

**Marginalization is not Marginal:
No Bad VAE Local Minima when Learning
Optimal Sparse Representations**

David Wipf

Amazon Web Services

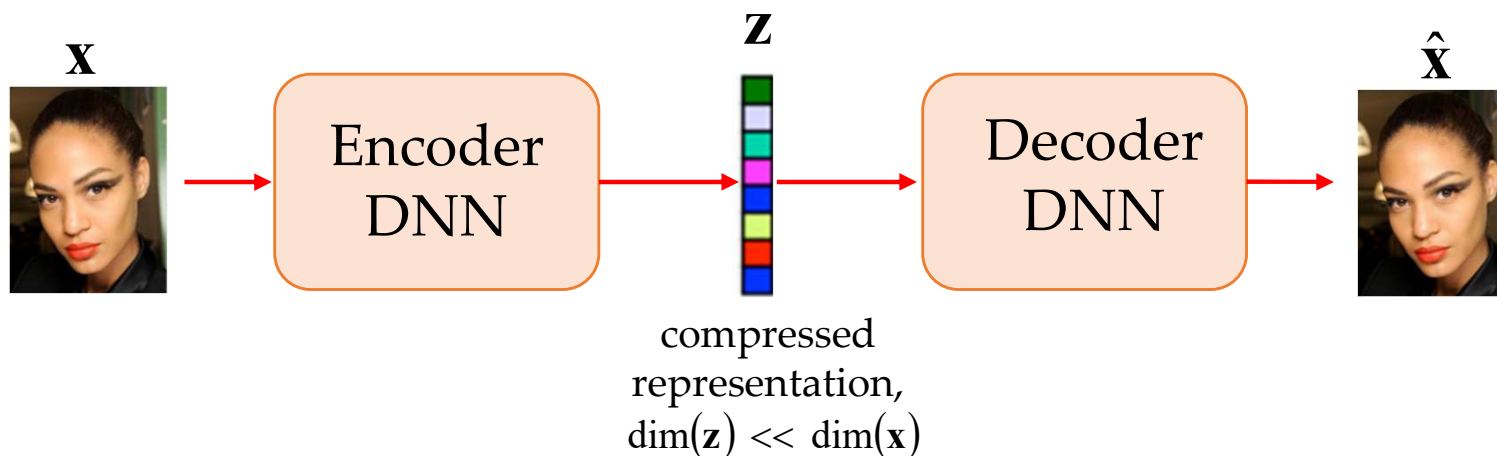
Variational Autoencoders (VAE)

- Popular generative model, probabilistic AE extension
- 33573+ citations [Kingma & Welling, 2014; Rezende et al., 2014]

Variational Autoencoders (VAE)

- Popular generative model, probabilistic AE extension
- 33573+ citations [Kingma & Welling, 2014; Rezende et al., 2014]

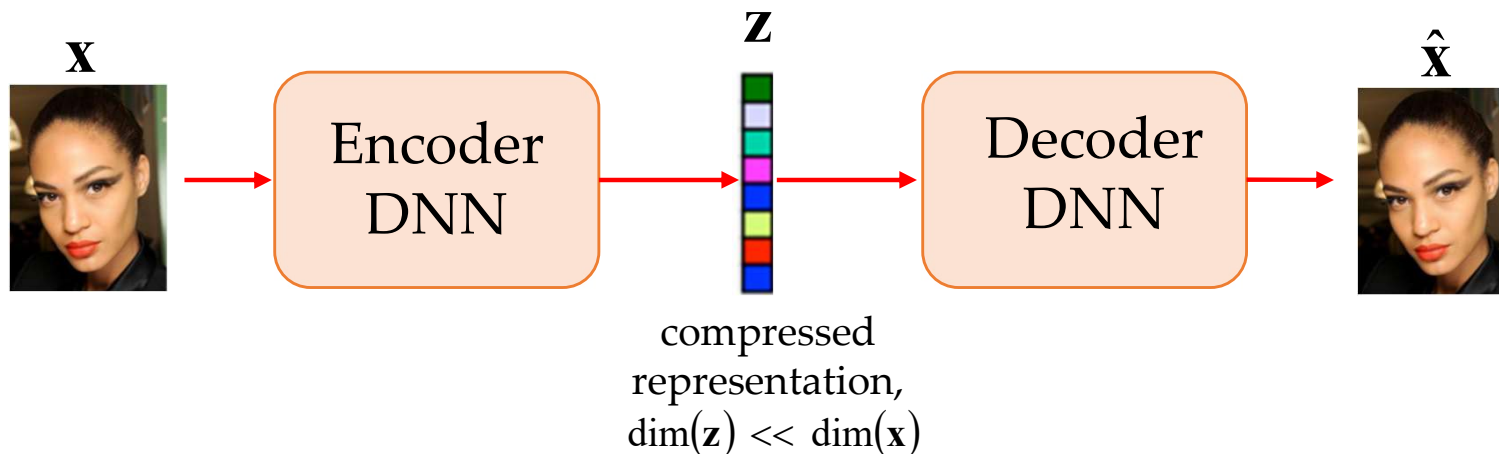
Autoencoder (AE):



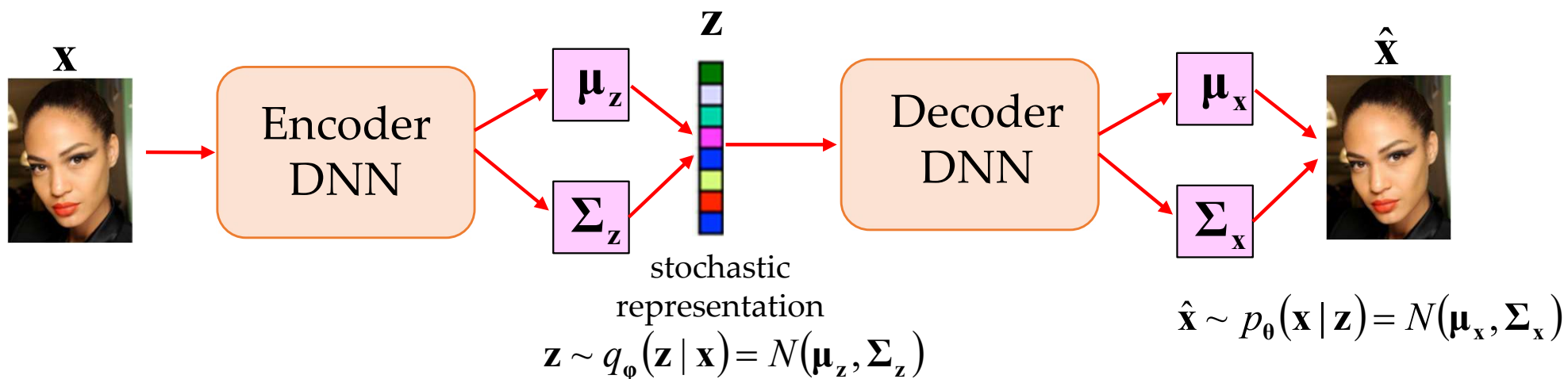
Variational Autoencoders (VAE)

- Popular generative model, probabilistic AE extension
- 33573+ citations [Kingma & Welling, 2014; Rezende et al., 2014]

Autoencoder (AE):

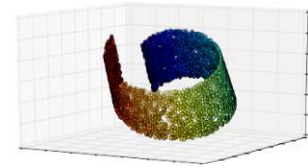


VAE (continuous data):



Why Variational Autoencoders?

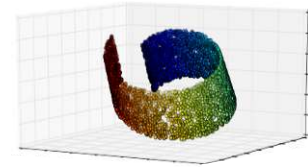
- Beyond generating samples, VAE **global minima** can be used to:
 - Compute the dimension of data manifolds.
 - Model/Exploit low-dimensional structure in data.



[Zheng et al., 2022]

Why Variational Autoencoders?

- Beyond generating samples, VAE **global minima** can be used to:
 - Compute the dimension of data manifolds.
 - Model/Exploit low-dimensional structure in data.



[Zheng et al., 2022]

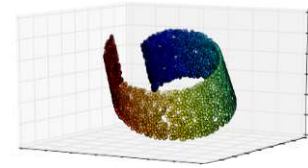
- **Simple representative example:**

With linear decoder, VAE minima learn PCA dimensions.

[Dai et al., 2019; Lucas et al., 2019]

Why Variational Autoencoders?

- Beyond generating samples, VAE **global minima** can be used to:
 - Compute the dimension of data manifolds.
 - Model/Exploit low-dimensional structure in data.



[Zheng et al., 2022]

- **Simple representative example:**

With linear decoder, VAE minima learn PCA dimensions.

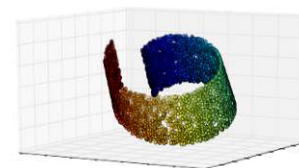
[Dai et al., 2019; Lucas et al., 2019]

- More broadly:
 - VAE global minima generalize many classical linear methods.
 - Useful for feature learning, interpretability, solving underdetermined inverse problems ...

[Dai et al., 2021]

Why Variational Autoencoders?

- Beyond generating samples, VAE **global minima** can be used to:
 - Compute the dimension of data manifolds.
 - Model/Exploit low-dimensional structure in data.



[Zheng et al., 2022]

- **Simple representative example:**

With linear decoder, VAE minima learn PCA dimensions.

[Dai et al., 2019; Lucas et al., 2019]

- More broadly:
 - VAE global minima generalize many classical linear methods.
 - Useful for feature learning, interpretability, solving underdetermined inverse problems ...
- [Dai et al., 2021]
- These VAE capabilities can be formalized through the notion of **optimal sparse representations**.

Optimal Sparse Representations

□ Two VAE criteria

1) Optimal reconstruction:
$$\sum_i \left\{ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[\left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] \right\|_2^2 \right] \right\} = 0$$

2) Maximal # of latent dimensions set to prior:
$$q_\phi(z_j | \mathbf{x}^{(i)}) = N(0,1), \quad \forall i$$
 no benefit to reconstructions

[Dai et al., 2021]

Optimal Sparse Representations

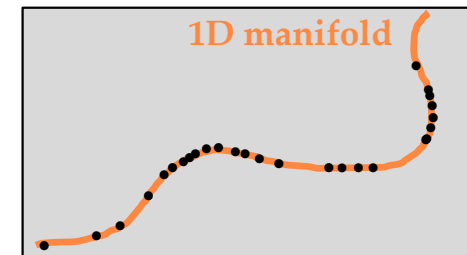
□ Two VAE criteria

1) Optimal reconstruction:
$$\sum_i \left\{ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[\left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] \right\|_2^2 \right] \right\} = 0$$

2) Maximal # of latent dimensions set to prior:
$$q_\phi(z_j | \mathbf{x}^{(i)}) = N(0,1), \quad \forall i$$
 no benefit to reconstructions

[Dai et al., 2021]

□ Note: Second criteria determines data manifold dimension.



● = observed data point

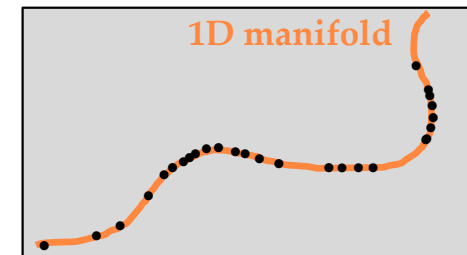
Optimal Sparse Representations

□ Two VAE criteria

1)	Optimal reconstruction:	$\sum_i \left\{ \mathbb{E}_{q_\phi(\mathbf{z} \mathbf{x}^{(i)})} \left[\left\ \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] \right\ _2^2 \right] \right\} = 0$	
2)	Maximal # of latent dimensions set to prior:	$q_\phi(z_j \mathbf{x}^{(i)}) = N(0,1), \quad \forall i$	no benefit to reconstructions

[Dai et al., 2021]

- Note: Second criteria determines data manifold dimension.



● = observed data point

- Under broad conditions, VAE **global minima** provably satisfy both criteria.

[Zheng et al., 2022]

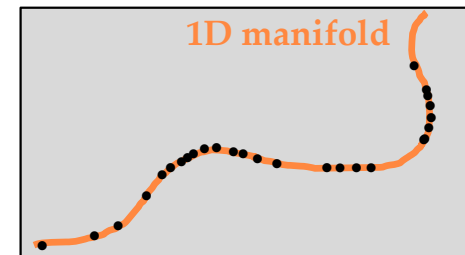
Optimal Sparse Representations

□ Two VAE criteria

1)	Optimal reconstruction:	$\sum_i \left\{ \mathbb{E}_{q_\phi(\mathbf{z} \mathbf{x}^{(i)})} \left[\left\ \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] \right\ _2^2 \right] \right\} = 0$	
2)	Maximal # of latent dimensions set to prior:	$q_\phi(z_j \mathbf{x}^{(i)}) = N(0,1), \quad \forall i$	no benefit to reconstructions

[Dai et al., 2021]

□ Note: Second criteria determines data manifold dimension.



● = observed data point

□ Under broad conditions, VAE **global minima** provably satisfy both criteria.

[Zheng et al., 2022]

□ Open question: What about bad VAE **local minima**?

Practical Example: Multiple-Response Sparse Regression

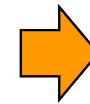
- Observed data: $\mathbf{X} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n \in \mathbb{R}^{d \times n}$

Practical Example: Multiple-Response Sparse Regression

□ Observed data: $\mathbf{X} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n \in \mathbb{R}^{d \times n}$

□ Assumed ground-truth model:

$$\mathbf{X} = \mathbf{\Phi} \mathbf{U}_0, \quad \underbrace{\mathbf{\Phi} \in \mathbb{R}^{d \times \kappa}}_{\text{known feature dictionary}}, \quad \underbrace{\mathbf{U}_0 \in \mathbb{R}^{\kappa \times n}}_{\text{unknown row-sparse coefficients}}$$



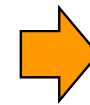
optimal sparse representation

Practical Example: Multiple-Response Sparse Regression

□ Observed data: $\mathbf{X} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n \in \mathbb{R}^{d \times n}$

□ Assumed ground-truth model:

$$\mathbf{X} = \mathbf{\Phi} \mathbf{U}_0, \quad \underbrace{\mathbf{\Phi} \in \mathbb{R}^{d \times \kappa}}_{\text{known feature dictionary}}, \quad \underbrace{\mathbf{U}_0 \in \mathbb{R}^{\kappa \times n}}_{\text{unknown row-sparse coefficients}}$$



optimal sparse representation

□ NP-hard inverse problem:

$$\mathbf{U}_0 = \arg \min_{\mathbf{U}} \rho_0(\mathbf{U}), \quad \text{s.t. } \mathbf{X} = \mathbf{\Phi} \mathbf{U}, \quad \rho_0(\mathbf{U}) \triangleq \sum_{j=1}^{\kappa} \mathbf{1} \left[\|\mathbf{u}_{(j)}\|_2 \neq 0 \right] \quad \left. \vphantom{\sum_{j=1}^{\kappa}} \right\} \text{counts \# of nonzero rows}$$

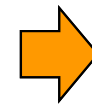
[Cotter et al., 2005; Lee et al., 2012; Adcock et al., 2019]

Practical Example: Multiple-Response Sparse Regression

□ Observed data: $\mathbf{X} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n \in \mathbb{R}^{d \times n}$

□ Assumed ground-truth model:

$$\mathbf{X} = \mathbf{\Phi} \mathbf{U}_0, \quad \underbrace{\mathbf{\Phi} \in \mathbb{R}^{d \times \kappa}}_{\text{known feature dictionary}}, \quad \underbrace{\mathbf{U}_0 \in \mathbb{R}^{\kappa \times n}}_{\text{unknown row-sparse coefficients}}$$



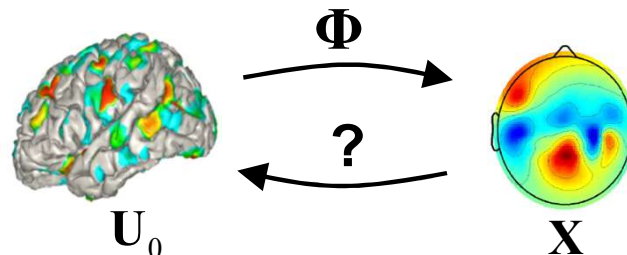
optimal sparse representation

□ NP-hard inverse problem:

$$\mathbf{U}_0 = \arg \min_{\mathbf{U}} \rho_0(\mathbf{U}), \quad \text{s.t. } \mathbf{X} = \mathbf{\Phi} \mathbf{U}, \quad \rho_0(\mathbf{U}) \triangleq \sum_{j=1}^{\kappa} \mathbf{1} \left[\|\mathbf{u}_{(j)}\|_2 \neq 0 \right] \quad \left. \vphantom{\sum} \right\} \text{counts \# of nonzero rows}$$

[Cotter et al., 2005; Lee et al., 2012; Adcock et al., 2019]

□ Applications:



[Banner et al., 2021; Bhutada et al., 2022; Cai et al., 2018]

Corresponding VAE Design

- VAE components for multiple sparse regression:

$$\text{Decoder: } \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z}, \quad \boldsymbol{\Sigma}_x[\mathbf{z}, \boldsymbol{\theta}] = \lambda \mathbf{I}, \quad \boldsymbol{\theta} = \{\mathbf{w}_x, \lambda\}$$

Corresponding VAE Design

- VAE components for multiple sparse regression:

$$\text{Decoder: } \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z}, \quad \boldsymbol{\Sigma}_x[\mathbf{z}, \boldsymbol{\theta}] = \lambda \mathbf{I}, \quad \boldsymbol{\theta} = \{\mathbf{w}_x, \lambda\}$$

$$\text{Encoder: } \boldsymbol{\mu}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{W}_z \mathbf{x}, \quad \boldsymbol{\Sigma}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{S} \mathbf{S}^T, \quad \boldsymbol{\varphi} = \{\mathbf{W}_z, \mathbf{S}\}$$

Corresponding VAE Design

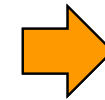
- VAE components for multiple sparse regression:

$$\text{Decoder: } \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z}, \quad \boldsymbol{\Sigma}_x[\mathbf{z}, \boldsymbol{\theta}] = \lambda \mathbf{I}, \quad \boldsymbol{\theta} = \{\mathbf{w}_x, \lambda\}$$

$$\text{Encoder: } \boldsymbol{\mu}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{W}_z \mathbf{x}, \quad \boldsymbol{\Sigma}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{S} \mathbf{S}^T, \quad \boldsymbol{\varphi} = \{\mathbf{W}_z, \mathbf{S}\}$$

- Resulting VAE energy function (i.e., -ELBO):

$$L_{VAE}(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \sum_{i=1}^n \left\{ \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\frac{1}{\lambda} \left\| \mathbf{x}^{(i)} - \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z} \right\|_2^2 \right] + d \log \lambda \right. \\ \left. + \text{tr}[\mathbf{S} \mathbf{S}^T] - \log |\mathbf{S} \mathbf{S}^T| + \left\| \mathbf{W}_z \mathbf{x}^{(i)} \right\|_2^2 \right\}$$



nonconvex,
potentially many
local minima

Corresponding VAE Design

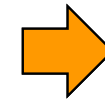
- VAE components for multiple sparse regression:

$$\text{Decoder: } \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z}, \quad \boldsymbol{\Sigma}_x[\mathbf{z}, \boldsymbol{\theta}] = \lambda \mathbf{I}, \quad \boldsymbol{\theta} = \{\mathbf{w}_x, \lambda\}$$

$$\text{Encoder: } \boldsymbol{\mu}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{W}_z \mathbf{x}, \quad \boldsymbol{\Sigma}_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{S} \mathbf{S}^T, \quad \boldsymbol{\varphi} = \{\mathbf{W}_z, \mathbf{S}\}$$

- Resulting VAE energy function (i.e., -ELBO):

$$L_{VAE}(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \sum_{i=1}^n \left\{ \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\frac{1}{\lambda} \left\| \mathbf{x}^{(i)} - \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z} \right\|_2^2 \right] + d \log \lambda \right. \\ \left. + \text{tr}[\mathbf{S} \mathbf{S}^T] - \log |\mathbf{S} \mathbf{S}^T| + \left\| \mathbf{W}_z \mathbf{x}^{(i)} \right\|_2^2 \right\}$$



nonconvex,
potentially many
local minima

- Analogous AE for multiple sparse regression:

$$L_{AE}(\mathbf{w}_x, \mathbf{W}_z) = \sum_{i=1}^n \frac{1}{\lambda} \left\| \mathbf{x}^{(i)} - \boldsymbol{\Phi} \text{diag}[\mathbf{w}_x] \mathbf{z}^{(i)} \right\|_2^2 + \underbrace{g(\|\mathbf{w}_x\|_2)}_{\text{for avoiding scaling ambiguity}} + \sum_{j=1}^{\kappa} \underbrace{h(\|\mathbf{z}_{(j)}\|_2)}_{\text{promotes row sparsity}}$$

s.t. $\mathbf{Z} = \mathbf{W}_z \mathbf{X} \in \mathbb{R}^{\kappa \times n}$

Main Result

Assume: $\mathbf{X} = \Phi \mathbf{U}_0 \in \mathbb{R}^{d \times n}$, $\rho_0(\mathbf{U}_0) < d$ (+ other minor tech. cond. on Φ)

Main Result

Assume: $\mathbf{X} = \Phi \mathbf{U}_0 \in \mathbb{R}^{d \times n}$, $\rho_0(\mathbf{U}_0) < d$ (+ other minor tech. cond. on Φ)

Denote: $\{\boldsymbol{\theta}^*, \boldsymbol{\varphi}^*\} \equiv \{\mathbf{w}_x^*, \mathbf{W}_z^*, \mathbf{S}^*\} =$ any local minimum of $L_{VAE}(\boldsymbol{\theta}, \boldsymbol{\varphi})$ as $\lambda \rightarrow 0$

Main Result

Assume: $\mathbf{X} = \Phi \mathbf{U}_0 \in \mathbb{R}^{d \times n}$, $\rho_0(\mathbf{U}_0) < d$ (+ other minor tech. cond. on Φ)

Denote: $\{\boldsymbol{\theta}^*, \boldsymbol{\varphi}^*\} \equiv \{\mathbf{w}_x^*, \mathbf{W}_z^*, \mathbf{S}^*\} =$ any local minimum of $L_{VAE}(\boldsymbol{\theta}, \boldsymbol{\varphi})$ as $\lambda \rightarrow 0$

Theorem (informal version)

If \mathbf{U}_0 has nonzero row norms with sufficiently different *scales*, then we have that:

- i) $\mathbf{U}_0 \Rightarrow$ unique optimal sparse representation
 - ii) $\{\mathbf{w}_x^*, \mathbf{W}_z^*, \mathbf{S}^*\}$ is a VAE global minimum
 - iii) $\text{diag}[\mathbf{w}_x^*] \mathbf{W}_z^* \mathbf{X} = \mathbf{U}_0$
 - iv) No analogous AE can satisfy equivalent recovery result
- bad local and/or global minima cannot be ruled out

Main Result

Assume: $\mathbf{X} = \Phi \mathbf{U}_0 \in \mathbb{R}^{d \times n}$, $\rho_0(\mathbf{U}_0) < d$ (+ other minor tech. cond. on Φ)

Denote: $\{\boldsymbol{\theta}^*, \boldsymbol{\varphi}^*\} \equiv \{\mathbf{w}_x^*, \mathbf{W}_z^*, \mathbf{S}^*\} =$ any local minimum of $L_{VAE}(\boldsymbol{\theta}, \boldsymbol{\varphi})$ as $\lambda \rightarrow 0$

Theorem (informal version)

If \mathbf{U}_0 has nonzero row norms with sufficiently different *scales*, then we have that:

- i) $\mathbf{U}_0 \Rightarrow$ unique optimal sparse representation
 - ii) $\{\mathbf{w}_x^*, \mathbf{W}_z^*, \mathbf{S}^*\}$ is a VAE global minimum
 - iii) $\text{diag}[\mathbf{w}_x^*] \mathbf{W}_z^* \mathbf{X} = \mathbf{U}_0$
 - iv) No analogous AE can satisfy equivalent recovery result
- bad local and/or global minima cannot be ruled out

How is this possible? \Rightarrow **Marginalization** over VAE latent space introduces selective smoothing effect

Empirical Examples

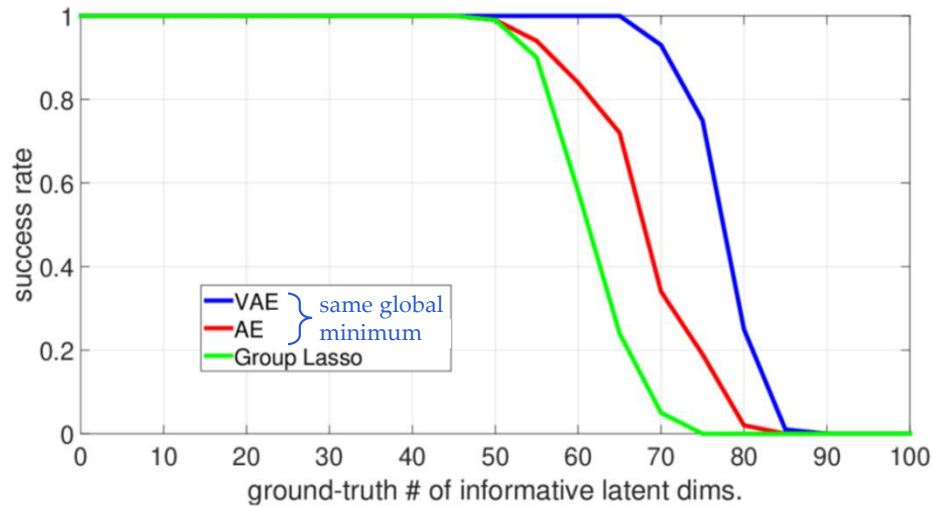
Goal: Recover ground-truth \mathbf{U}_0 from observations $\mathbf{X} = \Phi\mathbf{U}_0$

Empirical Examples

Goal: Recover ground-truth \mathbf{U}_0 from observations $\mathbf{X} = \Phi\mathbf{U}_0$

poor approx to theory

$\Phi \sim N(0,1)$, \mathbf{U}_0 has unit-norm nonzero rows

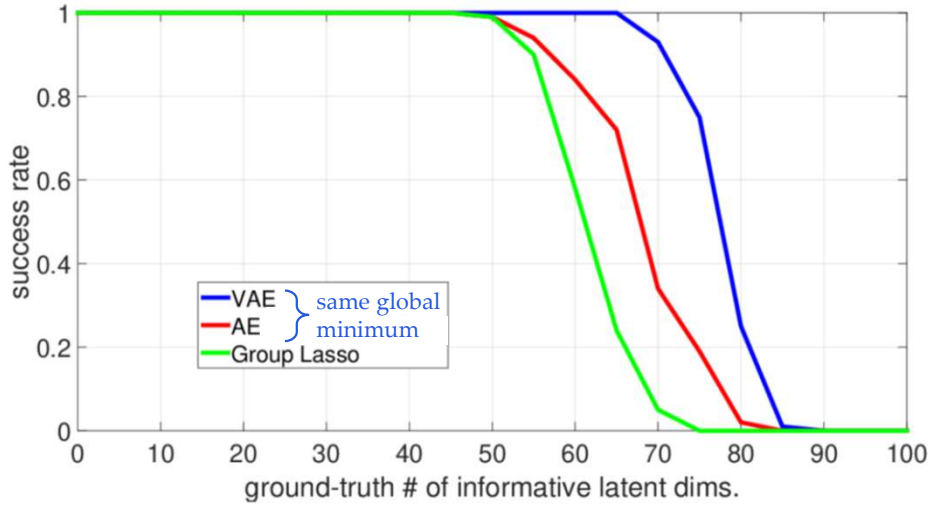


Empirical Examples

Goal: Recover ground-truth \mathbf{U}_0 from observations $\mathbf{X} = \Phi\mathbf{U}_0$

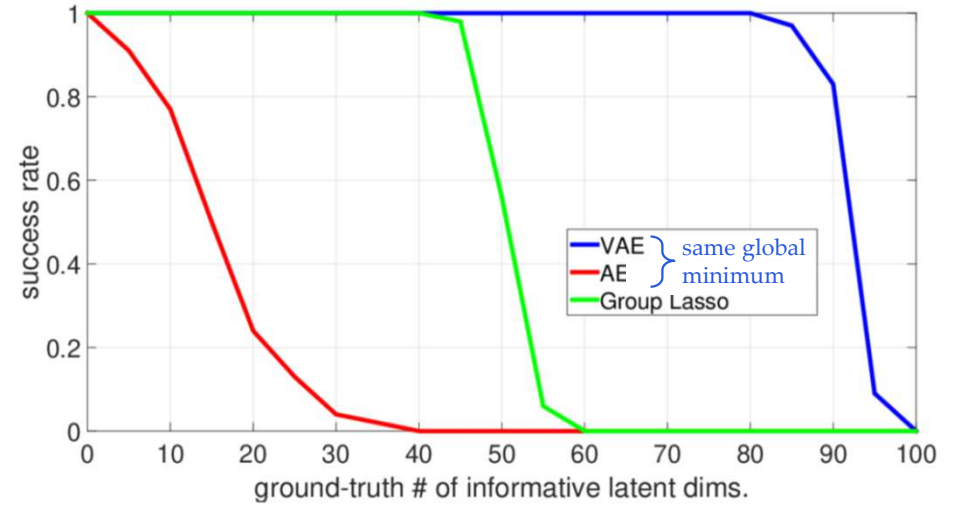
poor approx to theory

$\Phi \sim N(0,1)$, \mathbf{U}_0 has unit-norm nonzero rows



close approx to theory

$\Phi \sim N(0,1)$, \mathbf{U}_0 has scaled nonzero rows



Empirical Examples

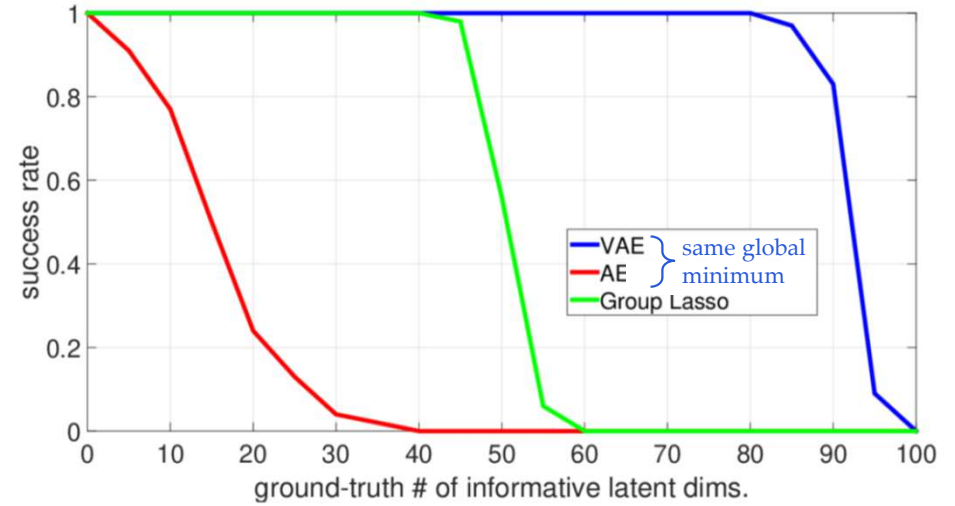
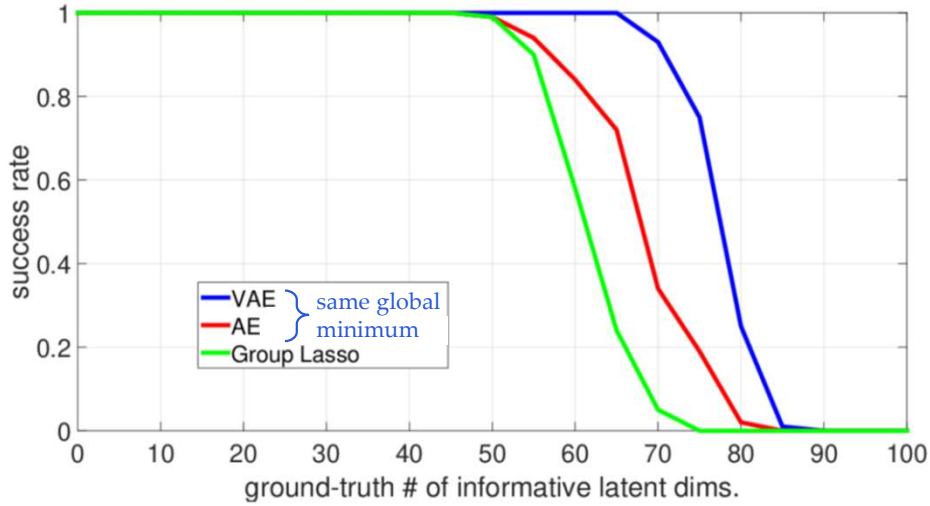
Goal: Recover ground-truth \mathbf{U}_0 from observations $\mathbf{X} = \Phi\mathbf{U}_0$

poor approx to theory

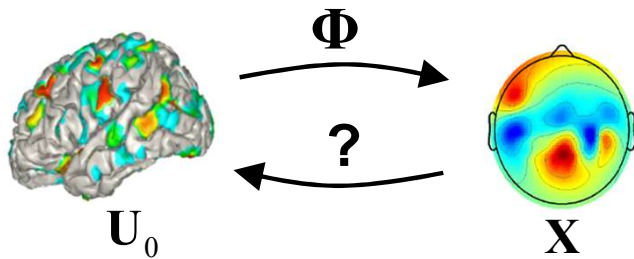
close approx to theory

$\Phi \sim N(0,1)$, \mathbf{U}_0 has unit-norm nonzero rows

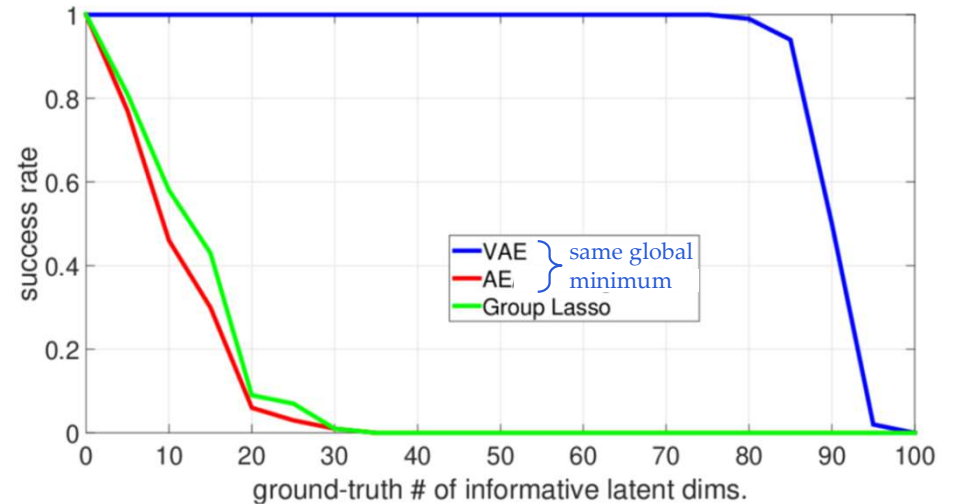
$\Phi \sim N(0,1)$, \mathbf{U}_0 has scaled nonzero rows



Neuromagnetic inverse modeling



Φ = leadfield matrix, highly correlated columns



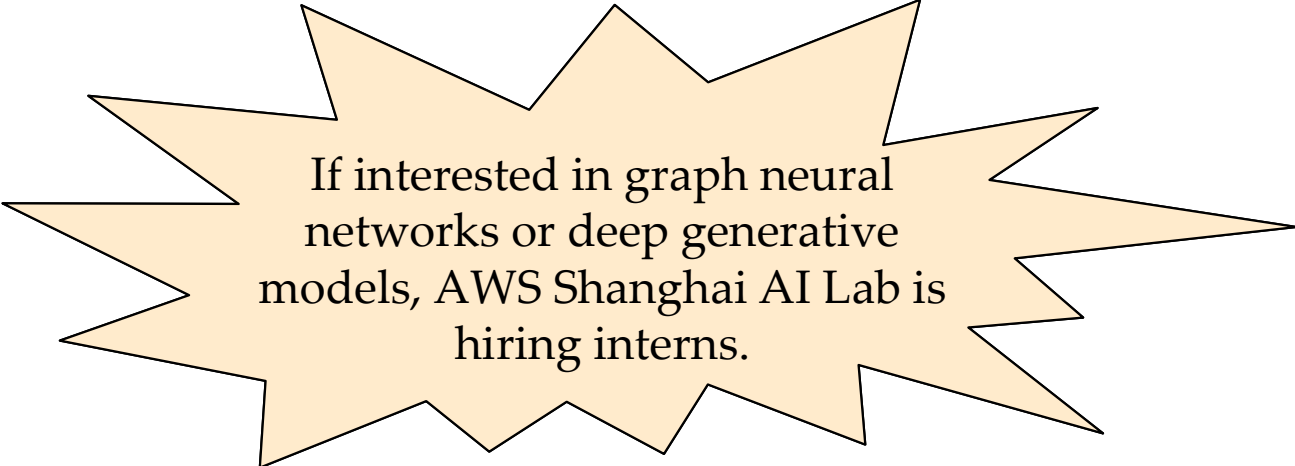
Summary

- ❑ Prior work has shown that the VAE **global** minima can provably recover low-dimensional structure in data.
- ❑ But previously no strict guarantees (outside of cases that reduce to PCA) regarding bad **local** minima.
- ❑ We demonstrate a challenging regime whereby **all** VAE local minima produce optimal sparse representations.
- ❑ Made possible by VAE **marginalization**.
- ❑ Practical relevance:
 - ❑ Helps to explain the effectiveness of VAEs in modeling low-dimensional structure in data.
 - ❑ Motivates diverse VAE use cases beyond generating samples.
 - ❑ Theory makes accurate predictions regarding empirical VAE behavior in broader regimes of interest, e.g., more complex decoders.

Thank You

Links:

- ❑ <http://www.davidwipf.com/>
- ❑ <http://www.dgl.ai/>



If interested in graph neural networks or deep generative models, AWS Shanghai AI Lab is hiring interns.