

Semi-supervised Bayesian Deep Multi-modal Emotion Recognition

Changde Du¹, Changying Du², Jinpeng Li¹, Wei-long Zheng³, Bao-liang Lu³, Huiguang He¹

¹Research Center for Brain-Inspired Intelligence,

Institute of Automation, Chinese Academy of Sciences (CAS), Beijing, China

²Laboratory of Parallel Software and Computational Science, Institute of Software, CAS, Beijing, China

³Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China
{duchangde2016, huiguang.he}@ia.ac.cn

Abstract

In emotion recognition, it is difficult to recognize human's emotional states using just a single modality. Besides, the annotation of physiological emotional data is particularly expensive. These two aspects make the building of effective emotion recognition model challenging. In this paper, we first build a multi-view deep generative model to simulate the generative process of multi-modality emotional data. By imposing a mixture of Gaussians assumption on the posterior approximation of the latent variables, our model can learn the shared deep representation from multiple modalities. To solve the labeled-data-scarcity problem, we further extend our multi-view model to semi-supervised learning scenario by casting the semi-supervised classification problem as a specialized missing data imputation task. Our semi-supervised multi-view deep generative framework can leverage both labeled and unlabeled data from multiple modalities, where the weight factor for each modality can be learned automatically. Compared with previous emotion recognition methods, our method is more robust and flexible. The experiments conducted on two real multi-modal emotion datasets have demonstrated the superiority of our framework over a number of competitors.

1 Introduction

With the development of human-computer interaction, emotion recognition has become increasingly important. Since human's emotion contains many nonverbal cues, various modalities ranging from facial expressions, voice, Electroencephalogram (EEG), eye movements to other physiological signals can be used as the indicators of emotional states [Calvo and D'Mello, 2010]. In real-world applications, it is difficult to recognize the emotional states using just a single modality, because signals from different modalities represent different aspects of emotion and provide complementary information. Recent studies show that integrating multiple modalities can significantly boost the emotion recognition accuracy [Verma and Tiwary, 2014; Pang *et al.*, 2015;

Lu *et al.*, 2015; Liu *et al.*, 2016; Soleymani *et al.*, 2016; Zhang *et al.*, 2016].

The most successful approach to fuse the information from multiple modalities is based on deep multi-view representation learning [Ngiam *et al.*, 2011; Andrew *et al.*, 2013; Srivastava and Salakhutdinov, 2014; Wang *et al.*, 2015; Chandar *et al.*, 2016]. For example, [Pang *et al.*, 2015] proposed to learn a joint density model for emotion analysis with a multi-modal Deep Boltzmann Machine (DBM) [Srivastava and Salakhutdinov, 2014]. This multi-modal DBM is exploited to model the joint distribution over visual, auditory, and textual features. [Liu *et al.*, 2016] proposed a multi-modal emotion recognition method by using multi-modal Deep Autoencoders (DAE) [Ngiam *et al.*, 2011], in which the joint representations of EEG and eye movement signals were extracted. Nevertheless, there are still limitations with these deep multi-modal emotion recognition methods, e.g., their performances depend on the amount of labeled data.

By using the modern sensor equipments, we can easily collect massive physiological signals, which are closely related to people's emotional states. Despite the convenience of data acquisition, the data labeling procedure requires lots of manual efforts. Therefore, in most cases only a small set of labeled samples is available, while the majority of whole dataset is left unlabeled. Traditional emotion recognition approaches only utilized the limited amount of labeled data, which may result in severe overfitting. The most attractive way to deal with this issue is based on Semi-supervised Learning (SSL), which builds more robust model by exploiting both labeled and unlabeled data [Schels *et al.*, 2014; Jia *et al.*, 2014; Zhang *et al.*, 2016]. Though multi-modal approaches have been widely implemented for emotion recognition, very few of them explored SSL simultaneously. To the best of our knowledge, only [Zhang *et al.*, 2016] proposed an enhanced multi-modal co-training algorithm for semi-supervised emotion recognition, but its shallow structure is hard to capture the high-level correlation between different modalities.

Amongst existing SSL approaches, the most competitive one is based on deep generative models, which employs the Deep Neural Networks (DNNs) to learn discriminative features and casts the semi-supervised classification problem as a specialized missing data imputation task. [Kingma *et al.*, 2014] and [Maaløe *et al.*, 2016] have shown that deep generative models and approximate Bayesian inference exploiting

recent advances in scalable variational methods [Kingma and Welling, 2014; Rezende *et al.*, 2014] can provide state-of-the-art performance for semi-supervised classification. Though the Variational Autoencoder (VAE) framework [Kingma and Welling, 2014] has shown great advantages in SSL, its potential merits remain under-explored. For example, until recently, there was no successful multi-view extension for it. The main difficulty lies in its inherent assumption that the posterior approximation should be conditioned on the data point, which is natural to single-view data but becomes problematic for multi-view case.

In this paper, we propose a novel semi-supervised multi-view deep generative framework for multi-modal emotion recognition. Our framework combines the advantages of deep multi-view representation learning and Bayesian modeling, thus it has sufficient flexibility and robustness in learning joint features and classifier. Our main contributions can be summarized as follows.

- We propose a multi-view extension for VAE by imposing a mixture of Gaussians assumption on the posterior approximation of the latent variables. For multi-view learning, this is critical for fully exploiting the information from multiple views.
- We introduce a semi-supervised multi-modal emotion recognition framework based on multi-view VAE. Our framework can leverage both labeled and unlabeled samples from multiple modalities and the weight factor for each modality can be learned automatically, which is critical for building a robust emotion recognition system.
- We demonstrate the superiority of our framework and provide insightful observations on two real multi-modal emotion datasets.

2 Multi-view Variational Autoencoder for Semi-supervised Emotion Recognition

The VAE framework has recently been introduced as a robust model for latent feature learning [Kingma and Welling, 2014; Rezende *et al.*, 2014]. However, the single-view architecture in VAE can't effectively deal with multi-view data. In this section, we first build a multi-view VAE, which can learn the shared deep representation from multi-view data. And then, we extend it to the semi-supervised scenario. Assume we are faced with multi-view data that appears as pairs $(\mathcal{X}, y) = (\{\mathbf{x}^{(v)}\}_{v=1}^V, y)$, with observation $\mathbf{x}^{(v)}$ from the v -th view and the corresponding class label y .

2.1 Multi-view Variational Autoencoder

DNN-parameterized Likelihoods

We assume the latent variable \mathbf{z} can generate multi-view features $\{\mathbf{x}^{(v)}\}_{v=1}^V$. Specifically, we assume \mathbf{z} generates $\mathbf{x}^{(v)}$ for any $v \in \{1, \dots, V\}$, with the following generative model (cf. Fig. 1a):

$$p_{\theta^{(v)}}(\mathbf{x}^{(v)}|\mathbf{z}) = f(\mathbf{x}^{(v)}; \mathbf{z}, \theta^{(v)}), \quad (1)$$

where $f(\mathbf{x}^{(v)}; \mathbf{z}, \theta^{(v)})$ is a suitable likelihood function (e.g. a Gaussian for continuous observation or Bernoulli for binary

observation), which is formed by a non-linear transformation of the latent variable \mathbf{z} . This non-linear transformation is essential to allow for higher moments of the data to be captured by the density model, and we choose these non-linear functions to be DNNs, referred to as the generative networks, with parameters $\{\theta^{(v)}\}_{v=1}^V$. Note that, the likelihoods for different data views are assumed to be independent of each other, with different nonlinear transformations.

The Bayesian Canonical Correlation Analysis (CCA) model [Klami *et al.*, 2013] can be seen as a special case of our model, where linear shallow transformations were used to generate each data view and only two different views were considered. [Wang *et al.*, 2016] used a similar deep non-linear generative process as ours to construct deep Bayesian CCA model, but during inference they construct the variational posterior approximation from just one view and ignore the rest one. Such a choice is convenient for inference and computation, but only seeks suboptimal solutions as it doesn't fully exploit the data. As shown in the following, we assume the variational approximation to the posterior of latent variables to be a mixture of Gaussians, utilizing information from multiple views.

Gaussian Prior and Mixture of Gaussians Posterior

Typically, both the prior $p(\mathbf{z})$ and the approximate posterior $q_{\phi}(\mathbf{z}|\mathcal{X})$ are assumed to be Gaussian distributions [Kingma and Welling, 2014; Rezende *et al.*, 2014] in order to maintain mathematical and computational tractability. Although this assumption has led to favorable results on several tasks, it is clearly a restrictive and often unrealistic assumption. Specifically, the choice of a Gaussian distribution for $p(\mathbf{z})$ and $q_{\phi}(\mathbf{z}|\mathcal{X})$ imposes a strong uni-modal structure assumption on the latent space. However, for data distributions that are strongly multi-modal, the uni-modal Gaussian assumption inhibits the model's ability to extract and represent important structure in the data. To improve the flexibility of the model, one way is to impose a mixture of Gaussians assumption on $p(\mathbf{z})$. However, it has the risk of creating separate "islands" of discontinuous manifolds that may break the meaningfulness of the representation in the latent space.

To learn more powerful and expressive models – in particular, models with multi-modal latent variable structures for multi-modal emotion recognition applications – we seek a mixture of Gaussians for $q_{\phi}(\mathbf{z}|\mathcal{X})$, while preserving $p(\mathbf{z})$ as a standard Gaussian. Thus (cf. Fig. 1b),

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}),$$

$$q_{\phi}(\mathbf{z}|\mathcal{X}) = \sum_{v=1}^V \lambda^{(v)} \mathcal{N}\left(\mathbf{z}|\boldsymbol{\mu}_{\phi^{(v)}}(\mathbf{x}^{(v)}), \boldsymbol{\Sigma}_{\phi^{(v)}}(\mathbf{x}^{(v)})\right), \quad (2)$$

where the mean $\boldsymbol{\mu}_{\phi^{(v)}}$ and the covariance $\boldsymbol{\Sigma}_{\phi^{(v)}}$ are nonlinear functions of the observation $\mathbf{x}^{(v)}$, with variational parameter $\phi^{(v)}$. As in our generative model, we choose these nonlinear functions to be DNNs, referred to as the inference networks. $\lambda^{(v)}$ is the non-negative normalized weight factor for the v -th view, i.e., $\lambda^{(v)} > 0$ and $\sum_{v=1}^V \lambda^{(v)} = 1$. By conditioning the posterior approximation on the data point, we avoid variational parameters per data point, instead only requiring to fit global variational parameters. Note that, our mixed

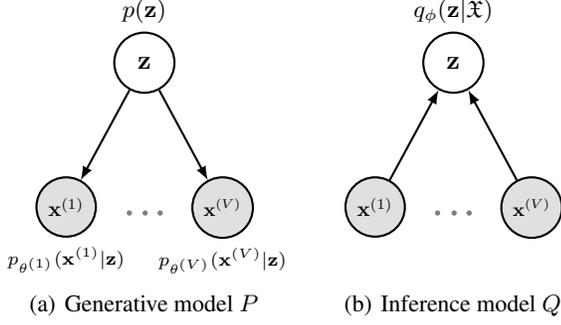


Figure 1: Graphical model of the multi-view VAE, where $\mathfrak{X} = \{\mathbf{x}^{(v)}\}_{v=1}^V$.

Gaussian assumption on the variational approximation distinguishes our method from all existing ones using the auto-encoding variational framework [Kingma and Welling, 2014; Wang *et al.*, 2016; Burda *et al.*, 2016; Kingma *et al.*, 2016; Serban *et al.*, 2016; Maaløe *et al.*, 2016]. For multi-view learning, this is critical for fully exploiting the information from multiple views.

2.2 Semi-supervised Emotion Recognition

In semi-supervised classification, only a subset of the samples have corresponding class labels, and we focus on using the multi-view VAE to build a model (semiMVAE) that learns classifier from both labeled and unlabeled multi-view data. Since the emotional data is continuous, we choose the Gaussian likelihoods. Then the generative model P is defined as $p(y)p(\mathbf{z})\prod_{v=1}^V p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})$ (cf. Fig. 2a):

$$\begin{aligned} p(y) &= \text{Cat}(y|\boldsymbol{\pi}), \\ p(\mathbf{z}) &= \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}), \\ p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}) &= \mathcal{N}(\boldsymbol{\mu}_{\theta^{(v)}}(y, \mathbf{z}), \text{diag}(\boldsymbol{\sigma}_{\theta^{(v)}}^2(y, \mathbf{z}))), \end{aligned} \quad (3)$$

where $\text{Cat}(\cdot)$ denotes the categorical distribution, y is treated as a latent variable for the unlabeled data points, and the mean $\boldsymbol{\mu}_{\theta^{(v)}}$ and variance $\boldsymbol{\sigma}_{\theta^{(v)}}^2$ are nonlinear functions of y and \mathbf{z} , with parameter $\theta^{(v)}$. The inference model Q is defined as $q_{\varphi}(y|\mathfrak{X})q_{\phi}(\mathbf{z}|\mathfrak{X}, y)$ (cf. Fig. 2b):

$$\begin{aligned} q_{\varphi}(y|\mathfrak{X}) &= \text{Cat}(y|\boldsymbol{\pi}_{\varphi}(\mathfrak{X})), \\ q_{\phi}(\mathbf{z}|\mathfrak{X}, y) &= \sum_{v=1}^V \lambda^{(v)} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi^{(v)}}(\mathbf{x}^{(v)}, y), \boldsymbol{\Sigma}_{\phi^{(v)}}(\mathbf{x}^{(v)}, y)), \end{aligned} \quad (4)$$

where $q_{\phi}(\mathbf{z}|\mathfrak{X}, y)$ is assumed to be a mixture of Gaussians to combine the information from multiple data views. Intuitively, $q_{\phi}(\mathbf{z}|\mathfrak{X}, y)$, $p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})$ and $q_{\varphi}(y|\mathfrak{X})$ correspond to the encoder, the decoder and the classifier, respectively.

For brevity, we omit the explicit dependencies on $\mathbf{x}^{(v)}$, y and \mathbf{z} for the moment variables mentioned above hereafter. In principle, $\boldsymbol{\mu}_{\theta^{(v)}}$, $\boldsymbol{\sigma}_{\theta^{(v)}}^2$, $\boldsymbol{\pi}_{\varphi}$, $\boldsymbol{\mu}_{\phi^{(v)}}$ and $\boldsymbol{\Sigma}_{\phi^{(v)}}$ can be implemented by various DNN models, e.g., Multiple Layer Perceptrons (MLP) and Convolutional Neural Networks (CNN).

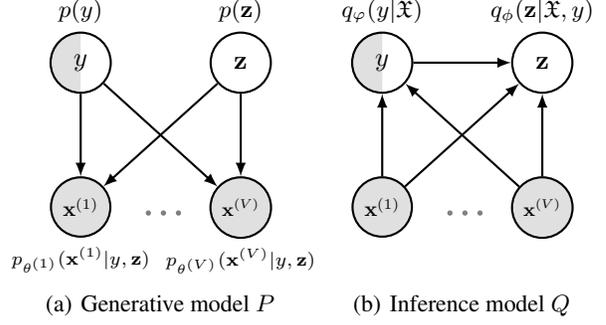


Figure 2: Graphical model of the semiMVAE for semi-supervised multi-view learning, where $\mathfrak{X} = \{\mathbf{x}^{(v)}\}_{v=1}^V$.

2.3 Variational Lower Bound

The variational lower bound on the marginal likelihood for a single labeled data point is

$$\begin{aligned} \log p_{\theta}(\mathfrak{X}, y) &= \log \int_{\mathbf{z}} p_{\theta}(\mathfrak{X}, y, \mathbf{z}) dz \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)} \left[\log \frac{p_{\theta}(\mathfrak{X}, y, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)} \left[\sum_{v=1}^V \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}) + \log p(y) \right. \\ &\quad \left. + \log p(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathfrak{X}, y) \right] \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)} \left[\sum_{v=1}^V \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}) + \log p(y) \right. \\ &\quad \left. + \log p(\mathbf{z}) \right] - \sum_{v=1}^V \lambda^{(v)} \cdot \log \left(\sum_{l=1}^V \lambda^{(l)} \cdot \omega_{v,l} \right) \\ &\equiv -\mathcal{L}(\mathfrak{X}, y), \end{aligned} \quad (5)$$

where $\omega_{v,l} = \mathcal{N}(\boldsymbol{\mu}_{\phi^{(v)}}|\boldsymbol{\mu}_{\phi^{(l)}}, \boldsymbol{\Sigma}_{\phi^{(v)}} + \boldsymbol{\Sigma}_{\phi^{(l)}})$. Note that, the Shannon entropy $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)}[-\log q_{\phi}(\mathbf{z}|\mathfrak{X}, y)]$ is hard to compute analytically, and we have used the Jensen's inequality to derive a lower bound of it:

$$\begin{aligned} &\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathfrak{X}, y)}[-\log q_{\phi}(\mathbf{z}|\mathfrak{X}, y)] \\ &= -\sum_{v=1}^V \lambda^{(v)} \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi^{(v)}}, \boldsymbol{\Sigma}_{\phi^{(v)}}) \log q_{\phi}(\mathbf{z}|\mathfrak{X}, y) dz \\ &\geq -\sum_{v=1}^V \lambda^{(v)} \log \left(\sum_{l=1}^V \lambda^{(l)} \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi^{(v)}}, \boldsymbol{\Sigma}_{\phi^{(v)}}) \right. \\ &\quad \left. \cdot \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi^{(l)}}, \boldsymbol{\Sigma}_{\phi^{(l)}}) dz \right) \\ &= -\sum_{v=1}^V \lambda^{(v)} \log \left(\sum_{l=1}^V \lambda^{(l)} \mathcal{N}(\boldsymbol{\mu}_{\phi^{(v)}}|\boldsymbol{\mu}_{\phi^{(l)}}, \boldsymbol{\Sigma}_{\phi^{(v)}} + \boldsymbol{\Sigma}_{\phi^{(l)}}) \right). \end{aligned}$$

For unlabeled data, we further introduce the variational dis-

tribution $q_\varphi(y|\mathfrak{X})$ for y :

$$\begin{aligned}\log p_\theta(\mathfrak{X}) &= \log \int_{\mathbf{z}} \int_y p_\theta(\mathfrak{X}, y, \mathbf{z}) dy d\mathbf{z} \\ &\geq \mathbb{E}_{q_{\varphi, \phi}(y, \mathbf{z}|\mathfrak{X})} \left[\log \frac{p_\theta(\mathfrak{X}, y, \mathbf{z})}{q_{\varphi, \phi}(y, \mathbf{z}|\mathfrak{X})} \right] \\ &= \mathbb{E}_{q_\varphi(y|\mathfrak{X})} [-\mathcal{L}(\mathfrak{X}, y) - \log q_\varphi(y|\mathfrak{X})] \\ &\equiv -\mathcal{U}(\mathfrak{X}),\end{aligned}\quad (6)$$

with $q_{\varphi, \phi}(y, \mathbf{z}|\mathfrak{X}) = q_\varphi(y|\mathfrak{X})q_\phi(\mathbf{z}|\mathfrak{X}, y)$. The objective function for the entire dataset is now:

$$\mathcal{J} = \sum_{(\mathfrak{X}, y) \in S_l} \mathcal{L}(\mathfrak{X}, y) + \sum_{\mathfrak{X} \in S_u} \mathcal{U}(\mathfrak{X}), \quad (7)$$

where S_l and S_u are labeled and unlabeled dataset, respectively. The classification accuracy can be improved by introducing an explicit classification loss for labeled data. The extended objective function is:

$$\mathcal{F} = \mathcal{J} + \alpha \cdot \sum_{(\mathfrak{X}, y) \in S_l} [-\log q_\varphi(y|\mathfrak{X})], \quad (8)$$

where the hyper-parameter α is a weight between generative and discriminative learning. We set $\alpha = \beta \cdot (N_l + N_u)$, where β is a scaling constant, and N_l and N_u are the numbers of labeled and unlabeled data points in one minibatch, respectively. Note that, the classifier $q_\varphi(y|\mathfrak{X})$ is also used at test phase for the prediction of unseen data.

2.4 Optimization

Eq. (8) provides a unified objective function for optimizing the parameters of encoder, decoder and classifier networks. This optimization can be done jointly, without resort to the variational EM algorithm, using the stochastic backpropagation technique [Kingma and Welling, 2014; Rezende *et al.*, 2014].

Reparameterization Trick

The reparameterization trick is a vital component of the algorithm, because it allows us to easily take the derivative of $\mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})]$ with respect to the variational parameters ϕ . However, the use of a mixture of Gaussians for the variational distribution $q_\phi(\mathbf{z}|\mathfrak{X}, y)$ makes the application of reparameterization trick challenging. It can be shown that, for any $v \in \{1, \dots, V\}$, $\mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})]$ can be rewritten, using the location-scale transformation for the Gaussian distribution, as:

$$\begin{aligned}\mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})] \\ = \sum_{l=1}^V \lambda^{(l)} \mathbb{E}_{\mathcal{N}(\epsilon^{(l)}|\mathbf{0}, \mathbf{I})} [\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \boldsymbol{\mu}_{\phi^{(l)}} + \mathbf{R}_{\phi^{(l)}} \epsilon^{(l)})],\end{aligned}\quad (9)$$

where $\mathbf{R}_{\phi^{(l)}} \mathbf{R}_{\phi^{(l)}}^\top = \Sigma_{\phi^{(l)}}$ and $l \in \{1, \dots, V\}$.

Gradients of the Objective

While the expectations on the right hand side of Eq. (9) still cannot be solved analytically, their gradients w.r.t. $\theta^{(v)}$,

$\phi^{(l)}$ and $\lambda^{(l)}$ can be efficiently estimated using the following Monte-Carlo estimators,

$$\begin{aligned}\frac{\partial}{\partial \theta^{(v)}} \mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})] \\ = \sum_{l=1}^V \lambda^{(l)} \mathbb{E}_{\mathcal{N}(\epsilon^{(l)}|\mathbf{0}, \mathbf{I})} \left[\frac{\partial}{\partial \theta^{(v)}} \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l)}) \right] \\ \approx \frac{\lambda^{(l)}}{T} \sum_{t=1}^T \sum_{l=1}^V \frac{\partial}{\partial \theta^{(v)}} \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l,t)}),\end{aligned}\quad (10)$$

$$\begin{aligned}\frac{\partial}{\partial \phi^{(l)}} \mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})] \\ = \lambda^{(l)} \frac{\partial}{\partial \phi^{(l)}} \mathbb{E}_{\mathcal{N}(\epsilon^{(l)}|\mathbf{0}, \mathbf{I})} \left[\frac{\partial}{\partial \mathbf{z}^{(l)}} \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l)}) \right. \\ \left. \cdot \left(\frac{\partial \boldsymbol{\mu}_{\phi^{(l)}}}{\partial \phi^{(l)}} + \frac{\partial \mathbf{R}_{\phi^{(l)}}}{\partial \phi^{(l)}} \epsilon^{(l)} \right) \right] \\ \approx \frac{\lambda^{(l)}}{T} \sum_{t=1}^T \frac{\partial}{\partial \mathbf{z}^{(l,t)}} \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l,t)}) \\ \left. \cdot \left(\frac{\partial \boldsymbol{\mu}_{\phi^{(l)}}}{\partial \phi^{(l)}} + \frac{\partial \mathbf{R}_{\phi^{(l)}}}{\partial \phi^{(l)}} \epsilon^{(l,t)} \right),\end{aligned}\quad (11)$$

$$\begin{aligned}\frac{\partial}{\partial \lambda^{(l)}} \mathbb{E}_{q_\phi(\mathbf{z}|\mathfrak{X}, y)}[\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z})] \\ = \mathbb{E}_{\mathcal{N}(\epsilon^{(l)}|\mathbf{0}, \mathbf{I})} [\log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l)})] \\ \approx \frac{1}{T} \sum_{t=1}^T \log p_{\theta^{(v)}}(\mathbf{x}^{(v)}|y, \mathbf{z}^{(l,t)}),\end{aligned}\quad (12)$$

where $\mathbf{z}^{(l)}$ is evaluated at $\mathbf{z}^{(l)} = \boldsymbol{\mu}_{\phi^{(l)}} + \mathbf{R}_{\phi^{(l)}} \epsilon^{(l)}$ and $\mathbf{z}^{(l,t)} = \boldsymbol{\mu}_{\phi^{(l)}} + \mathbf{R}_{\phi^{(l)}} \epsilon^{(l,t)}$ with $\epsilon^{(l,t)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In practice, it suffices to use a small T (e.g. $T = 1$) and then estimate the gradient using minibatches of data points. We use the same random numbers $\epsilon^{(l,t)}$ for all estimators to have lower variance. The gradient w.r.t. φ is omitted here, since it can be derived straightforwardly by using traditional reparameterization trick [Kingma *et al.*, 2014].

The gradients of the loss for semiMVAE (Eq. (8)) can then be computed by a direct application of the chain rule and estimators presented above. During optimization we can use the estimated gradients in conjunction with standard stochastic gradient based optimization methods such as SGD, RM-Sprop or Adam [Kingma and Ba, 2014]. Overall, the model can be trained with reparameterization trick for backpropagation through the mixed Gaussian latent variables.

3 Experiments

In this section, we present extensive experimental results to demonstrate the effectiveness of the proposed semi-supervised multi-view framework for emotion recognition.

3.1 Experimental Testbed and Setup

Data description Two multi-modal emotion datasets, the SEED dataset¹ [Lu *et al.*, 2015] and the DEAP

¹<http://bcmi.sjtu.edu.cn/%7Eseed/index.html>

dataset² [Koelstra *et al.*, 2012], were used in our experiments.

The SEED dataset contains EEG and eye movement signals from 15 subjects during watching 15 movie clips, where each movie clip lasts about 4 minutes long. The EEG signals were recorded from 62 channels and the eye movement signals contained information about blink, saccade fixation and so on. In our experiment, we used the data from 9 subjects across 3 sessions, totally 27 data files. For each data file, data from watching the 1-9 movie clips were used as training set, while data from watching the 10-12 movie clips were used as validation set and the rest (13-15) were used as testing set.

The DEAP dataset contains EEG and peripheral physiological signals of 32 participants. Signals were recorded when they were watching 40 one-minute duration music videos. The EEG signals were recorded from 32 channels, whereas the peripheral physiological signals were recorded from 8 channels. The participants, using values from 1 to 9, rated each music video in terms of the levels of valence, arousal and so on. In our experiment, the valence-arousal space was divided into four quadrants according to the ratings. The threshold we used was 5, leading to four classes of data. Considering the fuzzy boundary of emotions and the variations of participants' ratings possibly associated with individual difference in rating scale, we discarded the samples whose ratings of arousal and valence are between 3 and 6. The dataset was randomly divided into 10-folds, where 8 folds for training, one fold for validation and the last fold for testing. The size of testing set is relative small, because some graph-based semi-supervised baselines are hard to deal with large dataset.

Feature selection For SEED dataset, Lu *et al.* have extracted the Differential Entropy (DE) features and 33 eye movement features from EEG and eye movement signals [Lu *et al.*, 2015]. We also used these features in our experiments. For DEAP dataset, we extracted the DE features from EEG and peripheral physiological signals. The DE features can be calculated in four frequency bands: theta (4-8Hz), alpha (8-14Hz), beta (14-31Hz), and gamma (31-45Hz), and we used all band's features. The details of the data used in our experiments were summarized in Table 1.

Table 1: The details of the datasets used in our experiments.

Datasets	#Instances	#Features	#Training	#Validation	#Testing	#Classes
SEED	22734	310(EEG), 33(Eye)	13473	4725	4536	3
DEAP	21042	128(EEG), 32(Phy.)	16834	2104	2104	4

Compared methods We compared our semiMVAE with a broad range of solutions, including supervised learning, transductive and inductive semi-supervised learning. We briefly summarize the various baselines in the following.

- **MAE**: the multi-view extension of deep autoencoders, which can be used to extract the joint representations from multiple modalities [Ngiam *et al.*, 2011].
- **DCCA**: the full deep neural network extension of Canonical Correlation Analysis (CCA). DCCA can learn deep nonlinear mappings of two views, which are maximally correlated [Andrew *et al.*, 2013].
- **DCCAE**: a deep multi-view representation learning model which combines the advantages of the DCCA

and deep autoencoders. In particular, DCCAE consists of two autoencoders and optimizes the combination of canonical correlation between the learned bottleneck representations and the reconstruction errors of the autoencoders [Wang *et al.*, 2015].

- **AMMSS**: a graph-based multi-view semi-supervised classification algorithm, which can integrate heterogeneous features from both labeled and unlabeled data [Cai *et al.*, 2013].
- **AMGL**: a latest auto-weighted multiple graph learning framework, which can be applied to multi-view semi-supervised classification task [Nie *et al.*, 2016].
- **semiVAE**: a single-view semi-supervised deep generative model proposed in [Kingma *et al.*, 2014]. We evaluate semiVAE's performance for each modality and the concatenation of all modalities, respectively.

For MAE, DCCA and DCCAE, we used the Support Vector Machines³ (SVM) and transductive SVM⁴ (TSVM) for supervised learning and transductive semi-supervised learning, respectively.

Parameter setting For semiMVAE, we considered multiple layer perceptrons as the type of inference and generative networks. On both datasets, we set the structures of the inference and generative networks for each view as '100-50-30' and '30-50-100', respectively. We used the Adam optimizer [Kingma and Ba, 2014] with a learning rate $\eta = 3 \times 10^{-4}$ in training. The scaling constant β was selected from {0.1, 0.5, 1} throughout the experiments. The weight factor for each view was initialized with $\lambda^{(v)} = 1/V$, where V is the number of views. For MAE, DCCA and DCCAE, we considered the same setups (network structure, learning rate, etc.) as our semiMVAE. For AMMSS, we tuned the parameters as suggested in [Cai *et al.*, 2013]. For AMGL and semiVAE, we used their default settings.

3.2 Performance Evaluation

To simulate semi-supervised learning scenario, on both datasets, we randomly labeled different proportions of samples in the training set, and remained the rest samples in the training set unlabeled. For transductive semi-supervised learning, we trained models on the dataset consisting of the testing data and labeled data belonging to training set. For inductive semi-supervised learning, we trained models on the entire training set consisting of the labeled and unlabeled data. For supervised learning, we trained models on the labeled data belonging to training set, and test their performance on the testing set. Table 2 presents the classification accuracies of all methods on SEED and DEAP datasets. The proportions of labeled samples in the training set vary from 1% to 3%. Several observations can be drawn as follows.

First, the average accuracy of semiMVAE significantly surpasses the baselines in all cases. Second, by examining semiMVAE against supervised learning approaches trained on very limited labeled data, we can find that semiMVAE always outperforms them. This encouraging result shows that

²<http://www.eecs.qmul.ac.uk/mmv/datasets/deap/download.html>

³<http://www.csie.ntu.edu.tw/%7Eejlin/liblinear/>.

⁴<http://svmlight.joachims.org/>.

Table 2: Comparison with several supervised and semi-supervised methods on SEED and DEAP with few labels. Results (mean \pm std) were averaged over 20 independent runs.

SEED data	Algorithms	1% labeled	2% labeled	3% labeled	
Supervised learning	MAE+SVM	.814 \pm .031	.896 \pm .024	.925 \pm .024	
	DCCA+SVM	.809 \pm .035	.891 \pm .035	.923 \pm .028	
	DCCAE+SVM	.819 \pm .036	.893 \pm .034	.923 \pm .027	
Transductive semi-supervised learning	AMMSS	.731 \pm .055	.839 \pm .036	.912 \pm .018	
	AMGL	.711 \pm .047	.817 \pm .023	.886 \pm .028	
	MAE+TSVM	.818 \pm .035	.910 \pm .025	.931 \pm .026	
	DCCA+TSVM	.811 \pm .031	.903 \pm .024	.928 \pm .021	
	DCCAE+TSVM	.823 \pm .040	.907 \pm .027	.929 \pm .023	
Inductive semi-supervised learning	semiMVAE	.861\pm.037	.931\pm.020	.960\pm.021	
	semiVAE (Eye)	.753 \pm .024	.849 \pm .055	.899 \pm .049	
	semiVAE (EEG)	.768 \pm .041	.861 \pm .040	.919 \pm .026	
	semiVAE (Concat.)	.803 \pm .035	.876 \pm .043	.926 \pm .044	
DEAP data	Algorithms	1% labeled	2% labeled	3% labeled	
	Supervised learning	MAE+SVM	.353 \pm .027	.387 \pm .014	.411 \pm .016
		DCCA+SVM	.359 \pm .016	.400 \pm .014	.416 \pm .018
		DCCAE+SVM	.361 \pm .023	.403 \pm .017	.419 \pm .013
Transductive semi-supervised learning	AMMSS	.303 \pm .029	.353 \pm .024	.386 \pm .014	
	AMGL	.291 \pm .027	.341 \pm .021	.367 \pm .019	
	MAE+TSVM	.376 \pm .025	.403 \pm .031	.417 \pm .026	
	DCCA+TSVM	.379 \pm .021	.408 \pm .024	.421 \pm .017	
	DCCAE+TSVM	.384 \pm .022	.412 \pm .027	.425 \pm .021	
Inductive semi-supervised learning	semiMVAE	.424\pm.020	.441\pm.013	.456\pm.013	
	semiVAE (Phy.)	.366 \pm .024	.389 \pm .048	.402 \pm .034	
	semiVAE (EEG)	.374 \pm .019	.397 \pm .013	.407 \pm .016	
	semiVAE (Concat.)	.383 \pm .019	.404 \pm .016	.416 \pm .012	
semiMVAE	.421\pm.019	.439\pm.025	.451\pm.022		

semiMVAE can effectively leverage the useful information from unlabeled data. Third, multi-view semi-supervised algorithms AMMSS and AMGL perform worst in all cases. We attribute this to the fact that graph-based shallow models AMMSS and AMGL can't extract the deep features from the original data. Fourth, the performances of three TSVM based semi-supervised methods are moderate. Although MAE+TSVM, DCCA+TSVM and DCCAE+TSVM can also integrate multi-modality information from unlabeled samples, their two-stage learning can't obtain the global optimal model parameters. Finally, compared with the single-view semi-supervised method semiVAE, our multi-view method is more effective in integrating multiple modalities.

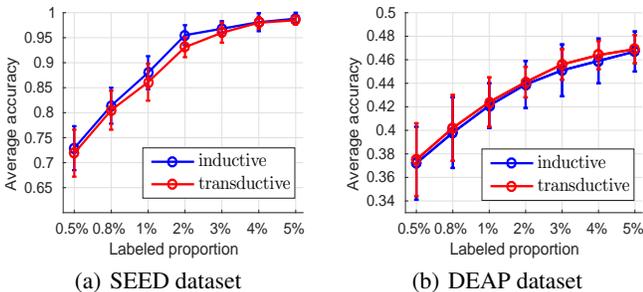


Figure 3: semiMVAE's performance with different proportions of labeled samples in the training set.

The proportion of labeled and unlabeled samples in the training set will affect the performance of semi-supervised

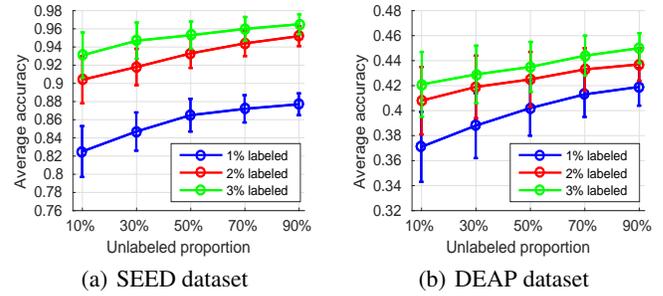


Figure 4: Inductive semiMVAE's performance with different proportions of unlabeled samples in the training set.

models. Figs. 3 and 4 show the changes of semiMVAE's average accuracy on both datasets with different proportions of labeled and unlabeled samples in the training set. We can observe that both labeled and unlabeled samples can effectively boost the classification accuracy of semiMVAE.

Instead of treating each modality equally, our semiMVAE can weight each modality and perform classification simultaneously. Fig. 5a shows the learned weight factors by inductive semiMVAE on SEED and DEAP datasets (1% labeled). From it, we can observe that EEG modality has the highest weight on both datasets, which is consistent with single modality's performance of semiVAE shown in Table 2 and the results in previous work [Lu *et al.*, 2015].

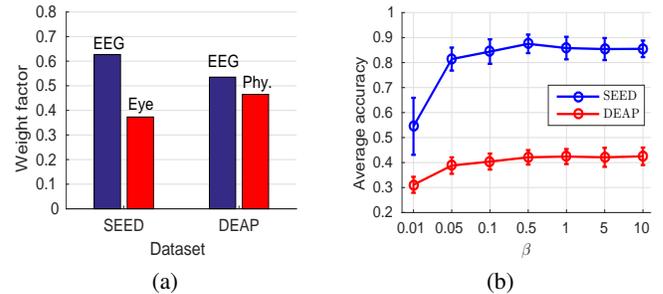


Figure 5: (a) Learned weight factors by inductive semiMVAE. (b) The impact of scaling constant β .

The scaling constant β controls the weight of discriminative learning in semiMVAE. Fig. 5b shows the performance of inductive semiMVAE with different β values (1% labeled). From it, we can find that the scaling constant β can be chosen from $\{0.1, 0.5, 1\}$, where semiMVAE achieves good results.

4 Conclusion

This paper proposes a semi-supervised multi-view deep generative framework for emotion recognition, which can leverage both labeled and unlabeled data from multiple modalities. The key to the framework are two parts: 1) multi-view VAE can fully integrate the information from multiple modalities and 2) semi-supervised learning can overcome the labeled-data-scarcity problem. Experimental results on two real multi-modal emotion datasets demonstrate the effectiveness of our approach.

References

- [Andrew *et al.*, 2013] Galen Andrew, Raman Arora, Jeff A Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *ICML*, pages 1247–1255, 2013.
- [Burda *et al.*, 2016] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In *ICLR*, 2016.
- [Cai *et al.*, 2013] Xiao Cai, Feiping Nie, Weidong Cai, and Heng Huang. Heterogeneous image features integration via multi-modal semi-supervised learning model. In *ICCV*, pages 1737–1744, 2013.
- [Calvo and D’Mello, 2010] Rafael A Calvo and Sidney D’Mello. Affect detection: An interdisciplinary review of models, methods, and their applications. *IEEE Transactions on Affective Computing*, 1(1):18–37, 2010.
- [Chandar *et al.*, 2016] Sarath Chandar, Mitesh M Khapra, Hugo Larochelle, and Balaraman Ravindran. Correlation neural networks. *Neural computation*, 28(2):257–285, 2016.
- [Jia *et al.*, 2014] Xiaowei Jia, Kang Li, Xiaoyi Li, and Aidong Zhang. A novel semi-supervised deep learning framework for affective state recognition on EEG signals. In *International Conference on Bioinformatics and Bioengineering (BIBE)*, pages 30–37. IEEE, 2014.
- [Kingma and Ba, 2014] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [Kingma and Welling, 2014] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014.
- [Kingma *et al.*, 2014] Diederik P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *NIPS*, pages 3581–3589, 2014.
- [Kingma *et al.*, 2016] Diederik P Kingma, Tim Salimans, and Max Welling. Improving variational inference with inverse autoregressive flow. In *NIPS*, 2016.
- [Klami *et al.*, 2013] Arto Klami, Seppo Virtanen, and Samuel Kaski. Bayesian canonical correlation analysis. *Journal of Machine Learning Research*, 14(1):965–1003, 2013.
- [Koelstra *et al.*, 2012] Sander Koelstra, Christian Muhl, Mohammad Soleymani, Jong-Seok Lee, Ashkan Yazdani, Touradj Ebrahimi, Thierry Pun, Anton Nijholt, and Ioannis Patras. Deap: A database for emotion analysis; using physiological signals. *IEEE Transactions on Affective Computing*, 3(1):18–31, 2012.
- [Liu *et al.*, 2016] Wei Liu, Wei-Long Zheng, and Bao-Liang Lu. Multimodal emotion recognition using multimodal deep learning. *arXiv preprint arXiv:1602.08225*, 2016.
- [Lu *et al.*, 2015] Yifei Lu, Wei-Long Zheng, Binbin Li, and Bao-Liang Lu. Combining eye movements and EEG to enhance emotion recognition. In *IJCAI*, pages 1170–1176, 2015.
- [Maaløe *et al.*, 2016] Lars Maaløe, Casper Kaae Sønderby, Søren Kaae Sønderby, and Ole Winther. Auxiliary deep generative models. In *ICML*, pages 1445–1453, 2016.
- [Ngiam *et al.*, 2011] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. Multimodal deep learning. In *ICML*, pages 689–696, 2011.
- [Nie *et al.*, 2016] Feiping Nie, Jing Li, Xuelong Li, et al. Parameter-free auto-weighted multiple graph learning: A framework for multiview clustering and semi-supervised classification. In *IJCAI*, pages 1881–1887, 2016.
- [Pang *et al.*, 2015] Lei Pang, Shiai Zhu, and Chong-Wah Ngo. Deep multimodal learning for affective analysis and retrieval. *IEEE Transactions on Multimedia*, 17(11):2008–2020, 2015.
- [Rezende *et al.*, 2014] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *NIPS*, pages 1278–1286, 2014.
- [Schels *et al.*, 2014] Martin Schels, Markus Kächele, Michael Glodek, David Hrabal, Steffen Walter, and Friedhelm Schwenker. Using unlabeled data to improve classification of emotional states in human computer interaction. *Journal on Multimodal User Interfaces*, 8(1):5–16, 2014.
- [Serban *et al.*, 2016] Iulian V Serban, II Ororbia, G Alexander, Joelle Pineau, and Aaron Courville. Multimodal variational encoder-decoders. *arXiv preprint arXiv:1612.00377*, 2016.
- [Soleymani *et al.*, 2016] Mohammad Soleymani, Sadjad Asghari-Esfeden, Yun Fu, and Maja Pantic. Analysis of EEG signals and facial expressions for continuous emotion detection. *IEEE Transactions on Affective Computing*, 7(1):17–28, 2016.
- [Srivastava and Salakhutdinov, 2014] Nitish Srivastava and Ruslan Salakhutdinov. Multimodal learning with deep boltzmann machines. *Journal of Machine Learning Research*, 15:2949–2980, 2014.
- [Verma and Tiwary, 2014] Gyanendra K Verma and Uma Shanker Tiwary. Multimodal fusion framework: A multiresolution approach for emotion classification and recognition from physiological signals. *NeuroImage*, 102:162–172, 2014.
- [Wang *et al.*, 2015] Weiran Wang, Raman Arora, Karen Livescu, and Jeff A Bilmes. On deep multi-view representation learning. In *ICML*, pages 1083–1092, 2015.
- [Wang *et al.*, 2016] Weiran Wang, Xinchen Yan, Honglak Lee, and Karen Livescu. Deep variational canonical correlation analysis. *arXiv: 1610.03454*, 2016.
- [Zhang *et al.*, 2016] Zixing Zhang, Fabien Ringeval, Bin Dong, Eduardo Coutinho, Erik Marchi, and Björn Schüller. Enhanced semi-supervised learning for multimodal emotion recognition. In *ICASSP*, pages 5185–5189. IEEE, 2016.