

# Fairness

[NIPS 2010: Discriminative Clustering by Regularized Information Maximization (RIM)]

Problem: learn a probabilistic discriminative classifier from an unlabeled dataset

$$X = (x_1, \dots, x_N), \text{ where } x_i = (x_{i1}, \dots, x_{iD})^T \in \mathbb{R}^D \xrightarrow{\text{learn}} p(y|x, W)$$

RIM:  $F(p(y|x, W); X; \lambda) \Rightarrow$  evaluate the suitability of  $p(y|x, W)$

① cluster assumption (decision boundaries  $\times$  dense)  $\rightarrow$  confidence level  
conditional entropy  $\frac{1}{N} \sum_i H\{p(y|x_i, W)\}$

(On unsupervised assumption, it can be reduced by removing decision boundaries)

② class balance (avoid degenerate solutions)  $\rightarrow$  distribution level

$$\text{empirical label distribution: } \hat{p}(y; W) = \frac{1}{N} \sum_i p(y|x_i, W)$$

entropy  $H\{\hat{p}(y; W)\}$

Combine ① + ②  $I_W\{y; x\} = H\{\hat{p}(y; W)\} - \frac{1}{N} \sum_i H\{p(y|x_i, W)\}$   
mutual information

( $I_W\{y; x\}$  may be trivially maximized by a conditional model that classifies each data point  $x_i$  into its own category  $y_i$ )

③ classifier complexity (penalty)

$$F(W; X, \lambda) = I_W\{y; x\} - R(W; \lambda)$$

Learning a conditional distributional distribution for  $y$  that preserves information from the data set, subject to a complexity penalty.

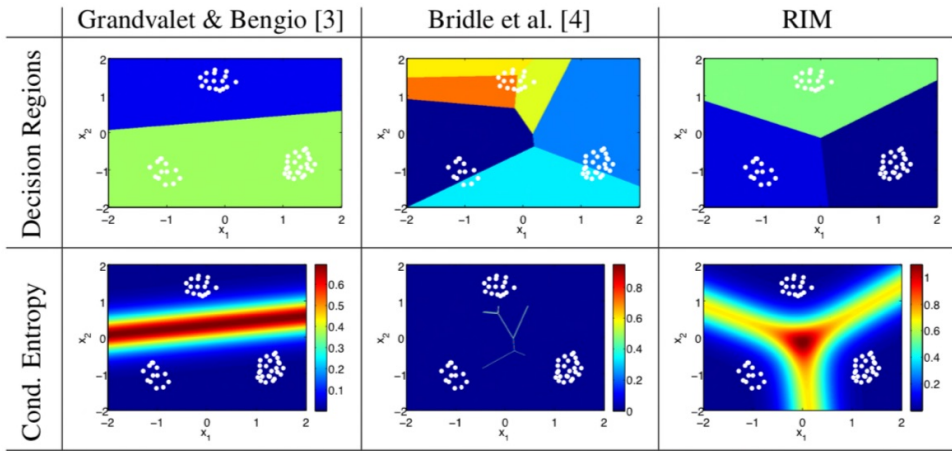


Figure 1: Example unsupervised multilogit regression solutions on a simple dataset with three clusters. The top and bottom rows show the category label  $\arg \max_y p(y|x, \mathbf{W})$  and conditional entropy  $H\{p(y|x, \mathbf{W})\}$  at each point  $\mathbf{x}$ , respectively. We find that both class balance and regularization terms are necessary to learn unsupervised classifiers suitable for multi-class clustering.

Since  $H\{\hat{p}(y; \mathbf{W})\} = \log K - \text{KL}\{\hat{p}(y; \mathbf{W}) \| \mathcal{U}\}$ , then

$$F(\mathbf{W}; \mathbf{X}, \lambda) = -\frac{1}{N} \sum_i H\{p(y|x_i, \mathbf{W})\} - \underbrace{\text{KL}\{\hat{p}(y; \mathbf{W}) \| \mathcal{U}\}}_{\text{class balance} \rightarrow D(y; \tau)} - R(\mathbf{W}; \lambda)$$

Others

$$F(\mathbf{W}; \mathbf{X}, \lambda) = I_{\mathbf{W}}(x_i; y) - H\{\hat{p}(y; \mathbf{W}) \| \mathcal{D}(y; \tau)\} - R(\mathbf{W}; \lambda)$$

In SSL,  $S(\mathbf{W}; \tau, \lambda) = \underbrace{\tau I_{\mathbf{W}}(y; x_i)}_{D_U} - R(\mathbf{W}; \lambda) + \underbrace{\sum_i \log p(y_i | x_i^L, \mathbf{W})}_{D_L}$

[IJCNN 2020: Pseudo-Labeling and Confirmation Bias in Deep SSL]

Two Reg' to improve convergence.

①  $R_A = \sum_{c=1}^C p_c \log\left(\frac{p_c}{h_c}\right)$  distribution level  $D_{KL}(\mathcal{U} \| \Gamma_{\mathbf{W}})$   
prior distribution of  $c$  mean softmax probability of model for  $c$

$$\textcircled{2} R_H = -\frac{1}{N} \sum_{i=1}^N \sum_{c=0}^C h_0^c(x_i) \log(h_0^c(x_i)) \quad \text{confidence level (entropy regularization)}$$

→ mix up loss

$$l = l^* + \lambda_A R_A + \lambda_H R_H$$

[ICLR 2023: FreeMatch]

Self-adaptive fairness (distribution level)

$\tilde{p}_t(c) \rightarrow$  estimate of the expectation of prediction distribution over  $D_V$ .

In RIM, we use  $H(\hat{p}(y|w))$  to realize class balance. ( $H(\mathbb{E}_u[p_m(y|u)])$  max)

We optimize CE of  $\tilde{p}_t$  and  $\bar{p} = \mathbb{E}_{p_B}[p_m(y|\Omega(u_b))]$  as an estimate.

(We expect  $D_{KL}(\tilde{p}_t \| \bar{p}) \leq 0$ , that is  $H(\tilde{p}_t, \bar{p}) = H(\tilde{p}_t) = H(\bar{p})$ )

? The underlying pseudo-label distribution may not be uniform ( $\tilde{p}_t, \bar{p} \neq U$ )

(max  $H(\tilde{p}_t, \bar{p}) \Rightarrow$  tend to be uniform distribution)

★ modulate the fairness objective in a self-adaptive way (normalize)

$$\bar{p} = \frac{1}{\mu_B} \sum_{b=1}^{\mu_B} \mathbb{1}(\max(q_b) \geq \tau_t(\arg \max(q_b))) Q_b$$

$$\tilde{h} = \text{Hist}_{\mu_B}(\mathbb{1}(\max(q_b) \geq \tau_t(\arg \max(q_b))) \hat{Q}_b)$$

$$\tilde{h}_t = \lambda \tilde{h}_{t-1} + (1-\lambda) \text{Hist}_{\mu_B}(\hat{q}_b)$$

The self-adaptive fairness (SAF)  $L_f$  at the  $t$ -th iteration is:

$$L_f = - \mathcal{H} \left( \underbrace{\text{SumNorm} \left( \frac{\tilde{p}_t}{\tilde{h}_t} \right)}_{\downarrow}, \underbrace{\text{SumNorm} \left( \frac{\bar{p}}{\tilde{h}} \right)}_{\downarrow} \right)$$

Nearly Uniform