

---

# Deep Neural Network With Massive Learned Knowledge

By Abhilasha

---

**Teacher :**

$$q^*(\mathbf{Y}) \propto p_{\theta}(\mathbf{Y}|\mathbf{X}) \exp \left\{ C \sum_l \lambda_l f_l(\mathbf{X}, \mathbf{Y}) \right\},$$

**Student :**

$$\begin{aligned} \boldsymbol{\theta}^{(t+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{n=1}^N & (1 - \pi) \ell(\mathbf{y}_n, \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}_n)) \\ & + \pi \ell(\mathbf{s}_n^{(t)}, \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}_n)), \end{aligned}$$

1.  $f_l(\mathbf{X}, \mathbf{Y}) = r_l(\mathbf{X}, \mathbf{Y}) - 1$ , where  $r_l(\mathbf{X}, \mathbf{Y})$  in  $[0, 1]$
2.  $r_l(\mathbf{X}, \mathbf{Y})$  is fully specified a priori and fixed throughout training.
3.  $\lambda_l$  has to be manually set.

1. To substantially extend the scope of knowledge used in the framework, we introduce learnable modules  $\varphi$  in the knowledge expression denoted as  $f\varphi$ .
2. We aim to learn the knowledge by determining  $\varphi$  from data.
3. As any meaningful knowledge is expected to be consistent with the observations, a straightforward way is then to directly optimize against the training data:

$$\Phi^* = \operatorname{argmax}_{\phi} \left( \frac{1}{N} \sum_n f_{\phi}(x_n, y_n) \right)$$

$$\begin{aligned}\phi^{(t+1)} = \arg \max_{\phi \in \Phi} \frac{1}{N} \sum_{n=1}^N (1 - \pi') h_{\phi}(\mathbf{x}_n, \mathbf{y}_n) \\ + \pi' \mathbb{E}_{q^{(t)}(\mathbf{y})} [h_{\phi}(\mathbf{x}_n, \mathbf{y})]\end{aligned}$$

$$\lambda^{(t+1)} = \arg \max_{\lambda \geq 0} \frac{1}{N} \sum_{n=1}^N q_{\lambda}(\mathbf{y}_n)$$

---

**Algorithm 1** Mutual Distillation

---

**Input:** Training data  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ ,  
Initial knowledge constraints  $\mathcal{F} = \{f_{\phi,l}\}_{l=1}^L$ ,  
Initial neural network  $p_\theta$ ,  
Parameters:  $\pi, \pi'$  – imitation parameters  
 $C$  – regularization parameters

- 1: Initialize neural network parameters  $\theta$
- 2: Initialize knowledge parameters  $\phi$  and weights  $\lambda$
- 3: **while** not converged **do**
- 4:   Sample a minibatch  $(\mathbf{X}, \mathbf{Y}) \subset \mathcal{D}$
- 5:   Build the teacher model  $q$  with Eq.(2) and Eq.(6)
- 6:   Update  $p_\theta$  with distillation objective Eq.(3)
- 7:   Update  $f_l$  ( $l = 1, \dots, L$ ) with distillation objective Eq.(5)
- 8: **end while**

**Output:** Learned network  $p$ , knowledge modules  $\mathcal{F}$ , and the joint teacher network  $q$

---