

1. Problem statement for biased datasets

Using definitions of $q(\mathbf{x}, \mathbf{y})$ and $p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$, (2) can be analytically derived as

$$\begin{aligned} D_{KL}(Q_{\mathbf{x}, \mathbf{y}} \| P_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})) &= \\ &= \int \int q(\mathbf{y}|\mathbf{x})q(\mathbf{x}) \log \frac{q(\mathbf{y}|\mathbf{x})q(\mathbf{x})}{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})q(\mathbf{x})} d\mathbf{y}d\mathbf{x} = \\ &= \int q(\mathbf{x}) \int q(\mathbf{y}|\mathbf{x}) \log \frac{q(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})} d\mathbf{y}d\mathbf{x} = \\ &= \mathbb{E}_{Q_{\mathbf{x}}} [D_{KL}(Q_{\mathbf{y}|\mathbf{x}} \| P_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\theta}))]. \end{aligned}$$

Assuming that $Q_{\mathbf{y}|\mathbf{x}}$ can be replaced by empirical $\hat{Q}_{\mathbf{y}|\mathbf{x}}$ and $\mathbf{y} = \mathbf{1}_d \in \mathbb{R}^D$ is one-hot vector with only d th class not equal to zero, (4) can be derived as

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{N^b} \sum_{i \in \mathbb{N}^b} [D_{KL}(Q_{\mathbf{y}_i|\mathbf{x}_i} \| P_{\mathbf{y}_i|\mathbf{x}_i}(\boldsymbol{\theta}))] = \\ &= \frac{1}{N^b} \sum_{i \in \mathbb{N}^b} \sum_{d=1}^D \mathbf{1}_d(i) \log \frac{\mathbf{1}_d(i)}{p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta})} = \\ &= -\frac{1}{N^b} \sum_{i \in \mathbb{N}^b} \log p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}). \end{aligned}$$

2. Relationship between $D_{KL}(P_z^v \| P_z)$ and Fisher information

Using the sufficiency property [1], we approximate our optimal acquisition function (5) using the distributions of learned representations \mathbf{z} as

$$\mathcal{R}_{opt}(b, P) = \arg \min_{\mathcal{R}(b, P)} D_{KL}(\hat{P}_z^v \| \hat{P}_z),$$

Then, a *connection* between the main task (2) and $D_{KL}(P_z^v \| P_z)$ minimization in (7) via Fisher information can be derived with respect to small perturbations in $\boldsymbol{\theta}$. Assuming that the task model minimizes distribution shift in (2) every backward pass as

$$p^v(\mathbf{z}|\boldsymbol{\theta}) = p(\mathbf{z}|\boldsymbol{\theta}) + \Delta p,$$

where $\Delta p = \Delta \boldsymbol{\theta} \frac{\partial p(\mathbf{z}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ and $\Delta \rightarrow 0$.

By substituting (8), the expanded form of $D_{KL}(P_z^v \| P_z)$ can be written as

$$\begin{aligned} D_{KL}(P_z^v \| P_z) &= \int (p(\mathbf{z}|\boldsymbol{\theta}) + \Delta p) \log \frac{p(\mathbf{z}|\boldsymbol{\theta}) + \Delta p}{p(\mathbf{z}|\boldsymbol{\theta})} dz = \\ &= \int (p(\mathbf{z}|\boldsymbol{\theta}) + \Delta p) \log \left(1 + \frac{\Delta p}{p(\mathbf{z}|\boldsymbol{\theta})} \right) dz. \end{aligned}$$

Using the Taylor series of natural logarithm, this can be

approximated by

$$\begin{aligned} D_{KL}(P_z^v \| P_z) &\approx \int (p(\mathbf{z}|\boldsymbol{\theta}) + \Delta p) \times \\ &\left(\frac{\Delta p}{p(\mathbf{z}|\boldsymbol{\theta})} - \frac{(\Delta p)^2}{2(p(\mathbf{z}|\boldsymbol{\theta}))^2} \right) dz = \int \Delta p dz + \\ &\frac{1}{2} \int \left(\frac{\Delta p}{p(\mathbf{z}|\boldsymbol{\theta})} \right)^2 p(\mathbf{z}|\boldsymbol{\theta}) dz - \int \frac{(\Delta p)^3}{2p(\mathbf{z}|\boldsymbol{\theta})^2} dz, \end{aligned}$$

where the first term using the definition of Δp is equal to zero and the third $\mathcal{O}(\Delta \boldsymbol{\theta}^3) \rightarrow 0$.

By substituting Δp and rewriting vector $\boldsymbol{\theta}$ as a discrete sum, the term

$$\frac{\Delta p}{p(\mathbf{z}|\boldsymbol{\theta})} \approx \sum_i \frac{\partial \log p(\mathbf{z}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i} \Delta \boldsymbol{\theta}_i.$$

Using this approximation, the final form of (7) can be obtained as

$$\begin{aligned} \mathcal{R}_{opt}(b, P) &= \arg \min_{\mathcal{R}(b, P)} D_{KL}(P_z^v \| P_z) \\ &\approx \arg \min_{\mathcal{R}(b, P)} \sum_{m, n} \mathcal{I}_{m, n} \Delta \boldsymbol{\theta}_m \Delta \boldsymbol{\theta}_n \approx \arg \min_{\mathcal{R}(b, P)} \Delta \boldsymbol{\theta}^T \mathcal{I} \Delta \boldsymbol{\theta}, \end{aligned}$$

where $\mathcal{I} = \mathbb{E}_{P_z} [\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})^T]$ is a Fisher information matrix and $\mathbf{g}(\boldsymbol{\theta}) = \frac{\partial \log p(\mathbf{z}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ is a Fisher score with respect to $\boldsymbol{\theta}$.

3. Practical Fisher kernel for DNNs

Using the chain rule for a DNN layer ($\tilde{z}_i^j = \boldsymbol{\theta}^T \mathbf{z}_i^j = \boldsymbol{\theta}^T \sigma(\tilde{z}_i^{j-1})$) with $\sigma(\cdot)$ nonlinearity, Jacobian of interest can be simplified as follows

$$\frac{\partial L(\mathbf{y}_i, \hat{\mathbf{y}}_i)}{\partial \boldsymbol{\theta}} = \frac{\partial L(\mathbf{y}_i, \hat{\mathbf{y}}_i)}{\partial \tilde{z}_i} \frac{\partial \tilde{z}_i}{\partial \boldsymbol{\theta}} = \frac{\partial L(\mathbf{y}_i, \hat{\mathbf{y}}_i)}{\partial \tilde{z}_i} \mathbf{z}_i^T = \mathbf{g}_i \mathbf{z}_i^T,$$

where $\boldsymbol{\theta} \in \mathbb{R}^{L \times L}$, $\mathbf{z}_i \in \mathbb{R}^{L \times 1}$, and $\mathbf{g}_i \in \mathbb{R}^{L \times 1}$.

Then, approximation of FK in (11) for $\mathbf{g}_i(\boldsymbol{\theta}) = \text{vec}(\partial L(\mathbf{y}_i, \hat{\mathbf{y}}_i)/\partial \boldsymbol{\theta}) \in \mathbb{R}^{L^2 \times 1}$ can be derived as

$$\begin{aligned} R_{z, g}(\mathbf{z}_m, \mathbf{z}_n) &= \mathbf{g}_m(\boldsymbol{\theta})^T \mathcal{I}^{-1} \mathbf{g}_n(\boldsymbol{\theta}) \stackrel{\text{PFK}}{\approx} \mathbf{g}_m(\boldsymbol{\theta})^T \mathbf{g}_n(\boldsymbol{\theta}) = \\ &= \text{vec} \left(\frac{\partial L(\mathbf{y}_m, \hat{\mathbf{y}}_m)}{\partial \tilde{z}_m} \mathbf{z}_m^T \right)^T \text{vec} \left(\frac{\partial L(\mathbf{y}_n, \hat{\mathbf{y}}_n)}{\partial \tilde{z}_n} \mathbf{z}_n^T \right) = \\ &= \text{vec}(\mathbf{g}_m \mathbf{z}_m^T)^T \text{vec}(\mathbf{g}_n \mathbf{z}_n^T) = [\mathbf{g}_m^1 \mathbf{z}_m, \mathbf{g}_m^2 \mathbf{z}_m, \dots, \mathbf{g}_m^L \mathbf{z}_m]^T \times \\ &= [\mathbf{g}_n^1 \mathbf{z}_n, \mathbf{g}_n^2 \mathbf{z}_n, \dots, \mathbf{g}_n^L \mathbf{z}_n] = \mathbf{z}_m^T \mathbf{z}_n \sum_l \mathbf{g}_m^l \mathbf{g}_n^l = \mathbf{z}_m^T \mathbf{z}_n \mathbf{g}_m^T \mathbf{g}_n. \end{aligned}$$

References

- [1] Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep representations. *Journal of Machine Learning Research*, pages 1947–1980, 2018.