# A General Formulation for Safely Exploiting Weakly Supervised Data

## Lan-Zhe Guo    Yu-Feng Li

LAMDA Group,
National Key Lab for Novel Software Technology
Nanjing University, China

{guolz, liyf}@lamda.nju.edu.cn

# What is this paper about

Weakly supervised data is one important machine learning data

It suffers one serious issue

➤ The usage of weakly supervised data may even degenerate performance, which means, it could be outperformed by its supervised counterpart using only a small number of labeled data

Contribution of this work

In this work, we consider to learn a <u>safe</u> prediction for weakly supervised learning, where safe means it will not be worse than its supervised counterpart. We propose a general formulation and give theoretical analysis. The experiments also show quite promising results

# Outline

☐ **Introduction**

☐ Proposed Approach

☐ Experiments

☐ Conclusion

# What is Weakly Supervised Learning

- Weakly supervised learning use the data that does not require a large amount of precise label information

- For Example:
  - Label Noise Learning  [Fr´enay and Verleysen 2014]
  - Semi-Supervised Learning  [Chapelle et al. 2006)]
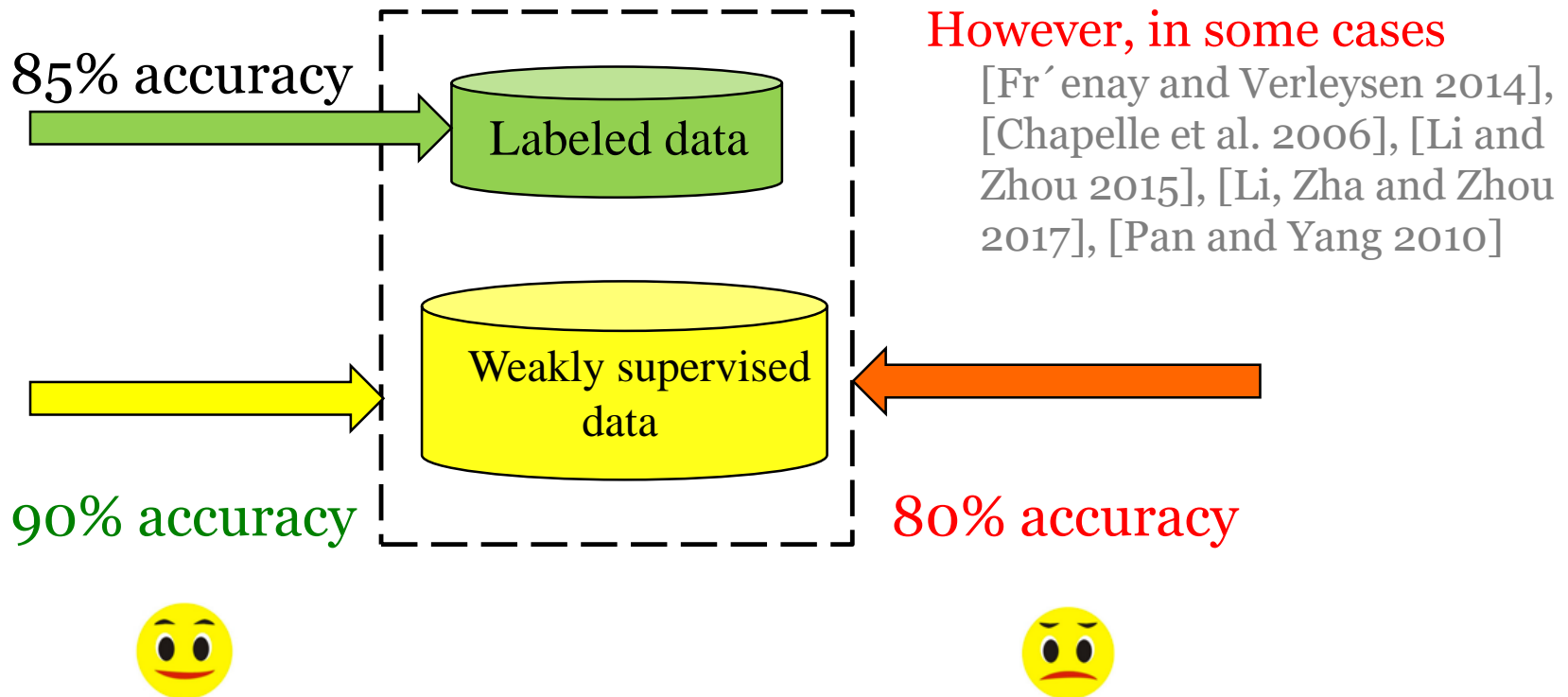  - Domain Adaptation  [Pan and Yang 2010]
  - ...

# Examples

- Label Noise Learning
  - We have only a small number of high-quality labeled data and a lot of noisy labeled data

- Semi-Supervised Learning
  - We have only limited labeled data and need to leverage a number of unlabeled data

- Domain Adaptation
  - Label information in target domain is not sufficient and we need to exploit further label information from other domains

# Weakly supervised learning is not safe

It is often expected that weakly supervised data can help improve performance since more data are used. However, it sometimes fails.

85% accuracy

Labeled data

90% accuracy

Weakly supervised data

80% accuracy

However, in some cases
[Fr´enay and Verleysen 2014], [Chapelle et al. 2006], [Li and Zhou 2015], [Li, Zha and Zhou 2017], [Pan and Yang 2010]

# Outline

☐ Introduction

☐ **Proposed Approach**

☐ Experiments

☐ Conclusion

# The Basic Setup

➢ Suppose we have a set of weakly supervised learning predictions $\{y_i\}_{i=1}^n$

➢ These base predictions can be obtained in various ways, e.g., by different type of algorithms

➢ Moreover, we can easily train a supervised method with the use of only limited labeled data and let $y_0$ denote the prediction

The goal: to learn a safe prediction $g((y_1, \cdots y_n), y_0)$, which often outperform, and will not be worse than $y_0$

http://lamda.nju.edu.cn/guolz

# A Direct Approach

Suppose we know the ground-truth $\boldsymbol{y}^*$, we can directly maximize the performance gain

$$\max_{\boldsymbol{y} \in \mathbb{H}^u} \ell(\mathbf{y}_0, \mathbf{y}^*) - \ell(\mathbf{y}, \mathbf{y}^*)$$

Trained on high-quality labeled data only

Ground Truth

http://lamda.nju.edu.cn/guolz

# However

Obviously, $y^*$ is unknown

We can construct $y^*$ with $\{y_i\}_{i=1}^n$, and consider the worst case for the requirement of safeness.

SafeW

$$\max_{\mathbf{y} \in \mathbb{H}^u} \min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell\left(\mathbf{y}_0, \sum_{i=1}^n \alpha_i \mathbf{y}_i\right) - \ell\left(\mathbf{y}, \sum_{i=1}^n \alpha_i \mathbf{y}_i\right)$$

Worst case consideration [Li and Zhou, ICML2011/TPAMI2015; Balsubramani and Freund, COLT2015]

# Three Questions about the formulation

$$\max_{\mathbf{y} \in \mathbb{H}^u} \min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell\left(\mathbf{y}_0, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i\right) - \ell\left(\mathbf{y}, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i\right)$$

- Is this formulation reasonable?

- How to setup the set of weights $\mathcal{M}$?

- How to solve it efficiently?

# Is this formulation reasonable?

$$\max_{\mathbf{y} \in \mathbb{H}^u} \min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell(\mathbf{y}_0, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i) - \ell(\mathbf{y}, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i)$$

Theoretical analysis:

**Theorem 1.** *Suppose the ground-truth $\mathbf{y}^*$ can be constructed by the base learners, i.e., $\mathbf{y}^* \in \{\mathbf{y} | \sum_{i=1}^{b} \alpha_i \mathbf{y}_i, \boldsymbol{\alpha} \in \mathcal{M}\}$. Let $\hat{\mathbf{y}}$ and $\hat{\boldsymbol{\alpha}}$ be the optimal solution to Eq.(1), we then have $\ell(\hat{\mathbf{y}}, \mathbf{y}^*) \leq \ell(\mathbf{y}_0, \mathbf{y}^*)$ and $\hat{\mathbf{y}}$ has already achieved the maximal performance gain against $\mathbf{y}_0$.*

# If the assumption is not satisfied?

If $\ell(\cdot,\cdot)$ is $\eta$-Lipschitz, i.e., $|\ell(y_1, y_2) - \ell(y_1, y_3| \leq \eta||y_2 - y_3||_1$

Let $\boldsymbol{\beta}^* = \arg\min\limits_{\boldsymbol{\beta} \in \mathcal{M}} \ell(\sum\limits_{i=1}^{n} \beta_i \mathbf{y}_i, \mathbf{y}^*)$ and $\boldsymbol{\epsilon} = \mathbf{y}^* - \sum_{i=1}^{n} \beta_i^* \mathbf{y}_i$

We have,

**Theorem 4.** *The performance gain of $\hat{\mathbf{y}}$ against $\mathbf{y}_0$, i.e.,* $\ell(\mathbf{y}_0, \mathbf{y}^*) - \ell(\hat{\mathbf{y}}, \mathbf{y}^*)$, *has a lower-bound* $-2\eta||\boldsymbol{\epsilon}||_1$.

http://lamda.nju.edu.cn/guolz

# Setup $\mathcal{M}$

For regression $C_{ij} = (y_i - u_i)^\top (y_j - u_j)$, for classification $C_{ij} = y_i^\top y_j$

We prove that $\mathbf{C\alpha}$ indicates the performance of the base learner

And we can think the base learner has a lower-bound performance
[Balsubramani and Freund 2015]

Hence, we can setup $\mathcal{M} = \{\boldsymbol{\alpha} | \mathbf{C\alpha} \geq \boldsymbol{\delta}, \mathbf{1}^\top \boldsymbol{\alpha} = 1, \boldsymbol{\alpha} \geq \mathbf{0}\}$

Moreover, if we have prior knowledge, we can setup $\mathcal{M}$ more flexible

# Optimization

Original Form:

$$\max_{\mathbf{y} \in \mathbb{H}^u} \min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell(\mathbf{y}_0, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i) - \ell(\mathbf{y}, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i)$$

Usually non-convex and not easy to solve

For regression task, we can get a convex optimization: $\min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell(\mathbf{y}_0, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i)$

For classification task and hinge loss, we can get a linear programming:

$$\min_{\boldsymbol{\alpha} \in \mathcal{M}} \ell(\mathbf{y}_0, \sum_{i=1}^{n} \alpha_i \mathbf{y}_i) + \frac{1}{u} \| \sum_{i=1}^{n} \alpha_i \mathbf{y}_i \|_1 - 1$$

Become much easier to solve

# Outline

☐ Introduction

☐ Proposed Approach

☐ **Experiments**

☐ Conclusion

# Label Noise Learning

**❑Setup**

  – 8 frequently-used datasets

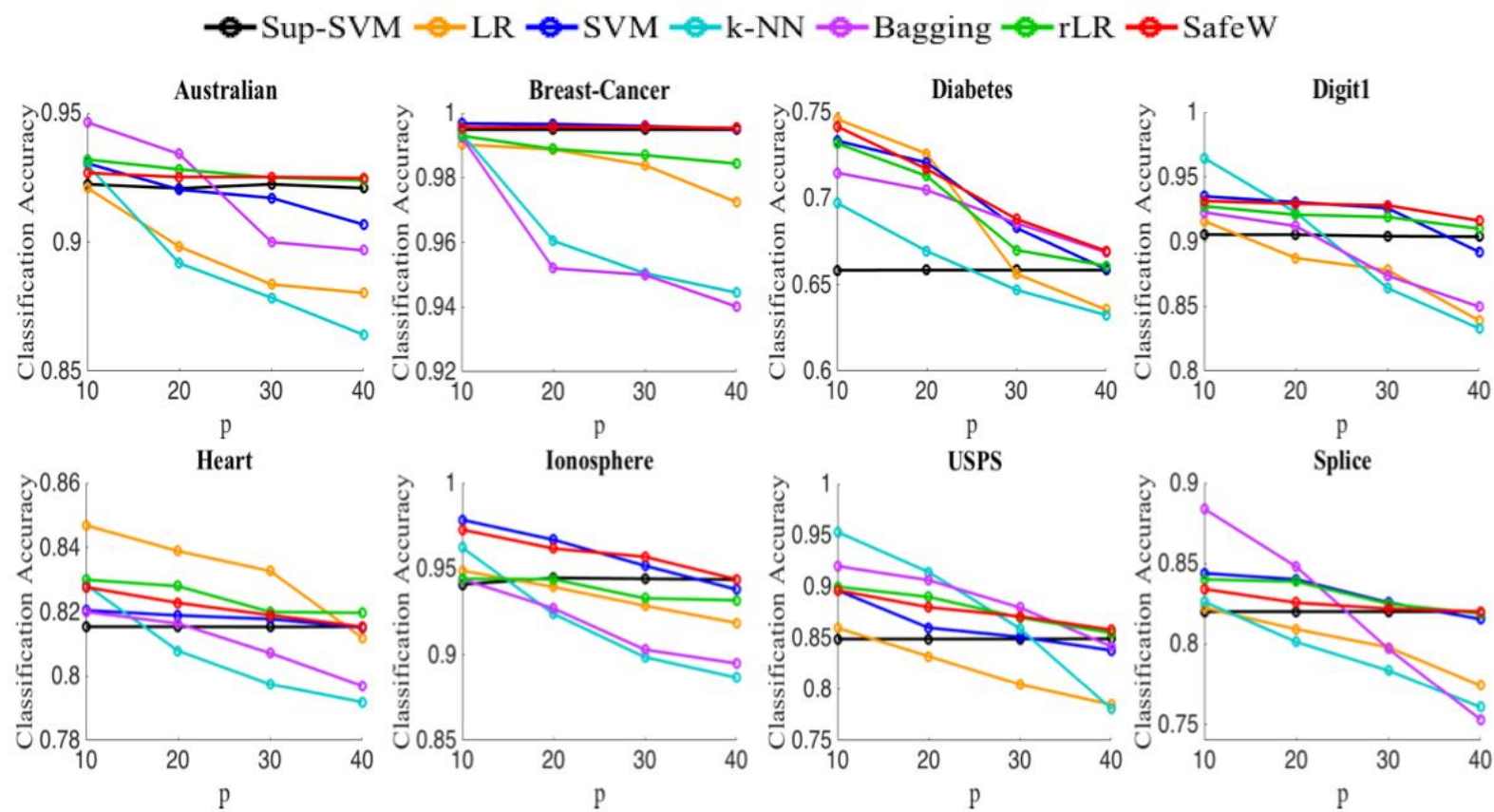  – 30% high-quality labeled data and 70% noisy data which their labels are random reversed with a probability p%.

**❑Compared Methods**

  – Baseline

   ➢ Sup-SVM

  – 2 state-of-art noisy robust methods

   ➢ Bagging  [Fr´enay and Verleysen 2014],

   ➢ rLR  [Bootkrajang and Kab´an 2012],

  – 3 traditional methods

   ➢ SVM

   ➢ k-NN

   ➢ Logistic Regression(LR)

**❑SafeW: adopted SVM, k-NN and LR as base learners.**

# Label Noise Learning

Classification accuracy on 8 datasets with p range from 10% to 40%

# Domain Adaptation

❑Setup
– 2 benchmark datasets: 20newgroup, landmine

❑Compared Methods
– Baseline
➢ Logistic Regression
– 4 transfer learning methods
➢ Original
➢ MIDA [Yan, Kou, and Zhang 2016]
➢ TCA [Pan et al. 2011]
➢ TrAdaBoost [Dai et al. 2007]

❑SafeW: adopted Original, MIDA and TCA as base learners.

# Domain Adaptation

20newsgroup: 19,997 news and 20 groups

| Dataset | Logistic Regression | Original | MIDA | TCA | TrAdaBoost | SAFEW |
|---|---|---|---|---|---|---|
| Comp vs Rec | .7028 ± .0091 | **.7492 ± .0135** | **.7961 ± .0197** | **.7940 ± .0162** | **.8077 ± .0155** | **.7956 ± .0170** |
| Comp vs Sci | .8225 ± .0662 | .7985 ± .0194 | **.8946 ± .0188** | .8255 ± .0172 | **.8583 ± .0201** | **.8925 ± .0212** |
| Comp vs Talk | .8423 ± .0685 | .8022 ± .0182 | .8231 ± .0164 | .8434 ± .0110 | .8247 ± .0143 | .8451 ± .0158 |
| Sci vs Talk | .7294 ± .1045 | .7100 ± .0121 | **.7456 ± .0164** | .7022 ± .0092 | .7166 ± .0213 | **.7468 ± .0153** |
| Rec vs Sci | .8006 ± .0758 | .7754 ± .0161 | .8033 ± .0151 | **.8440 ± .0118** | .8016 ± .0151 | **.8435 ± .0157** |
| Rec vs Talk | .8278 ± .0446 | .8276 ± .0115 | **.8566 ± .0105** | **.8580 ± .0128** | **.8415 ± .0113** | **.8579 ± .0105** |
| Average | .7876 | .7805 | .8199 | .8112 | .8084 | .8302 |
| Win/Tie/Loss against LR | | 1/2/3 | 4/1/1 | 3/2/1 | 3/2/1 | **5/1/0** |

Directly using weakly supervised data often degenerate performance while SafeW does not suffer this problem.

http://lamda.nju.edu.cn/guolz

# Domain Adaptation

Landmine: 29 domain, 9 features

| Dataset | Logistic Regression | Original | MIDA | TCA | TrAdaBoost | SAFEW |
|---|---|---|---|---|---|---|
| Domain 20 | .9215 ± .0173 | .9237 ± .0034 | **.9265 ± .0039** | .9255 ± .0045 | .9183 ± .0029 | **.9271 ± .0035** |
| Domain 21 | .9360 ± .0095 | .9310 ± .0047 | .9384 ± .0045 | .9304 ± .0051 | .9261 ± .0033 | .9396 ± .0038 |
| Domain 22 | .9594 ± .0051 | .9555 ± .0038 | .9506 ± .0065 | **.9650 ± .0017** | .9095 ± .0026 | **.9648 ± .0016** |
| Domain 23 | .9361 ± .0095 | .9310 ± .0041 | **.9424 ± .0045** | .9314 ± .0051 | **.9627 ± .0043** | **.9426 ± .0038** |
| Domain 24 | .9535 ± .0052 | .9524 ± .0029 | .9447 ± .0025 | .9432 ± .0029 | .9535 ± .0034 | .9550 ± .0024 |
| Average | .9413 | .9387 | .9405 | .9391 | .9340 | .9458 |
| Win/Tie/Loss against LR | | 0/3/2 | 2/1/2 | 1/1/3 | 1/2/2 | **3/2/0** |

We have similar observations.

http://lamda.nju.edu.cn/guolz

# Semi-Supervised Learning

❑Setup
- – 8 commonly used regression datasets
- – 10 labeled data are chosen for each dataset

❑Compared Methods
- – Baseline
  - ➢ 1NN
- – Semi-Supervised Regressor
  - ➢ Self-kNN  [Yarowsky 1995]
  - ➢ Self-LS  [Hastie, Tibshirani, and Friedman 2001]
- – Ensemble methods
  - ➢ Average
  - ➢ Safer  [Li, Zha and Zhou 2017]

❑For Average, Safer and SafeW: adopted Self-kNN(Euclidean), Self-kNN(Cosine) and Self-LS as base learners.

# Semi-Supervised Learning

Mean Square Error on 8 datasets

| Dataset | 1NN | Self-$k$NN(Euclidean) | Self-$k$NN(Cosine) | Self-LS | Average | Safer | SAFEW |
|---|---|---|---|---|---|---|---|
| abalone | $.020 \pm .010$ | $.014 \pm .005$ | $.014 \pm .003$ | $.013 \pm .004$ | $.012 \pm .003$ | $.013 \pm .005$ | $.013 \pm .005$ |
| bodyfat | $.019 \pm .005$ | $.018 \pm .006$ | $.019 \pm .005$ | $.041 \pm .013$ | $.023 \pm .009$ | $.018 \pm .007$ | $.017 \pm .005$ |
| cadata | $.083 \pm .029$ | $.063 \pm .012$ | $.058 \pm .009$ | $.056 \pm .007$ | $.057 \pm .009$ | $.060 \pm .013$ | $.057 \pm .005$ |
| cpusmall | $.024 \pm .012$ | $.027 \pm .011$ | $.028 \pm .009$ | $.025 \pm .010$ | $.024 \pm .005$ | $.025 \pm .011$ | $.024 \pm .009$ |
| housing | $.039 \pm .010$ | $.036 \pm .009$ | $.033 \pm .006$ | $.036 \pm .009$ | $.034 \pm .008$ | $.034 \pm .009$ | $.033 \pm .005$ |
| mg | $.051 \pm .009$ | $.039 \pm .006$ | $.038 \pm .006$ | $.035 \pm .015$ | $.038 \pm .014$ | $.038 \pm .006$ | $.038 \pm .006$ |
| mpg | $.022 \pm .007$ | $.020 \pm .006$ | $.018 \pm .006$ | $.021 \pm .008$ | $.020 \pm .006$ | $.019 \pm .004$ | $.018 \pm .004$ |
| pyrim | $.023 \pm .006$ | $.021 \pm .005$ | $.022 \pm .005$ | $.052 \pm .014$ | $.020 \pm .007$ | $.020 \pm .006$ | $.020 \pm .006$ |
| Ave. Mse. | $.035$ | $.030$ | $.029$ | $.035$ | $.029$ | $.030$ | $.028$ |
| Win/Tie/Loss against 1NN | | 4/3/1 | 3/4/1 | 3/3/2 | 5/2/1 | 6/2/0 | 6/2/0 |

SafeW also obtain safe predictions

# Outline

☐ Introduction

☐ Proposed Approach

☐ Experiments

☐ **Conclusion**

# Conclusion

- We propose a general formulation for safely exploiting weakly supervised data

- It has three advantages

  ➢ Has safeness guarantee for commonly used loss functions in both regression and classification tasks

  ➢ Can setup the weight of base learner flexibly

  ➢ Can be solved globally in an efficient manner

**Thanks!**

http://lamda.nju.edu.cn/guolz