

---

# Deep Generative Models for Signal Processing and Beyond

---

David Wipf  
Microsoft Research

**Note:** Updated version of slides available at <http://www.davidwipf.com/>

# Part I: Introduction

# Generative Models

- Given: Training data

$$\{\mathbf{x}^{(i)}\}_{i=1}^n, \quad \mathbf{x}^{(i)} \sim p_{gt}(\mathbf{x})$$

Example: MNIST digits



- Goal: Learn a parametric model capable of producing **new** samples

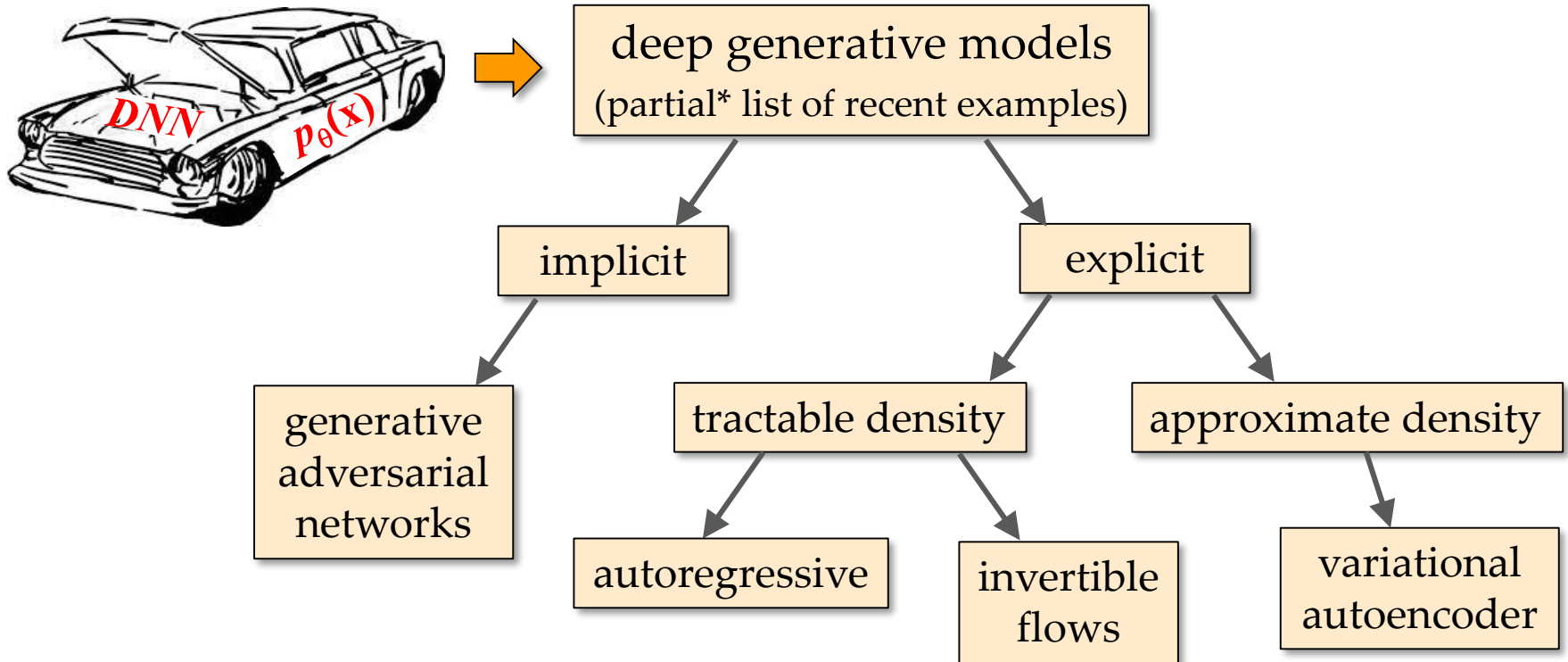
$$\{\mathbf{x}^{(j)}\}_{j=1}^m, \quad \mathbf{x}^{(j)} \sim p_{\theta}(\mathbf{x}) \approx p_{gt}(\mathbf{x})$$

Similar to real data



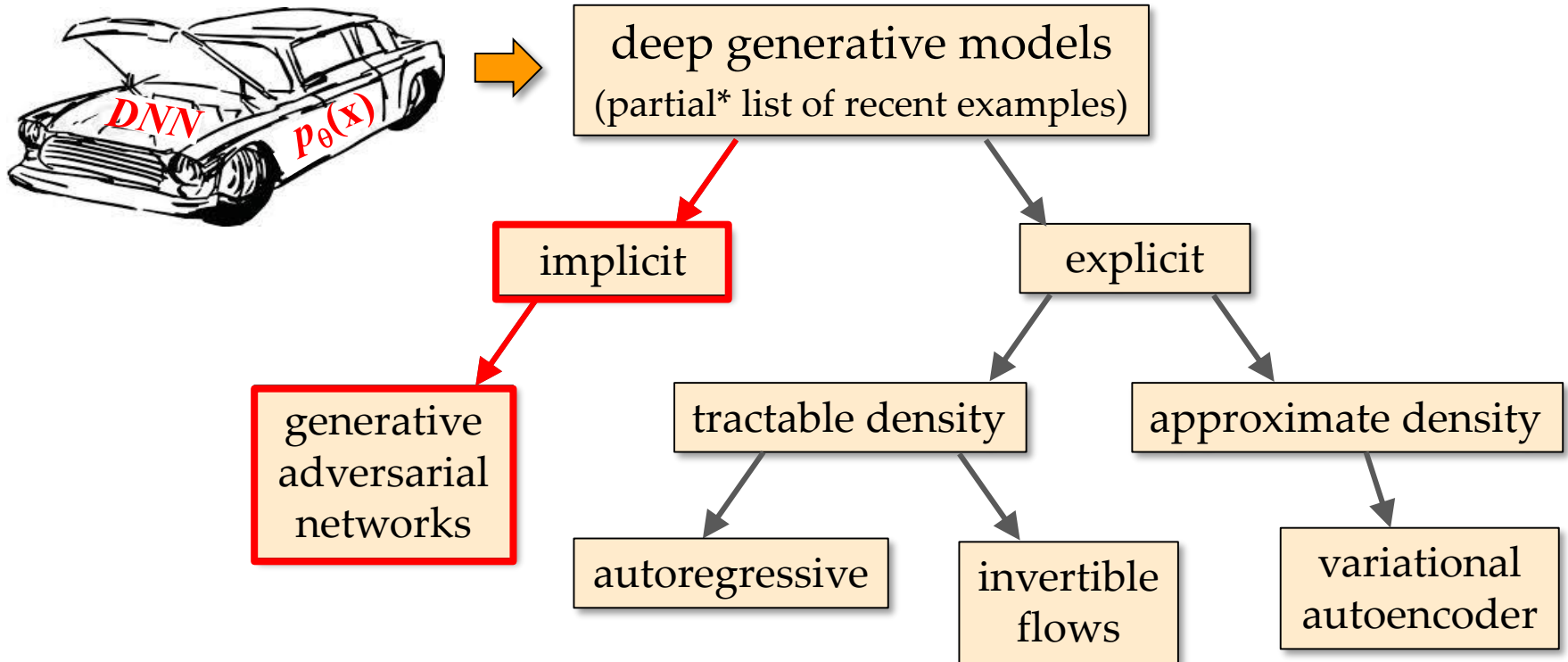
- **Deep generative models** use neural networks for implicitly or explicitly defining the density  $p_{\theta}(\mathbf{x})$

# Types of Deep Generative Models



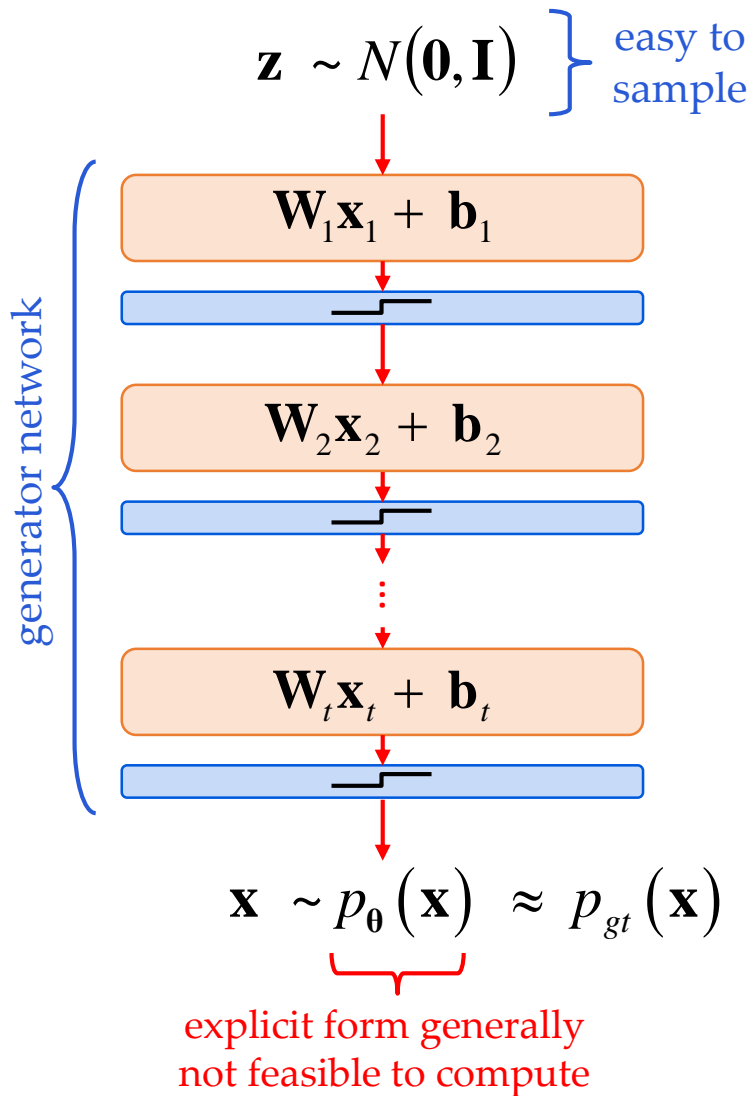
\*For additional examples, see tutorial [Goodfellow, 2016]

# Implicit Deep Generative Models

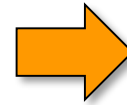
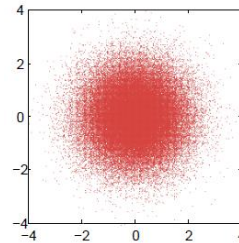


\*For additional examples, see tutorial [Goodfellow, 2016]

# Implicit Deep Generative Modeling



## Example



## Popular Example

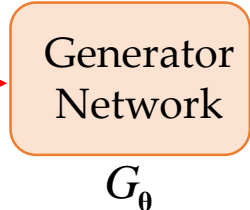
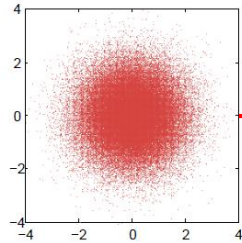
### Generative Adversarial Networks:

- Based on game theory, Nash equilibrium
- [8679 citations](#) [Goodfellow et al., 2014]

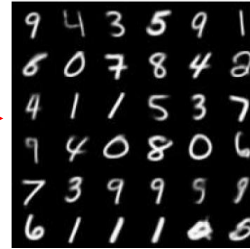
# Generative Adversarial Networks (GANs)

noise samples

$$\{\mathbf{z}^{(j)}\}_{j=1}^m, \mathbf{z}^{(j)} \sim N(\mathbf{z} | \mathbf{0}, \mathbf{I})$$



fake samples



real samples

$$\{\mathbf{x}^{(i)}\}_{i=1}^n, \mathbf{x}^{(i)} \in \mathbb{R}^d \sim p_{gt}(\mathbf{x})$$



real/fake

$$D_\phi : \mathbb{R}^d \rightarrow (0,1)$$

Binary classification:

$$D_\phi(\underbrace{\mathbf{x}^{(i)}}_{\text{real}}) \approx 1, \quad D_\phi\left[\underbrace{G_\theta(\mathbf{z}^{(j)})}_{\text{fake}}\right] \approx 0$$

**Basic GAN objective (cross-entropy-based):**

$$\min_{\theta} \max_{\phi} \mathbb{E}_{p_{gt}(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{N(\mathbf{z}|\mathbf{0},\mathbf{I})} [\log(1 - D_\phi[G_\theta(\mathbf{z})])]$$

expectations approximated with samples

# GAN Strengths

State-of-the-art GAN models generate highly realistic samples, e.g., StyleGAN [Karras et al, 2019]:



**real**



**fake**

Examples from <http://www.whichfaceisreal.com/>

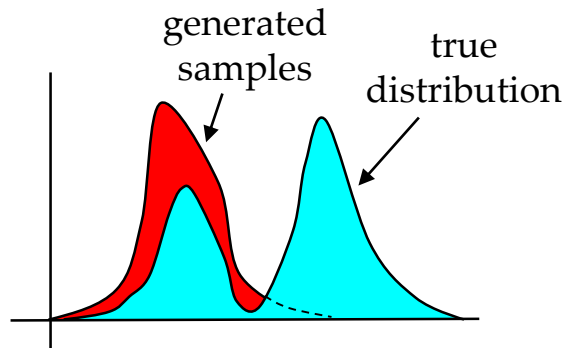


# GAN Weaknesses

- Training involves potentially unstable minimax problem, iterations may diverge, be sensitive to tuning.

[Lucic et al., 2018]

- Can be susceptible to mode collapse:



low sample diversity



[Arora and Zhang, 2017]

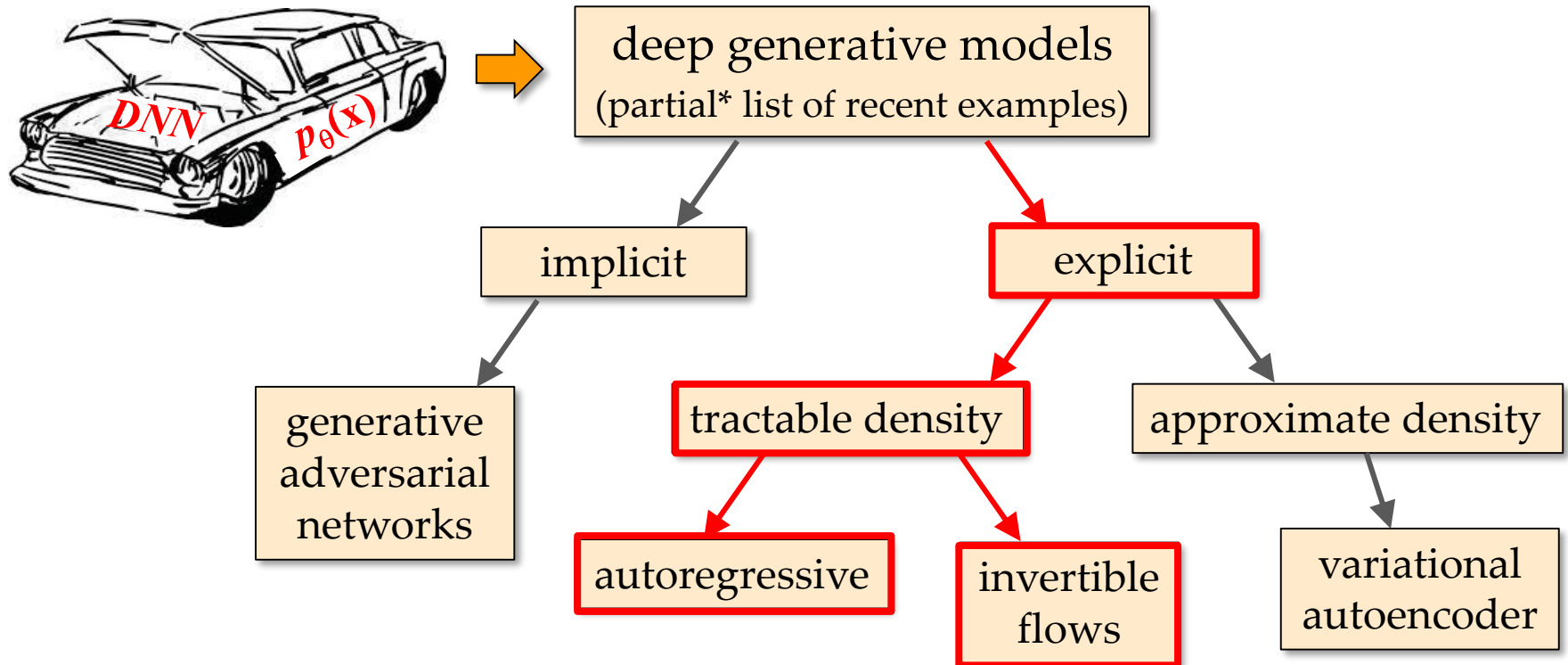
- No explicit density estimate  $p_{\theta}(\mathbf{x}) \approx p_{gt}(\mathbf{x})$ , cannot infer the latent code that produced a sample:



$$p_{\theta}(\mathbf{z} | \mathbf{x}) ?$$

cannot compute low-dimensional representation

# Explicit Deep Generative Modeling w/ a Tractable Density



\*For additional examples, see tutorial [Goodfellow, 2016]

# Explicit Deep Generative Modeling w/ a Tractable Density

- Density  $p_{\theta}(\mathbf{x})$  and gradients  $\nabla p_{\theta}(\mathbf{x})$  can be computed exactly
- Given training data  $\{\mathbf{x}^{(i)}\}_{i=1}^n$ ,  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ , can solve via SGD:

$$\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} - \sum_i \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \quad \Rightarrow \quad \text{maximum likelihood estimator}$$

- **Key advantage:** Closed-form test data likelihood  $p_{\boldsymbol{\theta}_*}(\mathbf{x}^{test})$
- **Disadvantages:**
  - Generated samples arguably inferior to GANs
  - No dimensionality reduction, representation learning (mostly)

# Examples

## □ Autoregressive methods:

Apply chain rule to form:

$$p_{\theta}(\mathbf{x}) = \prod_{j=1}^d p_{\theta}(x_j | x_1, \dots, x_{j-1}) \quad \left. \vphantom{\prod_{j=1}^d} \right\} \begin{array}{l} \text{conditionals parameterized} \\ \text{as RNN or CNN} \end{array}$$

[Larochelle & Murray, 2011; van den Oord et al., 2016]

## □ Invertible flows:

Assumptions:

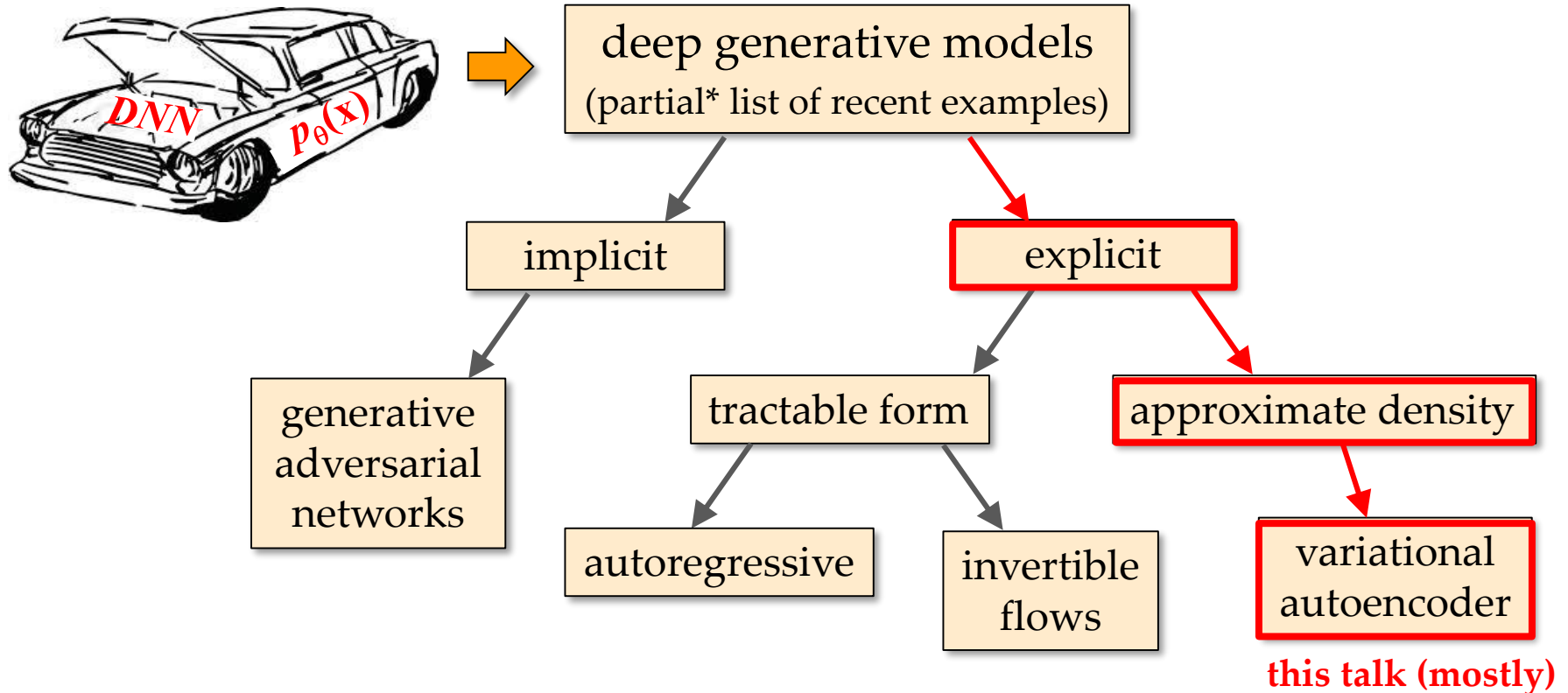
$$p(\mathbf{z}) = N(\mathbf{z} | \mathbf{0}, \mathbf{I}), \quad \dim(\mathbf{z}) = \dim(\mathbf{x}), \quad \mathbf{z} = f_{\theta}(\mathbf{x}), \quad \mathbf{x} = f_{\theta}^{-1}(\mathbf{z})$$

Change of variables formula:

$$p_{\theta}(\mathbf{x}) = N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \left| \det \left( \frac{\partial f_{\theta}(\mathbf{x})}{\partial \mathbf{x}^T} \right) \right| \quad \left. \vphantom{\det} \right\} \begin{array}{l} \text{tractable determinant because} \\ \text{of special DNN structure} \end{array}$$

[Dinh et al., 2016; Kingma & Dhariwal, 2018]

# Explicit Deep Generative Modeling Using a Density Approximation/Bound



\*For additional examples, see tutorial [Goodfellow, 2016]

# Explicit Deep Generative Modeling Using a Density Approximation/Bound

- Often interested in densities of the form:

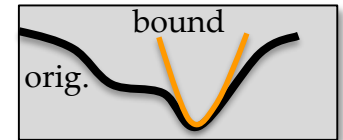
$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

low-dimensional latent factors

- Required integral is intractable ...



Optimize upper bound on  $-\sum_i \log p_{\theta}(\mathbf{x}^{(i)})$

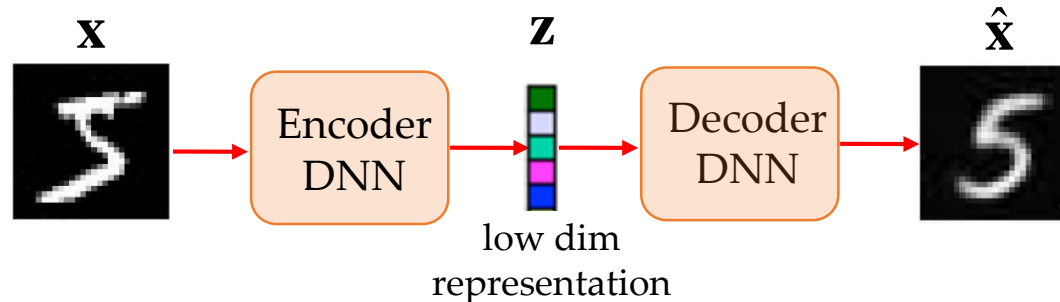


- Popular example: **The variational autoencoder (VAE)**

[Kingma and Welling, 2014; Rezende et al., 2014]

Upper bound based on autoencoder-like structure

6079 citations



# Variational Autoencoders

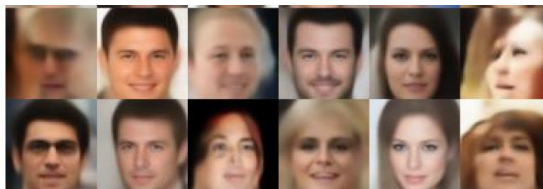
## (details in Part II)

### Advantages:

- ❑ Less prone to mode collapse than GANs, more stable training.
- ❑ Provides explicit estimate of latent distribution  $p_{\theta}(\mathbf{z} | \mathbf{x})$ ; many applications in representation learning.
- ❑ Natural generalization of dimensionality reduction tools in common use for signal processing (Part III).

### Disadvantages:

- ❑ Optimizes a bound on the data likelihood, not exact likelihood (but conditions for when bound is tight discussed in Part IV).
- ❑ Generated samples usually inferior to GANs ...

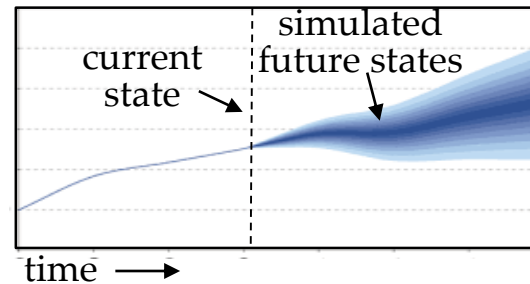


... although improvements possible (Part IV).

# Representative Applications

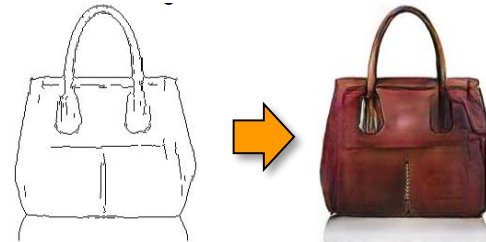
## Generative models in general:

- Model-based reinforcement learning:



[Finn et al., 2016]

- Image-to-image translation:

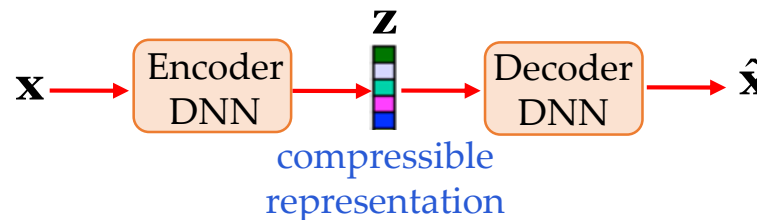


[Isola et al., 2016]

- Many more, a generic unsupervised learning tool

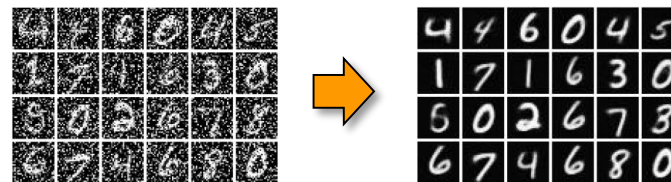
## VAEs in particular:

- Compression:



[Ballé et al., 2018]

- Data cleaning, outlier removal:



[Dai et al., 2018]



# Caveat



- ❑ Deep generative modeling is a rapidly changing field.
- ❑ Strengths and weaknesses of various methods frequently need recalibration in accordance with new developments.
- ❑ Also, important to differentiate:
  - 1) General-purpose improvements in DNN architectures
  - 2) Advances in specific generative modeling paradigms

# Remainder of Tutorial

- ❑ Part II: Details of the variational autoencoder
- ❑ Part III: Connections with existing signal processing methods for finding low-dimensional structure in data
- ❑ Part IV: From signal reconstruction to generative modeling
- ❑ Part V: Practical usage issues and examples

**Questions?**

## Part II: Details of the Variational Autoencoder

**Note:** Updated version of slides available at <http://www.davidwipf.com/>

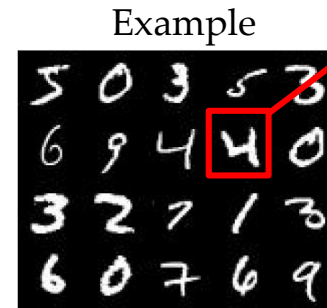
# Latent Variable Model

Observed data:  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n$ ,  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ ,  $\forall i$

Assumed latent factors:

$$\mathbf{Z} = \{\mathbf{z}^{(i)}\}_{i=1}^n, \quad \mathbf{z}^{(i)} \in \mathbb{R}^k, \quad \forall i, \quad k \ll d$$

low-dimensional  
representation of significant  
factors of variation



28x28 = 784 dim  
MNIST digit

candidate latent factors:  
digit type, stroke width,  
slant angle, etc.

$$k < 20 \ll d = 784$$

sufficient in practice

Ground-truth generative process:

$$\mathbf{z}^{(i)} \sim p_{gt}(\mathbf{z}), \quad \mathbf{x}^{(i)} \sim p_{gt}(\mathbf{x} | \mathbf{z}^{(i)})$$

prior on  
latent factors  
of variation

mapping to  
observation space;  
could be delta function

$$\mathbf{x}^{(i)} \sim p_{gt}(\mathbf{x}) = \int p_{gt}(\mathbf{x} | \mathbf{z}) p_{gt}(\mathbf{z}) d\mathbf{z}$$

# Parameterized Latent-Variable Model

Without loss of generality, assume:  $p_{gt}(\mathbf{z}) = N(\mathbf{z} | \mathbf{0}, \mathbf{I})$

Also assume parameterized family:

$$p_{\theta}(\mathbf{x} | \mathbf{z}), \quad \theta \in \Omega$$

$$\text{s.t. } p_{\theta_*}(\mathbf{x} | \mathbf{z}) \approx p_{gt}(\mathbf{x} | \mathbf{z}), \quad \text{for some } \theta_* \in \Omega$$

High-level goal:

$$\text{Given } \mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n, \quad \mathbf{x}^{(i)} \sim p_{gt}(\mathbf{x})$$

$$\text{Solve } \min_{\theta} \underbrace{-\sum_{i=1}^n \log p_{\theta}(\mathbf{x}^{(i)})}_{\text{equivalent to maximum likelihood}} \equiv \min_{\theta} -\sum_{i=1}^n \log \int p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}$$

equivalent to maximum likelihood

$$\text{Key problem: } p_{\theta}(\mathbf{x}^{(i)}) = \int p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}$$

$$p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)}) = p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) / p_{\theta}(\mathbf{x}^{(i)})$$



intractable

# Naïve Approximation

Finite-sample approximation to intractable integral for each  $i$ :

$$\text{sample } \mathbf{z}^{(i,j)} \sim N(\mathbf{z} | \mathbf{0}, \mathbf{I}), \quad j = 1, \dots, m$$

$$\Rightarrow \int p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z} \equiv \mathbb{E}_{N(\mathbf{z} | \mathbf{0}, \mathbf{I})} [p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z})] \approx \frac{1}{m} \sum_{j=1}^m p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,j)})$$

Revised tractable objective:

$$\min_{\theta} - \sum_{i=1}^n \log \left[ \frac{1}{m} \sum_{j=1}^m p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,j)}) \right]$$

Lingering problem:

$$\text{for most } \mathbf{z}^{(i,j)} \sim N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad \Rightarrow \quad p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,j)}) \approx 0$$

**Need huge number of samples for reasonable approximation ...**

# A Useful Variational Bound

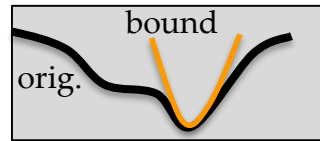
Define an approximate distribution as

$$q_{\varphi}(\mathbf{z} | \mathbf{x}^{(i)}) \approx p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)}) = p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) / p_{\theta}(\mathbf{x}^{(i)})$$

$\underbrace{\hspace{10em}}$ 
 $\underbrace{\hspace{10em}}$

tractable intractable

Variational upper bound:



$$-\sum_i \log p_{\theta}(\mathbf{x}^{(i)}) \leq L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \triangleq \sum_i \left\{ \underbrace{\text{KL}[q_{\varphi}(\mathbf{z} | \mathbf{x}^{(i)}) \| p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)})]}_{\geq 0} - \log p_{\theta}(\mathbf{x}^{(i)}) \right\}$$

After standard manipulations ...

$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \equiv \sum_i \left\{ \underbrace{\text{KL}[q_{\varphi}(\mathbf{z} | \mathbf{x}^{(i)}) \| N(\mathbf{z} | \mathbf{0}, \mathbf{I})] - \mathbb{E}_{q_{\varphi}(\mathbf{z} | \mathbf{x}^{(i)})}[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z})]}_{\geq 0} \right\}$$

Does not depend on intractable  $p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)})$  or  $p_{\theta}(\mathbf{x}^{(i)})$



# Basic VAE Energy Function Decomposition

$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \underbrace{\sum_i \left\{ \text{KL} \left[ q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right] \right\}}_{\text{regularization factor}} - \underbrace{\mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right]}_{\text{data-fit term}}$$

# Handling the Regularization Term

$\text{KL}\left[q_\phi(\mathbf{z}|\mathbf{x}^{(i)})\|N(\mathbf{z}|\mathbf{0},\mathbf{I})\right]$  is still intractable in general

Simplifying Gaussian approximate posterior assumption:

$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = N\left(\mathbf{z}|\boldsymbol{\mu}_z[\mathbf{x}^{(i)},\phi], \boldsymbol{\Sigma}_z[\mathbf{x}^{(i)},\phi]\right) \quad \left. \vphantom{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \right\} \begin{array}{l} \text{encoder} \\ \text{distribution} \end{array}$$

Encoder moments computed by deep networks:



KL term now satisfies:

$$2 \text{KL}\left[q_\phi(\mathbf{z}|\mathbf{x}^{(i)})\|N(\mathbf{z}|\mathbf{0},\mathbf{I})\right] \equiv \left\|\boldsymbol{\mu}_z[\mathbf{x}^{(i)},\phi]\right\|_2^2 + \text{tr}\left(\boldsymbol{\Sigma}_z[\mathbf{x}^{(i)},\phi]\right) - \log\left|\boldsymbol{\Sigma}_z[\mathbf{x}^{(i)},\phi]\right|$$

**Differentiable, suitable for minimization via SGD**

# Basic VAE Energy Function Decomposition

$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \sum_i \left\{ \begin{array}{l} \text{regularization factor} \\ \text{KL}[q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) \| N(\mathbf{z} | \mathbf{0}, \mathbf{I})] \\ \text{data-fit term} \\ - \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z})] \end{array} \right\}$$

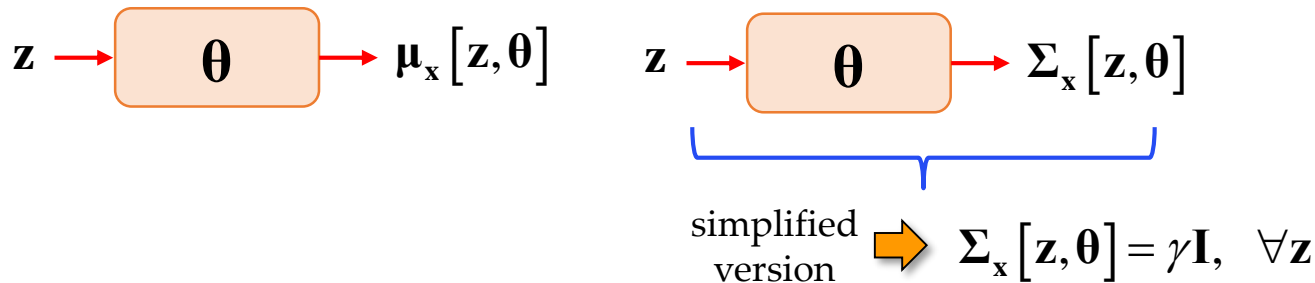
# Handling the Data-Fit Term

$-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right]$  is also generally intractable

For continuous data, typical assumption is

$$p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) = N(\mathbf{x}^{(i)} | \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}], \boldsymbol{\Sigma}_x[\mathbf{z}, \boldsymbol{\theta}]) \quad \left. \vphantom{p_\theta} \right\} \text{decoder distribution}$$

Decoder moments computed by deep networks:



But ...

$$-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right] \equiv \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \frac{1}{2\gamma} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}]\|_2^2 \right] \Rightarrow \text{still intractable}$$

# Revisiting Finite-Sample Approximations

From before:  $\mathbf{z}^{(i,j)} \sim N(\mathbf{z} | \mathbf{0}, \mathbf{I}), \quad j = 1, \dots, m$

$$\mathbb{E}_{N(\mathbf{z} | \mathbf{0}, \mathbf{I})} \left[ p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \neq \frac{1}{m} \sum_{j=1}^m p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,j)}) \quad \Rightarrow \quad \text{bad approximation unless } m \text{ is huge}$$

But what about the present situation:  $\mathbf{z}^{(i,j)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}), \quad j = 1, \dots, m$

$$\mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}] \right\|_2^2 \right] \stackrel{?}{\iff} \frac{1}{m} \sum_{j=1}^m \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}^{(i,j)}, \boldsymbol{\theta}] \right\|_2^2$$

Unlike the prior  $N(\mathbf{z} | \mathbf{0}, \mathbf{I})$ , during training the encoder  $q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})$ :

- Confines mass to narrow region of  $\mathbf{z}$ -space
  - *Excludes* regions that are unlikely to have produced  $\mathbf{x}^{(i)}$
- } much better for sampling

In practice, can use just  $m = 1$  sample at each training iteration:

$$\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}), \quad \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}] \right\|_2^2 \right] \approx \underbrace{\frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}^{(i)}, \boldsymbol{\theta}] \right\|_2^2}_{\text{unbiased estimator}}$$

# Reparameterization Trick

Data-term approximation:

$$\frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}} \left[ \mathbf{z}^{(i)}, \boldsymbol{\theta} \right] \right\|_2^2 \quad \Rightarrow \quad \begin{array}{l} \text{easy to minimize} \\ \text{over } \boldsymbol{\theta} \text{ via SGD} \end{array}$$

But what about sampling operator  $\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})$ ?

**Problem:** Cannot directly propagate gradients w.r.t.  $\boldsymbol{\phi}$  through sampling operator ...

Equivalent sampling procedures:

$$\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}) \quad \Leftrightarrow \quad \begin{array}{l} \boldsymbol{\varepsilon}^{(i)} \sim N(\boldsymbol{\varepsilon} | \mathbf{0}, \mathbf{I}) \\ \mathbf{z}^{(i)} = \boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] + \boldsymbol{\Sigma}_{\mathbf{z}}^{1/2}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] \boldsymbol{\varepsilon}^{(i)} \end{array}$$

Revised approximation:

$$\begin{aligned} \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}] \right\|_2^2 \right] &= \mathbb{E}_{N(\boldsymbol{\varepsilon} | \mathbf{0}, \mathbf{I})} \left[ \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}} \left( \boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] + \boldsymbol{\Sigma}_{\mathbf{z}}^{1/2}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] \boldsymbol{\varepsilon}, \boldsymbol{\theta} \right) \right\|_2^2 \right] \\ &\approx \frac{1}{2\gamma} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}} \left( \boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] + \boldsymbol{\Sigma}_{\mathbf{z}}^{1/2}[\mathbf{x}^{(i)}, \boldsymbol{\phi}] \boldsymbol{\varepsilon}^{(i)}, \boldsymbol{\theta} \right) \right\|_2^2 \end{aligned} \quad \left. \vphantom{\mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})}} \right\} \text{SGD friendly}$$

differentiable sample from encoder

# VAE Optimization Summary

Basic energy function:

$$\begin{aligned} L(\boldsymbol{\theta}, \boldsymbol{\varphi}) &\equiv \sum_i \left\{ \text{KL} \left[ q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right] - \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \right\} \\ &\geq -\sum_i \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \end{aligned}$$

Solve:

$$\begin{aligned} \boldsymbol{\theta}_*, \boldsymbol{\varphi}_* &= \operatorname{argmin}_{\boldsymbol{\theta}, \boldsymbol{\varphi}} L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \\ \text{s.t. } & \left. \begin{aligned} q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) &= N\left(\mathbf{z} \mid \boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\varphi}], \boldsymbol{\Sigma}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\varphi}]\right) \\ p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) &= N\left(\mathbf{x}^{(i)} \mid \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}], \boldsymbol{\Sigma}_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}]\right) \end{aligned} \right\} \begin{array}{l} \text{approximate via} \\ \text{reparameterization} \\ \text{trick + SGD} \end{array} \end{aligned}$$

# Generating New Samples

Simple hierarchical sampling:

$$\left. \begin{aligned} \mathbf{z}^{(j)} &\sim N(\mathbf{z} | \mathbf{0}, \mathbf{I}), \quad j = 1, \dots, m \\ \mathbf{x}^{(j)} &\sim p_{\theta_*}(\mathbf{x} | \mathbf{z}^{(j)}), \quad j = 1, \dots, m \end{aligned} \right\} \begin{array}{l} \text{only decoder with} \\ \text{optimized parameters} \\ \text{is needed} \end{array}$$

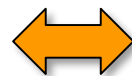
Typical to ignore decoder noise variance, i.e.,

$$\text{replace } \mathbf{x}^{(j)} \sim p_{\theta_*}(\mathbf{x} | \mathbf{z}^{(j)}) \text{ with } \mathbf{x}^{(j)} = \boldsymbol{\mu}_{\mathbf{x}}(\mathbf{z}^{(j)}, \boldsymbol{\theta}_*) \quad \Rightarrow \quad \text{cleaner samples}$$

Ideal scenario:

new samples

$$\left\{ \mathbf{x}^{(j)} \right\}_{j=1}^m$$



similar in  
distribution

training data

$$\left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n \sim p_{gt}(\mathbf{x})$$



# MNIST Examples

Ground-truth samples



VAE-generated samples  
with  $\Sigma_x[\mathbf{z}, \boldsymbol{\theta}] = \mathbf{I}, \forall \mathbf{z}$



(better VAE options available; Part IV)

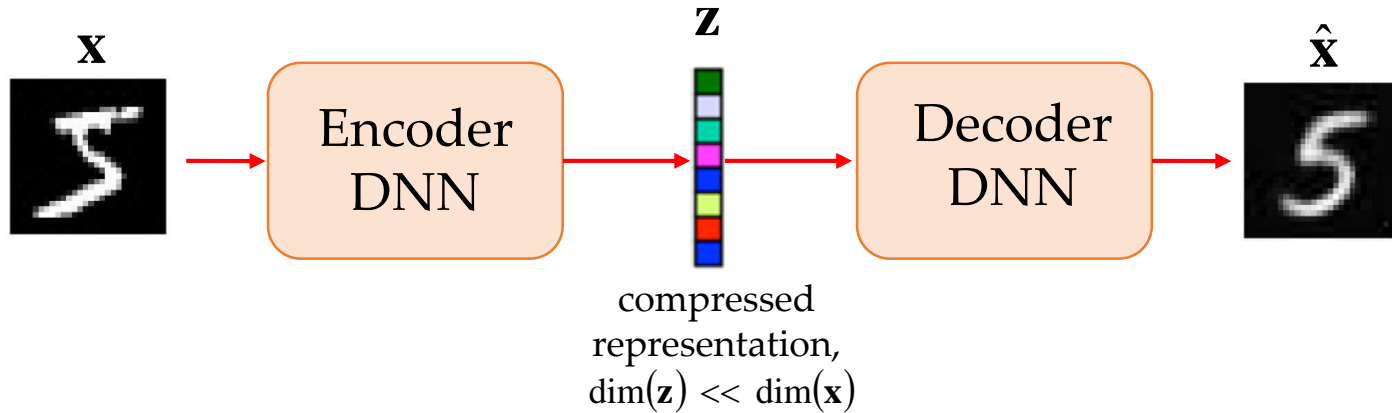
# Computing Negative Log-Likelihood (NLL) Estimates

Can apply unbiased estimate of VAE bound:

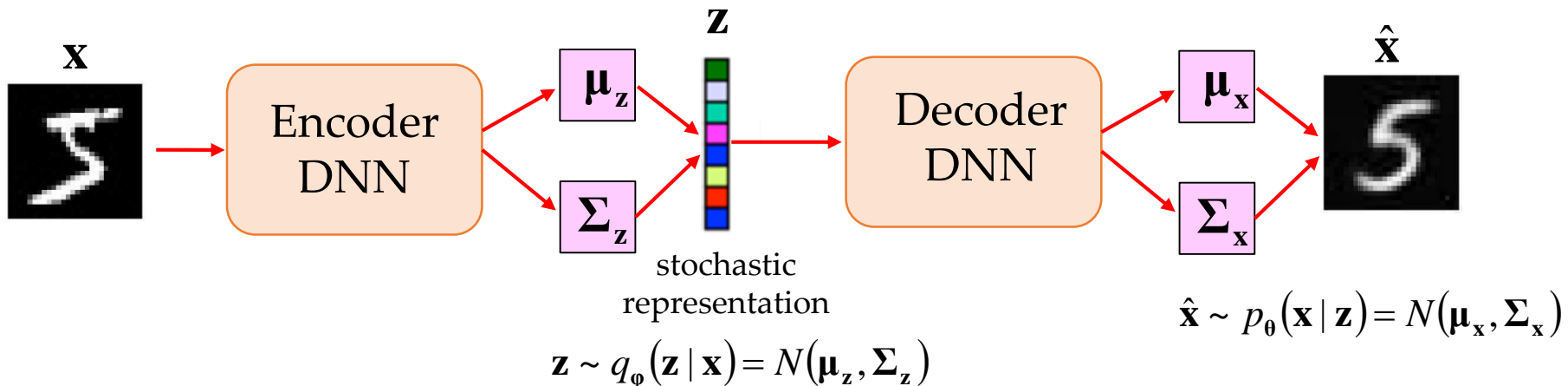
$$\begin{aligned} -\log p_{\theta}(\mathbf{x}^{test}) &\leq \underbrace{\text{KL}[q_{\phi}(\mathbf{z} | \mathbf{x}^{test}) \| N(\mathbf{z} | \mathbf{0}, \mathbf{I})]}_{\text{exact, closed-form}} - \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{test})}[\log p_{\theta}(\mathbf{x}^{test} | \mathbf{z})]}_{\text{use unbiased estimate}} \\ &\approx \frac{1}{m} \sum_{j=1}^m \frac{1}{2\gamma} \left\| \mathbf{x}^{test} - \boldsymbol{\mu}_{\mathbf{x}}[\mathbf{z}^{(j)}, \boldsymbol{\theta}] \right\|_2^2 \\ &\quad \mathbf{z}^{(j)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}^{test}) \end{aligned}$$

# Comparison with an Autoencoder

## Autoencoder (AE):



## VAE:

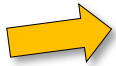


**Questions?**

**Part III:**  
**Connections with Existing Signal Processing  
Models for Finding Low-Dimensional Structure**

**Note:** Updated version of slides available at <http://www.davidwipf.com/>

# Outline

- ❑ Finding low-dimensional structure in high-dimensional data, possibly corrupted with outliers
- ❑ NP-hard decompositions into inlier and sparse outlier components
- ❑ Weaknesses of existing methods and useful VAE-based alternatives
- ❑ Case Study: Robust PCA  representative example
  - Connections with restricted class of VAE models
  - Advantages of the VAE in finding low-dimensional structure

# Context

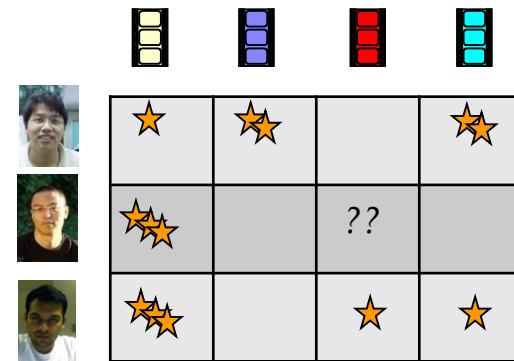
- Data is increasingly massive, high-dimensional



Images



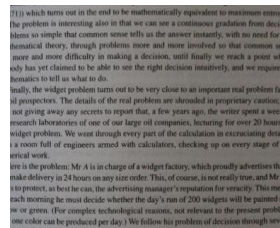
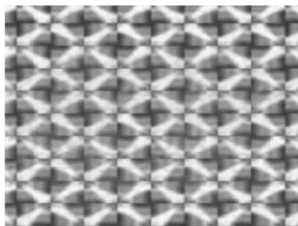
Videos



User data

- Blessing of dimensionality:**

Real data often concentrate on low-dimensional or degenerate structures in high-dimensional ambient space



local regularities, global symmetries, repetitive patterns, redundant sampling ...

# Robust Estimation

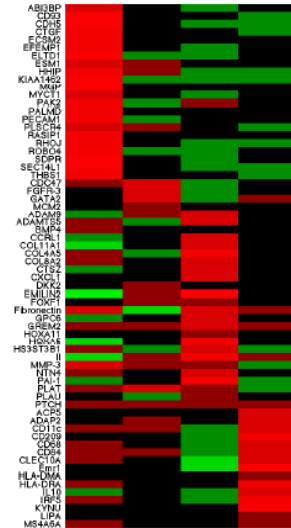
- But real-world data also frequently contain extraneous features, missing observations, or corruptions/outliers



face recognition  
[Wright et al., 2009]

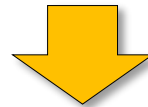


3D reconstruction  
[Zhang et al., 2011]



gene expression  
[Wang et al., 2012]

- Traditional methods (e.g., PCA, least squares regression) break down ...



**Replacements:** Robust PCA, sparse representations/regression, and many others



# Building Blocks for Robust Estimation

- Sparse representations:

$$\mathbf{y} = \begin{bmatrix} -4 \\ -5 \\ 3 \end{bmatrix}, \quad \Phi = \begin{bmatrix} 1 & 4 & 1 & 1 & 6 \\ -2 & 1 & -4 & 2 & -3 \\ 3 & 3 & 2 & -2 & 1 \end{bmatrix}$$

feasible solutions to  $\mathbf{y} = \Phi\mathbf{x}$

$$\mathbf{u} = \begin{bmatrix} 4 \\ -1 \\ 3 \\ 5 \\ -2 \end{bmatrix} \quad \mathbf{u}_0 = \begin{bmatrix} 0 \\ 0 \\ 2 \\ 0 \\ -1 \end{bmatrix}$$

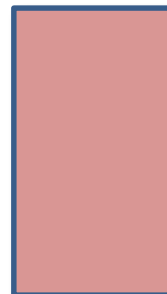
non-sparse

sparse

- Low-Rank matrices:



=



\*



defines low-dimensional  
subspaces

# High-Level Data Decomposition

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

Basic building blocks can be combined in various ways to construct models of the form:

$$\mathbf{X} = \underbrace{\text{inlier component } (\mathbf{L})}_{\mathbf{L} = \left\{ \mathbf{l}^{(i)} \right\}_{i=1}^n = \mathbf{U}\mathbf{Z}, \mathbf{l}^{(i)} = \mathbf{U}\mathbf{z}^{(i)}} + \underbrace{\text{outlier/noise component } (\mathbf{E})}_{\mathbf{E} = \left\{ \mathbf{e}^{(i)} \right\}_{i=1}^n}$$

- low-dimensional latent structure, e.g.,  
 $\mathbf{U}$  defines a low-dim inlier subspace
- sometimes not fully observable, e.g.,  
have measurement operator  $\mathcal{A}(\mathbf{L})$
- sparse corruptions (possibly large)
- other errors or model mismatch

Background detection example:



observed video  
frames

=



1D subspace  
background  
component

+




sparse foreground  
component

# Typical Objective for Signal Recovery

Challenging ill-posed inverse problem to recover low-dimensional representation  $\mathbf{L}$  :

$$\min_{\mathbf{L}, \mathbf{E}} \left\| \mathbf{X} - \mathcal{A}(\mathbf{L}) - \mathbf{E} \right\|_2^2 + \lambda_1 g_1(\mathbf{L}) + \lambda_2 g_2(\mathbf{E})$$

$$\min_{\mathbf{L}, \mathbf{E}} g_1(\mathbf{L}) + \lambda g_2(\mathbf{E}), \quad \text{s.t. } \mathbf{X} = \mathcal{A}(\mathbf{L}) + \mathbf{E} \quad (\text{constrained version})$$

$\mathcal{A}$  : linear measurement operator  
 $g_1$  : favors low-dim representations  $\rightarrow$   $\mathbf{L} = \mathbf{U} * \mathbf{Z}$   
 $g_2$  : favors sparsity  $\rightarrow$  

**Example penalties:**  $g_1(\mathbf{L}) = \text{rank}(\mathbf{L}) \equiv \|\sigma(\mathbf{L})\|_0 \rightarrow$  # nonzero singular values of  $\mathbf{L}$

$g_2(\mathbf{E}) = \|\mathbf{E}\|_0 \rightarrow$  # nonzero elements in  $\mathbf{E}$

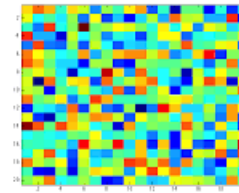
**Note:** Penalties are primarily used for limiting:

- 1) the intrinsic dimensionality of the inlier space
- 2) the cardinality of outliers

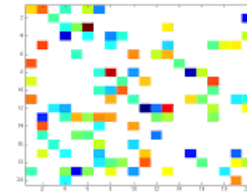
... not used for learning distribution within the inlier space

# Special Cases

- Matrix recovery/completion:



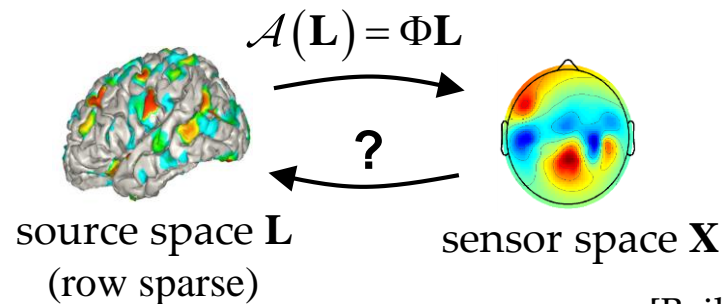
low rank  $\mathbf{L}$



partial observations  $\mathcal{A}(\mathbf{L})$

[Candès and Recht, 2008]

- Source localization:



[Baillet et al., 2001]

- Robust PCA:

observations  $\mathbf{X}$  = low rank  $\mathbf{L}$  + sparse  $\mathbf{E}$

**case study**

[Chandrasekaran et al., 2011; Candès et al. 2011]

- Many more ...

# Weakness of Traditional Pipeline

## Primary:

- Difficult nonconvex, NP-hard estimation process:

$$\min_{\mathbf{L}, \mathbf{E}} g_1(\mathbf{L}) + \lambda g_2(\mathbf{E}), \quad \text{s.t. } \mathbf{X} = \mathcal{A}(\mathbf{L}) + \mathbf{E}$$

(and convex relaxations often fail ...)

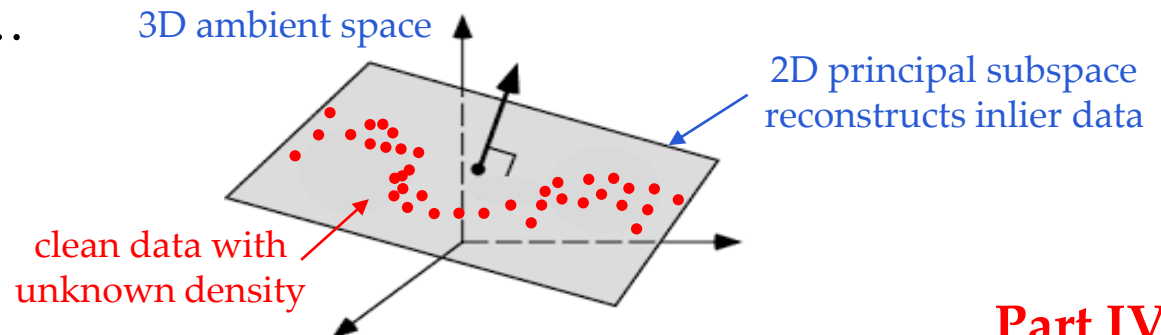
- Limited capacity inlier models, e.g.,

$$\mathbf{L} = \mathbf{U} * \mathbf{Z}$$

**remainder of Part III**

## Secondary:

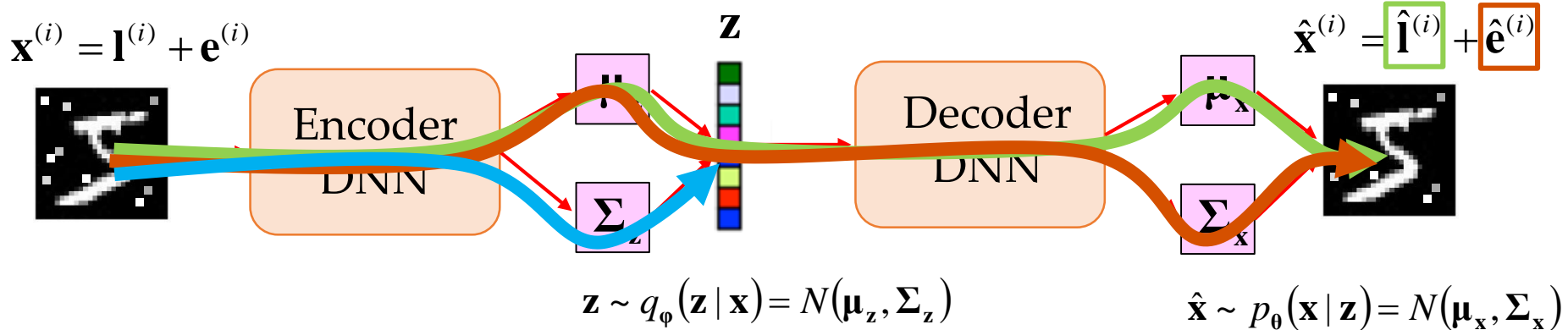
- Limited generative modeling capability, primarily used for data reconstruction ...



**Part IV**

# How might the VAE model help?

Basic VAE architecture:



## High-Level Picture

Correspondences between VAE components and signal recovery:

- Deterministic path provides nonlinear inlier model:



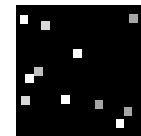
$$\boldsymbol{\mu}_x(\boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\phi}], \boldsymbol{\theta}) \approx \mathbf{l}^{(i)} =$$



- Decoder covariance path models sparse outliers:



$$\boldsymbol{\Sigma}_x(\boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\phi}], \boldsymbol{\theta}) \approx (\mathbf{e}^{(i)})^2 =$$



- VAE Encoder covariance:

- 1) At global optima: determines inlier dimensionality
- 2) Elsewhere: smooths bad local minima

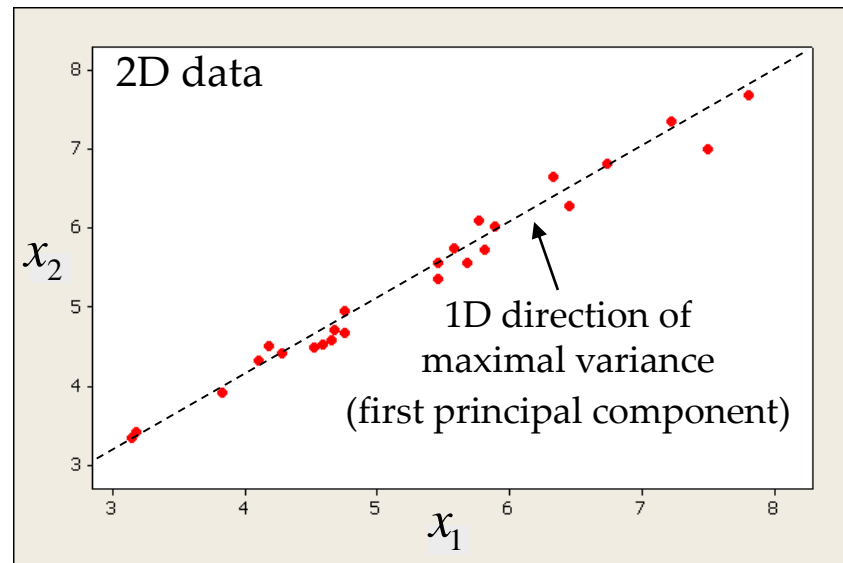
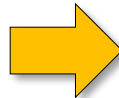
# Case Study: Robust PCA

## Why this is a good choice?

- Highly influential model e.g., [4358 citations](#) to [Candès et al, 2011].
- Exactly follows inlier + outlier data decomposition (common to many signal processing applications ...).
- Limited by
  - NP-hard estimation,
  - simple low-rank (bilinear inlier) model
- Correspondences with VAE can be explicitly quantified.
- Representative of connections between the VAE and other ill-posed inverse problems (e.g., compressive sensing, source localization, subspace clustering, matrix completion, ...).

# PCA Background

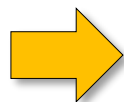
PCA finds directions of maximal variance:



Many different formulations given data  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n \in \mathbb{R}^{d \times n}$

Example (AE-like):

$$\underbrace{\mathbf{U}_*, \mathbf{V}_*}_{\text{defines } \kappa \text{ principal component directions}} = \arg \min_{\mathbf{U}, \mathbf{V}} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \underbrace{\mathbf{U} \mathbf{z}^{(i)}}_{\text{linear decoder}} \right\|_2^2, \quad \text{s.t. } \mathbf{z}^{(i)} = \underbrace{\mathbf{V} \mathbf{x}^{(i)}}_{\text{linear encoder}} \quad \forall i, \quad \mathbf{U} \in \mathbb{R}^{d \times \kappa}, \quad \mathbf{V} \in \mathbb{R}^{\kappa \times d}$$



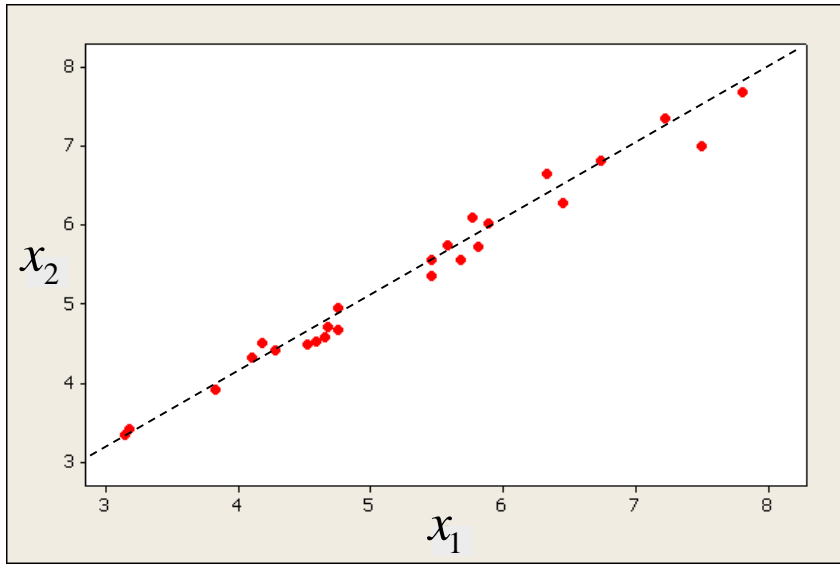
Simple AE can compute principal components

[Bourlard and Kamp, 1988]

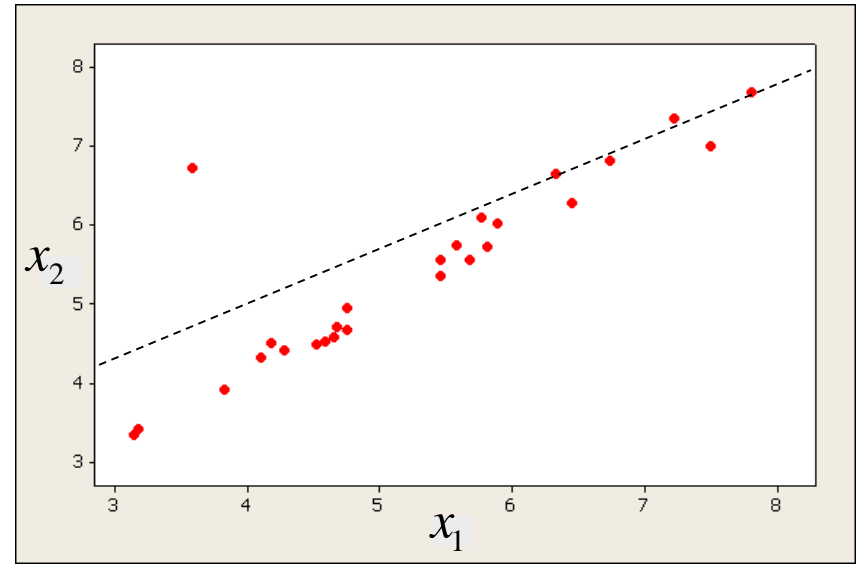


# PCA Sensitivity to Outliers

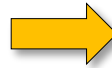
No Outliers



Single Outlier



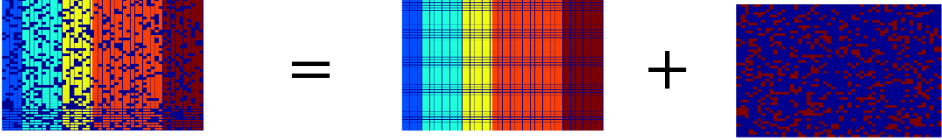
dashed line



first principal  
component direction

# Robust PCA

[Chandrasekaran et al., 2011; Candès et al. 2011]



observations ( $\mathbf{X}$ ) = low rank ( $\mathbf{L}$ ) + sparse outliers ( $\mathbf{E}$ )

RPCA inverse problem:

$$\mathbf{L}_*, \mathbf{E}_* = \arg \min_{\mathbf{L}, \mathbf{E}} \underbrace{\text{rank}[\mathbf{L}] + \frac{1}{n} \|\mathbf{E}\|_0}_{\text{NP-hard}} \quad \text{s.t. } \mathbf{X} = \mathbf{L} + \mathbf{E}$$

Convex relaxation:

$$\hat{\mathbf{L}}, \hat{\mathbf{E}} = \arg \min_{\mathbf{L}, \mathbf{E}} \|\mathbf{L}\|_* + \frac{1}{\sqrt{n}} \|\mathbf{E}\|_1 \quad \text{s.t. } \mathbf{X} = \mathbf{L} + \mathbf{E}$$

**Theory:**  $\{\hat{\mathbf{L}}, \hat{\mathbf{E}}\} = \{\mathbf{L}_*, \mathbf{E}_*\}$  in very specialized conditions, but these rarely hold in practice ...



**What about the VAE?**

# Illustrative Degenerate Case

Original VAE objective:

$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \sum_i \left\{ \text{KL} \left[ q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel p(\mathbf{z}) \right] - \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \right\}$$

Two assumptions:

1. Degenerate encoder covariance  $\Rightarrow \Sigma_{\mathbf{z}}[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{0}$
2. High capacity decoder covariance  $\Rightarrow \Sigma_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}]$  arbitrary

Can collapse VAE objective using *assumption 1*:

$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \equiv - \sum_i \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z} = \boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\varphi}])$$

Further simplification possible using *assumption 2*:

$$\min_{\Sigma_{\mathbf{x}}} L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \equiv \sum_i \sum_j \log |e_j^{(i)}| \quad \text{s.t.} \quad \mathbf{e}^{(i)} = \mathbf{x}^{(i)} - \boldsymbol{\mu}_{\mathbf{x}}[\boldsymbol{\mu}_{\mathbf{z}}[\mathbf{x}^{(i)}, \boldsymbol{\varphi}], \boldsymbol{\theta}], \quad \forall \mathbf{x}^{(i)}$$

Equivalent to a deterministic autoencoder with Gaussian entropy loss ...

# VAE and Induced AE Side-by-Side

## VAE model

$$\Sigma_z[\mathbf{x}, \boldsymbol{\phi}] \text{ arbitrary}$$

$$\Sigma_x[\mathbf{z}, \boldsymbol{\theta}] \text{ arbitrary}$$

## Induced AE

$$\Sigma_z[\mathbf{x}, \boldsymbol{\phi}] = \mathbf{0}$$

$$\Sigma_x[\mathbf{z}, \boldsymbol{\theta}] \text{ arbitrary}$$

Induced AE is like a typical AE but with an outlier robust loss function:

$$L_{\text{AE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) \triangleq \underbrace{\sum_i \sum_j \log |e_j^{(i)}|}_{\substack{\text{approximates} \\ \ell_0 \text{ norm}}} \quad \text{s.t. } \mathbf{e}^{(i)} = \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\phi}], \boldsymbol{\theta}], \quad \forall \mathbf{x}^{(i)}$$

$$\Rightarrow \sum_i \sum_j \log |e_j^{(i)}| = \lim_{p \rightarrow 0} \sum_i \sum_j \frac{1}{p} \left( |e_j^{(i)}|^p - 1 \right) \equiv \|\mathbf{E}\|_0$$

VAE is like a smoothed, regularized version of the induced AE:

$$L(\boldsymbol{\theta}, \boldsymbol{\phi}) = \sum_i \left\{ \text{KL} \left[ q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) \| p(\mathbf{z}) \right] - \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right] \right\}$$

$$\geq \underbrace{\sum_i \text{KL} \left[ q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) \| p(\mathbf{z}) \right]}_{\text{regularization}} + \underbrace{\sum_i \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \sum_j \log |\tilde{e}_j^{(i)}| \right]}_{\text{smoothing}} \Rightarrow \text{best possible w/ } \Sigma_x[\mathbf{z}, \boldsymbol{\theta}] \text{ arbitrary}$$

$$\text{s.t. } \tilde{\mathbf{e}}^{(i)} = \mathbf{x}^{(i)} - \boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}]$$

Both VAE and Induced AE have a distinct relationship with **Robust PCA** ...

# RPCA and the Induced AE

$$\text{RPCA: } \mathbf{L}_*, \mathbf{E}_* = \arg \min_{\mathbf{L}, \mathbf{E}} \text{rank}[\mathbf{L}] + \frac{1}{n} \|\mathbf{E}\|_0 \quad \text{s.t. } \mathbf{X} = \mathbf{L} + \mathbf{E}$$

Assume *affine* decoder mean (arbitrary encoder mean):

$$\boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \mathbf{W}\mathbf{z} + \mathbf{b}, \quad \boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}\}$$

$$\dim[\mathbf{z}] = \text{rank}[\mathbf{L}_*]$$

**Result:** Induced AE shares the same combinatorial constellation of local and global minima of the constrained RPCA problem

$$\min_{\mathbf{L}, \mathbf{E}} \|\mathbf{E}\|_0 \quad \text{s.t.} \quad \begin{array}{l} \mathbf{X} = \mathbf{L} + \mathbf{E} \\ \text{rank}[\mathbf{L}] \leq \text{rank}[\mathbf{L}_*] \end{array} \quad \rightarrow \quad \boxed{\text{local minima a huge issue}}$$

**Additional concern:** If  $\dim[\mathbf{z}] \neq \text{rank}[\mathbf{L}_*]$ , then even the global minimum need not be optimal ...

$\rightarrow$  superfluous dimensions can cause trouble

# RPCA and the VAE

$$\text{RPCA: } \mathbf{L}_*, \mathbf{E}_* = \arg \min_{\mathbf{L}, \mathbf{E}} \text{rank}[\mathbf{L}] + \frac{1}{n} \|\mathbf{E}\|_0 \quad \text{s.t. } \mathbf{X} = \mathbf{L} + \mathbf{E}$$

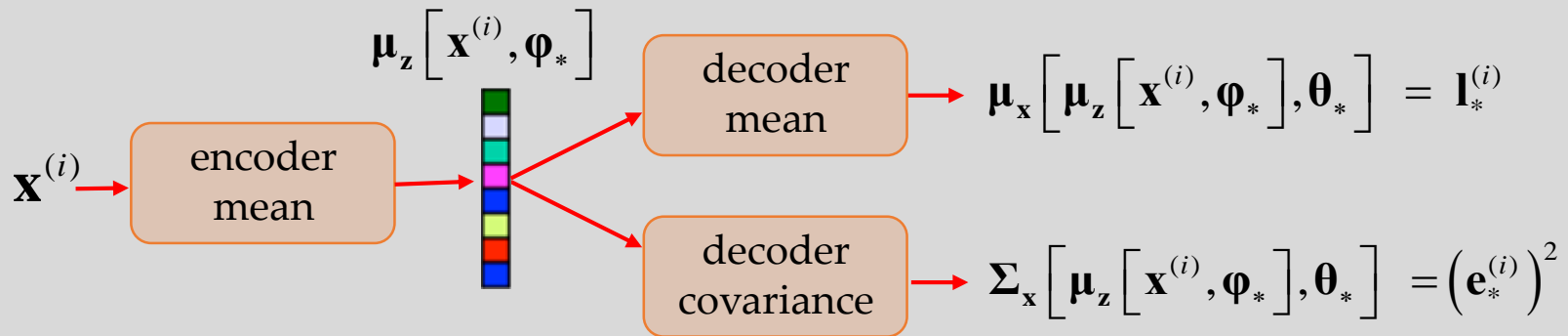
Assume *affine* decoder mean (arbitrary encoder mean):

$$\boldsymbol{\mu}_x[\mathbf{z}, \boldsymbol{\theta}] = \mathbf{W}\mathbf{z} + \mathbf{b}$$

$$\dim[\mathbf{z}] \geq \text{rank}[\mathbf{L}_*]$$

## Theorem (Perfect Recovery):

Any VAE global optimum  $\{\boldsymbol{\theta}_*, \boldsymbol{\varphi}_*\}$  is such that:



Matching global optima ... **even after smoothing!**

... true even if  $\dim[\mathbf{z}] > \text{rank}[\mathbf{L}_*]$

# Two Underappreciated Distinctions

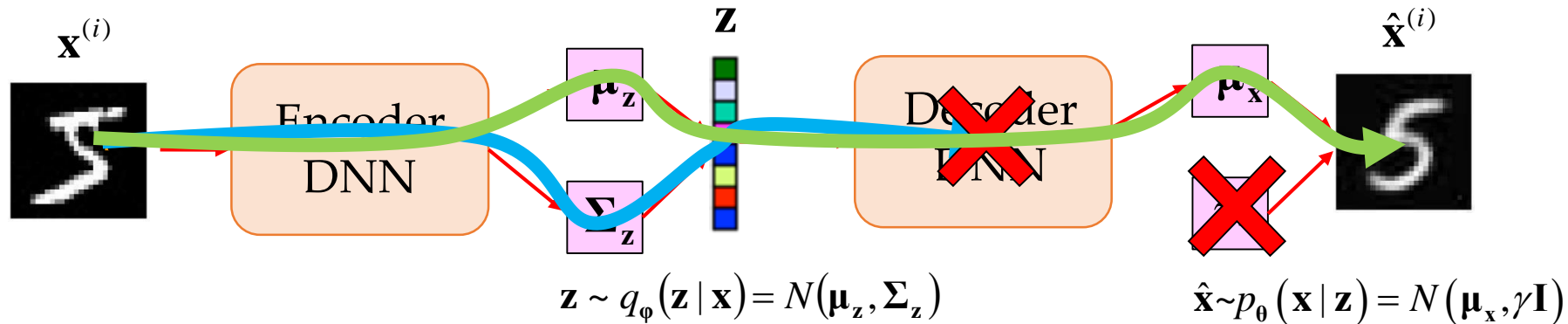
- 1) VAE can learn the optimal/minimal latent dimension of inlier model ... unnecessary dimensions **can be automatically discarded**.
- 2) VAE smoothing/KL regularization impacts bad local minimum, does **not** change the global optimum.

**Note:** VAE capabilities motivated by Robust PCA example, but also translate to more complex inlier models

# Discarding Unnecessary Latent Dimensions

Observed data:  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n$ ,  $\mathbf{x}^{(i)} = \mathbf{I}^{(i)} \in \mathcal{X}$   $\rightarrow$  arbitrary inlier manifold; for simplicity no outliers

Assumed VAE (arbitrary encoder/decoder networks):



## Theorem (Reconstruction Invariance):

Under some technical conditions, any VAE global optimum  $\{\boldsymbol{\theta}_*, \boldsymbol{\varphi}_*\}$  is such that  $\gamma \rightarrow 0$  and reconstructions are exact:

$$\boldsymbol{\mu}_x \left( \boldsymbol{\mu}_z \left[ \mathbf{x}^{(i)}, \boldsymbol{\varphi}_* \right] + \boldsymbol{\Sigma}_z^{1/2} \left[ \mathbf{x}^{(i)}, \boldsymbol{\varphi}_* \right] \boldsymbol{\varepsilon}, \boldsymbol{\theta}_* \right) = \boldsymbol{\mu}_x \left( \boldsymbol{\mu}_z \left[ \mathbf{x}^{(i)}, \boldsymbol{\varphi}_* \right], \boldsymbol{\theta}_* \right) = \mathbf{x}^{(i)}, \quad \forall \boldsymbol{\varepsilon}, \forall i$$

**Key Conclusion:** At global minimum, encoder randomness will not impact perfect reconstructions  $\rightarrow$  can be “pruned” with white noise



# Discarding Unnecessary Latent Dimensions Cont.

- Recall VAE KL term with Gaussian encoder satisfies

$$\text{KL}\left[q_\phi(\mathbf{z}|\mathbf{x}^{(i)})\|N(\mathbf{z}|\mathbf{0},\mathbf{I})\right] \propto \left\|\boldsymbol{\mu}_z[\mathbf{x}^{(i)},\boldsymbol{\phi}]\right\|_2^2 + \text{tr}\left(\boldsymbol{\Sigma}_z[\mathbf{x}^{(i)},\boldsymbol{\phi}]\right) - \log\left|\boldsymbol{\Sigma}_z[\mathbf{x}^{(i)},\boldsymbol{\phi}]\right|$$

- With diagonal covariance (common choice), further decouples to

$$\text{KL}\left[q_\phi(\mathbf{z}|\mathbf{x}^{(i)})\|N(\mathbf{z}|\mathbf{0},\mathbf{I})\right] \propto \sum_{j=1}^K \left\{ \mu_z[\mathbf{x}^{(i)},\boldsymbol{\phi}]_j^2 + \sigma_z^2[\mathbf{x}^{(i)},\boldsymbol{\phi}]_j - \log\left(\sigma_z^2[\mathbf{x}^{(i)},\boldsymbol{\phi}]_j\right) \right\}$$

- Reconstruction Invariance Theorem implies that certain dimensions will not influence VAE data term.

- Along these dimensions, KL term can be minimized independently:

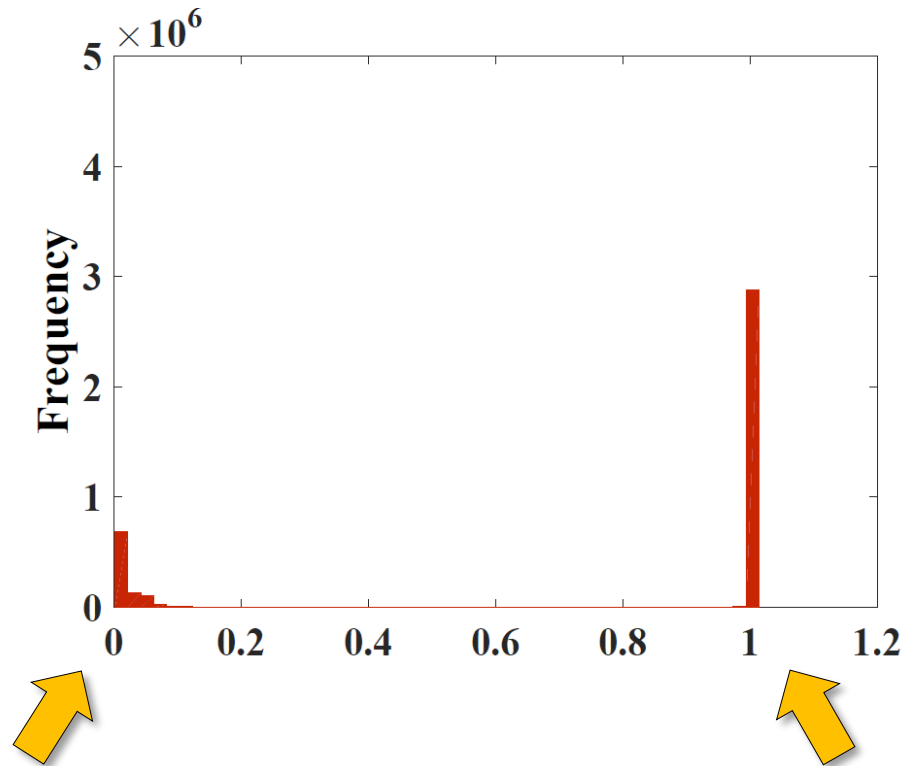
Optimal moments for these unnecessary dimensions:

$$\mu_z[\mathbf{x}^{(i)},\boldsymbol{\phi}]_j \rightarrow 0, \quad \sigma_z^2[\mathbf{x}^{(i)},\boldsymbol{\phi}]_j \rightarrow 1$$

- This non-informative white noise will be filtered out by the decoder.

# Empirical Example

Histogram of  $\sigma_z^2[\mathbf{x}^{(i)}, \boldsymbol{\phi}]_j$  values for VAE trained on MNIST data



for useful dimensions, encoder variance is near zero; facilitates good reconstructions

for unnecessary dimensions, encoder variance is near one; optimizes KL term

( Encoder noise will serve an important purpose in Part IV... )

# Filtering Unnecessary Dimensions

$$\sigma_z^2[\mathbf{x}^{(i)}, \boldsymbol{\phi}]_j = 1.0 \quad \rightarrow \quad \text{unnecessary dimension}$$

Reconstructions as we change latent code along this dimension (other dimensions fixed)



Image Variance = 0.000 **no changes**

---

$$\sigma_z^2[\mathbf{x}^{(i)}, \boldsymbol{\phi}]_j = 0.005 \approx 0 \quad \rightarrow \quad \text{necessary dimension}$$

Reconstructions as we change latent code along this dimension (other dimensions fixed)



Image Variance = 27.20 **large changes**

# Two Underappreciated Distinctions

- 1) VAE can learn the optimal/minimal latent dimension of inlier model ... unnecessary dimensions **can be automatically discarded**.

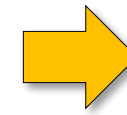
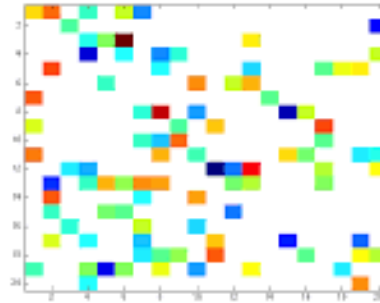
- 2) VAE smoothing/KL regularization impacts bad local minimum, does **not** change the global optimum.

**Note:** VAE capabilities motivated by Robust PCA example, but also translate to more complex inlier models

# Benefits of VAE smoothing

With induced AE (no smoothing), we enter a local minima at any outlier support pattern

$$\mathbf{E} = \left[ \mathbf{e}^{(1)}, \dots, \mathbf{e}^{(n)} \right] =$$
$$\mathbf{e}^{(i)} = \mathbf{x}^{(i)} - \boldsymbol{\mu}_x \left[ \boldsymbol{\mu}_z \left[ \mathbf{x}^{(i)}, \boldsymbol{\varphi} \right], \boldsymbol{\theta} \right]$$

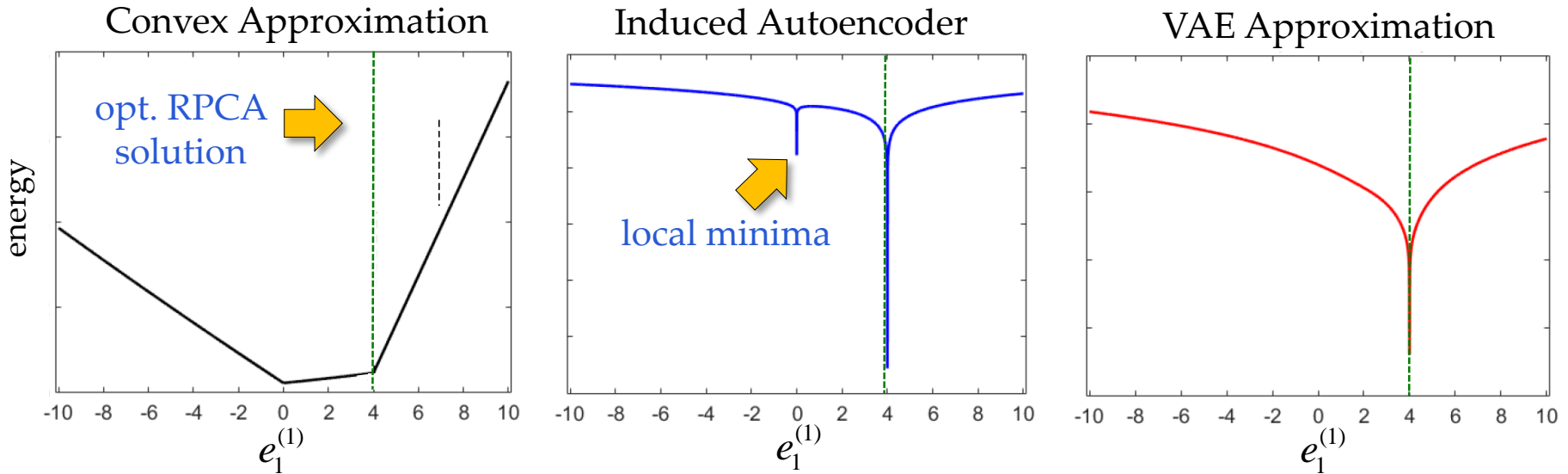
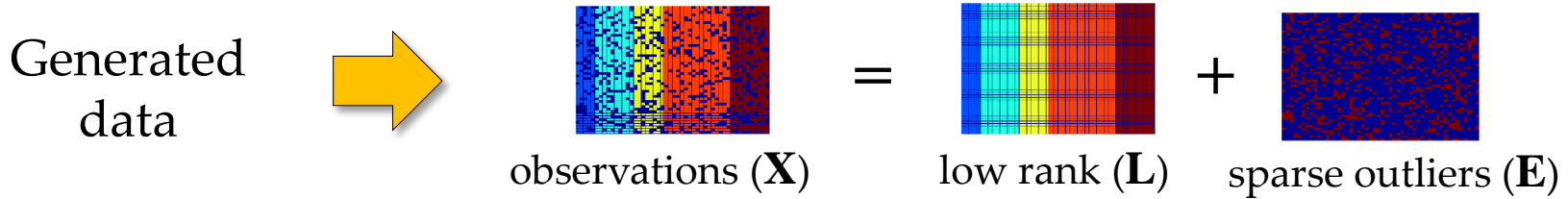


zero-valued elements can never change

But for the VAE, every support pattern need **not** be a local minimum because of **selective smoothing** ...

does not impact global minimum  
(unlike convex relaxations ...)

# Illustration of Selective Smoothing Effects



Representative 1D slice of energy functions while varying the coefficient

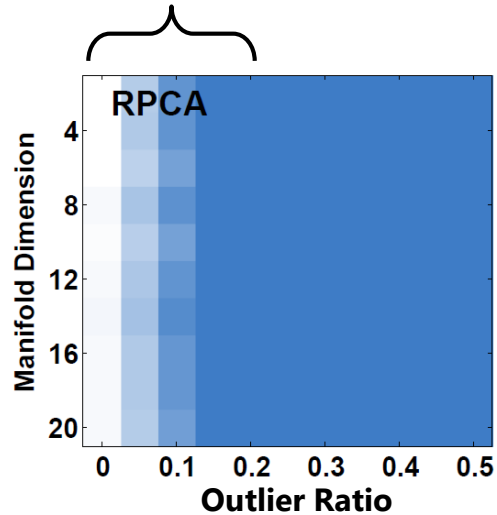
$$e_1^{(1)} = x_1^{(1)} - \mu_x \left[ \boldsymbol{\mu}_z \left[ \mathbf{x}^{(1)}, \boldsymbol{\phi} \right], \boldsymbol{\theta} \right]_1$$

# Non-Linear Manifold Recovery

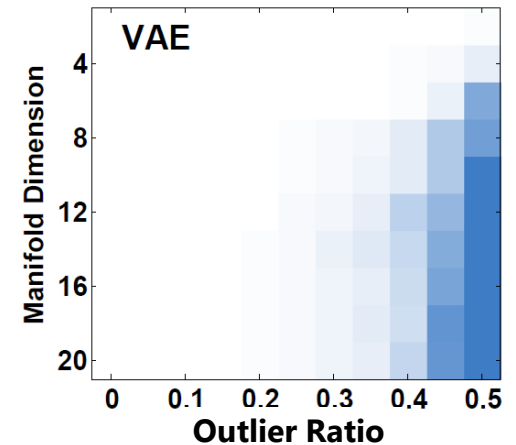
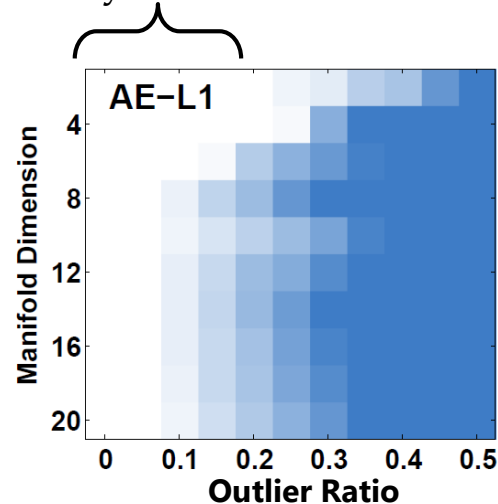
$\mathbf{X}$  = low-dimensional manifold component + sparse outlier component

↳  $\mu_{\mathbf{x}} \left[ \mu_{\mathbf{z}} \left[ \mathbf{x}^{(i)}, \boldsymbol{\varphi}' \right], \boldsymbol{\theta}' \right], \forall \mathbf{x}^{(i)}$

convex relaxation



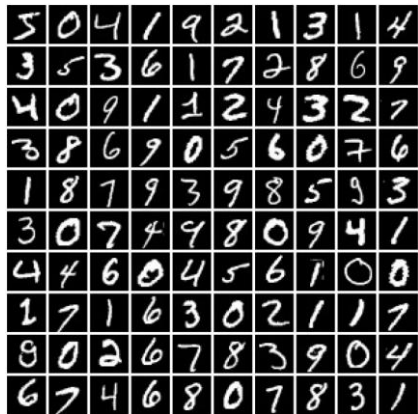
induced AE with extra  $\ell_1$  penalty on latent code



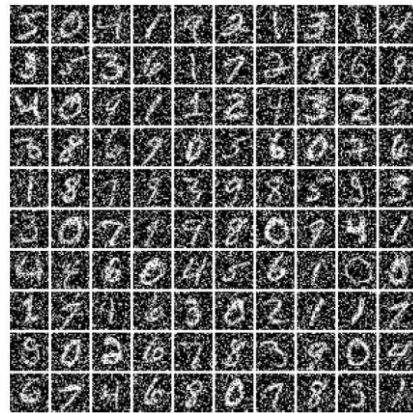
white = success (zero error), blue = failure (large error)

# MNIST Example

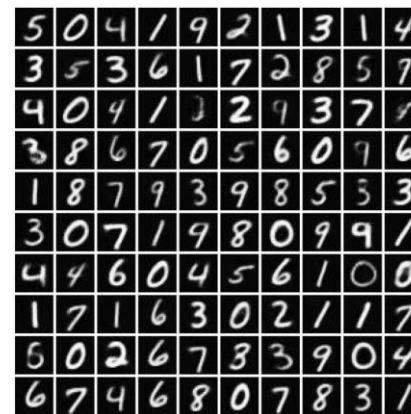
Original data



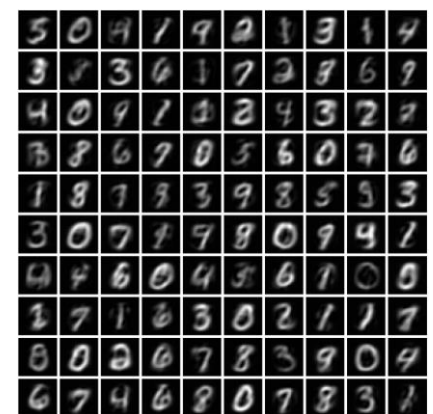
40% corrupted



VAE reconstructions



Convex RPCA reconstructions





# A Lingering Issue ...

A large training corpus  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n$  is required for learning complex manifolds with outliers

## Solution:

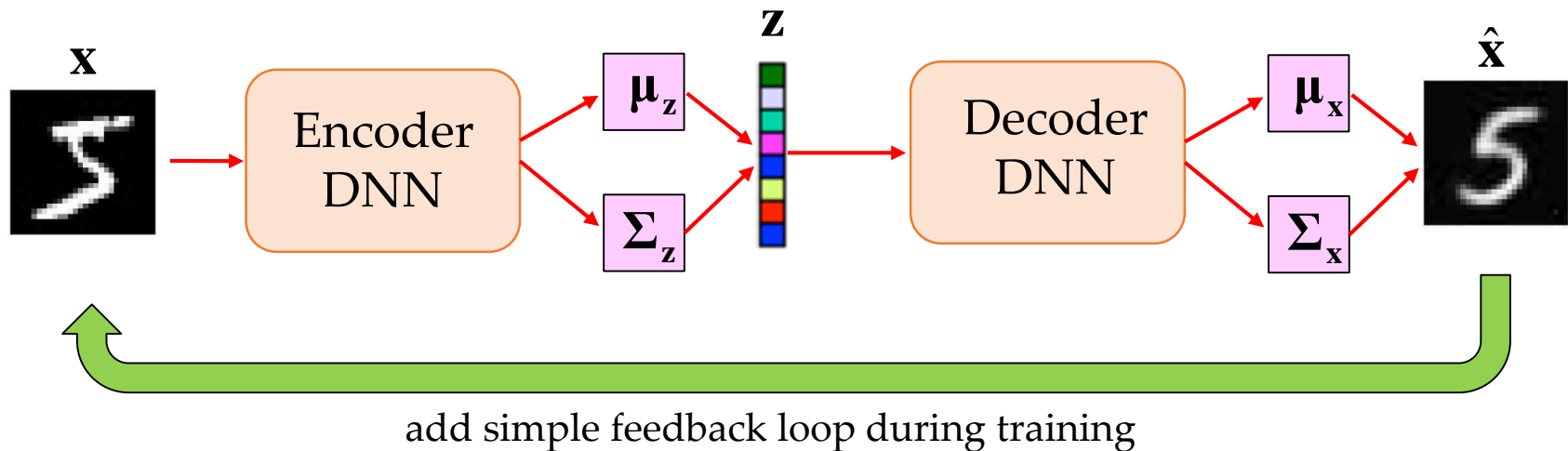
Recycle dirty samples via specialized recurrent connections ... [automated data augmentation](#)

[Wang et al., 2018]



# Recycled/Recurrent VAE

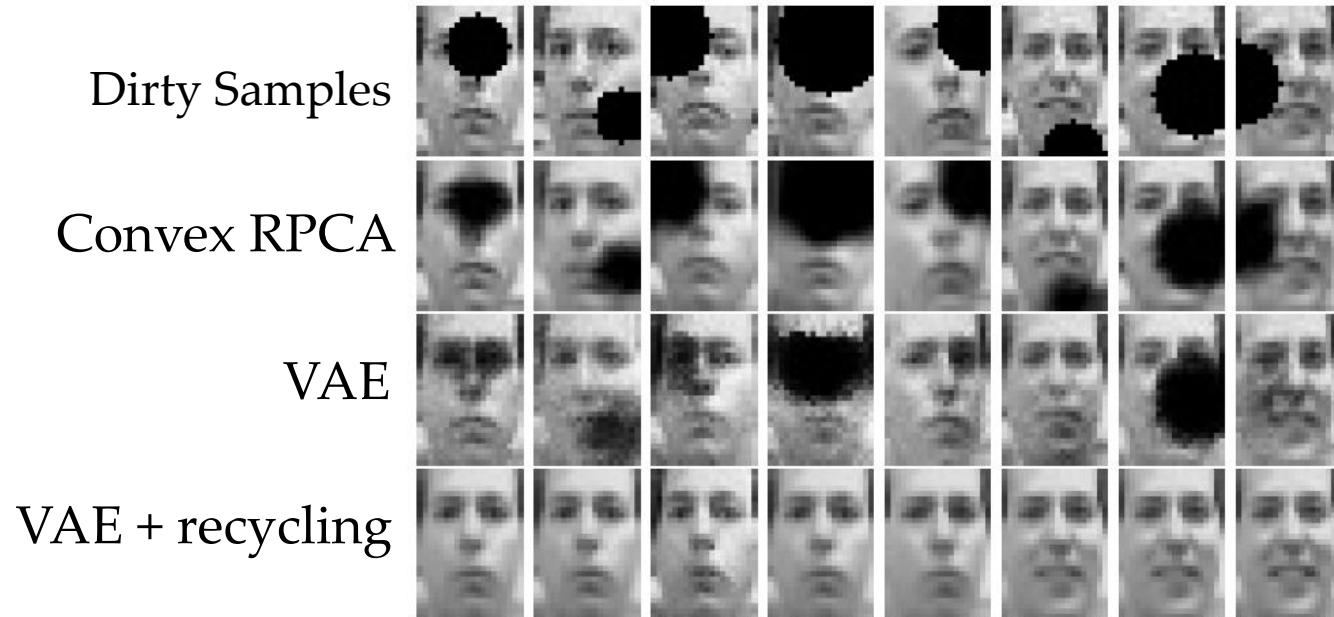
Given a single input sample, bootstrap virtual samples via recurrent connection



## Properties:

- ❑ No additional parameters required (simple SGD still works ...)
- ❑ Partially detected outliers can be removed in multiple passes
- ❑ Close connection to iterative reweighting algorithms

# Frey Face Data Recovery



# Summary of Robust PCA Case Study

- The VAE with an affine decoder mean collapses to a robust PCA variant with attractive properties.
  
- In broader regimes, can be viewed as powerful nonlinear extension.
  
- Analysis reveals underappreciated effects of VAE regularization:
  1. Can learn optimal latent dimensionality
  2. Can selectively smooth away bad local minima while preserving good global solutions
  3. Can potentially be useful for deterministic data cleaning tasks unrelated to generative modeling per se.
  4. Extra recurrent connections/recycling, can serve as a useful form of data augmentation.
  
- Representative of connections between the VAE and other ill-posed inverse problems.

**Questions?**

# Part IV: From Signal Reconstruction to Generative Modeling

**Note:** Updated version of slides available at <http://www.davidwipf.com/>

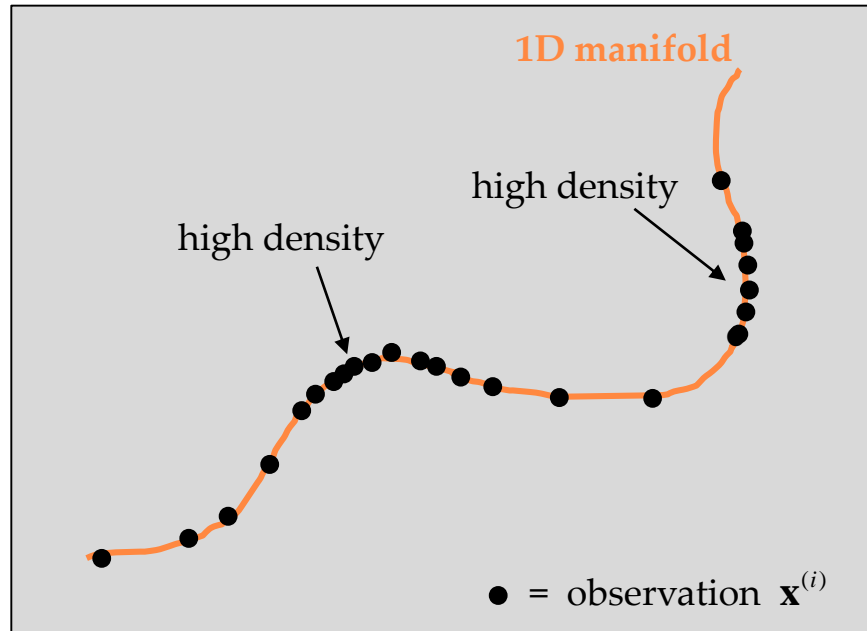
# Recap

- VAE can extend/enhance capabilities of traditional algorithms for finding low-dimensional structure.
- Low-dimensional structure could be an arbitrary manifold.
- Can reconstruct data (possibly corrupted by outliers) by fitting a parsimonious inlier model.
- But this is not sufficient for a full generative model ...

Note: Will mostly assume no outliers in Part IV for simplicity ... but the key concepts generalize.

# Illustration

2D ambient space



- ❑ Reconstructing data using parsimonious inlier model provides estimate of the 1D manifold (Part III).
- ❑ But it does **not** provide any information about the data distribution **within** the manifold.
- ❑ **Key question:** How can good reconstructions segue to a good generative model?



# Revisiting Original VAE Bound

Variational upper bound (from Part II):

$$\begin{aligned}
 -\sum_i \log p_{\theta}(\mathbf{x}^{(i)}) &\leq \sum_i \left\{ \text{KL} \left[ q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)}) \right] - \log p_{\theta}(\mathbf{x}^{(i)}) \right\} \\
 &\equiv \sum_i \left\{ \text{KL} \left[ q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right] - \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \right\}
 \end{aligned}$$

Equality iff:

$$\underbrace{q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)})}_{\text{encoder distribution}} = p_{\theta}(\mathbf{z} | \mathbf{x}^{(i)}) = \frac{p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I})}{\int p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}}$$

decoder distribution

**Consequence:** If encoder and decoder are sufficiently complex such that

Gaussian encoder/decoder  $\leftarrow$   $q_{\phi_*}(\mathbf{z} | \mathbf{x}) = p_{\theta_*}(\mathbf{z} | \mathbf{x})$   $\xrightarrow{\text{generally not Gaussian}}$  can estimate ground-truth distributions just by minimizing VAE cost

$$p_{\theta_*}(\mathbf{x}) = \int p_{\theta_*}(\mathbf{x} | \mathbf{z}) N(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z} = p_{gt}(\mathbf{x})$$

**Problem:** But typical VAEs for continuous data often involve Gaussian encoder and decoder distributions ... no match with true latent posterior.

# Impact of VAE Gaussian Assumptions

Assume for simplicity:

- model**
- Decoder covariance:  $\Sigma_{\mathbf{x}}[\mathbf{z}, \boldsymbol{\theta}] = \gamma \mathbf{I}, \quad \forall \mathbf{z}$  single learnable parameter  
(common in practice if no outliers)
  - Decoder mean, encoder mean/covariance all arbitrary functions
- 
- data**
- Asymptotic regime
 
$$L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \equiv \sum_{i=1}^n \left\{ \text{KL} \left[ q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right] - \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \right\}$$

$$\xrightarrow{n \rightarrow \infty} \int \left\{ \text{KL} \left[ q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x}) \parallel N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right] - \mathbb{E}_{q_{\boldsymbol{\varphi}}(\mathbf{z} | \mathbf{x})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z}) \right] \right\} \underbrace{\mu_{gt}(d\mathbf{x})}_{\text{ground-truth measure}}$$

$$\underbrace{\hspace{15em}}_{\vec{L}(\boldsymbol{\theta}, \boldsymbol{\varphi}) \Rightarrow \text{asymptotic loss}}$$
  - Potential low-dimensional structure in data (and no outliers):
 
$$\mu_{gt} \neq 0 \text{ on } r\text{-dimensional manifold } \mathcal{X} \Rightarrow \Pr(\mathbf{x} \notin \mathcal{X}) = 0$$

(Note: If  $r = \dim(\mathbf{x})$ , then no manifold structure)

# Impact of VAE Gaussian Assumptions Cont.

Notation:  $\dim(\mathbf{x}) = d$ ,  $\dim(\mathbf{z}) = \kappa$

i.e., no manifold, VAE  
latent dim large enough,  
and density exists

## Theorem (Exact Density Recovery):

Scenario:  $r = d$ ,  $\kappa \geq r$ , and  $\mu_{gt}(d\mathbf{x}) = p_{gt}(\mathbf{x})d\mathbf{x}$

Then any optimum  $\{\theta_*, \varphi_*\} \in \arg \min_{\theta, \varphi} \bar{L}(\theta, \varphi)$  will be such that\*

$$\text{KL}[q_{\varphi_*}(\mathbf{z}|\mathbf{x}) \| p_{\theta_*}(\mathbf{z}|\mathbf{x})] = 0 \quad \text{and} \quad p_{\theta_*}(\mathbf{x}) = \int p_{\theta_*}(\mathbf{x}|\mathbf{z})N(\mathbf{z}|\mathbf{0}, \mathbf{I})d\mathbf{z} = p_{gt}(\mathbf{x})$$

[Dai & Wipf, 2019]

### Positive:

- When there is no manifold, VAE global optimum exactly corresponds with recovery of ground-truth measure **even with Gaussian assumptions.**

### Negative Corollary:

- When there is a manifold, i.e.,  $r < d$ , cannot rule out globally optimal solutions that do **not** correspond with the ground-truth measure ...

Conclusion: VAE needs modifications to correctly handle manifolds

\*Some additional technical conditions apply

# One Candidate Solution: More Complex, Non-Gaussian Encoders

- A variety of non-Gaussian decoders have been proposed based on invertible flows (Part I) and related.

[Burda et al., 2015; Kingma et al., 2016; Rezende & Mohamed, 2016; van den Berg et al., 2018]

- This can improve non-negative likelihood (NLL) scores on test data:

Test set performance on the CIFAR-10 data.				
	$K = 0$	$K = 2$	$K = 5$	$K = 10$
$-\ln p(\mathbf{x})$	-293.7	-308.6	-317.9	-320.7

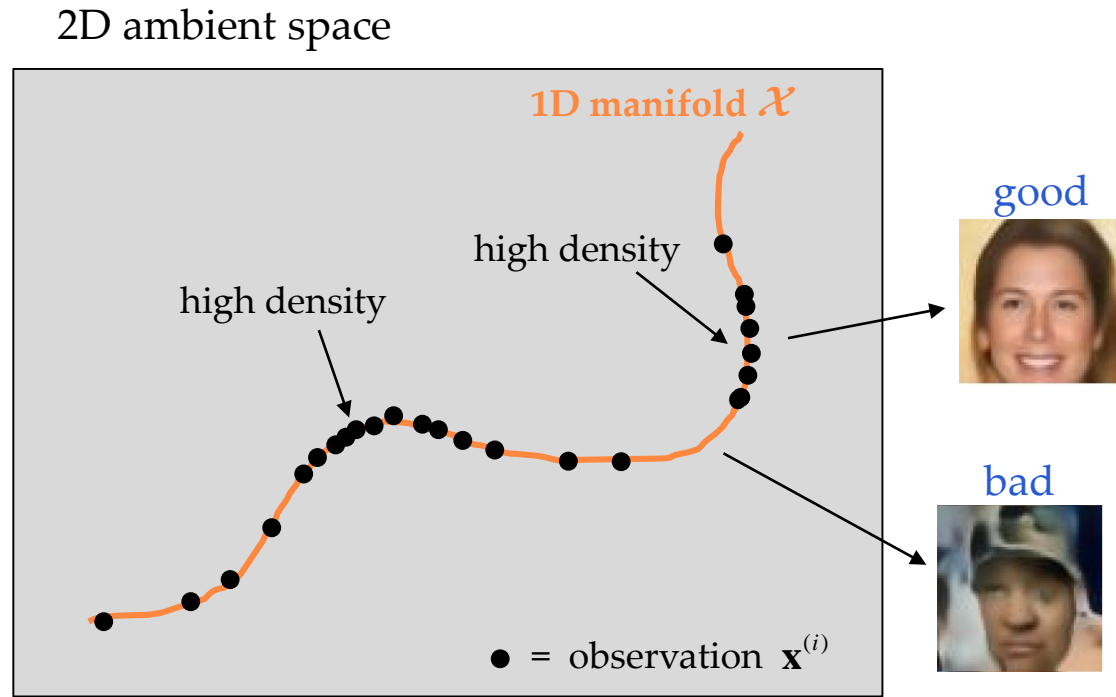
} flow length

[Rezende & Mohamed, 2016]

- Weaknesses:
  - 1) Does not solve the non-uniqueness issue with low-dim manifolds.
  - 2) Has not yet shown quantitative improvement generating new samples (... this is of course subject to change).
- Similar conclusions for non-Gaussian VAE latent priors

[(Tomczak & Welling, 2018; Zhao et al., 2018)]

# Potentially Misleading NLL Scores



Can have  $-\sum_i \log p_\theta(\mathbf{x}) \rightarrow -\infty$  (infinite density)  
with just a uniform measure on  $\mathcal{X}$  and  $\Pr(\mathbf{x} \notin \mathcal{X}) = 0$

---

But samples drawn from the low-density manifold regions might be bad ...

**NLL scores need not correlate with sample quality**

# Helpful Alternative Viewpoint

□ Fix:  $\Sigma_z[\mathbf{x}, \boldsymbol{\varphi}] = \mathbf{0}$ ,  $\Sigma_x[\mathbf{z}, \boldsymbol{\theta}] = \mathbf{I}$

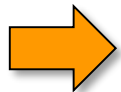
□ VAE energy collapses to a simple deterministic, induced AE:

$$L_{\text{AE}}(\boldsymbol{\theta}, \boldsymbol{\varphi}) \triangleq \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \underbrace{\boldsymbol{\mu}_x[\mathbf{z}^{(i)}, \boldsymbol{\theta}]}_{\text{decoder}} \right\|_2^2, \quad \text{s.t. } \mathbf{z}^{(i)} = \underbrace{\boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\varphi}]}_{\text{encoder}}$$

□ Compute:  $\boldsymbol{\theta}_*, \boldsymbol{\varphi}_* = \arg \min_{\boldsymbol{\theta}, \boldsymbol{\varphi}} L_{\text{AE}}(\boldsymbol{\theta}, \boldsymbol{\varphi})$

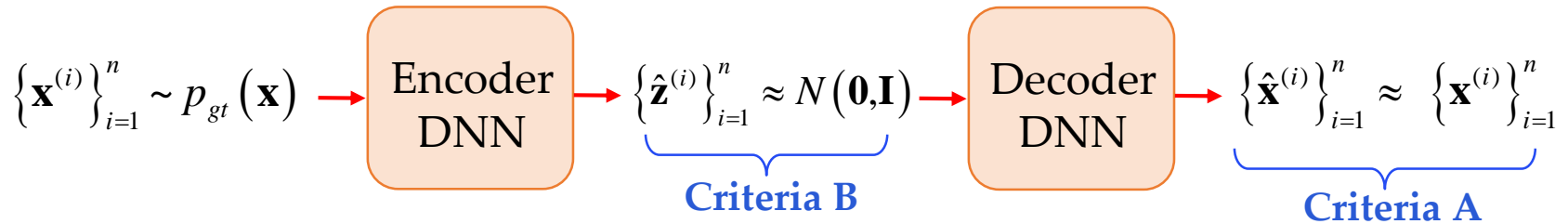
□ Collect corresponding latent variables:  $\{\mathbf{z}^{(i)}\}_{i=1}^n$ ,  $\mathbf{z}^{(i)} = \boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\varphi}_*]$

□ **Hypothetical:** Suppose  $\underbrace{L_{\text{AE}}(\boldsymbol{\theta}_*, \boldsymbol{\varphi}_*)}_{\text{Criteria A:}} \approx 0$  and  $\underbrace{\{\mathbf{z}^{(i)}\}_{i=1}^n}_{\text{Criteria B:}} \approx N(\mathbf{0}, \mathbf{I})$   
Good reconstruction of training data (like VAE from Part III)      Approximation to some known latent distribution

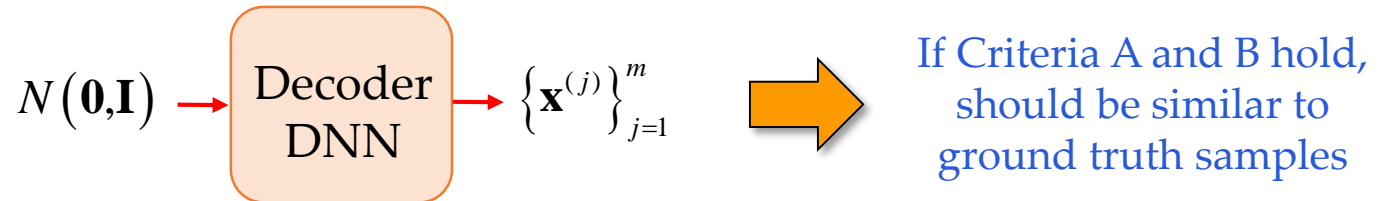


Can in principle apply an AE for generating new samples ...

# Illustration of AE Required Criteria



- Could generate new samples via:



- In practice, an AE can satisfy Criteria A, but will generally **not** satisfy Criteria B ...

## Practical workaround:

Can penalize some measure of the distance between samples  $\{\hat{\mathbf{z}}^{(i)}\}_{i=1}^n$  and  $N(\mathbf{0}, \mathbf{I})$

# Generic Form of AE-Based Generative Model

Enhanced AE energy:

$$L_{\text{AE}^+}(\boldsymbol{\theta}, \boldsymbol{\varphi}) \triangleq \underbrace{\sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_x(\mathbf{z}^{(i)}, \boldsymbol{\theta}) \right\|_2^2}_{\text{data fit term}} + \lambda \Delta \left[ \underbrace{\left\{ \mathbf{z}^{(i)} \right\}_{i=1}^n, N(\mathbf{0}, \mathbf{I})}_{\text{penalty favors latent samples "similar" to standardized Gaussian}} \right], \quad \text{s.t. } \mathbf{z}^{(i)} = \boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\varphi}], \quad \forall i$$

Candidate penalties based on Wasserstein distance measures



Wasserstein AE (WAE)

Two main variants incorporate:

- Maximum mean discrepancy (MMD)
- Generative adversarial network (GAN)



WAE-MMD, WAE-GAN

(Note: There also exists stochastic versions of the WAE encoder, but empirical results are not available)



# WAE Results

quantitative measure of  
perceptual quality; lower is better



Algorithm	FID
VAE	63
WAE-MMD	55
WAE-GAN	42

} significant improvement over the VAE


WAE-MMD generated samples:



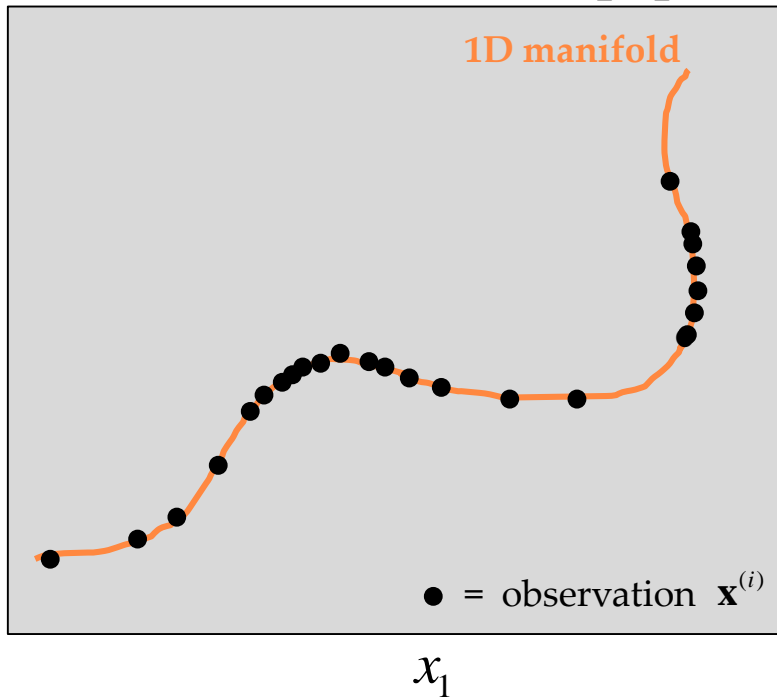
# Potential WAE Limitations

- Must tune trade-off parameter  $\lambda$
- If  $\dim(\mathbf{z})$  is too small  large reconstruction error  
(fails Criteria A)
- If  $\dim(\mathbf{z})$  is too high  large distribution mismatch  
(fails Criteria B)

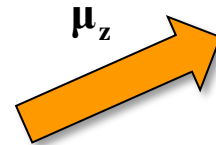
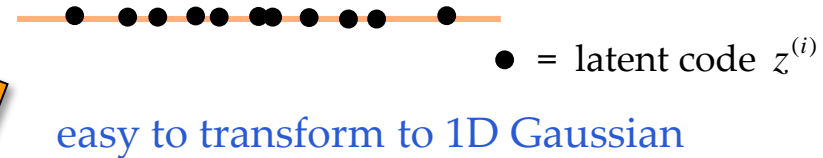
# The Problem of Excess Latent Dimensions

$\mathbf{z}^{(i)} = \boldsymbol{\mu}_z[\mathbf{x}^{(i)}, \boldsymbol{\phi}], \forall i$   deterministic encoder mapping between  $\mathbf{x}$  and  $\mathbf{z}$  space

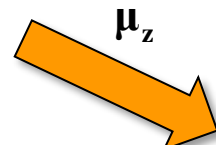
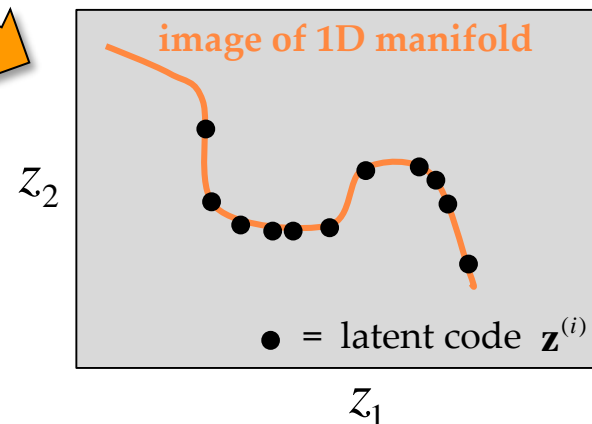
2D observation space,  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



1D latent space,  $\mathbf{z} = z$   
(optimal)



2D latent space,  $\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$   
(suboptimal)



difficult to transform samples on 1D manifold to a 2D Gaussian

# Returning to the VAE ...

## Critical Questions:

- ❑ How does the VAE behave w.r.t. Criteria A (perfect reconstructions) and B (latent space distribution match)?
- ❑ And can we use this information to make improvements?

# Perfect VAE Reconstructions (Criteria A)

Recall from Part III:

## Theorem (Reconstruction Invariance):

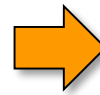
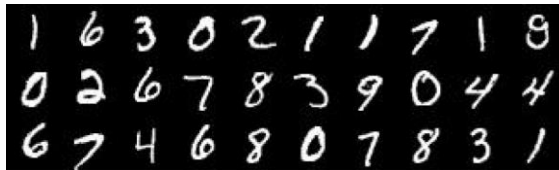
Under some technical conditions, any VAE global optimum  $\{\theta_*, \phi_*\}$  is such that  $\gamma \rightarrow 0$  and reconstructions are exact:

$$\mu_x\left(\mu_z\left[\mathbf{x}^{(i)}, \phi_*\right] + \Sigma_z^{1/2}\left[\mathbf{x}^{(i)}, \phi_*\right]\boldsymbol{\varepsilon}, \theta_*\right) = \mu_x\left(\mu_z\left[\mathbf{x}^{(i)}, \phi_*\right], \theta_*\right) = \mathbf{x}^{(i)}, \quad \forall \boldsymbol{\varepsilon}, \forall i$$

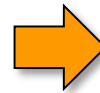
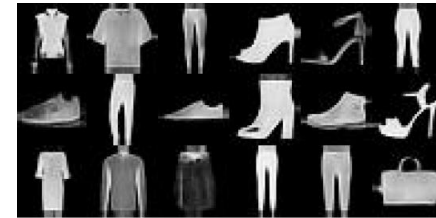
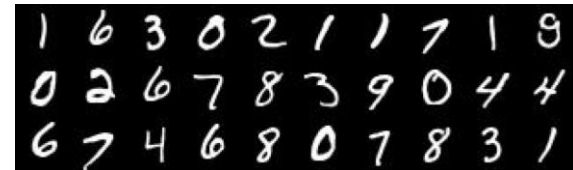
**Key (rephrased) Conclusion:** At global minimum, encoder randomness will not impact perfect reconstructions  VAE can satisfy Criteria A

# Example Reconstructions

Ground Truth Samples



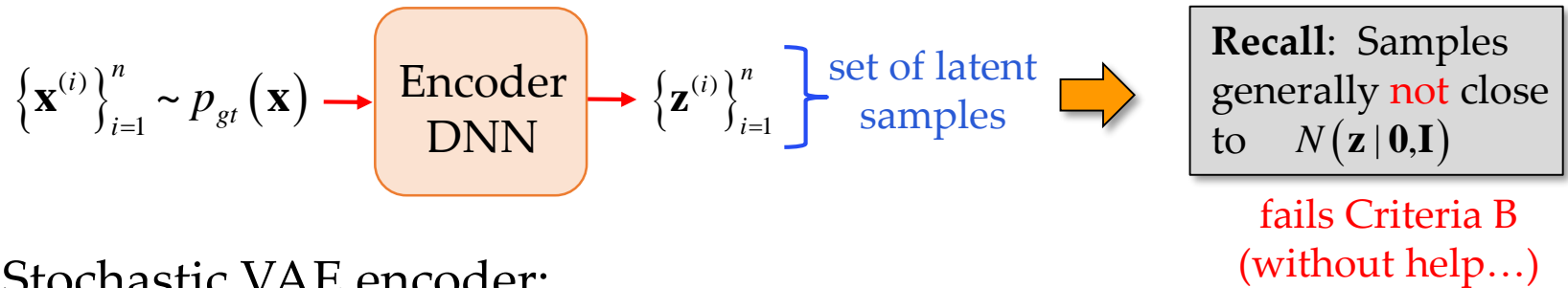
VAE Reconstructions



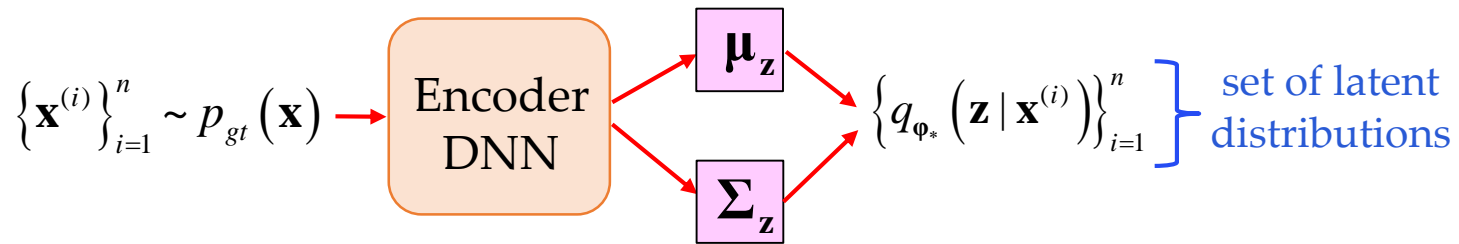
So poor VAE performance may be related to Criteria B

# Addressing the VAE Latent Space (Criteria B)

- Deterministic AE encoder:



- Stochastic VAE encoder:

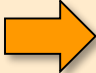
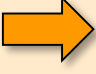


- Aggregated distribution of VAE latent space:

$$\underbrace{q_{\phi_*}(\mathbf{z})}_{\text{aggregated posterior}} \triangleq \int q_{\phi_*}(\mathbf{z} | \mathbf{x}) p_{gt}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n q_{\phi_*}(\mathbf{z} | \mathbf{x}^{(i)})$$

For generating good samples, should be close to  $N(\mathbf{z} | \mathbf{0}, \mathbf{I})$  } VAE version of Criteria B

# Properties of VAE Aggregated Posterior

- When data lies on a manifold ( $r < d$ ), at global minimum can have  $q_{\phi_*}(\mathbf{z}) \neq N(\mathbf{z}|\mathbf{0},\mathbf{I})$   fails Criteria B
- But under reasonable assumptions, VAE aggregated posterior  $q_{\phi_*}(\mathbf{z})$  will satisfy conditions of Exact Density Recovery Theorem. (Note: samples from an AE generally will **not**)
- This means that a **second** VAE could be trained to learn  $q_{\phi_*}(\mathbf{z})$ .  
 **implicitly addresses Criteria B**



# Matching the VAE Aggregated Posterior

- From Exact Density Recovery Theorem, when  $r = d$  we have

$$\text{KL}[q_{\phi_*}(\mathbf{z}|\mathbf{x}) \| p_{\theta_*}(\mathbf{z}|\mathbf{x})] = 0 \quad \text{and} \quad p_{\theta_*}(\mathbf{x}) = p_{gt}(\mathbf{x})$$

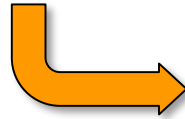
at any optimal solution, provided  $\kappa \geq r$ .

- This implies that:

$$q_{\phi_*}(\mathbf{z}) \triangleq \int q_{\phi_*}(\mathbf{z}|\mathbf{x}) p_{gt}(\mathbf{x}) d\mathbf{x} = \int p_{\theta_*}(\mathbf{z}|\mathbf{x}) p_{\theta_*}(\mathbf{x}) d\mathbf{x} = \int p_{\theta_*}(\mathbf{x}|\mathbf{z}) N(\mathbf{z}|\mathbf{0},\mathbf{I}) d\mathbf{x} = N(\mathbf{z}|\mathbf{0},\mathbf{I})$$

perfect match!

- But when the data lie on a manifold (i.e.,  $r < d$ ), this no longer need be the case, i.e.,  $q_{\phi_*}(\mathbf{z}) \neq N(\mathbf{z}|\mathbf{0},\mathbf{I})$



$\begin{aligned} \mathbf{z}^{(j)} &\sim q_{\phi_*}(\mathbf{z}) \\ \mathbf{x}^{(j)} &\sim p_{\theta_*}(\mathbf{x} \mathbf{z}^{(j)}) \end{aligned} \neq \begin{aligned} \mathbf{z}^{(j)} &\sim N(\mathbf{z} \mathbf{0},\mathbf{I}) \\ \mathbf{x}^{(j)} &\sim p_{\theta_*}(\mathbf{x} \mathbf{z}^{(j)}) \end{aligned}$
---

generates training data, but is intractable

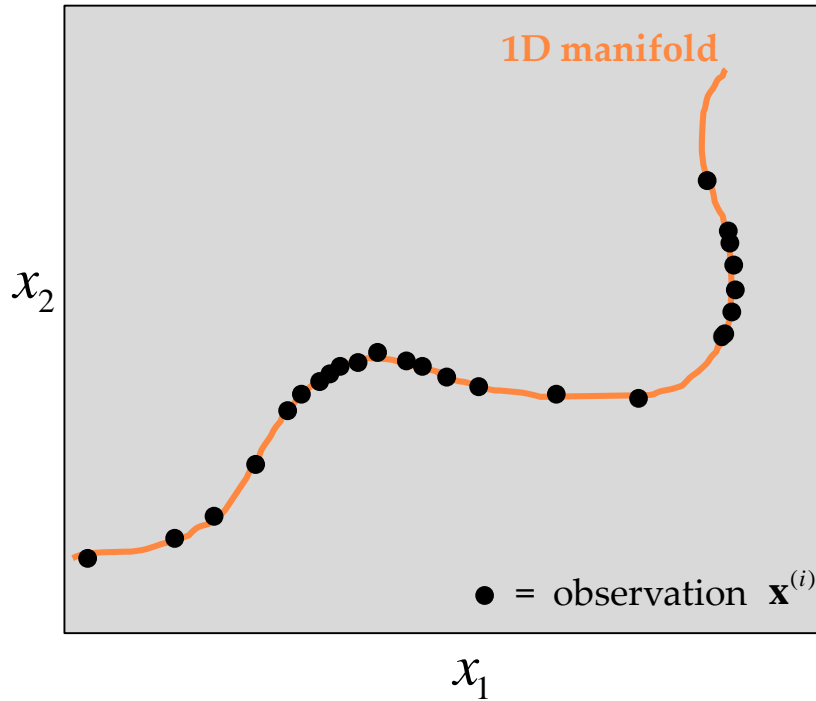
?

- But intrinsic VAE properties suggest a practical solution ...

# 2D Illustration

Optimal VAE encoder  
mapping to 2D latent space

2D observation space,  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

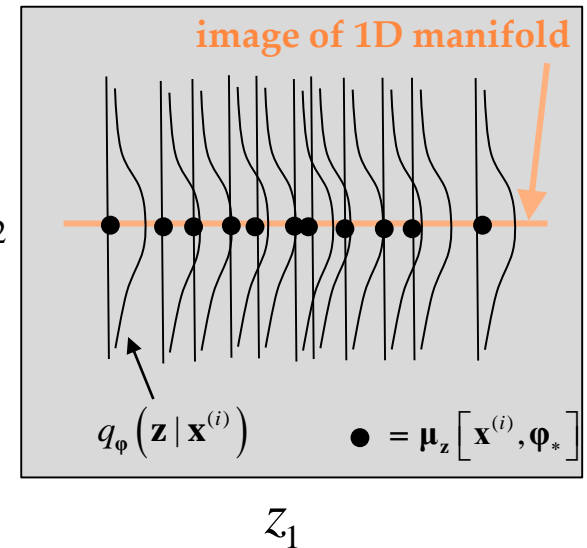


$$\mu_z[\mathbf{x}^{(i)}, \phi_*] \rightarrow \begin{bmatrix} \mu_z[\mathbf{x}^{(i)}, \phi_*]_1 \\ 0 \end{bmatrix}$$

$$\Sigma_z[\mathbf{x}^{(i)}, \phi_*] \rightarrow \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



2D VAE latent space,  $\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$



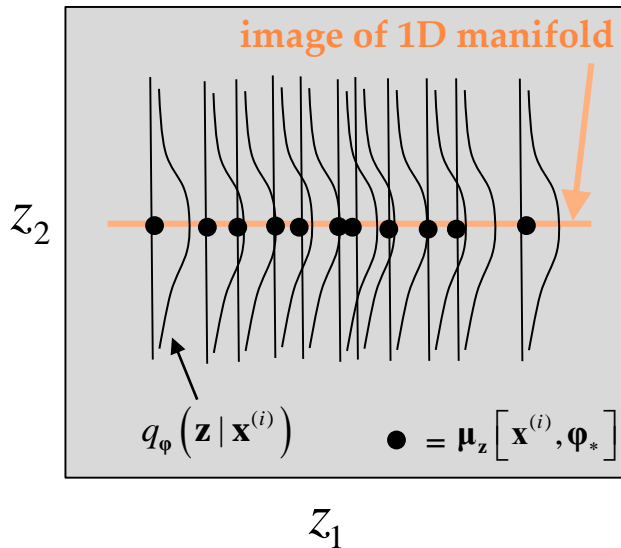
$$\underbrace{q_{\phi_*}(\mathbf{z})}_{\text{aggregated posterior}} \triangleq \int q_{\phi_*}(\mathbf{z} | \mathbf{x}) p_{gt}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n q_{\phi_*}(\mathbf{z} | \mathbf{x}^{(i)})$$

Aggregated posterior does not lie on a low-dim manifold as with deterministic AE

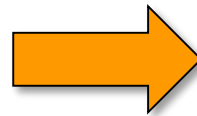
# Comparing Aggregated Posterior Samples

$$q_{\phi_*}(\mathbf{z}) \triangleq \int q_{\phi_*}(\mathbf{z}|\mathbf{x}) p_{gt}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n q_{\phi_*}(\mathbf{z}|\mathbf{x}^{(i)})$$

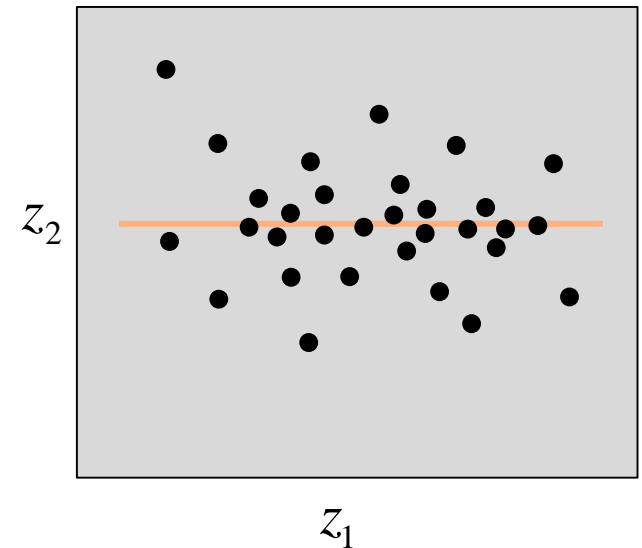
2D VAE latent space,  $\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$



$$\mathbf{z}^{(j)} \sim q_{\phi_*}(\mathbf{z})$$



Aggregated posterior samples:  $\{\mathbf{z}^{(j)}\}_{j=1}^m$



## Remarks:

- Still no guarantee that the aggregated posterior will be close to  $N(\mathbf{0}, \mathbf{I})$
- But samples will **not** lie on a low-dimensional manifold
- The VAE decoder “fills out” unnecessary dimensions with random noise (Part III)
- This leads to a simple 2-stage VAE enhancement based on Exact Density Recovery Theorem from earlier ...

# Two-Stage VAE Strategy

$$\mathbf{X} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^n, \quad \mathbf{x}^{(i)} \in \mathcal{X} \subset \mathbb{R}^d, \quad r\text{-dim manifold}, \quad r < d$$

- Choose  $\dim(\mathbf{z}) = \kappa$  sufficiently large, ensure  $\kappa \geq r$  (do not need exact value)

- Solve via SGD:  $\boldsymbol{\theta}_*, \boldsymbol{\varphi}_* = \arg \min_{\boldsymbol{\theta}, \boldsymbol{\varphi}} L(\boldsymbol{\theta}, \boldsymbol{\varphi}) \quad \Rightarrow \quad \text{first-stage VAE}$

- Form aggregated posterior approximation:  $q_{\boldsymbol{\varphi}_*}(\mathbf{z}) \approx \frac{1}{n} \sum_{i=1}^n q_{\boldsymbol{\varphi}_*}(\mathbf{z} | \mathbf{x}^{(i)})$

- Samples from this approximation form new data set:

$$\mathbf{Z} = \left\{ \mathbf{z}^{(j)} \right\}_{j=1}^m, \quad \mathbf{z}^{(j)} \sim \frac{1}{n} \sum_{i=1}^n q_{\boldsymbol{\varphi}_*}(\mathbf{z} | \mathbf{x}^{(i)}) \quad \left. \vphantom{\mathbf{z}^{(j)}} \right\} \text{latent codes associated with training data; no manifold structure}$$

- This is regime where Exact Density Recovery Theorem applies

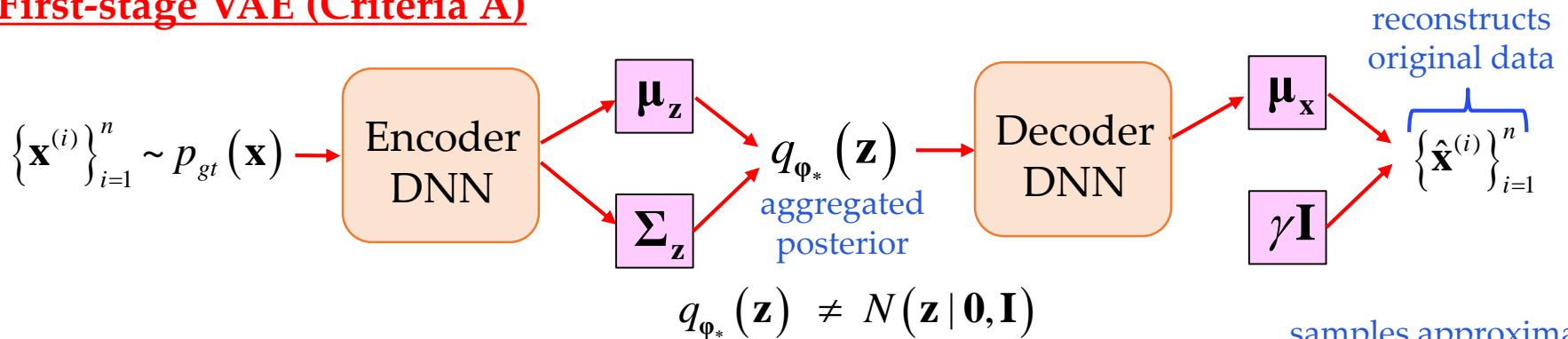
- Train a second VAE on data  $\mathbf{Z}$  with latent code  $\mathbf{u}$ , and  $\dim(\mathbf{u}) = \kappa$

$$\boldsymbol{\theta}_{2*}, \boldsymbol{\varphi}_{2*} = \arg \min_{\boldsymbol{\theta}_2, \boldsymbol{\varphi}_2} L(\boldsymbol{\theta}_2, \boldsymbol{\varphi}_2) \quad \Rightarrow \quad \text{second-stage VAE (much smaller)}$$

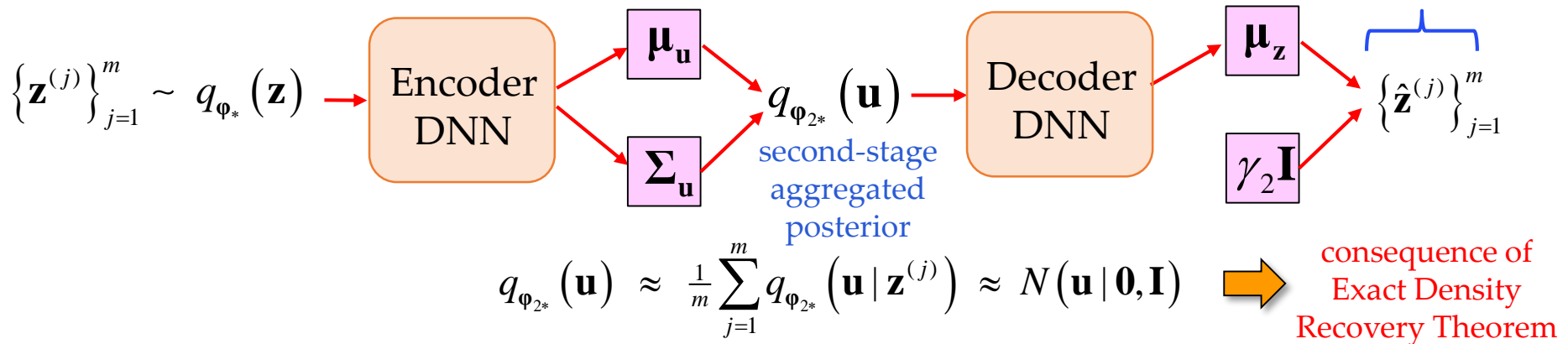
- By design, this VAE will (asymptotically) learn the exact aggregated posterior from the first-stage VAE

# Two-Stage VAE Visualization

## First-stage VAE (Criteria A)



## Second-stage VAE (Criteria B)

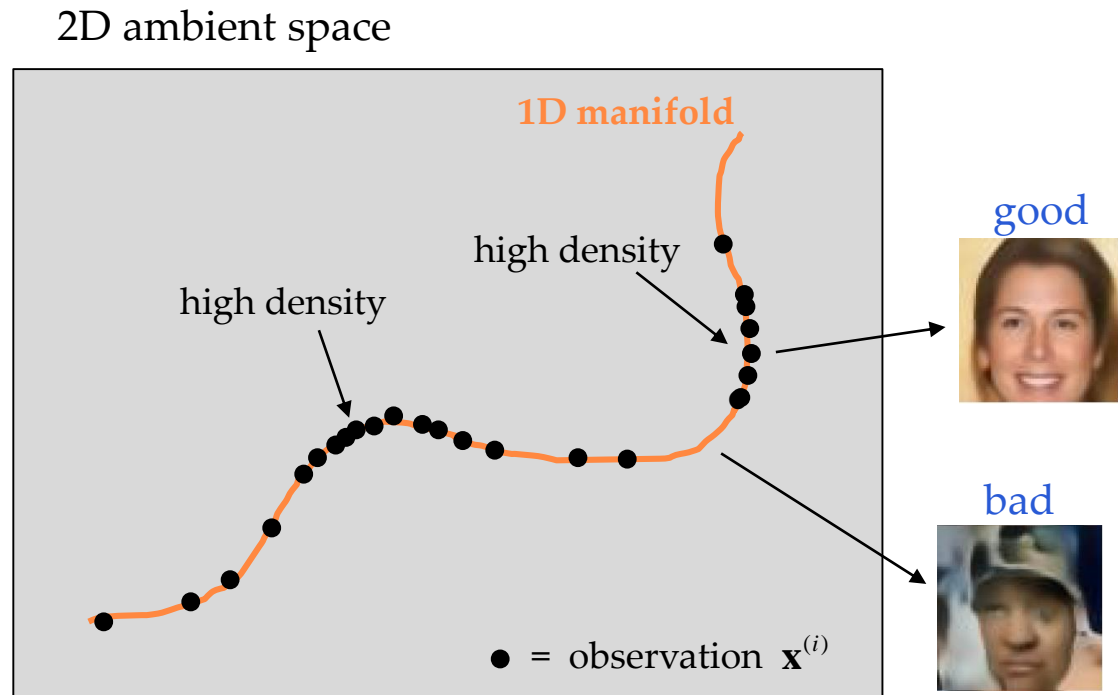


Generating new samples is now trivial:

$$\mathbf{u}^{new} \sim N(\mathbf{u} | \mathbf{0}, \mathbf{I}), \quad \mathbf{z}^{new} \sim p_{\theta_{2*}}(\mathbf{z} | \mathbf{u}^{new}), \quad \mathbf{x}^{new} \sim p_{\theta_*}(\mathbf{x} | \mathbf{z}^{new})$$

$$\underbrace{\hspace{15em}}_{\mathbf{z}^{new} \sim q_{\phi_*}(\mathbf{z})}$$

# Two-Stage VAE Intuition



- ❑ First-stage VAE learns manifold model by efficiently reconstructing samples (analogous to Criteria A).
- ❑ Second-stage VAE learns distribution **within** the manifold (analogous to Criteria B).
- ❑ Note: Joint training does **not** work in this context.

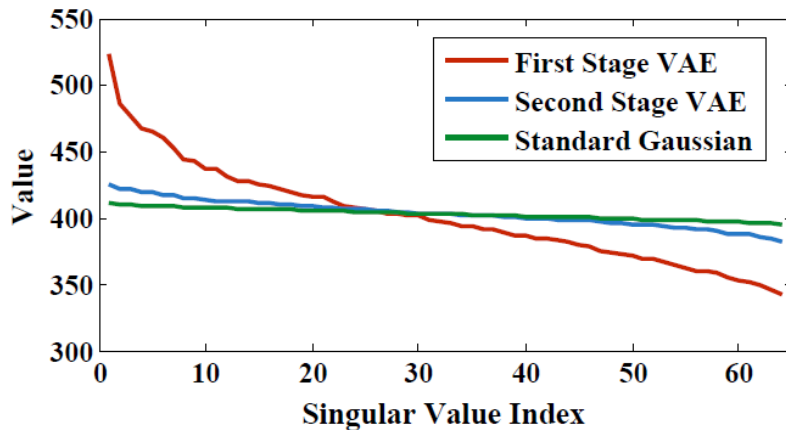
# Aggregated Posterior Comparisons

First-stage  
latent codes:  $\mathbf{Z} = \{\mathbf{z}^{(j)}\}_{i=1}^m$ ,  $\mathbf{z}^{(j)} \sim q_{\hat{\phi}}(\mathbf{z}) \approx \frac{1}{n} \sum_{i=1}^n q_{\hat{\phi}}(\mathbf{z} | \mathbf{x}^{(i)})$

---

Second-stage  
latent codes:  $\mathbf{U} = \{\mathbf{u}^{(j)}\}_{i=1}^m$ ,  $\mathbf{u}^{(j)} \sim q_{\hat{\phi}_2}(\mathbf{u}) \approx \frac{1}{m} \sum_{j=1}^m q_{\hat{\phi}_2}(\mathbf{u} | \mathbf{z}^{(j)})$

## Singular Value Comparison



## MMD from ideal $N(0, \mathbf{I})$

	First Stage	Second Stage
MNIST	2.85	0.43
Fashion	1.37	0.40
Cifar10	1.08	0.00
CelabA	7.42	0.29



low values; close to Gaussian

# Two-Stage VAE Results

quantitative measure of perceptual quality; lower is better

## Averaged **FID** Score Comparisons

Neutral testing conditions from [Lucic et al., 2018]

		MNIST	Fashion	CIFAR-10	CelebA
optimized, data-dependent settings	MM GAN	$9.8 \pm 0.9$	$29.6 \pm 1.6$	$72.7 \pm 3.6$	$65.6 \pm 4.2$
	NS GAN	$6.8 \pm 0.5$	$26.5 \pm 1.6$	$58.5 \pm 1.9$	$55.0 \pm 3.3$
	LSGAN	$7.8 \pm 0.6$	$30.7 \pm 2.2$	$87.1 \pm 47.5$	$53.9 \pm 2.8$
	WGAN	$6.7 \pm 0.4$	$21.5 \pm 1.6$	$55.2 \pm 2.3$	$41.3 \pm 2.0$
	WGAN GP	$20.3 \pm 5.0$	$24.5 \pm 2.1$	$55.8 \pm 0.9$	$30.3 \pm 1.0$
	DRAGAN	$7.6 \pm 0.4$	$27.7 \pm 1.2$	$69.8 \pm 2.0$	$42.3 \pm 3.0$
	BEGAN	$13.1 \pm 1.0$	$22.9 \pm 0.9$	$71.4 \pm 1.6$	$38.9 \pm 0.9$
default settings	Best GAN	$\sim 10$	$\sim 32$	$\sim 70$	$\sim 49$
	VAE (cross-entr.)	$16.6 \pm 0.4$	$43.6 \pm 0.7$	$106.0 \pm 1.0$	$53.3 \pm 0.6$
	VAE (fixed $\gamma$ )	$52.0 \pm 0.6$	$84.6 \pm 0.9$	$160.5 \pm 1.1$	$55.9 \pm 0.6$
	VAE (learned $\gamma$ )	$54.5 \pm 1.0$	$60.0 \pm 1.1$	$76.7 \pm 0.8$	$60.5 \pm 0.6$
	VAE + Flow	$54.8 \pm 2.8$	$62.1 \pm 1.6$	$81.2 \pm 2.0$	$65.7 \pm 2.8$
	WAE-MMD	$115.0 \pm 1.1$	$101.7 \pm 0.8$	$80.9 \pm 0.4$	$62.9 \pm 0.8$
	<b>2-Stage VAE</b>	$12.6 \pm 1.5$	$29.3 \pm 1.0$	$72.9 \pm 0.9$	$44.4 \pm 0.7$

similar to GANs,  
no tuning

WAE testing conditions from [Tolstikhin et al., 2018]

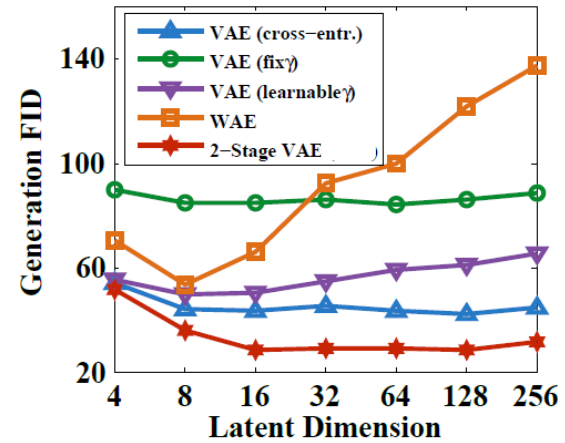
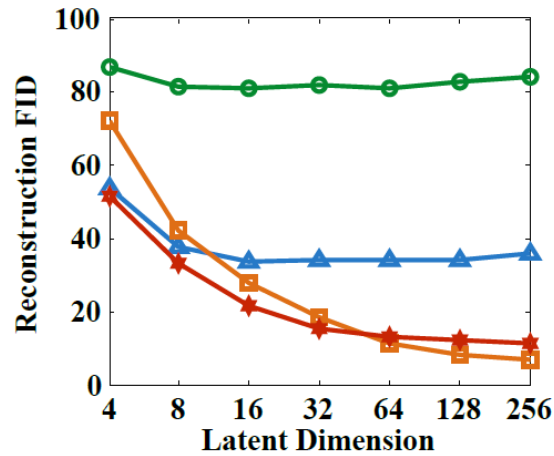
	VAE	WAE-MMD	WAE-GAN	2-Stage VAE
CelebA FID	63	55	42	<b>34</b>

improvement  
over WAE

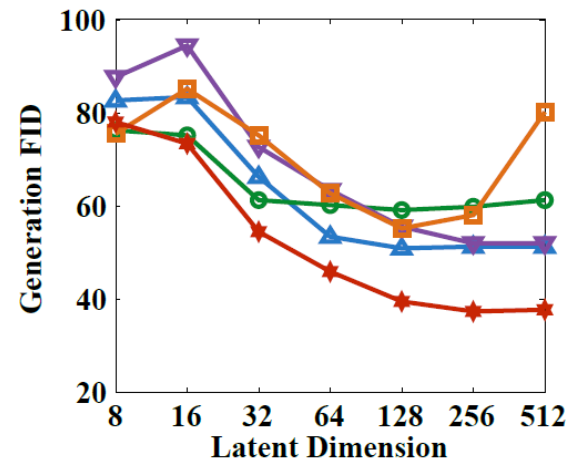
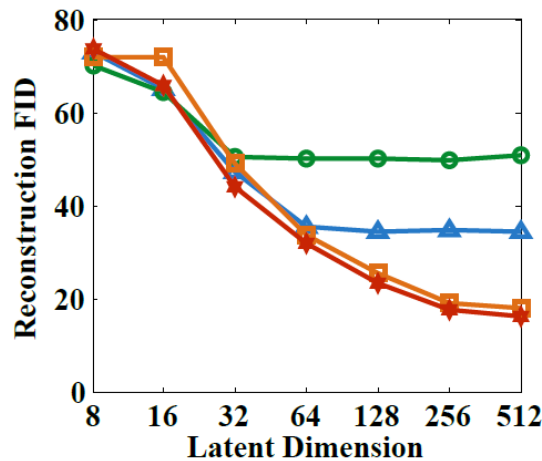


# Robustness to the Latent Space Dimension

Fashion  
MNIST

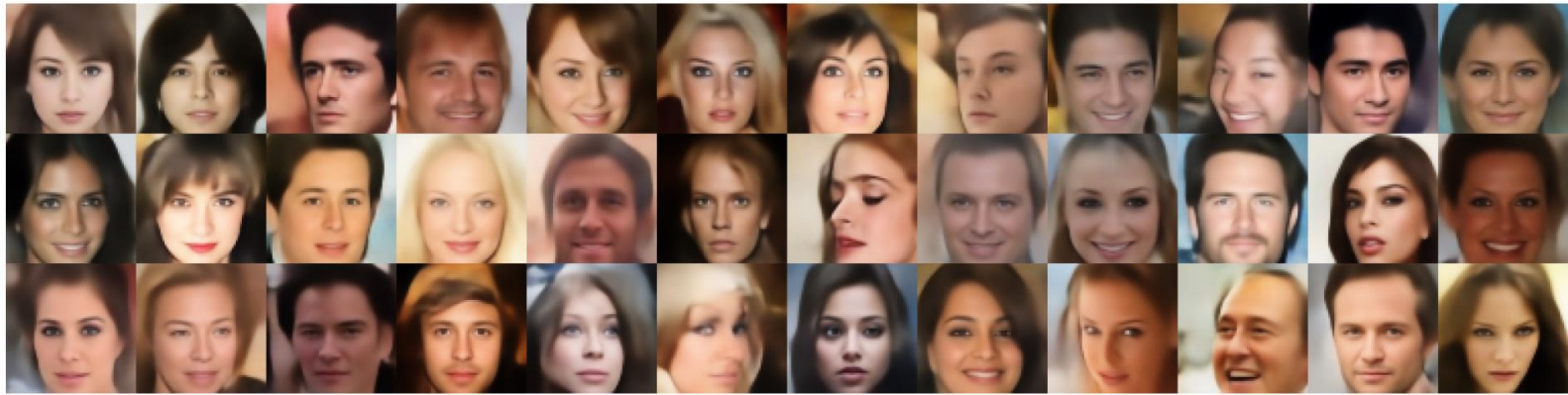


CelebA



# Comparison of Generated CelebA Samples

## Two-Stage VAE



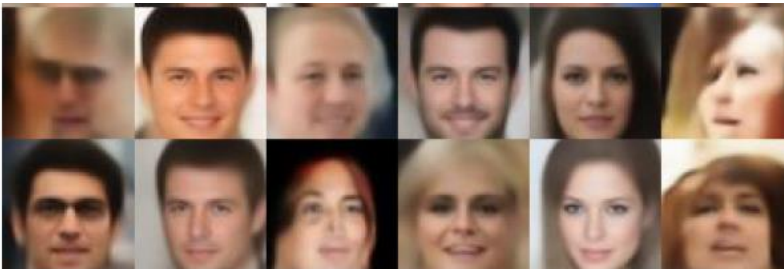
## Single-Stage VAE (learned $\gamma$ )

good \ clear  
reconstructions,  
poor new samples



## Single-Stage VAE (fixed $\gamma$ )

bad \ blurry  
reconstructions,  
poor new samples



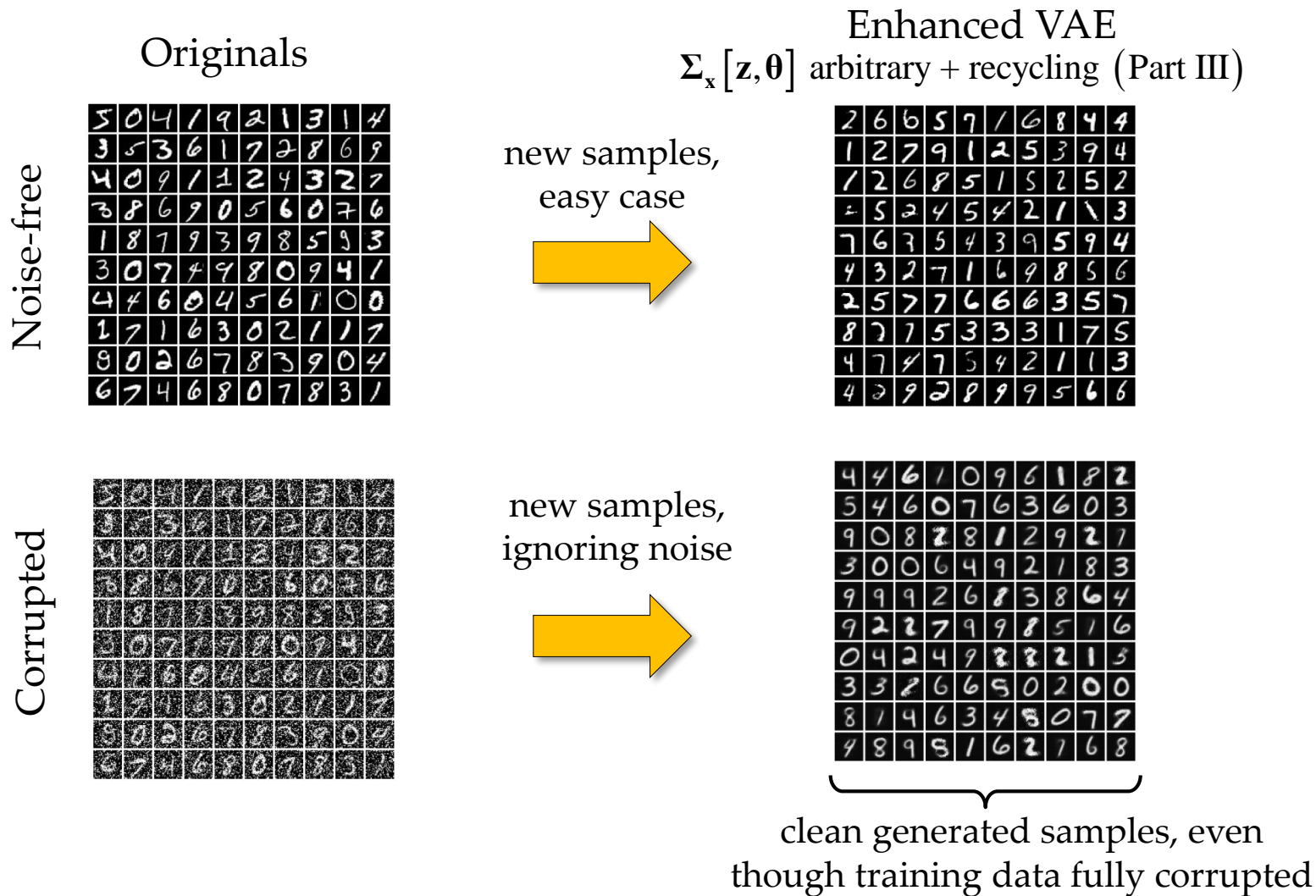
## WGAN-GP [Gulrajani et al., 2017]



<https://github.com/LynnHo/WGAN-GP-DRAGAN-Celeba-Pytorch>

# MNIST Example with Corruptions

Note: MNIST data is much simpler than CelebA, but with corruptions it is challenging to generate new clean samples



# Summary

- ❑ Standard VAE can reconstruct data lying on a low-dimensional manifold (first-stage VAE).
- ❑ But generated samples may not resemble training data.
- ❑ Fortunately, the distribution of the VAE latent codes can be successfully modeled and sampled from (second-stage VAE).
- ❑ Combined stages can produce more realistic samples, **comparable to many GANs** (w/ same neutral architecture).
- ❑ But two-stage model retains original VAE advantages (and additional complexity is minimal, second-stage can be small).
- ❑ Alternative VAE-inspired approaches like the WAE can also produce good results (but may be more sensitive to latent dimensions).

**Questions?**

# Part V: Practical Usage Issues and Examples

**Note:** Updated version of slides available at <http://www.davidwipf.com/>

# Outline

- Cases of over- and under-regularization
- Identifiability of semantically-meaningful latent factors
- Practical applications via Conditional VAEs

# Over-Regularized/Degenerate VAE Local Solutions

VAE Objective:

$$L(\theta, \phi) = \sum_i \left\{ \underbrace{\text{KL} \left[ q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) \| N(\mathbf{z} | \mathbf{0}, \mathbf{I}) \right]}_{\text{KL term has trivial minimum, only requires parameters of last encoder layer}} - \underbrace{\mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right]}_{\text{minimizing data term requires complex coordination of all parameters in encoder/decoder networks (hard)}} \right\}$$

KL term has trivial minimum, only requires parameters of last encoder layer

minimizing data term requires complex coordination of all parameters in encoder/decoder networks (hard)

$$\sum_{j=1}^K \left\{ \mu_z[\mathbf{x}^{(i)}, \phi]_j^2 + \sigma_z^2[\mathbf{x}^{(i)}, \phi]_j - \log \left( \sigma_z^2[\mathbf{x}^{(i)}, \phi]_j \right) \right\}$$

trivial  
minimum



$$\begin{aligned} \mu_z[\mathbf{x}^{(i)}, \phi]_j &\rightarrow 0 \\ \sigma_z^2[\mathbf{x}^{(i)}, \phi]_j &\rightarrow 1 \end{aligned}$$

Potential for convergence to bad, overregularized (local) solutions

Candidate workarounds:

- ❑ KL warm-start [Bowman et al., 2015; Sønderby et al., 2016]
- ❑ Skip connections [Cai et al., 2017; Dieng et al., 2018]
- ❑ Ladder networks [Sønderby et al., 2016; Maaløe et al., 2019]



# Under-Regularized VAE Global Optima

- In principle, VAE encoder can be arbitrarily complex; this just tightens the original upper bound

$$-\sum_i \log p_{\theta}(\mathbf{x}^{(i)}) \leq \sum_i \left\{ \underbrace{\text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) \| p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})]}_{\rightarrow 0 \text{ w/ complex encoder}} - \log p_{\theta}(\mathbf{x}^{(i)}) \right\}$$

- Likewise, decoder covariance can be arbitrarily complex to learn outlier locations (Part III).
- But the **decoder mean network** is more subtle ...

Problem: While the VAE cost does penalize excessive dimensions of  $\mathbf{z}$  (Part III), it cannot prevent overfitting from excessive **depth**.

## Theorem

Even with  $\dim(\mathbf{z}) = 1$ , VAE cost can be globally optimized by solution that just memorizes the training data if the decoder mean is too complex.

# Outline

- ❑ Cases of over- and under-regularization
- ❑ Identifiability of semantically-meaningful latent factors
- ❑ Practical applications via Conditional VAEs

# Interpretability of the VAE Latent Space

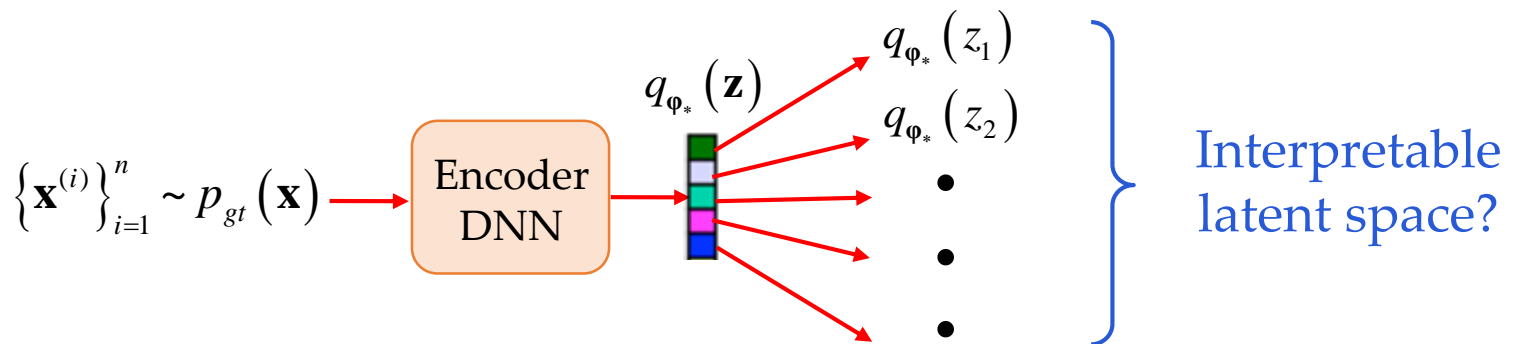
If VAE training is successful, then:

$$q_{\phi_*}(\mathbf{z}) \triangleq \int q_{\phi_*}(\mathbf{z}|\mathbf{x}) p_{gt}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^n q_{\phi_*}(\mathbf{z}|\mathbf{x}^{(i)}) \approx N(\mathbf{z}|\mathbf{0}, \mathbf{I})$$



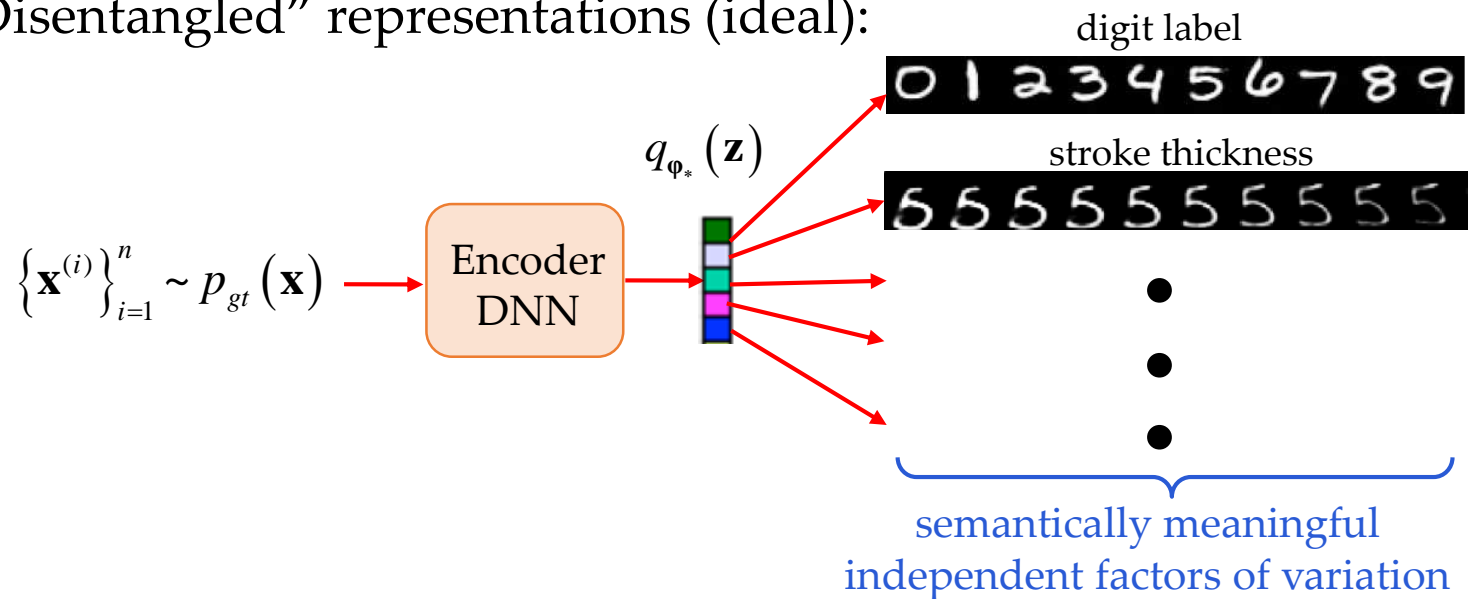
$$q_{\phi_*}(\mathbf{z}) \approx \prod_{j=1}^{\kappa} q_{\phi_*}(z_j)$$

Independent latent factors of variation:

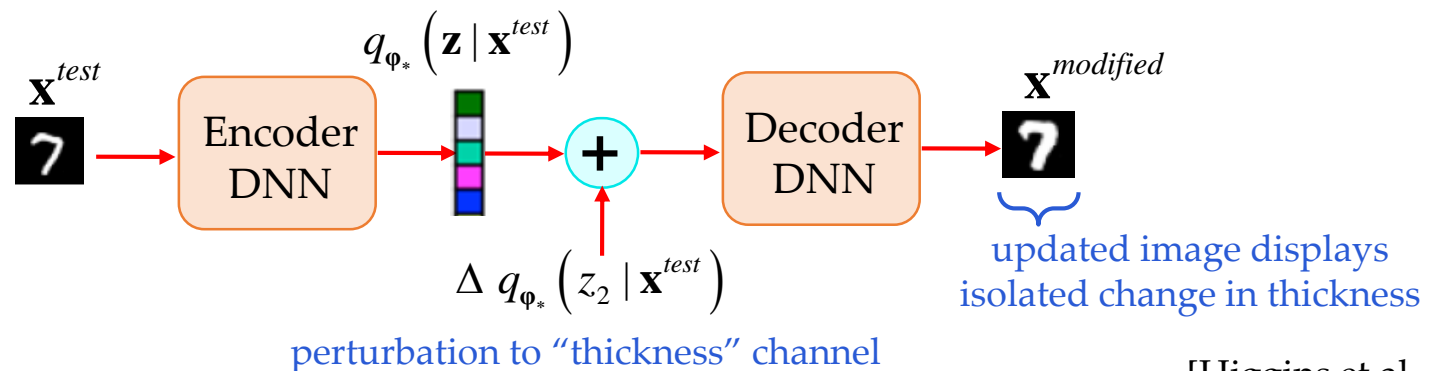


# Independence vs. Semantic Meaning

“Disentangled” representations (ideal):



Useful for numerous computer vision, image processing applications, e.g., photo editing/manipulation:



# Identifiability Issues

Assume the ground-truth generative process:

$$\mathbf{z}^{(i)} \sim p_{gt}(\mathbf{z}), \quad \mathbf{x}^{(i)} = \underbrace{f_{gt}(\mathbf{z}^{(i)})}_{\text{arbitrary deterministic decoder}}, \quad i = 1, \dots, n$$

Also assume a disentangled latent density:

$$p_{gt}(\mathbf{z}) = \prod_j p_{gt}(z_j) \quad \Rightarrow \quad \begin{array}{l} \text{semantically} \\ \text{meaningful factors} \\ \text{of variation} \end{array}$$

**Identifiability Problem:** Exact same samples can be generated using a simple transformed process ...

$$\mathbf{z}^{(i)} \sim p_{gt}(\mathbf{z}) \quad \Rightarrow \quad \tilde{\mathbf{z}}^{(i)} \sim \tilde{p}(\tilde{\mathbf{z}}) = \prod_j \tilde{p}(\tilde{z}_j)$$

new latent factorial distribution by design
mix of previous disentangled factors

$$\tilde{\mathbf{z}}^{(i)} \triangleq D_2 \left[ \mathbf{R} \cdot D_1(\mathbf{z}^{(i)}) \right]$$

transformation to arbitrary marginals  $\nearrow$        $\nwarrow$  rotation       $\nwarrow$  transformation to  $N(\mathbf{z} | \mathbf{0}, \mathbf{I})$

$$\tilde{\mathbf{z}}^{(i)} \sim \tilde{p}(\tilde{\mathbf{z}}) = \prod_j \tilde{p}(\tilde{z}_j)$$

$$\mathbf{x}^{(i)} = \underbrace{f_{gt} \left( D_1^{-1} \left[ \mathbf{R}^T \cdot D_2^{-1}(\tilde{\mathbf{z}}^{(i)}) \right] \right)}_{\text{composite decoder inverts transformation}} = f_{gt}(\mathbf{z}^{(i)})$$

Same samples as before, no unique disentangled representation

# Trivial Discrete Example

Data set of 4 equiprobable images:



Two latent attributes: gender {female, male}, age {young, old}

Candidate 2D latent codes

	Disentangled Representation	Entangled Representation
	$z_1$ $z_2$	$\tilde{z}_1$ $\tilde{z}_2$
$\mathbf{x}_{FY}$	0 0	0 0
$\mathbf{x}_{FO}$	0 1	1 1
$\mathbf{x}_{MY}$	1 0	0 1
$\mathbf{x}_{MO}$	1 1	1 0

one attribute changes      both attributes change

$$p_{gt}(\mathbf{z}) = p_{gt}(z_1) p_{gt}(z_2) \quad \tilde{p}(\tilde{\mathbf{z}}) = \tilde{p}(\tilde{z}_1) \tilde{p}(\tilde{z}_2)$$

both latent distributions factorize  $\Rightarrow$  **not identifiable**

# Workarounds

- Constraints on the ground-truth generative process, e.g.,

$$p_{gt}(\mathbf{z}) = \prod_{j=1}^r p_{gt}(z_j) \neq N(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p_{gt}(\mathbf{x} | \mathbf{z}) = N(\mathbf{x} | \mathbf{W}\mathbf{z}, \gamma\mathbf{I}) \quad \rightarrow \quad \text{linear ground-truth decoder}$$

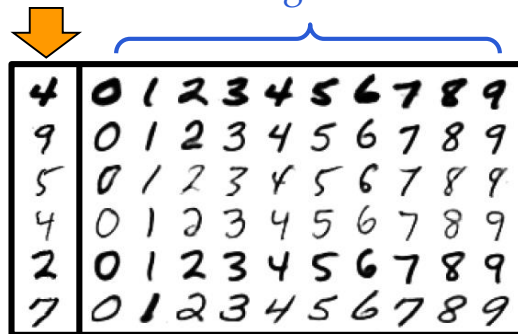
**VAE model:** Use linear decoder mean network and non-Gaussian (possibly parameterized) prior  $p_{\theta}(\mathbf{z})$

Leads to ICA-like model  $\rightarrow$  identifiable up to permutation and scaling

[Hyvärinen et al., 2001]

- Apply some form of weak supervision or semi-supervised learning to resolve ambiguity, e.g.,

test set images      generated images using same style



[Kingma et al., 2014]

# Outline

- Cases of over- and under-regularization
- Identifiability of semantically-meaningful latent factors
- Practical applications via Conditional VAEs

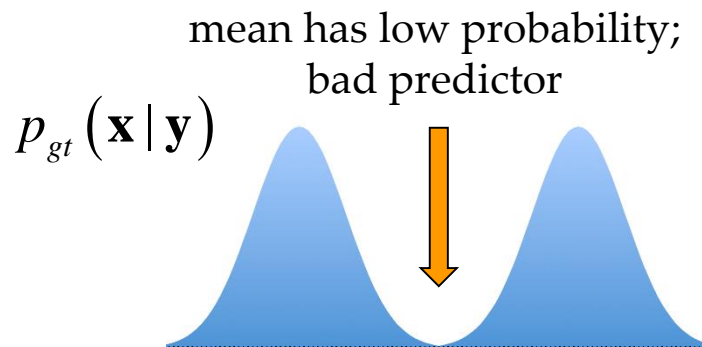


# Conditional VAEs

- Often want a generative model for data conditioned on some variable of interest, e.g.,

$$p_{gt}(\mathbf{x}|\mathbf{y}) = \int p_{gt}(\mathbf{x}|\mathbf{z},\mathbf{y}) \underbrace{p_{gt}(\mathbf{z})}_{\text{independent of } \mathbf{y}} d\mathbf{z}$$

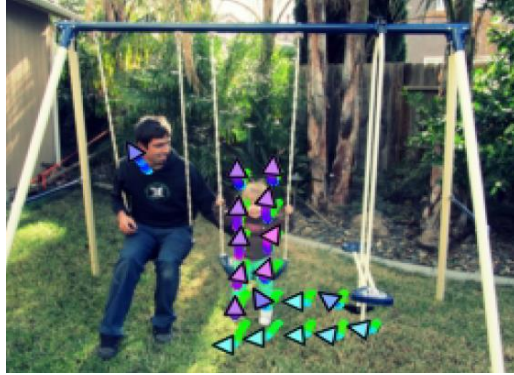
- Basic VAE derivations go through as before, with extra conditioning variable  $\mathbf{y}$ .
- Many applications, e.g., structured output prediction:



# Example Applications

- Forecasting possible motions from static images:

$$p_{gt}(\mathbf{x}|\mathbf{y}), \mathbf{y} = \text{static image}, \mathbf{x} = \text{dense motion trajectory}$$



[Walker et al., 2016]

- Image Captioning:

$$p_{gt}(\mathbf{x}|\mathbf{y}), \mathbf{y} = \text{image}, \mathbf{x} = \text{caption}$$



a woman sitting at a table with a cup of coffee  
a person sitting at a table with a cup of coffee  
a table with two plates of donuts and a cup of coffee  
a woman sitting at a table with a plate of coffee  
a man sitting at a table with a plate of food

[Wang et al., 2017]

# Final Thoughts

- ❑ The VAE represents a natural extension of many existing signal processing, dimensionality reduction tools
- ❑ This is complementary to its role capability as a generative model
- ❑ Many diverse applications, algorithmic variants, extensions ...

**Questions?**