

Semi-Supervised Learning from Crowds Using Deep Generative Model

Kyohei Atarashi
Hokkaido University
atarashi_k@complex.ist.hokudai.ac.jp

Satoshi Oyama
Hokkaido University/RIKEN AIP
oyama@ist.hokudai.ac.jp

Masahito Kurihara
Hokkaido University
kurihara@ist.hokudai.ac.jp

AAAI-18

Introduction

- Previous Methods: only crowd data
- A more practical problem setting: crowd data + unlabeled data
 - Semi-supervised learning from crowds

Previous Methods

- Infer true labels
 - Crowd label -> “True” label -> Supervised classification model
 - MV, DS,
 - Additional information are lost.
- Learning from crowds
 - Crowd label -> Classifier
 - True labels as latent variables
 - EM
- Semi-Supervised Learning
 - SSVAE-M2

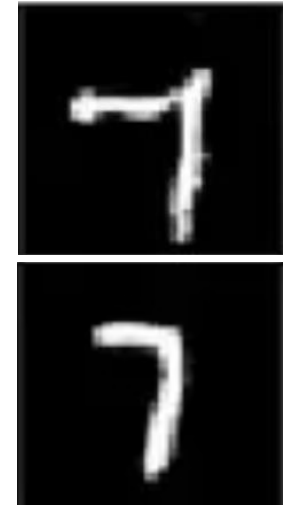
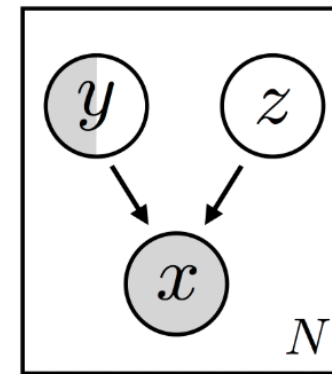
SSVAE-M2

- Model

- $p(x, y, z) = p(x|y, z)p(y)p(z)$
- $p(y) = \text{Cat}(y|\pi)$
- $p(z) = \mathcal{N}(z|0, I)$
- $p_\theta(x|z, y) = \mathcal{N}(x|\mu(z, y), \sigma^2(z, y))$ decoder

- Inference

- $q_\phi(\mathbf{y} | \mathbf{x}) = \text{Cat}(\mathbf{y} | \alpha_\phi(\mathbf{x}))$ classifier
- $q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{z} | \mu_\phi(\mathbf{x}, \mathbf{y}), \sigma_\phi^2(\mathbf{x}, \mathbf{y}))$ encoder



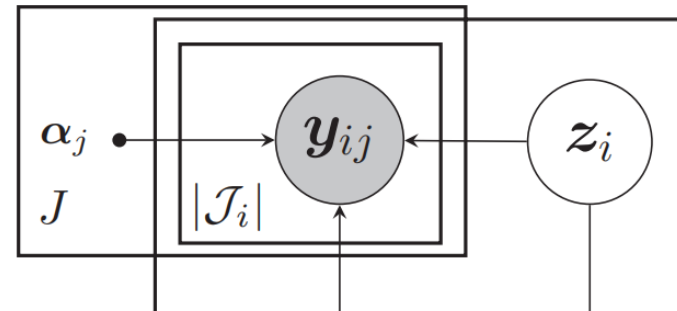
- Loss

- $$\mathcal{J}^\alpha = \mathcal{J} + \alpha \mathbb{E}_{\tilde{p}_l(\mathbf{x}, \mathbf{y})} [-\log(q_\phi(\mathbf{y} | \mathbf{x}))] = \mathcal{L}(\mathbf{x}, \mathbf{y}) + \mathcal{U}(\mathbf{x}) + \alpha \mathbb{E}_{\tilde{p}_l(\mathbf{x}, \mathbf{y})} [-\log(q_\phi(\mathbf{y} | \mathbf{x}))]$$

Model

$$\begin{aligned}
 & p(\mathcal{X}_c, \mathcal{Y}_c, \mathcal{T}_c, \mathcal{Z}_c, \mathcal{X}_u, \mathcal{T}_u, \mathcal{Z}_u) \\
 &= p(\mathcal{X}_c, \mathcal{Y}_c, \mathcal{T}_c, \mathcal{Z}_c) p(\mathcal{X}_u, \mathcal{T}_u, \mathcal{Z}_u) \\
 &= p(\mathcal{X}_c | \mathcal{T}_c, \mathcal{Z}_c) p(\mathcal{Y}_c | \mathcal{T}_c, \mathcal{Z}_c) p(\mathcal{T}_c) p(\mathcal{Z}_c) \\
 &\times p(\mathcal{X}_u | \mathcal{T}_u, \mathcal{Z}_u) p(\mathcal{T}_u) p(\mathcal{Z}_u) \\
 &= \prod_{i \in \Lambda_c} \prod_{j \in \mathcal{J}_i} p(\mathbf{x}_i | \mathbf{t}_i, \mathbf{z}_i) p(\mathbf{y}_{ij} | \mathbf{t}_i, \mathbf{z}_i) p(\mathbf{t}_i) p(\mathbf{z}_i) \\
 &\times \prod_{i \in \Lambda_u} p(\mathbf{x}_i | \mathbf{t}_i, \mathbf{z}_i) p(\mathbf{t}_i) p(\mathbf{z}_i),
 \end{aligned}$$

$$\begin{aligned}
 p(\mathbf{t}_i) &= \text{Cat}(\mathbf{t}_i | \boldsymbol{\pi}) \\
 p(\mathbf{z}_i) &= \mathcal{N}(\mathbf{z}_i | \mathbf{0}, \mathbf{I}), \\
 p(\mathbf{x}_i | \mathbf{t}_i, \mathbf{z}_i) &= \mathcal{N}(\mathbf{x}_i | \mu_{\boldsymbol{\theta}}(\mathbf{t}_i, \mathbf{z}_i), \text{diag}(\sigma_{\boldsymbol{\theta}}^2(\mathbf{t}_i, \mathbf{z}_i))) \text{ decoder} \\
 p(\mathbf{y}_{ij} | \mathbf{t}_i, \mathbf{z}_i) &= \text{Cat}(\mathbf{y}_{ij} | f_{\alpha_j}(\mathbf{t}_i, \mathbf{z}_i))
 \end{aligned}$$



of label-

Optimization and Inference

- Computing the posteriors of \mathbf{z} and \mathbf{t} are intractable, cannot use EM
- Cannot get a classifier directly

$$q(\mathbf{t}, \mathbf{z} \mid \mathbf{x}) = \underbrace{q(\mathbf{t} \mid \mathbf{x})}_{\text{classifier}} \underbrace{q(\mathbf{z} \mid \mathbf{t}, \mathbf{x})}_{\text{encoder}}$$

- The marginal log-likelihood

$$\log p(\mathcal{X}_c, \mathcal{Y}_c, \mathcal{X}_u) = \log p(\mathcal{X}_c, \mathcal{Y}_c) + \log p(\mathcal{X}_u)$$

Unlabeled data

$$\begin{aligned} \log p(\mathcal{X}_u) &= \log \int \int p(\mathcal{X}_u, \mathcal{T}_u, \mathcal{Z}_u) d\mathbf{t} d\mathbf{z} \\ &\geq \sum_{i \in \Lambda_u} \mathbb{E}_{q(\mathbf{t}_i, \mathbf{z}_i \mid \mathbf{x}_i)} \left[\log p(\mathbf{x}_i \mid \mathbf{t}_i, \mathbf{z}_i) \right] \\ &\quad - \sum_{i \in \Lambda_u} \text{KL}[q(\mathbf{t}_i, \mathbf{z}_i \mid \mathbf{x}_i) \parallel p(\mathbf{t}_i, \mathbf{z}_i)] \\ &:= -\mathcal{U}(\mathcal{X}_u). \end{aligned}$$

Labeled data

$$\begin{aligned} \log p(\mathcal{X}_c, \mathcal{Y}_c) &= \log \int \int p(\mathcal{X}_c, \mathcal{Y}_c, \mathcal{T}_c, \mathcal{Z}_c) d\mathbf{t} d\mathbf{z} \\ &\geq \mathbb{E}_{q(\mathcal{T}_c, \mathcal{Z}_c \mid \mathcal{X}_c)} \left[\log \frac{p(\mathcal{X}_c, \mathcal{Y}_c, \mathcal{T}_c, \mathcal{Z}_c)}{q(\mathcal{T}_c, \mathcal{Z}_c \mid \mathcal{X}_c)} \right] \\ &= \sum_{i \in \Lambda_c} \sum_{j \in \mathcal{J}_i} \mathbb{E}_{q(\mathbf{t}_i, \mathbf{z}_i \mid \mathbf{x}_i)} \left[\log p(\mathbf{x}_i, \mathbf{y}_{ij} \mid \mathbf{t}_i, \mathbf{z}_i) \right] \\ &\quad - \sum_{i \in \Lambda_c} \text{KL}[q(\mathbf{t}_i, \mathbf{z}_i \mid \mathbf{x}_i) \parallel p(\mathbf{t}_i, \mathbf{z}_i)] \\ &:= -\mathcal{C}(\mathcal{X}_c, \mathcal{Y}_c), \end{aligned}$$

Optimization and Inference

- Pseudo classification loss

$$-\sum_{i \in \Lambda_c} \log(q(\bar{\mathbf{y}}_i | \mathbf{x}_i))$$
$$\bar{\mathbf{y}}_i = \frac{1}{|\mathcal{J}_i|} \sum_{j \in \mathcal{J}_i} \mathbf{y}_{ij}$$

- The prior of \mathbf{t} for crowdsourced data

$$p(\mathbf{t}_i) = \text{Cat}(\mathbf{t}_i | \bar{\mathbf{y}}_i)$$

- Final objective function

$$\mathcal{J} = \mathcal{C}(\mathcal{X}_c, \mathcal{Y}_c) + \mathcal{U}(\mathcal{X}_u) - \alpha \cdot \sum_{i \in \Lambda_c} \log(q(\bar{\mathbf{y}}_i | \mathbf{x}_i))$$

Experiments

Table 1: Experimental results for the MNIST dataset.

	PC labeling	DS labeling
Proposed	0.9646	0.8647
MV-LR	0.7162	0.6705
MV-MLP	0.7128	0.6547
PC	0.7756	0.6063
PC-GP	0.7754	0.6055
M2	0.9421	0.8445

Table 2: Experimental results for the Rotten Tomatoes dataset. The accuracies of the MVhard (this model is same as the MV-LR), the MVsoft, the Raykar, the Rayker w/prior and the MA-LR are taken from Rodrigues, Pereira, and Ribeiro (2013).

Model	Accuracy	Model	Accuracy
Proposed	0.7290	MVhard	0.7027
MV-MLP	0.6979	MVsoft	0.7165
PC	0.7261	Raykar	0.4867
PC-GP	0.7263	Raykar w/prior	0.7078
M2	0.7096	MA-LR	0.7240

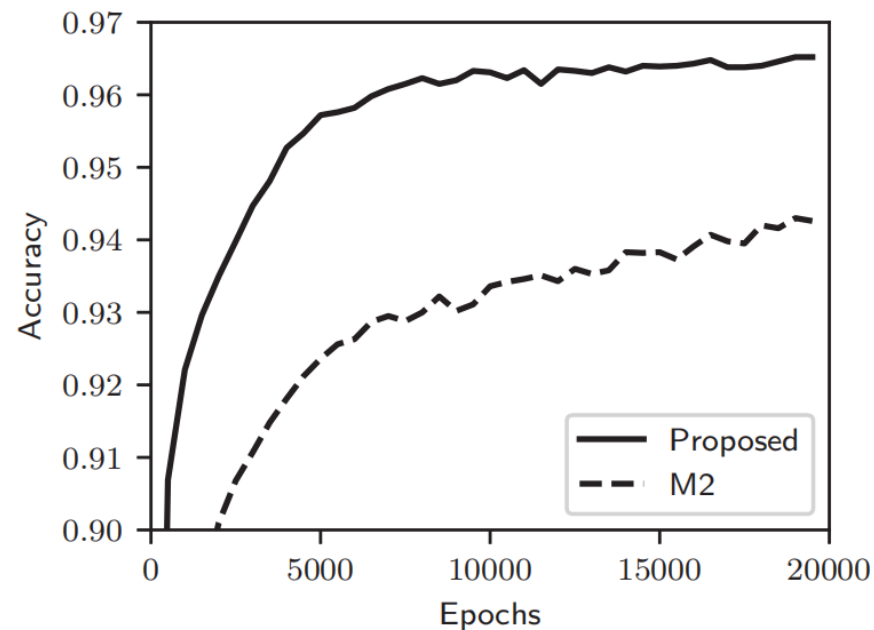


Figure 3: Accuracy for the PC labeling evaluated every 100 epochs between the proposed model and the M2 model.

Experiments

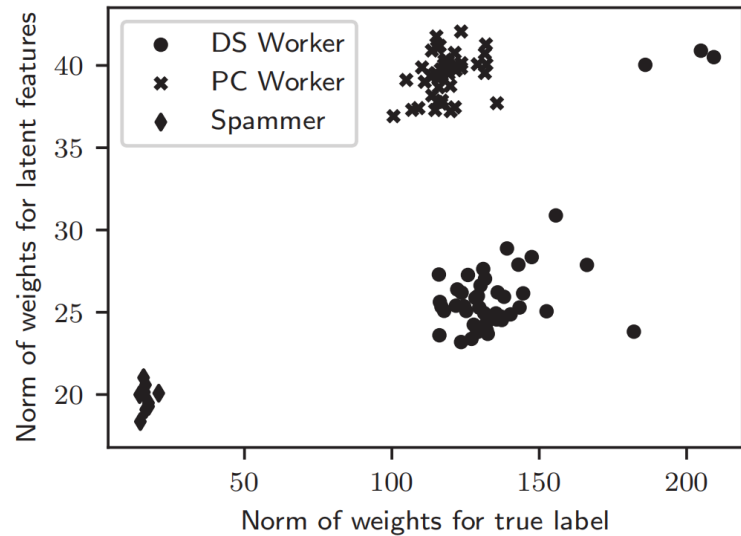


Figure 4: Estimated $\|\alpha^{(t)}\|_1$ and $\|\alpha^{(z)}\|_1$ of 45 DS workers, 45 PC workers and 10 spammers.

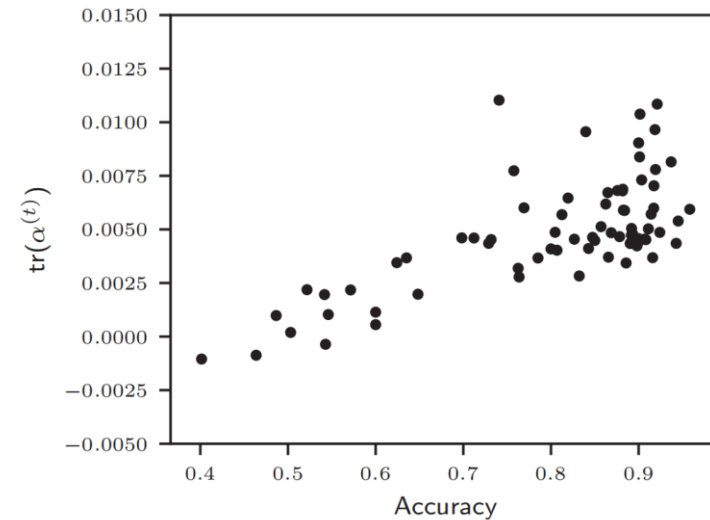


Figure 5: Estimated $\text{diag}(\alpha^{(t)})$ and the accuracy of each worker who labels more than 100 data.