

Homework 7- Saeideh Khadem- Paper summery

Variational Autoencoders for Cancer Data Integration: Design Principles and Computational Practice

Aim of the research:

- collecting multiple data sets at different genome-scales with the aim to identify novel cancer bio-markers and predict patient survival.
- there is a distinct lack of work on how to integrate diverse patient data and identify the optimal design best suited to the available data.
- Paper is performed extensive analyses of these approaches and provide a clear methodological and computational framework for designing systems that enable clinicians to investigate cancer traits and translate the results into clinical applications.
- Paper is demonstrated how these networks can be designed, built, and, in particular, applied to tasks of integrative analyses of heterogeneous breast cancer data.
-

There are several existing machine learning approaches that integrate diverse data. These can be classified into three different categories based on how the data is being utilized: (i) output (or late) integration, (ii) partial (or intermediate) integration, and (iii) full (or early) integration.

Partial integration refers to specifically designed and developed methods that produce a joint model learned from multiple data simultaneously.

Variational Autoencoders

Generally, an autoencoder consists of two networks, an *encoder* and a *decoder*, which broadly perform the following tasks:

- **Encoder:** Maps the high dimensional input data into a latent variable embedding which has lower dimensions than the input.
- **Decoder:** Attempts to reconstruct the input data from the embedding.

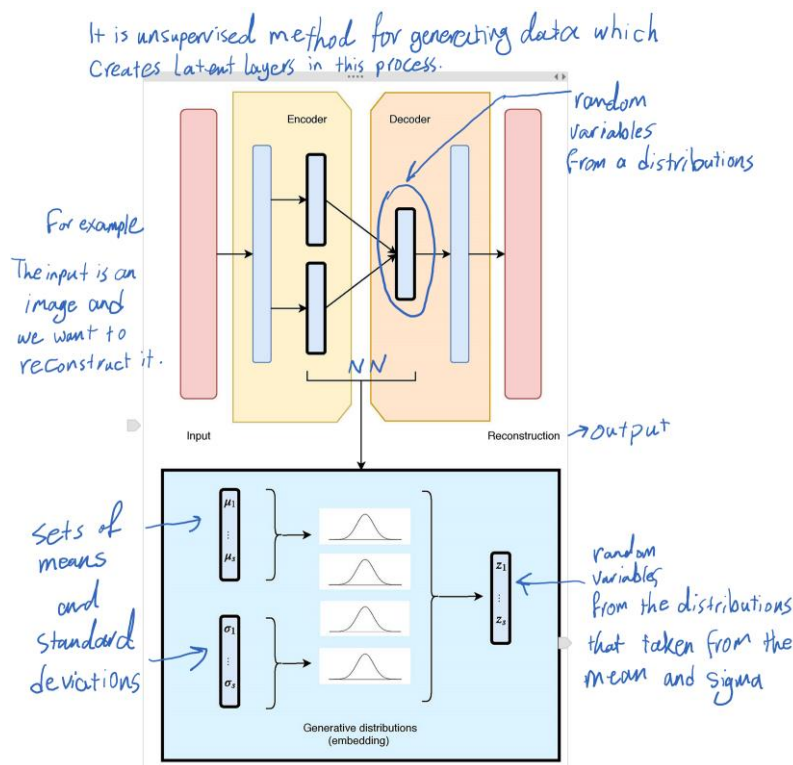


Figure 1 The unimodal Variational Autoencoder (VAE) architecture and latent embedding: the red layers correspond to the input and reconstructed data, given and generated by the model. The hidden layers are in blue, with the embedding framed in black.

Each latent component is made of two nodes (mean and standard deviation), which define a Gaussian distribution. The combination of all Gaussian constitutes the VAE generative embedding.

$$q_{\phi}(z|x^{(i)}) = N(z; \mu^{(i)}, \sigma^{(i)} I) \quad , \quad L_{AE}(\theta, \phi) = \frac{1}{n} \sum_{i=0}^n (x_i - f_{\theta}(g_{\phi}(x_i)))^2 \quad , \quad q_{\phi}(z|x^{(i)}) = N(z; \mu^{(i)}, \sigma^{(i)} I)$$

$$l_i(\theta, \phi) = \underbrace{-E_{q_{\phi}(z|x^{(i)})}[\log p_{\theta}(x|z)]}_{\text{Likelihood}} + \underbrace{KL(q_{\phi}(z|x^{(i)}) || p_{\theta}(z))}_{\text{Prior}}$$

encoder
output
input

How different are two distributions

where the first term corresponds to the reconstruction loss, which encourages the decoder to learn to reconstruct the data from the embedding space. The second term is regularization, and measures the divergence between the encoding distributions $q(z|x)$ and $p(z)$, and penalizes the entanglement between components in the latent space. It is typically estimated by the Kullback–Leibler (KL) divergence, a measure of discrepancy between two probability distributions, which in this case is applied between the prior and the representation.

Moreover, some approaches have experimented with different regularization terms, such as the InfoVAE where **Maximum Mean Discrepancy** (MMD) is employed as an alternative to KL divergence. MMD is based on the concept that two distributions are identical if, and only if, all their moments are identical. Therefore, by employing MMD *via* the kernel embedding trick, the divergence can be defined as the discrepancy between the moments of two distributions $p(z)$ and $q(z)$ as:

$$MMD(q(z)||p(z)) = E_{p(z), p(z')} [k(z, z')] + E_{q(z), q(z')} [k(z, z')] - 2E_{q(z), p(z')} [k(z, z')]$$

where $k(z, z')$ denotes any universal kernel. In this paper, we employ a Gaussian kernel $k(z, z') = e^{-\frac{\|z-z'\|^2}{2\sigma^2}}$ when considering MMD regularization in the objective function.

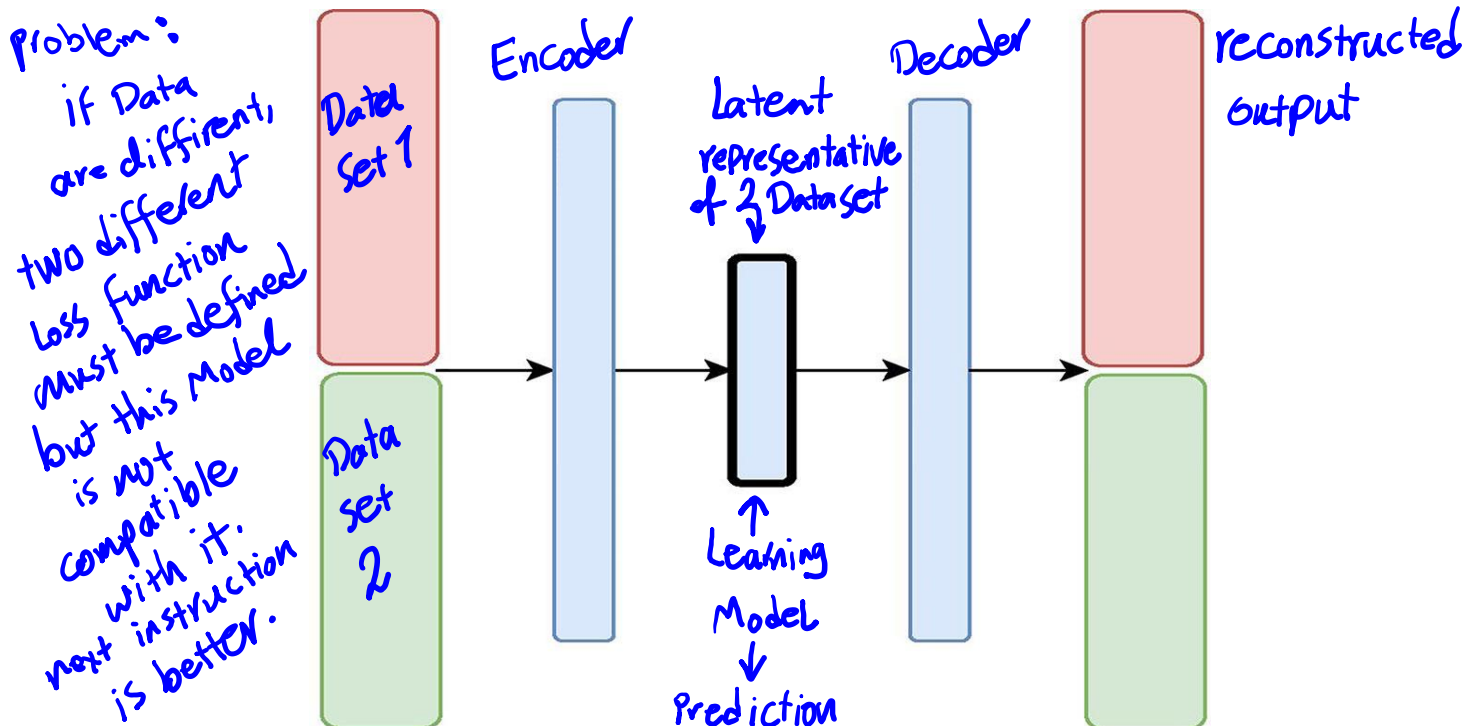


Figure 2 The Variational Autoencoder with Concatenated Inputs (CNC-VAE) Architecture: the red and green layers on the left correspond to two inputs from different data sources. The blue layers are shared, with the embedding being framed in black.

Problem: overfit
on the Data
with higher
Population
since the
Latent representative
is shared.

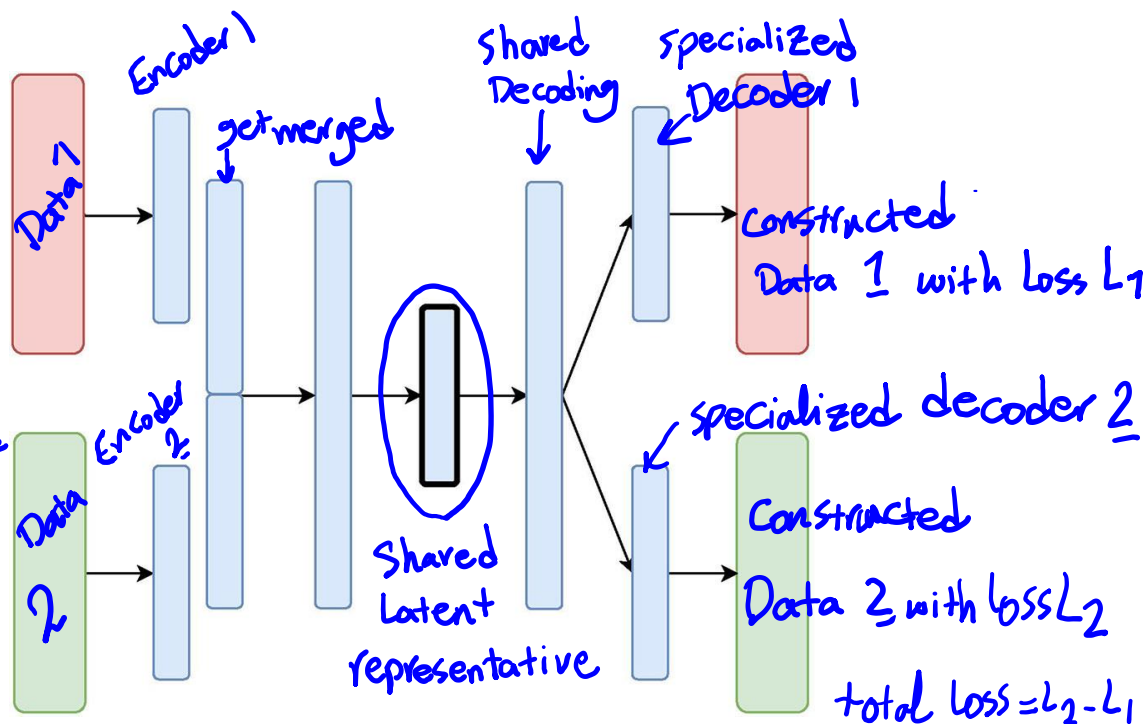


Figure 3 The X-shaped Variational Autoencoder (X-VAE) Architecture: the red and green layers on the left correspond to two inputs from different data sources. The blue layers are shared, with the embedding being framed in black.

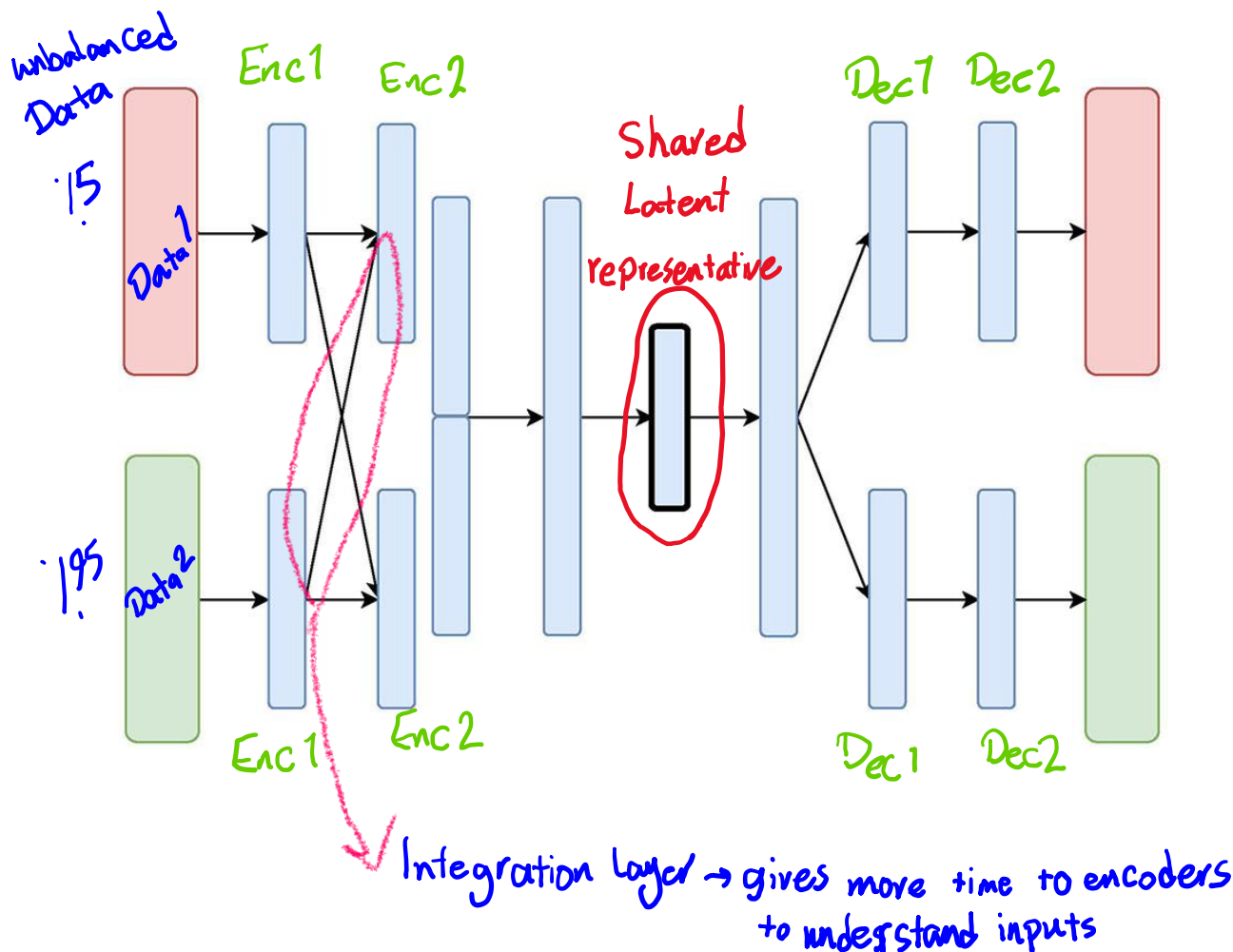


Figure 4 The Mixed-Modal Variational Autoencoder (MM-VAE) Architecture: the red and green layers on the left correspond to two inputs from different data sources. The blue layers are shared, with the embedding being framed in black.

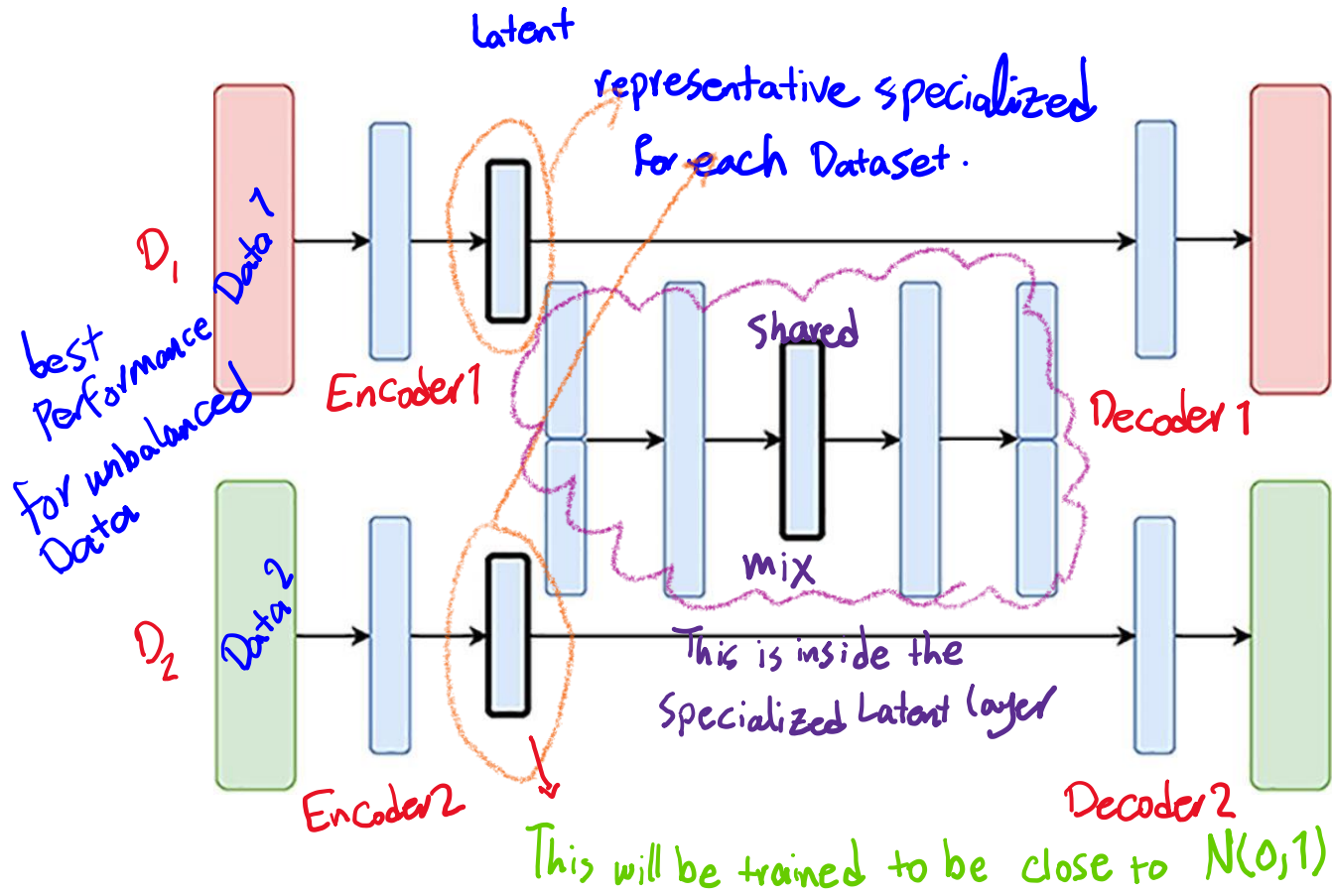


Figure 5 The Hierarchical Variational Autoencoder (H-VAE) Architecture: the red and green layers on the left correspond to two inputs from different data sources. The blue layers are shared, with the embedding being framed in black.

In conclusion, in this study we demonstrate the utility of VAEs for full data integration. The design and the analyses of different integrative VAE architectures and configurations, and in particular their application to the tasks of integrative modeling and analyzing heterogeneous breast cancer data, are the main contributions of this paper. The studied approaches have several distinguishing properties. First, they are able to produce representations that capture the structure (i.e., intrinsic relationships between the data variables) of the data and therefore allow for more accurate downstream analyses. Second, they are able to reduce the dimensionality of the input data without loss of quality or performance. Therefore, in the process of compressing the input data, they can reduce noise implicitly present in the data. Third, they are modular and easily extendable to handle integration of a multitude of heterogeneous data sets. Next, while the integrative VAEs can be used as a data pre-processing approach for learning representations, they can also be utilized in a more generative setting for producing surrogate data, which can be used for more in-depth analysis. Finally, we show that VAEs can be successfully applied to learn representations in complex integrative tasks, such as integrative analyses of breast cancer data, that ultimately lead to more accurate and stable diagnoses.