
Appendix for A Novel Approach for Effective Multi-View Clustering with Information-Theoretic Perspective

Anonymous Submission

1 A Formulas about Mutual Information

2 In this section we introduce several formulas about mutual information that are used in this work.

(P1) Positivity:

$$I(x; y) \geq 0, I(x; y|z) \geq 0$$

(P2) Chain rule:

$$I(xy; z) = I(y; z) + I(x; z|y)$$

(P3) Chain rule (Multivariate Mutual Information):

$$I(x; y; z) = I(y; z) - I(y; z|x)$$

3 B Theorems and Loss Computation

4 B.1 Decomposition of $I(x^i; z^i)$

Theorem B.1. Let z^i be a representation of x^i , then the mutual information between x^i and z^i can be decomposed as:

$$I(x^i, z^i) = I(z^i, x^i|x^j) + I(x^j, z^i).$$

5 **Proof.** Hypothesis (H1): z^i is a representation of x^i : $I(x^j; z^i|x^i) = 0$.

6 Based on the Hypothesis (H1) and Chain rule (P3), then we have

$$\begin{aligned} I(x^i; z^i) &\stackrel{(P3)}{=} I(x^i; z^i|x^j) + I(x^j; z^i; x^j) \\ &\stackrel{(P3)}{=} I(x^i; z^i|x^j) + I(x^j; z^i) - I(x^j; z^i|x^i) \\ &\stackrel{(H1)}{=} I(x^i; z^i|x^j) + I(x^j; z^i) \end{aligned}$$

7

□

8 B.2 Multi-View Redundancy and Sufficiency [2]

Theorem B.2. Let x^i, x^j , and y be random variables with joint distribution $p(x^i, x^j, y)$. Let z^i be a representation of x^i , then:

$$I(x^i; y|z^i) \leq I(x^i; x^j|z^i) + I(x^j; y|x^j).$$

9 **Proof.** Hypothesis (H1): z^i is a representation of x^i : $I(y; z^i | x^j x^i) = 0$.
 10 Based on the Hypothesis (H1), Positivity (P1) and Chain rule (P3), we have

$$\begin{aligned}
 I(x^i; y | z^i) &\stackrel{(P3)}{=} I(x^i; y | z^i x^j) + I(x^i; x^j; y | z^i) \\
 &\stackrel{(P3)}{=} I(x^i; y | x^j) - I(x^i; y; z^i | x^j) + I(x^i; x^j; y | z^i) \\
 &\stackrel{(P3)}{=} I(x^i; y | x^j) - I(y; z^i | x^j) + I(y; z^i | x^j x^i) + I(x^i; x^j; y | z^i) \\
 &\stackrel{(P1)}{\leq} I(x^i; y | x^j) + I(y; z^i | x^j x^i) + I(x^i; x^j; y | z^i) \stackrel{(H1)}{=} I(x^i; y | x^j) + I(x^i; x^j; y | z^i) \\
 &\stackrel{(P3)}{=} I(x^i; y | x^j) + I(x^i; x^j | z^i) - I(x^i; x^j | z^i y) \\
 &\stackrel{(P1)}{\leq} I(x^i; y | x^j) + I(x^i; x^j | z^i)
 \end{aligned}$$

11 **Theorem B.3.** Let x^i be a redundant view with respect to x^j for y . Any representation z^i of x^i that is sufficient for x^j is also sufficient for y , i.e.,

$$I(x^i; x^j | z^i) = 0 \rightarrow I(x^i; y | z^i) = 0.$$

12 **Proof.** Hypotheses (H1): z^i is a representation of x^i : $I(y; z^i | x^i) = 0$.

13 (H2): x^i is redundant with respect to x^j for y : $I(y; x^i | x^j) = 0$.

14 Using the results from Theorem B.2 $I(x^i; y | z^i) \stackrel{(Th.B.2)}{\leq} I(x^i; y | x^j) + I(x^i; x^j | z^i) \stackrel{(H2)}{=} I(x^i; x^j | z^i)$, we have
 15 $I(x^i; x^j | z^i) = 0 \rightarrow I(x^i; y | z^i) = 0$. \square

Theorem B.4. Let x^i , x^j and y be random variables with distribution $p(x^i, x^j, y)$. Let z^i be a representation of x^i , then

$$I(y; z^i) \geq I(y; x^i x^j) - I(x^i; x^j | z^i) - I(x^i; y | x^j) - I(x^j; y | x^i).$$

16 **Proof.** Hypothesis (H1): z^i is a representation of x^i : $I(y; z^i | x^i x^j) = 0$.

$$\begin{aligned}
 I(y; z^i) &\stackrel{(P3)}{=} I(y; z^i | x^i x^j) + I(y; x^i x^j; z^i) \\
 &\stackrel{(H1)}{=} I(y; x^i x^j; z^i) \\
 &\stackrel{(P3)}{=} I(y; x^i x^j) - I(y; x^i x^j | z^i) \\
 &\stackrel{(P2)}{=} I(y; x^i x^j) - I(y; x^i | z^i) - I(y; x^j | z^i x^i) \\
 &\stackrel{(P3)}{=} I(y; x^i x^j) - I(y; x^i | z^i) - I(y; x^j | x^i) + I(y; x^j; z^i | x^i) \\
 &\stackrel{(P3)}{=} I(y; x^i x^j) - I(y; x^i | z^i) - I(y; x^j | x^i) + I(y; z^i | x^i) - I(y; z^i | x^i x^j) \\
 &\stackrel{(H1)}{=} I(y; x^i x^j) - I(y; x^i | z^i) - I(y; x^j | x^i) + I(y; z^i | x^i) \\
 &\stackrel{(P1)}{\geq} I(y; x^i x^j) - I(y; x^i | z^i) - I(y; x^j | x^i) \\
 &\stackrel{(Th.B.3)}{\geq} I(y; x^i x^j) - I(x^i; y | x^j) - I(x^i; x^j | z^i) - I(y; x^j | x^i)
 \end{aligned}$$

17 **Corollary B.1.** Let x^i and x^j be mutually redundant views for y . Let z^i be a representation of x^i that is sufficient for x^j . Then we have:

$$I(y; z^i) = I(x^i x^j; y).$$

18 **Proof.** Hypotheses (H1): z^i and z^j are mutually redundant for y : $I(y; x^i | x^j) + I(y; x^j | x^i) = 0$.

19 (H2): z^i is sufficient for x^j : $I(x^j; x^i | z^i) = 0$.

20 Using the Theorem B.4 and Hypotheses (H1)-(H2), we have

$$\begin{aligned}
 I(y; z^i) &\stackrel{(Th.B.4)}{\geq} I(y; x^i x^j) - I(x^i; y | x^j) - I(x^i; x^j | z^i) - I(y; x^j | x^i) \\
 &\stackrel{(H1)}{=} I(y; x^i x^j) - I(x^i; x^j | z^i) \\
 &\stackrel{(H2)}{=} I(y; x^i x^j)
 \end{aligned}$$

21 Since $I(y; z^i) \leq I(y; x^i x^j)$ is a consequence of the data processing inequality, we conclude that
 22 $I(y; z^i) = I(y; x^i x^j)$. \square

23 B.3 Consistent Variational Lower Bound for SCMVC

24 Consistent variational lower bound provides an efficient approximate posterior inference of the latent
 25 variable \vec{z}, y^i given an observed value x^i . We maximize the following formulation:

$$\max \log q_\phi(x^i) = D_{KL}(p_\theta(\vec{z}, y^i | x^i) \| q_\phi(\vec{z}, y^i | x^i)) + L(\phi, \theta; x^i), \quad (1)$$

as $D_{KL}(p_\theta(\vec{z}, y^i | x^i) \| q_\phi(\vec{z}, y^i | x^i)) \geq 0$, then we have

$$\log q_\phi(x^i) \geq L(\phi^i, \theta^i; x^i) = \mathbb{E}_{p_\theta(\vec{z}, y^i | x^i)}[-\log p_\theta(\vec{z}^i, y^i | x^i) + \log q_\phi(x^i, \vec{z}^i, y^i)],$$

26 Inspired by [1], which is better at generalization due to its emphasis on disentangled representations
 27 and robustness during training, $L(\phi, \theta, x^i)$ can be expressed as follows:

$$\max L_2^i = \mathbb{E}_{p_\theta(\vec{z}, y^i | x^i)}[\log q_\phi(x^i | \vec{z}, y^i)] - \gamma D_{KL}(p_\theta(\vec{z}, y^i | x^i) \| q_\phi(\vec{z}, y^i)), \quad (2)$$

As we assume that

$$q_\phi(x^i | \vec{z}, y^i) = q_\phi(x^i | \vec{z}), p_\theta(\vec{z}, y^i | x^i) = p_\theta(y^i | \vec{z}) p_\theta(\vec{z} | x^i),$$

28 $\mathbb{E}_{p_\theta(\vec{z}, y^i | x^i)}[\log q_\phi(x^i | \vec{z}, y^i)]$ is equivalent as

$$\begin{aligned} \mathbb{E}_{p_\theta(\vec{z}, y^i | x^i)}[\log q_\phi(x^i | \vec{z}, y^i)] &= \iint \sum_y p_\theta(\vec{z}, y^i | x^i) \log q_\phi(x^i | \vec{z}) dz dx \\ &= \iint \sum_y p_\theta(y^i | \vec{z}) p_\theta(\vec{z} | x^i) \log q_\phi(x^i | \vec{z}) dz dx \\ &= \iint p_\theta(\vec{z} | x^i) \log q_\phi(x^i | \vec{z}) dz dx \end{aligned}$$

29 For $D_{KL}(p_\theta(\vec{z}, y^i | x^i) \| q_\phi(\vec{z}, y^i))$,

$$\begin{aligned} D_{KL}(p_\theta(\vec{z}, y^i | x^i) \| q_\phi(\vec{z}, y^i)) &= \sum_y \iint p_\theta(\vec{z}, y^i | x^i) \ln \frac{p_\theta(\vec{z}, y^i | x^i)}{q_\phi(\vec{z}, y^i)} dz dx \\ &= \sum_y \iint p_\theta(y^i | \vec{z}) p_\theta(\vec{z} | x^i) \ln \frac{p_\theta(y^i | \vec{z}) p_\theta(\vec{z} | x^i)}{q_\phi(\vec{z} | y^i) q_\phi(y^i)} dz dx \\ &= \mathbb{E}_{x \sim \tilde{p}(x)} [\mathbb{E}_{\vec{z} \sim \tilde{p}(\vec{z} | x^i)} [\sum_y p_\theta(y^i | \vec{z}) KL(p_\theta(\vec{z} | x^i) \| q_\phi(\vec{z} | y^i)) + KL(p_\theta(y^i | \vec{z}) \| q_\phi(y^i))]] \end{aligned}$$

Given

$$q(\vec{z} | x^i) = \frac{1}{\prod_{t=1}^d \sqrt{2\pi\sigma_t^2(x^i)}} \exp\{-\frac{1}{2} \|\frac{\vec{z}^i - \mu(x^i)}{\sigma(x^i)}\|^2\}$$

30 and

$$q(\vec{z} | y_j^i) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp(-\frac{1}{2} \|\vec{z} - \mu_j^i\|^2),$$

31 the objective function of SCMVC can be formulated as follows:

$$\max L_{con}^i \leftrightarrow \max \sum_{m=1}^n -|x_m^i - \tilde{x}_m^i|^2 + \sum_{j=1}^k \gamma(y_{m,j}^i) (\sum_{t=1}^d \log \sigma_t^2(x)) - \frac{1}{2} \|\vec{z} - \mu_j^i\|^2 + \log(y_{m,j}^i))$$

32 where d is the dimension of hidden layer features, y_{ij}^v (obtained from z_i^v through softmax layer) is
 33 the probability of instance i assigned to class j in view- v , and \tilde{x}^i represents the reconstruction item
 34 through the decoder. For μ_j^i , which is initialized randomly, is the mean vector of the class of i -th view.

35 B.4 Bayes Error Rates for Arbitrary Learned Representations

36 **Theorem B.5.** *The Bayes Error Rates for Arbitrary Learned Representations can be minimized by*
 37 *maximizing L_c on $+\beta L_{suf}^i$, formally:*

$$\max L_{con}^i + \beta L_{suf}^{ij} \leftrightarrow \min \bar{p}_e. \quad (3)$$

Proof. *Our reparameterization strategy transforms z^i into a normal distribution that maximizes $H_{\theta^*}^{(n)}(z^i)$. Moreover, we obtain consistent information among views, as expressed in consistent lower bound, which minimizes $\hat{H}_{\theta^*}^{(n)}(z^i|x^j)$. As we have $\hat{I}_{\theta^*}^{(n)}(z^i; x^j) = \hat{H}_{\theta^*}^{(n)}(z^i) - \hat{H}_{\theta^*}^{(n)}(z^i|x^j)$, the resulting $\hat{I}_{\theta^*}^{(n)}(z^i; x^j)$ is maximized. As for minimizing $I(z^i; x^i|x^j, T)$, given Proof of the sufficient representation lower bound in Section 3.3 the conclusion can be stated as follows:*

$$\max -D_{KL}(p(z^i|x^i), p(z^j|x^j)) \leftrightarrow \min I(z^i; x^i|x^j, T).$$

38 *Overall, the Bayes Error Rates for Arbitrary Learned Representations are minimized:*

$$\max L_{con}^i + \beta L_{suf}^{ij} \leftrightarrow \min \bar{p}_e. \quad (4)$$

39 □

40 C Sufficient Representation Loss

41 Based on the assumption that $\{z^i\}_{i=1}^v$ follows a Gaussian distribution, $D_{KL}(p(z^i|x^i), p(z^j|x^j))$ can be
 42 directly optimized via minimizing the KL divergence:

$$D_{KL}(p(z^i|x^i), p(z^j|x^j)) = \frac{1}{2} [\text{tr}((\Sigma^i)^{-1} \Sigma^j) + (\mu^i - \mu^j)^T (\Sigma^i)^{-1} (\mu^i - \mu^j) - d + \ln \frac{|\Sigma^i|}{|\Sigma^j|}], \quad (5)$$

43 where tr stands for the trace of a matrix defined as the sum of the diagonal elements of the matrix,
 44 Σ^i, Σ^j are the diagonal matrices where the diagonal elements are the variances of z^i and z^j in each
 45 dimension respectively, and d is the dimension of embeddings $\{z^i\}_{i=1}^v$. The mutual information
 46 between the two representations can be maximized by using any sample-based differentiable mutual
 47 information lower bound like InfoNCE. In our methods we found that the predictive information
 48 term didn't play a role in improving final clustering results. As we focus on discarding superfluous
 49 information to finetune the model, we discard the signal-relevant term during finetuning. It's interesting
 50 to further explore this discovery in the future work.

51 D More Details of the Proposed Methods

52 D.1 Implementation Details

53 We use the convolutional autoencoder (CAE) for each view to learn embedded features. The encoder is
 54 Input-Conv₄³²-Conv₄⁶⁴-Conv₄⁶⁴-Fc₂₀₀ and the decoders are symmetric with the corresponding encoders.
 55 For the comparing methods, we use open-source codes with the settings recommended by the authors.
 56 For all datasets an encoder and a decoder are shared for all views. The experiments were conducted
 57 utilizing a Windows PC installed with an Intel(R) Core(TM) i5-12600K CPU@3.69 GHz, 32.0
 58 GB RAM, and GeForce RTX 3070ti GPU (8 GB cache). All baselines were tuned to their highest
 59 performance according to the respective papers to ensure fair comparison. We set β and γ to 0.1
 60 for all experiments. To better understand our approach, we provide the framework of SCMVC and
 61 SUMVC in Figs. 1-2.

62 D.2 Visualization of Learned Generative Models

63 We show the reconstructed pictures and sampled pictures generated on the Multi-Mnist and Multi-
 64 COIL-20 respectively (Figs. 3-5). We utilize the trained model, set the batch size to 32, and display
 65 the images of each view for three epochs.

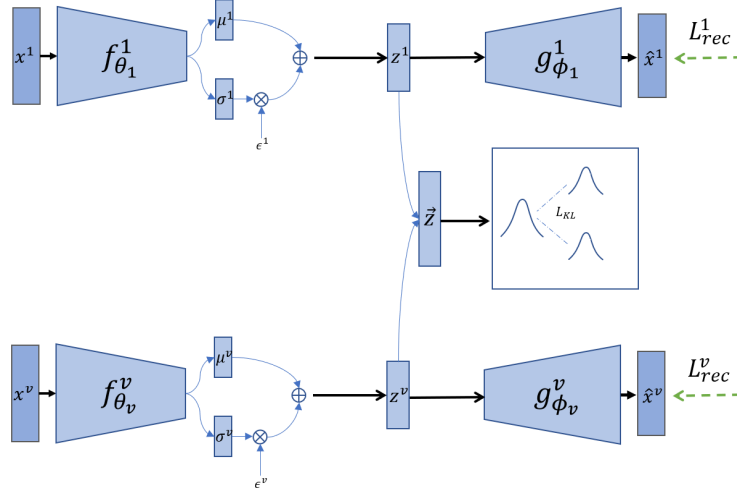


Figure 1: The framework of SCMVC. For the v -th view, x^v denotes the input data, z^v denotes the embedded features, $f_{\theta_v^v}$ and $g_{\phi_v^v}$ are the encoder and the decoder of v -th view respectively.

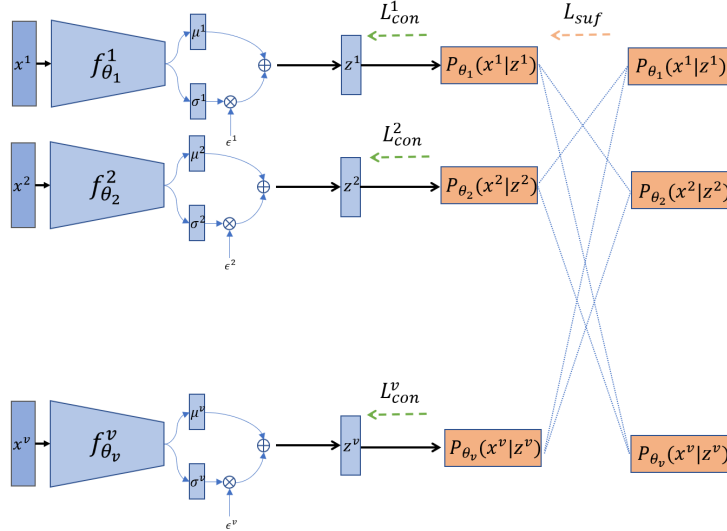


Figure 2: The framework of SUMVC. We utilize a collaborative training approach whereby each view sequentially serves as a guide to facilitate other views' learning process.

66 D.3 Limitation of the Proposed Methods

67 While our proposed methods have demonstrated promising results in addressing the multi-view
68 clustering research problem and have provided valuable insights into strengthening consistency
69 and decreasing view redundancy, it is important to acknowledge several limitations that should be
70 considered. In particular, our model does not perform well on datasets with particularly strong
71 heterogeneity between views, such as huge differences in dimensions of different views. This is
72 because variational inference assumes that all samples in the data come from a common probability
73 distribution. However, this assumption may not hold true when the data is highly heterogeneous,
74 and the distribution may exhibit a range of subgroups and variations. In such cases, the use of
75 variational inference may not be accurate or effective. To the best of our knowledge, this issue
76 remains unexplored within the existing literatures. In the future, we will consider employing
77 ensemble variational inference or gaussian process mixture model as a potential solution to address
78 this problem and achieve a more robust model.

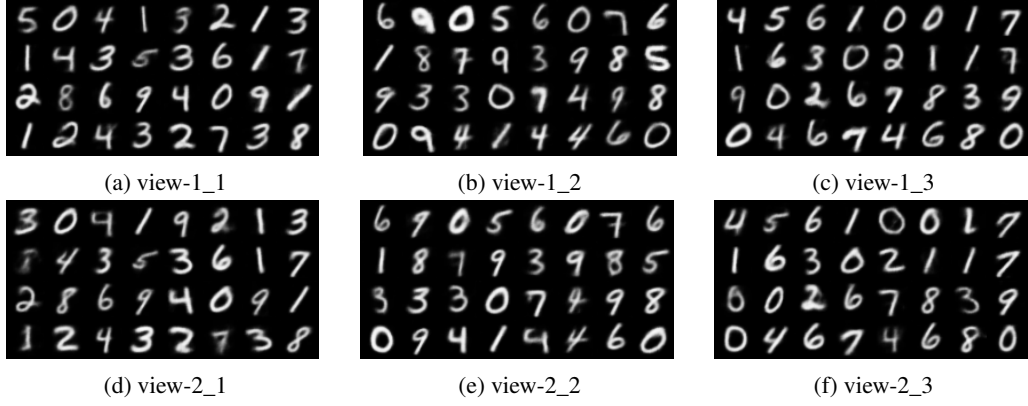


Figure 3: The reconstructed pictures generated on Multi-Mnist.

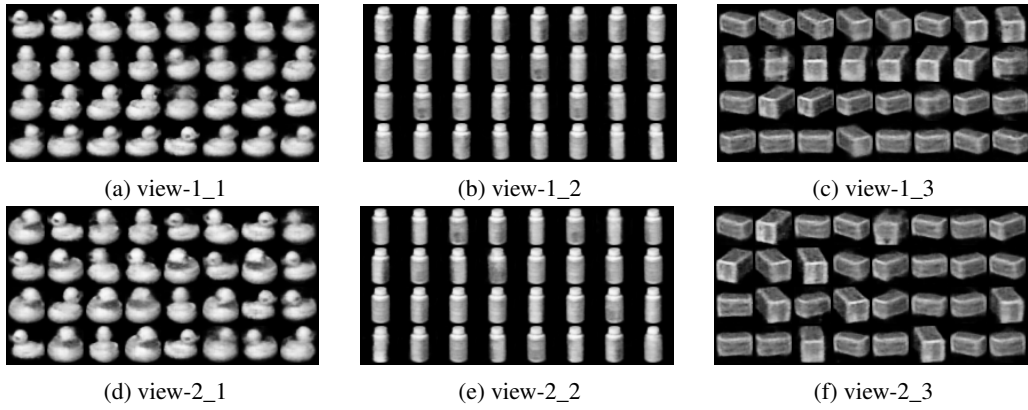


Figure 4: The reconstructed pictures generated on Multi-COIL-20.

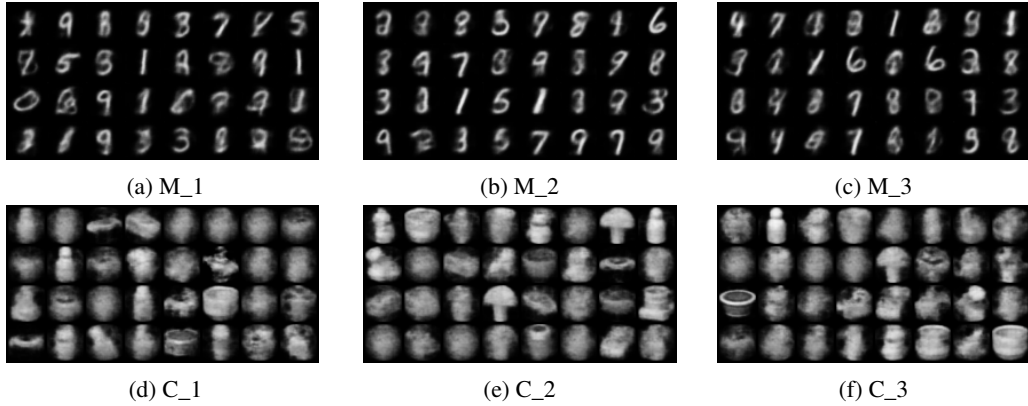


Figure 5: The sampled pictures generated on the Multi-Mnist and Multi-COIL-20.

References

- [1] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -vae. *arXiv preprint arXiv:1804.03599*, 2018.
- [2] Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. Learning robust representations via multi-view information bottleneck. In *ICLR*, 2020.