

# Consistency of Losses for Learning from Weak Labels

J. Cid-Sueiro<sup>1</sup>   D. García-García<sup>2</sup>   R. Santos-Rodríguez<sup>3</sup>

<sup>1</sup>Universidad Carlos III de Madrid, Spain

<sup>2</sup>Commonwealth Bank of Australia

<sup>3</sup>Universidad de Valencia, Spain

European Conf. on Machine Learning, Sept. 2014

# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - Mixing matrices
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - Mixing matrices
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

# Weak labels

- On some multiclass classification problems, we may have a training dataset with partial information about the class of the data.
- Weak label: a subset of labels, probably including the true class
  - Photo captions
  - Multiple specialized binary annotators.
  - Hierarchical classes and superclass labels



Alonso, Vettel and Webber in India

Alonso OR Vettel  
OR Webber

# Data model

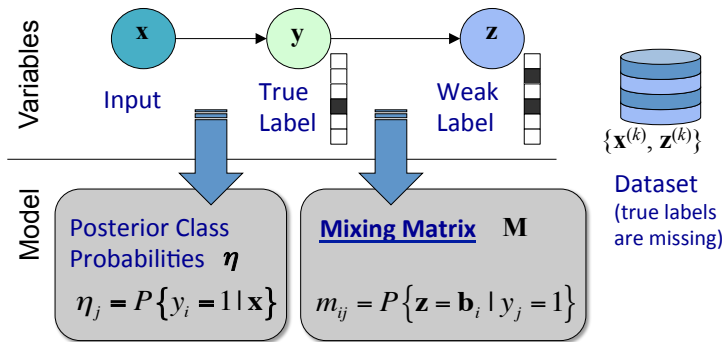


Figure: Data model

# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - **Mixing matrices**
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

# Mixing matrices.












- Supervised learning:

		True Label		
Weak Label		0	0	0
		1	0	0
		0	1	0
		0	0	1
		0	0	0
		0	0	0
		0	0	0
		0	0	0

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

# Mixing matrices.

- Unsupervised learning:

		True Label		
				
Weak Label		$\alpha$	$\beta$	$\gamma$
		$1-\alpha$	0	0
		0	$1-\beta$	0
		0	0	$1-\gamma$
		0	0	0
		0	0	0
		0	0	0
		0	0	0

# Mixing matrices.

- Noisy labels:

		True Label		
Weak Label		0	0	0
		$1-\alpha$	$\beta/2$	$\gamma/2$
		$\alpha/2$	$1-\beta$	$\gamma/2$
		$\alpha/2$	$\beta/2$	$1-\gamma$
		0	0	0
		0	0	0
		0	0	0
		0	0	0

# Mixing matrices.

- True label and one noisy label:

		True Label		
Weak Label		$1-\alpha$	0	0
		0	$1-\beta$	0
		0	0	$1-\gamma$
		$\alpha/2$	$\beta/2$	0
		$\alpha/2$	0	$\gamma/2$
		0	$\beta/2$	$\gamma/2$
		0	0	0

# Mixing matrices.

- True label and independent noisy labels:

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

# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - Mixing matrices
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

# Inference from partial labels.

- If  $\mathbf{M}$  is known, since

$$\mathbf{p} = \mathbf{M}\boldsymbol{\eta}$$

we can estimate  $\boldsymbol{\eta}$  from  $\mathbf{p}$  given  $\mathbf{M}$  as

$$\boldsymbol{\eta} = \mathbf{M}^+\mathbf{p}$$

- Two major problems:

- ① Even knowing  $\mathbf{M}$ , estimating (a high dimensional)  $\mathbf{p}$  from data may be infeasible
- ②  $\mathbf{M}$  is generally unknown, and cannot be estimated from  $\mathcal{S}$ .

# Inference from partial labels.

- If  $\mathbf{M}$  is known, since

$$\mathbf{p} = \mathbf{M}\boldsymbol{\eta}$$

we can estimate  $\boldsymbol{\eta}$  from  $\mathbf{p}$  given  $\mathbf{M}$  as

$$\boldsymbol{\eta} = \mathbf{M}^+ \mathbf{p}$$

- Two major problems:
  - 1 Even knowing  $\mathbf{M}$ , estimating (a high dimensional)  $\mathbf{p}$  from data may be infeasible
  - 2  $\mathbf{M}$  is generally unknown, and cannot be estimated from  $\mathcal{S}$ .

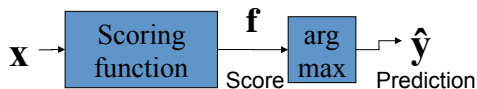
# Limitations of Learning from Weak Labels

- Any learning algorithm is effective for some weak labelling models only.
- Our goals:
  - Determine conditions or requirements for learnability
  - To analyze comparatively the impact of these requirements in probability estimation, class ranking and classification.
  - To propose general methods to design loss functions for learning.

# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - Mixing matrices
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

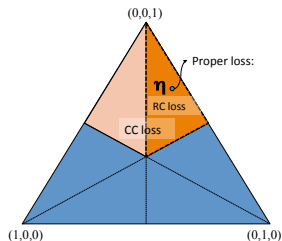
# Losses for weak labels.



- Label-based losses
  - A function of the true label and the score:  $\tilde{\Psi}(\mathbf{y}, \mathbf{f})$
- Weak losses
  - A function of the weak label and the score:  $\Psi(\mathbf{z}, \mathbf{f})$ .

# Proper, RC and CC losses

- Let  $\mathbf{f}^*$  be a minimizer of  $\mathbb{E}_{\mathbf{z}}\{\Psi(\mathbf{z}, \mathbf{f})\}$
- Three types of losses:
  - ① **M-Proper**: a loss for posterior class probability estimation
    - $\mathbf{f}^* = \boldsymbol{\eta}$
  - ② **M-Ranking Calibrated (RC)**: A loss for ranking classes
    - $f_i^* > f_j^* \Leftrightarrow \eta_i > \eta_j$
  - ③ **M-Classification Calibrated (CC)**: A loss to minimize errors
    - $f_i^* > \max_{j \neq i} f_j^* \Leftrightarrow \eta_i > \max_{j \neq i} \eta_j$



# Outline

- 1 Learning from Weak Labels
  - Weak labels
  - Mixing matrices
  - The inference problem
- 2 Weak Loss Functions
  - Proper, RC and CC losses
  - Consistency
- 3 Constructing weak losses

## Theorem

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

$$\tilde{\Psi}(\mathbf{f}) = \mathbf{M}^T \Psi(\mathbf{f}) \quad (1)$$

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

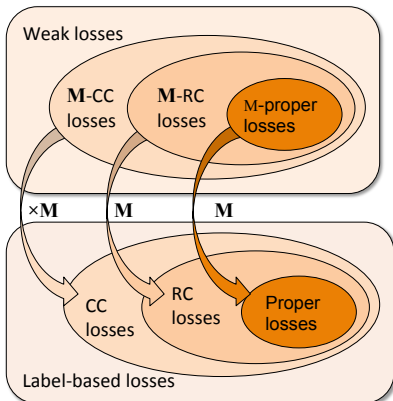


Figure: From Weak to labeled

- Let us proceed in a constructive way: we can design weak losses as a linear transformation of a label based loss

$$\Psi = \tilde{\mathbf{Y}}^T \tilde{\Psi}$$

- Maximal sets:

- $\mathcal{Q}_{\text{proper}}(\tilde{\mathbf{Y}})$  = the set of all mixing matrices,  $\mathbf{M}$  s.t.  $\Psi$  is  $\mathbf{M}$ -proper.
- $\mathcal{Q}_{rc}(\tilde{\mathbf{Y}})$  = the set of all mixing matrices,  $\mathbf{M}$  s.t.  $\Psi$  is  $\mathbf{M}$ -RC.
- $\mathcal{Q}_{cc}(\tilde{\mathbf{Y}})$  = the set of all mixing matrices,  $\mathbf{M}$  s.t.  $\Psi$  is  $\mathbf{M}$ -C.

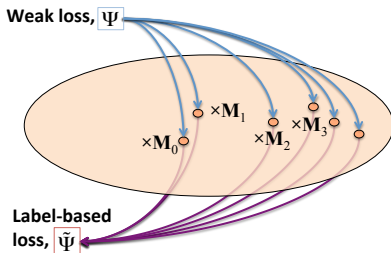
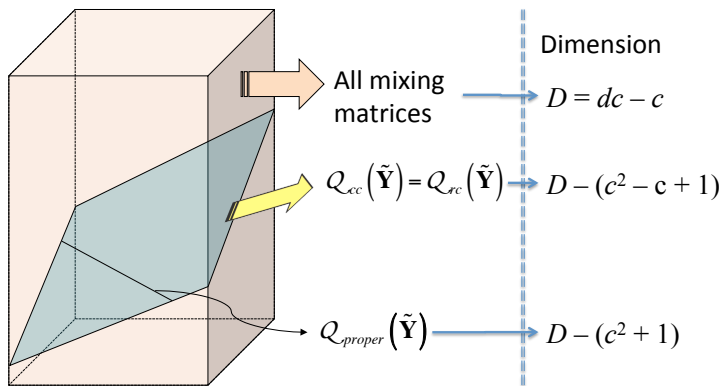


Figure: From Weak to labeled

## Theorem

- $\mathcal{Q}_{proper}(\tilde{\mathbf{Y}}) = \{\mathbf{M} | \tilde{\mathbf{Y}}\mathbf{M} = \lambda \mathbf{I}, \lambda > 0\}$
- $\mathcal{Q}_{cc}(\tilde{\mathbf{Y}}) = \mathcal{Q}_{rc}(\tilde{\mathbf{Y}}) = \{\mathbf{M} | \tilde{\mathbf{Y}}\mathbf{M} = \lambda \mathbf{I} + \mathbf{1}_c \mathbf{v}^T, \lambda \in \mathbb{R}^+, \mathbf{v} \in \mathbb{R}^c\}$



- Two be more practical: a proper loss recommendation:

- 1 Take any label-based proper loss, e.g. the cross entropy:

$$\tilde{\Psi}(\mathbf{f}) = \log(\mathbf{f}) \quad (2)$$

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

$$\tilde{y}_j = \begin{cases} 1 & z_j = 1 \\ -\frac{|z|-1}{c-|z|} & z_j = 0 \end{cases} \quad (3)$$

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

- $\Psi$  is proper for any quasi-independent mixing models.

- ... and a RC or CC loss recommendation:
  - 1 Take any label-based RC or CC loss (e.g. the cross entropy)
  - 2 Take  $\Psi(\mathbf{z}, \mathbf{f}) = \mathbf{z}^\top \tilde{\Psi}(\mathbf{f})$
- $\Psi$  is RC or CC for any mixing model given by

$$P(\mathbf{z}|y_c = 1) = \alpha^{z_m}(1 - \alpha)^{1-z_m}\beta^{|\mathbf{z}|-1}(1 - \beta)^{c-|\mathbf{z}|} \quad (4)$$

for any  $\alpha > \beta$ .

- $\Psi$  is convex