

Adapting supervised classification algorithms to weak label scenarios

Miquel Perelló-Nieto¹, Raúl Santos-Rodríguez¹, and Jesús Cid-Sueiro²

¹University of Bristol, UK

²Universidad Carlos III de Madrid, Spain

Email: ¹{Miquel.PerelloNieto, enrsr}@bristol.ac.uk, ²jcid@tsc.uc3m.es

October 26, 2017



University of
BRISTOL



Universidad
Carlos III de Madrid

Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



University of
BRISTOL

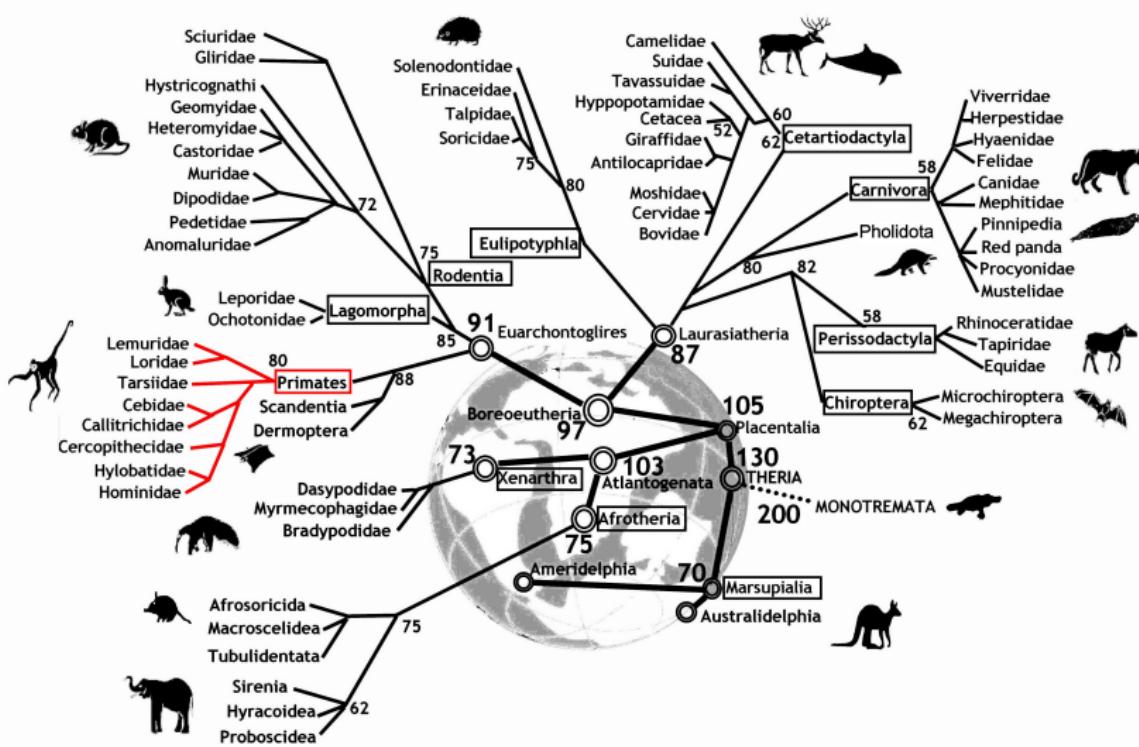
Universidad
Carlos III de Madrid

Weak labels: E.g. Captions

- One instance is one person
- Caption contains multiple labels (names)
- For every instance we have *multiple labels*
- Only one label is true



Weak labels: E.g. Supergroups



Weak labels: Other reasons

- Ambiguity
 - Lack of expertise
 - Cheaper annotations
 - Crowdsourcing
 - ▶ Fast to obtain labels
 - ▶ People may disagree



True label Weak label

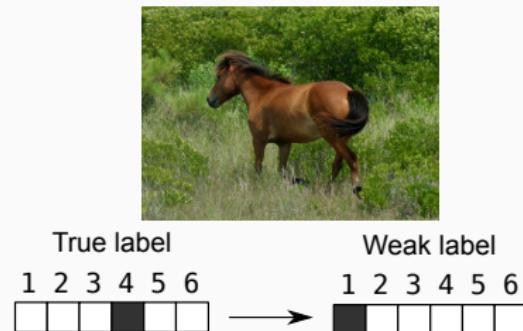
1	2	3	4	5	6
[]	[]	[]	[]	[]	[]

→

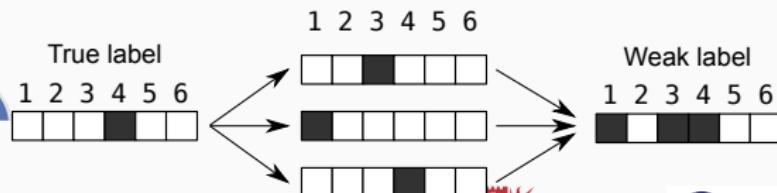
1	2	3	4	5	6
[]	[]	[]	[]	[]	[]

Weak labels: Other reasons

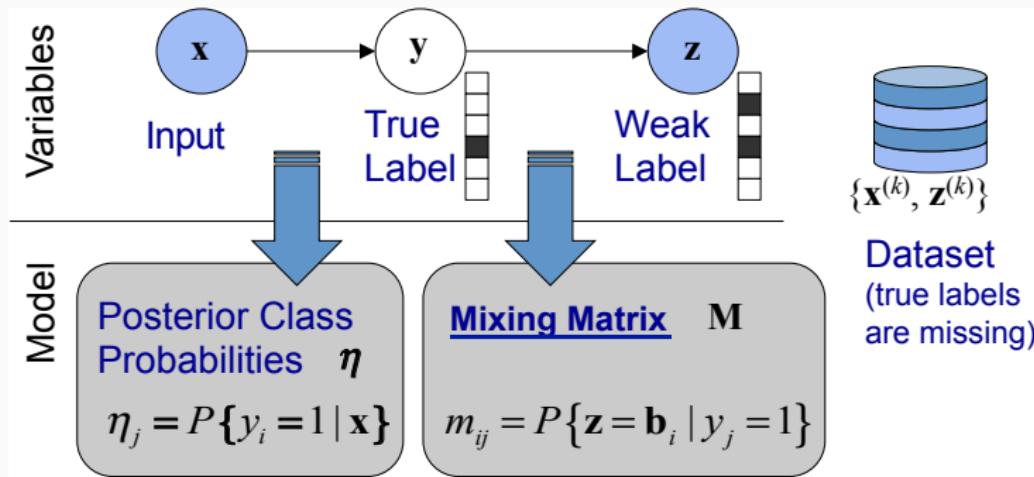
- Ambiguity
- Lack of expertise
- Cheaper annotations
- Crowdsourcing
 - ▶ Fast to obtain labels
 - ▶ People may disagree



Different Annotators



Data model



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

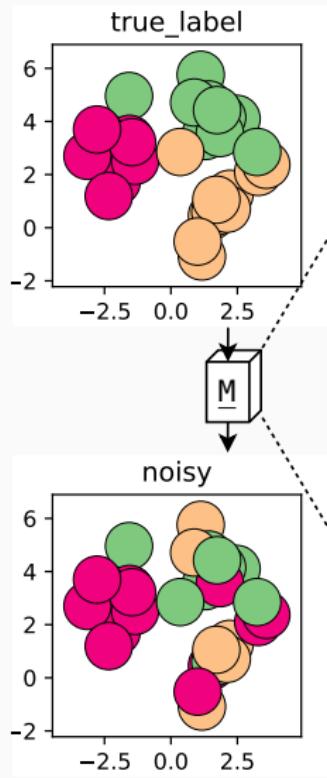
3 Experiments

- Description
- Results

4 Conclusion



Mixing matrices: Noisy labels

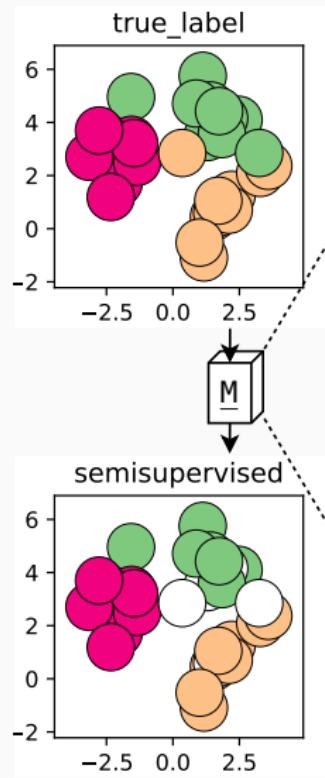


		True Label		
		0	$\beta/2$	$\beta/2$
Weak Label	0	1- β	$\beta/2$	$\beta/2$
	$\beta/2$	$\beta/2$	1- β	$\beta/2$
	$\beta/2$	$\beta/2$	$\beta/2$	1- β
	0	0	0	0
	0	0	0	0
	0	0	0	0
	0	0	0	0
	0	0	0	0

$\sum x_i = 1$

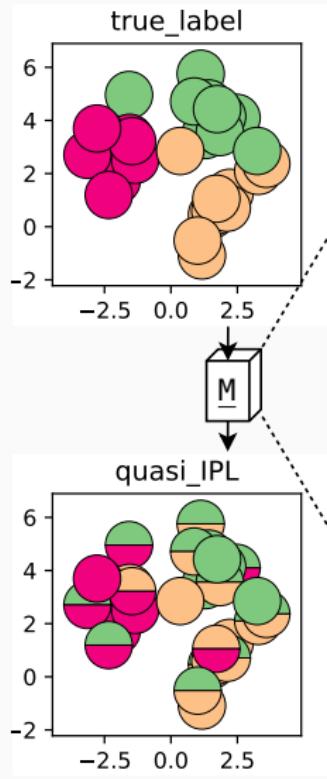
Probability of observing weak label $\blacksquare\blacksquare$ (classes 1 and 3) when the true class is $\square\blacksquare$ (class 3)

Mixing matrices: Semi-supervised learning



		True Label		
		α	β	γ
Weak Label	1 - α	0	0	0
	0	1 - β	0	0
0	0	0	1 - γ	0
1 - α	0	0	0	0
0	0	0	0	0
0	0	0	0	0
1 - β	0	0	0	0
0	0	0	0	0
1 - γ	0	0	0	0

Mixing matrices: True label + noisy label

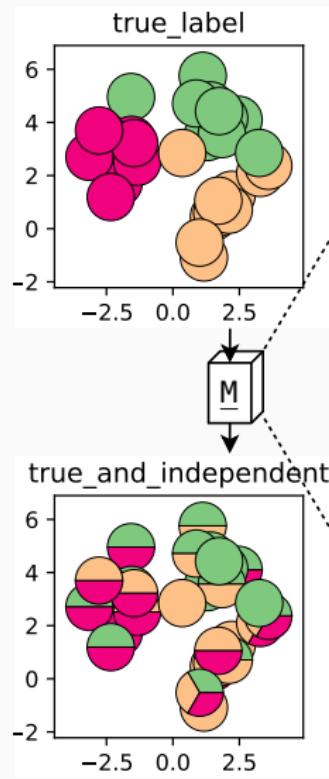


True Label

Weak Label

	0	1- β	$\beta/2$
0	0	0	0
1- β	$1-\beta$	0	0
$\beta/2$	0	$1-\beta$	0
$\beta/2$	$\beta/2$	$\beta/2$	0
$\beta/2$	$\beta/2$	0	$\beta/2$
$\beta/2$	0	$\beta/2$	$\beta/2$
1	0	0	0

Mixing matrices: True label + independent noisy label



		True Label		
		0	0	0
		$(1-\beta)^2$	0	0
Weak Label		0	$(1-\beta)^2$	0
		0	0	$(1-\beta)^2$
		$\beta(1-\beta)$	$\beta(1-\beta)$	0
		$\beta(1-\beta)$	0	$\beta(1-\beta)$
		0	$\beta(1-\beta)$	$\beta(1-\beta)$
		β^2	β^2	β^2

Index

1 Introduction

- Weak labels
- Generation of weak labels
- **Types of losses**

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

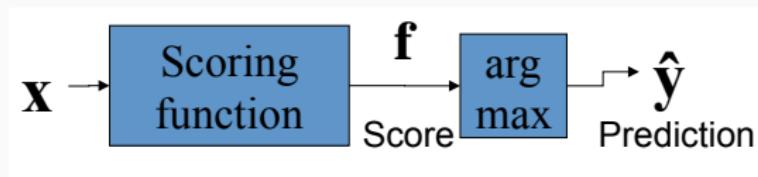
3 Experiments

- Description
- Results

4 Conclusion



Formulation



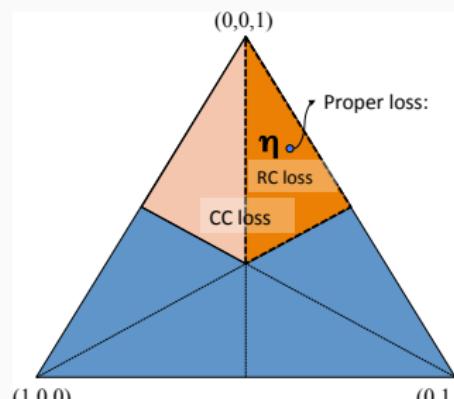
- *Conventional losses*
 - ▶ A function of the *true label* and the score: $\tilde{\Psi}(\mathbf{y}, \mathbf{f})$
in vector form, $\Psi(\mathbf{f}) = (\tilde{\Psi}(\mathbf{e}_0^c, \mathbf{f}), \dots, \tilde{\Psi}(\mathbf{e}_{c-1}^c, \mathbf{f}))$.
- *Weak losses*
 - ▶ A function of the *weak label* and the score: $\Psi(\mathbf{z}, \mathbf{f})$
in vector form, $\Psi(\mathbf{f}) = (\Psi(\mathbf{b}_0, \mathbf{f}), \dots, \Psi(\mathbf{b}_{d-1}, \mathbf{f}))$.
- We are interested in models based on Empirical Risk Minimization

$$\hat{R}_\Psi(\mathcal{S}) = \sum_{k=1}^K \Psi(\mathbf{z}_k, \mathbf{f}(\mathbf{x}_k)) \quad (1)$$



Types of losses

- Let \mathbf{f}^* be a minimizer of $\mathbb{E}_{\mathbf{y}}\{\tilde{\Psi}(\mathbf{y}, \mathbf{f})\}$
- Three types of losses:
 - Proper*: a loss for posterior class probability estimation
 $\mathbf{f}^* = \eta$
 - Ranking Calibrated (RC)*: A loss for ranking classes
 $f_i^* > f_j^* \Leftrightarrow \eta_i > \eta_j$
 - Classification Calibrated (CC)*: A loss to minimize errors
 $f_i^* > \max_{j \neq i} f_j^* \Leftrightarrow \eta_i > \max_{j \neq i} \eta_j$



Weak loss

Theorem

Consider a weak loss Ψ and a mixing matrix \mathbf{M} , and let the equivalent (conventional) loss $\tilde{\Psi}$ be given by

$$\tilde{\Psi}(\mathbf{f}) = \mathbf{M}^\top \Psi(\mathbf{f}) \quad (2)$$

- Ψ is (strictly) \mathbf{M} -proper iff $\tilde{\Psi}$ is (strictly) proper.
- Ψ is \mathbf{M} -RC iff $\tilde{\Psi}$ is RC.
- Ψ is \mathbf{M} -CC iff $\tilde{\Psi}$ is CC.



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



University of
BRISTOL



Universidad
Carlos III de Madrid

Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



University of
BRISTOL

Universidad
Carlos III de Madrid

Conventional loss into a weak loss

- We construct the weak loss in vector form as

$$\Psi(\mathbf{f}) = \tilde{\mathbf{Y}}^\top \tilde{\Psi}(\mathbf{f}) \quad (3)$$

- ▶ Where $\tilde{\mathbf{Y}}$ is a weight matrix (we call it *virtual label matrix*)
- The weak loss for a weak label \mathbf{b}_i can be written as

$$\Psi(\mathbf{f}, \mathbf{b}_i) = \tilde{\mathbf{y}}_i^\top \tilde{\Psi}(\mathbf{f}) \quad (4)$$

- ▶ where $\tilde{\mathbf{y}}_i$ is the i -th column of $\tilde{\mathbf{Y}}$
- For a conventional loss and a clean label \mathbf{y}

$$\tilde{\Psi}(\mathbf{f}, \mathbf{y}) = \mathbf{y}^\top \tilde{\Psi}(\mathbf{f}) \quad (5)$$

- i -th column of $\tilde{\mathbf{Y}}$ is a *virtual label vector*



Conventional loss into a weak loss

- We construct the weak loss in vector form as

$$\Psi(\mathbf{f}) = \tilde{\mathbf{Y}}^T \tilde{\Psi}(\mathbf{f}) \quad (3)$$

- ▶ Where $\tilde{\mathbf{Y}}$ is a weight matrix (we call it *virtual label matrix*)
- The weak loss for a weak label \mathbf{b}_i can be written as

$$\Psi(\mathbf{f}, \mathbf{b}_i) = \tilde{\mathbf{y}}_i^T \tilde{\Psi}(\mathbf{f}) \quad (4)$$

- ▶ where $\tilde{\mathbf{y}}_i$ is the i -th column of $\tilde{\mathbf{Y}}$
- For a conventional loss and a clean label \mathbf{y}

$$\tilde{\Psi}(\mathbf{f}, \mathbf{y}) = \mathbf{y}^T \tilde{\Psi}(\mathbf{f}) \quad (5)$$

- i -th column of $\tilde{\mathbf{Y}}$ is a *virtual label vector*



Conventional loss into a weak loss

- We construct the weak loss in vector form as

$$\Psi(\mathbf{f}) = \tilde{\mathbf{Y}}^T \tilde{\Psi}(\mathbf{f}) \quad (3)$$

- ▶ Where $\tilde{\mathbf{Y}}$ is a weight matrix (we call it *virtual label matrix*)
- The weak loss for a weak label \mathbf{b}_i can be written as

$$\Psi(\mathbf{f}, \mathbf{b}_i) = \tilde{\mathbf{y}}_i^T \tilde{\Psi}(\mathbf{f}) \quad (4)$$

- ▶ where $\tilde{\mathbf{y}}_i$ is the i -th column of $\tilde{\mathbf{Y}}$
- For a conventional loss and a clean label \mathbf{y}

$$\tilde{\Psi}(\mathbf{f}, \mathbf{y}) = \mathbf{y}^T \tilde{\Psi}(\mathbf{f}) \quad (5)$$

- i -th column of $\tilde{\mathbf{Y}}$ is a *virtual label vector*



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- **Weak labels into virtual labels**

3 Experiments

- Description
- Results

4 Conclusion



Mixing matrix is known

Theorem ([Cid-Sueiro et al., 2014])

Given a strictly proper loss $\tilde{\Psi}(\mathbf{f}, \mathbf{y})$ and a virtual label matrix $\tilde{\mathbf{Y}}$, the weak loss $\Psi(\mathbf{f}) = \tilde{\mathbf{Y}}^\top \tilde{\Psi}(\mathbf{f})$ is strictly \mathcal{M} -proper, for any $\mathbf{M} \in \mathcal{M}$ such that $\tilde{\mathbf{Y}}\mathbf{M} = \mathbf{I}$.

1. Compute a left inverse of the mixing matrix \mathbf{M} that we call *virtual label matrix* $\tilde{\mathbf{Y}}$
2. *Substitute* each weak label \mathbf{b}_i for its corresponding column $\tilde{\mathbf{y}}_i$
3. Train with the conventional loss $\tilde{\Psi}$



Mixing matrix is unknown: if assume quasi_IPL

- The following $\tilde{\Psi}$ loss is proper for any quasi-independent mixing matrix.
- Take any label-based proper loss, e.g. the cross entropy:

$$\tilde{\Psi}(\mathbf{f}) = \log(\mathbf{f}) \quad (6)$$

- Take $\tilde{\mathbf{y}}$ given by

$$\tilde{y}_j = \begin{cases} 1 & z_j = 1 \\ -\frac{|z|-1}{c-|z|} & z_j = 0 \end{cases} \quad (7)$$

- Take $\Psi(\mathbf{z}, \mathbf{f}) = \tilde{\mathbf{y}}^\top \tilde{\Psi}(\mathbf{f})$



Mixing matrix is unknown: if assume IPL

- Mixing matrix M of the form IPL doesn't have proper weak loss
- M-RC and M-CC losses if we train with the weak labels



Implementation

- Common implementations of Cross-entropy *ignore* all except the true class
- They can not deal with *negative labels*
- We used *Brier Score* as is proper and does not have the previous problems
- We use the *Moore-Penrose* pseudoinverse $\tilde{\mathbf{Y}} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$

Implementation

- Common implementations of Cross-entropy *ignore* all except the true class
- They can not deal with *negative labels*
- We used *Brier Score* as is proper and does not have the previous problems
- We use the *Moore-Penrose* pseudoinverse $\tilde{\mathbf{Y}} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$



Implementation

- Common implementations of Cross-entropy *ignore* all except the true class
- They can not deal with *negative labels*
- We used *Brier Score* as is proper and does not have the previous problems
- We use the *Moore-Penrose* pseudoinverse $\tilde{\mathbf{Y}} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$



Implementation

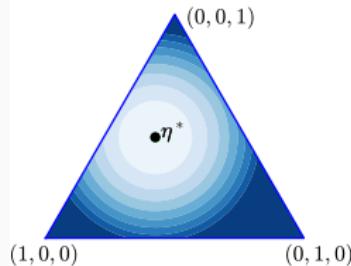
- Common implementations of Cross-entropy *ignore* all except the true class
- They can not deal with *negative labels*
- We used *Brier Score* as is proper and does not have the previous problems
- We use the *Moore-Penrose* pseudoinverse $\tilde{\mathbf{Y}} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T$



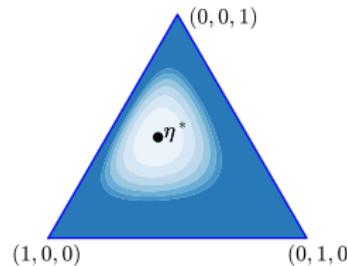
Convexity

- if the loss $\tilde{\Psi}(\mathbf{f}, \tilde{\mathbf{y}}_i)$ is **M-proper**, its conditional expectation is a proper loss

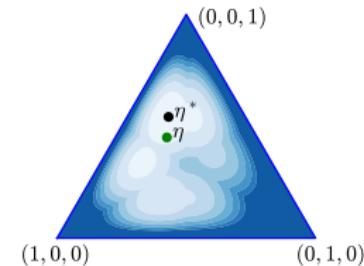
Brier



CE



OSL



- Other method: Optimistic Superset Learning

Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion

Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



University of
BRISTOL



Universidad
Carlos III de Madrid

Datasets

- 31 real-world datasets
 - ▶ 3 to 20 classes
 - ▶ No more than 11,000 instances
 - ▶ All standardized
 - ▶ Available from openml.org
- Generating synthetic weak labels
 - ▶ 5 types of random mixing matrices
 - ▶ 4 levels of noise (increasing β)
 - ▶ Every combination repeated 10 times



Datasets

- 31 real-world datasets
 - ▶ 3 to 20 classes
 - ▶ No more than 11,000 instances
 - ▶ All standardized
 - ▶ Available from openml.org
- Generating synthetic weak labels
 - ▶ 5 types of random mixing matrices
 - ▶ 4 levels of noise (increasing β)
 - ▶ Every combination repeated 10 times



Models and learning method

- Models
 - ▶ *LR*: Logistic Regression
 - ▶ *FNN*: Feed-Forward Neural Network
- Learning method
 - ▶ *Superv*: Using true labels
 - ▶ *Mproper*: Knowing M
 - ▶ *Weak*: Using weak labels (assuming IPL)
 - ▶ *quasiIPL*: Using virtual labels assuming quasi_IPL
 - ▶ *OSL*: Optimistic Superset Loss [Hüllermeier and Cheng, 2015]



Models and learning method

- Models
 - ▶ *LR*: Logistic Regression
 - ▶ *FNN*: Feed-Forward Neural Network
- Learning method
 - ▶ *Superv*: Using true labels
 - ▶ *Mproper*: Knowing **M**
 - ▶ *Weak*: Using weak labels (assuming IPL)
 - ▶ *quasiIPL*: Using virtual labels assuming quasi_IPL
 - ▶ *OSL*: Optimistic Superset Loss [Hüllermeier and Cheng, 2015]



Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

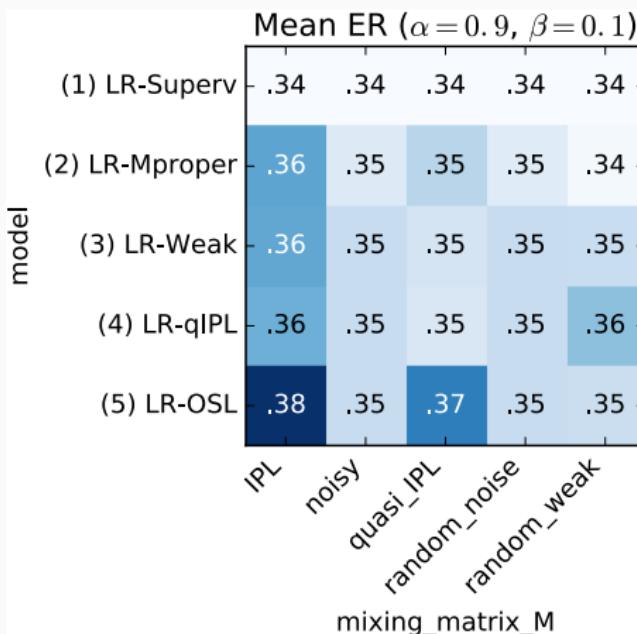
3 Experiments

- Description
- Results

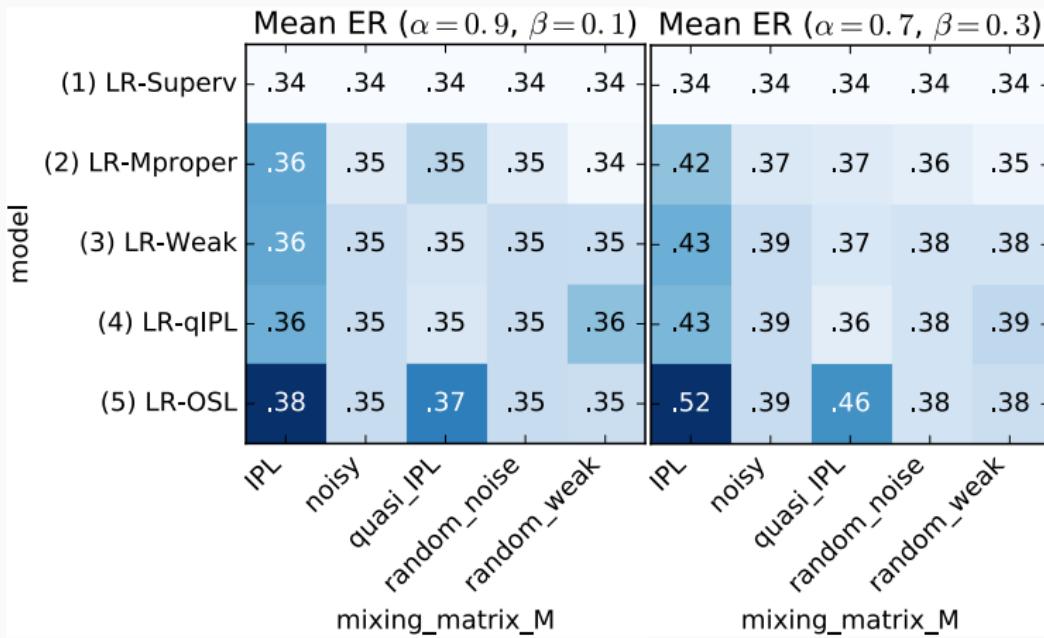
4 Conclusion



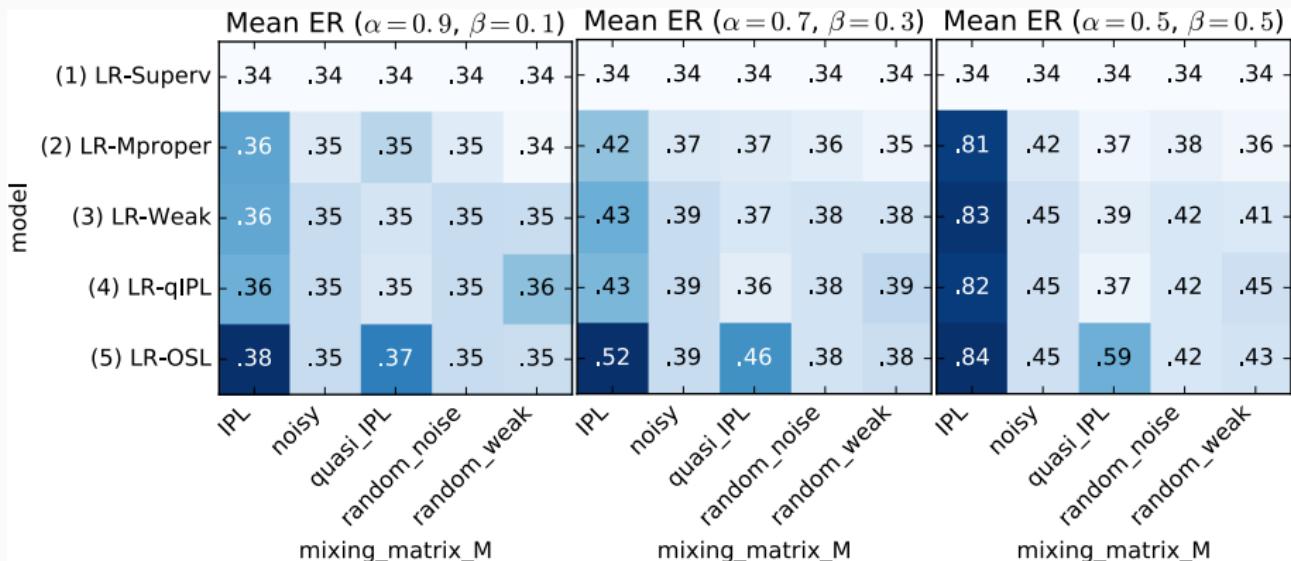
Error rate Logistic Regression



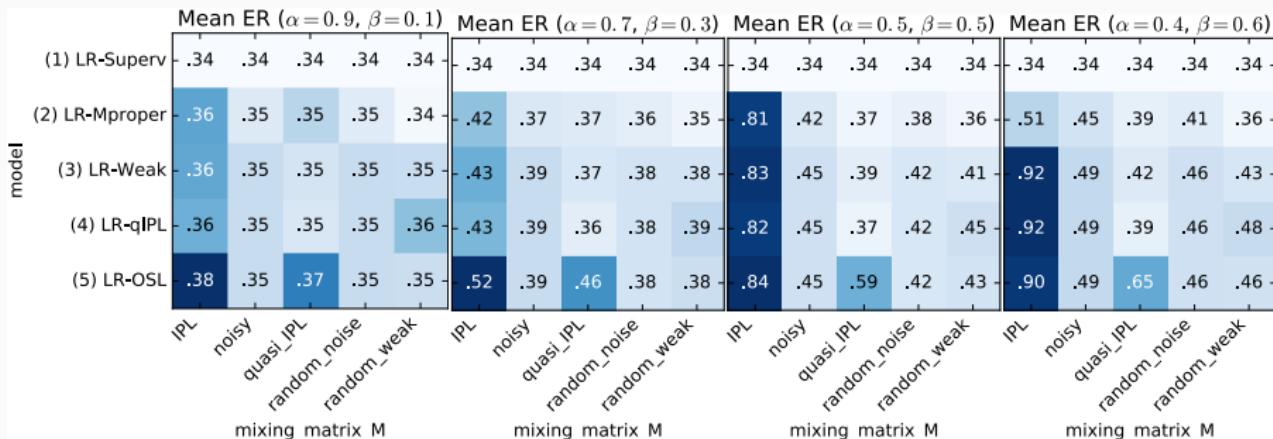
Error rate Logistic Regression



Error rate Logistic Regression



Error rate Logistic Regression



Error rate Feed-Forward Neural Network

	Mean ER ($\alpha=0.9, \beta=0.1$)					Mean ER ($\alpha=0.7, \beta=0.3$)					Mean ER ($\alpha=0.5, \beta=0.5$)					
model	(1) FNN-Superv	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30	.30
(2) FNN-Mproper	.32	.30	.31	.30	.30	.41	.34	.33	.32	.31	.83	.41	.33	.36	.32	
(3) FNN-Weak	.31	.31	.30	.30	.30	.38	.34	.32	.33	.32	.83	.41	.34	.38	.36	
(4) FNN-qIPL	.31	.31	.30	.30	.31	.39	.34	.32	.33	.35	.82	.41	.34	.38	.41	
(5) FNN-OSL	.33	.31	.32	.30	.30	.52	.34	.43	.33	.33	.85	.41	.56	.38	.37	

Index

1 Introduction

- Weak labels
- Generation of weak labels
- Types of losses

2 Method

- Conventional loss into a weak loss
- Weak labels into virtual labels

3 Experiments

- Description
- Results

4 Conclusion



Conclusion

- *Simple method* just changing the weak labels into virtual labels
- Adapts *existing classifiers* based on empirical minimization of losses
- Possible to add priors to improve results
- Best results if we can estimate the mixing matrix \mathbf{M}



Future work

- Real datasets with weak labels
- Merging datasets containing different types of noise
- Different noise depending on the input space





Cid-Sueiro, J., García-García, D., and Santos-Rodríguez, R. (2014).

Consistency of losses for learning from weak labels.

In *European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 197–210. Springer Berlin Heidelberg.



Hüllermeier, E. and Cheng, W. (2015).

Superset learning based on generalized loss minimization.

In *European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 260–275.



Adapting supervised classification algorithms to weak label scenarios

Miquel Perelló-Nieto¹, Raúl Santos-Rodríguez¹, and Jesús Cid-Sueiro²

¹University of Bristol, UK

²Universidad Carlos III de Madrid, Spain

Email: ¹{Miquel.PerelloNieto, enrsr}@bristol.ac.uk, ²jcid@tsc.uc3m.es

October 26, 2017



University of
BRISTOL



Universidad
Carlos III de Madrid