

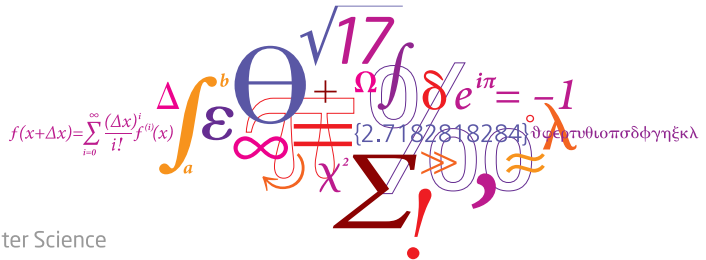
Towards Transferable Speech Emotion Representation: On Loss Functions For Cross-Lingual Latent Representations

Sneha Das¹, Nicole Nadine Lørfeldt², Anne Katrine Pagsberg^{2,3}, Line H. Clemmensen¹

¹Department of Applied Mathematics and Computer Science, Technical University of Denmark

²Child and Adolescent Mental Health Center, Copenhagen University Hospital, Capital Region

³Faculty of Health, Department of Clinical Medicine, Copenhagen University



Motivation

- Speech emotion recognition (SER): inferring emotional state from speech signals.
- Emotion recognition employed in healthcare, education sector, criminal justice system.
- SER: signal processing, machine learning, deep learning.
- Existing challenges: Generalizing over languages, corpora, recording condition (under low-resource conditions).

Objectives and Contributions

Objectives for transferability:

- ① Latent embedding with discrimination between emotion classes.
- ② Latent distribution that are consistent over corpora.

Contributions:

- ① Low-complexity DAE and VAE.
- ② VAE with KL-loss annealing: balancing KL-loss and reconstruction loss.
- ③ VAE with semi-supervision incorporating clustering in latent space.

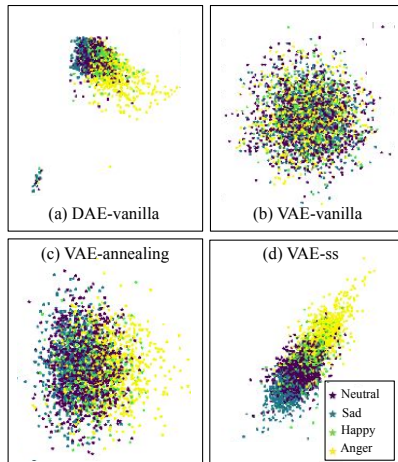
Formulation

- DAE:

$$\arg \min_{f_\theta, g_\phi} \mathcal{L}_{\text{rec}} = \mathbb{E} \|\mathbf{x} - g_\phi(f_\theta(\mathbf{x}_n))\|_2^2, \quad (1)$$

- VAE:

$$\begin{aligned} \arg \min_{\theta, \phi} \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{KL}} = & -\mathbb{E}_{\mathbf{z} \sim q_\theta(\mathbf{z}|\mathbf{x})} \log p_\phi(\mathbf{x}|\mathbf{z}) \\ & + D_{\text{KL}}(q_\theta(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})), \end{aligned} \quad (2)$$



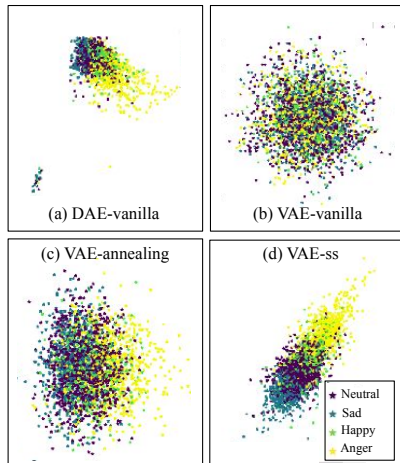
Formulation

- VAE with KL-annealing:

$$\arg \min_{\theta, \phi} \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{KL}} = -\mathbb{E}_{\mathbf{z} \sim q_{\theta}(\mathbf{z}|\mathbf{x})} \log p_{\phi}(\mathbf{x}|\mathbf{z}) + \beta_e D_{\text{KL}}(q_{\theta}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})), \quad (3)$$

where the standard formulation of β_e :

$$\beta_e = \begin{cases} f(\tau) = \frac{0.25}{R} \tau, & \tau \leq R \\ 0.25, & \tau > R \end{cases} \quad \text{where} \quad \tau = \frac{\text{mod}(e-1, \frac{T}{M})}{\frac{T}{M}}, \quad (4)$$

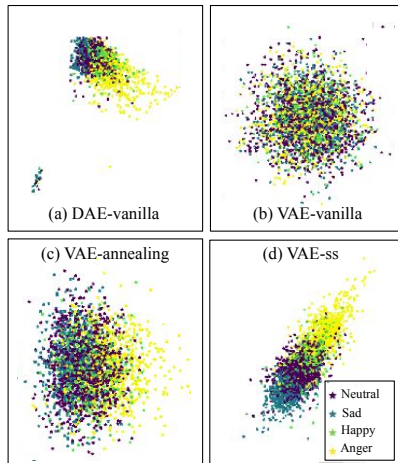


Formulation

- VAE with semi-supervision:

$$\arg \min_{\theta, \phi} \mathcal{L}_{\text{rec}} + \beta_e \mathcal{L}_{\text{KL}} + \gamma \mathcal{L}_{\text{clus}},$$

$$\mathcal{L}_{\text{clus}} = \frac{D_{\text{intra}}}{D_{\text{inter}}} = \frac{\sum_{k=1}^K \sum_{\forall i \in k} D(\mathbf{z}_i, \bar{\mathbf{z}}^k)}{\sum_{k=1}^{K-1} \sum_{j=k+1}^K D(\bar{\mathbf{z}}^k, \bar{\mathbf{z}}^j)}, \quad (5)$$



Architecture

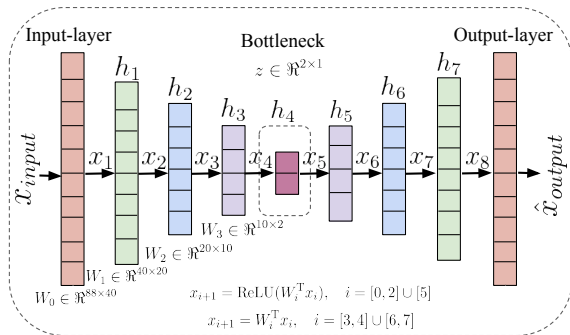


Figure: Illustration of the architecture employed for all the models explored in this work.

- Training: 50 epochs, batch size 64, Adam optimizer (learning rate: 1e-3).
- Latent embedding used as input features to a linear SVC.

Evaluation

- Datasets: IEMOCAP, SAVEE, Emo-DB, CaFE, URDU, AESD
- Input features: eGeMAPS using OpenSmile
- Preprocessing: remove outliers using z-score normalization ($-10 > z > 10$)
- 5-fold cross validation

Results: Classification performance

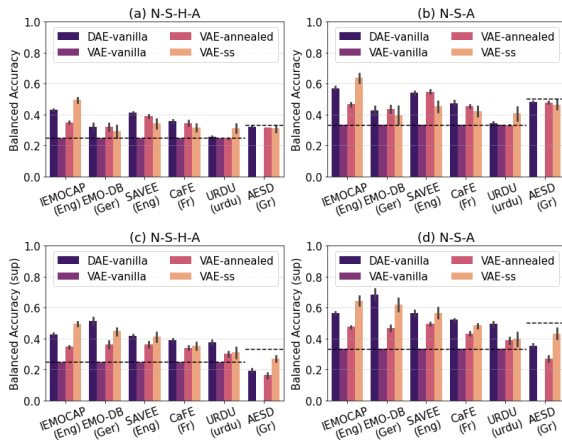


Figure: (1) Balanced accuracy on unseen transfer data sets using (a) 4 emotion classes, (b) 3 emotion classes; balanced accuracy with access to 20% of the unlabeled transfer data sets with (c) 4 emotions and (d) 3 emotion classes.

Results: Consistency of latent space

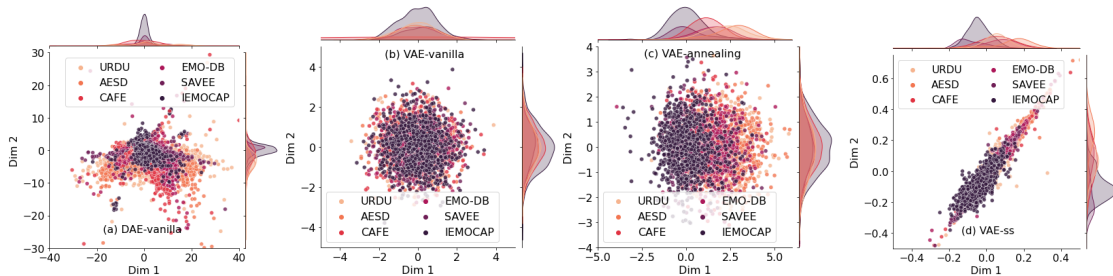
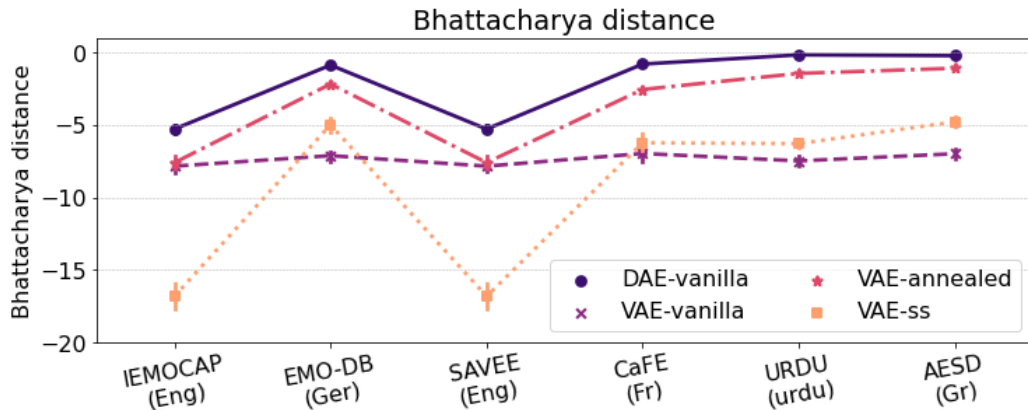


Figure: Scatter plots depicting the overlap between the latent embedding obtained from the methods investigated for all the transfer data sets.

Results: Consistency of latent space



Conclusions

- ① DAE: highest classification accuracy, worst distribution consistency.
- ② VAE-vanilla: best consistency, classification accuracy random.
- ③ VAE-ss: Classification accuracy similar to DAE and distribution consistency similar to VAE-vanilla.