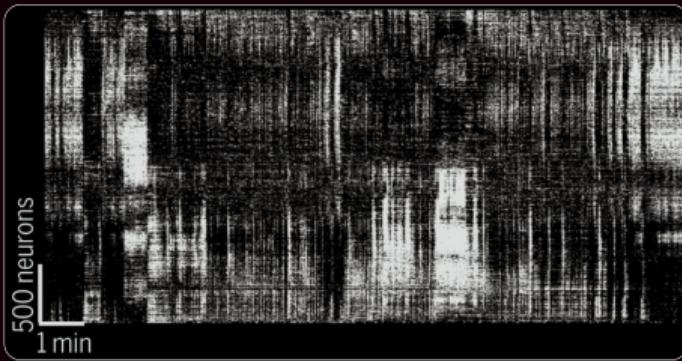
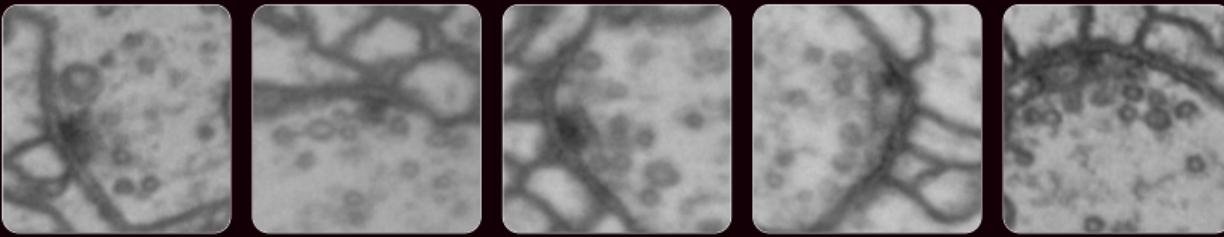


Learning of Disentangled Representations

Jan Funke

Data Analysis Journal Club

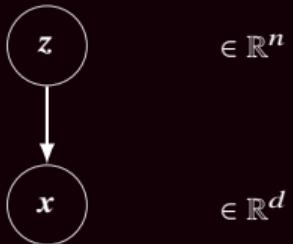
May 13, 2021



Neural recordings from: Stringer & Pachitariu et al., "Spontaneous behaviors drive multidimensional, brainwide activity", Science 2019

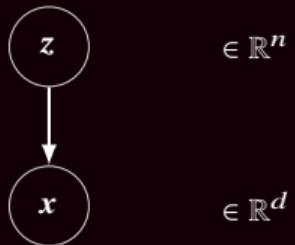
Formalization

latent variable:



Formalization

latent variable:



$\in \mathbb{R}^n$

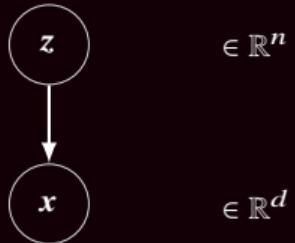
$$p(x, z) = p(x|z)p(z)$$

data:

$\in \mathbb{R}^d$

Formalization

latent variable:



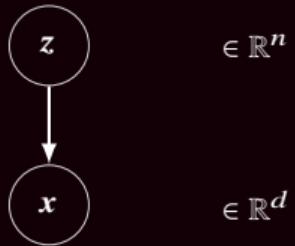
data:

$$p(x, z) = p(x|z)p(z)$$

$$p(x) = \int p(x|z)p(z)dz$$

Formalization

latent variable:



data:

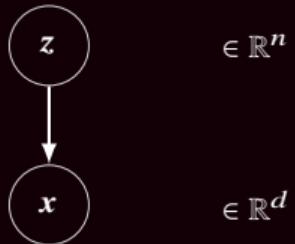
$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$$

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

$$p(\mathbf{z}) = \prod_i p(z_i)$$

Formalization

latent variable:



$$\in \mathbb{R}^n$$

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

data:

$$\in \mathbb{R}^d$$

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

$$p(\mathbf{z}) = \prod_i p(z_i)$$

Problem Statement

Given observations $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ with $\mathbf{x}^{(i)} \sim p(\mathbf{x})$, we aim to find a model $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})$.

Independent Component Analysis

Linear model:

$$\mathbf{x} = \mathbf{A}\mathbf{z} + \boldsymbol{\epsilon}$$

$$p(\mathbf{x}|\mathbf{z}) \sim \mathcal{N}(\mathbf{A}\mathbf{z}, \boldsymbol{\Sigma})$$

Search for $\mathbf{W} = \mathbf{A}^{-1}$, such that mutual information between components of $\mathbf{z} \approx \mathbf{W}\mathbf{x}$ is minimized.

Independent Component Analysis

Linear model:

$$\mathbf{x} = A\mathbf{z} + \boldsymbol{\epsilon}$$

$$p(\mathbf{x}|\mathbf{z}) \sim \mathcal{N}(A\mathbf{z}, \Sigma)$$

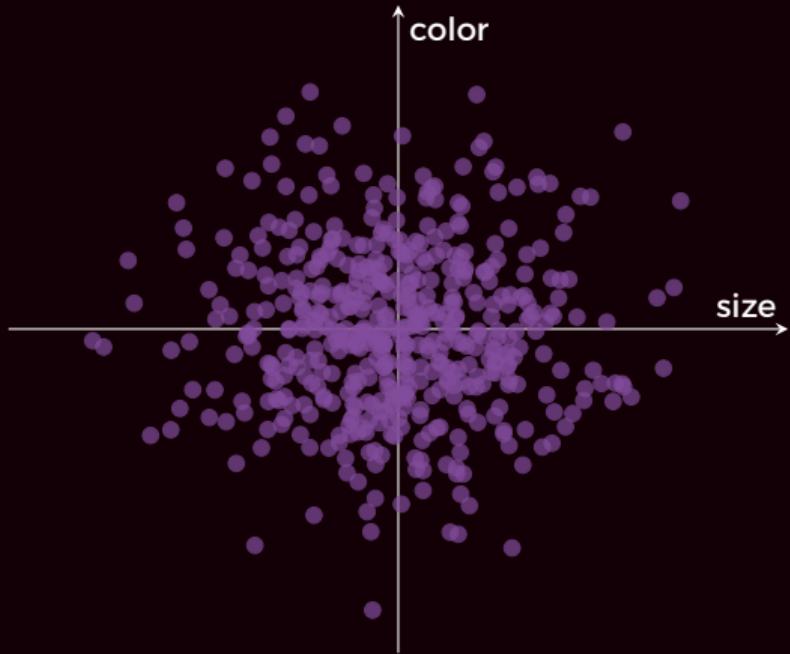
Search for $W = A^{-1}$, such that mutual information between components of $\mathbf{z} \approx W\mathbf{x}$ is minimized.

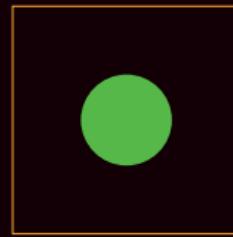
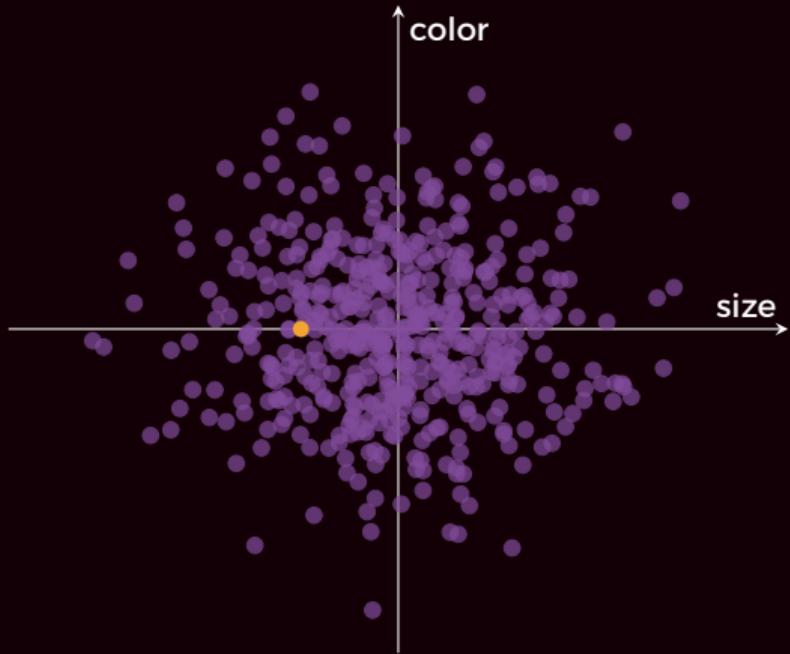
Disentanglement \approx nonlinear ICA

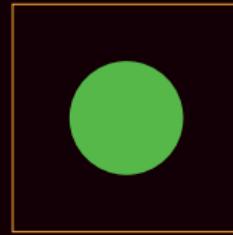
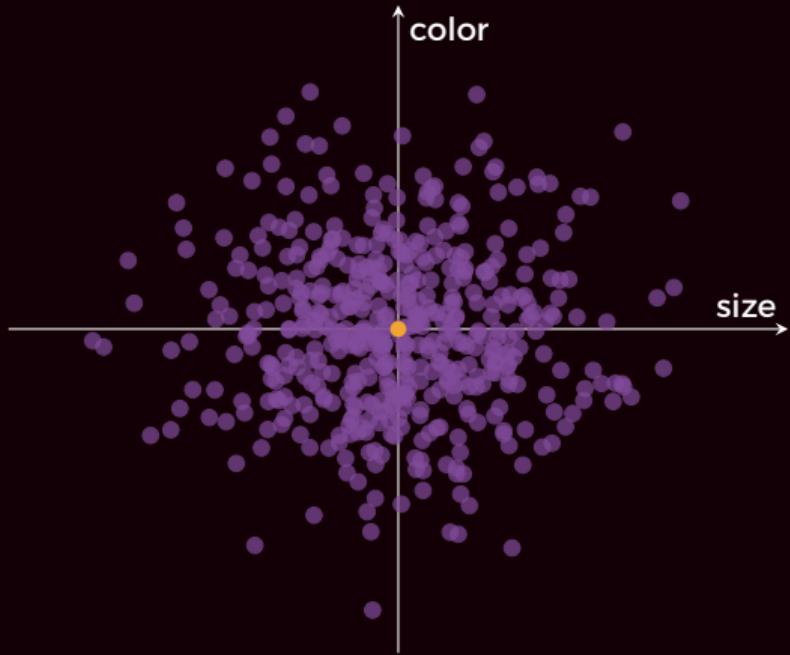
No particular restriction on $p(\mathbf{x}|\mathbf{z})$:

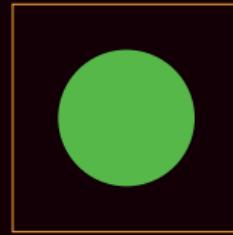
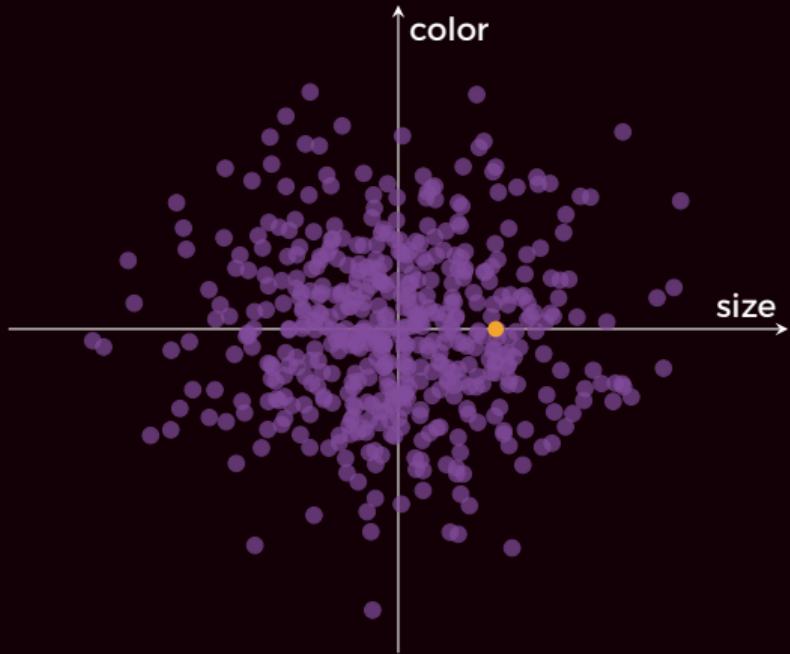
$$\mathbf{x} = f(\mathbf{z}) + \boldsymbol{\epsilon}$$

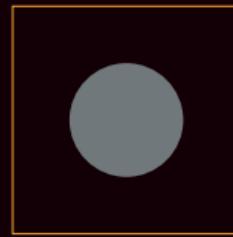
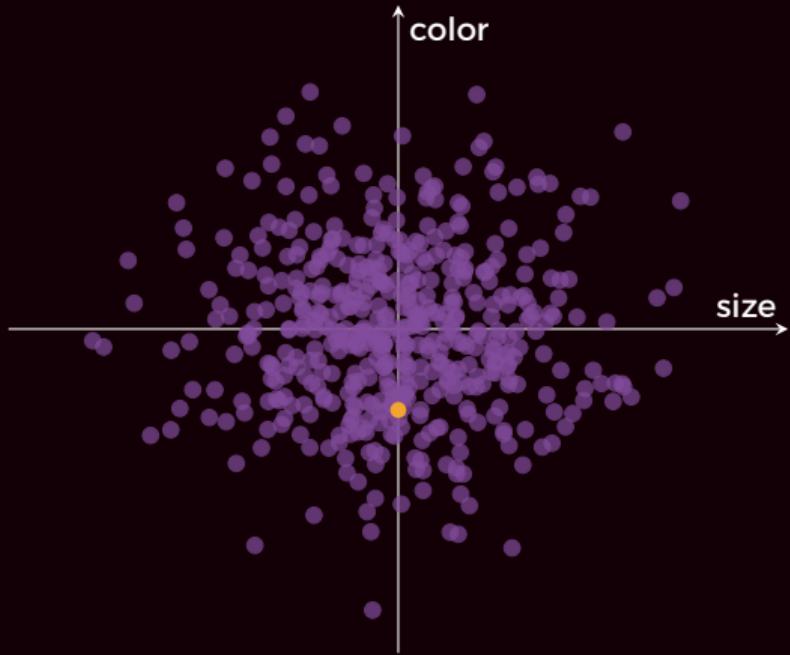
Search for f^{-1} , such that distribution of $\mathbf{z} = f^{-1}(\mathbf{x})$ factorizes: $p(\mathbf{z}) = \prod_i p(z_i)$.

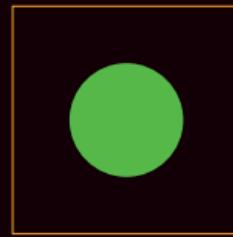
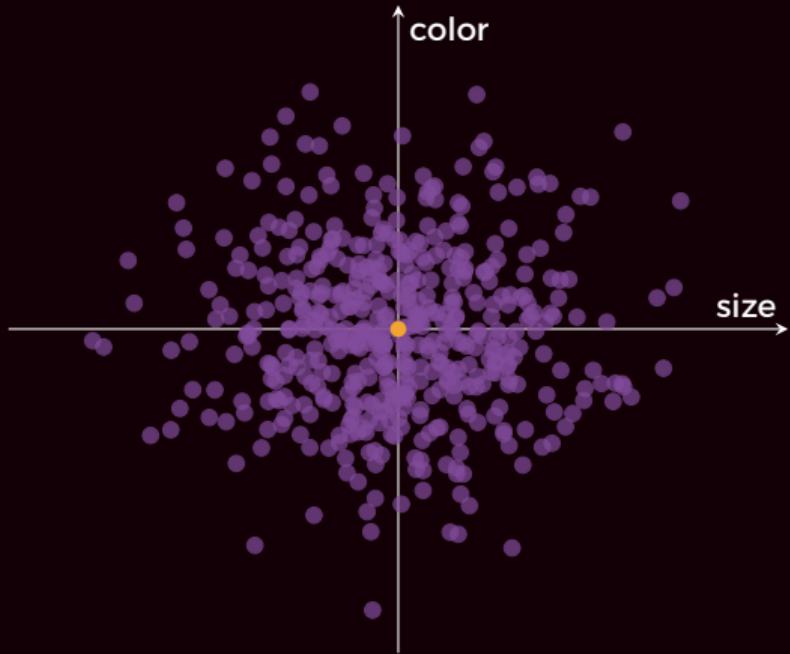


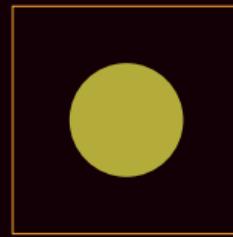
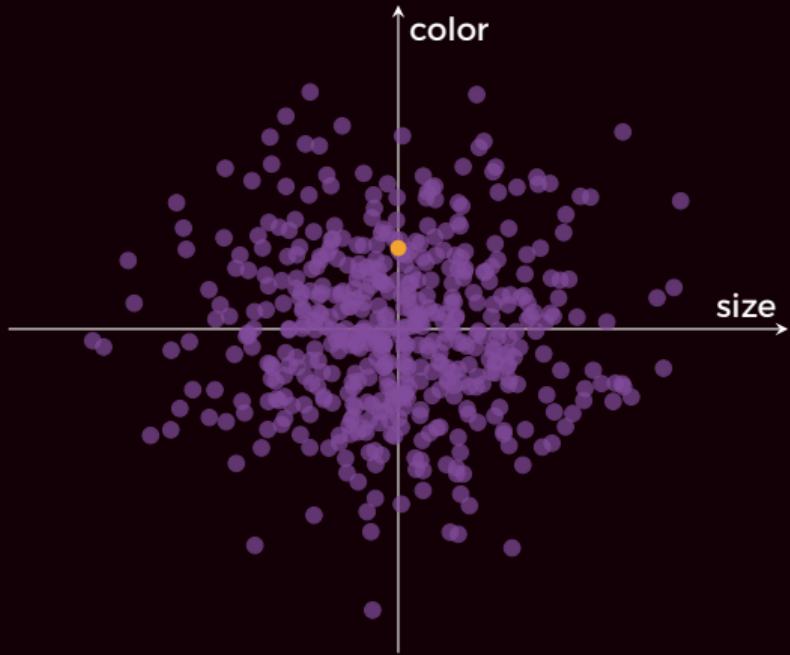








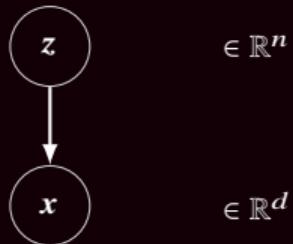




Chapter 1: Just do it!

Model Definition

latent variable:



$$\in \mathbb{R}^n$$

data:

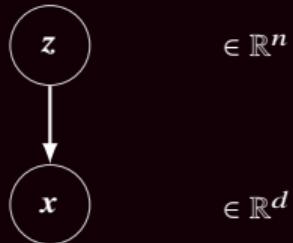
$$\in \mathbb{R}^d$$

Model Definition

- Model a factorizing prior over latent variables:

$$p_{\theta}(z) = \prod_i p_{\theta}(z_i)$$

latent variable:



$$\in \mathbb{R}^n$$

data:

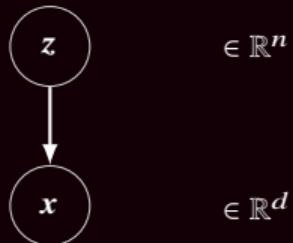
$$\in \mathbb{R}^d$$

Model Definition

- Model a factorizing prior over latent variables:

$$p_{\theta}(\mathbf{z}) = \prod_i p_{\theta}(z_i)$$

latent variable:



- Learn to generate data from latent variables:

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

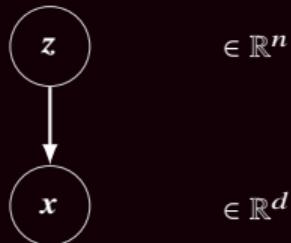
such that $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z} \approx p(\mathbf{x})$

Model Definition

- Model a factorizing prior over latent variables:

$$p_{\theta}(\mathbf{z}) = \prod_i p_{\theta}(z_i)$$

latent variable:



$$\in \mathbb{R}^n$$

- Learn to generate data from latent variables:

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

data:

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

- Extra Credit Learn to predict z from x :

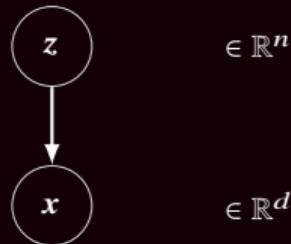
$$p_{\theta}(z|x)$$

Model Definition

- Model a factorizing prior over latent variables:

$$p_{\theta}(z) = \prod_i p_{\theta}(z_i)$$

latent variable:



$$\in \mathbb{R}^n$$

data:

$$\in \mathbb{R}^d$$

- Learn to generate data from latent variables:

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

such that $p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz \approx p(x)$

- Extra Credit Learn to predict z from x :

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

Variational Autoencoders (VAEs) are well suited for this task.

Variational Autoencoder

Ingredients

- Prior on z :

$$p(z) = \mathcal{N}(\mathbf{0}, I)$$

- Generator:

$$p_\theta(x|z) = \mathcal{N}(\mu_x, \Sigma_x)$$

$(\mu_x, \Sigma_x) = f_\theta(z)$ (typically a neural network)

- Encoder:

$$q_\theta(z|x) = \mathcal{N}(\mu_z, \Sigma_z)$$

$(\mu_z, \Sigma_z) = g_\theta(x)$ (typically a neural network)

$$q_\theta(z|x) \approx p_\theta(z|x)$$

Variational Autoencoder

Ingredients

- Prior on z :

$$p(z) = \mathcal{N}(\mathbf{0}, I)$$

- Generator:

$$p_\theta(x|z) = \mathcal{N}(\mu_x, \Sigma_x)$$

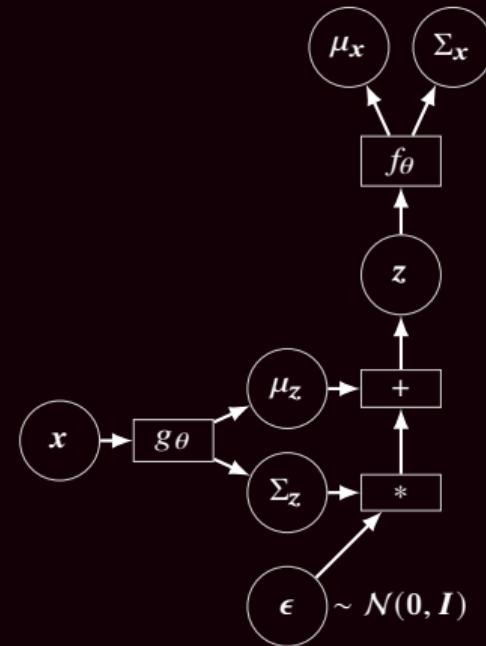
$(\mu_x, \Sigma_x) = f_\theta(z)$ (typically a neural network)

- Encoder:

$$q_\theta(z|x) = \mathcal{N}(\mu_z, \Sigma_z)$$

$(\mu_z, \Sigma_z) = g_\theta(x)$ (typically a neural network)

$$q_\theta(z|x) \approx p_\theta(z|x)$$



Variational Autoencoder

Ingredients

- Prior on z :

$$p(z) = \mathcal{N}(\mathbf{0}, I)$$

- Generator:

$$p_\theta(x|z) = \mathcal{N}(\mu_x, \Sigma_x)$$

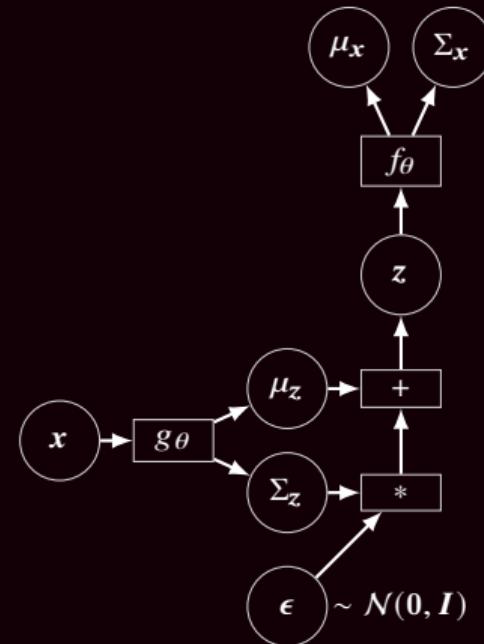
$(\mu_x, \Sigma_x) = f_\theta(z)$ (typically a neural network)

- Encoder:

$$q_\theta(z|x) = \mathcal{N}(\mu_z, \Sigma_z)$$

$(\mu_z, \Sigma_z) = g_\theta(x)$ (typically a neural network)

$$q_\theta(z|x) \approx p_\theta(z|x)$$



Training

$$\mathcal{L}(\theta, x) = \mathbb{E}_{q_\theta(z|x)} [-\log p_\theta(x|z)] + D_{\text{KL}}(q_\theta(z|x) || p(z)) \rightarrow \min$$

Variational Autoencoder

Ingredients

- Prior on z :

$$p(z) = \mathcal{N}(\mathbf{0}, I)$$

- Generator:

$$p_\theta(x|z) = \mathcal{N}(\mu_x, \Sigma_x)$$

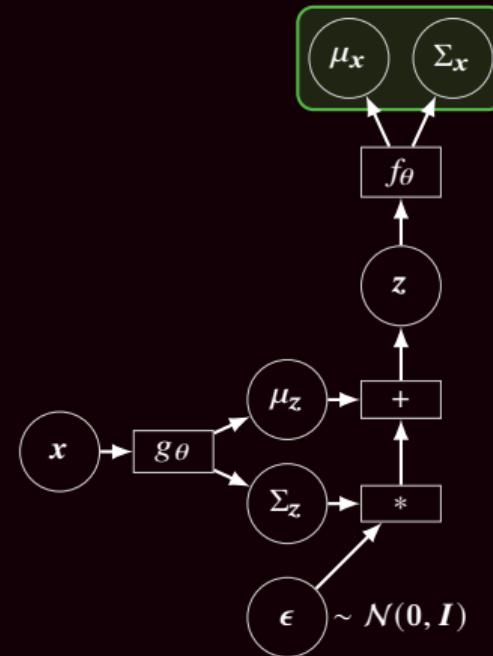
$(\mu_x, \Sigma_x) = f_\theta(z)$ (typically a neural network)

- Encoder:

$$q_\theta(z|x) = \mathcal{N}(\mu_z, \Sigma_z)$$

$(\mu_z, \Sigma_z) = g_\theta(x)$ (typically a neural network)

$$q_\theta(z|x) \approx p_\theta(z|x)$$



Training

$$\mathcal{L}(\theta, x) = \boxed{\mathbb{E}_{q_\theta(z|x)} [-\log p_\theta(x|z)]} + D_{\text{KL}}(q_\theta(z|x) || p(z)) \rightarrow \min$$

Variational Autoencoder

Ingredients

- Prior on z :

$$p(z) = \mathcal{N}(\mathbf{0}, I)$$

- Generator:

$$p_\theta(x|z) = \mathcal{N}(\mu_x, \Sigma_x)$$

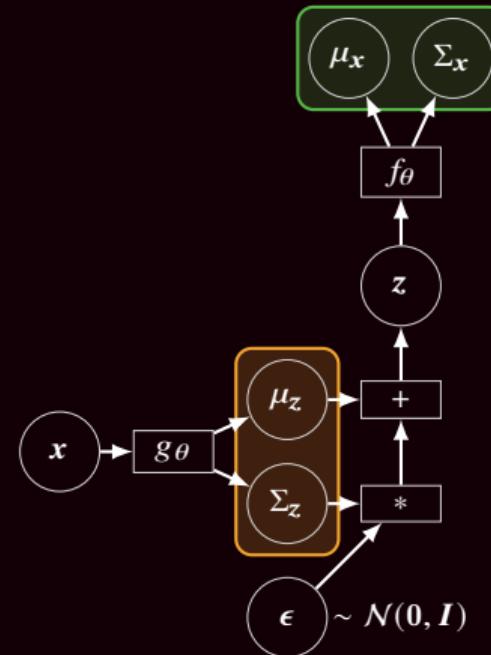
$$(\mu_x, \Sigma_x) = f_\theta(z) \quad (\text{typically a neural network})$$

- Encoder:

$$q_\theta(z|x) = \mathcal{N}(\mu_z, \Sigma_z)$$

$$(\mu_z, \Sigma_z) = g_\theta(x) \quad (\text{typically a neural network})$$

$$q_\theta(z|x) \approx p_\theta(z|x)$$



Training

$$\mathcal{L}(\theta, x) = \boxed{\mathbb{E}_{q_\theta(z|x)} [-\log p_\theta(x|z)]} + \boxed{D_{\text{KL}}(q_\theta(z|x) || p(z))} \rightarrow \min$$

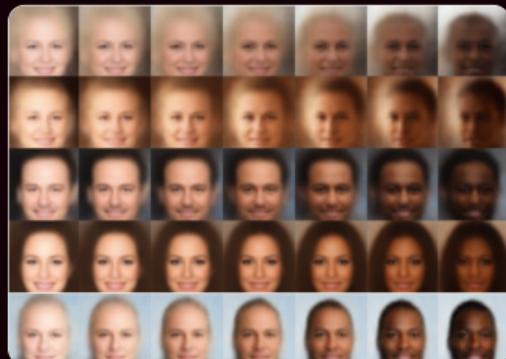
β -VAE

$$\mathcal{L}(\theta, \mathbf{x}) = \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})] + \beta D_{\text{KL}}(q_\theta(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \rightarrow \min$$

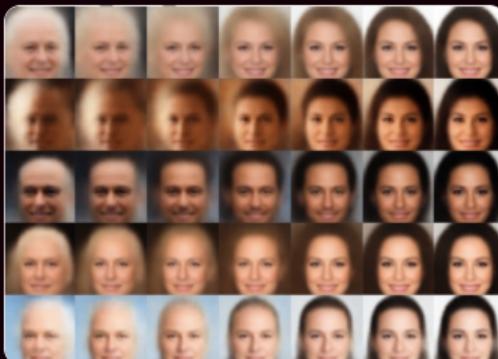
β -VAE

$$\mathcal{L}(\theta, \mathbf{x}) = \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})] + \beta D_{\text{KL}}(q_\theta(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \rightarrow \min$$

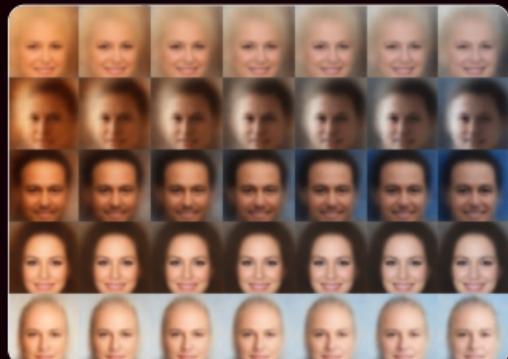
Examples on CelebA



skin color

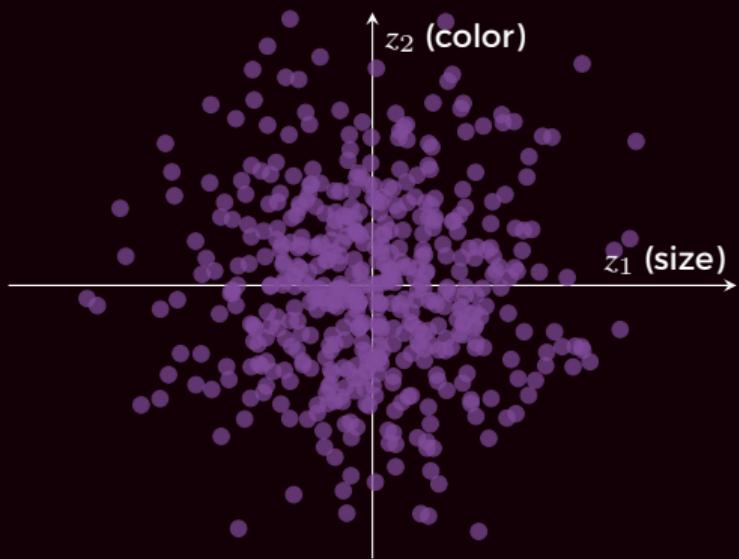


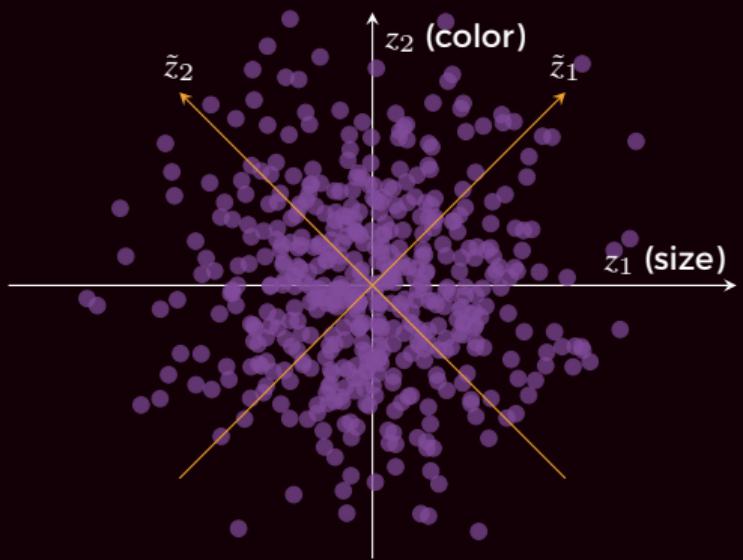
age/gender

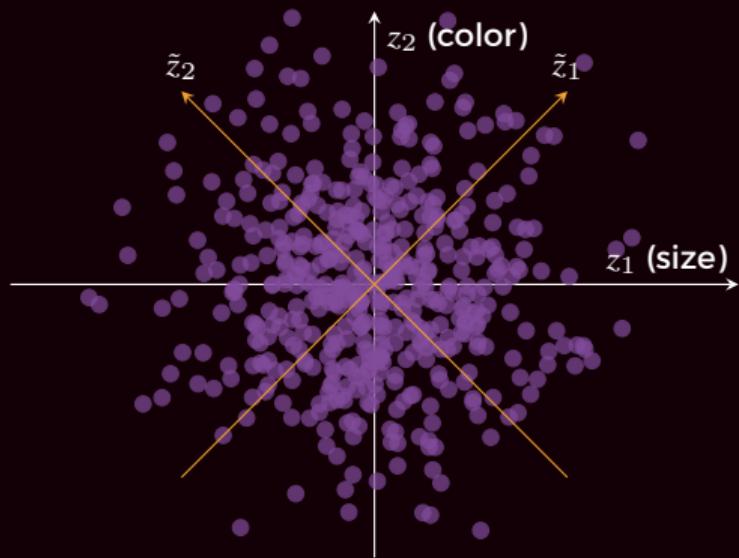


saturation

Chapter 2: Bad News



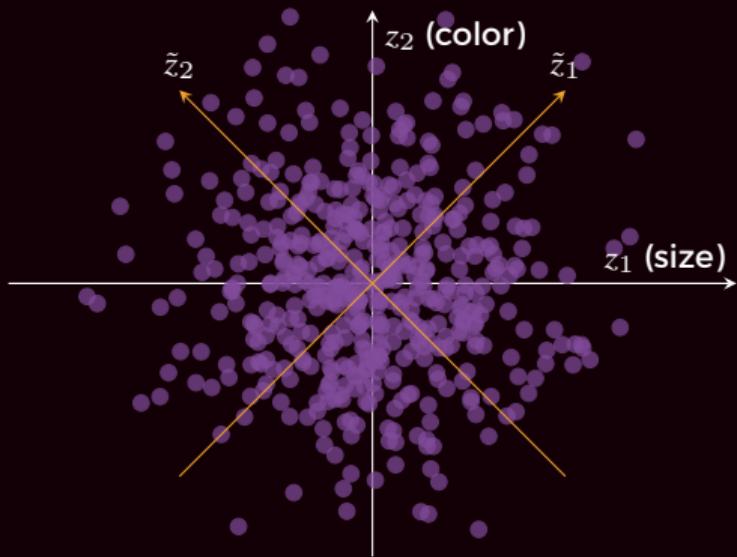




Gaussian Prior

$$p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$$

$$z \sim p_{\theta}(z) = \mathcal{N}(\mathbf{0}, I)$$



Problem

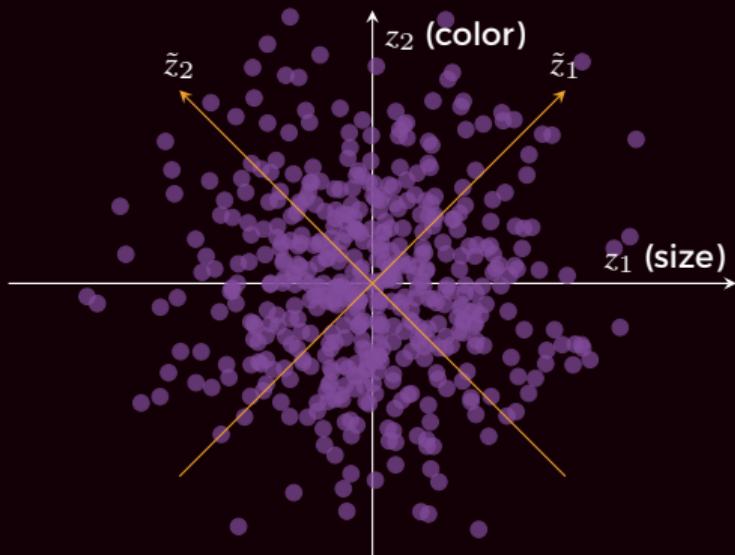
Let $\tilde{\mathbf{z}} = M\mathbf{z}$, where M is an orthogonal transformation:

$$\begin{aligned}
 p(\tilde{\mathbf{z}}) &= p(M^\top \mathbf{z}) |\det M| \\
 &= \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|M^\top \mathbf{z}\|^2\right) \\
 &= \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|\mathbf{z}\|^2\right) \\
 &= p(\mathbf{z})
 \end{aligned}$$

Gaussian Prior

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

$$\mathbf{z} \sim p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Problem

Let $\tilde{z} = Mz$, where M is an orthogonal transformation:

$$\begin{aligned} p(\tilde{z}) &= p(M^\top z) |\det M| \\ &= \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} ||M^\top z||^2\right) \\ &= \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} ||z||^2\right) \\ &= p(z) \end{aligned}$$

Gaussian Prior

$$p_\theta(x) = \int p_\theta(x|z)p_\theta(z)dz$$

$$z \sim p_\theta(z) = \mathcal{N}(\mathbf{0}, I)$$

Even Worse...

This applies to any factorial prior:

- transform z_i into uniform distribution via CDF $F_i(z_i)$
- further transform into Gaussian distribution: $\Phi^{-1}(F_i(z_i))$

Alternative Proof in Locatello et al.

“[...] unsupervised learning of disentangled representations is fundamentally impossible without inductive biases on both the models and the data.”

“[...] well-disentangled models seemingly cannot be identified without supervision.”

Chapter 3: Semi-Supervision to the Rescue!

Conditional Priors

Replace $p_\theta(z)$ with $p_\theta(z|u)$, where u is “an additionally observed variable”:

$$p_\theta(x, z|u) = p_\theta(x|z)p_\theta(z|u)$$

Conditional Priors

Replace $p_\theta(z)$ with $p_\theta(z|\mathbf{u})$, where \mathbf{u} is “an additionally observed variable”:

$$p_\theta(\mathbf{x}, \mathbf{z}|\mathbf{u}) = p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z}|\mathbf{u})$$

In particular, model the prior as an exponential family:

$$p_\theta(\mathbf{z}|\mathbf{u}) = \prod_i \frac{Q_i(z_i)}{Z_i(\mathbf{u})} \exp \left[\sum_{j=1}^k T_{i,j}(z_i) \lambda_{i,j}(\mathbf{u}) \right]$$

where $T_{i,j}$ are (learnable) sufficient statistics (k per component of \mathbf{z}) and $\lambda_{i,j}$ their (learnable) coefficients.

Identifiability Theorem 1

Under some mild conditions, the parameters θ of the true distribution are identifiable up to:

- a linear invertible transformation A
- point-wise nonlinearities

For $k = 1$:

$$(T_1^*(z_1^*), \dots, T_n^*(z_n^*)) = A(T_1(z_1), \dots, T_n(z_n))$$

Identifiability Theorem 1

Under some mild conditions, the parameters θ of the true distribution are identifiable up to:

- a linear invertible transformation A
- point-wise nonlinearities

For $k = 1$:

$$(T_1^*(z_1^*), \dots, T_n^*(z_n^*)) = A(T_1(z_1), \dots, T_n(z_n))$$

Identifiability Theorem 2

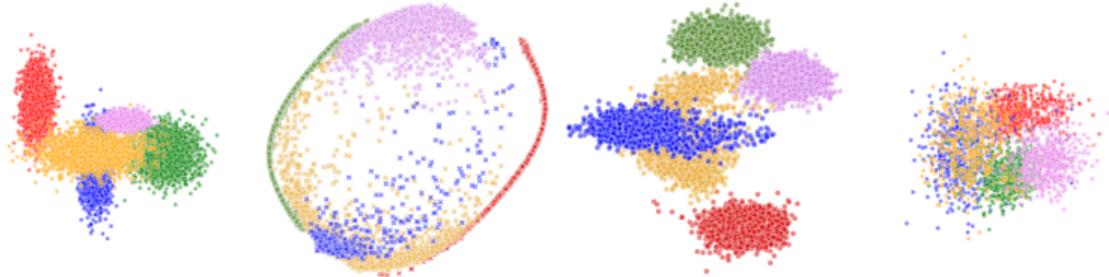
If further

- $T_{i,j}$ are twice differentiable
- f (the generator function) has all second order cross derivatives

then A reduces to a permutation matrix (for $k = 1$ some weaker conditions can be applied).

For $k = 1$:

$$T_i^*(z_i^*) = T_{i'}(z_{i'})$$



(a) $p_{\theta^*}(\mathbf{z}|\mathbf{u})$ (b) $p_{\theta^*}(\mathbf{x}|\mathbf{u})$ (c) $p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{u})$ (d) $p_{\text{VAE}}(\mathbf{z}|\mathbf{x})$

Discussion

Chapter 1

- VAEs are well suited to learn data distributions with handcrafted latent priors
- β -VAEs (and variants) put extra emphasis on “disentanglement”
- posterior estimate $q_\theta(z|x)$ encodes data, potentially useful for downstream tasks

Discussion

Chapter 1

- VAEs are well suited to learn data distributions with handcrafted latent priors
- β -VAEs (and variants) put extra emphasis on “disentanglement”
- posterior estimate $q_\theta(z|x)$ encodes data, potentially useful for downstream tasks

Chapter 2

- true factors of variation can not be discovered with nonlinear models
- this result does not depend on the design of $p_\theta(z)$
- supervision or inductive biases needed

Discussion

Chapter 1

- VAEs are well suited to learn data distributions with handcrafted latent priors
- β -VAEs (and variants) put extra emphasis on “disentanglement”
- posterior estimate $q_\theta(z|x)$ encodes data, potentially useful for downstream tasks

Chapter 2

- true factors of variation can not be discovered with nonlinear models
- this result does not depend on the design of $p_\theta(z)$
- supervision or inductive biases needed

Chapter 3

- conditional priors $p_\theta(z|u)$ allow identification of factors of variation up to simple transform
- under some conditions, the transform reduces to a permutation plus pointwise nonlinearity
- unclear what the requirements on u are

Conditions of Theorem 1

- The set $\{x \in X \mid \phi_\epsilon(x) = 0\}$ has measure zero.
- f is injective.
- $T_{i,j}$ are differentiable almost everywhere and $(T_{i,j})_{i \leq j \leq k}$ are linearly independent.
- There exist $nk + 1$ distinct points $\mathbf{u}^0, \dots, \mathbf{u}^{nk}$ such that

$$L = (\lambda(\mathbf{u}^1) - \lambda(\mathbf{u}^0), \dots, \lambda(\mathbf{u}^{nk} - \lambda(\mathbf{u}^0)))$$

is invertible.