

Generative Modeling

Variational Autoencoders

Denis Derkach, Artem Ryzhikov, Maxim Artemev

Laboratory for methods of big data analysis

Spring 2021



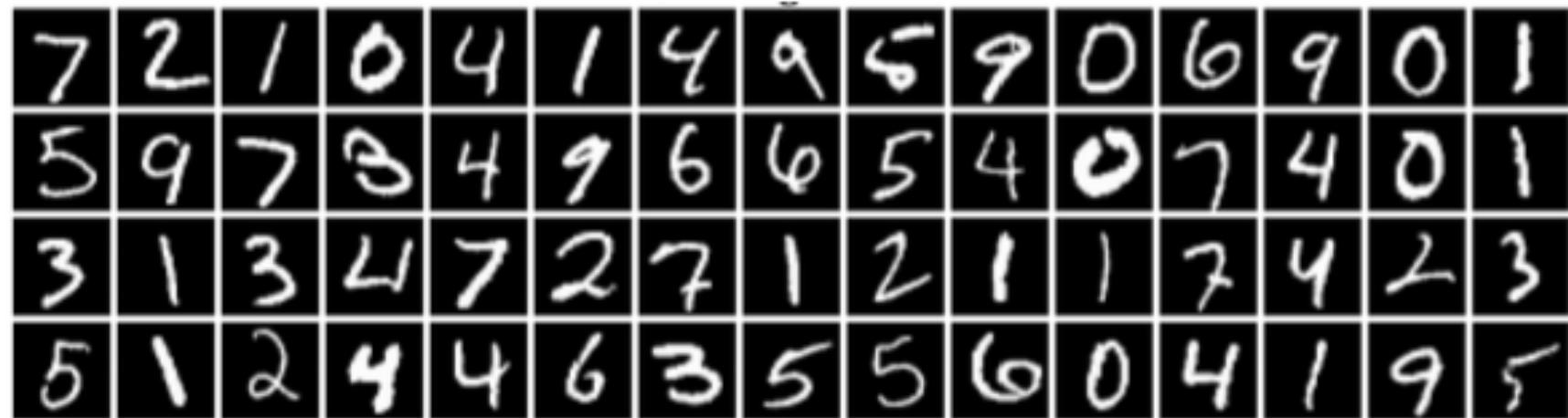
In this Lecture

- ▶ Autoencoders.
- ▶ Gaussian Mixture Models.
- ▶ Variational Autoencoders:
 - Construction.
 - ELBO.
 - Problems.

Idea

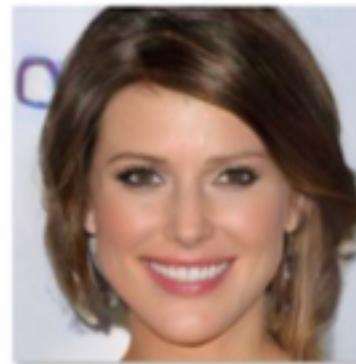
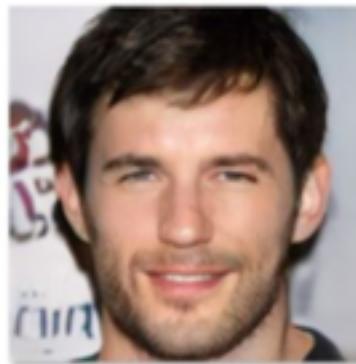
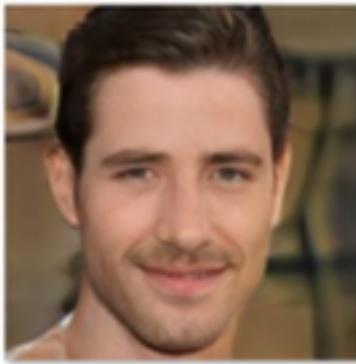


Latent variables - MNIST



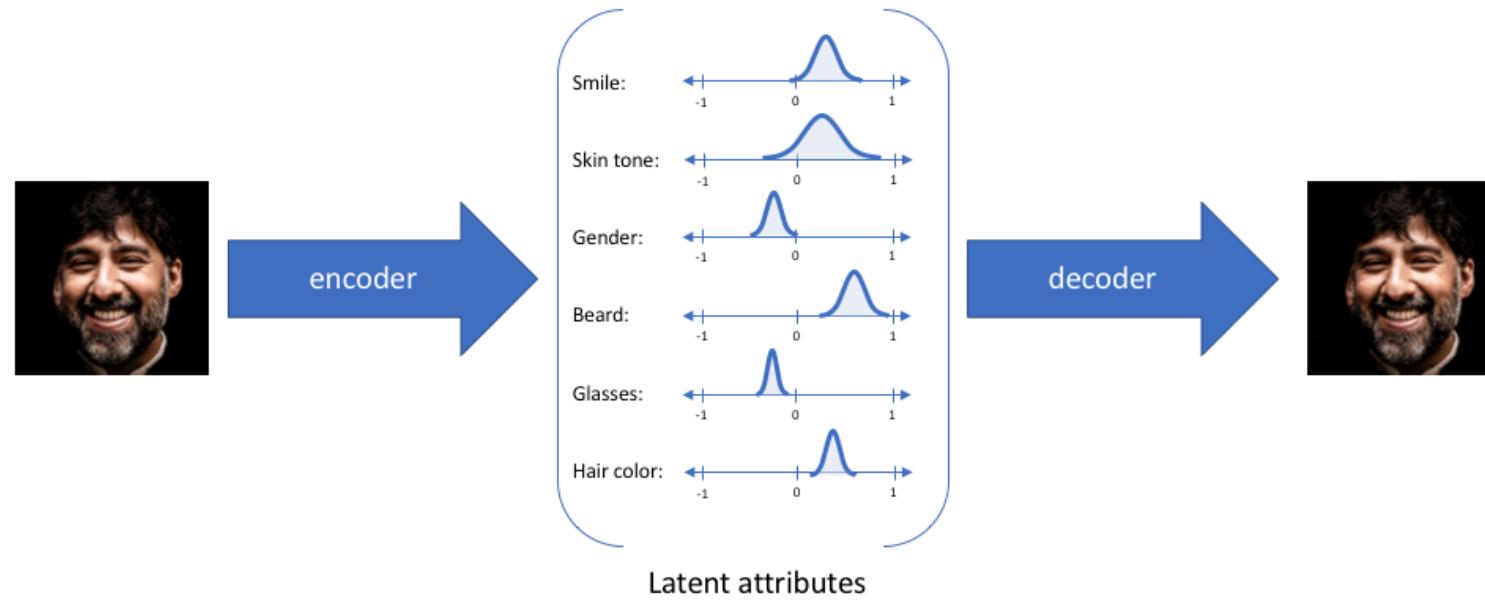
What is common in these images? How can they be characterized?

Latent Variables - Celeb



Variability due to eye and hair color, gender, race. However, unless images are annotated, these factors of variation are not explicitly available (latent).

Latent Variables - Celeb

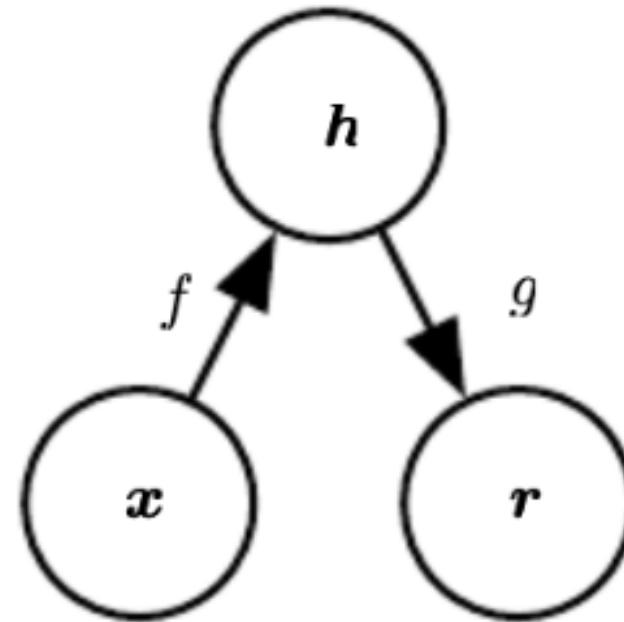


Variability due to eye and hair color, gender, race. However, unless images are annotated, these factors of variation are not explicitly available (latent).

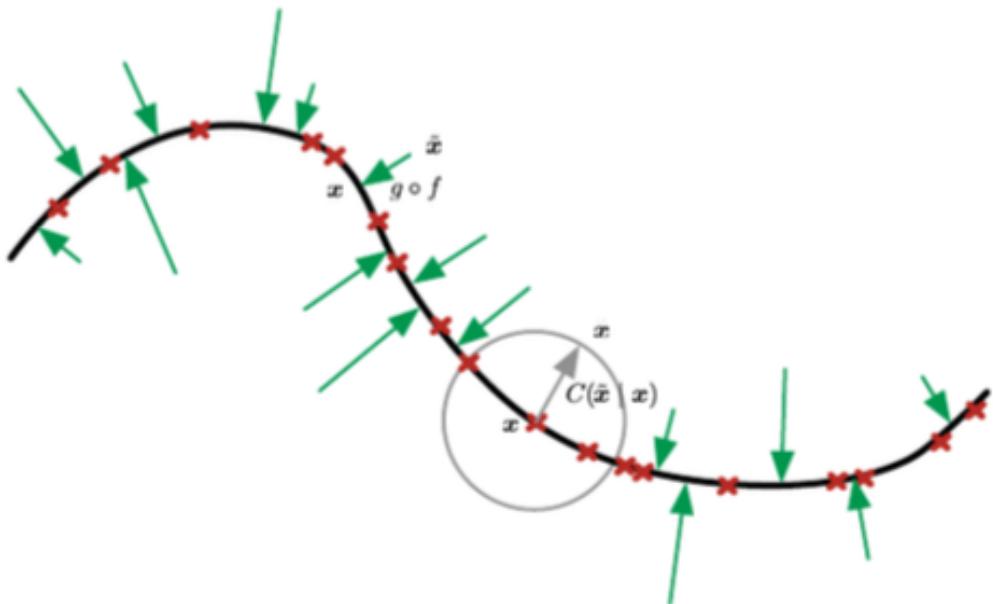
Reminder: Autoencoders

- ▶ Need a trick to reduce number of dimensions.
- ▶ Simplest idea: have a network with narrow hidden layer –
Autoencoder!
 - Encoder: $h = f(x)$;
 - Decoder: $r = g(h)$.
- ▶ Want to find transformation

$$g(f(x)) = x.$$



Contractive Denoising Autoencoders



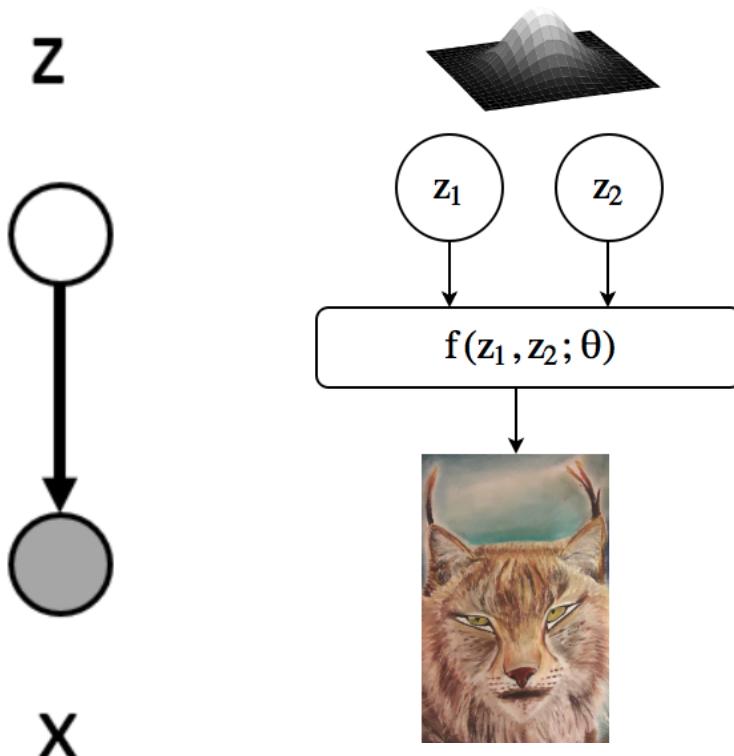
- ▶ The true signal is always situated on a manifold inside the R^D space.
- ▶ Denoising autoencoder is trained to map a corrupted data point \tilde{x} back to the original data point x .
- ▶ Still, produce lower-level results.

<https://arxiv.org/abs/1305.6663>

Latent Variable Models



Latent Variable Models: Motivation



- ▶ **Observe** X in the data sample.
Latent observables Z are **hidden**.
- ▶ Knowing latent variable space correctly can be beneficial as $p(x|z)$ can be described **simpler** than $p(x)$.
- ▶ **Hard** to find the real Z space manually, we need unsupervised learning.

Gaussian Mixture Model

We can use a mixture of distributions (like, Gaussians) to have a better estimate of true PDF. Gaussian Mixture Model proposes:

$$p_{GMM} = \sum_{l=1}^K \pi_l \phi(x; \mu_l, \sigma_l^2),$$

where $\pi_l \geq 0$ are the weights: $\sum \pi_l = 1$. Here the number K is a tuning parameter that specifies the number of Gaussians in our model.

Gaussian Mixture Model

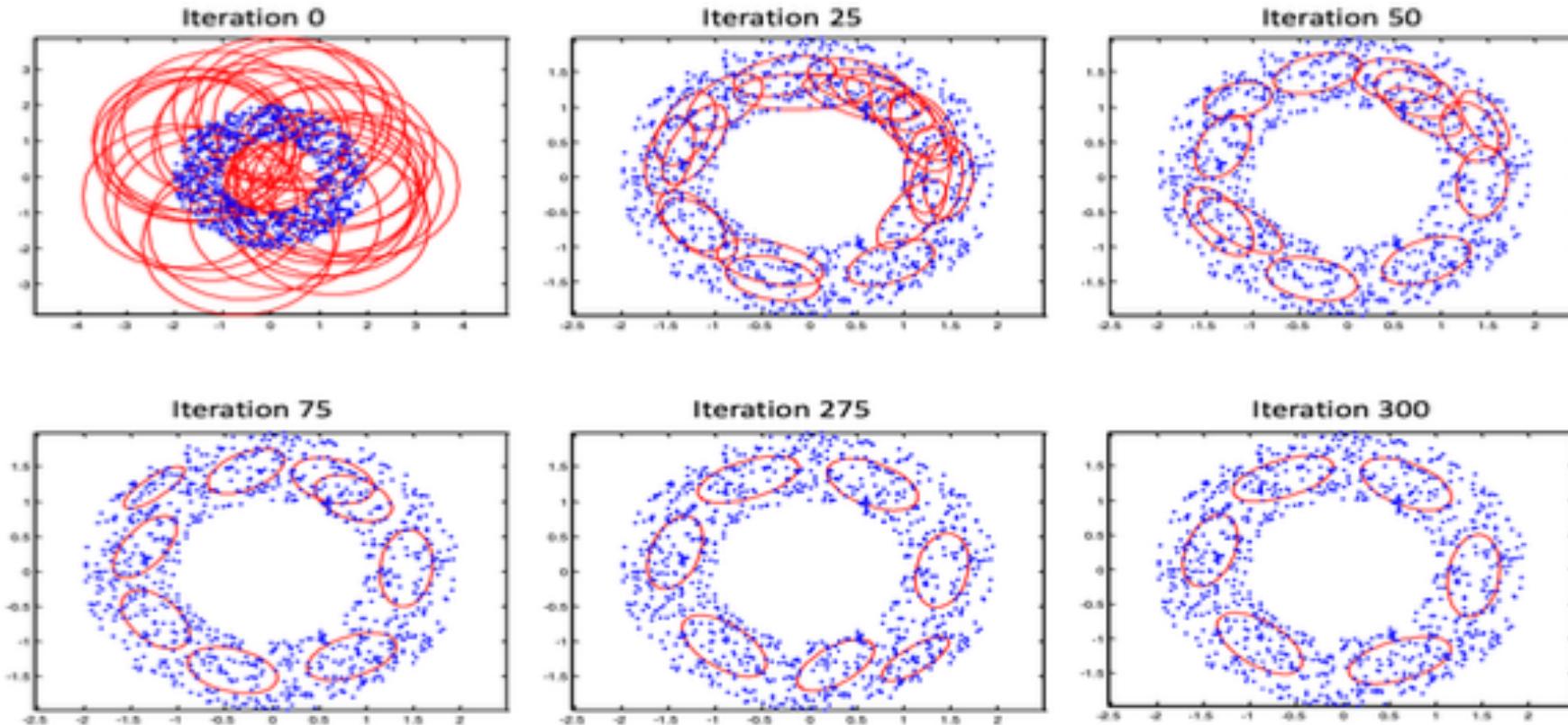
We thus have $3K - 1$ parameters. Estimation of them can be done using Maximum Likelihood Estimate:

$$\hat{\pi}_1, \hat{\mu}_1, \sigma_1^2, \dots, \sigma_K^2 = \arg \max_{\hat{\pi}_i, \hat{\mu}_i, \sigma_i^2} \sum_{i=1}^n \log \left(\sum_{i=1}^K \pi_i \phi(x; \mu_i, \sigma_i^2) \right)$$

Things get worth, when we go to multidimension, instead of σ^2 , we need to estimate matrix Σ .

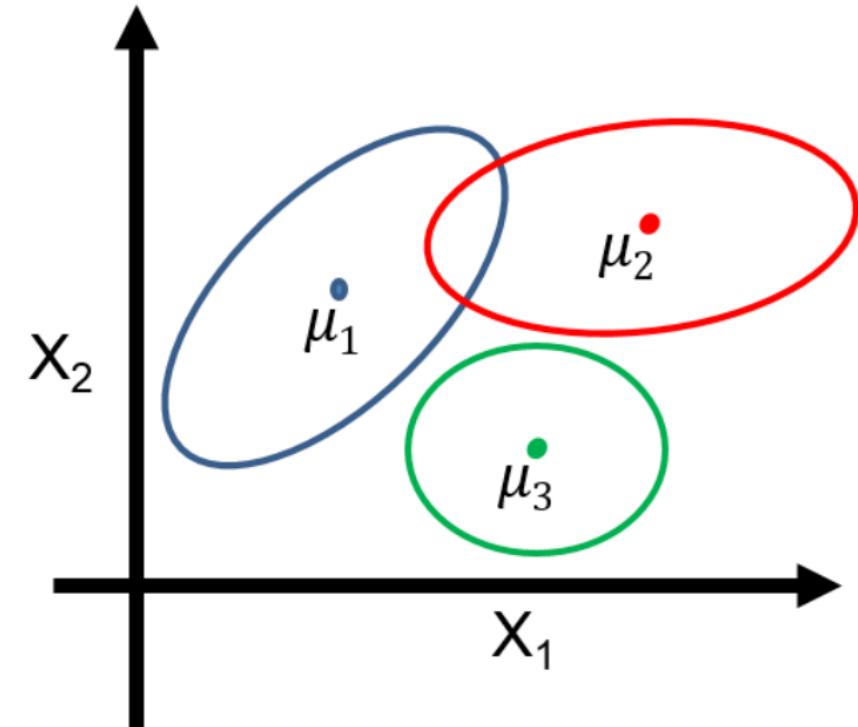
Gaussian Mixture Model

Training set: $n = 900$ examples from a uniform pdf inside an annulus,
model: GMM with $K = 30$ Gaussian components



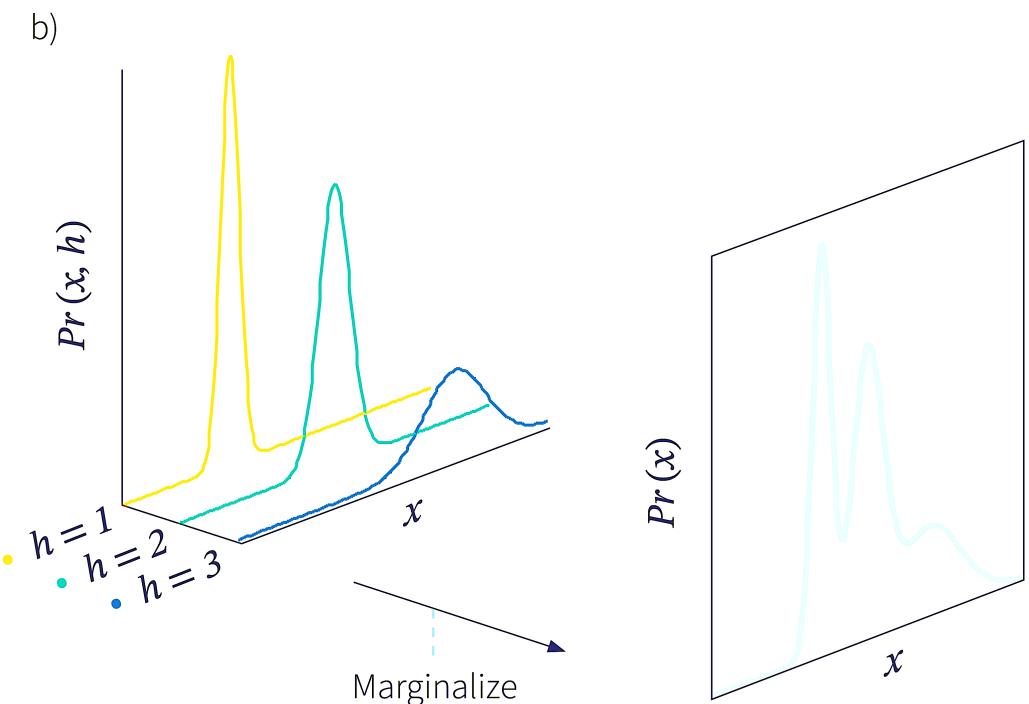
Mixture of Gaussians as Latent Variable Model

- ▶ GMM is a latent variable model.
 - $z \sim \text{Categorical}(1, K)$
 - $p(x|z = k) = N(\mu_k, \Sigma_k)$
- ▶ Generative process
 - Pick a mixture component k by sampling z .
 - Generate a data point by sampling from that Gaussian.



Combination of $p(x)$

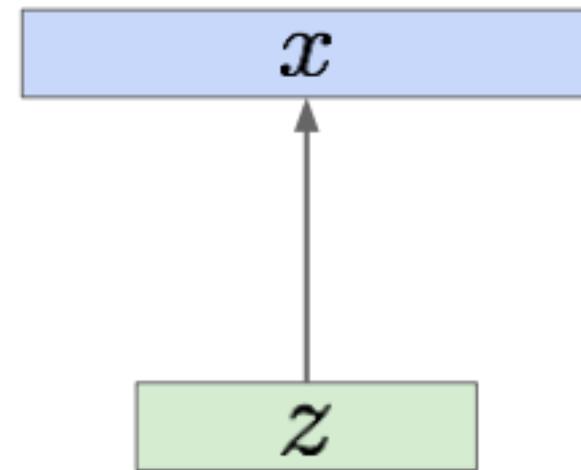
- ▶ Combine simple models into a more complex and expressive one.
- ▶ Obtain PDF using simple formula.



$$p(x) = \sum_{all z} p(x, z) = \sum_{all z} p(z)p(x|z) = \sum_{k=1}^K p(z = k)N(x; \mu_k, \Sigma_k)$$

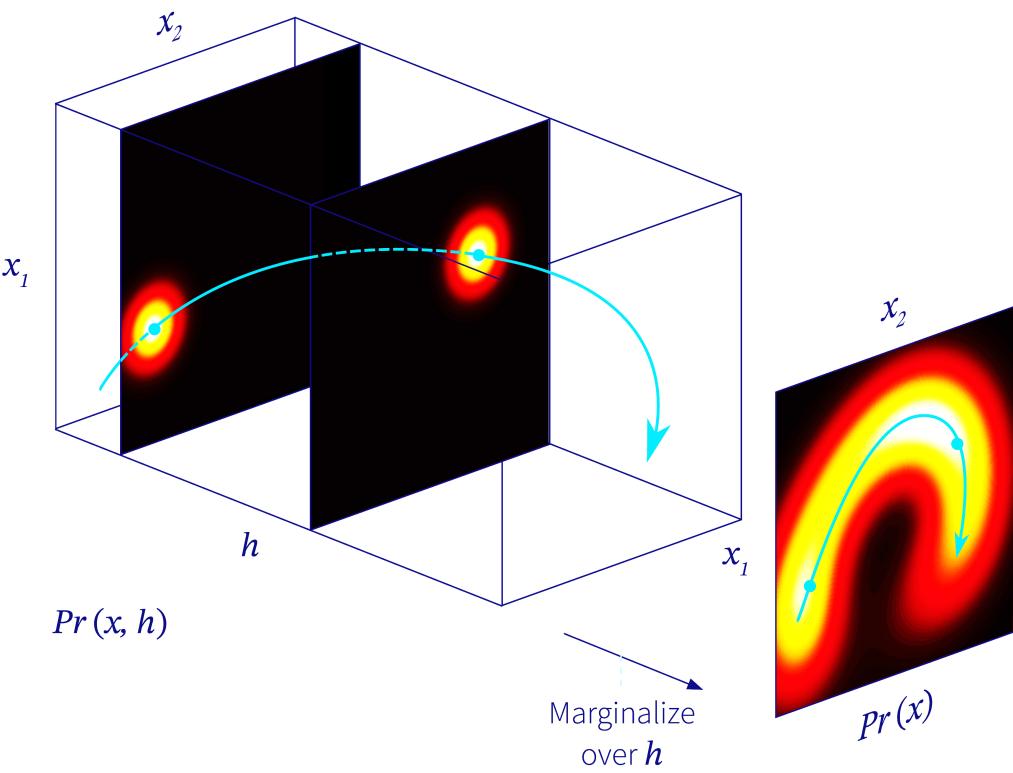
Continuous z

- ▶ Latent variable model.
 - $\mathbf{z} \sim p(\mathbf{z})$
 - $p(x|\mathbf{z} = \mathbf{z}^{(i)}) = N(\mu_{(i)}, \Sigma_{(i)})$
- ▶ Generative process
 - Pick z from prior.
 - Generate a data point by sampling from Gaussian.



Combination of $p(x)$

- ▶ Combine simple models into a more complex and expressive one.
- ▶ Obtain PDF using simple formula.
- ▶ What if we have a continuous observable?



$$p(x) = \sum_{\text{all } z} p(x, z) = \sum_{\text{all } z} p(z)p(x|z) = \sum_{k=1}^K p(z = k)N(x; \mu_k, \Sigma_k)$$

$$p(x; \theta) = \int_z p(x, z; \theta) dz = \int_z p(x; \theta|z)p(z) dz$$

Finding parameters θ

How do we find the parameters θ ? The most natural way is to maximise the log likelihood:

$$\log \prod_{i=1}^D p(x^i; \theta) = \sum_{x \in D} \log p(x; \theta)$$

In the following we will need to:

- › define a loss function to optimize
- › calculate gradients of the loss function with respect to the parameters θ .
- › apply a variant of stochastic gradient descent.

Maximizing the likelihood

Consider one sample, x . Then to maximize the likelihood, we need to calculate:

$$p(x; \theta) = \int_z p(x, z; \theta) dz = \int_z p(x; \theta|z)p(z) dz$$

- ▶ Easy to calculate in special cases:
 - in case of independent observables $x|z \perp z$.
 - special cases of conjugation.
- ▶ In general, x depends heavily on z (by construction).
- ▶ Integral is intractable.

Estimating likelihood: Naive Monte-Carlo

We can use the same algorithm as we have seen in GMM:

$$p(x; \theta) = \sum_{all z} p(x, z; \theta) = |Z| \mathbb{E}_{z \sim Uniform(z)}(p(x, z; \theta)).$$

For this, we can fix the distribution followed by z .

This is really only tractable when x is relatively low-dimensional. In case of image, we run into the curse of dimensionality, where we need to grab many samples to get an accurate view of x .

Wrap-up

- ▶ Problems:
 - We can't calculate $p(x)$.
 - We can't write the maximum-likelihood objective.
- ▶ We need a different approach or objective function. Idea:
 - if we can't write likelihood, let's instead derive a lower bound on it.
 - If the lower bound is tractable, then we can optimize the parameters with respect to the lower bound.
 - If we are making the lower bound larger, we are making the likelihood larger.

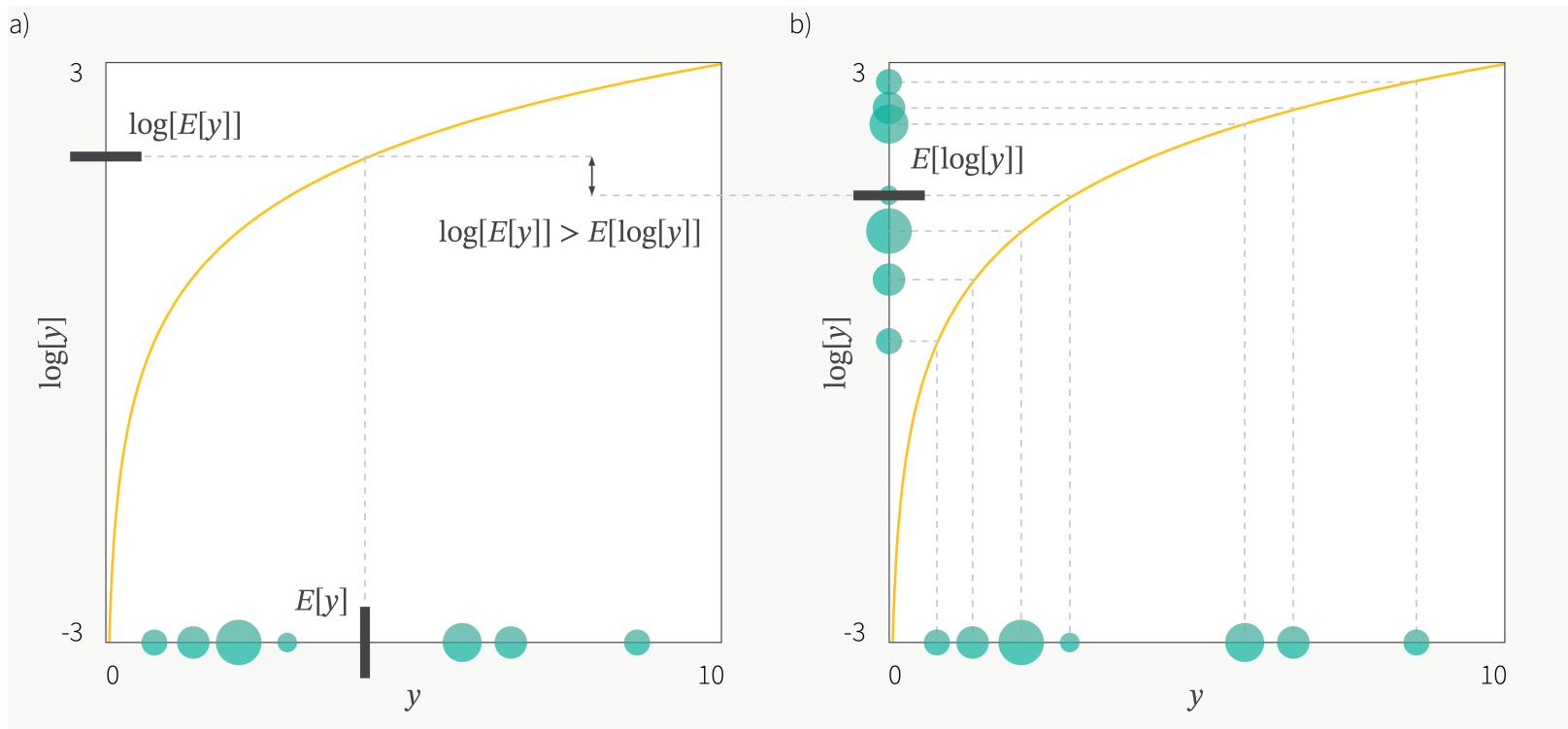
Evidence Lower Bound



Jensen's inequality

For concave function $g(\cdot)$:

$$g[\mathbb{E}[y]] \geq \mathbb{E}[g[y]].$$



$$\log \int p(y)y dy \geq \int p(y) \log y dy .$$

Figure: BorelisAI

Deriving the Evidence Lower BOund

- ▶ Start from arbitrary function multiplying and dividing $q(z)$:

$$\log p(x; \theta) = \log \int p(x, z; \theta) dz = \log \int q(z) \frac{p(x, z; \theta)}{q(z)} dz.$$

- ▶ Jensen's inequality for convex functions (\log):

$$\log \int q(z) \frac{p(x, z; \theta)}{q(z)} dz \geq \int q(z) \log \frac{p(x, z; \theta)}{q(z)} dz.$$

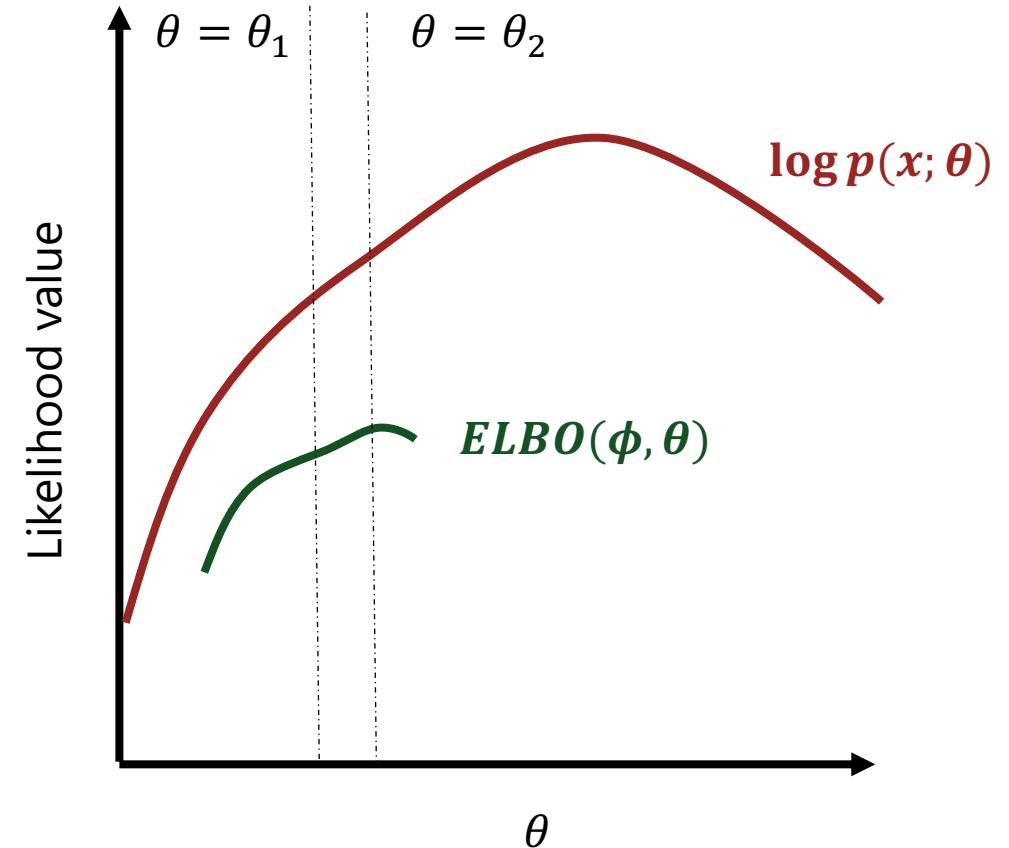
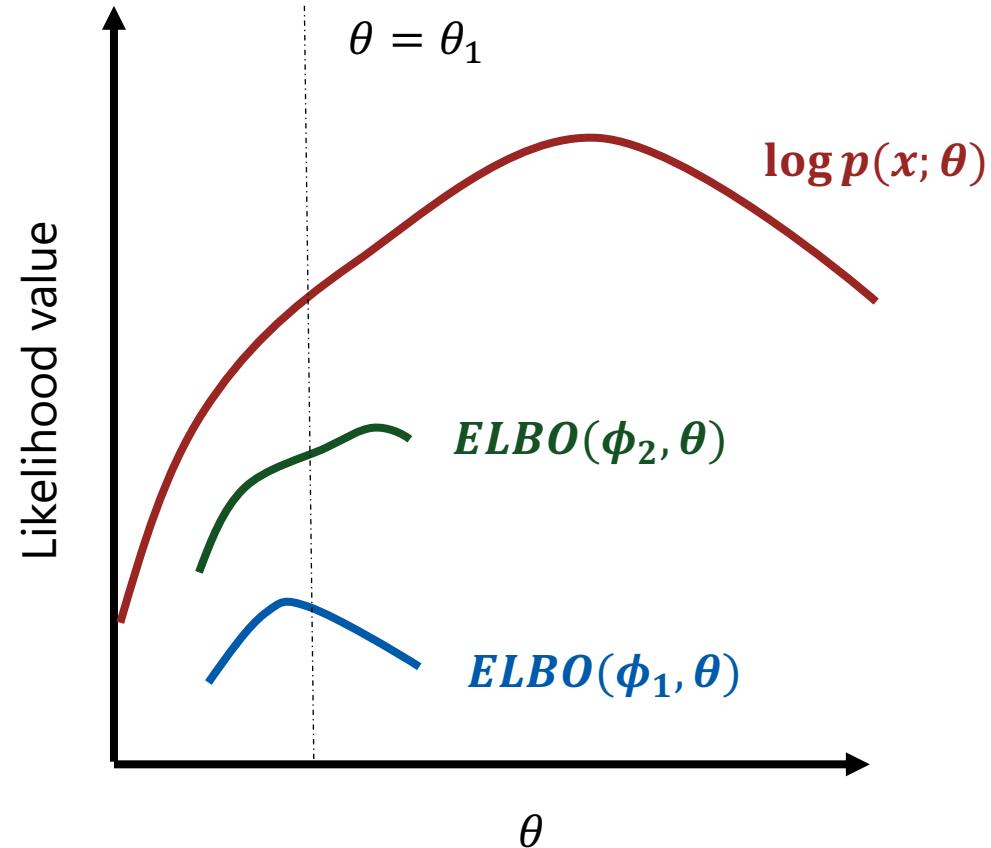
The rhs is **evidence lower bound**

- ▶ $q(z)$ has different set of parameters, ϕ

$$ELBO(\phi, \theta) = \int q(z; \phi) \log \frac{p(x, z; \theta)}{q(z; \phi)} dz$$

ELBO so far

$$ELBO(\phi, \theta) = \int q(z; \phi) \log \frac{p(x, z; \theta)}{q(z; \phi)} dz$$



Tightness of bound

$$ELBO(\phi, \theta) = \int q(z; \phi) \log \frac{p(x, z; \theta)}{q(z; \phi)} dz$$

- We use chain rule

$$ELBO(\phi, \theta) = \int q(z; \phi) \log \frac{p(x; \theta)p(z|x; \theta)}{q(z; \phi)} dz$$

- Separate part independent on z :

$$= \int q(z; \phi) \log p(x; \theta) dz + \int q(z; \phi) \log \frac{p(z|x; \theta)}{q(z; \phi)} dz =$$

$$= \log p(x; \theta) + \int q(z; \phi) \log \frac{p(z|x; \theta)}{q(z; \phi)} dz =$$

- Definition of KL divergence:

$$= \log p(x; \theta) - KL(q(z; \phi) || p(z|x; \theta)).$$

- Bound is tight iff $KL = 0$.

M. Hoffman, M. Johnson, ELBO surgery, NIPS2016

Variational Autoencoder



Dependence on prior on z

- ▶ We use chain rule

$$ELBO(\phi, \theta) = \int q(z; \phi) \log \frac{p(z)p(x|z; \theta)}{q(z; \phi)} dz =$$

- ▶ Separate KL divergence part:

$$= \int q(z; \phi) \log p(x|z; \theta) dz + \int q(z; \phi) \log \frac{p(z)}{q(z; \phi)} dz =$$

- ▶ Definition of KL

$$= \int q(z; \phi) \log p(x|z; \theta) dz - KL(q(z; \phi) || p(z)).$$

Reconstruction loss

Latent Parametrization agreement

Latent Parametrization Agreement

- ▶ Measures the degree to which the auxiliary distribution, $q(h, \theta)$, matches the prior $p(z)$:

$$KL(q(z; \phi) || p(z)).$$

- ▶ The KL divergence has closed form if both distributions are known parametric (e.g., Gaussian):

$$q(z; \phi) = N_z(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(z) = N_z(\mathbf{0}, \mathbf{I}).$$

- ▶ In this case:

$$KL[q(z|x, \theta), p(z)] = 12(Tr[\boldsymbol{\Sigma}] + \boldsymbol{\mu}^T \boldsymbol{\mu} - D - \log[\det[\boldsymbol{\Sigma}]])$$

D. Kingma et al. Auto-Encoding Variational Bayes, ICLR 2014

Latent Parametrization Agreement: training

- ▶ Measures the degree to which the auxiliary distribution, $q(h, \theta)$, matches the prior $p(z)$:

$$KL(q(z; \phi) || p(z)).$$

- ▶ The KL divergence has closed form if both distributions are known parametric (e.g., Gaussian):

$$q(z|x; \phi) = N_z(g_{\mu}(x; \phi), g_{\sigma}(x; \phi));$$

$$p(z) = N_z(\mathbf{0}, I).$$

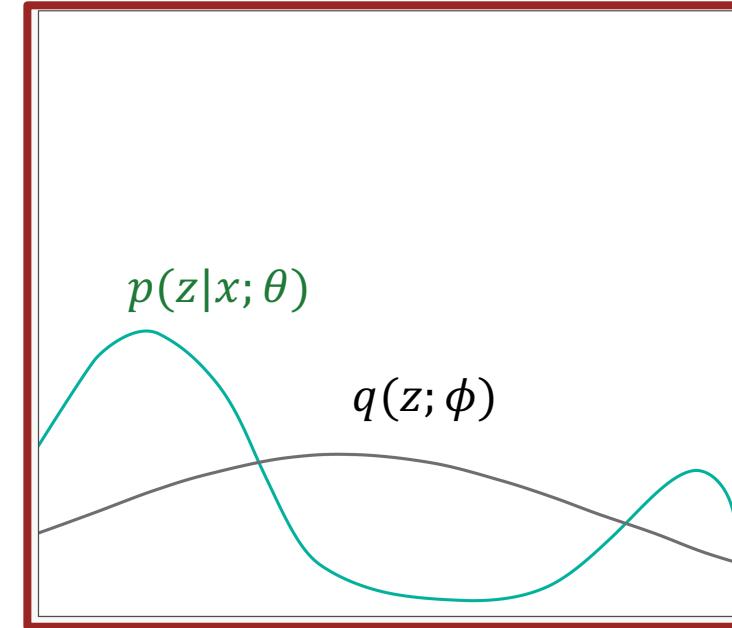
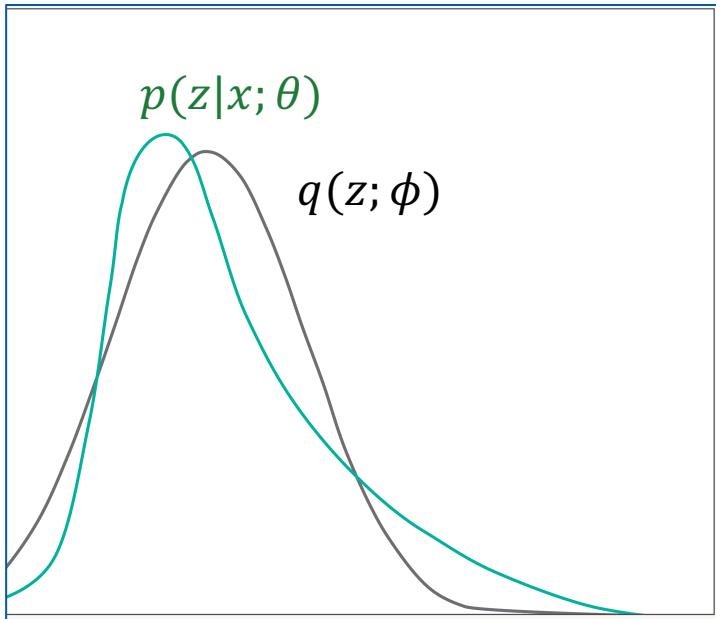
- ▶ In this case:

$$KL[q(z|x, \theta), p(z)] = \frac{1}{2} (Tr[\Sigma] + \mu^T \mu - D - \log[\det[\Sigma]])$$

Problems of parametrization choice

- ▶ For tightness of bound:

$$KL(q(z; \phi) || p(z|x; \theta)) \rightarrow 0.$$



- ▶ The tight bound might be unreachable.

Figure: BorelisAI

Reconstruction Loss

- ▶ Reconstruction loss:

$$\int q(z|x; \phi) \log p(x|z; \theta) dz.$$

- ▶ Still troubles to compute directly, we can approximate it with mathematical expectation (even with one sample z^*) and thus have:

$$\begin{aligned} ELBO(\theta, \phi) &= \int q(z; \phi) \log p(x|z; \theta) dz - KL(q(z; \phi) || p(z)) \approx \\ &\approx \log[p(x|z^*; \phi)] - KL[q(z|x, \theta), p(z)]. \end{aligned}$$

Variational Autoencoder

Maximize

$$\mathcal{L}_{VAE} = \mathbb{E}_z \log[p(x|z; \theta)] - KL[q(z|x, \phi), p(z)]$$

Autoencoder, as we reconstruct $x \rightarrow z \rightarrow x$.

Variational, as it constructs Gaussian approximation to the posterior distribution.

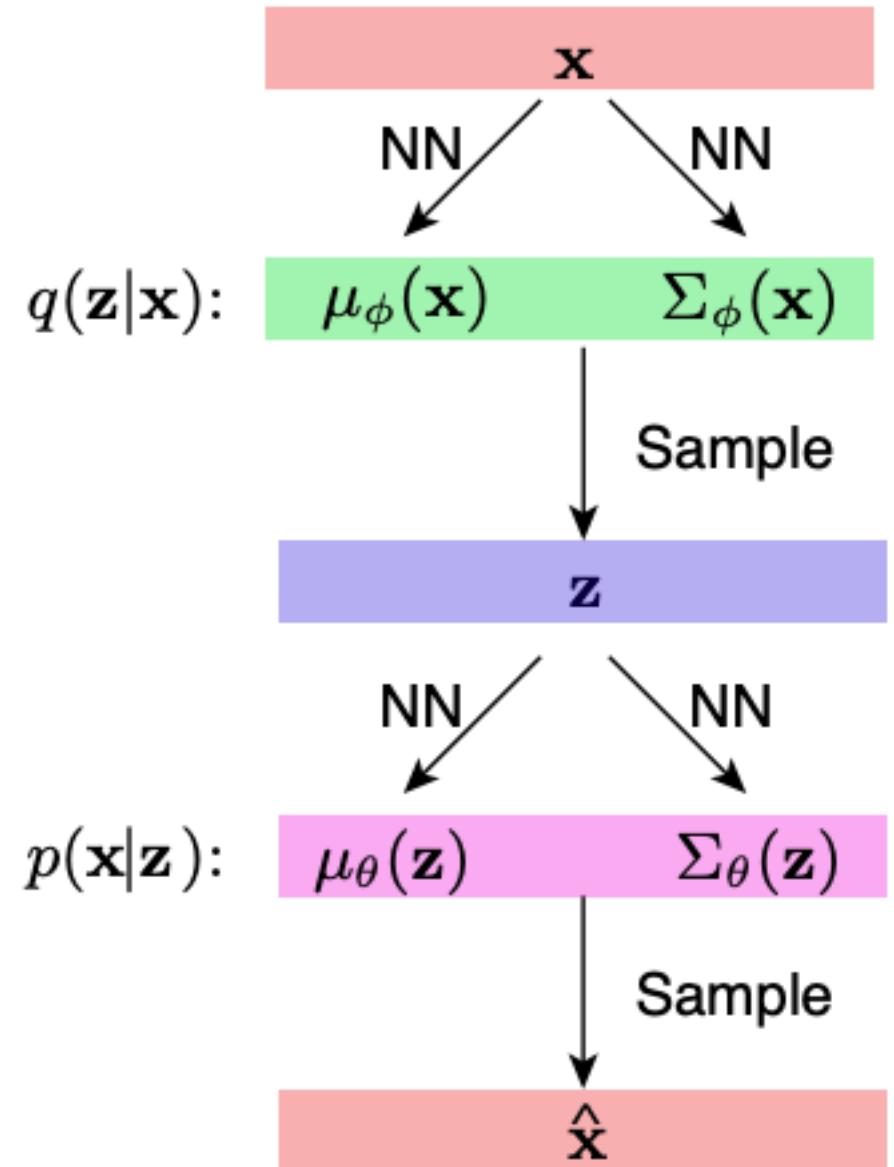


Fig: Kao UCLA DL lectures

Training VAE

We can estimate gradients for ELBO:

- › gradient wrt θ is easy:

$$\nabla_{\theta} \mathcal{L}(x; \theta, \phi) = \nabla_{\theta} \log p(x, z; \theta)$$

and can be done using a simple MC estimation.

- › for ϕ situation is much more tricky, since we calculate expectation over q :

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

luckily, we can use reparameterization.

Reparameterization trick

- ▶ Sampling:

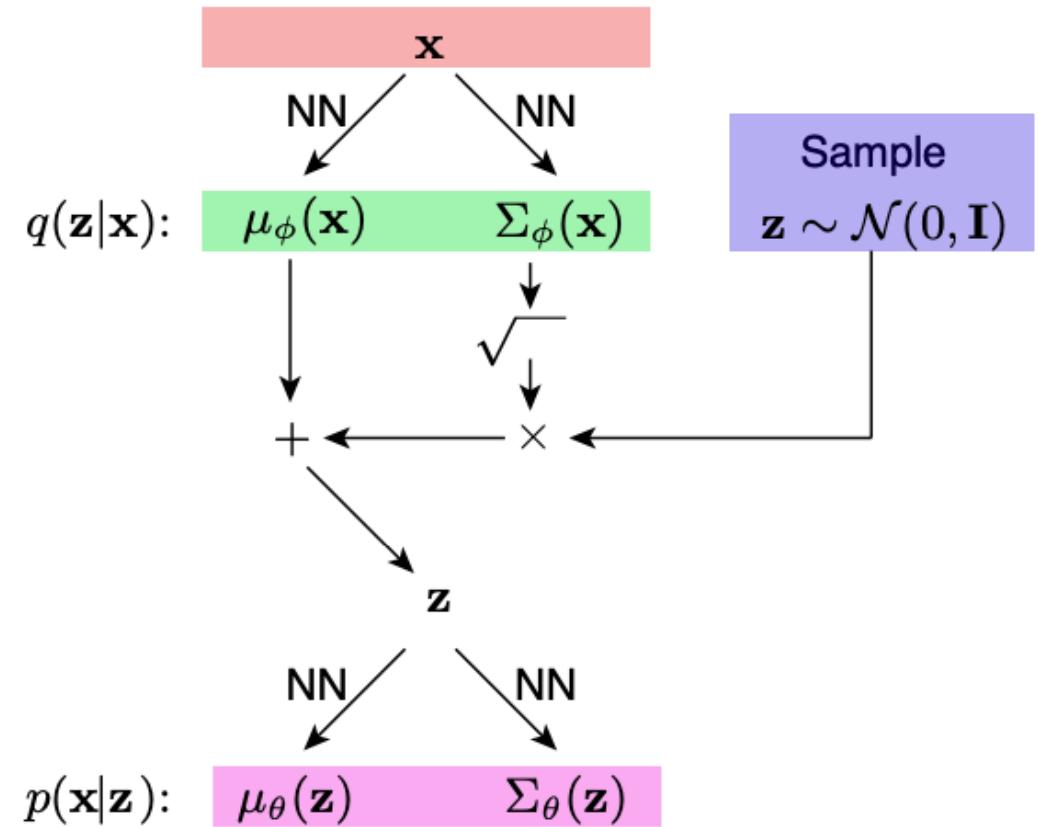
$$z \sim q(z|x)$$

- ▶ Can be done:

$$\varepsilon \sim N(0, I);$$

$$z = \mu_\phi(x) + \Sigma^{\frac{1}{2}}(x)\varepsilon.$$

- ▶ Backpropagation now does not pass through sampling.



Final Per-event Objective and Results

- ▶ Write out

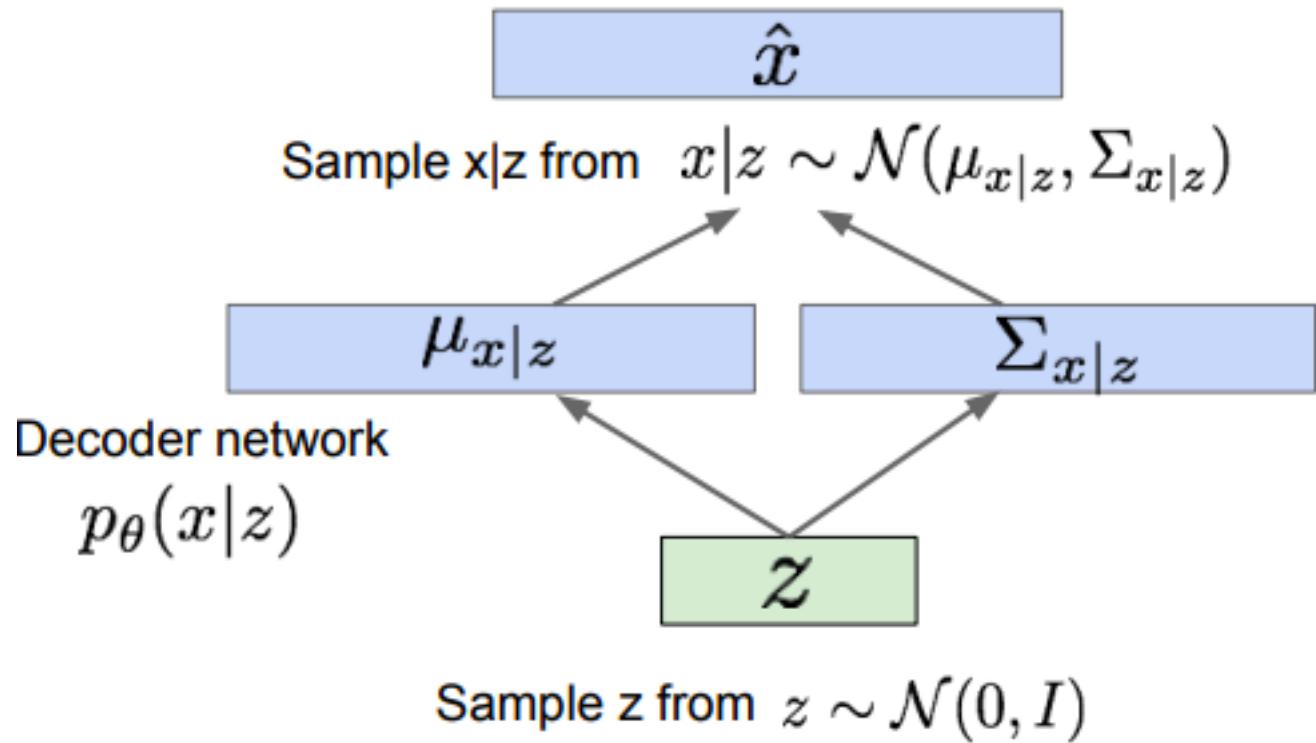
$$\mathcal{L}_{VAE}(x_i; \theta, \phi) = \mathbb{E}_{z \sim p(\varepsilon)} \log[p(x_i|z; \theta)] - KL[q(z|x_i, \phi), p(z)]$$

- ▶ We can even use a single noise sample:

$$\mathcal{L}_{VAE}(x_i; \theta, \phi) = \log[p(x_i|z; \theta)] - KL[q(z|x_i, \phi), p(z)]$$

- ▶ Note that the gradient estimates from per-event ELBO will be unbiased.

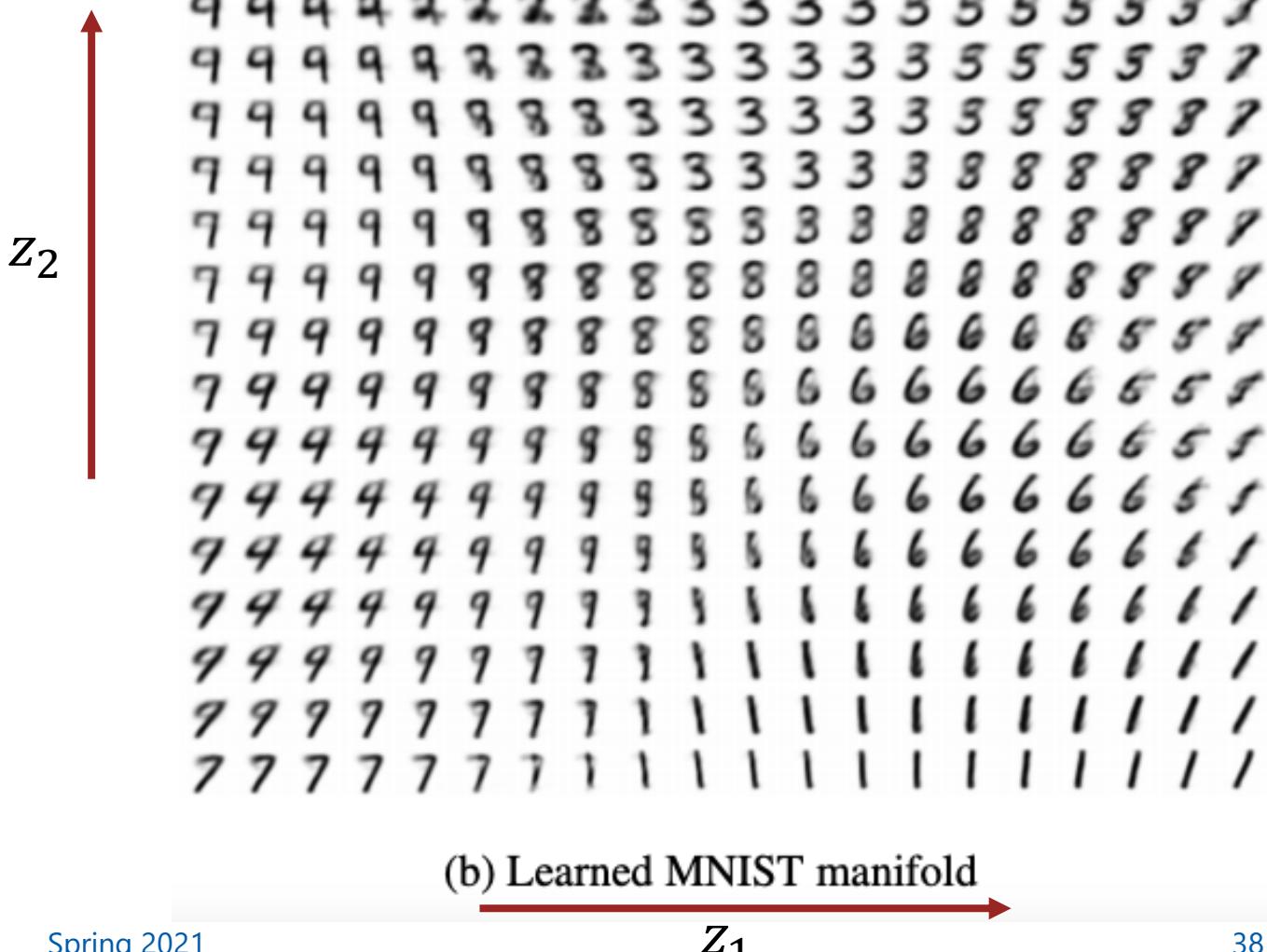
Sampling



- ▶ We now can use the same idea like in GMM.
- ▶ Sampling from prior.

Results

- ▶ The latent space is well organized.
- ▶ The results are good.



Results

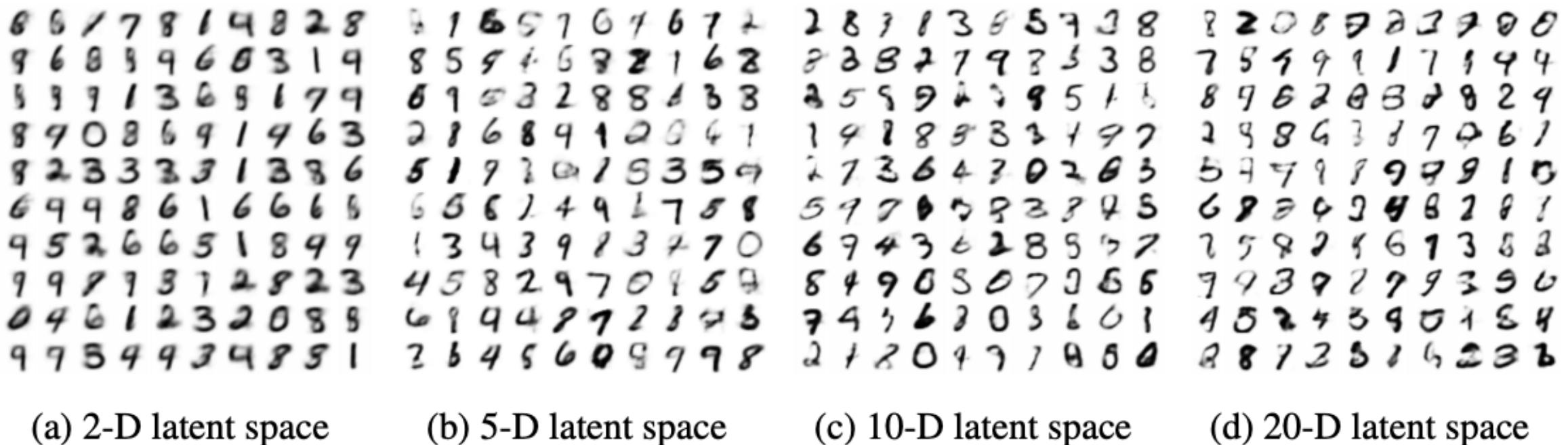
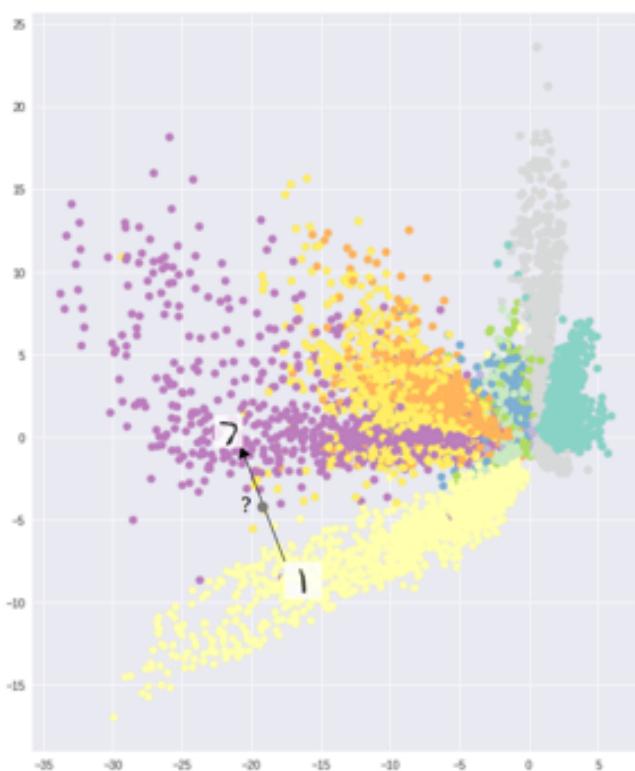


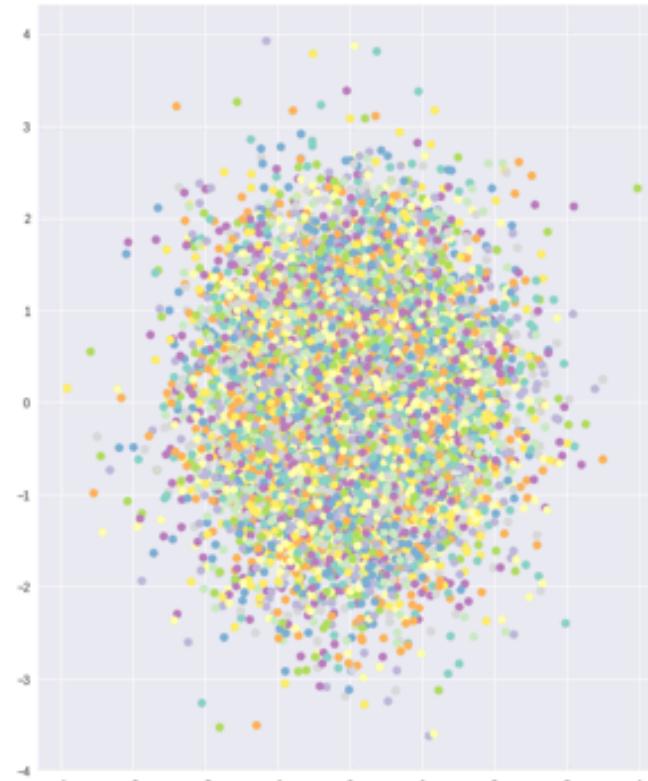
Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

Results Discussed

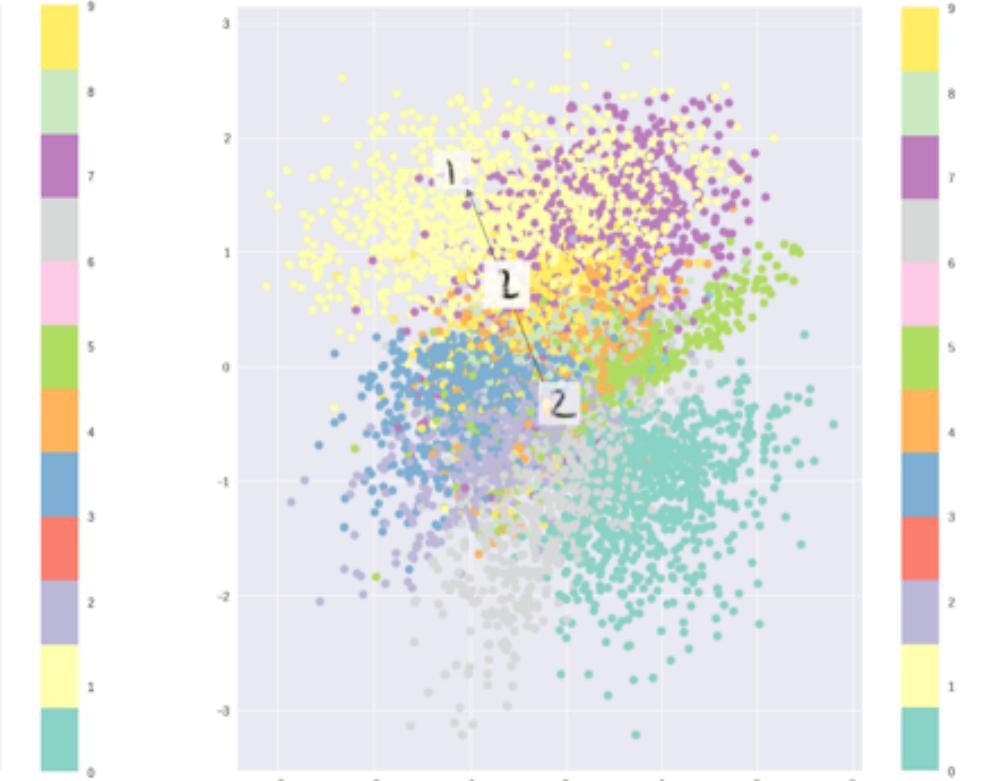
Only reconstruction loss



Only KL divergence



Combination



[Shafkat's blog](#)

Amortized Inference

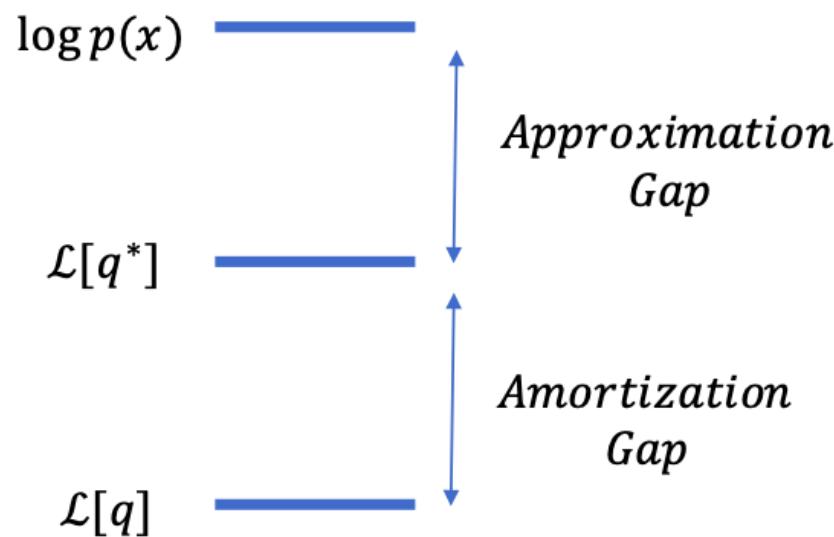


Figure 1. Gaps in Inference

- ▶ variational parameters ϕ_i should be calculated for every data point x_i ;
- ▶ for big datasets this is unsustainable;
- ▶ forget about it and make a single model.

Cremer et al Inference Suboptimality in Variational Autoencoders, ICML 2018

Images for Amortization Gap

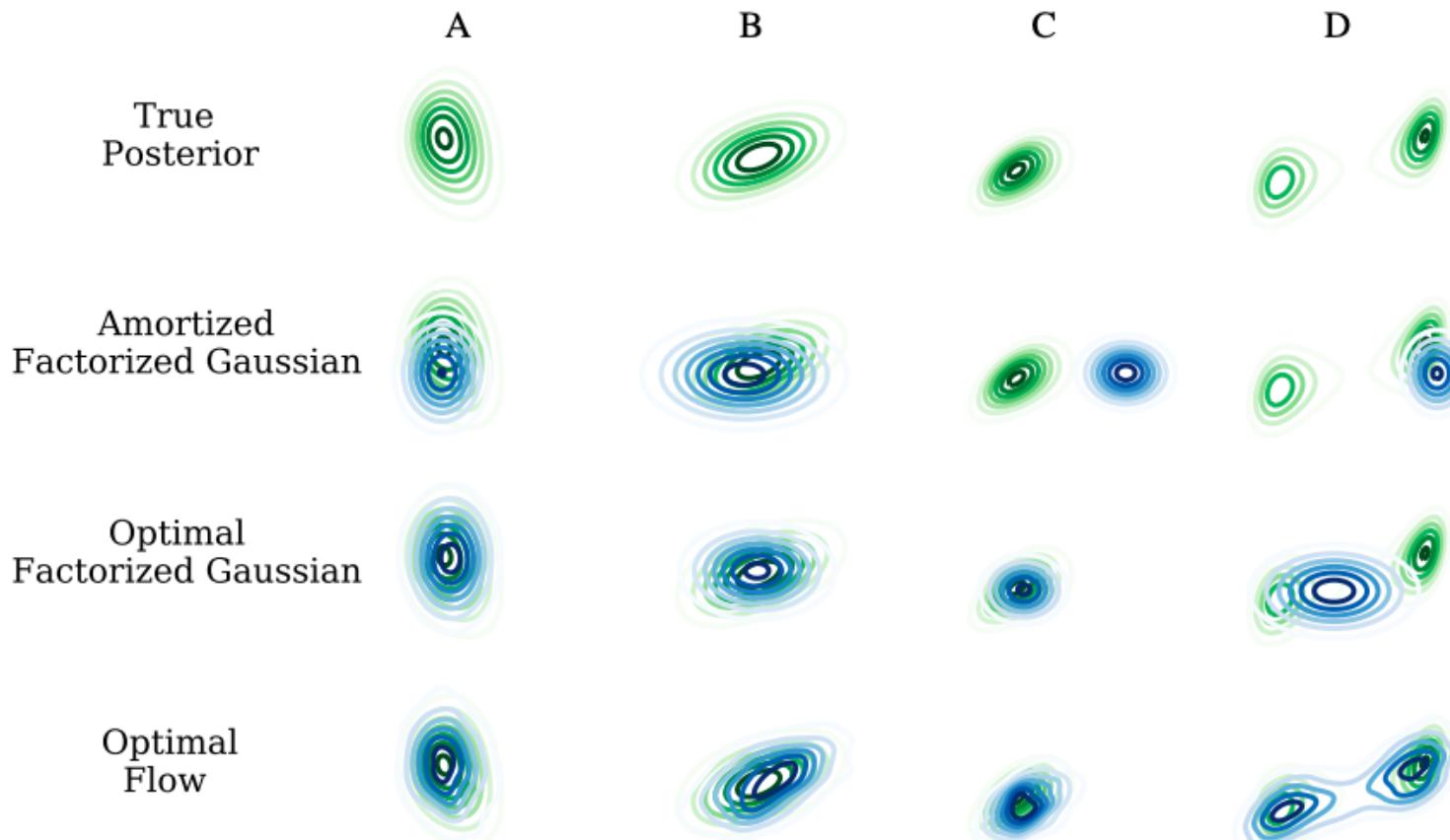


Figure 2. True Posterior and Approximate Distributions of a VAE with 2D latent space. The columns represent four different datapoints. The green distributions are the true posterior distributions, highlighting the mismatch with the blue approximations. Amortized: Variational parameters learned over the entire dataset. Optimal: Variational parameters optimized for each individual datapoint. Flow: Using a flexible approximate distribution.

VAE: starting optimization problem

- ▶ Stochastic optimization can get stuck:
 - in the beginning, likelihood term $\log p(x|z)$ is relatively weak, such that an initially attractive state is where $q(z|x) \approx p(z)$.
- ▶ Several ideas possible:
 - optimization schedule - gradually increase weight of KL.
 - free bits - ensure that on average, a certain minimum number of bits of information are encoded per latent variable.

VAE: posterior collapse



Can happen to weak encoders or noisy data.

The sampling does not depend on the latent model.

An averaged figures are produced.

KL-divergence term annealing is proposed.

Fig: J. Chou Generated Loss and Augmented Training of MNIST VAE

VAE: blurred problem

Due to minimization of KL
 $KL(q(z; \phi) \| p(z|x; \theta)) \rightarrow 0.$

Can be countered by choosing a sufficiently flexible inference model, and/or a sufficiently flexible generative model.



Some applications



VAE: Application

- ▶ Representation learning: learning better representations of the data.
Some uses of this are:
 - Data-efficient learning, such as semi-supervised learning;
 - Visualization of data as low-dimensional manifolds.
- ▶ Artificial creativity: plausible interpolation between data and extrapolation from data.

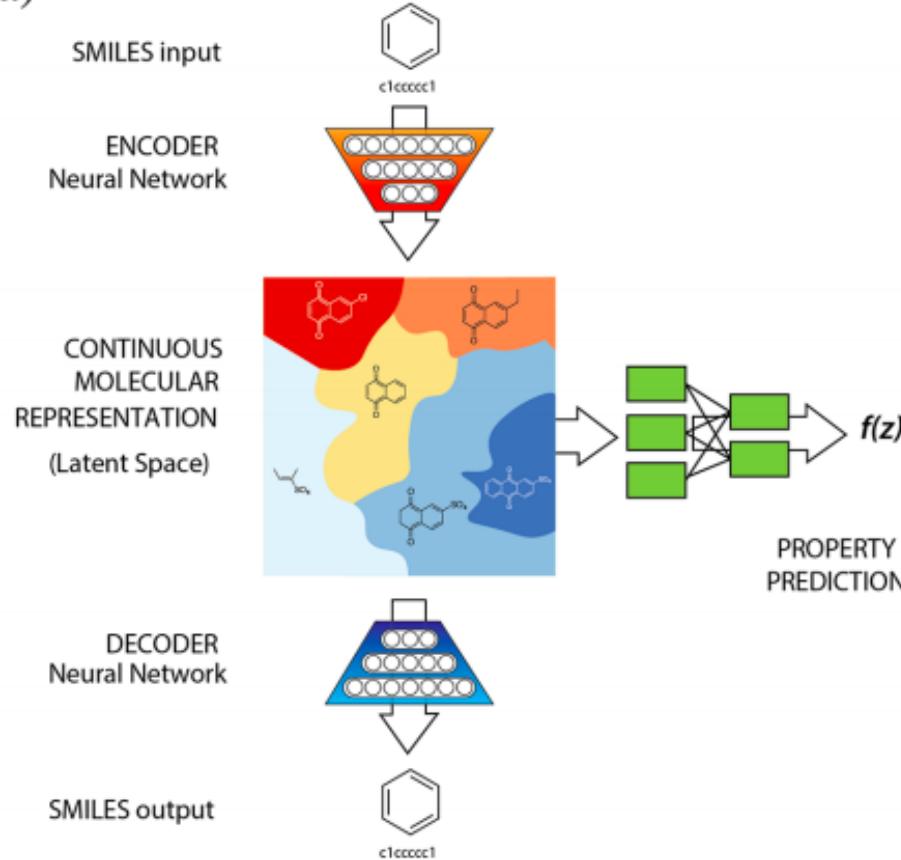
L. Maaløe et al. Auxiliary Deep Generative Models

R. Gómez-Bombarelli et al. Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

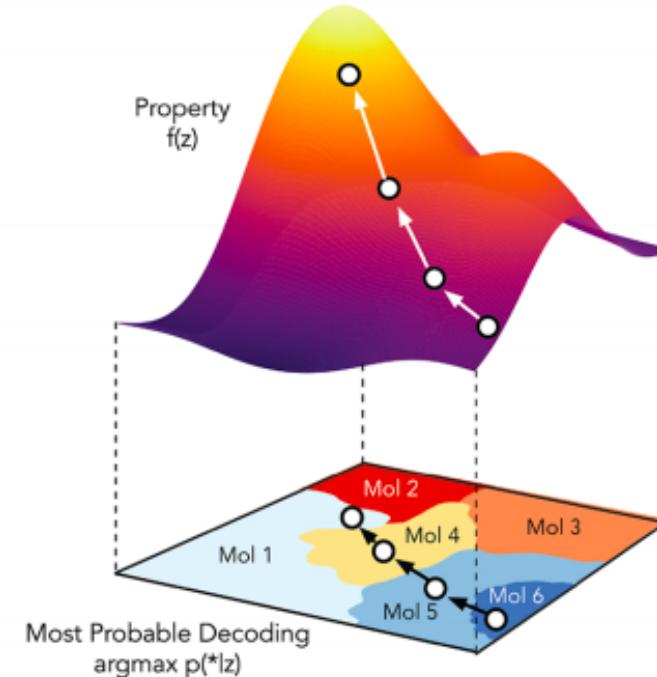
K. Gregor et al. DRAW: A Recurrent Neural Network For Image Generation

VAE: Application

(a)



(b)



- › interested in latent space;
- › look for the representation that has the best property $f(z)$;
- › generates a new molecule with best property.

Conclusions

- ▶ VAE is a powerful generation model
 - Individual images worse than GAN.
- ▶ Latent representations can be used as a separate instrument.
 - Many application in the real life.
- ▶ Next:
 - Merge good points of both approaches.