

Literature Review on VAEs with Correlated Latent Features

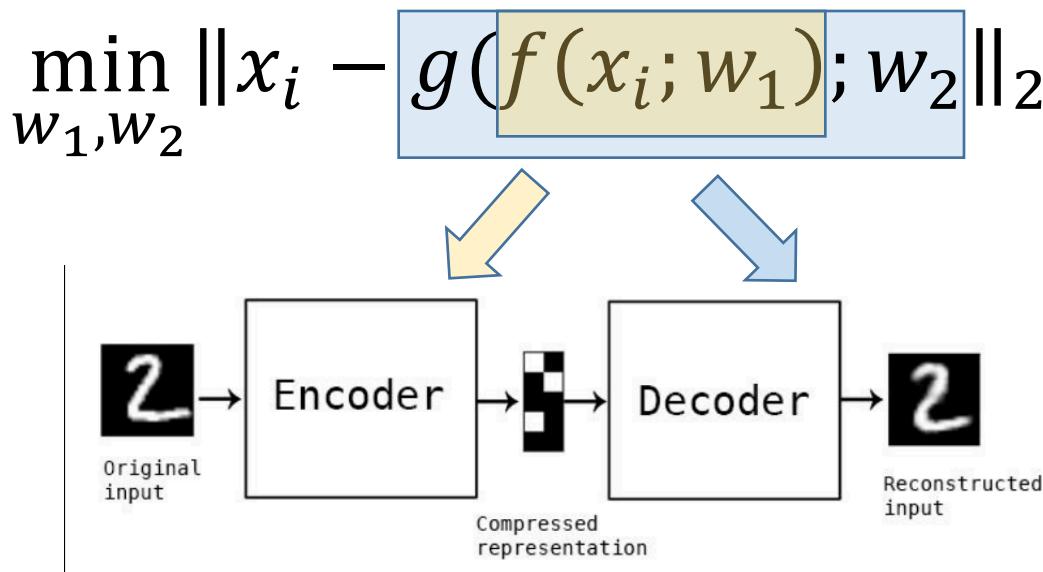
Javier Zapata

Today's plan

- VAE basics (mostly Kingma et al, 2014)
 - Autoencoders vs Variational Autoencoders
 - Variational Lower Bound
 - Reparameterization Trick
- VAE with correlated latent features (Casale et al, 2018)

Auto-Encoder (deterministic)

- A neural network that the output is the input itself.
- Intuition:
 - A good representation should keep the information well (reconstruction error)
 - Deep + nonlinearity might help enhance the representation power



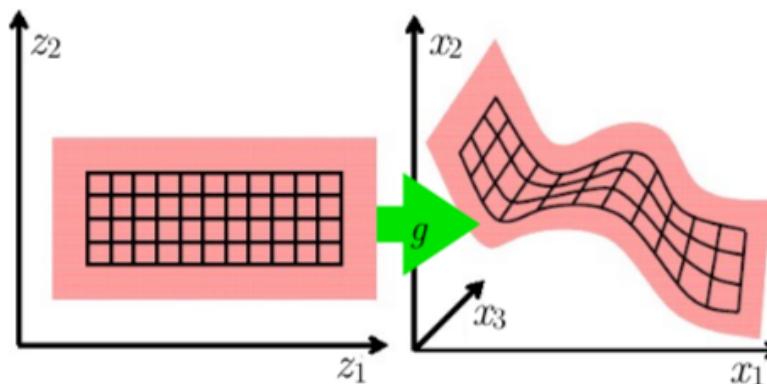
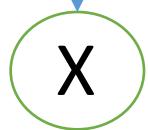
Variational Auto-Encoder (VAE) (probabilistic)

- Uses neural networks to learn a latent variable model
- A blend of variational inference & neural networks
- Seminal papers:
 - Kingma and Welling, *Auto-Encoding Variational Bayes, International Conference on Learning Representations (ICLR)* 2014.
 - Rezende, Mohamed and Wierstra, *Stochastic back-propagation and variational inference in deep latent Gaussian models*. ICML 2014.

Latent
variable



Observed
Sample
variable



Variational Auto-Encoder (VAE) (probabilistic)

Latent variable



- A blend of variational inference & neural networks
- **Latent Model:** Learn mapping from Z to X

where:

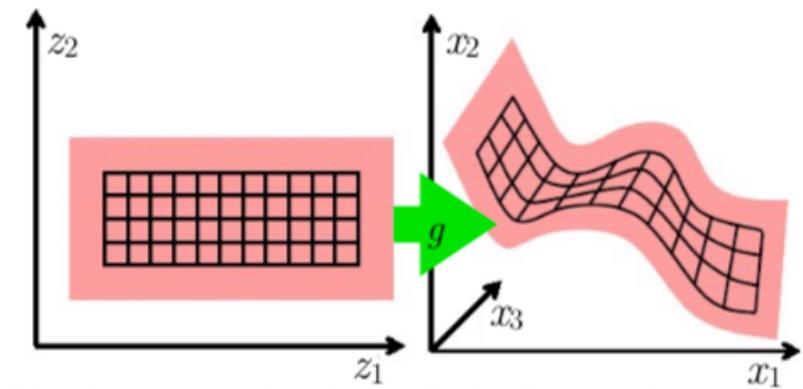
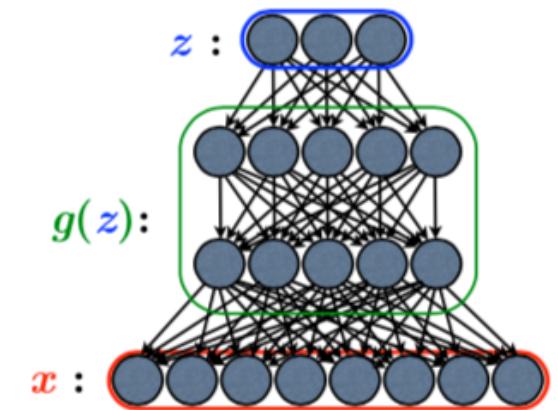
$$p(x) = \int p(x, z) dz = \int p(x|z) p(z) dz$$

$p(z)$: something simple

$p(x|z) = g(z)$ (a neural net)

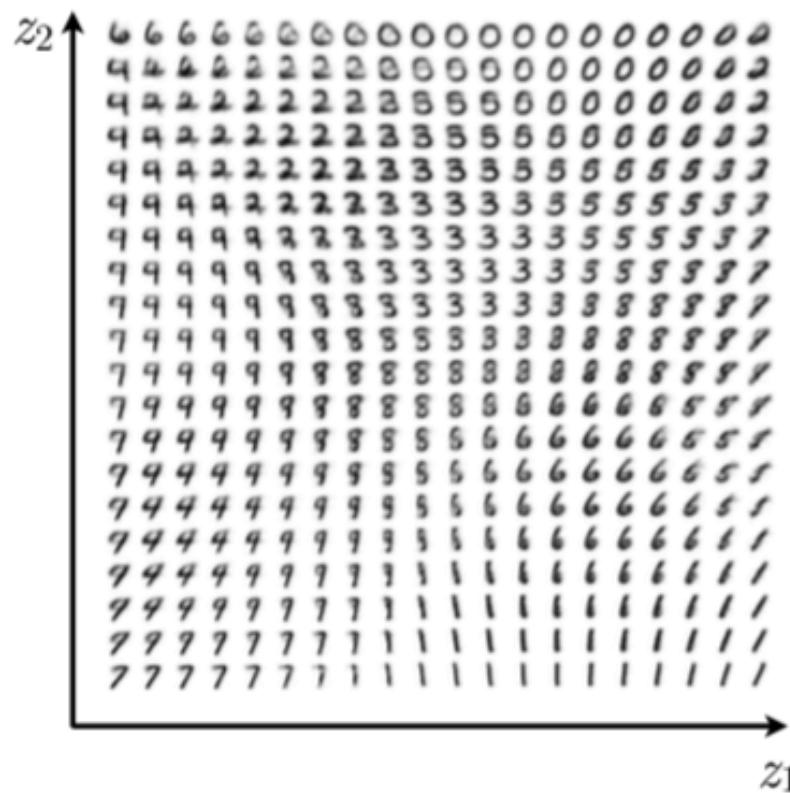
Observed Sample variable

- In principle we can think of:
- **Encoder:** $p(z|x)$
- **Decoder:** $p(x|z)$
- **Objective function:** ML or MAP on $p(x_1, \dots, x_n)$

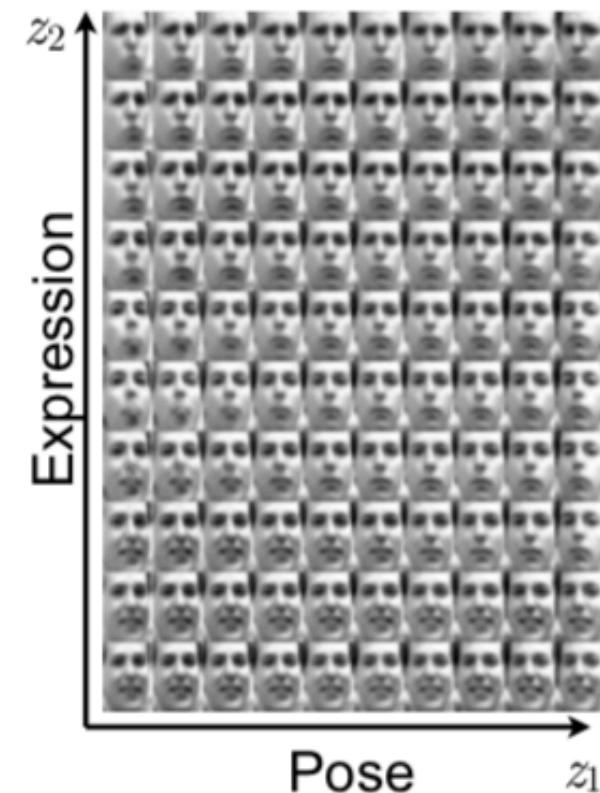


What VAE can do? Learn manifolds

Learned MNIST manifold

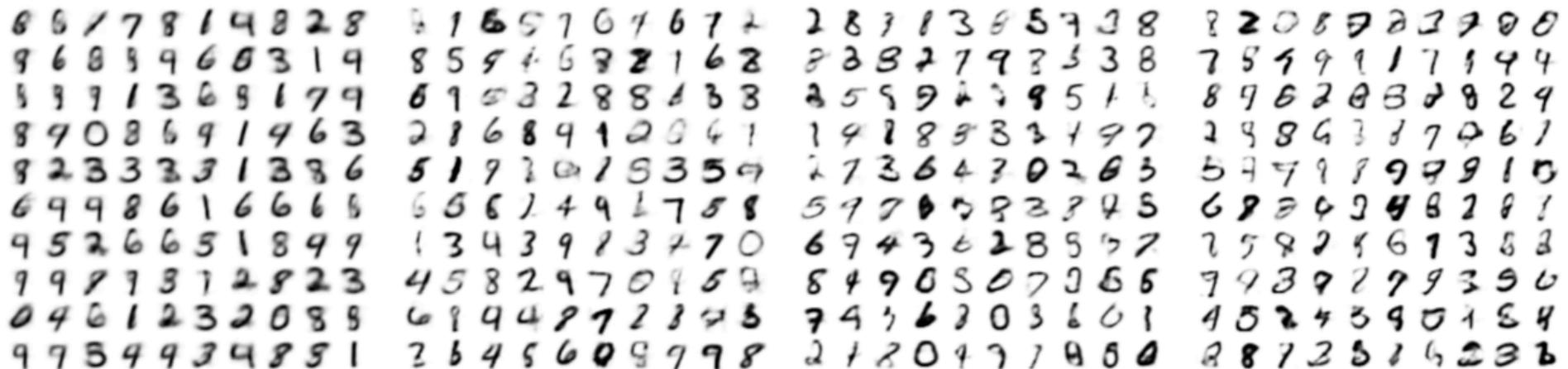


Learned face manifold



What VAE can do?

Sample from generative learned model



(a) 2-D latent space

(b) 5-D latent space

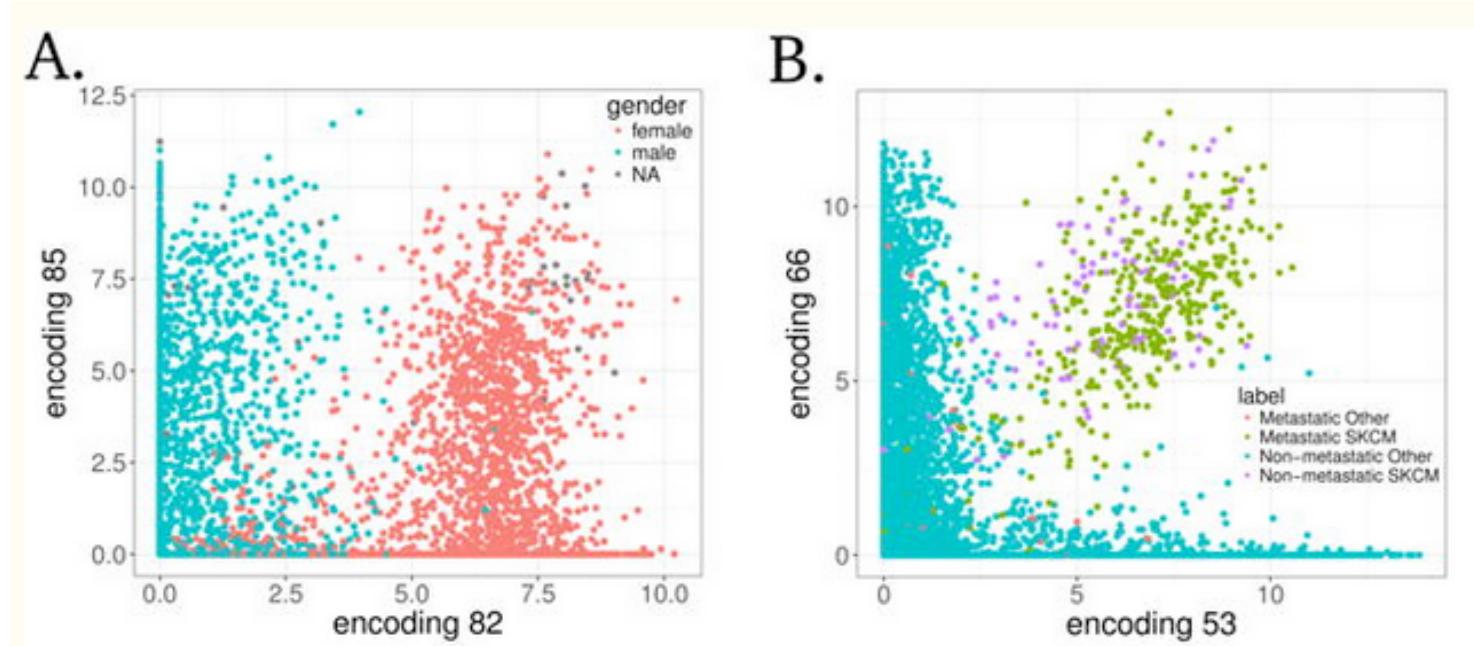
(c) 10-D latent space

(d) 20-D latent space

Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

What VAE can do? Interpret gene expressions

- Extracting a biologically relevant latent space: (Way et al, 2018)
 - Encoding 5,000 gene expression vectors into a vector of 100 dimensions
 - Distinguishing patients gender and metastatic activation



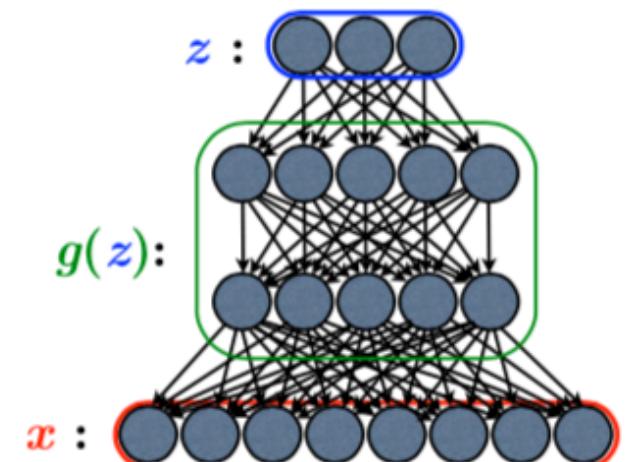
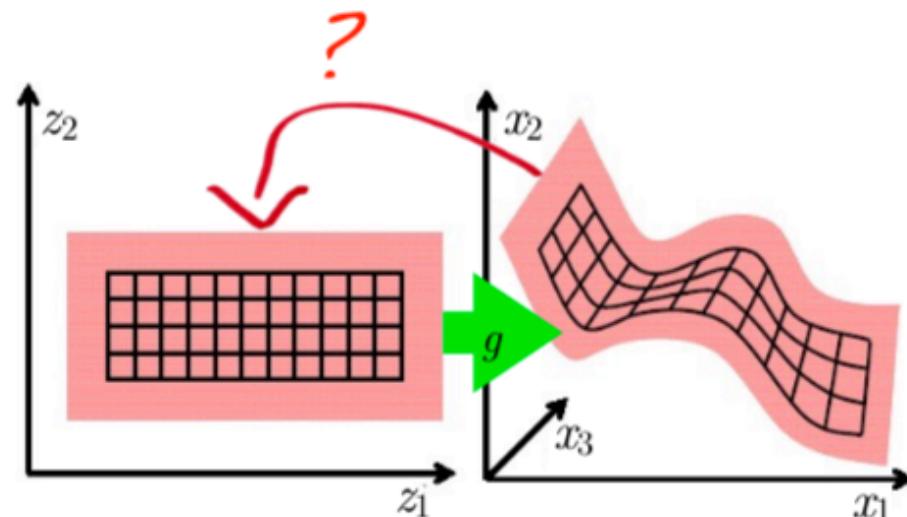
The inference / learning challenge

- **Where does z come from?** — The classic directed model dilemma.
- Computing the posterior $p(z | x)$ is intractable.
- We need it to train the directed model.

Latent
variable



Observed
Sample
variable



Example:

Decoder:

$$p_{\theta}(z): Z \sim N(0, I)$$

$$p_{\theta}(x|z): X|Z \sim No(\mu_{(z,\theta)}, \sigma_{(z,\theta)}^2 I)$$

where $\mu = W_1 h + b_1$

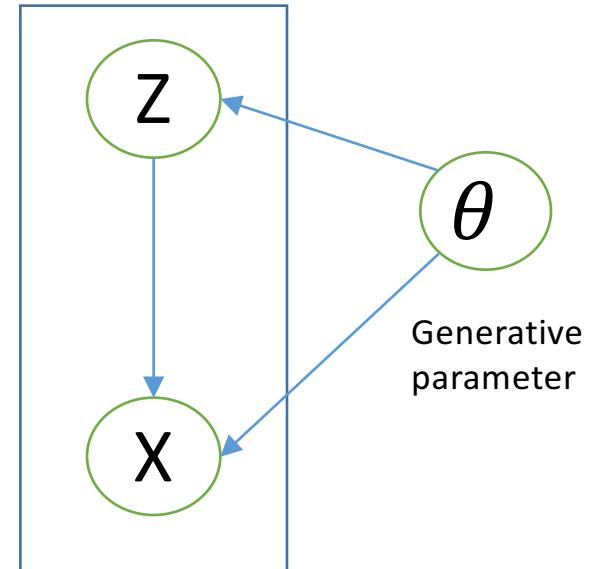
$\log \sigma^2 = W_2 h + b_2$ (componentwise)

$h = \tanh(W_3 z + b_3)$ (componentwise)

so $\theta := \{W_1, W_2, W_3, b_1, b_2, b_3\}$

Encoder:

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



Usually true posterior
is intractable ☹

- The VAE approach: introduce a simplified posterior $q_\phi(z|x)$ that **learns** to approximate the true posterior $p_\theta(z|x)$
- For this to happen we use a **variational lower bound (aka ELBO)**:

$$\log(p(x)) = \mathcal{L}(\theta, \phi, x) + KL(q||p)$$

where : $KL(q||p) = \mathbb{E}_{q(Z|x)} \left[\log \left(\frac{q(Z|x)}{p(Z|x)} \right) \right] \geq 0$

$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{q(Z|x)} [\log(p(x, Z)) - \log(q(Z|x))]$$

$$= \mathbb{E}_{q(Z|x)} [\log(p(x|Z)) + \log(p(Z)) - \log(q(Z|x))]$$

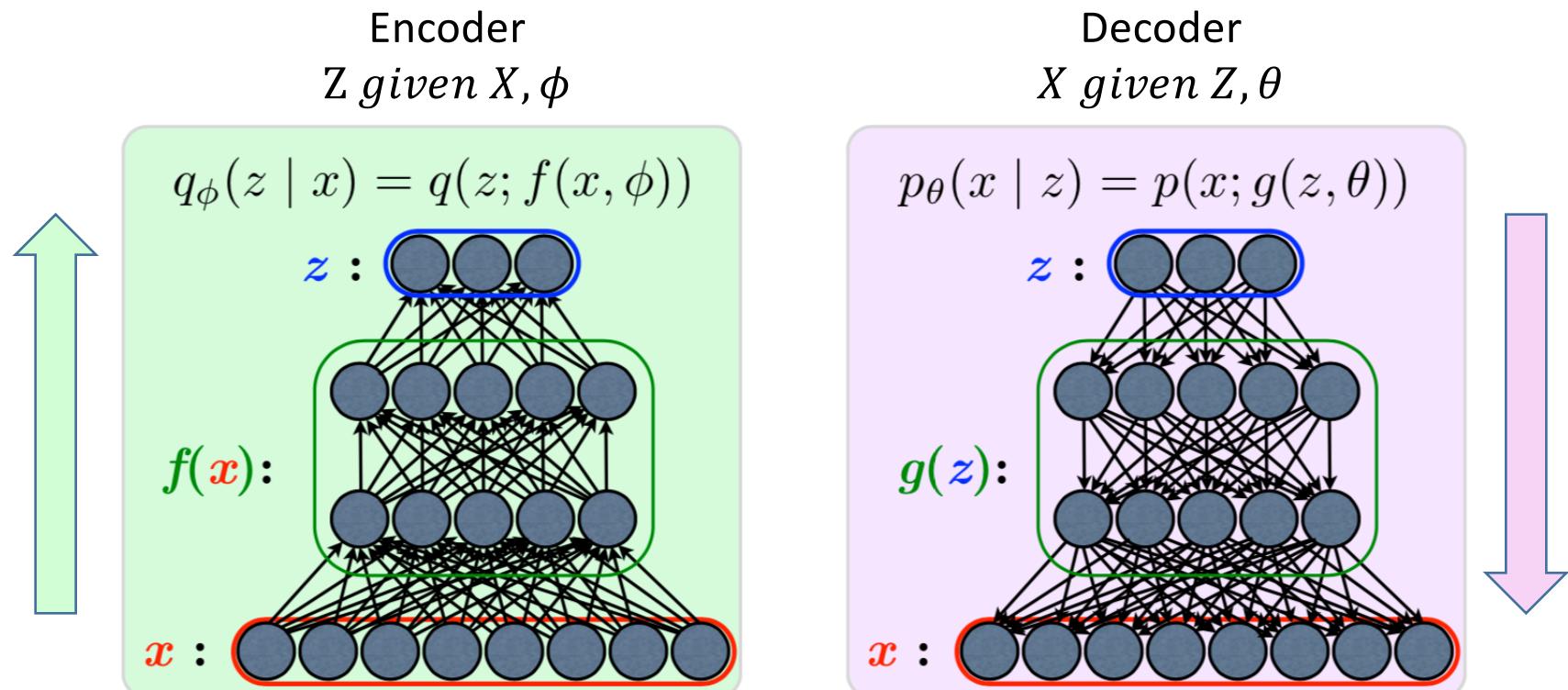
$$= \mathbb{E}_{q(Z|x)} [\log(p(x|Z))] - KL(q(Z|x)||p(Z))$$

reconstruction
term

regularization
term

- We include a neural network into $q_\phi(z|x)$ and now we solve:

$$\max_{\theta, \phi} \mathcal{L}(\theta, \phi, x) = \max_{\theta, \phi} \mathbb{E}_{q_\phi(Z|x)} \left[\log(p_\theta(x, Z)) - \log(q_\phi(Z|x)) \right]$$



Example:

Decoder:

$$p_{\theta}(z): Z \sim N(0, I)$$

$$p_{\theta}(x|z): X|Z = z \sim No(\mu_{(z,\theta)}, \sigma_{(z,\theta)}^2 I)$$

where $\mu = W_1 h + b_1$

$$\log \sigma^2 = W_2 h + b_2 \text{ (componentwise)}$$

$$h = \tanh(W_3 z + b_3) \text{ (componentwise)}$$

$$\text{so } \theta := \{W_1, W_2, W_3, b_1, b_2, b_3\}$$

Encoder:

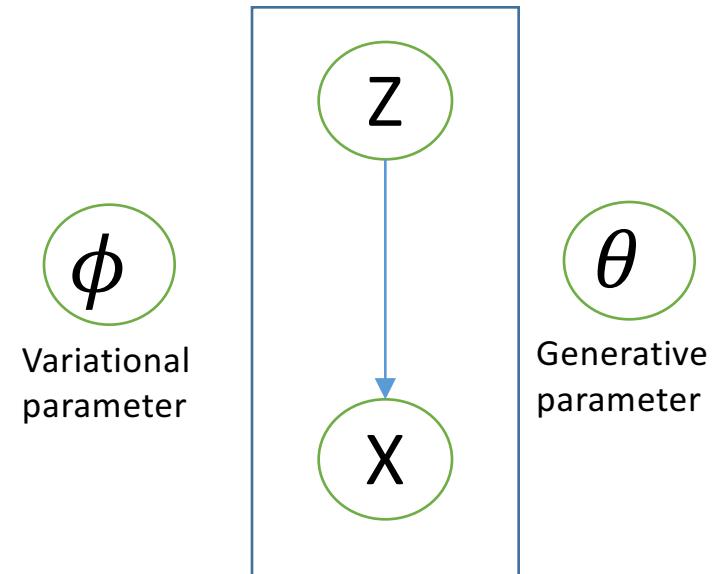
$$q_{\phi}(z|x): Z|X = x \sim No(\tilde{\mu}_{(x,\phi)}, \tilde{\sigma}_{(x,\phi)}^2 I)$$

where $\tilde{\mu} = \tilde{W}_1 h + \tilde{b}_1$

$$\log \tilde{\sigma}^2 = \tilde{W}_2 h + \tilde{b}_2 \text{ (componentwise)}$$

$$h = \tanh(\tilde{W}_3 x + \tilde{b}_3) \text{ (componentwise)}$$

$$\text{so } \phi := \{\tilde{W}_1, \tilde{W}_2, \tilde{W}_3, \tilde{b}_1, \tilde{b}_2, \tilde{b}_3\}$$



Reparameterization Trick

- We want to maximize the ELBO, wrt θ and ϕ using stochastic gradient descent

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log q_\phi(\mathbf{z}|\mathbf{x}) + \log p_\theta(\mathbf{x}, \mathbf{z})]$$

- Usual (naïve) Monte Carlo gradient exhibits high variance

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$$

- We use **Stochastic Gradient Variational Bayes (SGVB)** (Kingma 2014)

$q_\phi(z|x): Z|X=x \sim g_\phi(\epsilon, x)$ with $\epsilon \sim \text{known distribution}$

- Example:

$$q_\phi(z|x): Z|X=x \sim \text{No}(\tilde{\mu}_{(x,\phi)}, \tilde{\sigma}_{(x,\phi)}^2 I)$$

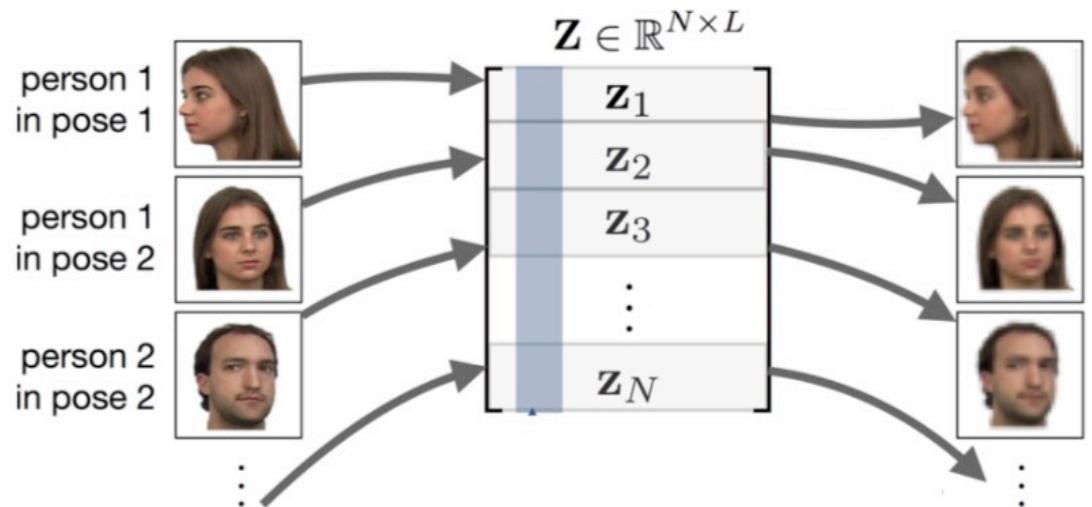
Sample using: $Z = \tilde{\mu} + \tilde{\sigma}^2 \epsilon$ with $\epsilon \sim N(0, I)$

VAE Pros

- Easy to train
- Blurry result due to minimizing the MSE based reconstruction error
- Nice probabilistic formulation, easy to introduce prior

VAE Cons

- In most VAE formulations latent variables Z are iid
- Ex: Multiple sequences of images from different cars, or medical image sequences from multiple patients.
- VAE prior should capture multiple levels of correlations at once, including time, object identities, etc.

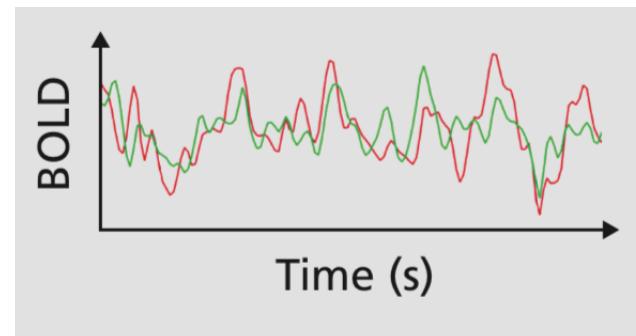
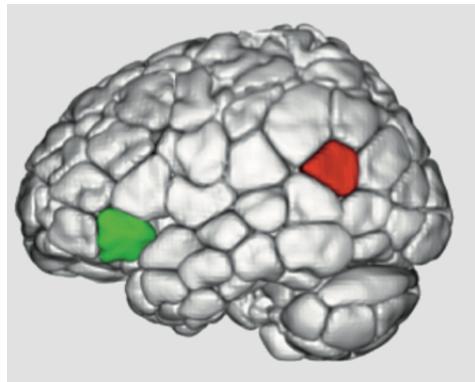


How to model the following:
“Object (person) p_n in view (pose) q_n ”???

We need some correlation in z_1, \dots, z_n

More examples

*Neuroimaging: fMRI signals (**sample**) of voxels (**objects**) under different motor tasks (**views**)*



- Right hand
- Left hand
- Right foot
- Left foot
- Tongue

*Art Images: J. Pollock paintings (**sample**), Number (**objects**) Series 1,5, 8 (**views**)*



Casale et al (2018)

Gaussian Process Prior Variational Autoencoders

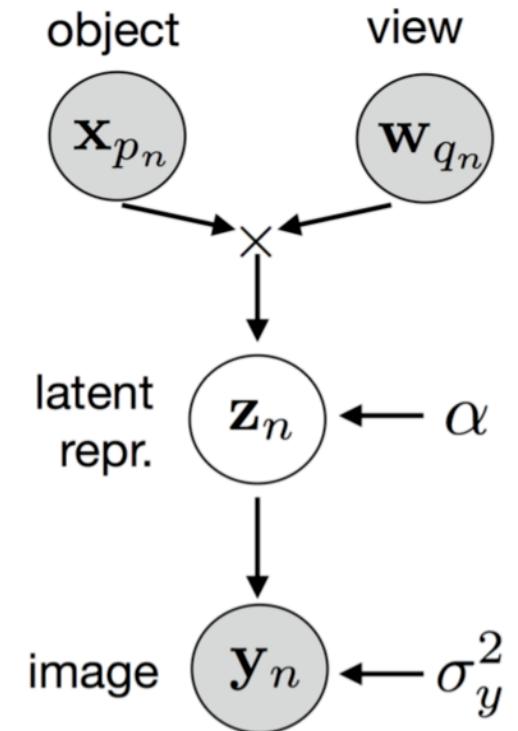
- y_1, \dots, y_n K-dimensional **samples** (not necessarily iid)
- z_1, \dots, z_n L-dimensional **latent variables**
- x_1, \dots, x_P M-dimensional **object feature vectors**
- w_1, \dots, w_Q R-dimensional **view feature vectors**

$N = \#$ images

$L = \#$ latent dims

$P = \#$ objects (people)

$Q = \#$ views (poses)



$$n = 1, \dots, N$$

Decoder $p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \phi, \sigma_y^2, \boldsymbol{\theta}, \alpha) = \int p(\mathbf{Y} | \mathbf{Z}, \phi, \sigma_y^2) p(\mathbf{Z} | \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha) d\mathbf{Z}$

$\mathbf{z}_n = f(\mathbf{x}_{p_n}, \mathbf{w}_{q_n}) + \boldsymbol{\eta}_n$, where $\boldsymbol{\eta}_n \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I}_L)$;

where $f(\cdot)$ generates \mathbf{z}_n with a GP prior:

$$p(\mathbf{Z} | \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha) = \prod_{l=1}^L \mathcal{N}(\mathbf{z}^l | \mathbf{0}, \mathbf{K}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{W}) + \alpha \mathbf{I}_N)$$

with covariance between samples n, m

$$\mathbf{K}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{W})_{nm} = \mathcal{K}_{\boldsymbol{\theta}}^{(\text{view})}(\mathbf{w}_{q_n}, \mathbf{w}_{q_m}) \mathcal{K}_{\boldsymbol{\theta}}^{(\text{object})}(\mathbf{x}_{p_n}, \mathbf{x}_{p_m})$$

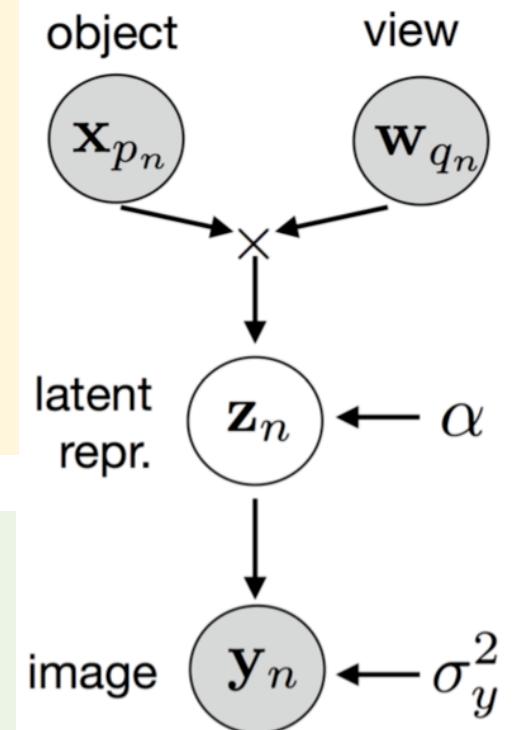
$\mathbf{y}_n = g(\mathbf{z}_n) + \boldsymbol{\epsilon}_n$, where $\boldsymbol{\epsilon}_n \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}_K)$

where $g(\cdot)$ projects \mathbf{z}_n to a high-dimensional space.

For images: a Convolutional Neural Net

For cntns variable:

activation function + linear transformation

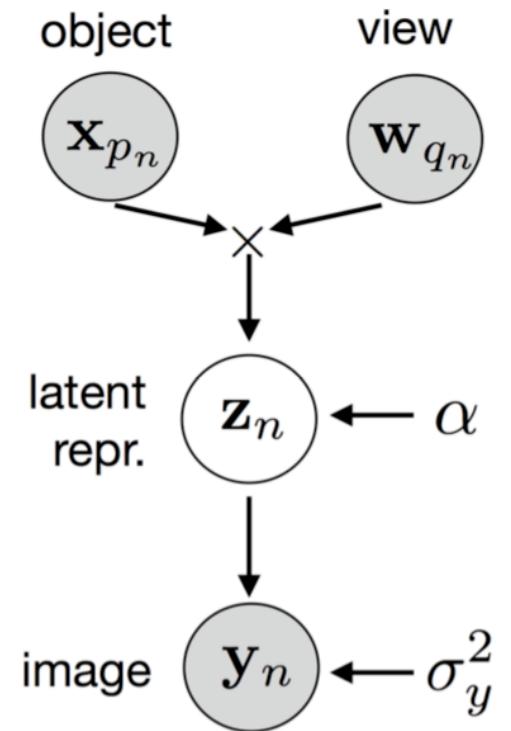


$$n = 1, \dots, N$$

Encoder

$$q_{\psi}(Z | Y) = \prod_n \mathcal{N} \left(z_n | \mu_{\psi}^z(y_n), \text{diag}(\sigma_{\psi}^{z^2}(y_n)) \right)$$

where μ_{ψ}^z and $\sigma_{\psi}^{z^2}$ are the variational parameters and are a neural network function of the observed data.



$$n = 1, \dots, N$$

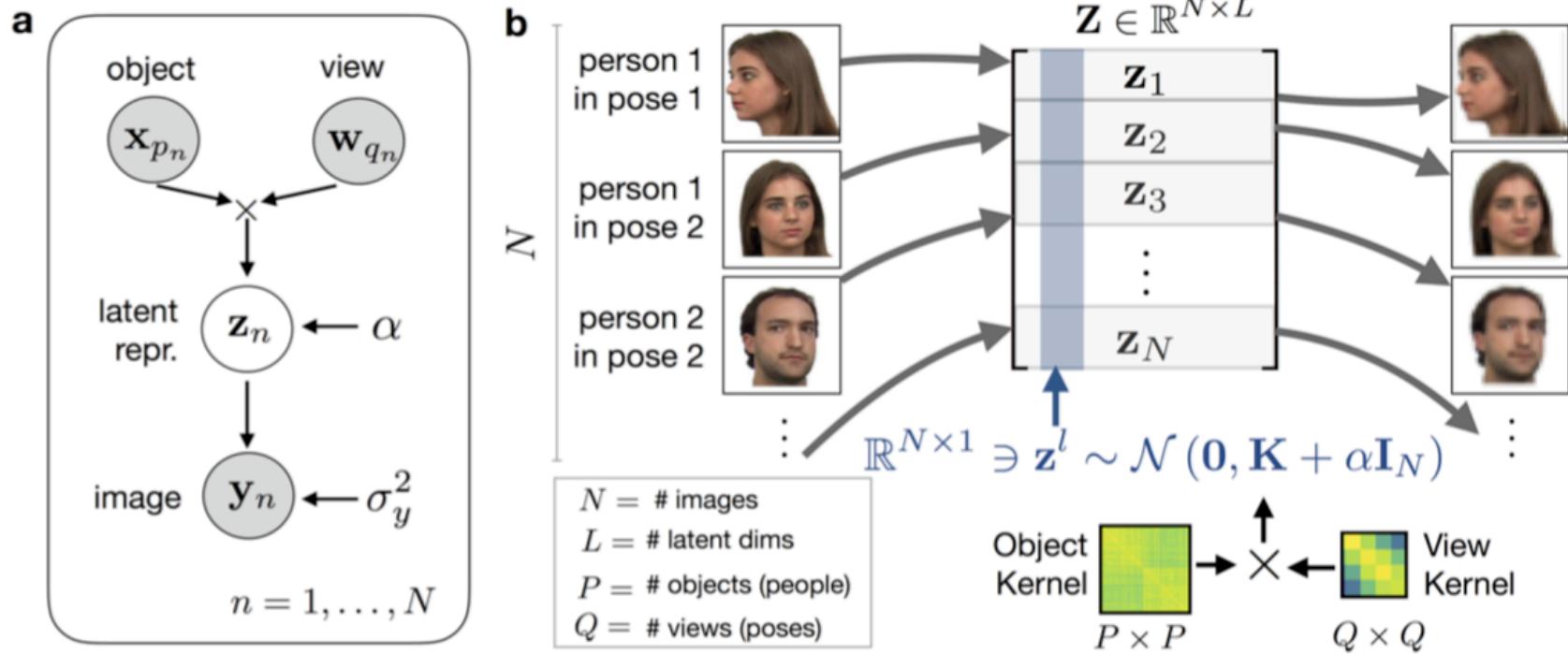


Figure 1: **(a)** Generative model underlying the proposed GPVAE. **(b)** Representation of the inference procedure in GPVAE. Each sample (here an image) is encoded in a low-dimensional space and then decoded to the original space. Covariances between samples modeled through a GP prior on each column of the latent representation matrix Z .

ELBO and Estimation

Variational lower bound:

$$\begin{aligned} \log p(\mathbf{Y} \mid \mathbf{X}, \mathbf{W}, \phi, \sigma_y^2, \boldsymbol{\theta}) &\geq \mathbb{E}_{\mathbf{Z} \sim q_{\psi}} \left[\sum_n \log \mathcal{N}(\mathbf{y}_n \mid g_{\phi}(\mathbf{z}_n), \sigma_y^2 \mathbf{I}_K) + \log p(\mathbf{Z} \mid \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha) \right] + \\ &+ \frac{1}{2} \sum_{nl} \log(\boldsymbol{\sigma}^{z^2}_{\psi}(\mathbf{y}_n)_l) + \text{const.} \end{aligned}$$

Stochastic Backpropagation:

$$\mathcal{L}(\phi, \psi, \boldsymbol{\theta}, \alpha, \sigma_y^2) = \frac{1}{K} \underbrace{\sum_n^N (\mathbf{y}_n - g_{\phi}(\mathbf{z}_{\psi_n}))^2}_{\text{reconstruction term}} - \frac{\lambda}{L} \left[\underbrace{\log p(\mathbf{Z}_{\psi} \mid \mathbf{X}, \mathbf{W}, \boldsymbol{\theta}, \alpha)}_{\text{latent-space GP model}} + \underbrace{\frac{1}{2} \sum_{nl} \log(\boldsymbol{\sigma}^{z^2}_{\psi}(\mathbf{y}_n)_l)}_{\text{regularization}} \right]$$

Where λ is a trade-off parameter between the data reconstruction term and the goodness of the latent space GP model

Reparameterization Trick

Remember that our posterior approximation is

$$q_{\psi}(\mathbf{Z} | \mathbf{Y}) = \prod_n \mathcal{N}\left(z_n | \mu_{\psi}^z(\mathbf{y}_n), \text{diag}(\sigma_{\psi}^{z2}(\mathbf{y}_n))\right)$$

and latent representations $\mathbf{Z}_{\psi} = [\mathbf{z}_{\psi_1}, \dots, \mathbf{z}_{\psi_N}] \in \mathbb{R}^{N \times L}$ are sampled as

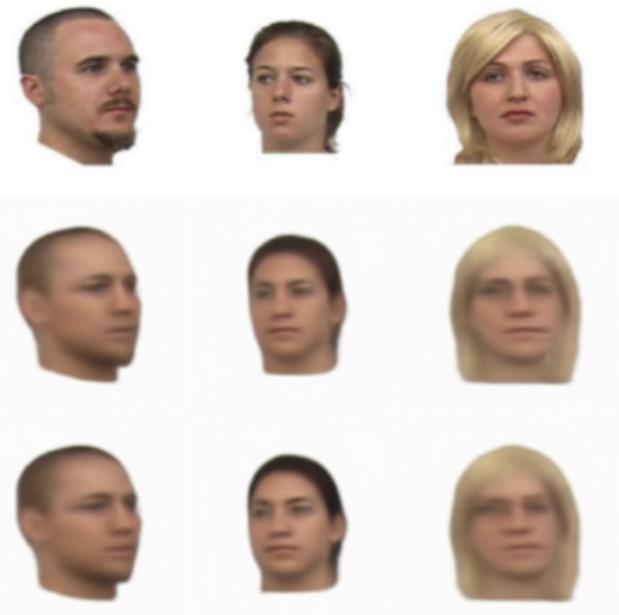
$$\mathbf{z}_{\psi_n} = \mu_{\psi}^z(\mathbf{y}_n) + \boldsymbol{\epsilon}_n \odot \sigma_{\psi}^z(\mathbf{y}_n), \quad \boldsymbol{\epsilon}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{L \times L}), \quad n = 1, \dots, N$$

where \odot represents the Hadamard product.

Predictive posterior sampling

$$p(\mathbf{y}_\star | \mathbf{x}_\star, \mathbf{w}_\star, \mathbf{Y}, \mathbf{X}, \mathbf{W}) \approx \int \underbrace{p(\mathbf{y}_\star | z_\star)}_{\text{decode GP prediction}} \underbrace{p(z_\star | \mathbf{x}_\star, \mathbf{w}_\star, \mathbf{Z}, \mathbf{X}, \mathbf{W})}_{\text{latent-space GP predictive posterior}} \underbrace{q(\mathbf{Z} | \mathbf{Y})}_{\text{encode training data}} dz_\star d\mathbf{Z}$$

this term replaces the true posterior $p_\theta(z|y)$



Compared Methods

1. **GPVAE-joint:** In each iteration
 - a) Fix kernel K_θ , and train all the other parameters
 - b) Fix all other parameters and train K_θ for 100 epochs
 - c) Train everything simultaneously for 1 epoch
2. **GPVAE-disjoint:** In each iteration iterate steps a) and b) of GPVAE-joint
3. **Conditional VAE (CVAE)**
4. **Linear Interpolation in VAE latents space (LIVAE)**

Example: Rotated MNIST (views known)

- **Data:** 400 handwritten versions of digit three with R= 16 rotation angles
- **Training Set:** 4,050 images spanning 15 rotations
- **Test Set:** 270 images with one rotation (out of sample)
- **Latent Space dimension:** L=16 **Object feature vector:** M=8

$$K_{\theta}(X, w)_{nm} = \underbrace{\beta \exp \left(-\frac{2\sin^2|w_{q_n} - w_{q_m}|}{\nu^2} \right)}_{\text{rotation kernel}} \cdot \underbrace{x_{p_n}^T x_{p_m}}_{\text{digit draw kernel}} \quad \text{where } \theta = \{ \beta, \nu \geq 0 \}$$

real (test set)	
GPVAE-joint	
GPVAE-dis	
CVAE	
LIVAE	

Example: Rotated MNIST (views known)

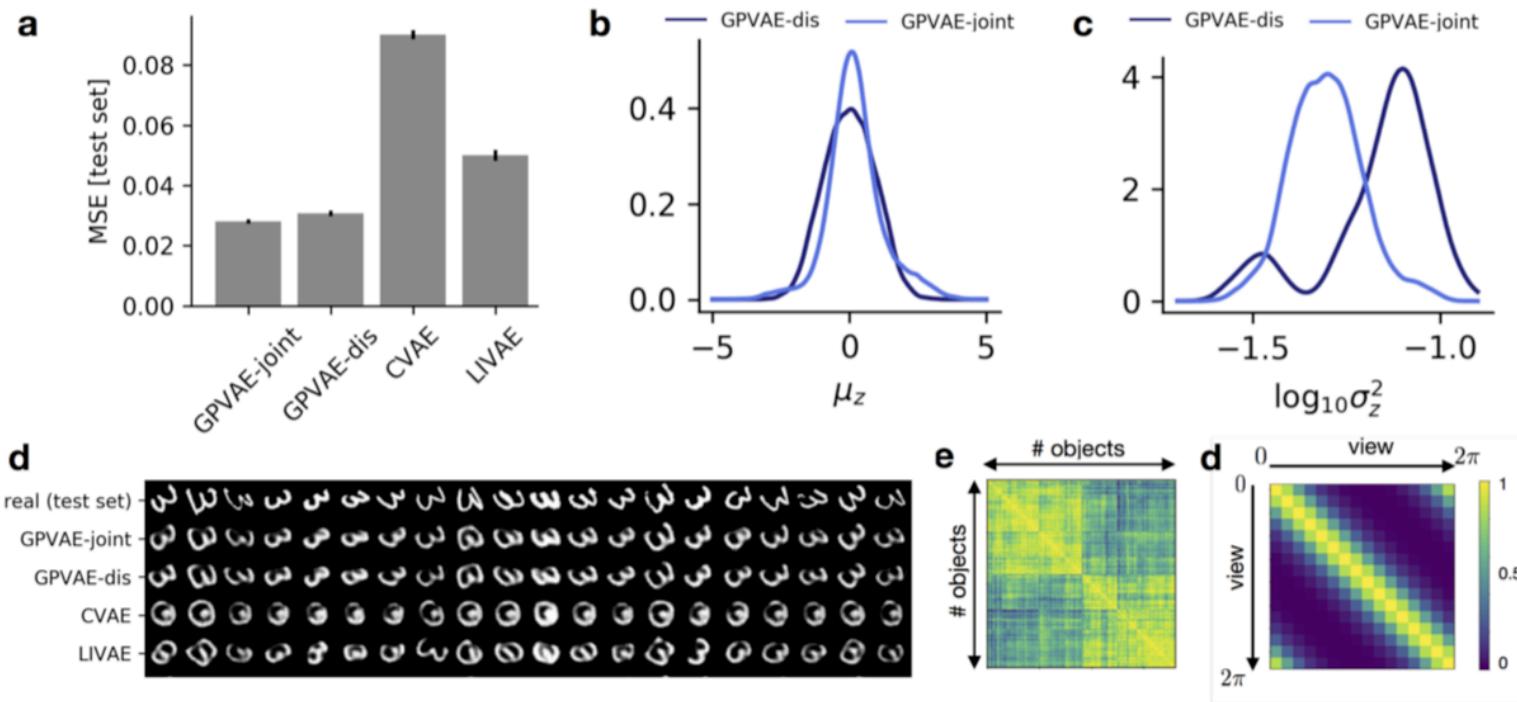


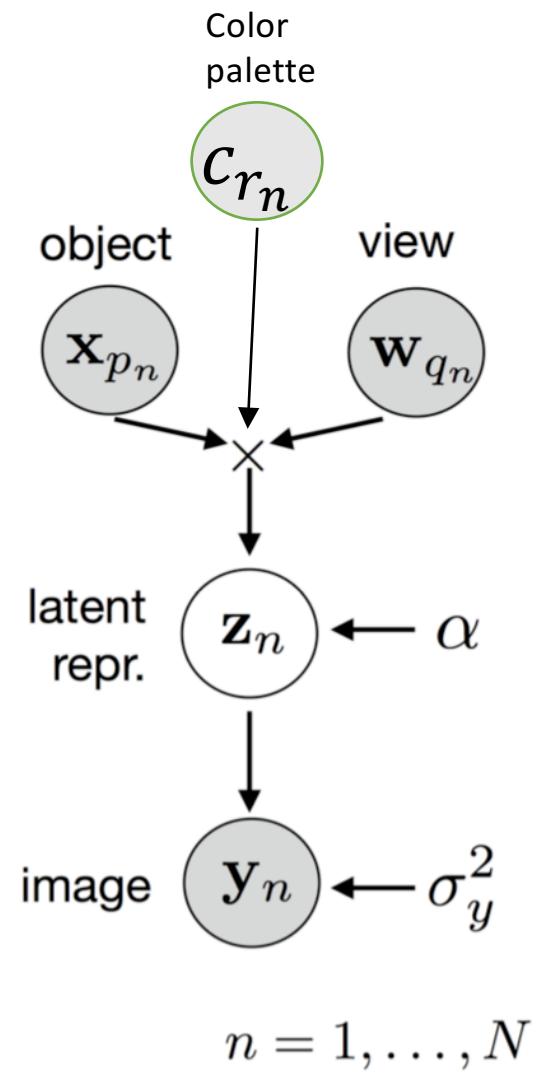
Figure 2: Results from experiments on rotated MNIST. (a) Mean squared error on test set. Error bars represent standard error of per-sample MSE. (b) Empirical density of estimated means of q_ψ , aggregated over all latent dimensions. (c) Empirical density of estimated log variances of q_ψ . (d) Out-of-sample predictions for ten random draws of digit "3" at the out-of-sample rotation state. (e, f) Object and view covariances learned through GPVAE-joint.

My Discussion: A hard example?

Andy Warhol's Mao Series #90-99



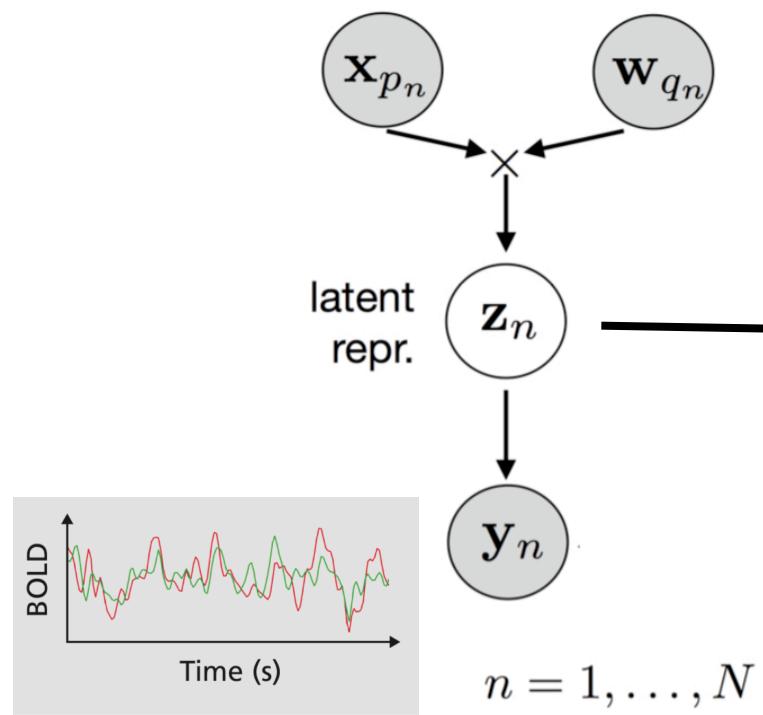
Andy Warhol's Marilyn Monroe Series #22-31



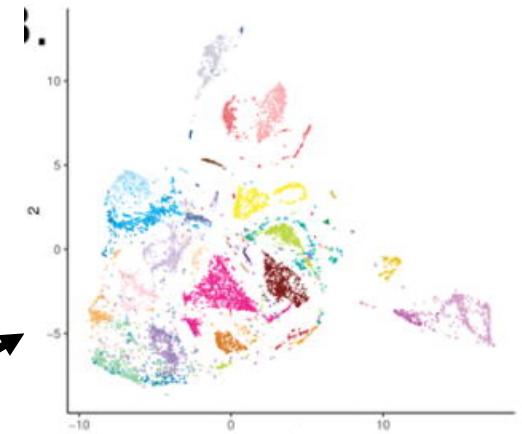
My Discussion: neuroimaging

signal (object)
(Functional PCA ??)

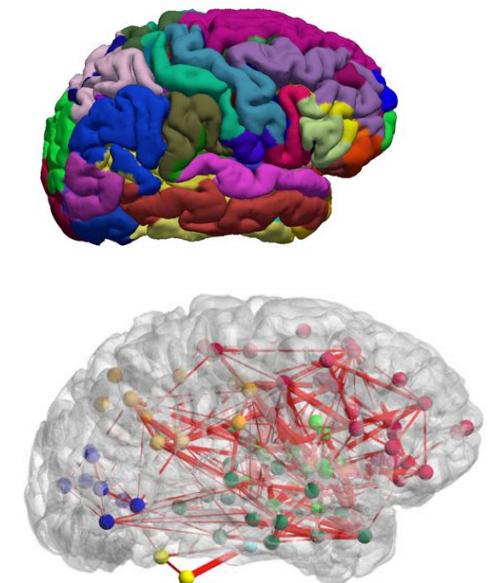
Motor Task
(view)



t-SNE Clustering
In latent space



Graphical
Lasso



Bibliography

- Courville - Variational Autoencoder and Extensions(http://videolectures.net/deeplearning2015_courville_autoencoder_extension/)
- Wang - Deep Generative Models (http://www.cs.toronto.edu/~slwang/generative_model.pdf)
- Way et. al. - Extracting a biologically relevant latent space from cancer transcriptomes with variational autoencoders (2018)
- Kingma et al - Auto-Encoding Variational Bayes (2014)
- Rezende et al - Stochastic back-propagation and variational inference in deep latent Gaussian models (2014)
- Casale et al – Gaussian Process Prior Variational Autoencoders (2018)