Extracting Speaker-Specific Information with a Regularized Siamese Deep Network

Ke Chen and Ahmad Salman School of Computer Science, The University of Manchester Manchester M13 9PL, United Kingdom {chen, salmana}@cs.manchester.ac.uk

Appendix

In this appendix, we derive the gradient of $L_D(X_1, X_2; \Theta)$ with respect to $\boldsymbol{u}_K(\boldsymbol{x}_{it})$ to obtain Eq. (5) in the main text.

To simplify the presentation, we first elucidate our notation system that is completely consistent to that used in the main text. We collectively denote the output of neurons in CS at the code layer or layer K of subnet i (i=1,2) as $CS(X_i) = \left\{ \left((CS(\boldsymbol{x}_{it}))_j \right)_{j=1}^{|CS|} \right\}_{t=1}^{T_B}$ for a speech segment of T_B frames, $X_i = \{\boldsymbol{x}_{it}\}_{t=1}^{T_B}$. Accordingly, we have $\boldsymbol{\mu}^{(i)} = (\mu_j^{(i)})_{j=1}^{|CS|}$ and $\boldsymbol{\Sigma}^{(i)} = \left[\sigma_{ln}^{(i)} \right]$ $(l, n = 1, \cdots, |CS|)$ where $\sigma_{ln}^{(i)} = \frac{1}{T_B - 1} \sum_{t=1}^{T_B} \left[(CS(\boldsymbol{x}_{it}))_l - \mu_l^{(i)} \right] \left[(CS(\boldsymbol{x}_{it}))_n - \mu_n^{(i)} \right]^T$. In the following derivation, we also drop all explicit parameters in $L_D(X_1, X_2; \Theta)$ and rewrite it into $L_D = L_m + L_S$ where $L_m = \mathcal{I}D_m + (1 - \mathcal{I})e^{-\frac{D_m}{\lambda_m}}$ and $L_S = \mathcal{I}D_S + (1 - \mathcal{I})e^{-\frac{D_S}{\lambda_S}}$.

Using our notation described above, we immediately achieve

$$\frac{\partial L_D}{\partial \boldsymbol{u}_K(\boldsymbol{x}_{it})} = \frac{\partial L_m}{\partial \boldsymbol{u}_K(\boldsymbol{x}_{it})} + \frac{\partial L_S}{\partial \boldsymbol{u}_K(\boldsymbol{x}_{it})} \\
= \left(\left(\left[\mathcal{I} - \lambda_m^{-1} (1 - \mathcal{I}) e^{-\frac{D_m}{\lambda_m}} \right] \frac{\partial D_m}{\partial \boldsymbol{u}_{Kj}(\boldsymbol{x}_{it})} \right)_{j=1}^{|\mathcal{C}S|}, \left(0 \right)_{j=|\mathcal{C}S|+1}^{|\boldsymbol{h}_K|} \right) + \left(\left(\left[\mathcal{I} - \lambda_S^{-1} (1 - \mathcal{I}) e^{-\frac{D_S}{\lambda_S}} \right] \frac{\partial D_S}{\partial \boldsymbol{u}_{Kj}(\boldsymbol{x}_{it})} \right)_{j=1}^{|\mathcal{C}S|}, \left(0 \right)_{j=|\mathcal{C}S|+1}^{|\boldsymbol{h}_K|} \right). \quad (A.1)$$

To facilitate the presentation, we define $\psi_j(\boldsymbol{x}_{it}) = \frac{\partial D_m}{\partial u_{K_j}(\boldsymbol{x}_{it})}$ and $\xi_j(\boldsymbol{x}_{it}) = \frac{\partial D_S}{\partial u_{K_j}(\boldsymbol{x}_{it})}$. Now we simply need to calculate $\psi_j(\boldsymbol{x}_{it})$ and $\xi_j(\boldsymbol{x}_{it})$ for $j = 1, \dots, |\mathcal{CS}|$ to obtain Eq. (5) in the main text.

As
$$D_m = || \boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)} ||_2^2 = \sum_{l=1}^{|\mathcal{CS}|} (\mu_l^{(1)} - \mu_l^{(2)})^2$$
, we have

$$\psi_{j}(\boldsymbol{x}_{it}) = \frac{\partial D_{m}}{\partial \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}} \frac{\partial \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}}{\partial u_{Kj}(\boldsymbol{x}_{it})}$$

$$= \frac{\partial \sum_{l=1}^{|\mathcal{CS}|} (\mu_{l}^{(1)} - \mu_{l}^{(2)})^{2}}{\partial \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}} \frac{\partial \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}}{\partial u_{Kj}(\boldsymbol{x}_{it})}$$

$$= p_{j}^{(i)} \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j} \left[1 - \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}\right], \qquad (A.2)$$

where

$$p_j^{(i)} = \frac{\partial \sum_{l=1}^{|\mathcal{CS}|} (\mu_l^{(1)} - \mu_l^{(2)})^2}{\partial (\mathcal{CS}(\boldsymbol{x}_{it}))_j} = \frac{2}{T_B} \operatorname{sign}(1.5 - i)(\mu_j^{(1)} - \mu_j^{(2)}),$$

and

$$\frac{\partial \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}}{\partial u_{Kj}(\boldsymbol{x}_{it})} = \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j} \left[1 - \left(\mathcal{CS}(\boldsymbol{x}_{it})\right)_{j}\right]$$

given the fact that the transfer function used in the code layer or layer K is the sigmoid function. Collectively, we have $p^{(i)} = \frac{2}{T_B} \text{sign}(1.5 - i)(\mu^{(1)} - \mu^{(2)})$.

Similarly, $D_S = ||\Sigma^{(1)} - \Sigma^{(2)}||_F^2 = \sum_{l=1}^{|\mathcal{CS}|} \sum_{n=1}^{|\mathcal{CS}|} (\sigma_{ln}^{(1)} - \sigma_{ln}^{(2)})^2$. Hence, we have

$$\xi_{j}(\boldsymbol{x}_{it}) = \frac{\partial D_{S}}{\partial (\mathcal{CS}(\boldsymbol{x}_{it}))_{j}} \frac{\partial (\mathcal{CS}(\boldsymbol{x}_{it}))_{j}}{\partial u_{Kj}(\boldsymbol{x}_{it})}$$

$$= \frac{\partial \sum_{l=1}^{|\mathcal{CS}|} \sum_{n=1}^{|\mathcal{CS}|} (\sigma_{ln}^{(1)} - \sigma_{ln}^{(2)})^{2}}{\partial (\mathcal{CS}(\boldsymbol{x}_{it}))_{j}} \frac{\partial (\mathcal{CS}(\boldsymbol{x}_{it}))_{j}}{\partial u_{Kj}(\boldsymbol{x}_{it})}$$

$$= q_{j}(\boldsymbol{x}_{it}) (\mathcal{CS}(\boldsymbol{x}_{it}))_{j} [1 - (\mathcal{CS}(\boldsymbol{x}_{it}))_{j}], \qquad (A.3)$$

where

$$q_j(\boldsymbol{x}_{it}) = \frac{4}{T_B - 1} \operatorname{sign}(1.5 - i) \sum_{n=1}^{|\mathcal{CS}|} (\sigma_{jn}^{(1)} - \sigma_{jn}^{(2)}) \left[(\mathcal{CS}(\boldsymbol{x}_{it}))_n - \mu_n^{(i)} \right]$$

and, collectively, we have $\boldsymbol{q}(\boldsymbol{x}_{it}) = \frac{4}{T_B - 1} \operatorname{sign}(1.5 - i)(\Sigma^{(1)} - \Sigma^{(2)})[\mathcal{CS}(\boldsymbol{x}_{it}) - \boldsymbol{\mu}^{(i)}].$ Inserting Eqs. (A.2) and (A.3) into Eq. (A.1), we obtain Eq. (5) in the main text as

$$\frac{\partial L_D}{\partial \boldsymbol{u}_K(\boldsymbol{x}_{it})} = \left(\left([\mathcal{I} - \lambda_m^{-1}(1-\mathcal{I})e^{-\frac{D_m}{\lambda_m}}]\psi_j(\boldsymbol{x}_{it}) \right)_{j=1}^{|\mathcal{CS}|}, \left(0\right)_{j=|\mathcal{CS}|+1}^{|\boldsymbol{h}_K|} \right) + \left(\left([\mathcal{I} - \lambda_S^{-1}(1-\mathcal{I})e^{-\frac{D_S}{\lambda_S}}]\xi_j(\boldsymbol{x}_{it}) \right)_{j=1}^{|\mathcal{CS}|}, \left(0\right)_{j=|\mathcal{CS}|+1}^{|\boldsymbol{h}_K|} \right).$$