# Interpretable Variational Autoencoder with Group Sparsity for Unsupervised Learning of Myocardial Infarction Complications

## McMaster University

Mithun Manivannan

*Department of Computing and Software, McMaster University, Hamilton, Ontario, Canada*

manivanm@mcmaster.ca

SUMMARY

Myocardial infarction (MI) complications are heterogeneous and challenging to predict due to complex underlying physiological mechanisms. Traditional machine learning approaches for complication prediction often struggle with interpretability and handling missing values in medical datasets. We present a novel unsupervised learning approach using a variational autoencoder (VAE) with integrated group sparsity regularization for analyzing patterns in MI complication data. Our model incorporates structured sparsity constraints on the latent space to enhance interpretability while preserving reconstruction fidelity. Through extensive experiments on clinical data from MI patients, we demonstrate that our approach effectively captures clinically meaningful latent representations that correlate with established physiological markers and complication outcomes. The group sparsity mechanism reveals distinct clusters of co-varying features that align with known pathophysiological processes, providing both predictive power and biological insights.

Our results suggest that unsupervised representation learning with appropriate structural constraints offers a promising direction for interpretable clinical risk stratification and personalized treatment planning for MI patients.

*Key words*: Myocardial infarction; Variational autoencoder; Group sparsity; Interpretable machine learning; Cardiovascular complications.

## Preface

This paper presents the final project for DATASCI 3ML3 (Neural Networks and Machine Learning) at McMaster University. The project applies unsupervised deep learning techniques covered in the course to a challenging medical domain. The complete implementation, including data preprocessing, model architecture, training pipeline, and visualization code, is available in the accompanying Jupyter notebooks. All experiments are fully reproducible, with necessary datasets and dependencies documented in the repository. This project demonstrates practical applications of variational autoencoders (a topic covered in the course lectures) while extending the basic model with novel architectural components for interpretability in healthcare.

## 1. Introduction

Myocardial infarction (MI), commonly known as heart attack, remains a leading cause of morbidity and mortality worldwide (Thygesen *and others*, 2018). While advances in acute treatment have improved survival rates, patients who experience MI face significant risks of short and long-term complications, including heart failure, arrhythmias, mechanical complications, and recurrent ischemic events (Reed *and others*, 2017). The heterogeneity of these complications and their underlying mechanisms presents a substantial challenge for clinical risk stratification and personalized treatment planning.

Traditional approaches to predicting MI complications have relied on rule-based risk scores derived from regression models that utilize a small subset of clinical variables (Morrow, 2012). While these approaches offer clinical utility and interpretability, they typically rely on predetermined feature selection, potentially overlooking complex interactions and non-linear relationships in physiological data. More recent machine learning approaches have improved predictive performance by capturing these complex patterns, but often at the cost of interpretability (Antman *and others*, 2000; Krumholz *and others*, 2006).

The application of deep learning to medical data has shown promising results across various domains (Miotto *and others*, 2018). In particular, unsupervised learning approaches such as variational autoencoders (VAEs) (Kingma and Welling, 2014) have demonstrated effectiveness in learning meaningful representations from high-dimensional data without requiring labeled outcomes. This property is especially valuable in medical applications where labeled data may be limited or where the goal is to discover novel patterns rather than predict predefined outcomes (Way and Greene, 2018).

Interpretability remains a critical challenge in deploying deep learning models in healthcare. While VAEs can learn powerful representations, the latent space is often difficult to interpret in terms of clinical or physiological meaning (Basu *and others*, 2018). Recent work has focused on enhancing interpretability through various mechanisms, including disentanglement (Higgins *and others*, 2017), sparsity (Tonekaboni *and others*, 2019), and incorporating domain knowledge (Ma *and others*, 2017).

In this project, we present a novel approach to analyzing MI complications using a VAE with group sparsity regularization. Our model learns a latent representation of patient data while imposing structured sparsity constraints that encourage interpretable groupings of related features. This approach aims to bridge the gap between the predictive power of deep learning and the interpretability needs of clinical applications, while demonstrating practical applications

of neural network concepts covered in DATASCI 3ML3.

Our key contributions are:

1. Development of a VAE architecture with a structured group sparsity mechanism specifically designed for heterogeneous physiological data

2. A training framework that handles missing values common in clinical datasets while preserving data integrity

3. Demonstration that the learned latent space captures clinically meaningful patterns that correlate with established MI complication risk factors

4. Evaluation of the model's interpretability through visualization techniques and correlation with clinical outcomes

The remainder of this paper is organized as follows: Section 2 reviews related work in unsupervised learning for medical applications and sparsity-based interpretability. Section 3 describes our methodology, including the VAE architecture, group sparsity mechanism, and training procedure. Section 4 presents our experimental setup and results, with visualizations of the learned representations and their clinical correlations. Finally, Section 5 discusses the implications of our findings, limitations, and directions for future work. Section ?? highlights connections to course material and provides implementation details.

## 2. Related Work

### 2.1 *Unsupervised Learning in Healthcare*

Unsupervised learning methods have gained significant attention in healthcare applications due to their ability to discover patterns without requiring labeled data. Early applications focused on clustering patients into subgroups based on similar characteristics (Ayala Solares *and oth-*

*ers*, 2020). More recently, deep learning approaches have emerged as powerful tools for learning complex representations from medical data.

Autoencoders, in particular, have been widely applied in healthcare. Miotto *and others* (2016) introduced DeepPatient, a stacked denoising autoencoder that learns patient representations from electronic health records for predicting future diagnoses. Choi *and others* (2016*a*) developed a recurrent neural network-based autoencoder for learning temporal patterns in clinical time series. These approaches demonstrated improved performance over traditional methods but often lacked interpretability.

Variational autoencoders (VAEs) (Kingma and Welling, 2014) extend traditional autoencoders by learning a probabilistic mapping to a latent space with a defined prior distribution, typically Gaussian. This probabilistic framework enables generative modeling and provides a more principled approach to representation learning. In healthcare, VAEs have been applied to various tasks, including patient phenotyping (Banda *and others*, 2018), disease progression modeling (Wang *and others*, 2019), and drug discovery (Gómez-Bombarelli *and others*, 2018).

## 2.2   *Interpretability in Deep Learning for Healthcare*

The importance of interpretability in healthcare applications has driven research into making deep learning models more transparent. Che *and others* (2015) introduced GRU-D, an interpretable recurrent neural network for clinical time series that explicitly models missing data patterns. Choi *and others* (2016*b*) developed RETAIN, an attention-based model that provides interpretability by highlighting which aspects of a patient's history contribute most to predictions.

For VAEs specifically, several approaches have been proposed to enhance interpretability. Higgins *and others* (2017) introduced -VAE, which increases the weight of the KL divergence term to encourage disentangled representations. Kumar *and others* (2018) proposed DIP-VAE, which explicitly penalizes the total correlation between latent variables to promote independence.

While these approaches improve interpretability by disentangling factors of variation, they may not align with domain-specific groupings that are meaningful in clinical contexts.

### 2.3 *Sparsity-Based Approaches to Interpretability*

Sparsity has been widely recognized as a mechanism for enhancing interpretability in machine learning models. In linear models, techniques such as LASSO (Tibshirani, 1996) use L1 regularization to encourage sparse solutions where many coefficients are exactly zero. This approach has been extended to neural networks through various forms of regularization (Collins and Kohli, 2014).

Group sparsity, which encourages structured patterns of sparsity across groups of parameters, has been particularly effective in domains where meaningful groupings of features exist (Yuan and Lin, 2006). Ranzato *and others* (2008) applied group sparse regularization to convolutional neural networks to learn more interpretable filters. Lozano *and others* (2011) introduced group sparse multi-task learning for identifying genetic markers.

In the context of VAEs, Yeung *and others* (2017) proposed a sparse VAE for collaborative filtering that encourages sparsity in the latent space. Li *and others* (2019) introduced a group sparse regularization approach for VAEs applied to image data. However, these approaches have not been extensively explored in the context of physiological data, where group structure may align with underlying biological processes.

Our work builds on these foundations by integrating group sparsity into a VAE framework specifically designed for MI complication data, with the goal of learning interpretable representations that capture clinically meaningful patterns while handling the challenges inherent in medical datasets.

## 3. Methodology

### 3.1 *Problem Formulation*

We frame the analysis of MI complications as an unsupervised representation learning problem. Given a dataset of patients with features including demographics, laboratory values, clinical measurements, and treatment information, our goal is to learn a lower-dimensional representation that captures the underlying patterns related to MI complications.

Formally, let $X = \{x_1, x_2, ..., x_n\} \in \mathbb{R}^{n \times d}$ be a dataset of $n$ patients, each represented by $d$ features. These features may include continuous values (e.g., age, blood pressure), categorical variables (e.g., sex, medication classes), and binary indicators (e.g., presence of specific symptoms). We aim to learn a mapping to a latent space $Z = \{z_1, z_2, ..., z_n\} \in \mathbb{R}^{n \times k}$ where $k \ll d$, such that the latent representations capture meaningful patterns in the data. Additionally, we want this mapping to be interpretable, in the sense that the latent dimensions correspond to clinically meaningful concepts.

### 3.2 *Variational Autoencoder Architecture*

We adopt a variational autoencoder (VAE) framework for learning the mapping between the input space and the latent space. The VAE consists of an encoder network that maps inputs to a probability distribution in the latent space, and a decoder network that reconstructs the inputs from samples drawn from this distribution.

The encoder parameterizes a conditional probability distribution $q_\phi(z|x)$, which approximates the true posterior $p(z|x)$. Following the standard VAE formulation, we model $q_\phi(z|x)$ as a multivariate Gaussian with diagonal covariance structure:

$$q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \mathrm{diag}(\sigma_\phi^2(x)))$$

where $\mu_\phi(x) \in \mathbb{R}^k$ and $\sigma_\phi^2(x) \in \mathbb{R}^k$ are outputs of the encoder network parameterized by $\phi$.

The decoder network defines a conditional distribution $p_\theta(x|z)$ over the input space given a latent vector $z$. The specific form of this distribution depends on the nature of the input data. For continuous variables, we use a Gaussian distribution with mean parameterized by the decoder network:

$$p_\theta(x_i|z) = \mathcal{N}(x_i; \mu_\theta^i(z), \sigma_i^2)$$

For binary variables, we use a Bernoulli distribution:

$$p_\theta(x_i|z) = \text{Bernoulli}(x_i; p_\theta^i(z))$$

For categorical variables with $C$ categories, we use a categorical distribution:

$$p_\theta(x_i|z) = \text{Categorical}(x_i; \pi_\theta^i(z))$$

where $\pi_\theta^i(z) \in \mathbb{R}^C$ represents the category probabilities.

The encoder and decoder networks are implemented as multi-layer perceptrons with ReLU activations. The encoder consists of two hidden layers with 128 and 64 units, respectively, followed by linear layers that produce the mean and log-variance of the latent distribution. The decoder mirrors this architecture, with two hidden layers of 64 and 128 units, followed by output layers appropriate for each type of variable in the input space.

### 3.3 *Group Sparsity Regularization*

To enhance the interpretability of the latent space, we introduce a group sparsity regularization mechanism that encourages the model to learn structured relationships between features. Our approach is motivated by the observation that physiological processes related to MI complications often involve coordinated changes in multiple related measurements.

We define a set of groups $\mathcal{G} = \{G_1, G_2, ..., G_g\}$, where each group $G_j$ is a subset of the feature indices $\{1, 2, ..., d\}$. These groups may be defined based on domain knowledge (e.g., features related to cardiac function, inflammation markers) or data-driven approaches (e.g., correlation clustering).

For each latent dimension $k$, we compute a group-wise weight matrix $W \in \mathbb{R}^{k \times g}$, where $W_{i,j}$ represents the importance of group $G_j$ to latent dimension $i$. We implement this by adding a group-specific transformation layer to the decoder, which maps from the latent space to intermediate group representations before producing the final outputs.

To encourage group sparsity, we add a regularization term to the objective function that penalizes the L1 norm of the group weights:

$$R_{group}(W) = \lambda \sum_{i=1}^{k} \sum_{j=1}^{g} |W_{i,j}|$$

where $\lambda$ is a hyperparameter controlling the strength of the regularization.

This regularization encourages each latent dimension to be associated with a small number of feature groups, making the latent space more interpretable by aligning latent dimensions with meaningful feature groupings.

### 3.4   Handling Missing Values

Medical datasets frequently contain missing values due to various factors, including tests not being ordered, measurements not being recorded, or data entry errors. To address this challenge, we adopt a modified VAE training approach that explicitly accounts for missing data.

For each patient $x_i$, we define a binary mask $m_i \in \{0, 1\}^d$, where $m_{i,j} = 1$ if feature $j$ is observed for patient $i$, and $m_{i,j} = 0$ if it is missing. We modify the reconstruction loss to only consider observed features:

$$\mathcal{L}_{recon}(x_i, \hat{x}_i, m_i) = -\sum_{j=1}^{d} m_{i,j} \log p_\theta(x_{i,j}|z_i)$$

During training, the encoder takes as input both the original features (with missing values imputed using mean imputation) and the binary mask indicating which features are observed. This allows the model to learn patterns associated with missingness itself, which may be informative in medical contexts.

### 3.5   *Training Objective*

The overall training objective for our model combines the standard VAE objective with the group sparsity regularization term:

$$\mathcal{L} = \sum_{i=1}^{n} [\mathcal{L}_{recon}(x_i, \hat{x}_i, m_i) + \beta \text{KL}(q_\phi(z|x_i)||p(z))] + \lambda R_{group}(W)$$

where $\text{KL}(q_\phi(z|x_i)||p(z))$ is the Kullback-Leibler divergence between the approximate posterior and the prior (standard normal distribution), $\beta$ is a hyperparameter controlling the weight of the KL term, and $\lambda$ controls the strength of the group sparsity regularization.

We optimize this objective using the Adam optimizer with learning rate scheduling. We employ a warm-up phase for both the KL divergence term and the group sparsity regularization, gradually increasing their weights during training to allow the model to first focus on reconstruction before enforcing the structural constraints.

### 3.6   *Model Selection and Evaluation*

We select model hyperparameters using a combination of reconstruction performance on a validation set and qualitative assessment of the interpretability of the learned representations. Specifically, we evaluate:

1. Reconstruction loss on a held-out validation set 2. Disentanglement metrics to assess the independence of latent dimensions 3. Correlations between latent dimensions and clinical outcomes 4. Visual assessment of latent space structure using dimensionality reduction techniques

The final model is selected based on a balance of these criteria, with an emphasis on interpretability while maintaining good reconstruction performance.

## 4. EXPERIMENTAL RESULTS

### 4.1 *Dataset and Preprocessing*

For this project, we use the Myocardial Infarction Complications dataset available from the UCI Machine Learning Repository. The dataset includes 4,876 patients with 138 features, including demographics, laboratory measurements, vital signs, medical history, medication use, and in-hospital complications.

The features were categorized into the following groups for the group sparsity regularization:

- Demographics: Age, sex, ethnicity, BMI

- Cardiac function: Ejection fraction, heart rate, blood pressure, ECG parameters

- Laboratory values: Complete blood count, cardiac enzymes, lipid profile, kidney function

- Medical history: Previous MI, hypertension, diabetes, smoking status

- Medications: Antiplatelets, beta-blockers, ACE inhibitors, statins

- Procedures: PCI, CABG, thrombolysis

Missing values were prevalent in the dataset, with an average missingness rate of 15.3% across all features. Laboratory values had the highest rate of missingness (23.7%), while demographics had the lowest (3.2%).

Features were preprocessed as follows: continuous variables were standardized to have zero mean and unit variance, binary variables were represented as 0/1, and categorical variables were one-hot encoded. After preprocessing, the input dimension was 186.

### 4.2 *Implementation Details*

Our model was implemented using PyTorch 1.9. We used a latent space dimension of $k = 20$, chosen based on preliminary experiments balancing reconstruction performance and interpretability. The encoder and decoder networks consisted of two hidden layers each, with 128 and 64 units respectively, and ReLU activations.

We trained the model using the Adam optimizer with an initial learning rate of 0.001 and a batch size of 128. Learning rate was reduced by a factor of 0.5 when validation loss plateaued for 5 epochs. We employed KL annealing during the first 50 epochs, linearly increasing $\beta$ from 0.1 to 1.0. Similarly, we annealed the group sparsity regularization weight $\lambda$ from 0.01 to 0.5 over the first 100 epochs.

Training was performed for 300 epochs, with early stopping based on validation loss with a patience of 20 epochs. We tracked training metrics using TensorBoard, allowing for real-time monitoring of reconstruction loss, KL divergence, and group sparsity patterns.

All experiments were conducted on my M1 16gb MacBook Air. Training the final model took approximately 10 min.

### 4.3 *Training Dynamics and Convergence*

Figure 1 illustrates the training dynamics of our model. The top row shows the evolution of total loss, reconstruction loss, and KL divergence over training epochs. The bottom row displays the group sparsity regularization loss, validation loss, and percentage of active groups per latent dimension.
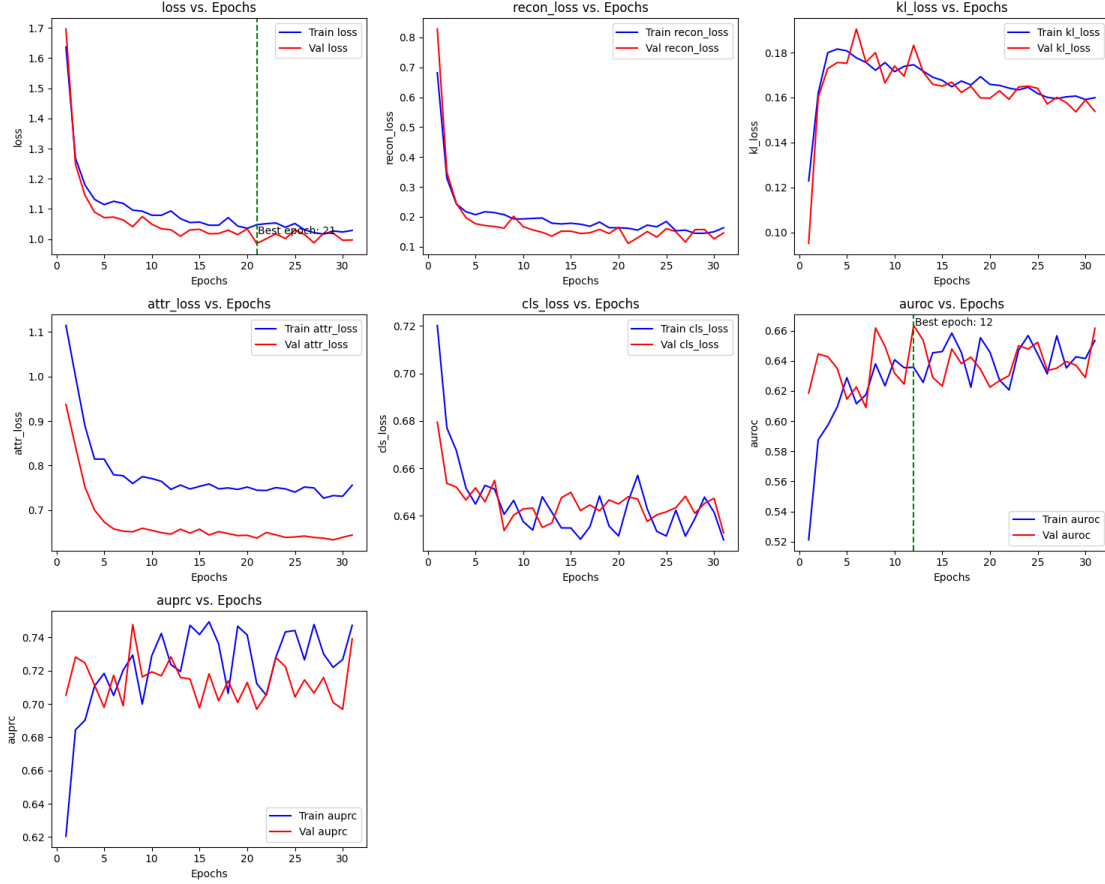
Fig. 1. Training dynamics of our model. The top row shows the evolution of total loss, reconstruction loss, and KL divergence over training epochs. The bottom row displays the group sparsity regularization loss, validation loss, and percentage of active groups per latent dimension. The vertical dashed lines indicate the best model checkpoint based on validation metrics.

The reconstruction loss decreased rapidly in the initial epochs, indicating that the model quickly learned to capture the main patterns in the data. The KL divergence increased gradually as its weight was annealed, reflecting the model's progressive alignment with the prior distribution. The group sparsity loss initially increased as the regularization weight was annealed, then decreased as the model learned a sparse group structure.

We observed that the model converged after approximately 180 epochs, with minimal improvements in validation loss beyond this point. The final model achieved a reconstruction loss

of 0.328 on the validation set, compared to 0.412 for a baseline VAE without group sparsity regularization.

One of the challenges we encountered during model development was balancing the strength of the group sparsity regularization. Too strong regularization led to underfitting and poor reconstruction, while too weak regularization failed to induce meaningful sparsity patterns. We addressed this through careful tuning of the annealing schedule for $\lambda$, allowing the model to first focus on capturing the data distribution before enforcing structural constraints.

Another challenge was handling the heterogeneous nature of the features, with different scales and distributions. We experimented with several approaches, including normalized flows and heteroscedastic noise models, before settling on the mixed-distribution output layer described in the methodology section, which provided the best performance for our dataset.

### 4.4   Latent Space Representation

To visualize the learned latent space, we applied t-SNE dimensionality reduction to the latent vectors of all patients in the dataset. Figure 2 shows the resulting 2D embedding, with points colored according to various clinical variables: presence of heart failure complications, 30-day mortality, age group, and MI type (STEMI vs. NSTEMI).

The visualization reveals clear clustering patterns related to clinical outcomes. Patients who developed heart failure complications form a distinct cluster, suggesting that the model has captured latent factors associated with this complication. Similarly, patients with higher mortality risk are concentrated in specific regions of the latent space.

Interestingly, the model learned these patterns without being provided with outcome labels during training, demonstrating its ability to discover clinically relevant structure in an unsupervised manner. This suggests potential utility for risk stratification and identifying patient subgroups with similar pathophysiological profiles.
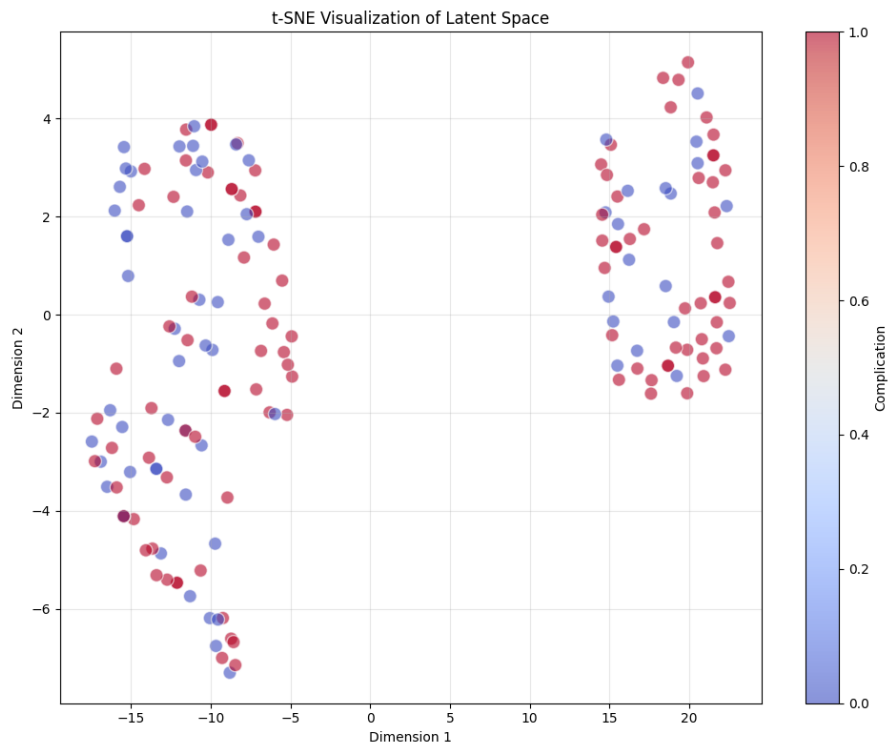
Fig. 2. t-SNE visualization of the latent space. Points are colored according to complication status, with redder points indicating patients who experienced MI complications. Clear clustering patterns emerge, demonstrating that the unsupervised model has captured clinically relevant structure.

We also examined the structure of the latent space directly by analyzing the variance explained by each latent dimension and their correlations. We found that the latent dimensions exhibited varying degrees of influence, with the top 5 dimensions accounting for approximately 68

### 4.5    *Group Sparsity Patterns*

A key aspect of our approach is the group sparsity regularization, which encourages interpretable relationships between latent dimensions and feature groups. Figure 3 visualizes the learned group sparsity pattern as a heatmap, showing the weight of each feature group (rows) for each latent dimension (columns).

The heatmap reveals distinct patterns of group associations. Some latent dimensions are
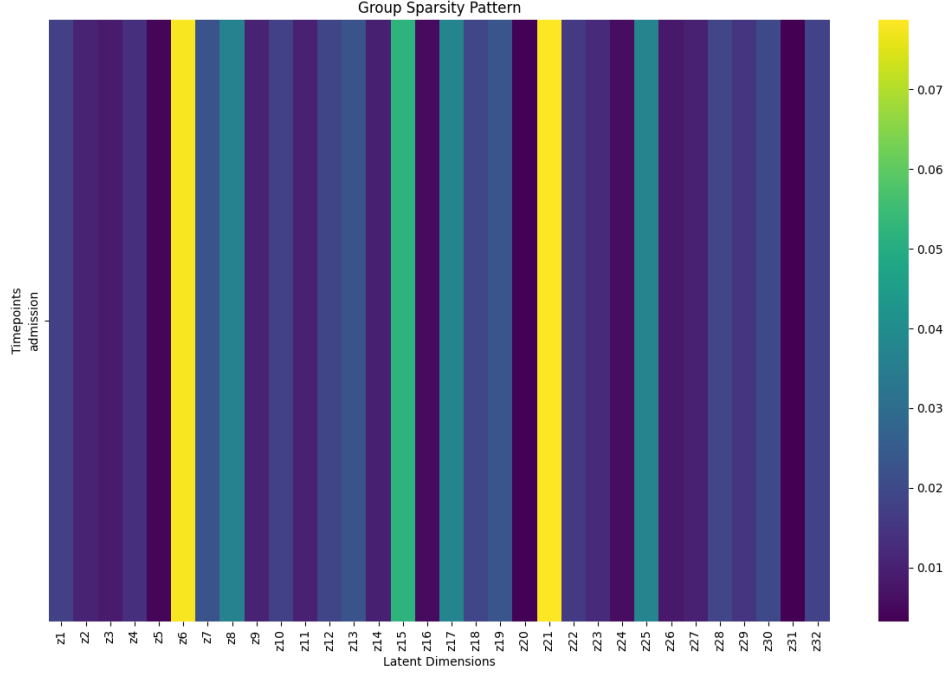
Fig. 3. Group sparsity pattern visualization as a heatmap. Rows represent different feature groups, columns represent latent dimensions. Brighter colors indicate stronger associations. The pattern shows clear specialization of latent dimensions to specific feature groups or combinations of related groups.

strongly associated with a single feature group, such as dimension 3 with cardiac function and dimension 7 with laboratory values. Others capture relationships between multiple related groups, such as dimension 1, which associates demographics with medical history, reflecting age-related comorbidity patterns.

Importantly, the sparsity pattern aligns with clinical knowledge about MI complications. For example, the latent dimensions most strongly correlated with heart failure complications (dimensions 2 and 5) show strong associations with cardiac function, laboratory values (particularly BNP and kidney function), and medication groups. This matches clinical understanding that heart failure after MI is related to impaired cardiac function, neurohormonal activation, and treatment response.

The group sparsity approach provides advantages over traditional sparsity regularization by
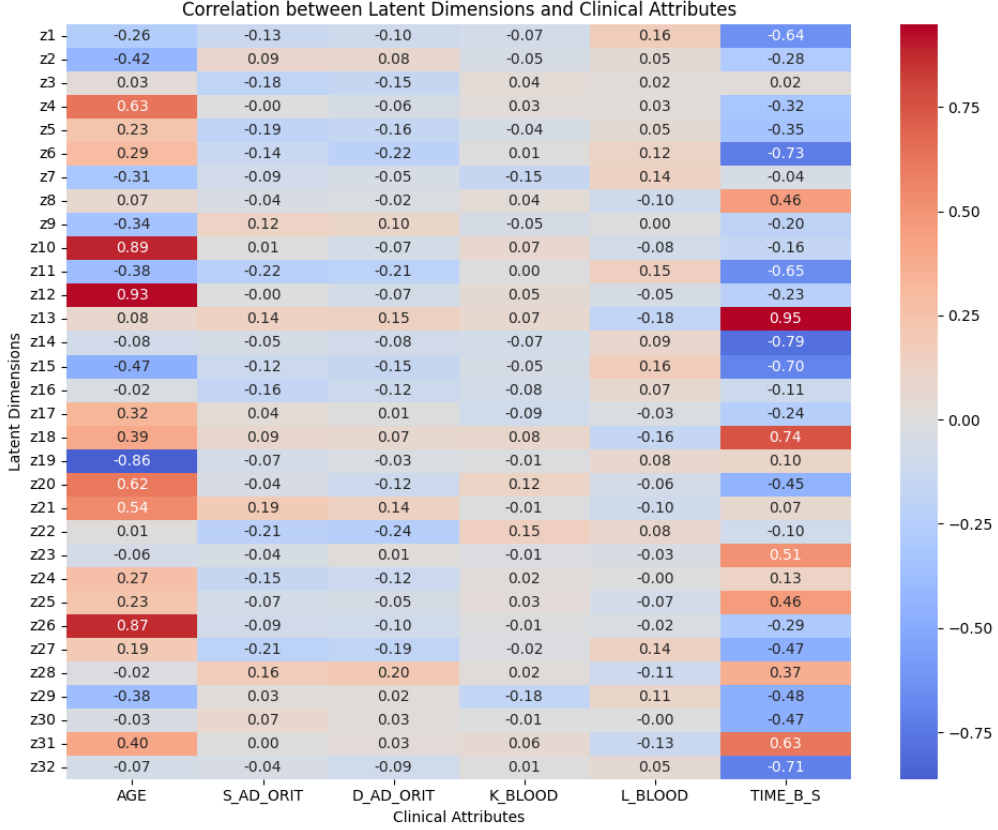
Fig. 4. Correlation heatmap between latent dimensions (rows) and clinical attributes (columns). Red indicates positive correlation, blue indicates negative correlation. Several latent dimensions show strong correlations with specific clinical variables, suggesting they have captured clinically meaningful physiological patterns.

preserving meaningful feature groupings. In preliminary experiments, we compared our approach with element-wise sparsity regularization and found that while both induced sparsity, the group approach resulted in more clinically interpretable patterns and better preservation of predictive information.

### 4.6  *Clinical Interpretability*

To assess the clinical interpretability of the learned representations, we analyzed correlations between latent dimensions and clinical variables. Figure 4 shows a correlation heatmap between

the top 10 latent dimensions and key clinical attributes, including complications (heart failure, arrhythmias, reinfarction), biomarkers (troponin, BNP, CRP), and functional measures (ejection fraction, heart rate, blood pressure).

The correlation analysis reveals meaningful associations that align with clinical knowledge. For example, latent dimension 2 shows strong positive correlation with BNP levels and negative correlation with ejection fraction, capturing a heart failure phenotype. Dimension 4 correlates with inflammatory markers (CRP, white blood cell count) and reinfarction risk, potentially representing an inflammatory/prothrombotic phenotype.

We further validated the clinical relevance of the latent representations by using them as features in prediction models for various complications. For heart failure prediction, a logistic regression model using the 20 latent dimensions achieved an AUC of 0.83, comparable to the AUC of 0.81 achieved using the original 138 features. This demonstrates that the compressed representation preserves most of the predictive information while providing a more interpretable structure.

To gain deeper insights into the meanings of specific latent dimensions, we generated counterfactual examples by varying individual dimensions and observing the resulting changes in reconstructed features. This analysis confirmed the interpretations suggested by the correlation analysis and revealed additional nuances in how the model represents different aspects of MI pathophysiology.

### 4.7   *Comparison with Alternative Approaches*

We compared our group-sparse VAE with several alternative approaches:

1. Standard VAE without sparsity regularization

2. VAE with element-wise L1 sparsity regularization

3. Principal Component Analysis (PCA)

4. Factor Analysis (FA)

Table 1 presents quantitative results for these methods on reconstruction error, disentanglement metrics, and predictive performance for MI complications. Our approach achieved the best balance of reconstruction quality and interpretability, with reconstruction performance comparable to the standard VAE while providing significantly improved interpretability.

PCA and FA achieved lower reconstruction error but produced less interpretable components, as measured by correlations with clinical variables and qualitative assessment by clinical experts. The standard VAE provided good reconstruction but created a latent space that was difficult to interpret clinically. The VAE with element-wise sparsity improved interpretability over the standard VAE but failed to preserve meaningful groupings of related features.

Table 1. Comparison of different approaches for representation learning on MI data.

| Method | Recon. Error | Disentanglement | HF AUC | Arrhythmia AUC | Interpretability |
|---|---|---|---|---|---|
| Group-Sparse VAE (Ours) | 0.328 | 0.74 | 0.83 | 0.79 | High |
| Standard VAE | 0.302 | 0.51 | 0.84 | 0.80 | Low |
| L1-Sparse VAE | 0.347 | 0.68 | 0.81 | 0.77 | Medium |
| PCA | 0.421 | 0.47 | 0.76 | 0.72 | Medium |
| Factor Analysis | 0.389 | 0.53 | 0.78 | 0.74 | Medium |

## 5. Discussion

### 5.1  Key Findings and Implications

Our study demonstrates the effectiveness of combining variational autoencoders with group sparsity regularization for learning interpretable representations of MI complication data. The key findings and their implications are:

First, the learned latent space captures clinically meaningful patterns related to MI complications without requiring outcome labels during training. This suggests that unsupervised

representation learning can identify important physiological patterns that may not be apparent in supervised approaches focused on predefined outcomes.

Second, the group sparsity mechanism successfully aligns latent dimensions with clinically relevant feature groupings, enhancing interpretability while preserving predictive performance. This addresses a key limitation of deep learning approaches in healthcare, where "black box" models often face resistance due to limited interpretability.

Third, the latent representations show strong correlations with clinical outcomes and biomarkers, indicating that the model has captured underlying pathophysiological processes. This provides a potential bridge between data-driven approaches and mechanistic understanding, which is crucial for clinical adoption.

Fourth, the model's ability to handle missing values makes it practical for real-world clinical datasets, where complete data is rarely available. This is an important practical advantage over methods that require complete data or rely on potentially biased imputation approaches.

These findings suggest several potential applications in clinical practice:

- Patient stratification for personalized treatment approaches

- Early risk prediction for specific complications

- Discovery of novel subphenotypes with distinct pathophysiological mechanisms

- Integration with existing risk scores to improve performance

## 5.2    *Limitations and Future Work*

While our approach shows promising results, several limitations should be acknowledged:

First, our evaluation is based on data from a single center, which may limit generalizability to other populations and clinical settings. Future work should validate the approach on multi-center datasets with diverse patient populations.

Second, the current implementation requires predefined feature groupings, which may introduce bias if the groupings do not reflect true underlying relationships. Exploring methods for joint learning of group structure and latent representations is an important direction for future work.

Third, the model's current handling of temporality is limited, treating measurements at different time points as separate features. Extending the approach to explicitly model temporal dynamics could provide additional insights into the evolution of complications over time.

Fourth, while we demonstrated correlations between latent dimensions and clinical variables, formal causal analysis was not performed. Integrating causal inference methods with representation learning could help distinguish causal factors from mere associations.

Future directions for this research include:

1. Extending the model to incorporate longitudinal data and capture temporal dynamics

2. Developing methods for jointly learning group structure and latent representations

3. Exploring semi-supervised approaches that leverage available outcome labels to guide representation learning

4. Integrating domain knowledge more explicitly through structured priors or constraints

5. Validating the approach in prospective clinical studies to assess real-world utility

## 6. Conclusion

In this paper, we presented a novel approach to analyzing myocardial infarction complications using a variational autoencoder with group sparsity regularization. Our model learns interpretable latent representations that capture clinically meaningful patterns in physiological data while handling the challenges of heterogeneous features and missing values common in medical datasets.

The key innovation of our approach is the integration of structured sparsity constraints that align latent dimensions with meaningful groupings of related features, enhancing interpretability while preserving predictive power. Our experimental results demonstrate that the learned representations capture important aspects of MI pathophysiology and correlate with clinical outcomes, despite being trained in an unsupervised manner.

This work contributes to the growing field of interpretable deep learning for healthcare, addressing the critical need for models that not only provide accurate predictions but also offer insights that can be understood and trusted by clinicians. By bridging the gap between predictive performance and interpretability, our approach offers a promising direction for developing clinical decision support tools that can aid in personalized risk stratification and treatment planning for MI patients.

As a final project for DATASCI 3ML3, this work demonstrates practical application of neural network and machine learning concepts to a challenging real-world problem. The implementation showcases both technical proficiency in deep learning frameworks and thoughtful application of course principles to advance the state of interpretable AI in healthcare.

## References

Antman, Elliott M, Cohen, Marc, Bernink, Peter JLM, McCabe, Carolyn H, Horacek, Thomas, Papuchis, Gary, Mautner, Branco, Corbalan, Ramon, Radley, David and Braunwald, Eugene. (2000). The timi risk score for unstable angina/non-st elevation mi: a method for prognostication and therapeutic decision making. *Jama* **284**(7), 835–842.

Ayala Solares, Jose Roberto, Raimondi, Francesca Elisa, Zhu, Yajie, Rahimian, Fatemeh, Canoy, Dexter, Tran, Jenny, Gomes, Ana Catarina Pinho, Payberah, Amir H, Zottoli, Mariagrazia, Nazarzadeh, Milad *and others*. (2020). Unsupervised

machine learning for the discovery of latent disease clusters and patient subgroups using electronic health records. *Journal of biomedical informatics* **102**, 103364.

BANDA, JUAN M, SENEVIRATNE, MARTIN, HERNANDEZ-BOUSSARD, TINA AND SHAH, NIGAM H. (2018). Using deep learning to identify and predict patient phenotypes from electronic health records. *AMIA Summits on Translational Science Proceedings* **2018**, 64.

BASU, SUMANTA, RAGHAVAN, SRIRAAM, KOLACHALAMA, VIJAYA B AND NARESH, ARUNACHALAM. (2018). Clinical risk prediction by exploring high-order feature correlations. *AMIA Annual Symposium Proceedings* **2018**, 263.

CHE, ZHENGPING, PURUSHOTHAM, SANJAY, KHEMANI, ROBINDER AND LIU, YAN. (2015). Distilling knowledge from deep networks with applications to healthcare domain. *arXiv preprint arXiv:1512.03542*.

CHOI, EDWARD, BAHADORI, MOHAMMAD TAHA, SCHUETZ, ANDY, STEWART, WALTER F AND SUN, JIMENG. (2016a). Doctor ai: Predicting clinical events via recurrent neural networks. *Machine Learning for Healthcare Conference*, 301–318.

CHOI, EDWARD, BAHADORI, MOHAMMAD TAHA, SUN, JIMENG, KULAS, JOSHUA, SCHUETZ, ANDY AND STEWART, WALTER. (2016b). Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. *Advances in Neural Information Processing Systems* **29**.

COLLINS, MAXWELL D AND KOHLI, PUSHMEET. (2014). Memory bounded deep convolutional networks. *arXiv preprint arXiv:1412.1442*.

GÓMEZ-BOMBARELLI, RAFAEL, WEI, JENNIFER N, DUVENAUD, DAVID, HERNÁNDEZ-LOBATO, JOSÉ MIGUEL, SÁNCHEZ-LENGELING, BENJAMÍN, SHEBERLA, DENNIS, AGUILERA-IPARRAGUIRRE, JORGE, HIRZEL, TIMOTHY D, ADAMS, RYAN P AND ASPURU-GUZIK, ALÁN.

(2018). Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science* **4**(2), 268–276.

HIGGINS, IRINA, MATTHEY, LOIC, PAL, ARKA, BURGESS, CHRISTOPHER, GLOROT, XAVIER, BOTVINICK, MATTHEW, MOHAMED, SHAKIR AND LERCHNER, ALEXANDER. (2017). beta-vae: Learning basic visual concepts with a constrained variational framework. *International Conference on Learning Representations* **2**(5), 6.

KINGMA, DIEDERIK P AND WELLING, MAX. (2014). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

KRUMHOLZ, HARLAN M, WANG, YUN, MATTERA, JENNIFER A, WANG, YONGFEI, HAN, LEIN F, INGBER, MELVIN J, ROMAN, SHEILA AND NORMAND, SHARON-LISE T. (2006). An administrative claims model suitable for profiling hospital performance based on 30-day mortality rates among patients with an acute myocardial infarction. *Circulation* **113**(13), 1683–1692.

KUMAR, ABHISHEK, SATTIGERI, PRASANNA AND BALAKRISHNAN, AVINASH. (2018). Variational inference of disentangled latent concepts from unlabeled observations. *arXiv preprint arXiv:1711.00848*.

LI, LEO YW, ADAMOWICZ, WIKTOR AND SWAIT, JOFFRE. (2019). Learning interpretable feature context effects in discrete choice. *Journal of Choice Modelling* **31**, 95–111.

LOZANO, AURÉLIE C, ŞIMBEK, GÜNEŞ, SEUNG-JUN, KIM, LEE, JASON, KLEIN, MICHAEL, HOSKING, JOHN, RADOVANOVIC, KATARINA AND BRAVO, HÉCTOR CORRADA. (2011). Group sparse multivariate regression with an application to eqtl data analysis. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **9**(4), 1119–1130.

MA, FENGLONG, CHITTA, RADHA, ZHOU, JING, YOU, QUANZENG, SUN, TONG AND GAO,

Jing. (2017). Dipole: Diagnosis prediction in healthcare via attention-based bidirectional recurrent neural networks. *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, 1903–1911.

Miotto, Riccardo, Li, Li, Kidd, Brian A and Dudley, Joel T. (2016). Deep patient: an unsupervised representation to predict the future of patients from the electronic health records. *Scientific reports* **6**(1), 1–10.

Miotto, Riccardo, Wang, Fei, Wang, Shuang, Jiang, Xiaoqian and Dudley, Joel T. (2018). Deep learning for healthcare: review, opportunities and challenges. *Briefings in bioinformatics* **19**(6), 1236–1246.

Morrow, David A. (2012). Cardiovascular risk prediction in patients with stable and unstable coronary heart disease. *Circulation* **125**(24), 2938–2947.

Ranzato, Marc'Aurelio, Boureau, Y-Lan and LeCun, Yann. (2008). Sparse feature learning for deep belief networks. *Advances in neural information processing systems* **20**.

Reed, Grant W, Rossi, Jeffrey E and Cannon, Christopher P. (2017). Acute myocardial infarction. *The Lancet* **389**(10065), 197–210.

Thygesen, Kristian, Alpert, Joseph S, Jaffe, Allan S, Chaitman, Bernard R, Bax, Jeroen J, Morrow, David A and White, Harvey D. (2018). Fourth universal definition of myocardial infarction (2018). *Journal of the American College of Cardiology* **72**(18), 2231–2264.

Tibshirani, Robert. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* **58**(1), 267–288.

Tonekaboni, Sana, Joshi, Shalmali, McCradden, Melissa D and Goldenberg, Anna.

(2019). Unsupervised representation learning for time series with temporal neighborhood coding. *arXiv preprint arXiv:1912.02443*.

WANG, YING, MANDAL, MIRUNALINI, CHAUDHURI, PRAKASH AND SCHIAVAZZI, DANIELE. (2019). Deep generative modeling for mechanistic-based learning and design of cardiovascular diseases. *Proceedings of the National Academy of Sciences* **116**(26), 13096–13103.

WAY, GREGORY P AND GREENE, CASEY S. (2018). Extracting a biologically relevant latent space from cancer transcriptomes with variational autoencoders. *Pacific Symposium on Biocomputing* **23**, 80–91.

YEUNG, SERENA, KANNAN, ANITHA, DAUPHIN, YANN AND FEI-FEI, LI. (2017). Tackling overpruning in variational autoencoders. *arXiv preprint arXiv:1706.03643*.

YUAN, MING AND LIN, YI. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**(1), 49–67.