## When Deep Learning Meets Weakly-Supervised Learning

## Gang Niu

Project Assistant Professor, Univ. of Tokyo Visiting Scientist, RIKEN AIP



Joint work with Masashi Sugiyama, Marthinus C. du Plessis, and many students



## About this talk

- 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)



## About the title

- There are various definitions of weakly-supervised learning
- Here, we mean (binary/multi-class) classification, such that
  - 1. The focus is still inductive learning but not transductive
  - 2. The performance measure is still the classification error
  - 3. Not all training data are equipped with (ordinary) labels
- Two types of weakly-supervised learning
  - Semi-supervised learning (Chapelle+, Semi-Supervised Learning, 2006) where we have a small set of fully labeled training data
  - Other learning problems where no such set is available

## Semi-supervised learning

Most popular form of learning objectives to be minimized:

Empirical risk (labeled data) + Regularization (unlabeled data)

- Empirical risk is defined exactly same as in supervised learning
- Regularization is based on the local smoothness or robustness

Explicit regularization in objective function

Manifold regularization (Belkin+, JMLR 2006)

Virtual adversarial training (Miyato+, ICLR'16) Implicit regularization in training algorithm

Temporal ensembling (Laine & Aila, *ICLR'17*)

Mean teacher (Tarvainen & Valpola, *NIPS'17*)

## Other weakly-supervised learning problems

Characteristic of labeled data for training



- Hence, we need to rewrite the true risk, if we want to follow ERM
- Fundamental questions:
  - 1. How to design **unbiased risk estimators**?
  - 2. When **deep learning** is involved, is this still the right way to go?

## Question 1

## How to design unbiased risk estimators?

## Problem settings in a nutshell

**PU** learning **PN** learning **PNU** learning (i.e., supervised learning) (i.e., semi-supervised learning) X  $\mathbf{O}$ 0 0 X P, N & U data are P & N data are P & U data are available for training available for training available for training

**O** : positive data

🗙 : negative data

: unlabeled data

## Notation

Random variable	Input $X \in \mathbb{R}^d$	Output $Y \in \{\pm 1\}$			
Density	Underlying joint density $p(x, y)$				
	$p(x) \qquad p_p(x) = p(x Y = +$	-1) $p_{n}(x) = p(x Y = -1)$			
Assumed known; can be	Class-prior probability $\pi_{\rm p} = p(2)$ estimated Ramaswamy+ ( <i>ICML'16</i> ); Jai	Y = +1) n+ ( <i>NIPS'16</i> ); du Plessis+ ( <i>MLJ 2017</i> )			

Expectation		$\mathbb{E}_{\mathbf{p}}[\cdot] = \mathbb{E}_{X \sim p_{\mathbf{p}}}[\cdot]$		$\mathbb{E}_{\mathbf{n}}[\cdot] = \mathbb{E}_{X \sim p_{\mathbf{n}}}[\cdot]$	
Dataset	$X_{\rm p} = \{x_i^{\rm p}\}$	$n_{\mathrm{p}}$ i.i.d. $\sim p_{\mathrm{p}}(x)$	$\boldsymbol{\mathcal{X}}_{\mathbf{n}} = \{\boldsymbol{x}_{i}^{\mathbf{n}}\}_{i=1}^{n_{\mathbf{n}}} \overset{\text{i.i.d.}}{\sim}$	$p_{\rm n}(x)$ $\chi_{\rm u}$	$= \{x_i^{u}\}_{i=1}^{n_u} \stackrel{\text{i.i.d.}}{\sim} p(x)$

## 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



## Question 2

# When deep learning is involved, is this still the right way to go?

## Thought experiment

- Assume g is fairly flexible (such as deep NNs) and  $\forall g, R(g) > 0$
- Consider when deep learning meets weakly-supervised learning:



## Real experiment

- $\hat{R}_{pu}(g)$  is nice for training **linear-in-parameter models**
- It cannot be used for training even the shallowest MLP



## 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)

• Safe to minimize  $\tilde{R}_{pu}(g)$  averaged over mini-batches  $\max\left\{0, \frac{1}{n_{u}}\sum_{x\in\mathcal{X}_{u}}\ell(-g(x)) - \frac{\pi_{p}}{n_{p}}\sum_{x\in\mathcal{X}_{p}}[\ell(-g(x))]\right\}$   $\leq \frac{1}{N}\sum_{i=1}^{N}\max\left\{0, \frac{1}{n_{u}/N}\sum_{x\in\mathcal{X}_{u}^{i}}\ell(-g(x)) - \frac{\pi_{p}}{n_{p}/N}\sum_{x\in\mathcal{X}_{p}^{i}}[\ell(-g(x))]\right\}$ Denote by  $\Delta$ 

Given *i*-th mini-batch  $(\chi_p^i, \chi_u^i)$ 

- Gradient according to *∆* otherwise ← Corr
- Correct overfitting
- Updates are done by external SGD-like algorithms

## 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)
   & Batch Normalization (loffe & Szegedy, ICML'15)

## Experiments on CIFAR-10



- 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 which is known as
 All Convolutional Net (Springenberg+, ICLR'15)

## When deep learning meets weakly-supervised learning

### PU classification is not a special case!

## Problems that suffer (similarly to PU classification)

PU learning while R of ERM is replaced with some other criteria

**AUC** maximization (Sakai+, *MLJ to appear*) **SMI** estimation & maximization (Sakai+, arXiv 2018) for dimensionality reduction & independence test

Learning binary classifiers from two datasets (neither PN nor PU)	
---	--

Two U having different class priors (du Plessis+, TAAI'13; Menon+, ICML'15)

Pairwise similarity dataset & U (Bao+, *arXiv 2018*)

Learning multi-class classifiers from extremely noisy labels

A complementary label specifies which class  $x_i$  is not from (Ishida+, NIPS'17)