrow$  Limit of NU will improve on that of PN
  - Exception:  $n_p/n_n \rightarrow \pi_p^2/\pi_n^2 \leftarrow$  All 3 limits are equally good

#### Take-home message

- Unbiased PU learning and simple models given small data
  - Logistic loss or double hinge loss for convexity
- Non-negative PU learning and complex models given big data
  - Sigmoid loss for a similar shape to zero-one loss
- When PU doesn't work → Label some N to see whether PN works
  - PN works and  $\pi_p/\sqrt{n_p} > \pi_n/\sqrt{n_n} \rightarrow$  Label more P (expensive) or switch to NU/PNU if you want
  - PN works and  $\pi_p / \sqrt{n_p} < \pi_n / \sqrt{n_n} \rightarrow \text{Collect more } \cup \text{ (cheap)}$
  - PN doesn't work as well → Find a more suitable model

#### Beyond this talk ...

■ We are **two-sample** PU (Ward, Hastie, Barry, Elith & Leathwick, *Biometrics 2009*)

- Fairly different from **one-sample** PU (Elkan & Noto, KDD'08)
- According to Menon, Van Rooyen, Ong, Williamson (ICML'15)
   ∈ mutually contaminated learning (Scott, AISTATS'15)
   ∉ class-conditional noise learning (Natarajan, Dhillon, Ravikumar & Tewari, NIPS'13)
- Following work in PU learning
  - Multi-label ranking (Kanehira & Harada, CVPR'16) Fancy application
  - AUC maximization (PU & SSL) (Sakai, Niu & Sugiyama, MLJ, to appear)
  - Multi-instance binary classification (Bao, Sakai, Sato & Sugiyama, arXiv 2017)
  - Multi-class classification (Konno, UTokyo bachelor thesis 2017)