

M²VAE – Derivation of a Multi-Modal Variational Autoencoder Objective from the Marginal Joint Log-Likelihood

Timo Korthals^{†*}

Abstract

This work gives an in-depth derivation of the trainable evidence lower bound (ELBO) obtained from the marginal joint log-Likelihood with the goal of training a multi-modal variational Autoencoder (M²VAE).

I. INTRODUCTION

Variational auto encoder (VAE) combine neural networks with variational inference to allow unsupervised learning of complicated distributions according to the graphical model shown in Fig. 1 (left). A D_a -dimensional observation a is modeled in terms of a D_z -dimensional latent vector z using a probabilistic decoder $p_{\theta_a}(z)$ with parameters θ . To generate the corresponding embedding z from observation a , a probabilistic encoder network with $q_{\phi_a}(z)$ is being provided which parametrizes the posterior distribution from which z is sampled. The encoder and decoder, given by neural networks, are trained jointly to bring a close to an a' under the constraint that an approximate distribution needs to be close to a prior $p(z)$ and hence inference is basically learned during training.

The specific objective of VAEs is the maximization of the marginal distribution $p(a) = \int p_{\theta}(a|z)p(z) da$. Because this distribution is intractable, the model is instead trained via *stochastic gradient variational Bayes* (SGVB) by maximizing the *evidence lower bound* (ELBO) \mathcal{L} of the marginal log-likelihood as described in Sec. II This approach proposed by [1] is used in settings where only a single modality a is present in order to find a latent encoding z (c.f. Fig. 1 (left)).

This work gives an in-depth derivation of the trainable *evidence lower bound* (ELBO) obtained from the marginal joint log-Likelihood, that satisfies all plate models as depicted in Fig. 1, we are with the goal of training a multi-modal variational Autoencoder (M²VAE).

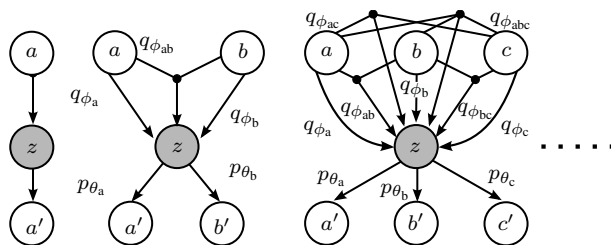


Fig. 1: Evolution of full uni-, bi-, and tri-modal VAEs comprising all modality permutations

[†]Bielefeld University, Cognitronics & Sensor Systems, Inspiration 1, 33619 Bielefeld, Germany

*tkorthals@cit-ec.uni-bielefeld.de

II. VARIATIONAL AUTOENCODER

First, the derivation of the vanilla *Variational Autoencoder* by [1] is recaped.

A. The Variational Bound

$$L = \log(p(a)) \quad (1)$$

$$= \sum_z q(z|a) \log(p(a)) \quad \text{Equation 89 w/o conditional} \quad (2)$$

$$= \sum_z q(z|a) \log\left(\frac{p(z,a)}{p(z|a)}\right) \quad \text{Equation 84} \quad (3)$$

$$= \sum_z q(z|a) \log\left(\frac{p(z,a) q(z|a)}{p(z|a) q(z|a)}\right) \quad \text{multiplied by 1} \quad (4)$$

$$= \sum_z q(z|a) \log\left(\frac{p(z,a) q(z|a)}{q(z|a) p(z|a)}\right) \quad \text{reordered} \quad (5)$$

$$= \sum_z q(z|a) \log\left(\frac{p(z,a)}{q(z|a)}\right) + \sum_z q(z|a) \log\left(\frac{q(z|a)}{p(z|a)}\right) \quad \text{Equation 86} \quad (6)$$

$$= \mathcal{L} + D_{\text{KL}}(q(z|a)||p(z|a)) \quad \text{Equation 92 \& 88} \quad (7)$$

$$\geq \mathcal{L} \quad D_{\text{KL}} \geq 0 \quad (8)$$

D_{KL} is the Kulbeck-Leibler divergence, with $D_{\text{KL}} \geq 0$, that depends on how good $q(z|a)$ can approximate $p(z|a)$. \mathcal{L} is the lower variational bound of the marginal log-likelihood, also called the *evidence lower bound (ELBO)*. If and only if the two distributions q and p are identical, D_{KL} becomes 0 ($q = p \Leftrightarrow D_{\text{KL}} = 0$). $\mathcal{L} = L$ means on the other hand therefore implicitly, that q perfectly approximates p . It is because \mathcal{L} and D_{KL} are in equilibrium so that minimizing D_{KL} is identical to the maximization of \mathcal{L} ($\min D_{\text{KL}} \Leftrightarrow \max \mathcal{L}$). Minimizing D_{KL} is not feasible, because we don't know the true posterior $p(z|a)$. Therefore, \mathcal{L} is further investigated.

B. Approximate Inference (i.e. rewriting \mathcal{L})

$$\mathcal{L} = \sum_z q(z|a) \log\left(\frac{p(z,a)}{q(z|a)}\right) \quad (9)$$

$$= \sum_z q(z|a) \log\left(\frac{p(a|z)p(z)}{q(z|a)}\right) \quad \text{Equation 84} \quad (10)$$

$$= \sum_z q(z|a) \log\left(\frac{p(z)}{q(z|a)}\right) + \sum_z q(z|a) \log(p(a|z)) \quad \text{Equation 86} \quad (11)$$

$$= -D_{\text{KL}}(q(z|a)||p(z)) + \mathbb{E}_{q(z|a)} \log(p(a|z)) \quad \text{Equation 92} \quad (12)$$

If the variable a is replaced by some real valued sample $a^{(i)}$ (e.g. image or LiDAR scan), two terms can be identified:

$$\mathcal{L} = \underbrace{-D_{\text{KL}}\left(q_\phi\left(z|a^{(i)}\right)||p(z)\right)}_{\text{Regularization}} + \underbrace{\mathbb{E}_{q_\phi\left(z|a^{(i)}\right)} \log\left(p_\theta\left(a^{(i)}|z\right)\right)}_{\text{Reconstruction}} \quad (13)$$

The first term is just a regularize that punishes the variational distribution q , that is the approximator of the posterior distribution, if it deviates from some prior $p(z)$. The reconstruction term on the other hand compares the difference between the data $a^{(i)}$ of $q_\phi(z|a^{(i)})$ wrt. the sampled data $a^{(i)}$ from the likelihood function $p_\theta(a^{(i)}|z)$.

That means if \mathcal{L} is going to be maximized, the posterior function has to be equal to some prior and the data that is used to sample the latent feature z from the posterior should be equal to the data from the likelihood function p_θ . The objective is now, to find a function q_ϕ and p_θ which own these properties. Luckily, if some parametrized functions q_ϕ and p_θ (with parameters ϕ and θ) are applied, every increase of Equation 13 means that the variational approximator $q_\phi(z|a)$ comes closer to the real posterior functions $p(z|a)$. Therefore, numerical optimization techniques like gradient descent can be applied to this issue. Commonly, two neuronal network, where each tries to find the optimal parameters ϕ and θ , are applied to approximate the functions q_ϕ (i.e. the encoder) and p_θ (i.e. the decoder). However, since the true value of L remains unknown, maximization of \mathcal{L} can only be done until convergence. Thus, the overall procedure has the property of finding local optima.

III. JOINT VARIATIONAL AUTOENCODER

Second, we expand the VAE from Sec. II to the marginal joint log-likelihood and derive the variational bound as follows:

$$L_J = \log(p(a,b)) \quad (14)$$

$$= \sum_z q(z|a,b) \log(p(a,b)) \quad \text{Equation 89 w/o conditional} \quad (15)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a,b)}{p(z|a,b)}\right) \quad \text{Equation 85} \quad (16)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a,b) q(z|a,b)}{p(z|a,b) q(z|a,b)}\right) \quad \text{multiplied by 1} \quad (17)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a,b) q(z|a,b)}{q(z|a,b) p(z|a,b)}\right) \quad \text{reordered} \quad (18)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a,b)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{q(z|a,b)}{p(z|a,b)}\right) \quad \text{Equation 86} \quad (19)$$

$$= \mathcal{L}_J + D_{\text{KL}}(q(z|a,b)||p(z|a,b)) \quad \text{Equation 92 \& 88} \quad (20)$$

$$\geq \mathcal{L}_J \quad (21)$$

Approximate Inference (i.e. rewriting \mathcal{L}_J):

$$\mathcal{L}_J = \sum_z q(z|a,b) \log\left(\frac{p(z,a,b)}{q(z|a,b)}\right) \quad (22)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(a,b|z)p(z)}{q(z|a,b)}\right) \quad \text{Equation 84} \quad (23)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log(p(a,b|z)) \quad \text{Equation 86} \quad (24)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log(p(a|z)) + \sum_z q(z|a,b) \log(p(b|z)) \quad \text{Equation 90} \quad (25)$$

$$= -D_{\text{KL}}(q(z|a,b)||p(z)) + \mathbb{E}_{q(z|a,b)} \log(p(a|z)) + \mathbb{E}_{q(z|a,b)} \log(p(b|z)) \quad \text{Equation 92} \quad (26)$$

Three different terms can be identified:

$$\mathcal{L}_J = \underbrace{-D_{\text{KL}}(q_{\phi_{ab}}(z|a,b)||p(z))}_{\text{Regularization}} + \underbrace{\mathbb{E}_{q_{\phi_{ab}}(z|a,b)} \log(p_{\theta_a}(a|z))}_{\text{Reconstruction wrt. } a} + \underbrace{\mathbb{E}_{q_{\phi_{ab}}(z|a,b)} \log(p_{\theta_b}(b|z))}_{\text{Reconstruction wrt. } b} \quad (27)$$

A regularization for the joint encoder $q_{\phi_{ab}}$ and two reconstruction terms, one for each decoder p_{θ_a} and p_{θ_b} .

IV. JOINT MULTI-MODAL VARIATIONAL AUTOENCODER VIA VARIATION OF INFORMATION

The issue with the joint VAE is the lacking possibility of encoding just one modality a or b . Thus, Suzuki et al. [2] exploit the *Variation of Information* (VI) and derive the evidence lower bound wrt. the VI.

First, the conditional probability is investigated

$$p(a|b) = \frac{p(z,a|b)}{p(z|a,b)} \quad \text{Equation 85} \quad (28)$$

$$= \frac{1}{p(z|a,b)} p(z,a|b) \quad (29)$$

$$= \frac{1}{p(z|a,b)} \frac{p(z,a,b)}{p(b)} \quad \text{Equation 84} \quad (30)$$

$$= \frac{1}{p(z|a,b)} \frac{p(a,b|z)p(z)}{p(b)} \quad \text{Equation 84} \quad (31)$$

$$= \frac{1}{p(z|a,b)} \frac{p(a|z)p(b|z)p(z)}{p(b)} \quad \text{Equation 90} \quad (32)$$

$$= \frac{1}{p(z|a,b)} \frac{p(a|z)p(z|b)\frac{p(b)}{p(z)}p(z)}{p(b)} \quad \text{Equation 84} \quad (33)$$

$$= \frac{p(a|z)p(z|b)}{p(z|a,b)} \quad (34)$$

Further, the marginal log-likelihood of a conditional distribution can be written as:

$$L_{M_a} = \log(p(a|b)) \quad (35)$$

$$= \sum_z q(z|a,b) \log(p(a|b)) \quad \text{Equation 89 w/o conditional} \quad (36)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a|b)}{p(z|a,b)}\right) \quad \text{Equation 85} \quad (37)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a|b)}{p(z|a,b)} \frac{q(z|a,b)}{q(z|a,b)}\right) \quad \text{multiplied by 1} \quad (38)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a|b)}{q(z|a,b)} \frac{q(z|a,b)}{p(z|a,b)}\right) \quad \text{reordered} \quad (39)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(z,a|b)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{q(z|a,b)}{p(z|a,b)}\right) \quad \text{Equation 86} \quad (40)$$

$$= \mathcal{L}_{M_a} + D_{\text{KL}}(q(z|a,b)||p(z|a,b)) \quad \text{Equation 92 \& 88} \quad (41)$$

$$\geq \mathcal{L}_{M_a} \quad (42)$$

Further, the log-likelihood of the VI can be written as

$$L_M = L_{M_a} + L_{M_b} \quad (43)$$

$$= \log(p(a|b)) + \log(p(b|a)) \quad (44)$$

$$= \mathcal{L}_{M_a} + \mathcal{L}_{M_b} + 2 D_{\text{KL}}(q(z|a,b)||p(z|a,b)) \quad \text{Equation 41} \quad (45)$$

$$\geq \mathcal{L}_{M_a} + \mathcal{L}_{M_b} \quad (46)$$

$$\mathcal{L}_{M_a} + \mathcal{L}_{M_b} = \sum_z q(z|a,b) \log\left(\frac{p(z,a|b)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{p(z,b|a)}{q(z|a,b)}\right) \quad (47)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(a|z)p(z|b)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{p(b|z)p(z|a)}{q(z|a,b)}\right) \quad \text{Equation 34} \quad (48)$$

$$= \sum_z q(z|a,b) \log\left(\frac{p(a|z)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{p(z|b)}{q(z|a,b)}\right) \quad (49)$$

$$+ \sum_z q(z|a,b) \log\left(\frac{p(b|z)}{q(z|a,b)}\right) + \sum_z q(z|a,b) \log\left(\frac{p(z|a)}{q(z|a,b)}\right) \quad \text{reordering} \quad (50)$$

$$= \mathbb{E}_{q(z|a,b)} \log(p(a|z)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|b)) \quad (51)$$

$$+ \mathbb{E}_{q(z|a,b)} \log(p(b|z)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a)) \quad \text{Equation 92} \quad (52)$$

$$= \mathbb{E}_{q(z|a,b)} \log(p(a|z)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|b)) \quad (53)$$

$$+ \mathbb{E}_{q(z|a,b)} \log(p(b|z)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a)) \quad (54)$$

$$+ \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a,b)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a,b)) \quad \text{added 0} \quad (55)$$

$$= \mathcal{L}_J - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|b)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a)) + \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a,b)) \quad \text{substitute eq. 26} \quad (56)$$

$$\geq \mathcal{L}_J - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|b)) - \text{D}_{\text{KL}}(q(z|a,b) \| p(z|a)) \quad (57)$$

$$=: \mathcal{L}_M \quad (58)$$

With respect to Equation 27, the following regularization terms can be identified:

$$\mathcal{L}_M = \mathcal{L}_J - \underbrace{\text{D}_{\text{KL}}(q_{\phi_{ab}}(z|a,b) \| q_{\phi_b}(z|b))}_{\text{Unimodal PDF fitting of encoder b}} - \underbrace{\text{D}_{\text{KL}}(q_{\phi_{ab}}(z|a,b) \| q_{\phi_a}(z|a))}_{\text{Unimodal PDF fitting of encoder a}} \quad (59)$$

A. Conclusion

The introduced KL regularization by Suzuki et al. [2] tries to find a mean representative of all parameters between clusters in the latent space. This is absolutely correct for the mean values but is insufficient for all other statistics of the distribution that are estimated by q_{ϕ^*} . The approach would be sufficient, iff the true latent joint probability would be known, because only then the KL divergence is able to adapt to it in a correct manner. Unfortunately, this is not the case and thus all the estimated statistic values of the uni-modal encoders q_{ϕ^*} , except the mean, are questionable.

V. PROPOSED JOINT MULTI-MODAL AUTOENCODER

A. The Variational Bound

$$L_{PJ} = \log(p(a,b)) \tag{60}$$

$$= \frac{2}{2} \log(p(a,b)) \quad \text{multiplied by 1} \tag{61}$$

$$= \frac{1}{2} \log(p(a,b)^2) \quad \text{Equation 87} \tag{62}$$

$$= \frac{1}{2} \log(p(a,b)p(a,b)) \tag{63}$$

$$= \frac{1}{2} \log(p(b)p(a|b)p(b|a)p(a)) \quad \text{Equation 84} \tag{64}$$

$$= \frac{1}{2} (\log(p(a)) + \log(p(b|a)) + \log(p(a|b)) + \log(p(b))) \quad \text{Equation 86} \tag{65}$$

$$\geq \frac{1}{2} (\mathcal{L}_a + \mathcal{L}_{M_a} + \mathcal{L}_{M_b} + \mathcal{L}_b) \quad \text{Equation 8 \& 46} \tag{66}$$

$$\geq \frac{1}{2} (\mathcal{L}_a + \mathcal{L}_M + \mathcal{L}_b) \quad \text{Equation 59} \tag{67}$$

$$:= \mathcal{L}_{PJ} \tag{68}$$

B. Approximating Inference (i.e. rewriting \mathcal{L}_{PJ})

$$2\mathcal{L}_{PJ} = \mathcal{L}_a + \mathcal{L}_M + \mathcal{L}_b \tag{69}$$

$$= -D_{\text{KL}}(q(z|a)||p(z)) + \mathbb{E}_{q(z|a)} \log(p(a|z)) \quad \text{Equation 12} \tag{70}$$

$$- D_{\text{KL}}(q(z|a,b)||p(z)) + \mathbb{E}_{q(z|a,b)} \log(p(a|z)) + \mathbb{E}_{q(z|a,b)} \log(p(b|z)) \quad \text{Equation 26} \tag{71}$$

$$- D_{\text{KL}}(q(z|a,b)||p(z|b)) - D_{\text{KL}}(q(z|a,b)||p(z|a)) \quad \text{Equation 113} \tag{72}$$

$$- D_{\text{KL}}(q(z|b)||p(z)) + \mathbb{E}_{q(z|b)} \log(p(b|z)) \quad \text{Equation 12} \tag{73}$$

Applying the corresponding function approximators, the formula can be written as:

$$2\mathcal{L}_{PJ} = \mathcal{L}_a + \mathcal{L}_M + \mathcal{L}_b \tag{74}$$

$$= -D_{\text{KL}}(q_{\phi_a}(z|a)||p(z)) + \mathbb{E}_{q_{\phi_a}(z|a)} \log(p_{\theta_a}(a|z)) \quad \text{Equation 13} \tag{75}$$

$$- D_{\text{KL}}(q_{\phi_{ab}}(z|a,b)||p(z)) + \mathbb{E}_{q_{\phi_{ab}}(z|a,b)} \log(p_{\theta_a}(a|z)) + \mathbb{E}_{q_{\phi_{ab}}(z|a,b)} \log(p_{\theta_b}(b|z)) \quad \text{Equation 27} \tag{76}$$

$$- D_{\text{KL}}(q_{\phi_{ab}}(z|a,b)||q_{\phi_b}(z|b)) - D_{\text{KL}}(q_{\phi_{ab}}(z|a,b)||q_{\phi_a}(z|a)) \quad \text{Equation 59} \tag{77}$$

$$- D_{\text{KL}}(q_{\phi_b}(z|b)||p(z)) + \mathbb{E}_{q_{\phi_b}(z|b)} \log(p_{\theta_b}(b|z)) \quad \text{Equation 13} \tag{78}$$

Investigating every line of the formula, the following properties can be identified: Equation 76 is the common multi-modal VAE loss derived from the joint probability, while Equation 77 adds the feature introduced by Suzuki et al. [2]. It introduces the KL regularization that brings the posterior distribution of an uni-modal encoder close to the distribution of the multi-modal case. The drawback of this approach is discussed in Section IV-A. The new lines, i.e. Equation 75 and 78, introduce the regularization of the uni-modal encoders wrt. the common prior and the reconstruction loss. The regularizer cares about the fact, that the uni-modal distribution does not deviate to much from the common prior while the reconstruction term shapes the remaining statistics including the mean. However, the last fact is very important, while the mean value in latent space might not be the best representative of the likelihood (i.e. the decoded data). This property cannot be respected by the KL divergence, but by the introduced reconstruction term.

1) *of General Expression for Arbitrary Number of Modalities:* Comprising the applied steps to derive \mathcal{L}_{PJ} from the former section, we can identify that by successively applying logarithm and Bayes rules, we derive the ELBO for the proposed multi-modal VAE as follows: First, given the independent set of observable modalities $\mathcal{M} = \{a,b,c,\dots\}$, its marginal log-likelihood $\log p(\mathcal{M}) =: L_{M^2}$ is multiplied by the cardinality of the set as the neutral element $1 = \frac{|\mathcal{M}|}{|\mathcal{M}|}$. Second, applying logarithm

multiplication rule, the nominator is written as the argument's exponent. Third, Bayes rule is applied to each term wrt. the remaining observable modalities to derive their conditionals. Therefore, we can write

$$L_{M^2_{\mathcal{M}}} = \log p(\mathcal{M}) \stackrel{\text{mul. 1}}{=} \frac{|\mathcal{M}|}{|\mathcal{M}|} \log p(\mathcal{M}) \stackrel{\text{log. mul.}}{=} \frac{1}{|\mathcal{M}|} \log p(\mathcal{M})^{|\mathcal{M}|} \quad (79)$$

$$\stackrel{\text{Bayes}}{=} \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \log p(\mathcal{M} \setminus m) p(m | \mathcal{M} \setminus m) \quad (80)$$

$$\stackrel{\text{log. add}}{=} \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \log p(\mathcal{M} \setminus m) + \log p(m | \mathcal{M} \setminus m). \quad (81)$$

The expression $\sum_{m \in \mathcal{M}} \log p(m | \mathcal{M} \setminus m)$ is the general form of the marginal log-likelihood for the *variation of information* (VI), as introduced by [2] for the JMVAE, for any set \mathcal{M} . Thus, it can be directly substituted with $L_{M_{\mathcal{M}}}$. The expression $\sum_{m \in \mathcal{M}} \log p(\mathcal{M} \setminus m)$ is the combination of all joint log-likelihoods of the subsets of \mathcal{M} which have one less element. Therefore, this term can be rewritten as

$$\sum_{m \in \mathcal{M}} \log p(\mathcal{M} \setminus m) = \sum_{\tilde{m} \in \tilde{\mathcal{M}}} \log p(\tilde{m}) \quad (82)$$

with $\tilde{\mathcal{M}} = \{m | m \in \mathcal{P}(\mathcal{M}), |m| = |\mathcal{M}| - 1\}$. Finally, $\log p(\tilde{m})$ can be substituted by $L_{M^2_{\tilde{m}}}$ without loss of generality. However, it is worth noticing that substitution stops at the end of recursion and therefore, all final expressions $\log p(\tilde{m}) \forall |\tilde{m}| \equiv 1$ remain. \square

This results in the final recursive log-likelihood expression from which the ELBO can be directly derived as follows:

$$L_{M^2_{\mathcal{M}}} = \frac{1}{|\mathcal{M}|} \left(L_{M_{\mathcal{M}}} + \sum_{\tilde{m} \in \tilde{\mathcal{M}}} L_{M^2_{\tilde{m}}} \right) \geq \frac{1}{|\mathcal{M}|} \left(\mathcal{L}_{M_{\mathcal{M}}} + \sum_{\tilde{m} \in \tilde{\mathcal{M}}} \mathcal{L}_{M^2_{\tilde{m}}} \right) =: \mathcal{L}_{M^2_{\mathcal{M}}}. \quad (83)$$

VI. APPENDIX

Variants of Bayes equation:

$$p(a) = \frac{p(z,a)}{p(z|a)}, \quad p(z|a) = \frac{p(z,a)}{p(a)}, \quad p(z,a) = \frac{p(a)}{p(z,a)} \quad (84)$$

$$p(a|b,c) \stackrel{\text{eq. 84}}{=} \frac{p(a,b|c)}{p(b|c)} \stackrel{\text{eq. 84}}{=} \frac{p(a,b,c)}{p(b|c)p(c)} \stackrel{\text{eq. 84}}{=} \frac{p(a,b,c)}{p(b,c)} \quad (85)$$

Logarithm rules:

$$\log(ab) = \log(a) + \log(b) \quad (86)$$

$$a \log(b) = \log(b^a) \quad (87)$$

Evidence lower bound:

$$\mathcal{L} = \sum_z q(z|a) \log\left(\frac{p(z,a)}{q(z|a)}\right) \quad (88)$$

Marginal likelihood:

$$p(a|b) = \sum_z p(a|z)p(z|b) \quad (89)$$

Independent and identically distributed random variables (i.i.d. or iid or IID):

$$p(a,b,c) = p(a)p(b)p(c) \quad (90)$$

A. Kulbeck-Leibler Divergence

$$D_{\text{KL}}(q(z|a)||p(z|a)) = \sum_z q(z|a) \log\left(\frac{q(z|a)}{p(z|a)}\right) \quad (91)$$

$$D_{\text{KL}}(\mathcal{N}_1(\mu_1, \sigma_1)||\mathcal{N}_2(\mu_2, \sigma_2)) = \log(\sigma_2) - \log(\sigma_1) + \frac{\sigma_1^2}{2\sigma_2^2} + \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \quad (92)$$

more tbd.

B. Variation of Information

1) Operator Names:

- $\text{VI}(A,B)$: Variation of Information between some properties A and B
- $\text{I}(A)$: Information of A (or mutual information)
- $\text{I}(A,B)$: *Mutual Information* (MI) of A and B
- $\text{I}(A,B|C)$: *Mutual Conditional Information* (MCI) of A and B given C
- $\text{H}(A)$: Entropy of A
- $\text{H}(A,B)$: *Joint Entropy* (JE) of A and B
- $\text{H}(A|B)$: *Conditional Entropy* (CE) of A given B

The *Variation of Information* (VI) between some random variables can be written as

$$\text{VI}(A,B) = \text{H}(A) + \text{H}(B) - 2\text{I}(A,B) = \text{H}(A|B) + \text{H}(B|A) \quad (93)$$

and

$$\text{VI}(A,B,C) = \text{H}(A) + \text{H}(B) + \text{H}(C) - 3\text{I}(A,B) = \text{H}(A|B,C) + \text{H}(B|A,C) + \text{H}(C|A,B) \quad (94)$$

and so on

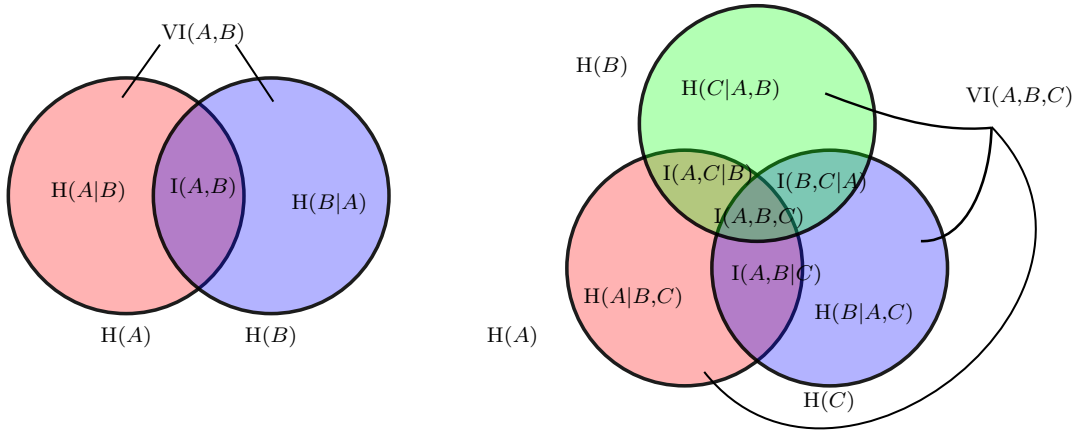


Fig. 2: Visualization of VI as Venn digram.

C. Extension to three Modalities

It should be clear that both approaches, proposed and Suzuki's [2], can be extended to multiple modalities. In the following, an example for three modalities is given.

First, the conditional probability is investigated:

$$p(a|b,c) = \frac{p(a,b,c,z)}{p(a,b,c)} \quad \text{multiplied by 1 \& Equation 85} \quad (95)$$

$$= \frac{p(z,a|b,c)}{p(z|a,b,c)} \quad \text{Equation 85} \quad (96)$$

$$= \frac{1}{p(z|a,b,c)} p(z,a|b,c) \quad \text{reorder} \quad (97)$$

$$= \frac{1}{p(z|a,b,c)} \frac{p(z,a,b,c)}{p(b,c)} \quad \text{Equation 85} \quad (98)$$

$$= \frac{1}{p(z|a,b,c)} \frac{p(a,b,c|z)p(z)}{p(b,c)} \quad \text{Equation 90} \quad (99)$$

$$= \frac{1}{p(z|a,b,c)} \frac{p(a|z)p(b,c|z)p(z)}{p(b,c)} \quad \text{Equation 85} \quad (100)$$

$$= \frac{1}{p(z|a,b,c)} \frac{p(a|z)p(z|c,b)\frac{p(c,b)}{p(z)}p(z)}{p(b,c)} \quad \text{Equation 85} \quad (101)$$

$$= \frac{p(a|z)p(z|c,b)}{p(z|a,b,c)} \quad (102)$$

The log-likelihood of a single joint distribution can be written as:

$$\log(p(a|b,c)) = \sum_z q(z|a,b,c) \log\left(\frac{p(z,a|b,c)}{q(z|a,b,c)}\right) + \sum_z q(z|a,b,c) \log\left(\frac{q(z|a,b,c)}{p(z|a,b,c)}\right) \quad (103)$$

$$= \mathcal{L}_{\tilde{M}_a} + \text{D}_{\text{KL}}(q(z|a,b,c) \| p(z|a,b,c)) \quad (104)$$

$$\geq \mathcal{L}_{\tilde{M}_a} \quad (105)$$

1) *JMVAE for three Modalities*: The log-likelihood of the VI between three distributions can be written as:

$$L_{3M} = \log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|b,c)) \quad (106)$$

$$= \mathcal{L}_{\tilde{M}_a} + \mathcal{L}_{\tilde{M}_b} + \mathcal{L}_{\tilde{M}_c} + 3 \text{D}_{\text{KL}}(q(z|a,b,c) \| p(z|a,b,c)) \quad \text{Equation 41} \quad (107)$$

$$\geq \mathcal{L}_{\tilde{M}_a} + \mathcal{L}_{\tilde{M}_b} + \mathcal{L}_{\tilde{M}_c} \quad (108)$$

The combined ELBO can then be rewritten as

$$\mathcal{L}_{\tilde{M}_a} + \mathcal{L}_{\tilde{M}_b} + \mathcal{L}_{\tilde{M}_c} = \mathbb{E}_{q(z|a,b,c)} \log(p(a|z)) - D_{\text{KL}}(q(z|a,b,c) \| p(z|c,b)) \quad (109)$$

$$+ \mathbb{E}_{q(z|a,b,c)} \log(p(b|z)) - D_{\text{KL}}(q(z|a,b,c) \| p(z|a,c)) \quad (110)$$

$$+ \mathbb{E}_{q(z|a,b,c)} \log(p(c|z)) - D_{\text{KL}}(q(z|a,b,c) \| p(z|a,b)) \quad \text{Equation 92} \quad (111)$$

$$\geq \mathcal{L}_{\tilde{J}} - D_{\text{KL}}(q(z|a,b,c) \| p(z|b,c)) \quad (112)$$

$$- D_{\text{KL}}(q(z|a,b,c) \| p(z|a,c)) - D_{\text{KL}}(q(z|a,b,c) \| p(z|b,c)) \quad (113)$$

$$:= \mathcal{L}_{\tilde{M}} \quad (114)$$

$\mathcal{L}_{\tilde{J}}$ is the joint ELBO of a joint probability distribution having three arguments (i.e. a, b, c). The derivation is analog to Section III. The next steps are the application of encoders and decoders for this network which is straight forward and should be clear to the reader.

However, if we investigate the last equations, the following properties can be identified: There are the common reconstruction terms (\mathbb{E}) for each decoder $p(*|z)$ wrt. the full multi-modal decoder $q(z|a,b,c)$. The KL terms show the **drawback** of the VI approach. As before, these regularizer tend to bring the encoders' distribution to match each other. But only pairwise encoders (e.g. $p(z|a,b)$) remain and thus, uni-modal encoders are neglected.

This means from a practical point of view, that when we have N modalities in a setup, we only can build derived setups having $N - 1$ modalities.

2) *Proposed JMVAE for three Modalities*: The derivation from the joint log likelihood can be written analogously:

$$\log(p(a,b,c)) = \frac{3}{3} \log(p(a,b,c)) \quad (115)$$

$$= \frac{1}{3} \log(p(a,b,c)^3) \quad (116)$$

$$= \frac{1}{3} \log(p(a,b,c)p(a,b,c)p(a,b,c)) \quad (117)$$

$$= \frac{1}{3} \log(p(a,b)p(b,c)p(a,c)p(a|b,c)p(b|a,c)p(c|a,b)) \quad (118)$$

$$= \frac{1}{3} (\log(p(a,b)) + \log(p(b,c)) + \log(p(a,c)) + \log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|a,b))) \quad (119)$$

$$= \frac{1}{3} \left(\frac{2}{2} (\log(p(a,b)) + \log(p(b,c)) + \log(p(a,c))) + \log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|a,b)) \right) \quad (120)$$

$$= \frac{1}{6} (\log(p(a,b)^2) + \log(p(b,c)^2) + \log(p(a,c)^2)) \quad (121)$$

$$+ \frac{1}{3} (\log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|a,b))) \quad (122)$$

$$= \frac{1}{6} (\log(p(a)p(b)p(a|b)p(b|a)) + \log(p(c)p(b)p(c|b)p(b|c)) + \log(p(a)p(c)p(a|c)p(c|a))) \quad (123)$$

$$+ \frac{1}{3} (\log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|a,b))) \quad (124)$$

$$= \frac{1}{6} (\log(p(a|b)) + \log(p(b|a)) + \log(p(c|b)) + \log(p(b|c)) + \log(p(a|c)) + \log(p(c|a))) \quad (125)$$

$$+ \frac{1}{3} (\log(p(a)) + \log(p(b)) + \log(p(c)) + \log(p(a|b,c)) + \log(p(b|a,c)) + \log(p(c|a,b))) \quad (126)$$

It is now straight forward, by applying all former mentioned equations, to derive the ELBO for the above marginal log-likelihood. As one can imagine, the above equation results in a pretty heavy loss term but with the big advantage of respecting all permutations of modalities.

This means again from a practical point of view, in comparison to the approach by Suzuki et al. [2], that we can build arbitrary sensor setups having 1 to N modalities.

REFERENCES

- [1] D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes," *CoRR*, vol. abs/1312.6, 2013.
- [2] M. Suzuki, K. Nakayama, and Y. Matsuo, "Joint multimodal learning with deep generative models," pp. 1–12, 2017.