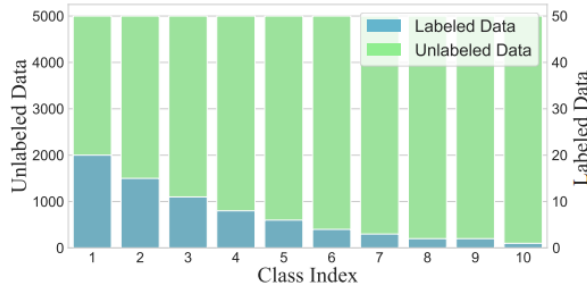# Towards Semi-supervised Learning with Non-random Missing Labels [ICLR 2023]

**Challenge**: label Missing Not At Random (MNAR)



Training under MNAR, the model increasingly favors some classes, seriously affecting the pseudo-rectifying procedure.

_Pseudo-rectifying_ : the change of the label assignment decision made by the SSL model for the same sample according to the knowledge learned at each new epoch
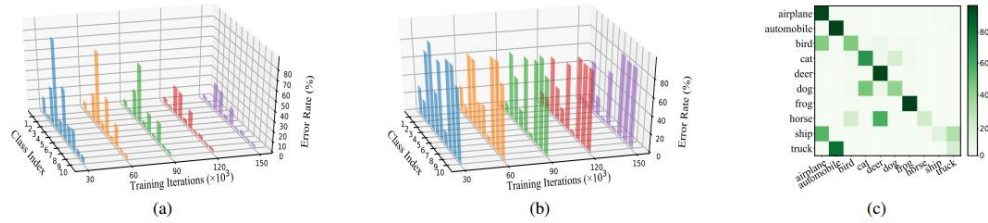


Figure 2: Results of FixMatch [27] in MNAR and the conventional SSL setting (_i.e._, balanced labeled and unlabeled data). The models are trained on CIFAR-10 with WRN-28-2 backbone [37]. (a) and (b): Class-wise pseudo-label error rate. (c): Confusion matrix of pseudo-labels. In (b) and (c), experiments are conducted with the setting of Fig. 1, whereas in (a) with the conventional setting. The label amount used in (a) is the same as that in (b) and (c).

**Motivation**: The mispredicted pseudo-labels for each class are often concentrated in a few classes, rather than scattered across all other classes. Inspired by this, we argue that it is feasible to guide pseudo-rectifying from the **class level**, i.e., pointing out the latent direction of class transition based on its current class prediction only.
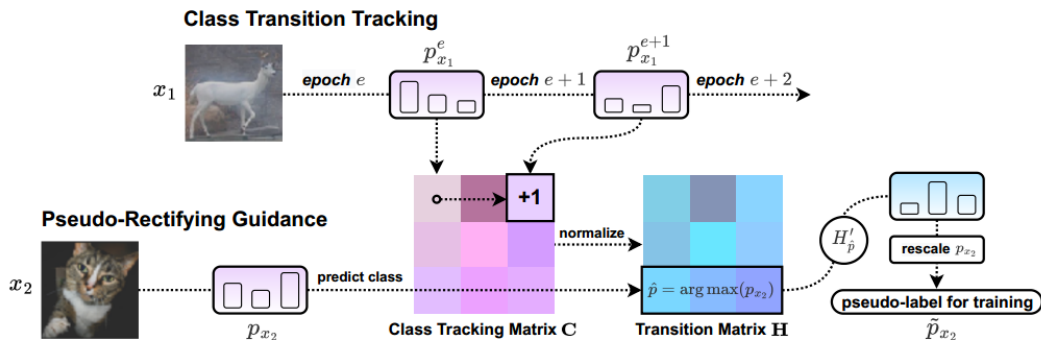
**Overview of Pseudo-Rectifying Guidance (PRG):**



Figure 3: Overview of PRG. Class tracking matrix $\mathbf{C}$ is obtained by tracking the class transitions of pseudo-labels (_e.g._, $p_{x_1}$ for sample $x_1$) between epoch $e$ and epoch $e+1$ caused by pseudo-rectifying procedure (Eq. (5)). The Markov random walk defined by transition matrix $\mathbf{H}$ (each row $H_i$ represents the transition probability vector corresponding to class $i$) is modeled on the graph constructed over $\mathbf{C}$. Generally, given a pseudo-label, _e.g._, $p_{x_2}$ for sample $x_2$, class- and batch-rescaled $\mathbf{H}$ (_i.e._, $\mathbf{H}'$) is utilized to provide the class-level pseudo-rectifying guidance for $p_{x_2}$ according to its class prediction $\hat{p} = \arg\max(p_{x_2})$ (Eqs. (6)~(7)). Finally, the rescaled pseudo-label $\tilde{p}_{x_2}$ is used for the training.

**Method:**

X: input space

Y: label space ( k classes)

M: label missing indicator ( m = 1 missing )

p = f(x,θ): pseudo-label

*Pseudo-Rectifying Guidance:*

Pseudo-rectifying process: the change on p by the next epoch $p^{e+1} = g_\theta(p^e)$

In general, the Pseudo-Rectifying Guidance (PRG) can be described as:

$$\tilde{p}^{e+1} = \text{Normalize}(\eta \circ g_\theta(p^e)), \qquad (3)$$

where $\circ$ is Hadamard product, scaling weight vector $\eta \in \mathbb{R}^k_+$

It is also feasible to guide pseudo-rectifying at class level. Hence, we define rectifying weight matrix $\mathbf{A} \in \mathbb{R}^{k \times k}_+$, where each row $A_i$ is representing the rectifying weight vector corresponding to class i. The class-level pseudo-rectifying guidance:

$$\tilde{p}^{e+1} = \text{Normalize}(A_{\hat{p}^{e+1}} \circ g_\theta(p^e)). \qquad (4)$$

Next, we will introduce a simple and feasible way to obtain an effective A for PRG to improve the pseudo-labels.

*Class Transition Tracking:*

Firstly, we consider building a fully connected graph G in class space Y. This graph is constructed by adjacency matrix C, where each element $C_{ij}$ represents the frequency of class transitions that occur from class i to class j.

$$C_{ij}^{(n)} =$$
$$\left| \left\{ \hat{p}^{(b)} \mid \hat{p}^{(b),e} = i, \hat{p}^{(b),e+1} = j, i \neq j, b \in \{1, ..., B_U\} \right\} \right|, \qquad (5)$$

$$n \in \{1, ..., N_B\} \text{ and } C_{ii}^{(n)} = 0$$

$$C_{ij} = \sum_{n=1}^{N_B} C_{ij}^{(n)} / N_B$$

Hereafter, we define the Markov random walk along the nodes of G, which is characterized by its transition matrix H. and $H_{ij}$ represents the transition probability for the class prediction transits from class i at epoch e to class j at epoch e + 1.

Since $H_{ii} = 0$ is unreasonable, we control the probability that does not transition class by setting $H_{ii} = \frac{\alpha}{k-1}$.

In addition, to better provide class-level guidance, we scale each element in H by

$$H'_{ij} = \frac{H_{ij}}{\frac{L_j}{\sum_{d=1}^{k} L_d}} \qquad (6)$$

where $L \in R^k_+$ and $L_i$ records the number of class predictions belonging to class i averaged on last $N_B$ batches.

Finally, we have

$$\tilde{p}^{e+1} = \text{Normalize}(H'_{\hat{p}^{e+1}} \circ g_\theta(p^e)), \qquad (7)$$

It is also feasible to use the class transition driven by $\widehat{p^e}$ to revise $p^{e+1}$.

---

**Algorithm 1:** PRG: Class Transition Tracking Based **P**seudo-**R**ectifying Guidance

**Input:** class tracking matrices $\mathcal{C} = \{\mathbf{C}^{(i)}; i \in (1, ..., N_B)\}$, labeled training dataset $D_L$, unlabeled training dataset $D_U$, model $\theta$, label bank $\{l^{(i)}; i \in (1, ..., n_T - n_L)\}$

1 **for** $n = 1$ **to** MaxIteration **do**
2    From $D_L$, draw a mini-batch $\mathcal{B}_L = \{(x_L^{(b)}, y_L^{(b)}); b \in (1, ..., B)\}$
3    From $D_U$, draw a mini-batch $\mathcal{B}_U = \{(x_U^{(b)}); b \in (1, ..., B_U)\}$
4    $\mathbf{H} = \text{RowWiseNormalize}(\text{Average}(\mathcal{C}))$                             // Construct transition matrix
5    $H'_{ij} = \frac{H_{ij}}{\frac{L_j}{\sum_{d=1}^{k} L_d}}$                           // Rescale **H** at class-level
6    **for** $b = 1$ **to** $B_U$ **do**
7       $p^{(b)} = f_\theta(x_U^{(b)})$                          // Compute model prediction
8       $\text{idx} = \text{Index}(x_U^{(b)})$               // Obtain the index of $x_U^{(b)}$ in $D_U$
9       $\hat{p}^{(b)} = \arg\max(p^{(b)})$               // Compute class prediction
10      **if** $l^{(\text{idx})} \neq \hat{p}^{(b)}$ **then**
11         $C^{(n)}_{l^{(\text{idx})} \hat{p}^{(b)}} = C^{(n)}_{l^{(\text{idx})} \hat{p}^{(b)}} + 1$      // Perform class transition tracking
12         $l^{(\text{idx})} = \hat{p}^{(b)}$
13      **end**
14       $\tilde{p}^{(b)} = \text{Normalize}(H'_{\hat{p}^{(b)}} \circ p^{(b)})$     // Perform pseudo-rectifying guidance
15    **end**
16    $\mathcal{L}_L, \mathcal{L}_U = \text{FixMatch}\left(\mathcal{B}_L, \mathcal{B}_U, \{\tilde{p}^{(b)}; b \in (1, ..., B_U)\}\right)$    // Run an applicable SSL learner
17    $\theta = \text{SGD}(\mathcal{L}_L + \mathcal{L}_U, \theta)$             // Update model parameters $\theta$
18 **end**

---

**for** $b = 1$ **to** $B_U$ **do**
  $p^{(b)} = f_\theta(x_U^{(b)})$                    // Compute model prediction
  $\text{idx} = \text{Index}(x_U^{(b)})$         // Obtain the index of $x_U^{(b)}$ in $D_U$
  $\tilde{p}^{(b)} = \text{Normalize}(H'_{l^{(\text{idx})}} \circ p^{(b)})$    // Perform pseudo-rectifying guidance
  $\hat{p}^{(b)} = \arg\max(p^{(b)})$           // Compute class prediction
  **if** $l^{(\text{idx})} \neq \hat{p}^{(b)}$ **then**
    $C^{(n)}_{l^{(\text{idx})} \hat{p}^{(b)}} = C^{(n)}_{l^{(\text{idx})} \hat{p}^{(b)}} + 1$    // Perform class transition tracking
    $l^{(\text{idx})} = \hat{p}^{(b)}$
  **end**
**end**