

Lecture 4: Variational AutoEncoder

Deep Generative Models

Sajjad Amini

Department of Electrical Engineering
Sharif University of Technology

Contents

- 1 Gaussian Mixture Model
- 2 Model Specification
- 3 Casting Likelihood Calculation as Expectation
- 4 Importance Sampling
- 5 Evidence Lower BOund
- 6 ELBO Tightening
- 7 ELBO Optimization
- 8 Reparameterization Trick
- 9 Learning
- 10 Amortization
- 11 Revisiting Objective Function
- 12 Attribute Vectors in Code Space

Section 1

Guassian Mixture Model

Example - Iris Flower



petal sepal

(a) Setosa



petal sepal

(b) Versicolor



petal sepal

(c) Virginica

Figure: Different types of Iris flower

Example - Iris Flower Data set

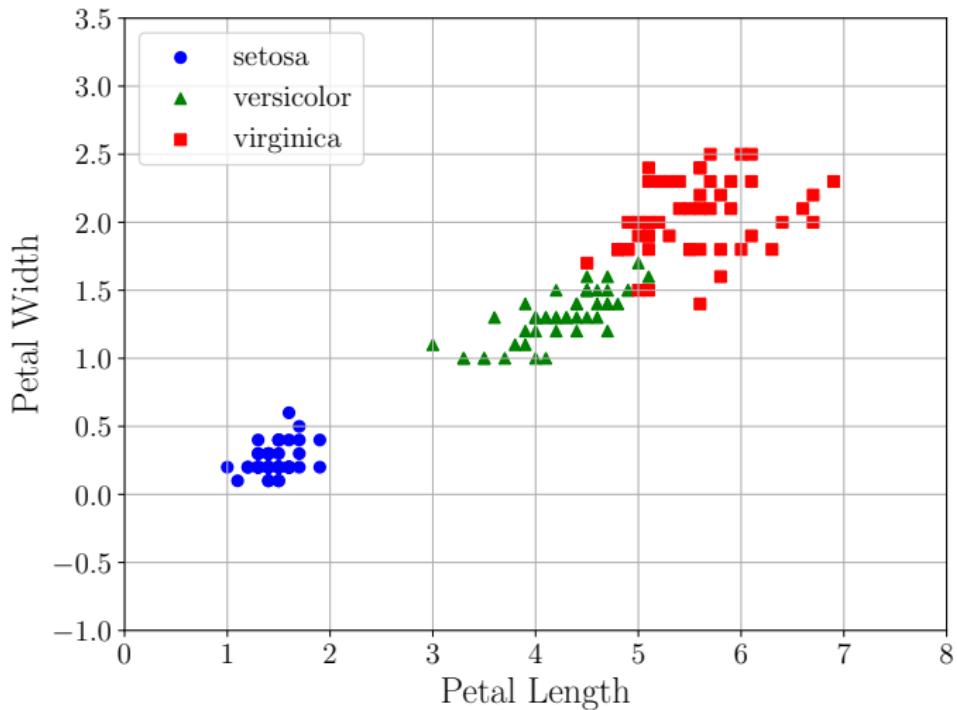


Figure: Labeled Iris dataset (in practice, you don't have the labels)

Example - Iris Flower Data set

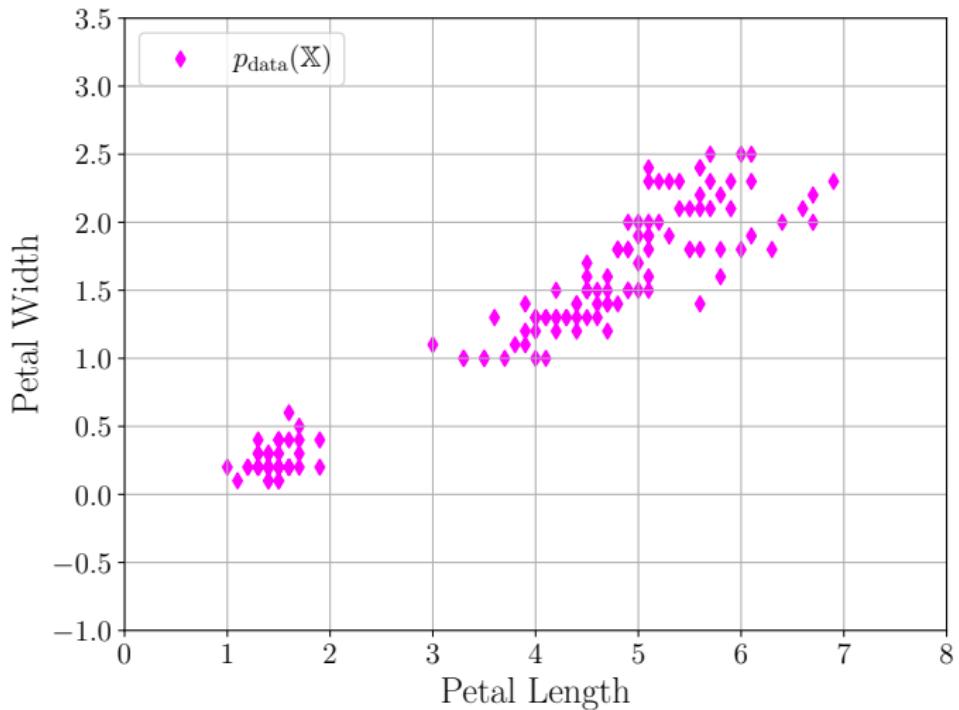


Figure: Iris dataset \mathcal{D}

Example - Iris Flower Data set

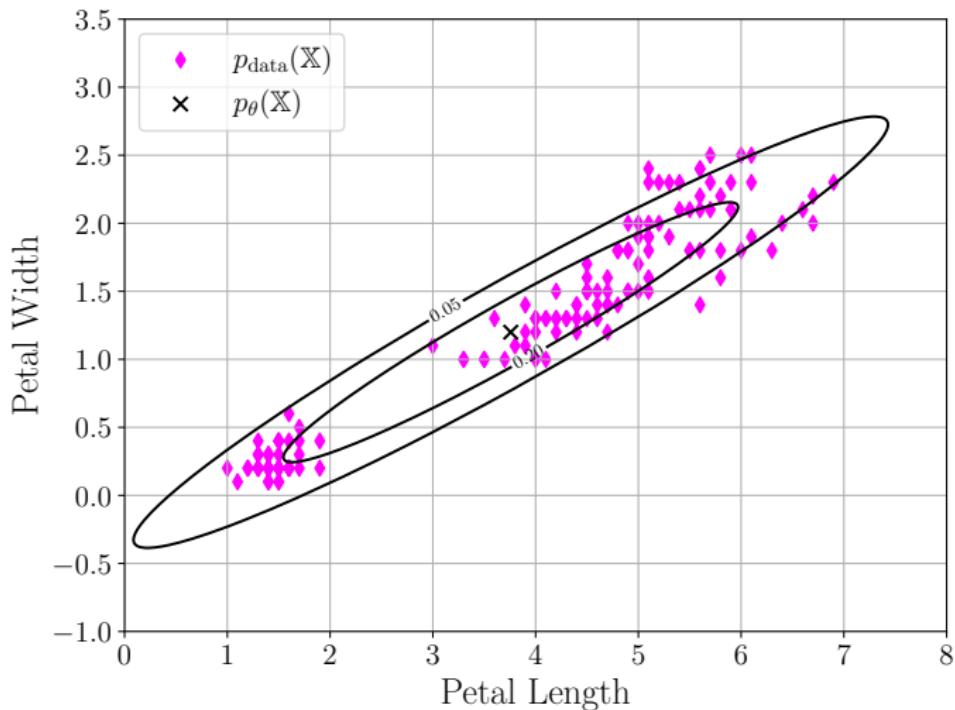


Figure: Modeling distribution using single Gaussian

Example - Iris Flower Data set

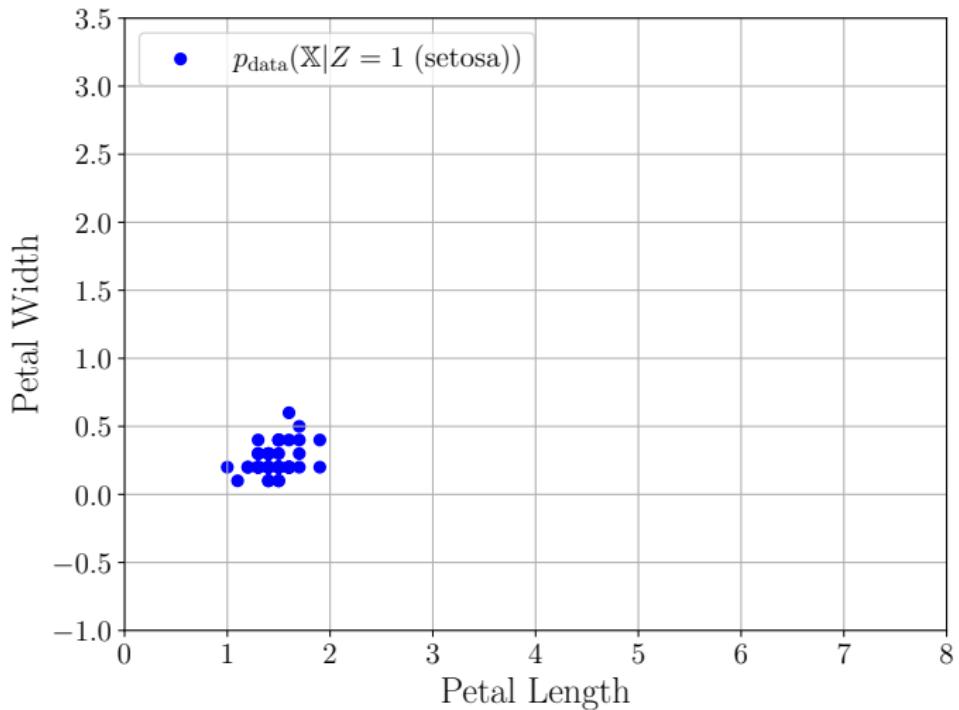


Figure: Data from $p_{\text{data}}(\mathbb{X}|Z = 1)$

Example - Iris Flower Data set

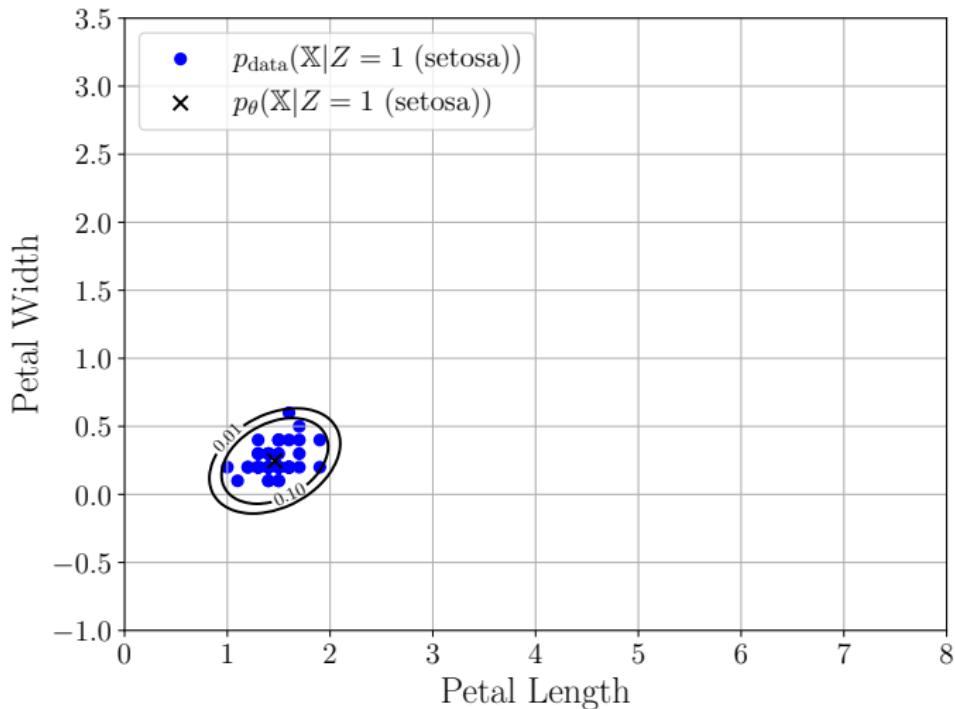


Figure: Modeling conditional distribution $p_{\text{data}}(\mathbb{X}|Z = 1)$ using signlne Gaussian

Example - Iris Flower Data set

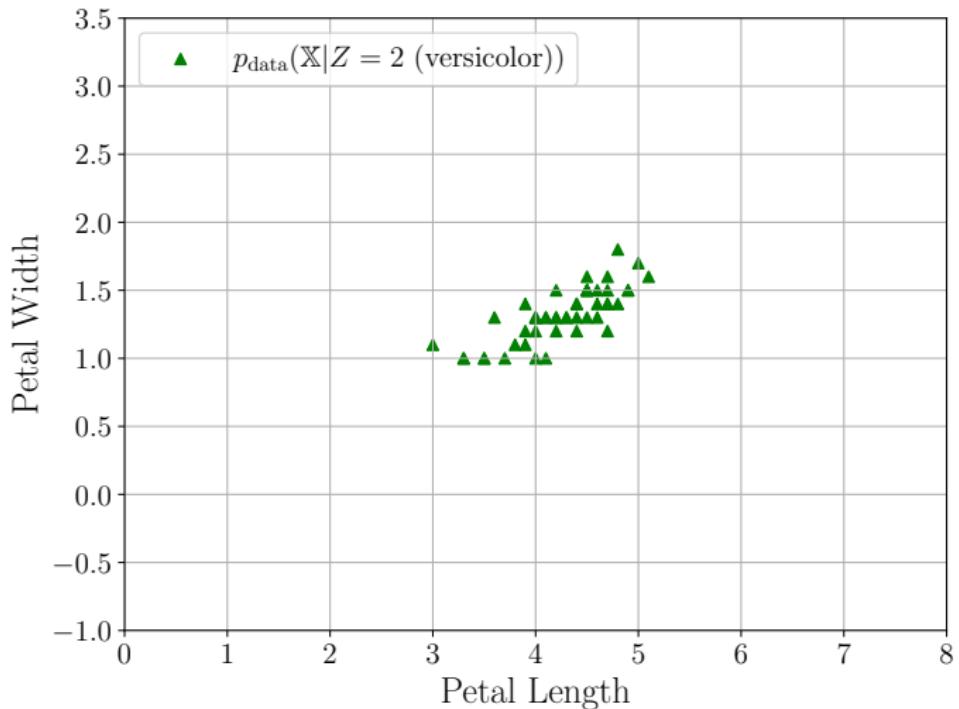


Figure: Data from $p_{\text{data}}(\mathbb{X}|Z = 2)$

Example - Iris Flower Data set

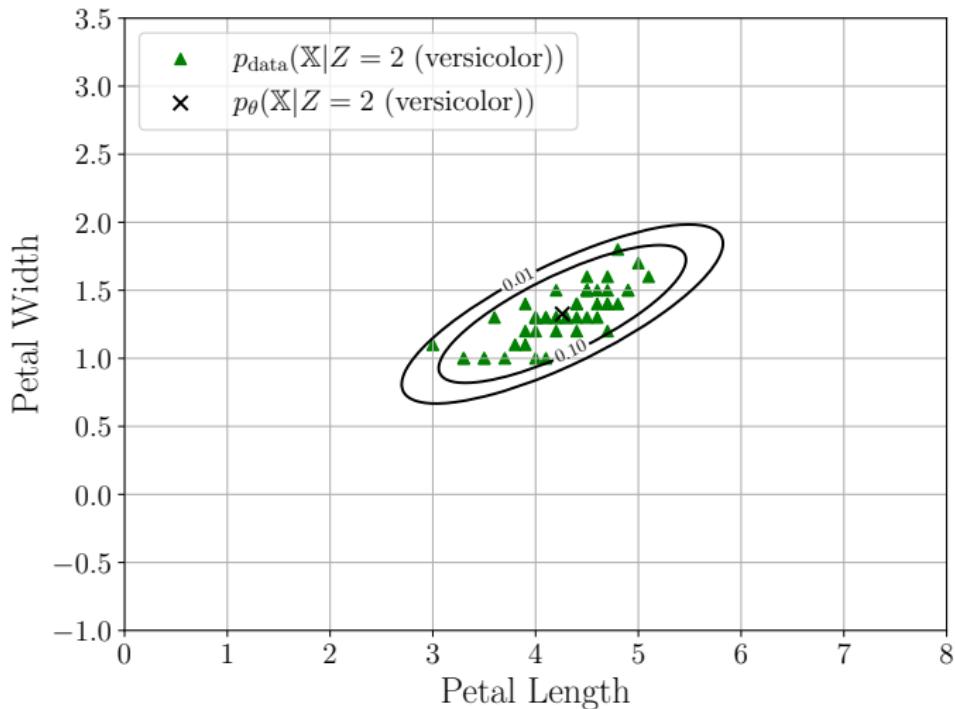


Figure: Modeling conditional distribution $p_{\text{data}}(\mathbb{X}|Z = 2)$ using signl Gaussian

Example - Iris Flower Data set

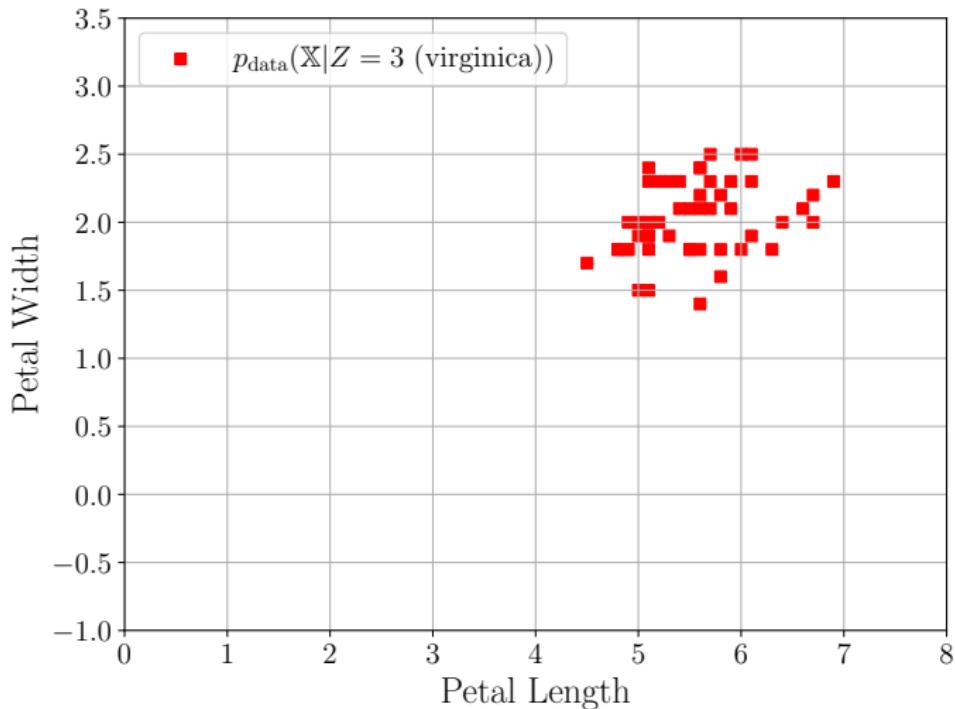


Figure: Data from $p_{\text{data}}(\mathbb{X}|Z=3)$

Example - Iris Flower Data set

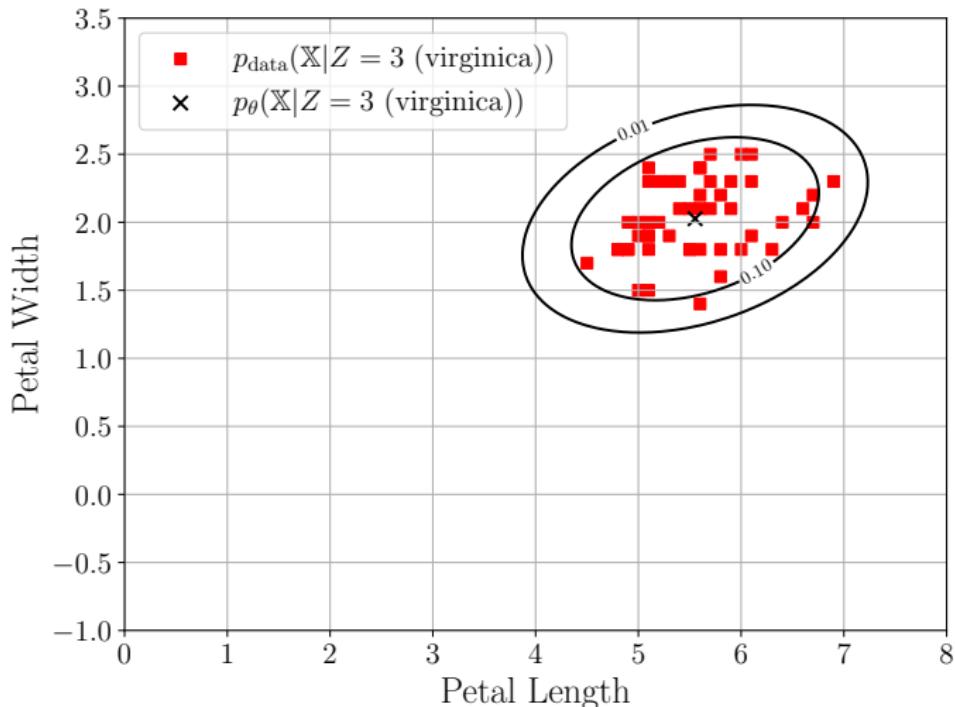


Figure: Modeling conditional distribution $p_{\text{data}}(\mathbb{X}|Z = 3)$ using signle Gaussian

Example - Iris Flower Data set

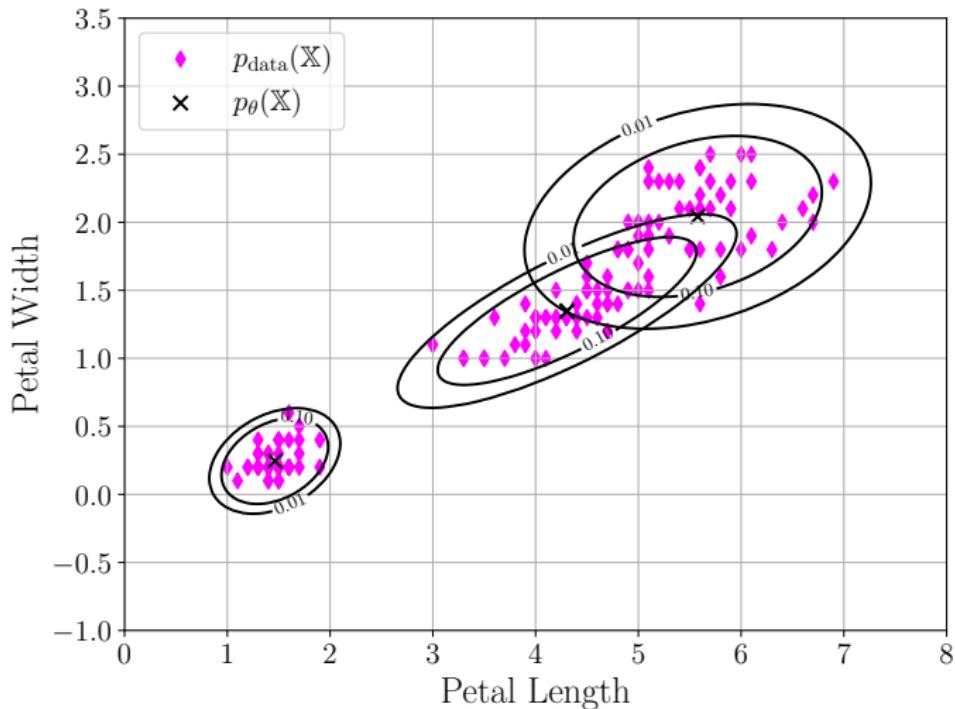


Figure: Modeling unconditional distribution using a Gaussian Mixture Model

Gaussian Mixture Model

Intuition

Intelligently combining simple distributions to create complex ones.

GMM

Assume you have access to Categorical and Gaussian distribution, then you generate each sample as:

- Generate sample z from Categorical distribution $\text{Cat}(\boldsymbol{\pi})$.
- Generate sample \mathbf{x} from conditional distribution $p(\mathbb{X}|z) = \mathcal{N}(\mathbb{X}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$

In the above process, the distribution over \mathbb{X} is:

$$p(\mathbb{X}) = \sum_z p(\mathbb{X}, Z = z) = \sum_z p(\mathbb{X}|z)p(z) = \sum_z \pi_z \mathcal{N}(\mathbb{X}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$$

while the Gaussian distribution is unimodal, the above distribution is multi-modal.

GMM

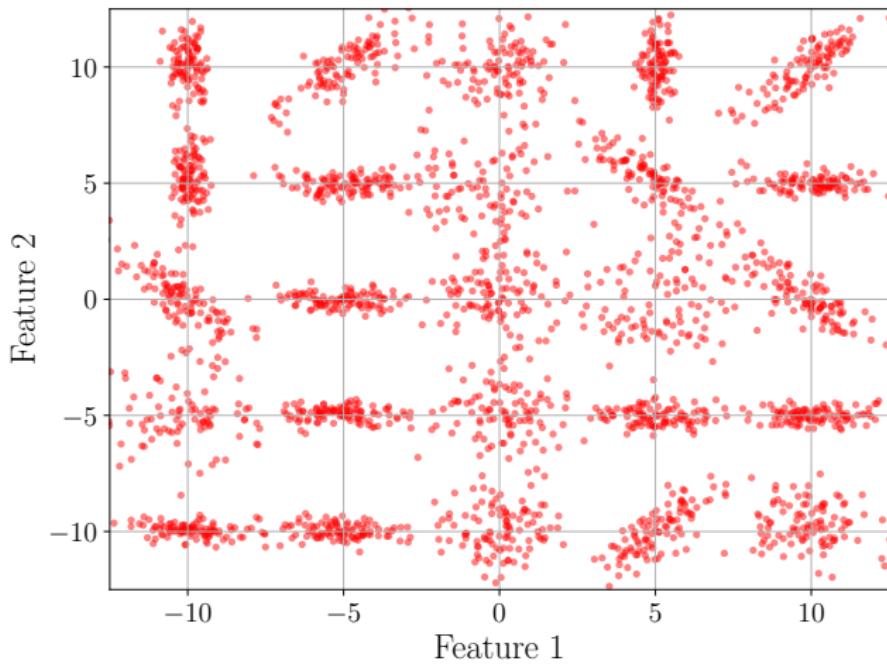


Figure: Synthesis dataset \mathcal{D} samples

GMM

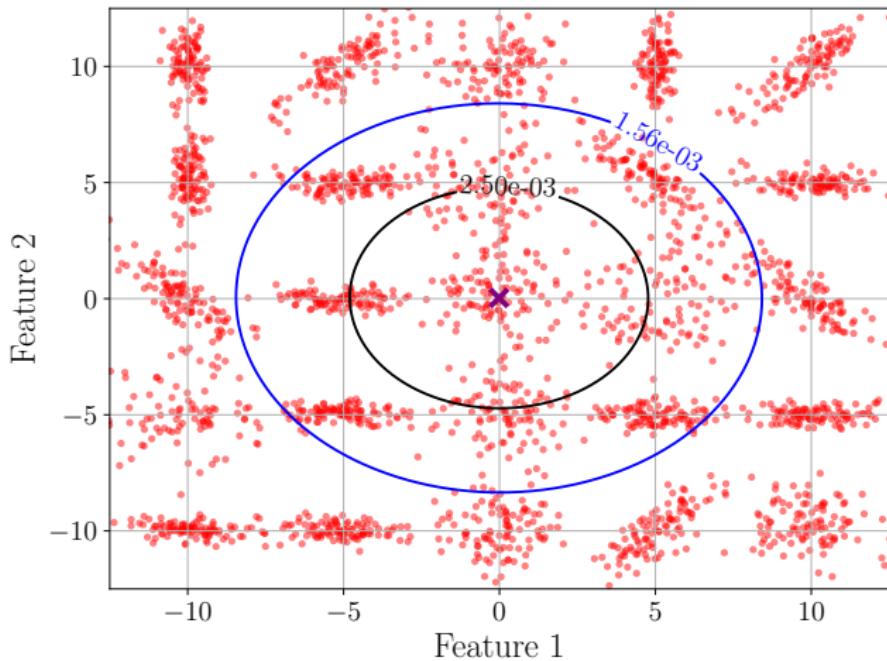


Figure: Gaussian distribution fitted to dataset \mathcal{D} (the purple 'x' symbol shows the mean location)

GMM

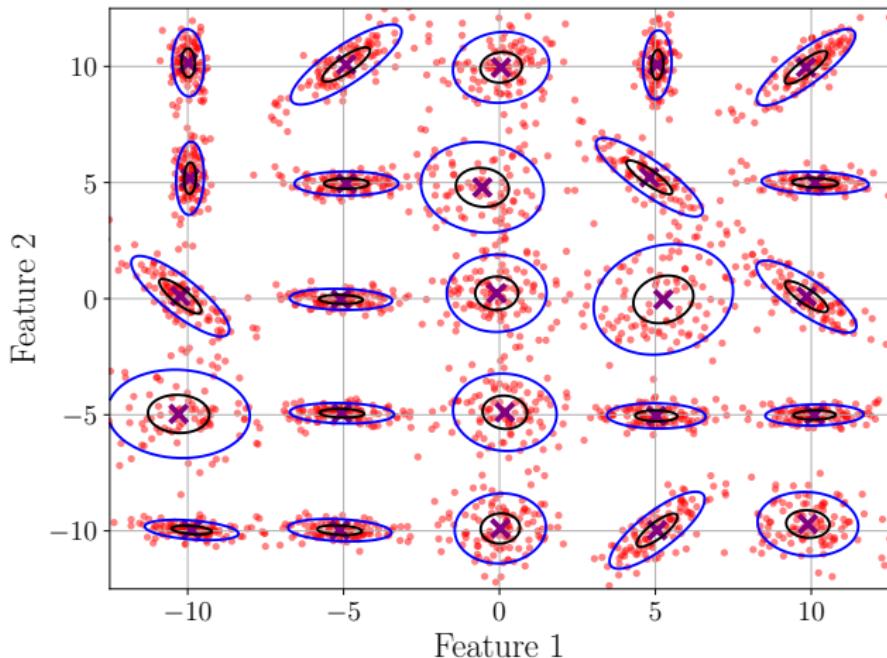


Figure: GMM distribution with four component fitted to dataset \mathcal{D} with 25 components (the purple 'x' symbols shows each component mean location)

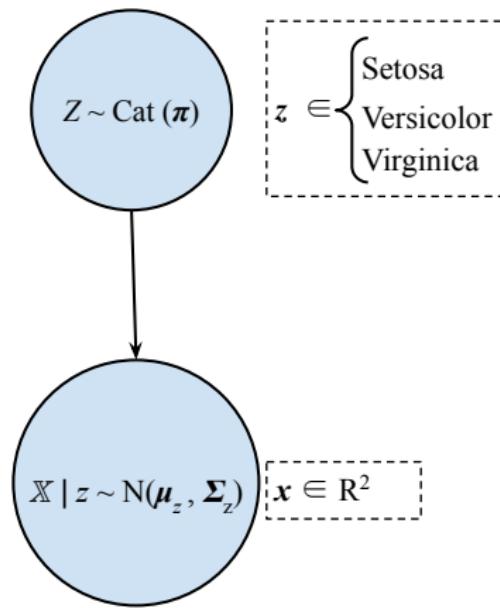
Pros

- A flexible distribution based on simple distribution functions
- Ancestral sampling based on simple distributions
 - ① Sample $z \sim \text{Cat}(\boldsymbol{\pi})$
 - ② Sample $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$
- Latent representation \mathbf{z} (For a given data point \mathbf{x} in Iris dataset, if we can calculate $p(Z|\mathbf{x})$ then we can predict the Iris flower type)

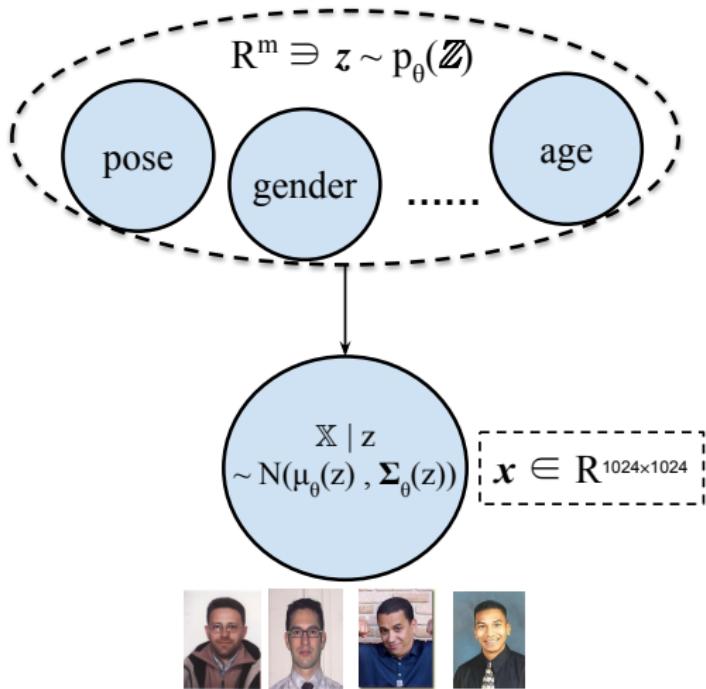
Cons

- No exact density estimation for new samples because of hidden variable \mathbf{z}
- Sensitivity to initialization
- Number of component

Extending GMM



(a) GMM Concept



(b) VAE Concept (images source: [1])

Section 2

Model Specification

Sample Model

Model Specification

Assume we have the following simple distributions:

- Latent variable distribution:

$$p_{\theta}(\mathbb{Z}) = \mathcal{N}(\mathbb{Z} | \boldsymbol{\mu}_z, \sigma_z^2 \mathbf{I}), \quad \mathbb{Z} \in \mathbb{R}^m$$

- Conditional observed variable distribution:

$$p_{\theta}(\mathbb{X} | \mathbf{z}) = \mathcal{N}(\mathbb{X} | \boldsymbol{\mu}_{\theta}(\mathbf{z}), \boldsymbol{\Sigma}_{\theta}(\mathbf{z})), \quad \begin{cases} \mathbb{X} \in \mathbb{R}^n \\ \boldsymbol{\mu}_{\theta}(\mathbf{z}) = \text{NN}_{\boldsymbol{\alpha}}(\mathbf{z}) \\ \boldsymbol{\Sigma}_{\theta}(\mathbf{z}) = \text{diag} \left(\sigma \left(\text{NN}_{\boldsymbol{\beta}}(\mathbf{z}) \right) \right) \end{cases}$$

thus $\boldsymbol{\theta} = \{\boldsymbol{\mu}_z, \sigma_z^2, \boldsymbol{\alpha}, \boldsymbol{\beta}\}$.

Task 1: Generation

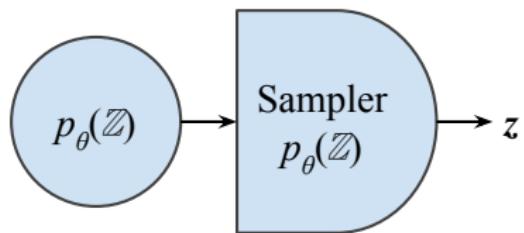


Figure: Sampling latent vector

Task 1: Generation

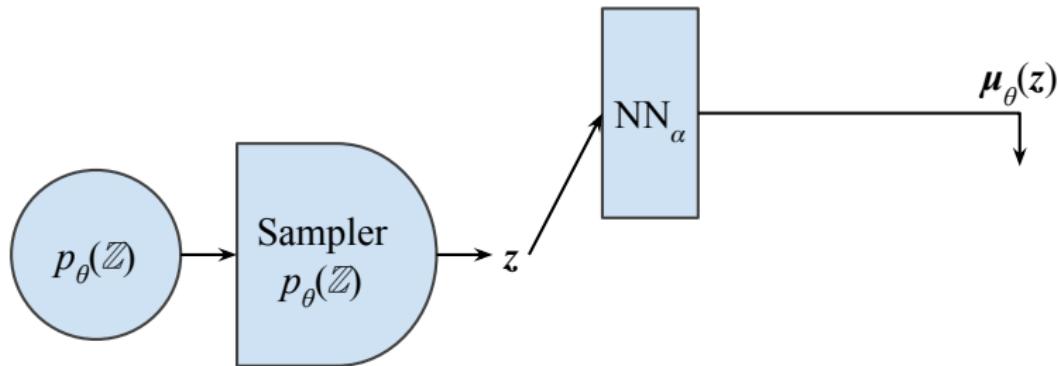


Figure: Calculating conditional distribution mean vector

Task 1: Generation

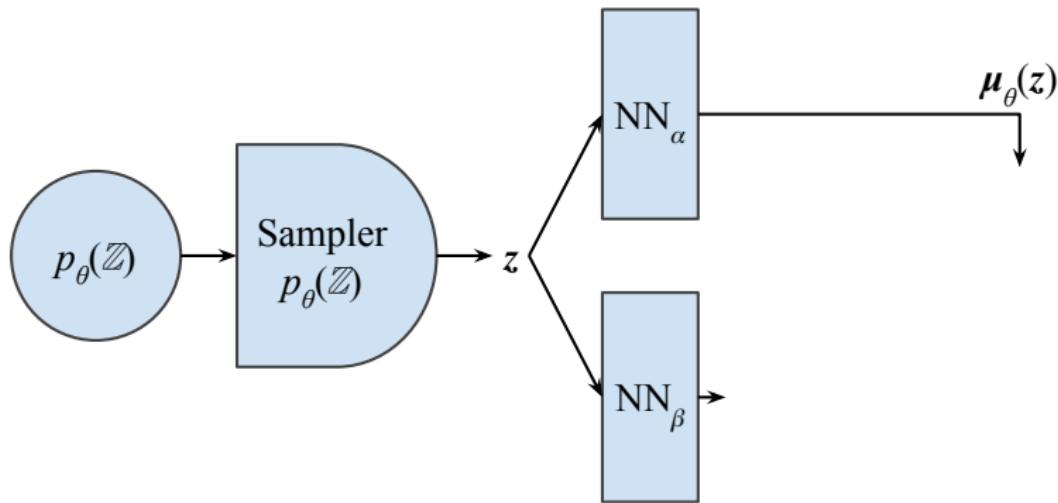


Figure: Calculating conditional distribution covariance matrix

Task 1: Generation

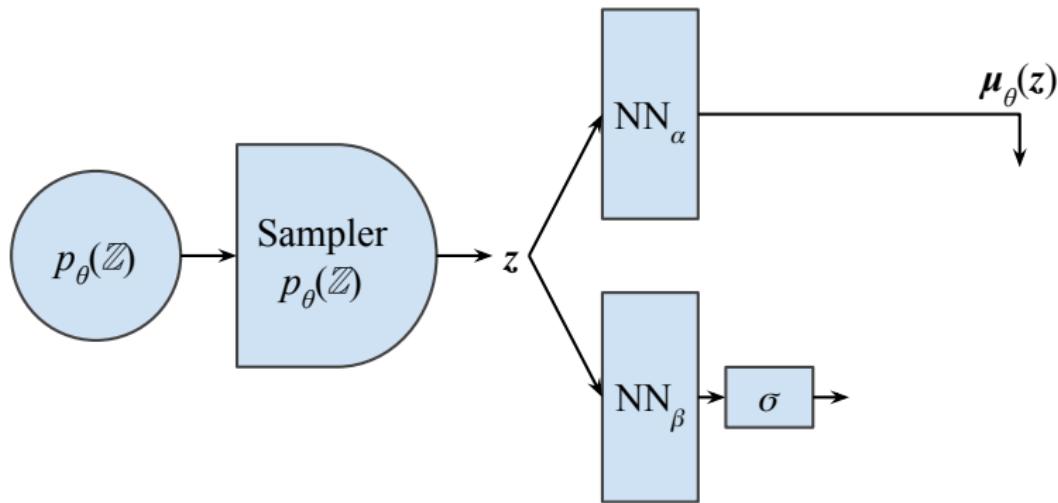


Figure: Calculating conditional distribution covariance matrix

Task 1: Generation

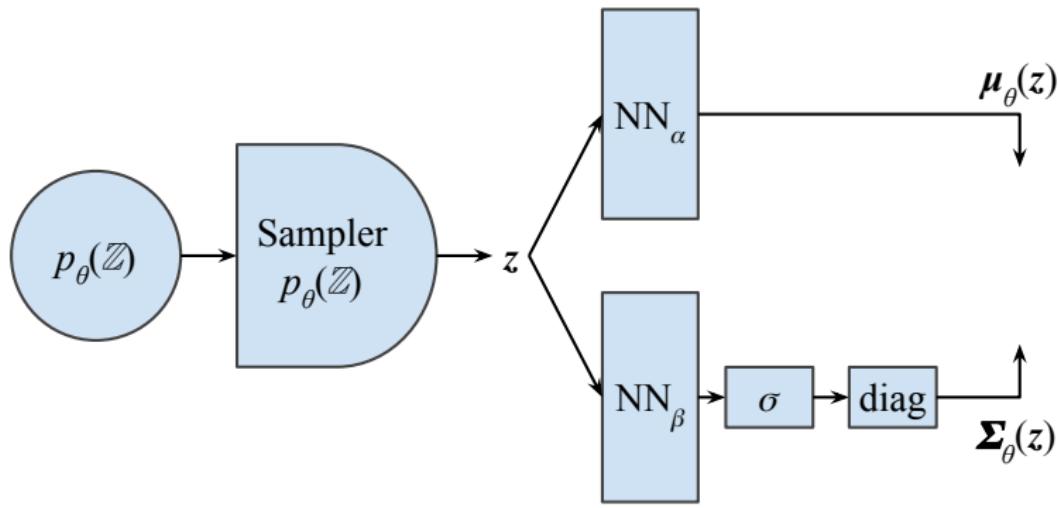


Figure: Calculating conditional distribution covariance matrix

Task 1: Generation

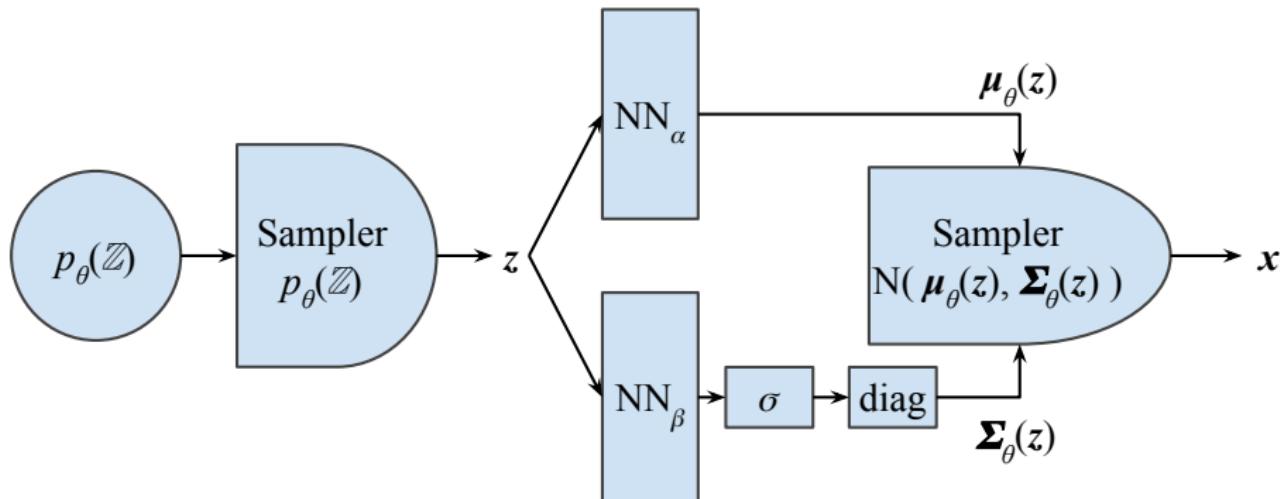


Figure: Sampling conditional distribution

Task 2: High-level Representation

High-level Representation

Remember the example of the Iris dataset. Accessing the latent random variable Z can reveal information about the type of flower. So Z can be considered as *High-level Representation*.

Challenges

For estimating \mathbf{z} corresponding to an observed vector \mathbf{x} , we have the following challenges:

- Estimating the parameters of conditional distribution $\mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{z}), \boldsymbol{\Sigma}_\theta(\mathbf{z}))$ with only one sample \mathbf{x} has a high variance.
- Assume you have access to both $\boldsymbol{\mu}_\theta(\mathbf{z})$ and $\boldsymbol{\Sigma}_\theta(\mathbf{z})$. Then calculating \mathbf{z} is impossible because Deep Neural Networks are non-invertible.

Task 3: Density Estimation

Density Estimation

We have access to both $p_\theta(\mathbb{Z})$ and $p_\theta(\mathbb{X}|\mathbb{Z})$, thus:

$$p_\theta(\mathbb{X}, \mathbb{Z}) = p_\theta(\mathbb{X}|\mathbb{Z})p_\theta(\mathbb{Z})$$

To calculate the density, we need to marginalize out the hidden random vector \mathbb{Z} , thus:

$$\begin{aligned} p_\theta(\mathbb{X}) &= \int_{\mathbb{Z}} p_\theta(\mathbb{X}, \mathbb{Z}) d\mathbb{Z} \\ &= \int_{\mathbb{Z}} p_\theta(\mathbb{X}|\mathbb{Z})p_\theta(\mathbb{Z}) d\mathbb{Z} \\ &= \int_{\mathbb{Z}} \mathcal{N}(\mathbb{X}|\boldsymbol{\mu}_\theta(\mathbb{Z}), \boldsymbol{\Sigma}_\theta(\mathbb{Z})) p_\theta(\mathbb{Z}) d\mathbb{Z} \end{aligned}$$

Challenge 1

Calculating the likelihood and thus training the model seems to be intractable.

Section 3

Casting Likelihood Calculation as Expectation

No Free Lunch: Challenges Start

Calculating the Likelihood for Sample \mathbf{x}

As we see before, we are interested in the following optimization:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} -\frac{1}{|\mathcal{D}|} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

To compute $p_{\boldsymbol{\theta}}(\mathbf{x})$ in a latent variable model we have:

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \int_{\mathbf{z}} p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

The above integral is hard to compute and in each update of the optimization problem, we want:

- Evaluate the integral as a function of $\boldsymbol{\theta}$
- Evaluate the gradient vector of the log-likelihood

Likelihood as Expectation

Likelihood as Expectation

We can write the data likelihood as:

$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbb{Z})} [p_{\theta}(\mathbf{x}|\mathbf{z})]$$

Thus we have an expectation. How can we estimate it?

Monte Carlo Estimate

We can estimate the expectation using the Monte Carlo estimate by k samples from the distribution as:

$$p_{\theta}(\mathbf{x}) \simeq \frac{1}{k} \sum_{i=1}^k p_{\theta}(\mathbf{x}|\mathbf{z}_i), \quad \mathbf{z}_i \sim p_{\theta}(\mathbb{Z}) \text{ for } i = 1, \dots, k$$

The above estimate, while unbiased, exhibits significant variance. Why?

High Variance Estimation - Toy Example

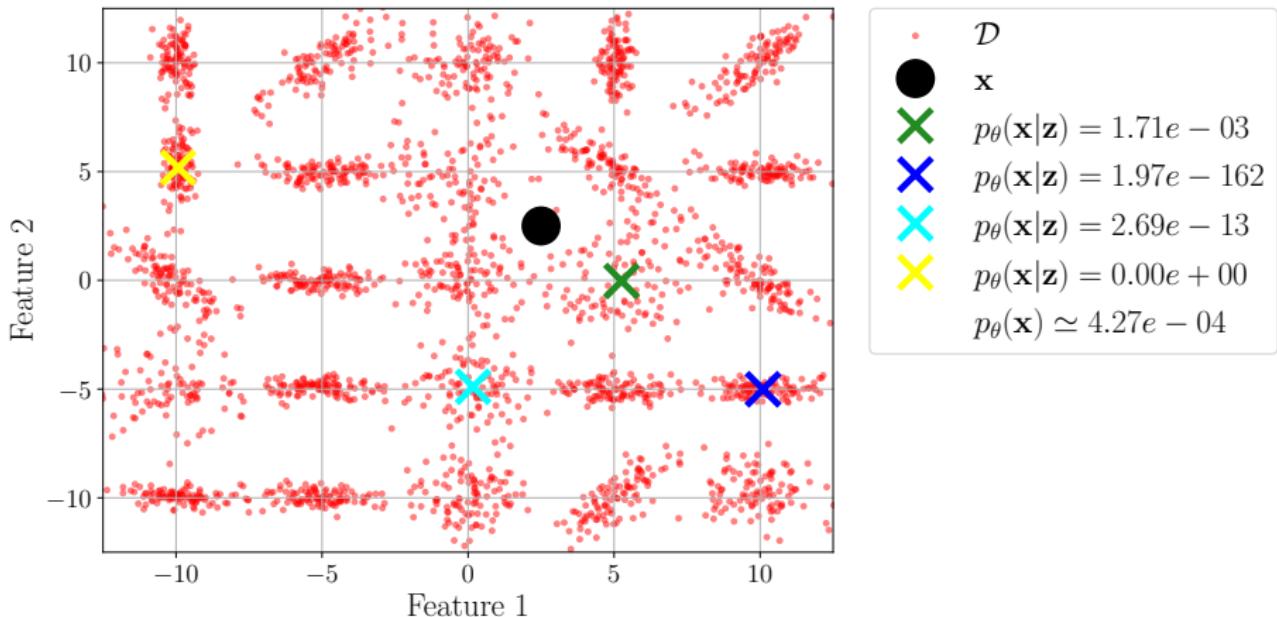


Figure: $[4.27 \times 10^{-4}]$

High Variance Estimation - Toy Example

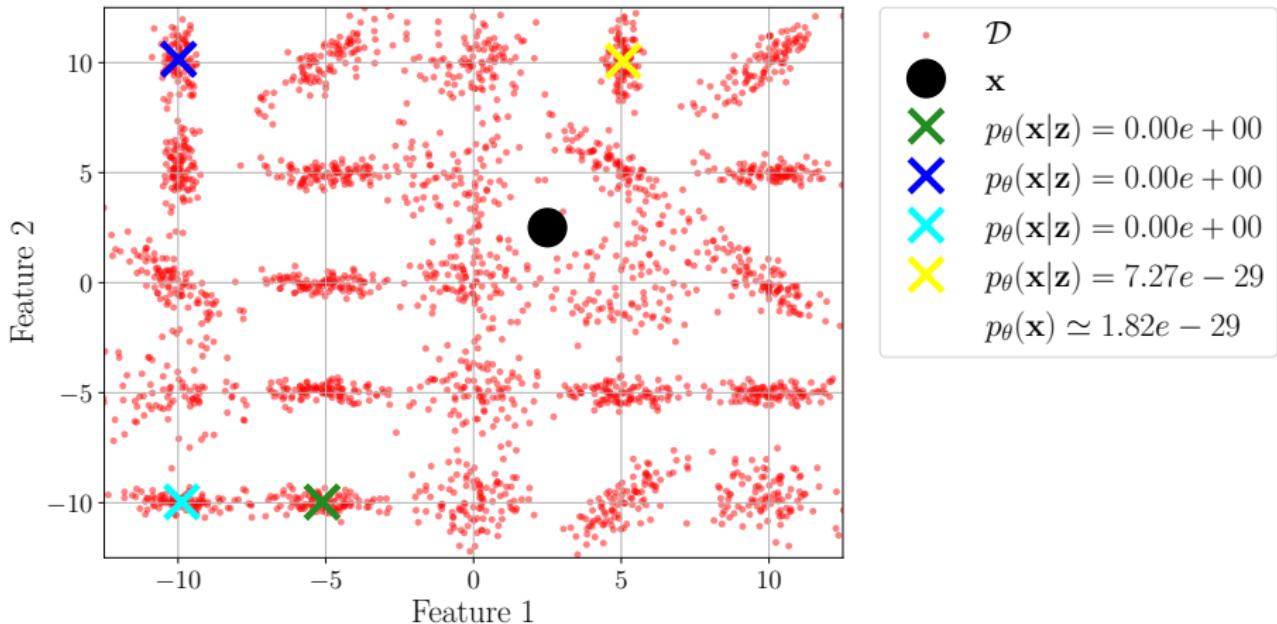


Figure: $[4.27 \times 10^{-4}, 1.82 \times 10^{-29}]$

High Variance Estimation - Toy Example

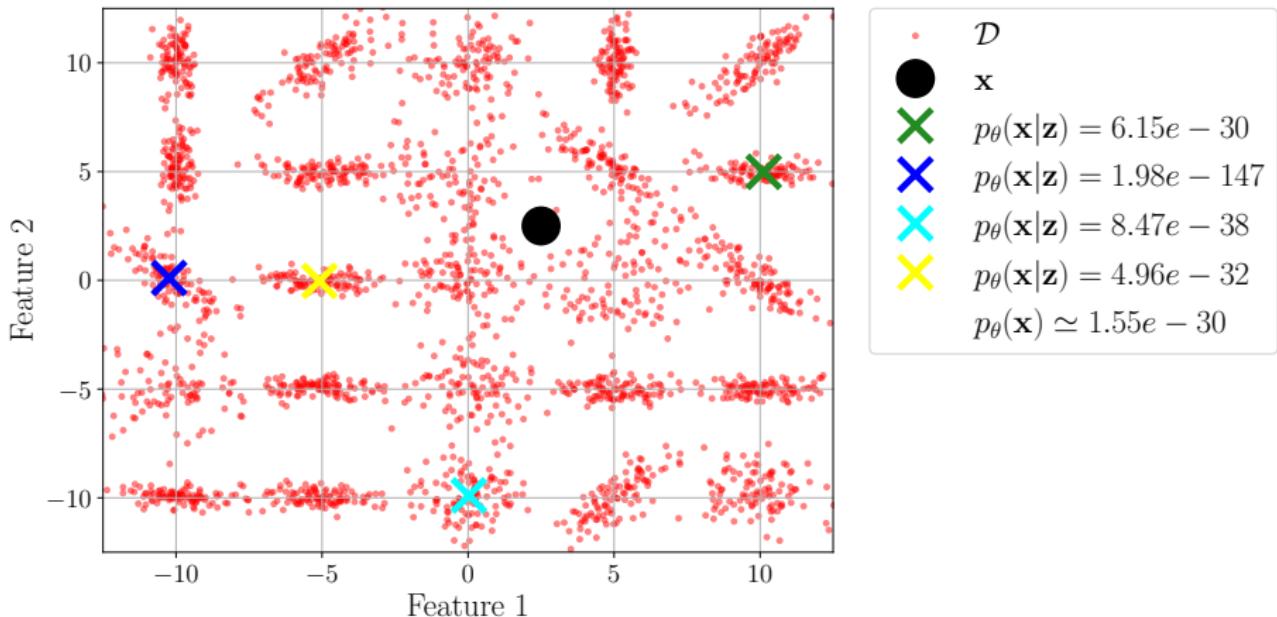


Figure: $[4.27 \times 10^{-4}, 1.82 \times 10^{-29}, 1.55 \times 10^{-30}]$

High Variance Estimation - Toy Example

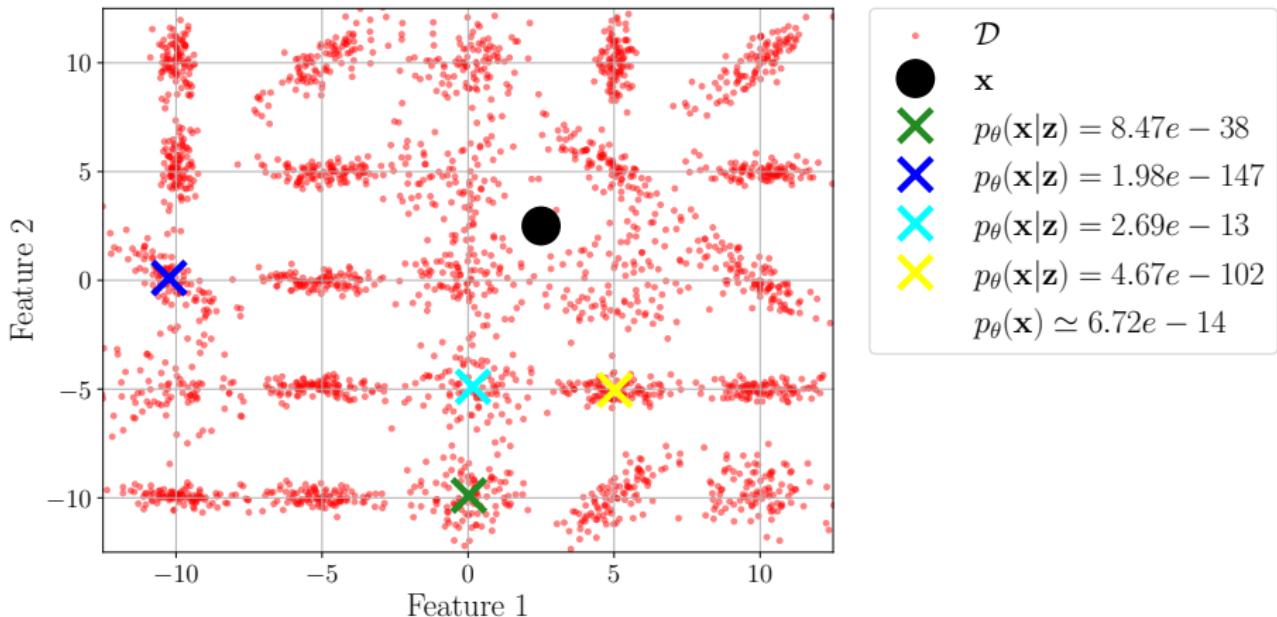
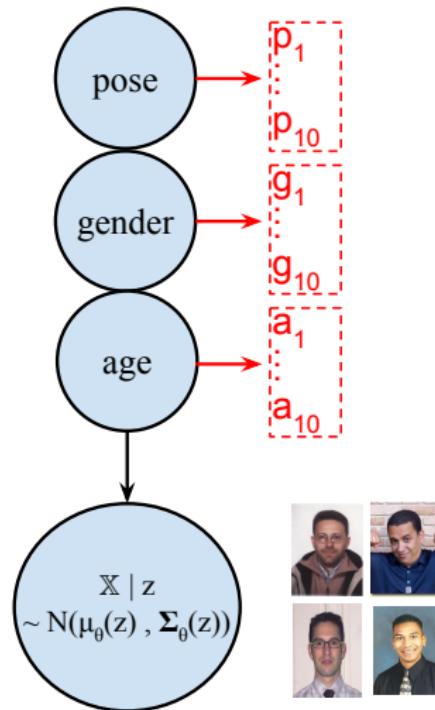


Figure: $[4.27 \times 10^{-4}, 1.82 \times 10^{-29}, 1.55 \times 10^{-30}, 6.72 \times 10^{-14}]$

High Variance Estimate - Real World Example



Scenario

As we can see $|\mathcal{Z}| = 10^3$. Assume we want to compute $p_\theta(\bar{\mathbf{x}})$, the true latent is \bar{z} and:

- $k = 10$
- \mathcal{S} is the set of generated latent vectors

Then:

- $p\{\bar{z} \in \mathcal{S}\} = 0.01 \Rightarrow p_\theta(\bar{\mathbf{x}}) > 0$
- $p\{\bar{z} \notin \mathcal{S}\} = 0.99 \Rightarrow p_\theta(\bar{\mathbf{x}}) \simeq 0$

while the estimation has a high variance and cannot be used in practice.

☞ So, what can we do?

Figure: (images source: [1])

Update on our Challenges

Challenge 1

Calculating the likelihood and thus training the model seems to be intractable.

- ☞ We can cast the likelihood as expectation and estimate the expectation using a Monte Carlo Estimation

Challenge 2

Estimating the likelihood using prior distribution, while unbiased, represents high variance.

Section 4

Importance Sampling

Importance Sampling

Importance Sampling

To calculate the likelihood, we can use importance sampling:

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} q(\mathbf{z}) \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \end{aligned}$$

where $q(\mathbf{z})$ can be any arbitrary distribution function over \mathbb{Z} . We can estimate the distribution using the Monte Carlo estimate:

$$p_{\theta}(\mathbf{x}) \simeq \frac{1}{k} \sum_{i=1}^k \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_i)}{q(\mathbf{z}_i)}, \quad \mathbf{z}_i \sim q(\mathbb{Z}) \text{ for } i = 1, \dots, k$$

Challenge

What is a good option for $q(\mathbb{Z})$. Maybe we need optimization over q too!

Importance Sampling in Toy Example

Estimating Latent Posterior

Assume the toy example of GMM, we can calculate $p_\theta(\mathbb{Z}|\mathbf{x})$ as:

$$p_\theta(\mathbb{Z} = i|\mathbf{x}) = \frac{p_\theta(\mathbf{x}, \mathbb{Z} = i)}{p_\theta(\mathbf{x})} \quad \# \text{Bayes Rule}$$

$$= \frac{p_\theta(\mathbf{x}, \mathbb{Z} = i)}{\sum_j p_\theta(\mathbf{x}, \mathbb{Z} = j)} \quad \# \text{Marginalization}$$

$$= \frac{p_\theta(\mathbf{x}|\mathbb{Z} = i)p_\theta(\mathbb{Z} = i)}{\sum_j p_\theta(\mathbf{x}|\mathbb{Z} = j)p_\theta(\mathbb{Z} = j)} \quad \# \text{Chain Rule}$$

$$= \frac{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\pi_i}{\sum_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\pi_j} \quad \text{Model Parameters}$$

Let's use importance sampling technique to calculate $p_\theta(\mathbf{x})$ with $q(Z) = p_\theta(Z|\mathbf{x})$

GMM - Posterior

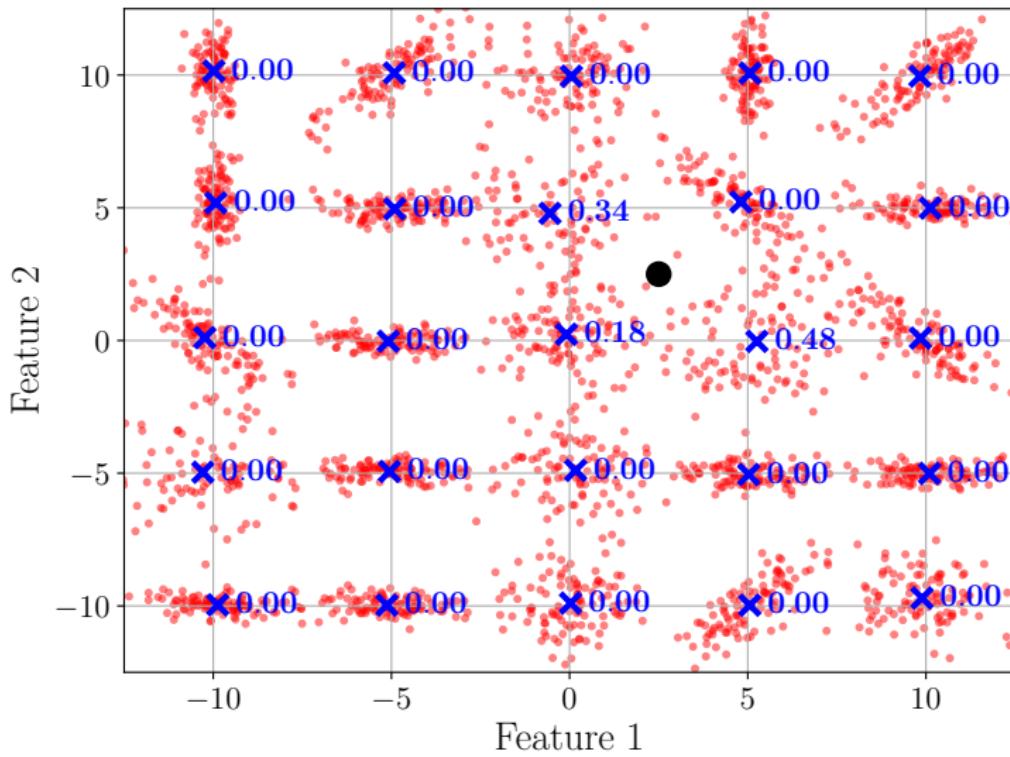


Figure: Posterior distribution $p(Z|x)$ (black circle represents x)

High Variance Estimation - Toy Example

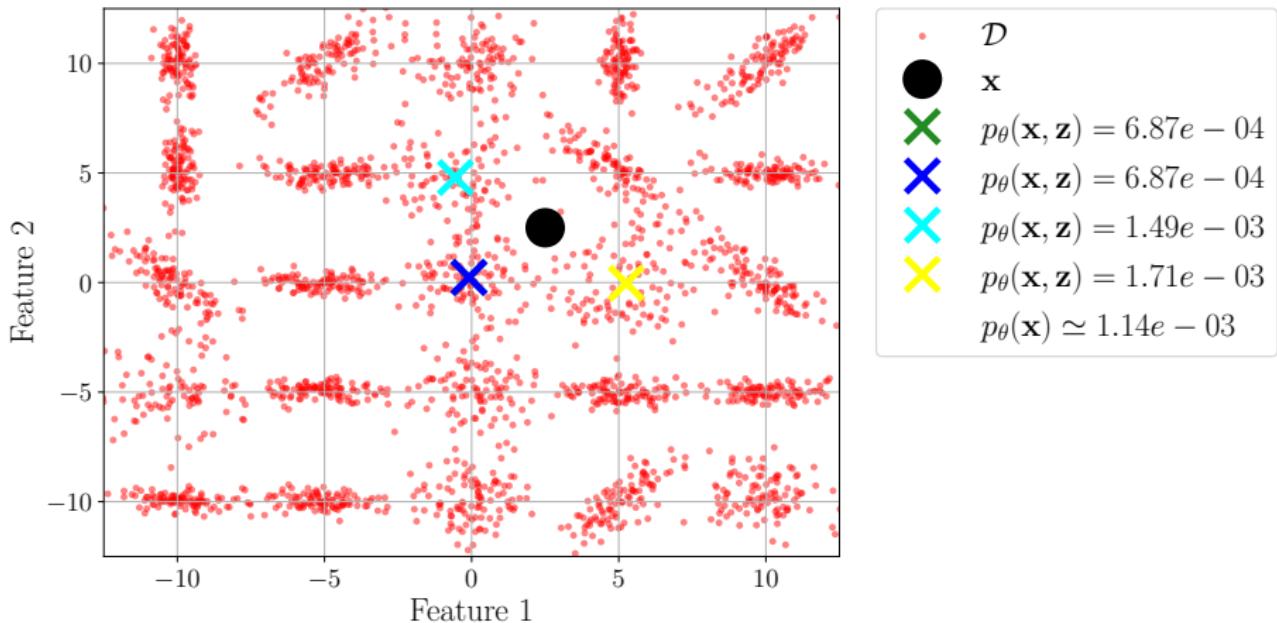


Figure: $[1.14 \times 10^{-3}]$

High Variance Estimation - Toy Example

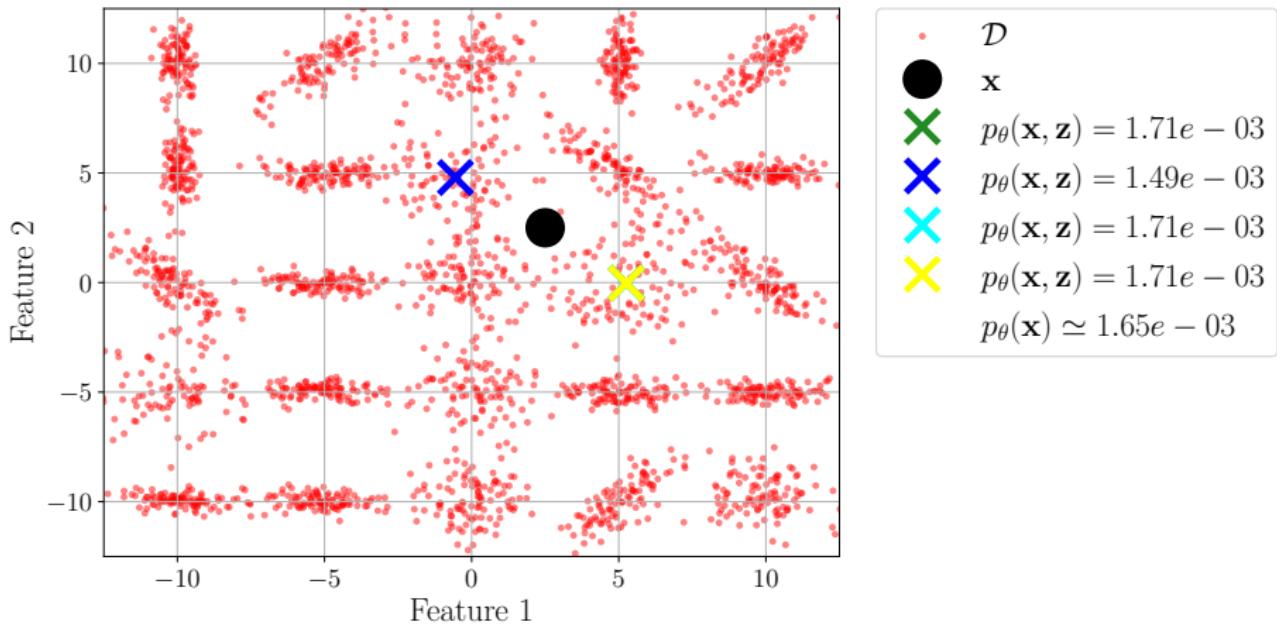


Figure: $[1.14 \times 10^{-3}, 1.65 \times 10^{-3}]$

High Variance Estimation - Toy Example

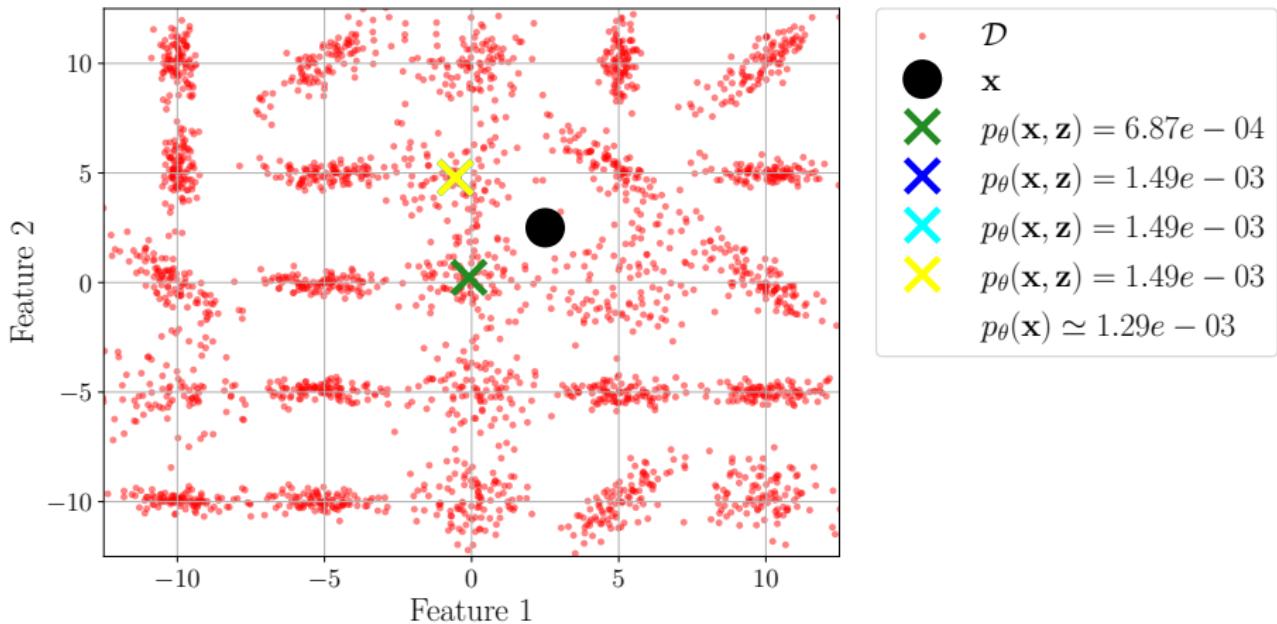


Figure: $[1.14 \times 10^{-3}, 1.65 \times 10^{-3}, 1.29 \times 10^{-3}]$

High Variance Estimation - Toy Example

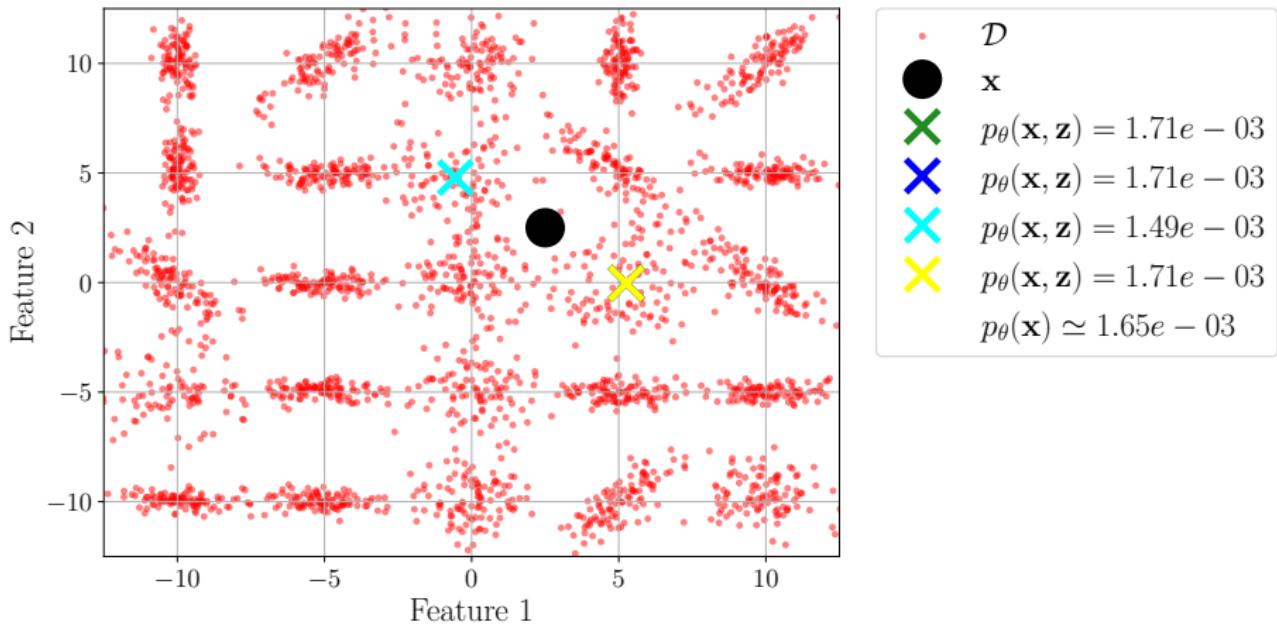


Figure: $[1.14 \times 10^{-3}, 1.65 \times 10^{-3}, 1.29 \times 10^{-3}, 1.65 \times 10^{-3}]$

Lowering Estimation Variance

Real-World Case

Although for the case of GMM latent variable posterior is tractable, in real-world application we have (assuming the latent prior is standard Gaussian):

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

and we know that calculating $p_{\theta}(\mathbf{x})$ is challenging (indeed we are looking for it!). So you can't use $p_{\theta}(\mathbf{Z}|\mathbf{x})$ as $q(\mathbf{Z})$ in the importance sampling technique.

Update on our Challenges

Challenge 2

Estimating the likelihood using prior distribution, while unbiased, represents high variance.

- ☞ Importance sampling may help. We observe $q(\mathbb{Z}) = p_\theta(\mathbb{Z}|\mathbf{x})$ is a good option at least in the case of GMM using the importance sampling technique.

Challenge 3

We have no formal derivation for the suitability of $p_\theta(\mathbb{Z}|\mathbf{x})$ to be used for importance sampling in the general real-world case (other than GMM). So we don't know how to select $q(\mathbb{Z})$

Section 5

Evidence Lower BOund

Working on Log Likelihood

Log Likelihood

In the previous section, we have seen that:

$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$$

But what we need in practice for training:

$$\log p_{\theta}(\mathbf{x}) = \log \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$$

For two important reasons (among others), we need to change the order of expectation and log on RHS.

Working on Log Likelihood

Reason 1: Disentangling Parameters

If we can change the order then on RHS we have:

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{Z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{Z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{Z})} [\log q(\mathbf{z})]$$

- The optimization over model parameters θ is done using the first term.
- The prospective optimization over q functions should still be done using both terms.

Working on Log Likelihood

Reason 2: Two Separate Learning Signals

For the log-likelihood, we have:

$$\log p_{\theta}(\mathbf{x}) = \log \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \simeq \log \left[\frac{1}{k} \sum_k \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z}_i)}{q(\mathbf{z}_i)} \right] \right], \quad \mathbf{z}_i \sim q(\mathbb{Z})$$

If we can change the order then on RHS we have:

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] &= \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log q(\mathbf{z})] \\ &\simeq \frac{1}{k} \sum_k \log p_{\theta}(\mathbf{x}, \mathbf{z}_i) - \frac{1}{k} \sum_k \log q(\mathbf{z}_i), \quad \mathbf{z}_i \sim q(\mathbb{Z}) \end{aligned}$$

Because we are estimating the expectations with MCE, then in the second case, the learning signal from each of the terms is separated leading to a better approximation.

Working on Log Likelihood

Challenge

The challenge is to relate the following terms:

$$\text{We have } \Rightarrow \log \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \begin{array}{c} \leq \\ = \\ \geq \end{array} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \Leftarrow \text{We want}$$

Jensen's Inequality

Assume random vector \mathbb{X} and concave function $\psi(\cdot)$. Then:

$$\psi(\mathbb{E}[\mathbb{X}]) \geq \mathbb{E}[\psi(\mathbb{X})]$$

Evidence Lower BOund

Evidence Lower BOund (ELBO)

The Log function is concave. So using Jensen's inequality, we have:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \\ &\geq \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right) \right]\end{aligned}$$

So by changing the order of log function and expectation, we get a lower bound for the log-likelihood known as *Evidence Lower BOund* or abbreviateley ELBO.

Notation

ELBO is defined for one sample \mathbf{x} , while the model parameters $\boldsymbol{\theta}$ and q distributions can be changed. Thus we use the following notation:

$$\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right) \right]$$

ELBO

ELBO

We can work on ELBO as:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q) \\&= \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log q(\mathbf{z})] \\&= \mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + \mathbb{H}(q)\end{aligned}$$

where $\mathbb{H}(q)$, called *distribution entropy*, is a well-known quantity in information theory.

ELBO Instead of Data Log-likelihood

Replacing Data Log-likelihood with ELBO

To train a model, we need to solve the following optimization:

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

We can replace the above objective with its lower bound, $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$. Hopefully maximizing the lower bound leads to maximizing the data log-likelihood. So:

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$$

Tightness Matters

As we are replacing $\log p_{\boldsymbol{\theta}}(\mathbf{x})$ with $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$, the tightness of the lower bound becomes important.

ELBO Tightness

Reframing ELBO

Based on KL divergence, we have:

$$\begin{aligned}\text{KL}(q(\mathbb{Z}) \| p_{\theta}(\mathbb{Z} | \mathbf{x})) &= \int_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \\ &= \int_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})}} d\mathbf{z} \\ &= \int_{\mathbf{z}} q(\mathbf{z}) [\log p_{\theta}(\mathbf{x}) - \log p_{\theta}(\mathbf{x}, \mathbf{z}) + \log q(z)] d\mathbf{z} \\ &= \log p_{\theta}(\mathbf{x}) \int_{\mathbf{z}} q(\mathbf{z}) d\mathbf{z} - \int_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} + \int_{\mathbf{z}} q(\mathbf{z}) \log q(z) d\mathbf{z} \\ &= \underbrace{\log p_{\theta}(\mathbf{x})}_{\text{Loh-likelihood}} - \underbrace{\left[\mathbb{E}_{\mathbf{z} \sim q(\mathbb{Z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + \mathbb{H}(q) \right]}_{\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)}\end{aligned}$$

ELBO Tightness

Tightest ELBO

So we have the following interesting equality:

$$\text{KL}(q(\mathbb{Z})\|p_{\theta}(\mathbb{Z}|\mathbf{x})) = \log p_{\theta}(\mathbf{x}) - \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$$

Thus the gap is $\text{KL}(q(\mathbb{Z})\|p_{\theta}(\mathbb{Z}|\mathbf{x}))$.

- As the distance between $q(\mathbb{Z})$ and $p_{\theta}(\mathbb{Z}|\mathbf{x})$ increases, the ELBO becomes a less tight bound.
- When you select $q(\mathbb{Z}) = p_{\theta}(\mathbb{Z}|\mathbf{x})$, ELBO touches model log Likelihood (as we see before).

So for $q(\mathbb{Z}) = p_{\theta}(\mathbb{Z}|\mathbf{x})$, we have:

$$\text{KL}(q(\mathbb{Z})\|p_{\theta}(\mathbb{Z}|\mathbf{x})) = 0 \Rightarrow \log p_{\theta}(\mathbf{x}) = \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$$

Consequently, we have found the $q(\mathbb{Z})$ distribution resulting in the tightest ELBO.

Update on our Challenges

Challenge 3

We have no formal derivation for the suitability of $p_\theta(\mathbb{Z}|\mathbf{x})$ to be used for importance sampling in the general real-world case (other than GMM). So we don't know how to select $q(\mathbb{Z})$

- ☞ $p_\theta(\mathbb{Z}|\mathbf{x})$ is the best option which leads ELBO touch data log-likelihood.

Challenge 4

Calculating $p_\theta(\mathbb{Z}|\mathbf{x})$ needs the data likelihood which is intractable as we observe in Slide 28. On the other hand, we need to decrease $\text{KL}(q(\mathbb{Z})\|p_\theta(\mathbb{Z}|\mathbf{x}))$ to tighten the ELBO.

Section 6

ELBO Tightening

ELBO Tightening

Parameterizing q

To have control over distribution $q(\mathbb{Z})$, we should parameterize this distribution. An example can be:

- $q_\lambda(\mathbb{Z}) = \mathcal{N}(\mathbb{Z} | \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ and $\boldsymbol{\lambda} = \{\boldsymbol{\mu}, \sigma^2\}$

Note that we are currently focused on the likelihood of one sample point \mathbf{x} and thus $\boldsymbol{\lambda}$ is the parameter of the distribution corresponding to \mathbf{x} .

Notation Update

By parameterizing q , instead of $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, q)$ we use:

$$\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbb{E}_{\mathbf{z} \sim q_\lambda(\mathbb{Z})} \left[\log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\lambda(\mathbf{z})} \right) \right]$$

And again pay attention that everything is just about one data point \mathbf{x} .

ELBO Tightening

Optimization Problem

For ELBO tightening in a model with parameter θ , we need to solve the following optimization problem:

$$\lambda^*(\theta) = \operatorname{argmin}_{\lambda} \text{KL}(q_\lambda(\mathbb{Z}) \parallel p_\theta(\mathbb{Z}|\mathbf{x}))$$

But pay attention to the following note:

- If you can find $\lambda^*(\theta)$, then it is a function of θ . Thus by changing the model parameter θ , λ^* is not optimum anymore and you should update it.

Challenge 4

Solving the above optimization seems intractable as calculating $p_\theta(\mathbb{Z}|\mathbf{x})$ is intractable in general.

ELBO Tightening

Intuition

Pay attention to the following equality:

$$\log p_{\theta}(\mathbf{x}) = \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) + \text{KL}(q_{\lambda}(\mathbb{Z}) \| p_{\theta}(\mathbb{Z} | \mathbf{x}))$$

Two important points:

- LHS is independent of $\boldsymbol{\lambda}$.
 - Both terms on the RHS are functions of $\boldsymbol{\lambda}$.
- ☞ How can we decrease the KL divergence between $q_{\lambda}(\mathbb{Z})$ and $p_{\theta}(\mathbb{Z} | \mathbf{x})$?

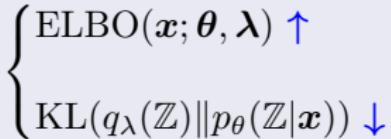
ELBO Tightening

Tightening Solution

We know:

$$\log p_\theta(\mathbf{x}) = \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) + \text{KL}(q_\lambda(\mathbb{Z}) \| p_\theta(\mathbb{Z}|\mathbf{x}))$$

Thus:

Because $\log p_\theta(\mathbf{x})$ is independent of $\boldsymbol{\lambda}$ \Rightarrow 

In other words for a fixed model $\boldsymbol{\theta}$:

$$\operatorname{argmax}_{\boldsymbol{\lambda}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) \equiv \operatorname{argmin}_{\boldsymbol{\lambda}} \text{KL}(q_\lambda(\mathbb{Z}) \| p_\theta(\mathbb{Z}|\mathbf{x}))$$

Thus we can tighten the ELBO by maximizing it with respect to $\boldsymbol{\lambda}$.

ELBO Tightening

ELBO Tightening

Assume the case of simple GMM with two components as:

$$p_{\theta}(\mathbf{x}) = \sum_{z \in \{0,1\}} p_{\theta}(Z = z) p_{\theta}(\mathbf{x}|Z = z), \quad \begin{cases} p_{\theta}(Z = z) = \pi_z \\ p_{\theta}(\mathbf{x}|Z = z) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) \end{cases}$$

So the model parameters are $\boldsymbol{\theta} = \{\pi_0, \pi_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0, \mu_1, \boldsymbol{\Sigma}_1\}$ and given.

In this example, we want to show the following problems are equivalent:

$$\operatorname{argmax}_{\lambda} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \lambda) \equiv \operatorname{argmin}_{\lambda} \text{KL}(q_{\lambda}(Z) \| p_{\theta}(Z|\mathbf{x}))$$

We will follow by solving each optimization problem separately.

ELBO Tightening

RHS Optimization

We are interested in the following optimization:

$$\operatorname{argmin}_{\lambda} \text{KL}(q_{\lambda}(Z) \| p_{\theta}(Z|\boldsymbol{x}))$$

- Because Z is binary thus:
 - $q_{\lambda}(Z) = \text{Ber}(Z|\lambda)$
 - $p_{\theta}(Z|\boldsymbol{x}) = \text{Ber}(Z|\beta(\boldsymbol{\theta}))$ (we write $\beta(\boldsymbol{\theta})$ to emphasize that $p_{\theta}(Z|\boldsymbol{x})$ has no new parameter and its parameter is a function of $\boldsymbol{\theta}$)
- The optimum λ based on RHS is:

$$\lambda^* = \beta(\boldsymbol{\theta})$$

So we need to find $\beta(\boldsymbol{\theta}) = p_{\theta}(Z=1|\boldsymbol{x})$

ELBO Tightening

RHS Optimization

$$\beta(\boldsymbol{\theta}) = p_{\theta}(Z = 1 | \mathbf{x})$$

$$= \frac{p_{\theta}(\mathbf{x}, Z = 1)}{p_{\theta}(\mathbf{x})} \quad \text{\#Bayes Rule}$$

$$= \frac{p_{\theta}(\mathbf{x}, Z = 1)}{\sum_{z \in \{0,1\}} p_{\theta}(\mathbf{x}, Z = z)} \quad \text{\#Marginalization}$$

$$= \frac{p_{\theta}(\mathbf{x}|Z = 1)p_{\theta}(Z = 1)}{\sum_{z \in \{0,1\}} p_{\theta}(\mathbf{x}|Z = z)p_{\theta}(Z = z)} \quad \text{\#Conditioning}$$

$$= \frac{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)\pi_1}{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)\pi_0 + \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)\pi_1} = \lambda^*$$

ELBO Tightening

RHS Optimization

Assume the model parameters are:

$$\boldsymbol{\theta} = \left\{ \begin{array}{l} \pi_1 = 0.33 \quad \mu_1 = \begin{bmatrix} 4.04 \\ 3.83 \end{bmatrix} \quad \boldsymbol{\Sigma}_1 = \begin{bmatrix} 1.79 & -0.10 \\ -0.10 & 2.00 \end{bmatrix} \\ \pi_0 = 0.67 \quad \mu_0 = \begin{bmatrix} 1.10 \\ 0.86 \end{bmatrix} \quad \boldsymbol{\Sigma}_0 = \begin{bmatrix} 1.20 & -0.97 \\ -0.97 & 1.15 \end{bmatrix} \end{array} \right\} \beta(\boldsymbol{\theta}) = 0.76$$

then using the equation in Slide 48 we have:

$$\lambda^* = \beta(\boldsymbol{\theta}) = 0.76$$

ELBO Tightening

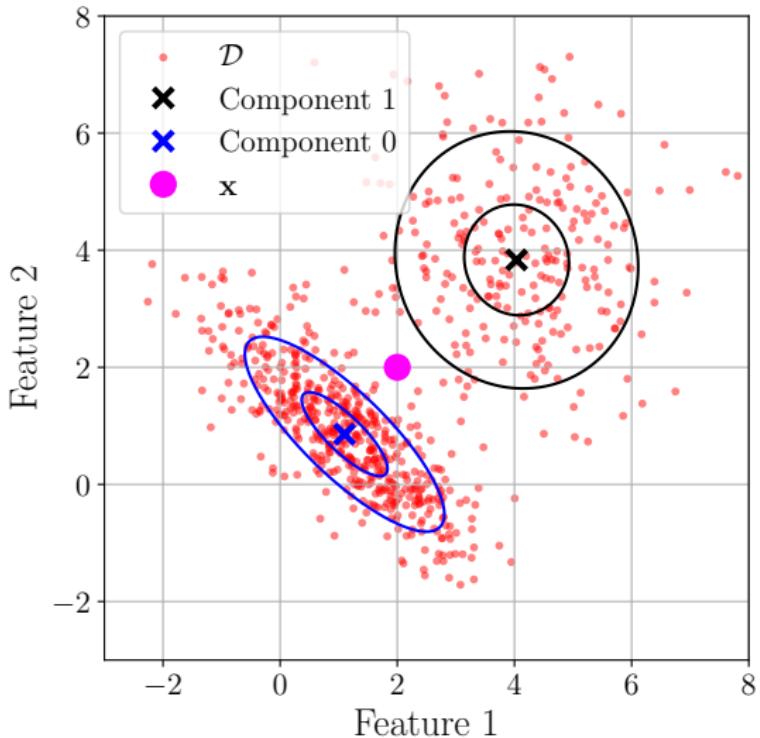


Figure: Dataset, GMM components and query point x

ELBO Tightening

LHS Optimization

Assume distribution $q_\lambda(Z) = \text{Ber}(Z|\lambda)$, then:

$$\begin{aligned}\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \lambda) &= \mathbb{E}_{z \sim q_\lambda(Z)} \log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, z)}{q_\lambda(z)} \right] \\ &= q_\lambda(1) \log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, 1)}{q_\lambda(1)} \right] + q_\lambda(0) \log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, 0)}{q_\lambda(0)} \right] \\ &= \lambda \log \left[\frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \pi_1}{\lambda} \right] + (1 - \lambda) \log \left[\frac{\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \pi_0}{1 - \lambda} \right]\end{aligned}$$

Assuming fixed model parameter $\boldsymbol{\theta}$, we can plot $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \lambda)$ as a function of λ .

ELBO Tightening

LHS Optimization

We know $p_\theta(Z|\mathbf{x}) = \text{Ber}(Z|0.76)$ and $q_\lambda(Z) = \text{Ber}(Z|\lambda)$, so:

$$\begin{aligned}\text{KL}(q_\lambda(Z)\|p_\theta(Z|\mathbf{x})) &= \mathbb{E}_{z \sim q_\lambda(Z)} \log \frac{q_\lambda(z)}{p_\theta(z|\mathbf{x})} \\ &= q_\lambda(1) \log \frac{q_\lambda(1)}{p_\theta(1|\mathbf{x})} + q_\lambda(0) \log \frac{q_\lambda(0)}{p_\theta(0|\mathbf{x})} \\ &= \lambda \log \frac{\lambda}{0.76} + (1 - \lambda) \log \frac{1 - \lambda}{(1 - 0.76)}\end{aligned}$$

In Slide 48, we have also calculate $p_\theta(\mathbf{x})$.

ELBO Tightening

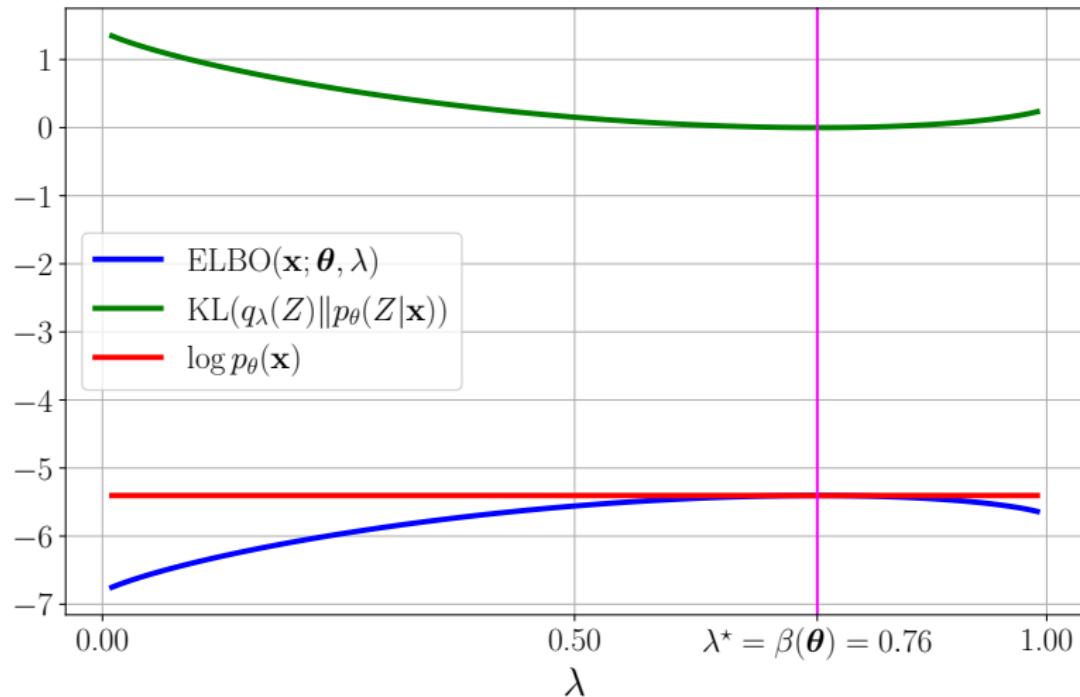


Figure: $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \lambda)$, $\text{KL}(q_\lambda(Z) \| p_\theta(Z|\mathbf{x}))$ and $\log p_\theta(\mathbf{x})$ as a function of λ

Altogether

Challenge 1 

$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{\theta}(\mathbf{z}), \boldsymbol{\Sigma}_{\theta}(\mathbf{z})) p_{\theta}(\mathbf{z}) d\mathbf{z}$$

Altogether

$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z})} [p_{\theta}(\mathbf{x}|\mathbf{z})]$$

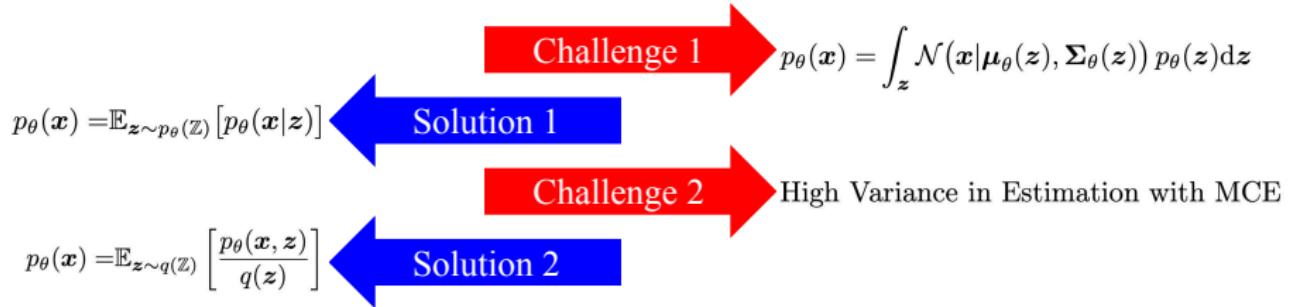
Challenge 1 → $p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\theta}(\mathbf{z}), \boldsymbol{\Sigma}_{\theta}(\mathbf{z})) p_{\theta}(\mathbf{z}) d\mathbf{z}$

← Solution 1

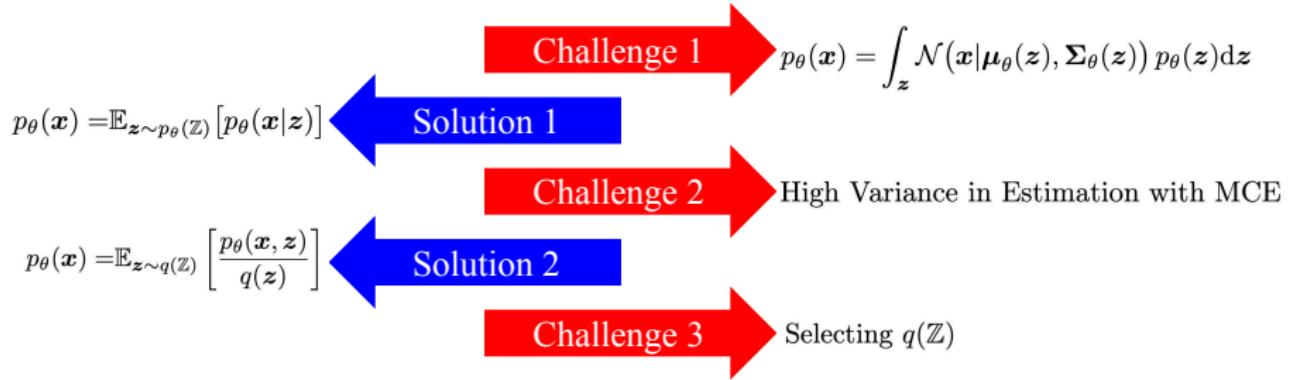
Altogether



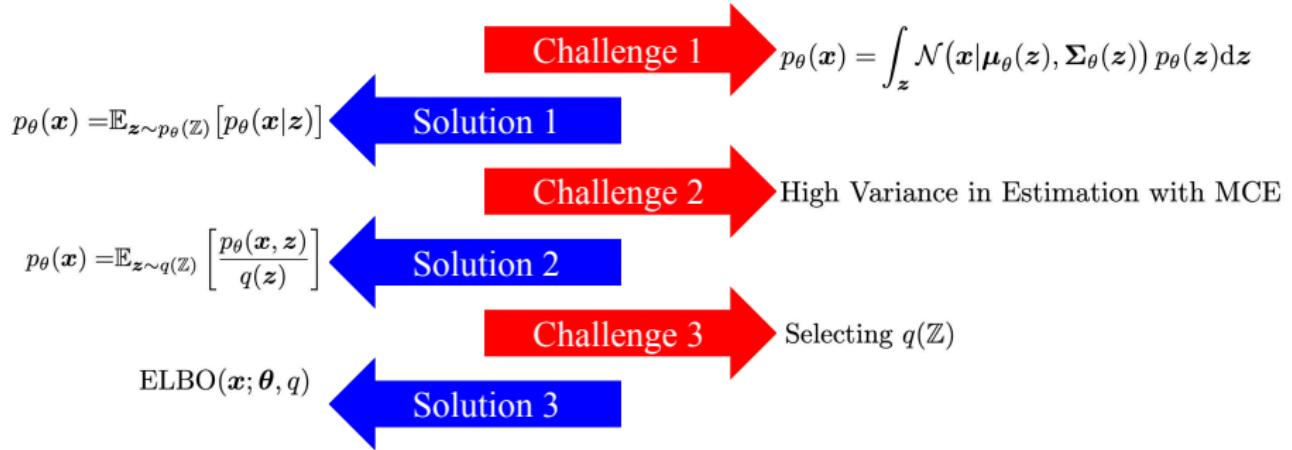
Altogether



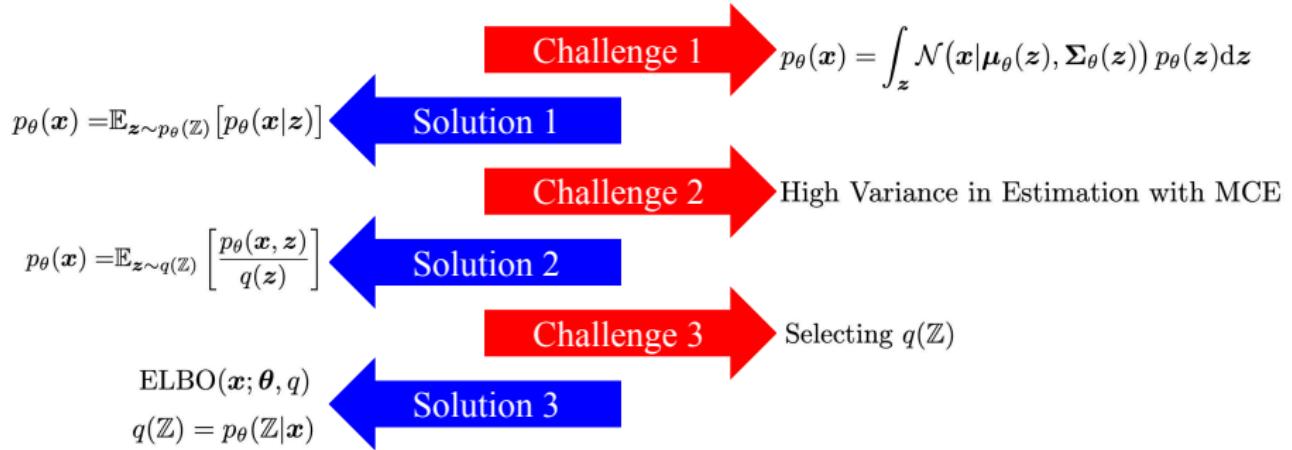
Altogether



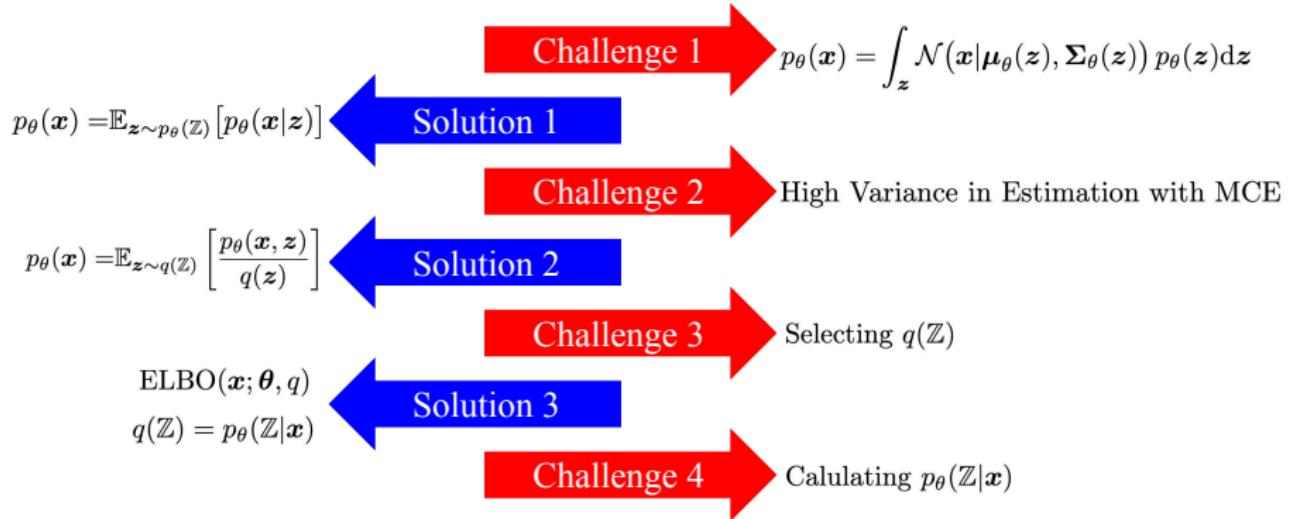
Altogether



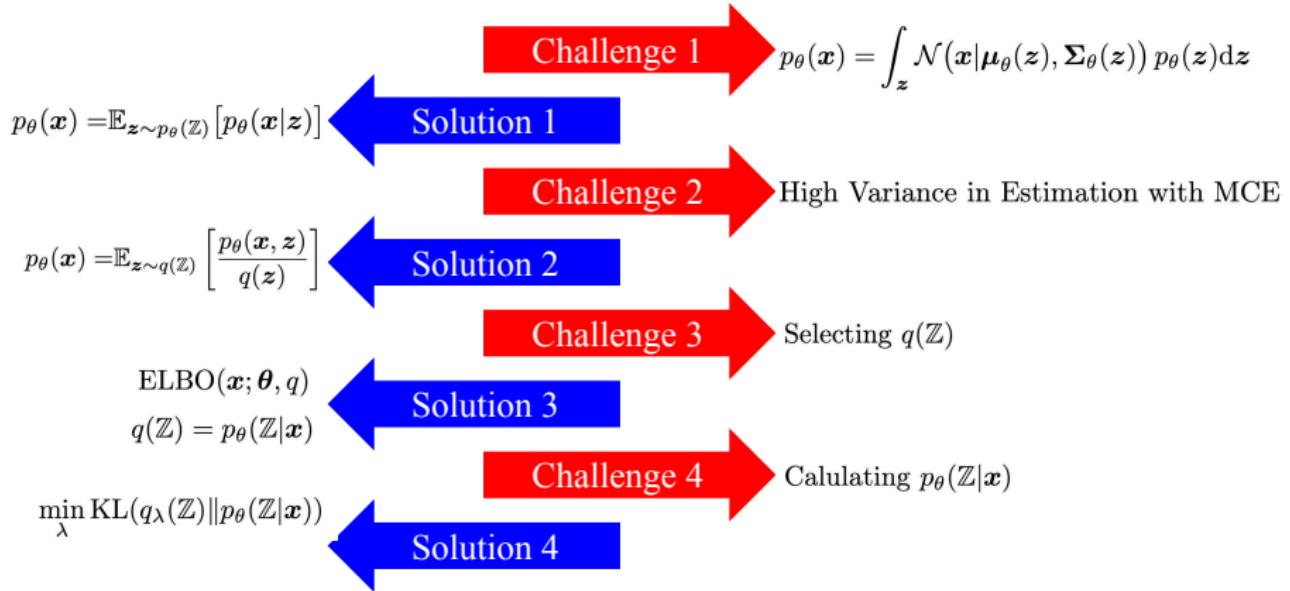
Altogether



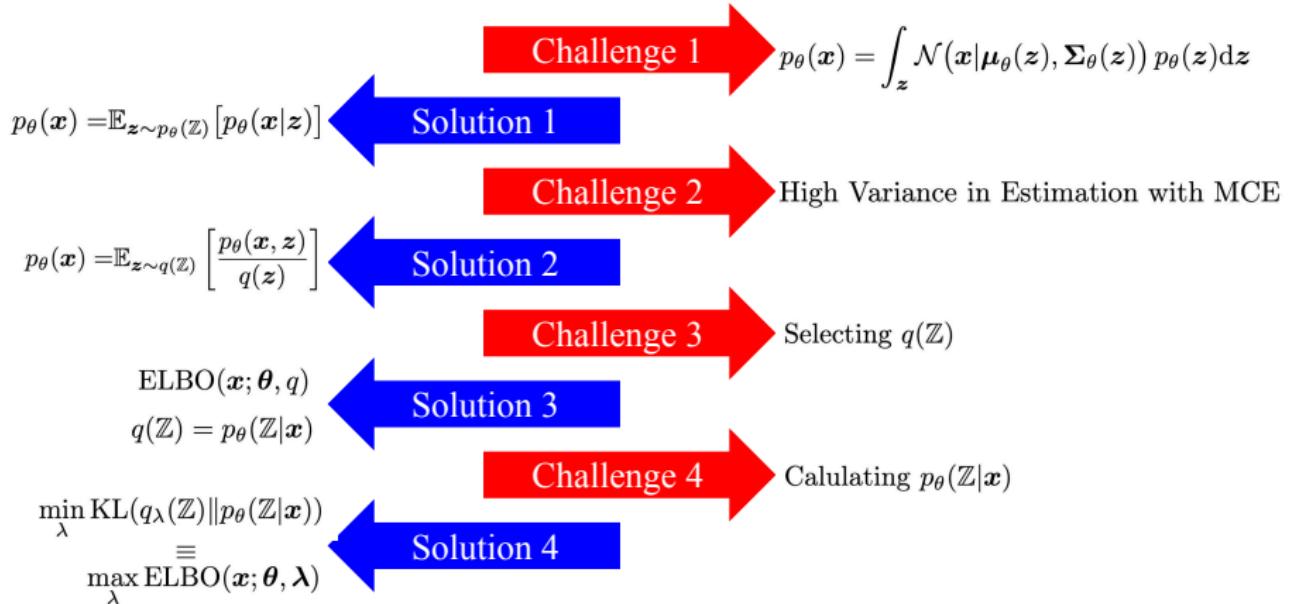
Altogether



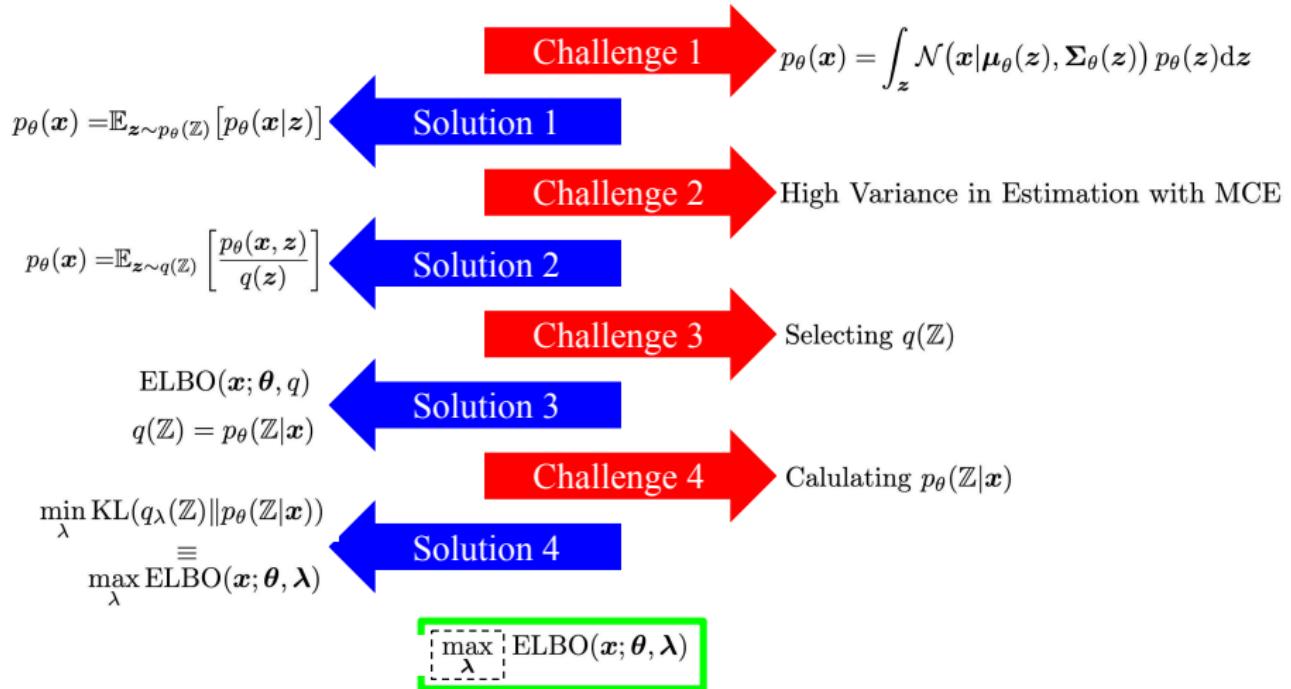
Altogether



Altogether



Altogether



Update on our Challenges

Challenge 4

Calculating $p_\theta(\mathbb{Z}|\mathbf{x})$ needs the data likelihood which is intractable as we observe in Slide 28. On the other hand, we need to decrease $\text{KL}(q(\mathbb{Z})\|p_\theta(\mathbb{Z}|\mathbf{x}))$ to tighten the ELBO.

☞ We can tighten the ELBO using the following optimization problem:

$$\max_{\boldsymbol{\lambda}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$$

Challenge 5

We need to maximize ELBO with respect to both $\boldsymbol{\theta}$ and $\boldsymbol{\lambda}$. Thus we need the following gradient vectors:

- $\nabla_{\boldsymbol{\theta}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$
- $\nabla_{\boldsymbol{\lambda}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$

Section 7

ELBO Optimization

ELBO Gradients

With Respect to θ

For the gradient w.r.t to θ , we have:

$$\begin{aligned}\nabla_{\theta} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) &= \nabla_{\theta} \left(\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\lambda}}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\lambda}}(\mathbf{z})] \right) \\ &\stackrel{a}{=} \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\lambda}}(\mathbf{z})} [\nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\lambda}}(\mathbf{z}))] \\ &\stackrel{b}{=} \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\lambda}}(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})] \\ &\stackrel{c}{\approx} \frac{1}{K} \sum_{k=1}^K \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z}_k)\end{aligned}$$

where:

- $a \Rightarrow$ Expectation is with respect to a distribution independent of θ
- $b \Rightarrow q_{\boldsymbol{\lambda}}(\mathbf{z})$ is independent of θ
- $c \Rightarrow$ Monte-Carlo estimation

ELBO Gradients

With Respect to λ

For the gradient w.r.t to λ , we have:

$$\nabla_{\lambda} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) = \nabla_{\lambda} \left(\mathbb{E}_{q_{\lambda}(Z)} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\lambda}(\mathbf{z})] \right)$$

Note that here we can not change the order of gradient and expectation as the distribution of expectation is a function of λ .

Update on our Challenges

Challenge 5

Calculating $p_\theta(\mathbb{Z}|\mathbf{x})$ needs the data likelihood which is intractable as we observe in Slide 28. On the other hand, we need to decrease $\text{KL}(q(\mathbb{Z})\|p_\theta(\mathbb{Z}|\mathbf{x}))$ to tighten the ELBO.

- ☞ We can tighten the ELBO using the following optimization problem:

$$\max_{\boldsymbol{\lambda}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$$

Challenge 6

We need to maximize ELBO with respect to $\boldsymbol{\lambda}$. Thus we need the following gradient vector:

- $\nabla_{\boldsymbol{\lambda}} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$

Section 8

Reparameterization Trick

Reparameterization Trick

Reparameterization Trick

Again pay attention to the following gradient.

$$\nabla_{\lambda} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) = \nabla_{\lambda} (\mathbb{E}_{q_{\lambda}(Z)} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\lambda}(\mathbf{z})])$$

To handle the above challenge we use the following steps:

- Limit the $q_{\lambda}(\mathbb{Z})$ to Gaussian distribution as:

$$q_{\lambda}(\mathbb{Z}) = \mathcal{N}(\mathbb{Z} | \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \Rightarrow \boldsymbol{\lambda} = \{\boldsymbol{\mu}, \sigma^2\}$$

- Reparameterize the random vector as:

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \Rightarrow \mathbf{z} = \boldsymbol{\mu} + \sigma \boldsymbol{\epsilon} = g(\boldsymbol{\epsilon}; \boldsymbol{\lambda}) \sim \mathcal{N}(\mathbb{Z} \boldsymbol{\mu}, \sigma^2 \mathbf{I}) = q_{\lambda}(\mathbb{Z})$$

ELBO Gradients

With Respect to λ

For the gradient w.r.t to λ , we have:

$$\begin{aligned}\nabla_{\lambda} \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda}) &= \nabla_{\lambda} \left(\mathbb{E}_{\mathbf{z} \sim q_{\lambda}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\lambda}(\mathbf{z})] \right) \\ &\stackrel{a}{=} \nabla_{\lambda} \left(\mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log p_{\theta}(\mathbf{x}, g(\boldsymbol{\epsilon}; \boldsymbol{\lambda})) - \log q_{\lambda}(g(\boldsymbol{\epsilon}; \boldsymbol{\lambda}))] \right) \\ &\stackrel{b}{=} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\lambda} (\log p_{\theta}(\mathbf{x}, g(\boldsymbol{\epsilon}; \boldsymbol{\lambda})) - \log q_{\lambda}(g(\boldsymbol{\epsilon}; \boldsymbol{\lambda})))] \\ &\stackrel{c}{=} \frac{1}{K} \sum_{k=1}^K \nabla_{\lambda} (\log p_{\theta}(\mathbf{x}, g(\boldsymbol{\epsilon}_k; \boldsymbol{\lambda})) - \log q_{\lambda}(g(\boldsymbol{\epsilon}_k; \boldsymbol{\lambda})))\end{aligned}$$

where:

- $a \Rightarrow$ Reparameterization trick
- $b \Rightarrow q_{\lambda}(\mathbf{z})$ Expectation is with respect to a distribution independent of λ
- $c \Rightarrow$ Monte-Carlo estimation

ELBO Optimization

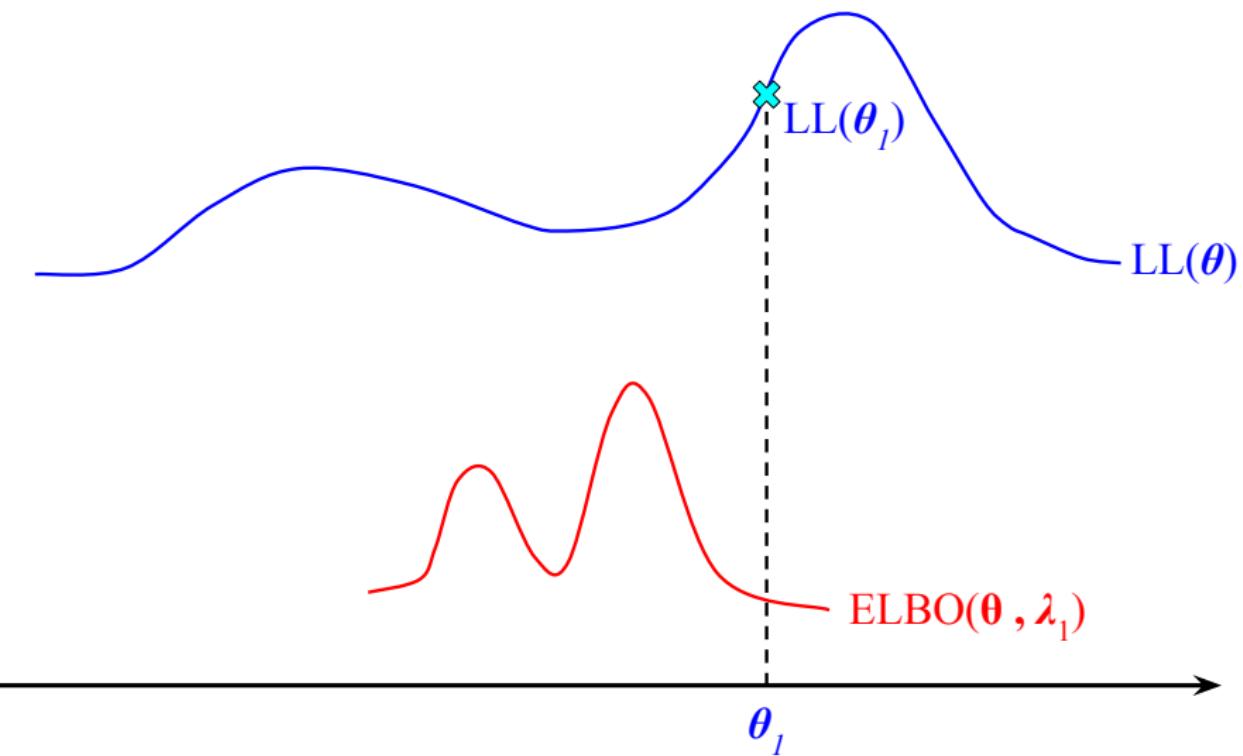


Figure: Starting point θ_1 and λ_1

ELBO Optimization

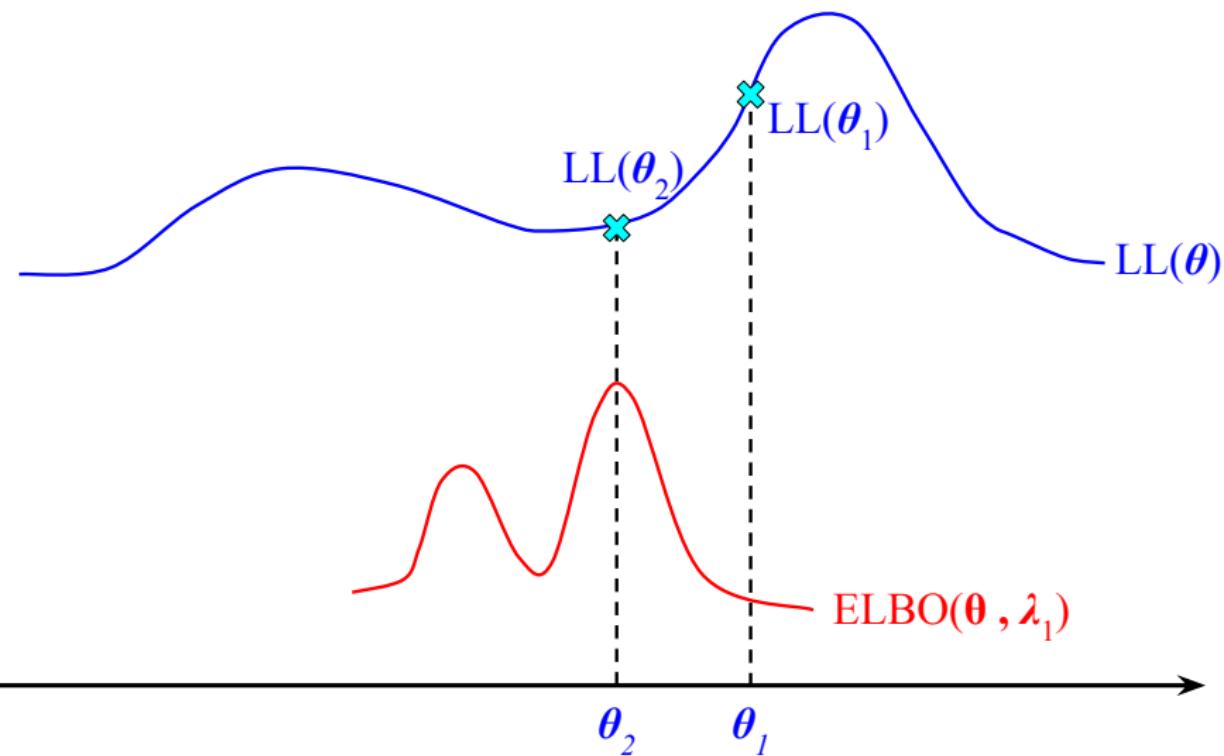


Figure: $\theta_2 = \operatorname{argmax}_{\theta} \text{ELBO}(\theta, \lambda_1)$

ELBO Optimization

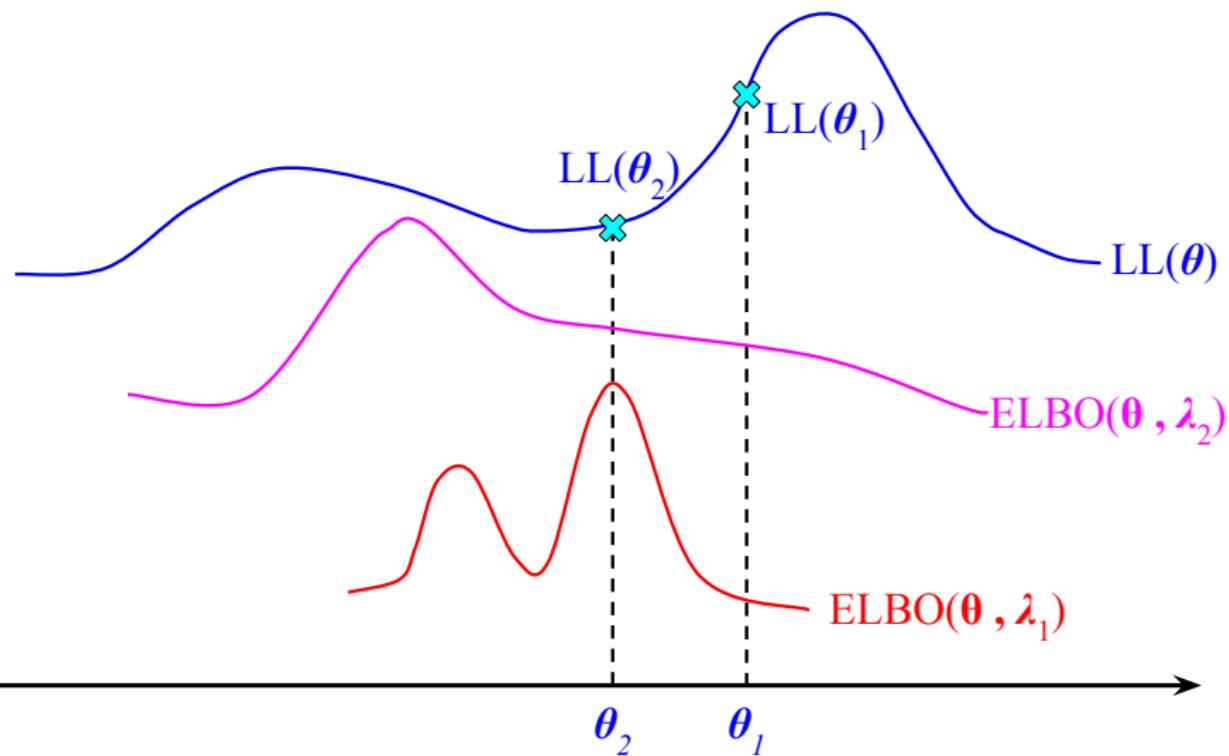


Figure: $\lambda_2 = \operatorname{argmax}_\lambda \text{ELBO}(\theta_2, \lambda)$

ELBO Optimization

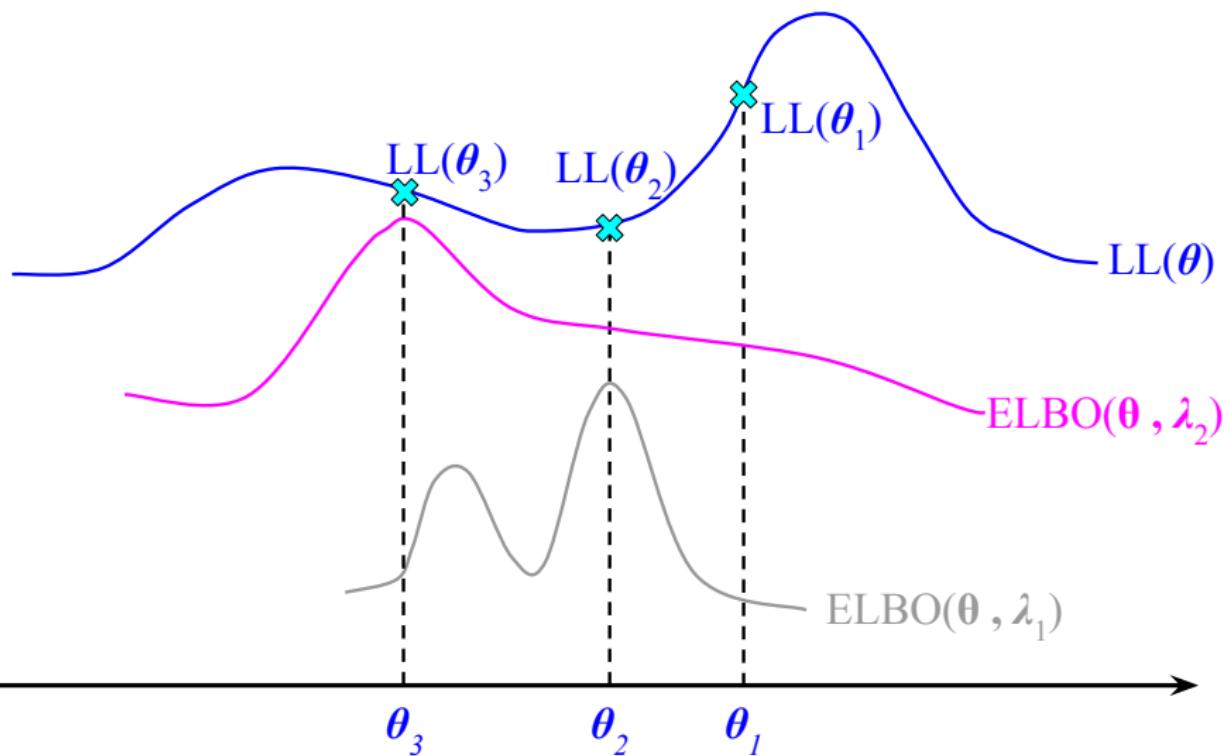


Figure: $\theta_3 = \operatorname{argmax}_{\theta} \text{ELBO}(\theta, \lambda_2)$

Update on our Challenges

Challenge 6

We need to maximize ELBO with respect to λ . Thus we need the following gradient vector:

- $\nabla_\lambda \text{ELBO}(\mathbf{x}; \theta, \lambda)$
- ☞ We can use the reparameterization trick to calculate the above gradient.

Challenge 7

To train a model, we need to maximize the log-likelihood on the dataset \mathcal{D} not an isolated sample \mathbf{x} .

Section 9

Learning

From Data Point to Dataset

Till now, we found a practical approach to increase the ELBO for a typical sample \mathbf{x} as:

$$\log p_{\theta}(\mathbf{x}) \geq \text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$$

Now assume we have a complete dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. We have the following optimization problem:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

Learning

Lower Bounding the Objective

We have:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}_i) &\geq \operatorname{argmax}_{\boldsymbol{\lambda}_i} \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i) \\ \Rightarrow \sum_i \log p_{\theta}(\mathbf{x}_i) &\geq \sum_i \operatorname{argmax}_{\boldsymbol{\lambda}_i} \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i) \\ &= \operatorname{argmax}_{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N} \sum_i \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)\end{aligned}$$

So:

$$\frac{1}{N} \sum_i \log p_{\theta}(\mathbf{x}_i) \geq \operatorname{argmax}_{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N} \frac{1}{N} \sum_i \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$$

Learning

Learning Problem

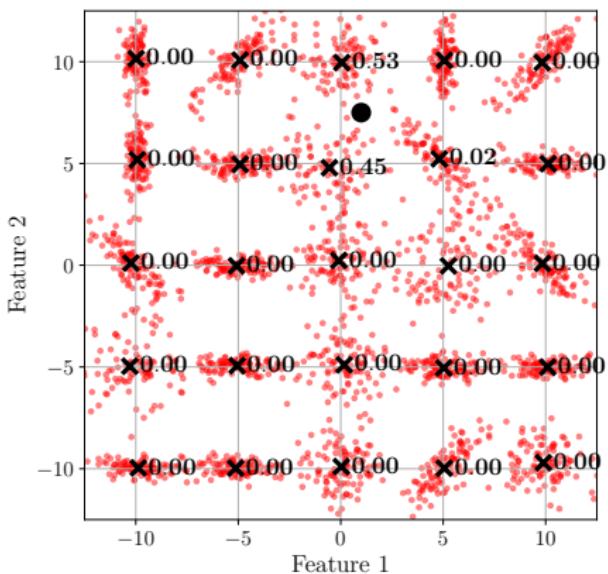
We see:

$$\frac{1}{N} \sum_i \log p_{\theta}(\mathbf{x}_i) \geq \operatorname{argmax}_{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N} \frac{1}{N} \sum_i \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$$

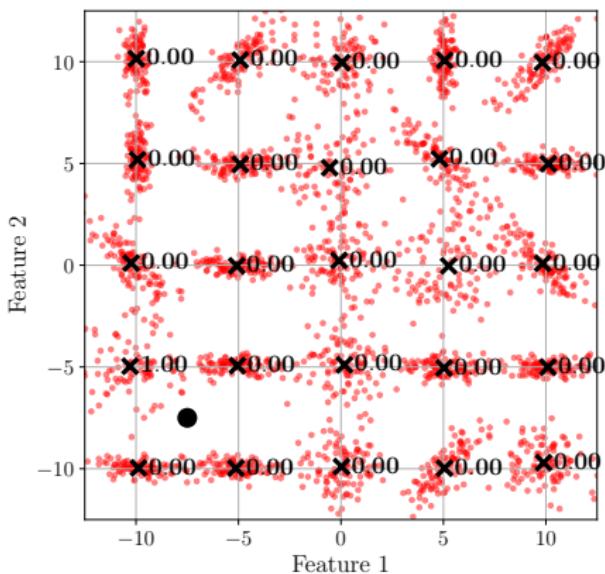
As a result, the learning problem is:

$$\operatorname{argmax}_{\boldsymbol{\theta}} \operatorname{argmax}_{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N} \sum_i \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$$

Different Sampling Distributions for Different Data Samples



(a) $q(Z|x_1)$



(b) $q(Z|x_2)$

Figure: The model posterior over latent variables for two different samples x_1 and x_2 (so we need one distribution for each sample)

Stochastic Variational Inference (SVI) Learning

Algorithm 1: Stochastic Variational Learning

Input : Dataset $\mathcal{D} = \{\mathbf{x}_i\}$
Initialization: $\boldsymbol{\theta}, \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N$
begin
 for $j = 1 : D$ **do**
 Select a Sample \mathbf{x}_i from \mathcal{D} randomly
 for $k = 1 : K$ **do**
 $\boldsymbol{\lambda}_i \leftarrow \boldsymbol{\lambda}_i + \eta \nabla_{\boldsymbol{\lambda}_i} \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$
 end
 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mu \nabla_{\boldsymbol{\theta}} \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$
 end
end
Output : $\boldsymbol{\theta}$

Update on our Challenges

Challenge 7

To train a model, we need to maximize the log-likelihood on the dataset \mathcal{D} , not an isolated sample \boldsymbol{x} .

☞ We derive the objective for the complete dataset as:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \underset{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N}{\operatorname{argmax}} \sum_i \text{ELBO}(\boldsymbol{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$$

Challenge 8

As we can see in the above optimization problem, the number of parameters is not scalable with the dataset size. Each training sample adds one $\boldsymbol{\lambda}$ vector to the optimization problem.

Section 10

Amortization

Amortization

Scalability of SVI

Using the SVI, we can solve the following maximization problem:

$$\boldsymbol{\theta}^*, \boldsymbol{\lambda}_1^*, \boldsymbol{\lambda}_2^*, \dots, \boldsymbol{\lambda}_N^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \underset{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N}{\operatorname{argmax}} \sum_i \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\lambda}_i)$$

Amortization is a type of parameter sharing to achieve scalability with respect to dataset size.

Amortization

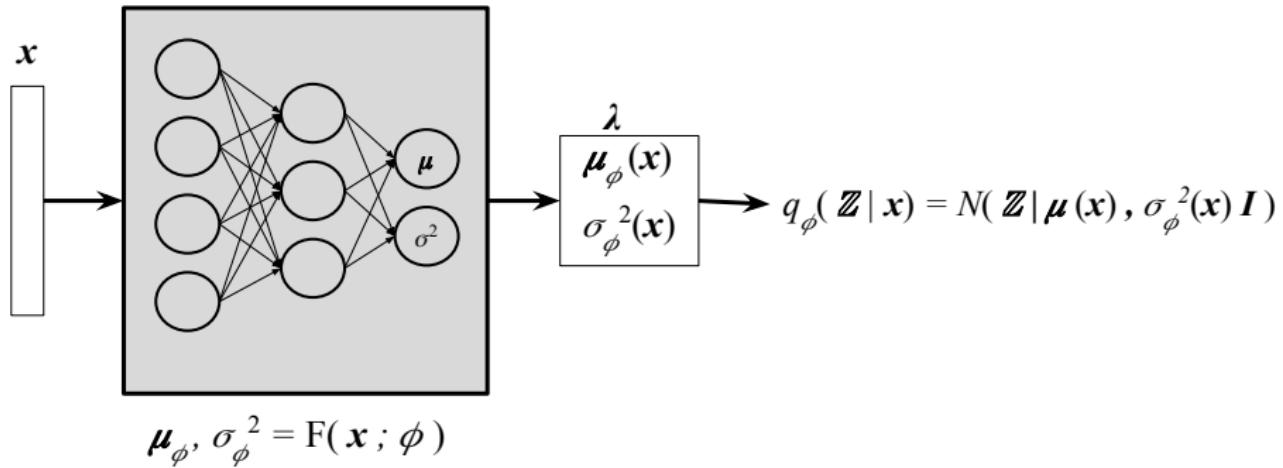


Figure: Amortization: The parameters of neural network $F(\cdot; \cdot)$ is shared between all the training samples.

Amortization

Notation Update

By using the Amortization technique, ϕ the parameters of F are shared between all training data samples. So instead of $\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\lambda})$ we use:

$$\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbb{Z}|\mathbf{x})} \left[\log \left(\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \right) \right]$$

Amortized Inference

Using Amortization, the ELBO for sample \mathbf{x}_i can be written as:

$$\text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbb{Z}|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}_i)]$$

Amortization

Algorithm 2: Amortized Stochastic Variational Inference

Input : Dataset $\mathcal{D} = \{\mathbf{x}_i\}$

Initialization: θ, ϕ

begin

for $j = 1 : D$ **do**

 Select a Sample \mathbf{x}_i from \mathcal{D} randomly

$\phi \leftarrow \phi + \eta \nabla_{\phi} \text{ELBO}(\mathbf{x}_i; \theta, \phi)$

$\theta \leftarrow \theta + \mu \nabla_{\theta} \text{ELBO}(\mathbf{x}_i; \theta, \phi)$

end

end

Output : θ, ϕ

Section 11

Revisiting Objective Function

Variational Autoencoder

Revisiting ELBO

We can reformulate ELBO as:

$$\begin{aligned}\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) &= \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[\underbrace{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log p_{\boldsymbol{\theta}}(\mathbf{z})}_{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})} + \underbrace{\log p_{\boldsymbol{\theta}}(\mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}_{-\log \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{z})}} \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{z}))\end{aligned}$$

Variational Autoencoder Overview

Encoder

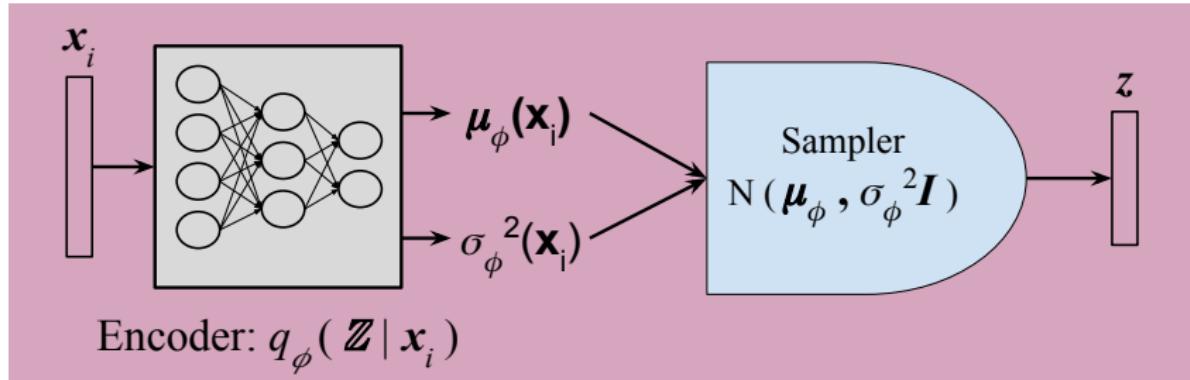


Figure: Encoding

Variational Autoencoder Overview

Encoder

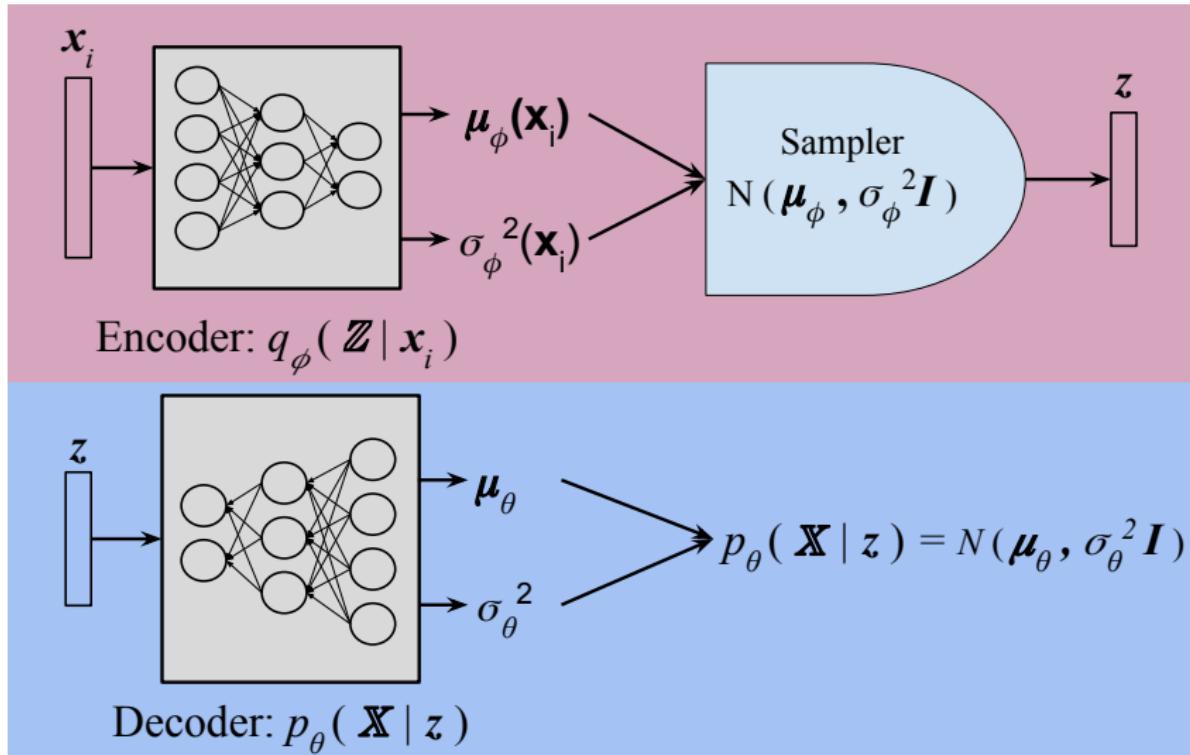


Figure: Encoding and Decoding

Variational Autoencoder

Revisiting ELBO

Let's focus on the first term:

$$\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbb{Z} | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | \mathbf{z})]$$

Conceptually, we have:

- Select a training sample \mathbf{x}_i .
- Generate the approximate posterior over the latent vector given \mathbf{x}_i
- Sample \mathbf{z} from the approximate posterior
- Generate the model posterior over \mathbb{X} given \mathbf{z} or $p_{\theta}(\mathbb{X} | \mathbf{z})$

A good model is one where $\log p_{\theta}(\mathbf{x}_i | \mathbf{z})$ is maximized.

Variational Autoencoder

Revisiting ELBO

Let's focus on the first term:

$$\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | \mathbf{z})]$$

In terms of formulations, we have:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}_i | \mathbf{z}) &= \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{\theta}, \sigma_{\theta}^2 \mathbf{I}) \\ &= -\frac{D}{2} \log(2\pi\sigma_{\theta}^2) - \frac{1}{2\sigma_{\theta}^2} \|\mathbf{x}_i - \boldsymbol{\mu}_{\theta}\|_2^2\end{aligned}$$

Note that:

- The optimization w.r.t. $\boldsymbol{\theta}$ is straightforward.
- The optimization w.r.t. $\boldsymbol{\phi}$ needs reparameterization trick.

Variational Autoencoder with Reparameterization Trick

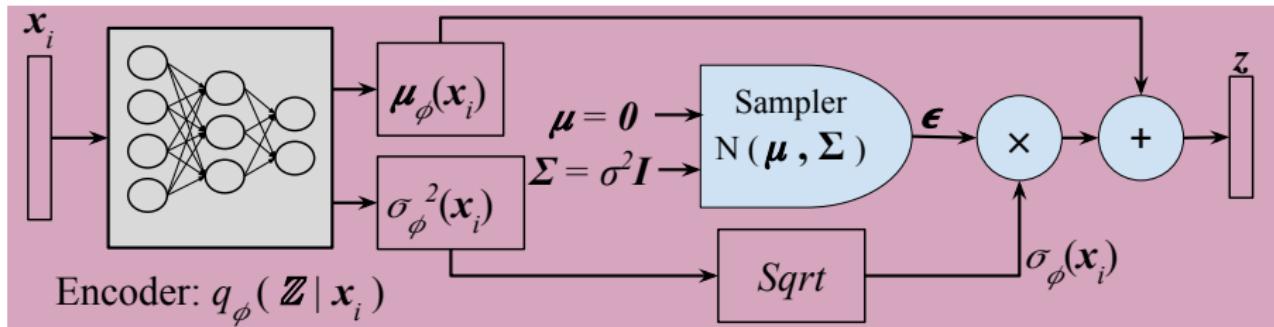


Figure: Encoding with Reparameterization Trick

Variational Autoencoder with Reparameterization Trick

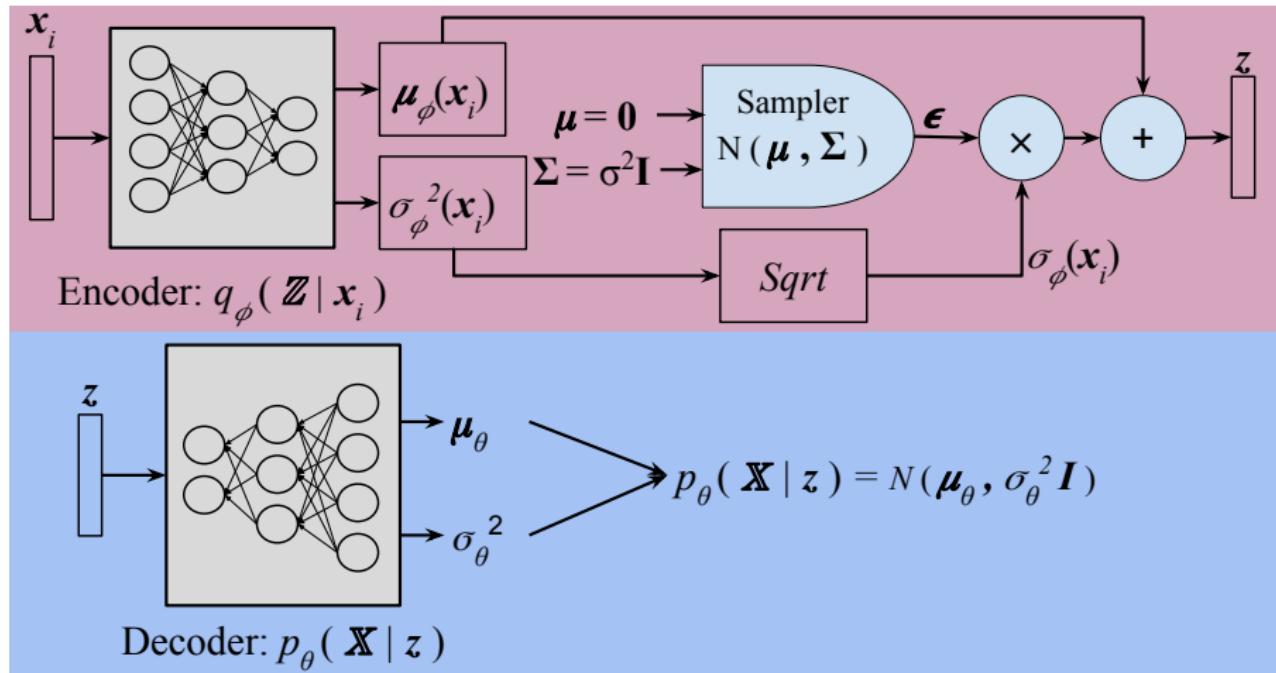


Figure: Encoding and Decoding

Revisiting ELBO

Now let's move on to the second term:

$$-\text{KL}(q_\phi(\mathbb{Z}|\mathbf{x}_i) \| p_\theta(\mathbb{Z}))$$

After training and at the end of the day:

- $q_\phi(\mathbb{Z}|\mathbf{x}_i), i = 1, \dots, N$ be near the prior distribution $p_\theta(\mathbb{Z})$ in KLD sense.

What does it imply?

Variational Autoencoder

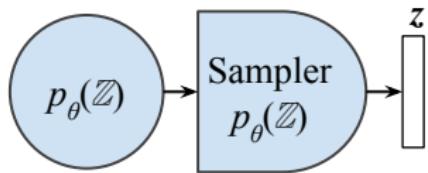


Figure: Encoding with Reparameterization Trick

Variational Autoencoder

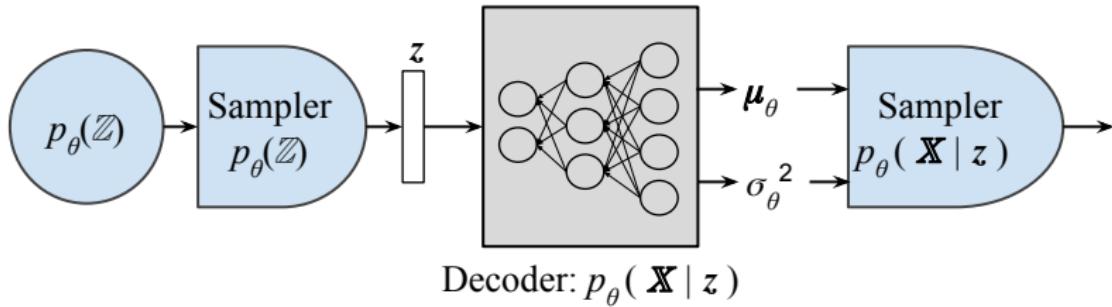


Figure: Encoding with Reparameterization Trick

Variational Autoencoder

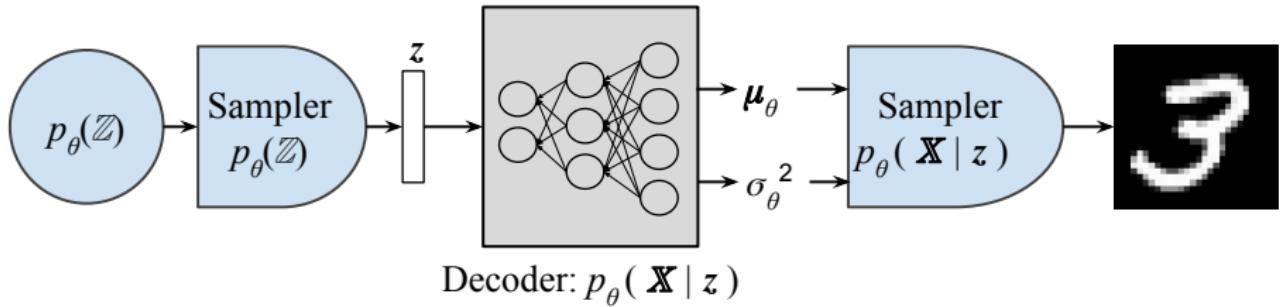


Figure: Encoding with Reparameterization Trick

Section 12

Attribute Vectors in Code Space

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

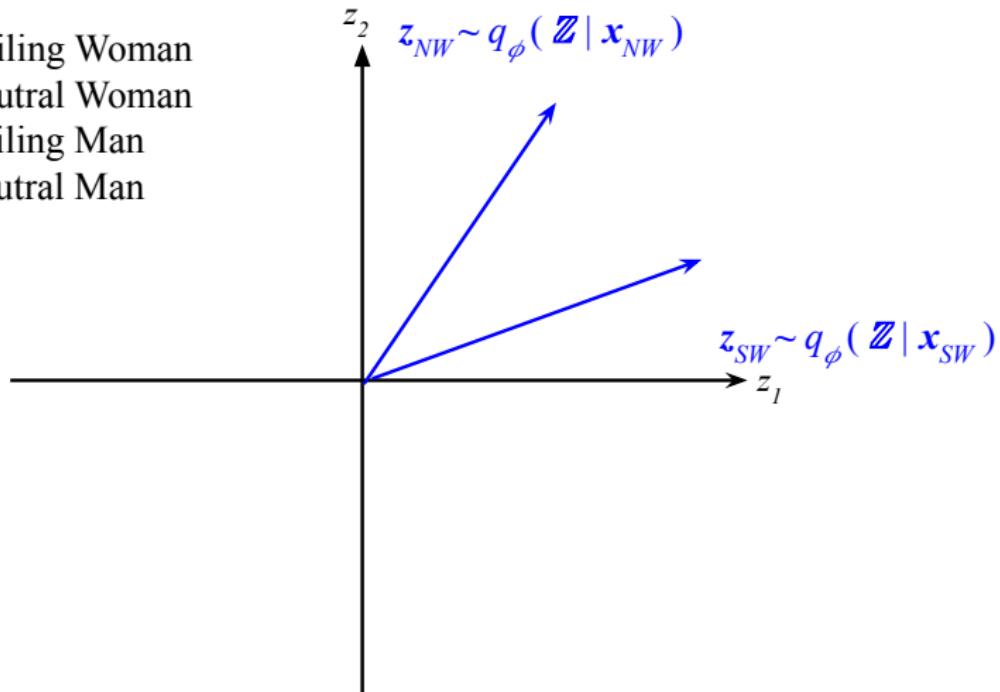


Figure: Encoding of *Neutral Woman* and *Smiling Woman* images

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

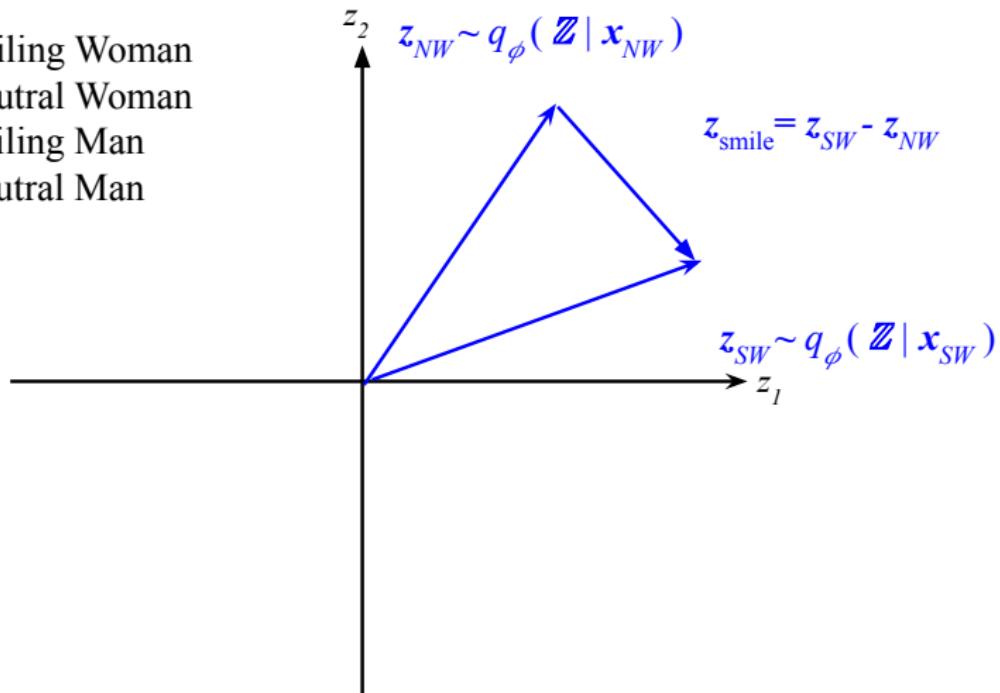


Figure: Calculating the attribute vector of *Smile* in latent space

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

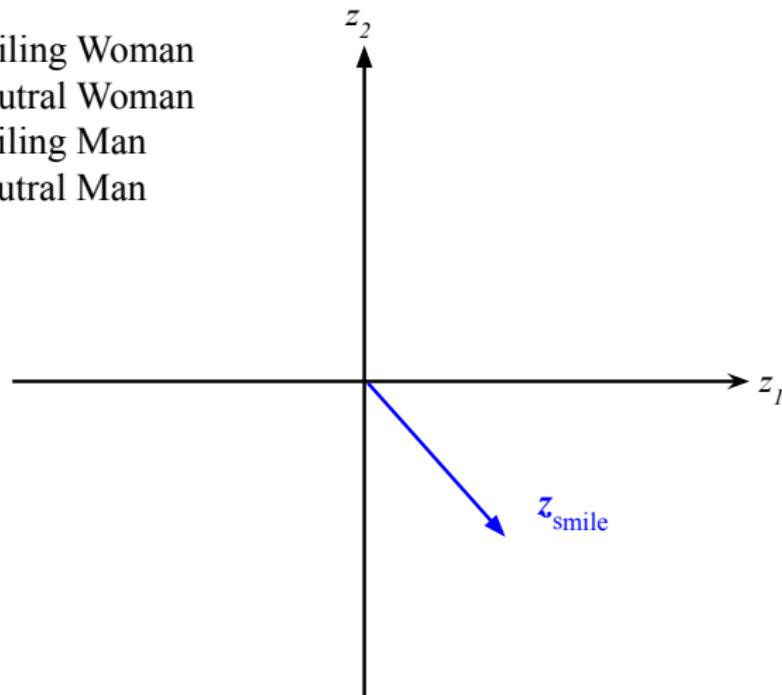


Figure: *Smile* attribute vector

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

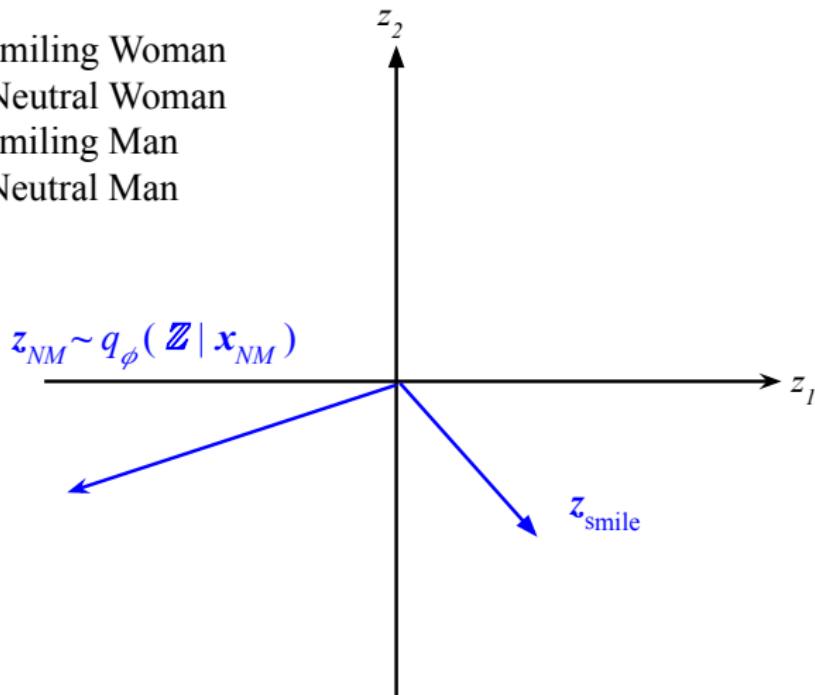


Figure: Encoding of Neutral Man image

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

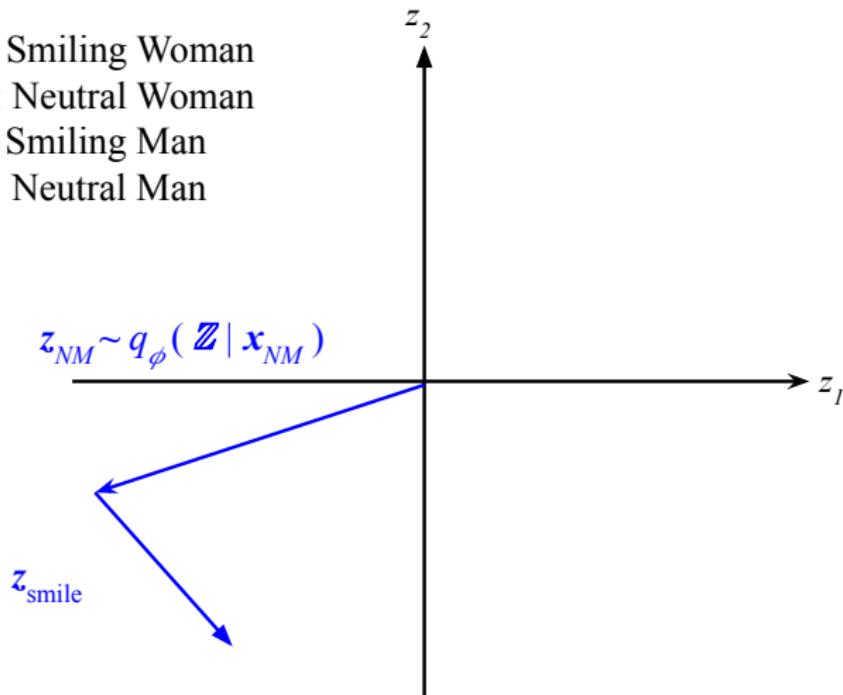


Figure: Adding *Smile* attribute vector to *Neutral Man* latent vector

Concept

SW: Smiling Woman
NW: Neutral Woman
SM: Smiling Man
NM: Neutral Man

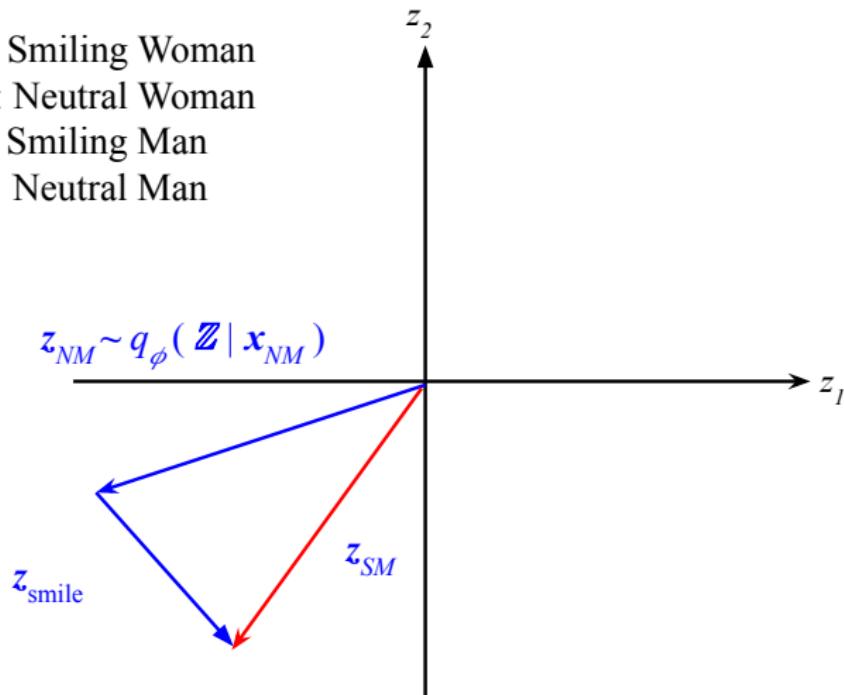


Figure: The latent vector of *Smiling Man*

Sample of Using Attribute Vectors

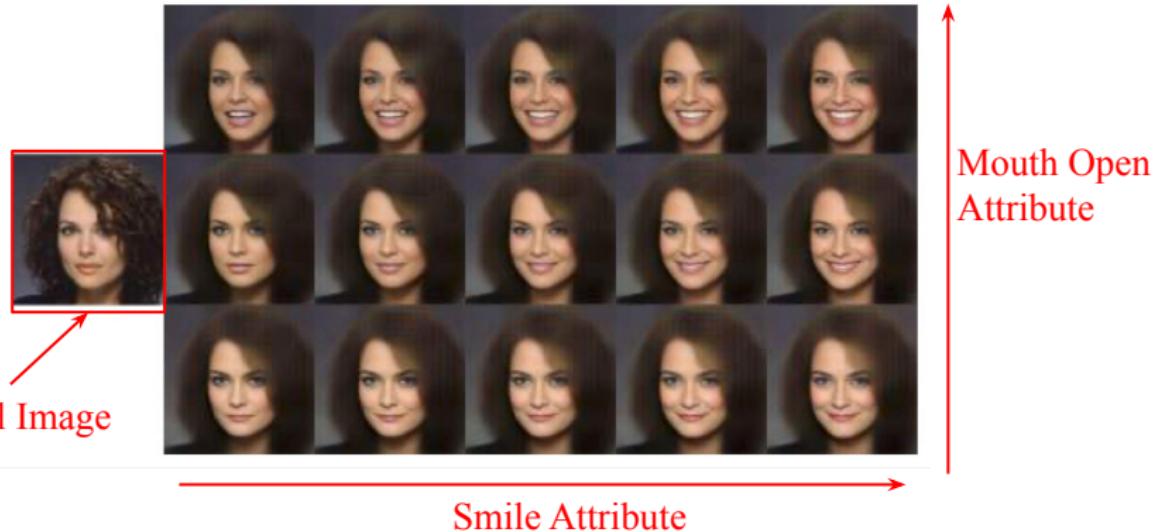


Figure: Application of attribute vectors in real worl application of face images
(source: [2])

List of Abbreviations

| Complete | Abbreviation |
|----------------------------------|--------------|
| Evidence Lower-BOund | ELBO |
| Gaussian Mixture Model | GMM |
| Kullback–Leibler | KL |
| Left-Hand Side | LHS |
| Log-Likelihood | LL |
| Right-Hand Side | RHS |
| Stochastic Variational Inference | SVI |
| Variational AutoEncoder | VAE |

References I

-  M. Fink and P. Perona,
“Caltech 10k web faces (1.0) [data set].,” <https://doi.org/10.22002/D1.20132>.
-  T White,
“Sampling generative networks: notes on a few effective techniques corr (2016),”
arXiv preprint arXiv:1609.04468.