

Machine Learning 2

**Variational Inference
- Variational AutoEncoders (VAE)**

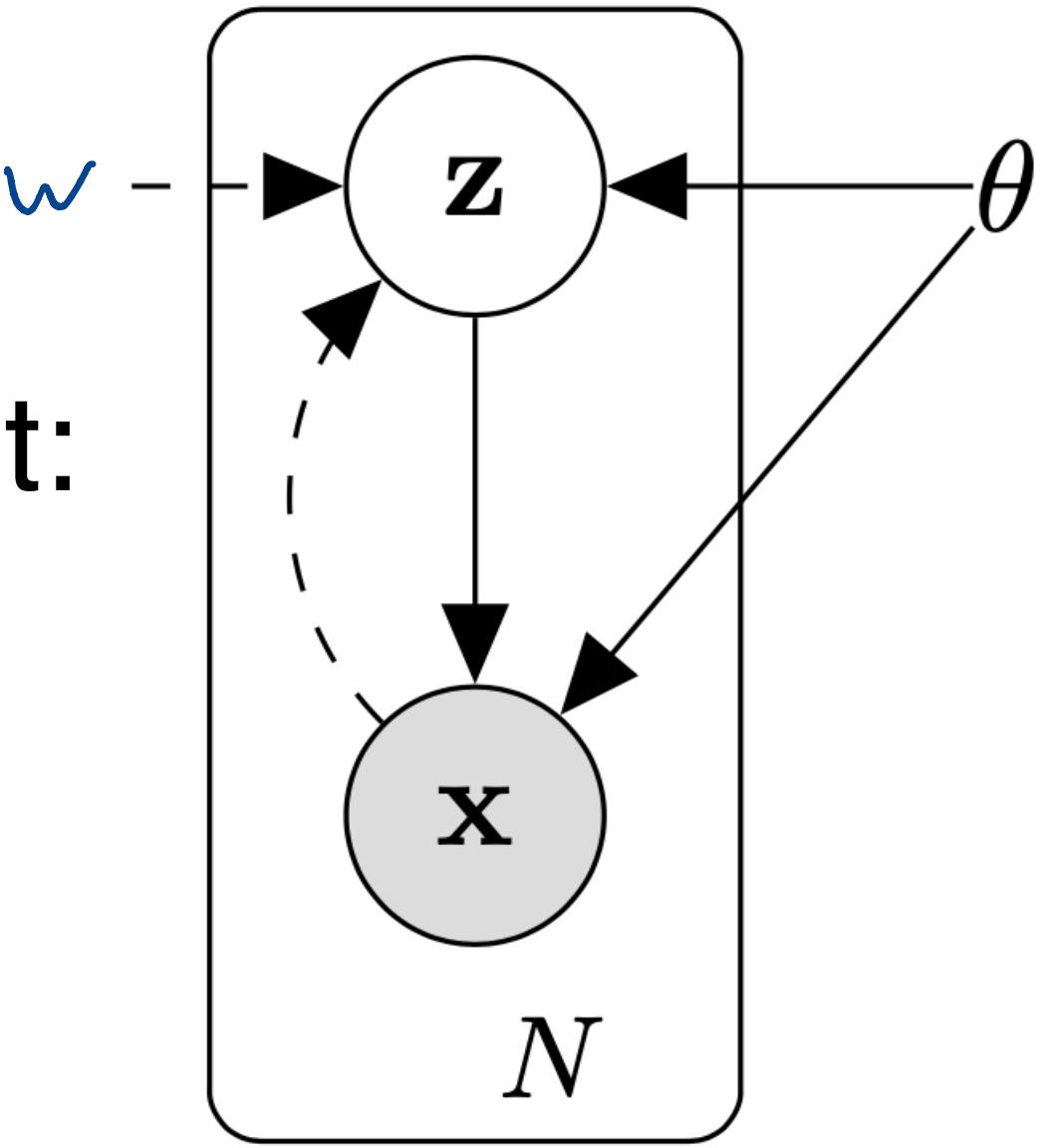
Patrick Forré

Variational AutoEncoders - Setup

- Let $q(x)$ be an underlying ‘true’ distribution in a **high dimensional space** \mathcal{X} and \hat{X} data sampled from it, e.g. **images**.
- We consider the ‘support’ of $q(x)$ to lie close to a **lower dimensional ‘manifold’**, which we consider given by a **latent variable**.
- So we assume a **lower dimensional latent** space \mathcal{Z} that ‘parameterizes’ the ‘data manifold’.
- Since **deep neural networks** are universal approximators we hope that they can learn the ‘data manifold’.
- We also introduce an **inference network** (the ‘encoder’) such that we can use Variational Inference for training: **maximizing ELBO**.

Variational AutoEncoders

- To maximize ELBO: $L_{\hat{X}}(w, \theta) = \mathbb{E}_{\hat{q}(X)} \mathbb{E}_{q_w(Z|X)} \left[\log \frac{p_\theta(Z, X)}{q_w(Z|X)} \right]$ we put:
 - $p(z) = \mathcal{N}(z | 0, I)$
 - $p_\theta(x | z) = \mathcal{N}(x | \mu_{\theta_1}(z), \sigma_{\theta_2}^2(z))$ *generator / decoder*
 - $q_w(z | x) = \mathcal{N}(z | \mu_{w_1}(x), \sigma_{w_2}^2(x))$ *encoder*
 - where $\mu_{\theta_1}(z), \sigma_{\theta_2}(z), \mu_{w_1}(x), \sigma_{w_2}(x)$ are parameterized by **neural networks.**
 - The ELBO is then maximized via **backpropagation** on mini-batches.



Reparameterization Trick

- To evaluate the expectation $\mathbb{E}_{q_w(Z|X)}[f(Z)]$ in the ELBO one approximates it with an empirical version based on samples:
 - For each $x^{(n)}$ sample M i.i.d. $\epsilon^{(n,1)}, \dots, \epsilon^{(n,M)} \sim \mathcal{N}(0,I)$ and put $z^{(n,m)} := \mu_{w_1}(x^{(n)}) + \sigma_{w_2}(x^{(n)}) \cdot \epsilon^{(n,m)}$, which are then i.i.d. $\sim \mathcal{N}(z | \mu_{w_1}(x), \sigma_{w_2}^2(x))$,
- separating parameters from random noise for back-propagation. We get:

$$\mathbb{E}_{q_w(Z|x^{(n)})}[f(Z)] \approx \frac{1}{M} \sum_{m=1}^M f(z^{(n,m)}) = \frac{1}{M} \sum_{m=1}^M f\left(\mu_{w_1}(x^{(n)}) + \sigma_{w_2}(x^{(n)}) \cdot \epsilon^{(n,m)}\right)$$

Training the VAE

- To maximize ELBO: $L_{\hat{X}}(w, \theta) = \mathbb{E}_{q(X)} \mathbb{E}_{q_w(Z|X)} \left[\log \frac{p_\theta(Z, X)}{q_w(Z|X)} \right]$ do until convergence:

- Sample mini-batch from data: $x^{(1)}, \dots, x^{(N)} \sim q(X)$
- For each n sample mini-batch:

$$z^{(n,1)}, \dots, z^{(n,M)} \sim \mathcal{N}(z | \mu_{w_1}(x^{(n)}), \sigma_{w_2}^2(x^{(n)}))$$

using the reparameterization trick.

- Forward pass: $\hat{L}(\theta, w) = \frac{1}{NM} \sum_{n=1}^N \sum_{m=1}^M \left[\log \frac{p_\theta(z^{(n,m)}, x^{(n)})}{q_w(z^{(n,m)} | x^{(n)})} \right]$

- Update parameters via Back-propagation (α learning rate)
 - $[\theta, w]^\top := [\theta, w]^\top + \alpha \cdot \nabla_{\theta, w} \hat{L}(\theta, w)$

Interpreting the ELBO

- The ELBO can be written as:

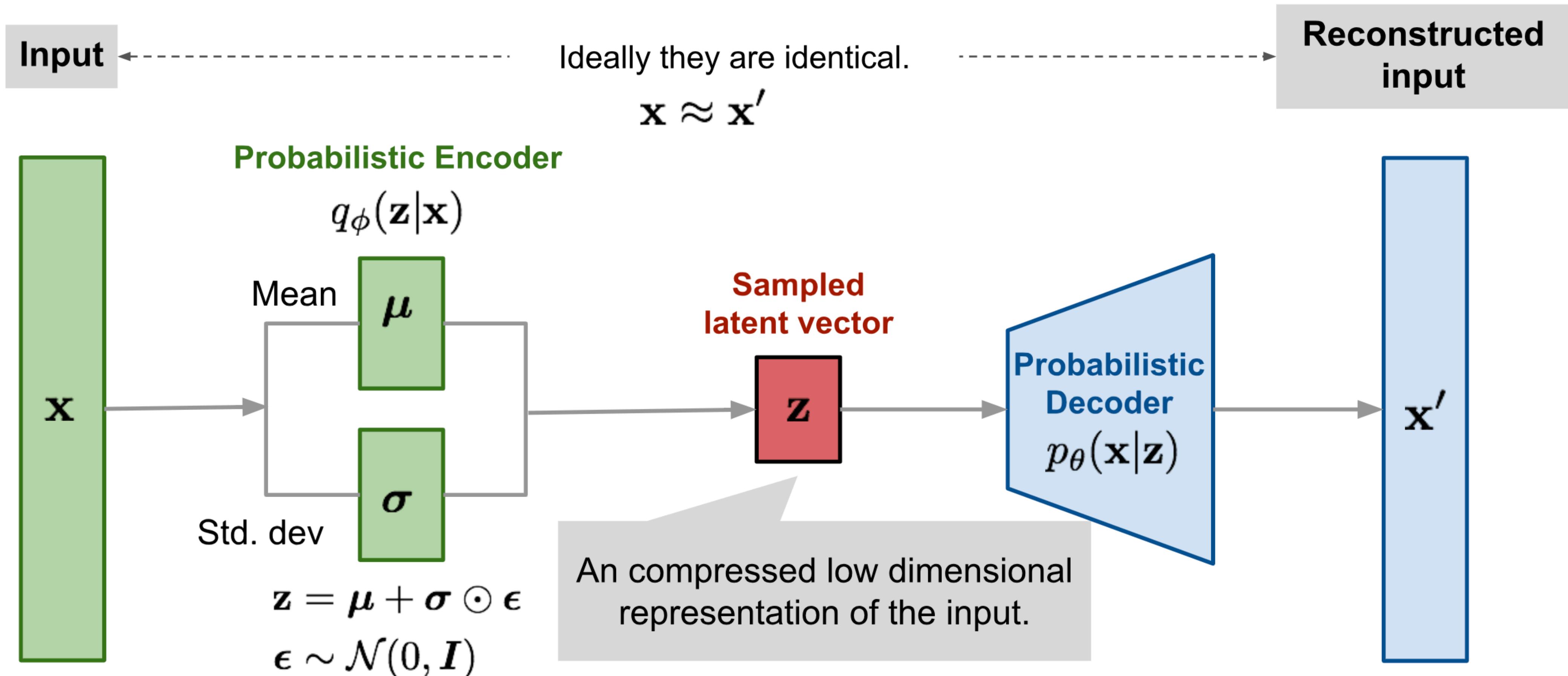
$$L_{\hat{X}}(w, \theta) = \mathbb{E}_{q(x)} \left[\mathbb{E}_{q_w(z|x)} \left[\log p_\theta(x|z) \right] - \mathbb{E}_{q(x)} [\text{KL}(q_w(z|x) || p(z))] \right]$$

$x \sim z \sim \hat{x}$
 $\sim -\|x - \hat{x}\|_2^2$
reconstruction

regularization.

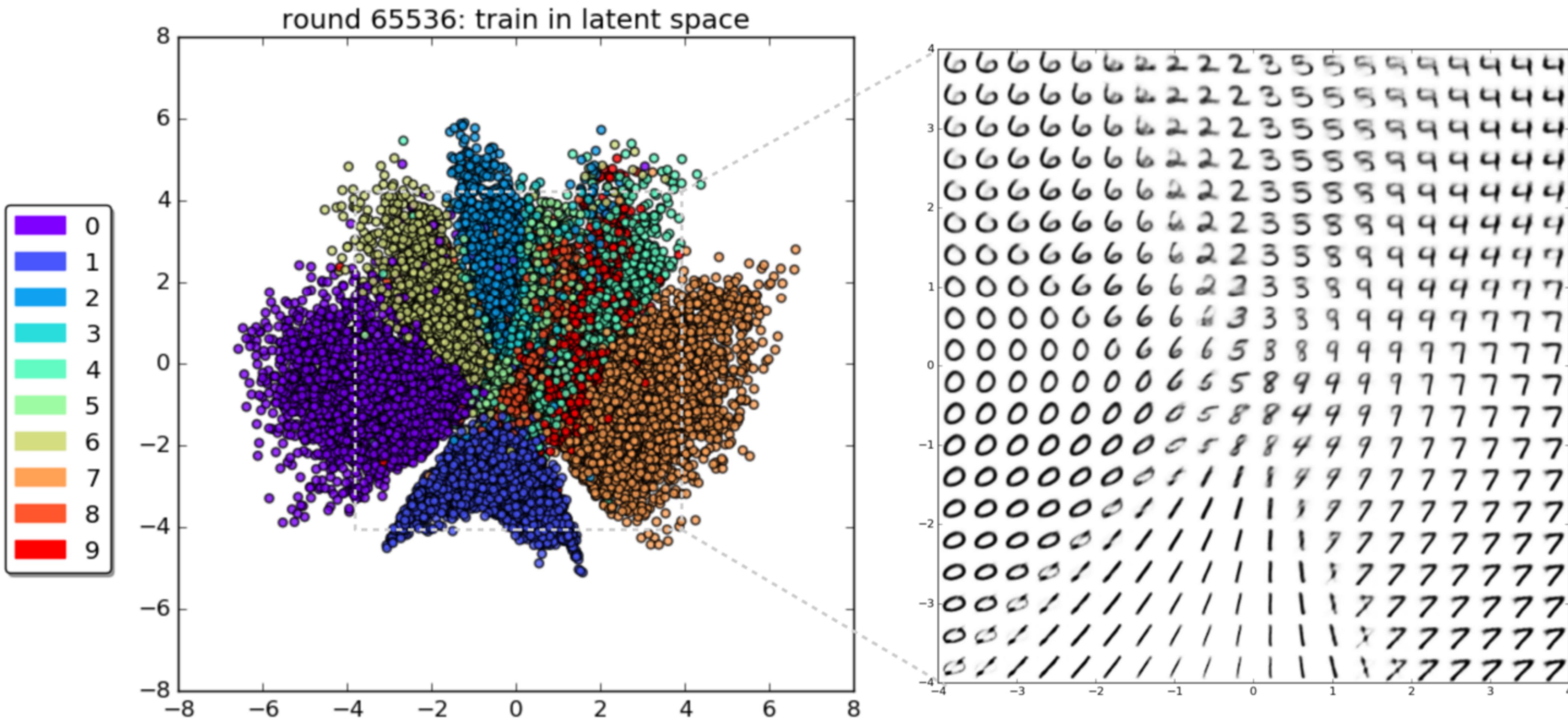
- splitting in reconstruction term, and
- regularization term

Training the VAE



- Source: <https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>

VAE - Latent Space



- <https://blog.fastforwardlabs.com/2016/08/12/introducing-variational-autoencoders-in-prose-and-code.html>

Machine Learning 2

Variational Inference
- **Variational Bayes /**
/ Mean Field Approximation

Patrick Forré

Recap: The Framework of Variational Inference (VI)

- Let $q(x)$ be an underlying ‘true’ distribution of interest and \hat{X} data sampled from it.
Let $\{p_\theta(x, z) \mid \theta \in \Theta\}$ be a **latent variable** model for $q(x)$ (or a Bayesian setting).

- Instead of **maximizing the evidence/log-likelihood** $\log p_\theta(\hat{X})$ w.r.t. θ ,
- in Variational Inference (VI) one introduces a variational **family of inference distributions** $\{q_w(z \mid x) \mid w \in \mathcal{W}\}$ and **Maximizes the Evidence Lower BOund** (ELBO):

$$L_{\hat{X}}(w, \theta) := \mathbb{E}_{q_w(Z \mid \hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z \mid \hat{X})} \right] \approx \leq \log p_\theta(\hat{X})$$

w.r.t. both distributions / parameters w, θ .

Variational Bayes - Mean Field Approximation

- Goal: Maximize ELBO wrt w, θ : $L_{\hat{X}}(w, \theta) = \mathbb{E}_{q_w(Z|\hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z|\hat{X})} \right]$
- Latent variables $Z = (Z_1, \dots, Z_K)$ with several components or component blocks.
- **Variational Bayes / Mean Field Approximation**: Family of inference distribution to be taken as **product distributions**:

$$q_w(z_1, \dots, z_K | \hat{X}) = \prod_{k=1}^K q_{w_k}(z_k | \hat{X}) \quad \approx \quad p_\theta^{(z | \hat{X})}$$

- This is said to be an “approximation”, as $q_w(z | \hat{X})$ is considered an approximation for the ‘true’ posterior: $p_\theta(z | \hat{X})$ in the Bayesian setting and $L_{\hat{X}}(w, \theta)$ after fitting for evidence: $\log p_\theta(\hat{X})$.
- Note that this is a ‘**high** entropy approximation’ as: $H(Z_1, \dots, Z_K) \leq H(Z_1) + \dots + H(Z_K)$.

Optimizing w.r.t. \mathbf{q}

dependent on w
and \hat{x}

- Abbrev. inference distributions with: $q(z) = q_1(z_1) \cdots q_K(z_K)$

- Langrangian (ELBO plus constraints, θ fixed):

$$\begin{aligned}
 \mathcal{L}(q) &= \mathbb{E}_{q(z)} \left[\log \frac{p(z, \hat{x})}{q(z)} \right] + \sum_{n=1}^N \lambda_n \cdot \left(\int q_n(z_n) dz_n - 1 \right) \\
 &= \mathbb{E}_{q(z)} [\log p(z, \hat{x})] - \mathbb{E}_{q(z)} [\log q(z)] + \sum_{n=1}^N \lambda_n \cdot \left(\int q_n(z_n) dz_n - 1 \right) \\
 &= \int q_1(z_1) \cdots q_N(z_N) \cdot \log p(z, \hat{x}) dz - \sum_{n=1}^N \int q_n(z_n) \cdot \log q_n(z_n) dz_n + \sum_{n=1}^N \lambda_n \left(\int q_n(z_n) dz_n - 1 \right)
 \end{aligned}$$

Functional derivatives w.r.t. q

$$0 \stackrel{!}{=} \frac{\partial \mathcal{L}}{\partial q_k(z_k)} = \int_{j \neq k} \prod_j q_j(z_j) \cdot \log p(z, \hat{x}) dz_{\neq k} - \log q_k(z_k) - 1 + \lambda_k$$

$$\mathcal{L}(q) = \int q_1(z_1) \cdots q_n(z_n) \cdot \log p(z, \hat{x}) dz - \sum_{u=1}^n \int q_u(z_u) \cdot \log q_u(z_u) dz_u + \sum_{u=1}^n \lambda_u \left(\int q_u(z_u) - 1 \right)$$

$$\begin{aligned} q_k(z_k) &= \exp(\lambda_k - 1) \cdot \exp \left(\int_{j \neq k} \prod_j q_j(z_j) \log p(z, \hat{x}) dz_{\neq k} \right) \\ &\quad \times \exp \left(\underbrace{\mathbb{E}_{q_{\neq k}(z_{\neq k})} \left[\log p(z_{\neq k}, z_k, \hat{x}) \right]}_{\text{red wavy line}} \right) \end{aligned}$$

VB-MF Fixpoint Equations

- We arrive at the following fixpoint equations for $k = 1, \dots, K$:

$$q_k(z_k) \propto \exp \left(\mathbb{E}_{q(Z_{\neg k})} \left[\log p(\hat{X}, z_k, Z_{\neg k}) \right] \right)$$

- Note that the rhs is dependent on the all other $q_j(z_j)$ for $j \neq k$.
- We are left to solve this system of fixpoint equation:
 - Often from this an iterative update rule for corresponding parameters can be derived that converges.

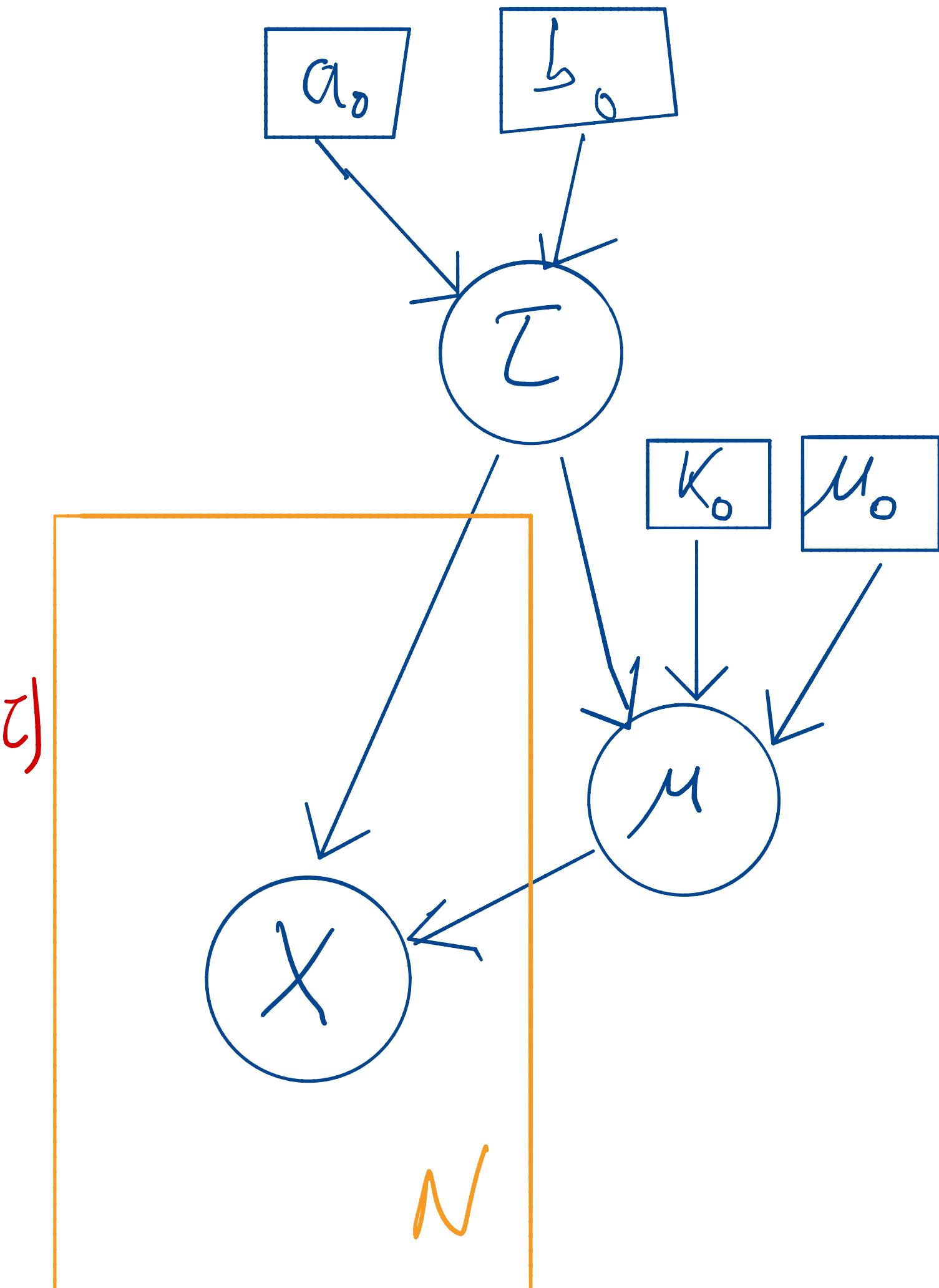
Machine Learning 2

Variational Inference
- Variational Bayes /
/ Mean Field Approximation
- Example

Patrick Forré

$$\begin{aligned}
 & \underbrace{x_1, \dots, x_N}_{\hat{x}} \stackrel{iid}{\sim} p(x|\mu, \tau) = N(x|\mu, \tau^{-1}) \\
 & p(\mu|\mu_0, K_0, \tau) = N(\mu|\mu_0, (K_0 \cdot \tau)^{-1}) \\
 & p(\tau|a_0, b_0) = F(\tau|a_0, b_0)
 \end{aligned}$$

$$\begin{aligned}
 \log p(\hat{x}, \mu, \tau) &= \log p(\hat{x}|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau) \\
 &= \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \\
 &\quad + \frac{1}{2} \log(K_0 \tau) - \frac{K_0 \tau}{2} (\mu - \mu_0)^2 \\
 &\quad + (a_0 - 1) \cdot \log \tau - b_0 \tau + \text{const.}
 \end{aligned}$$



Goal: $p(\mu, \tau | \hat{x}) \approx q_1(\mu) \cdot q_2(\tau)$

VB-MF Fixpoint equations:

$$\log q_1(\mu) = \mathbb{E}_{q_2(\tau)} [\log p(\hat{x}, \mu, \tau)] + \text{const.}$$

$$\log q_2(\tau) = \mathbb{E}_{q_1(\mu)} [\log p(\hat{x}, \mu, \tau)] + \text{const.}$$

$$\log q_1(\mu) = \mathbb{E}_{q_2(\tau)} [\log p(\hat{x}, \mu, \tau)] + \text{const.}$$

$$\log p(\hat{x}, \mu, \tau) = \frac{N}{2} \cancel{\log \tau} - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 + \frac{1}{2} \cancel{\log(K_0 \tau)} - \frac{x_0 \tau}{2} (\mu - \mu_0)^2 + (q_0 - 1) \cancel{\log \tau} - \cancel{b_0 \tau} + \text{const.}$$

$$\begin{aligned} \log q_1(\mu) &= \mathbb{E}_{q_2(\tau)} \left[-\frac{\tau}{2} \cdot \left(\sum_{n=1}^N (x_n - \mu)^2 + K_0 (\mu - \mu_0)^2 \right) \right] + \text{const} \\ &= -\frac{\mathbb{E}[\tau]}{2} \cdot \left(K_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right) + \text{const.} \end{aligned}$$

$$= -\frac{\tau_N}{2} (\mu - \mu_N)^2 + \text{const.}$$

$$\mu_N = \frac{K_0 \cdot \mu_0 + \sum_{n=1}^N x_n}{K_0 + N} \quad \begin{matrix} \leftarrow \text{known} \\ | \end{matrix} \quad \tau_N = (K_0 + N) \cdot \mathbb{E}_{q_2(\tau)} [\tau] \quad \begin{matrix} \leftarrow \text{unknown} \\ | \end{matrix}$$

$$\log q_1(\mu) = -\frac{\tau_N}{2} (\mu - \mu_N)^2 + \text{const.}$$

$$\Rightarrow q_1(\mu) = N(\mu | \mu_N, \tau_N^{-1})$$

where : $\mu_N = \frac{\lambda_0 \mu_0 + \sum_{n=1}^N x_n}{\lambda_0 + N}$

$$\tau_N = (\lambda_0 + N) \cdot \mathbb{E}_{q_2(c)}[c]$$

$$\log q_2(\tau) = \mathbb{E}_{q_1(\mu)} [\log p(\hat{x}, \mu, \tau)] + \text{const.}$$

$$\log p(\hat{x}, \mu, \tau) = \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 + \frac{1}{2} \log (k_0 \tau) - \frac{k_0 \tau}{2} (\mu - \mu_0)^2 \\ + (a_0 - 1) \cdot \log \tau - b_0 \cdot \tau + \text{const.}$$

$$\log q_2(\tau) = \frac{N}{2} \log \tau + \frac{1}{2} \log \tau + (a_0 - 1) \log \tau - b_0 \cdot \tau \\ - \frac{\tau}{2} \mathbb{E}_{q_1(\mu)} \left[\sum_{n=1}^N (x_n - \mu)^2 + k_0 (\mu - \mu_0)^2 \right] + \text{const.} \\ = \left(a_0 + \frac{N+1}{2} - 1 \right) \log \tau - \underbrace{\left[b_0 + \frac{1}{2} \mathbb{E}_{q_1(\mu)} \left[\sum_n (x_n - \mu)^2 + k_0 (\mu - \mu_0)^2 \right] \right]}_{b_N} \cdot \tau + \text{const.} \\ = (a_N - 1) \cdot \log \tau - b_N \cdot \tau + \text{const}$$

$$\log q_2(\tau) = (a_N - 1) \log \tau - b_N \cdot \tau + \text{const.}$$

$$\Rightarrow q_2(\tau) = P(\tau | a_N, b_N)$$

where

$$a_N = a_0 + \frac{N+1}{2} \quad \leftarrow \text{known}$$

$$\begin{aligned} b_N &= b_0 + \frac{1}{2} E_{\mu}[\mu] \left[K_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right] \\ &= b_0 + \frac{K_0}{2} \left(E_{\mu}[\mu^2] + \mu_0^2 - 2 E_{\mu}[\mu] \cdot \mu_0 \right) \\ &\quad + \frac{1}{2} \sum_{n=1}^N (x_i^2 + E_{\mu}[\mu^2] - 2 E_{\mu}[\mu] \cdot x_i) \end{aligned}$$

$$\rightarrow E_{\tau}[\tau] = \frac{a_N}{b_N}$$

unknown

$$q_2(\mu) = N(\mu | \mu_N, \sigma_N^{-2}), \quad E_{\mu}[\mu] = \mu_N, \quad E_{\mu}[\mu^2] = \frac{1}{\sigma_N^2} + \mu_N^2$$

Known

$$\mu_N$$

to be determined

$$c_N = (k_0 + N) \cdot \frac{a_N}{b_N} \Rightarrow c_N^{-1} = \frac{b_N}{(k_0 + N) \cdot a_N}$$

$$a_N$$

$$b_N = b_0 + \frac{k_0}{2} \left(\frac{1}{c_N} + \mu_N^2 + \mu_0^2 - 2 \cdot \mu_N \cdot \mu_0 \right)$$

$$+ \frac{1}{2} \sum_{n=1}^N \left[x_n^2 + \frac{1}{c_N} + \mu_N^2 - 2 \mu_N \cdot x_n \right]$$

$$= c_N^{-1} \left(\frac{k_0 + N}{2} \right) + \text{rest}$$

$$\Rightarrow b_N = \frac{1}{2} \frac{b_N}{a_N} + b_0 + \frac{k_0}{2} (\mu_N - \mu_0)^2 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu_N)^2$$

$$\Leftrightarrow b_N = \left(1 - \frac{a_N}{2} \right)^{-1} \left[b_0 + \frac{k_0}{2} (\mu_N - \mu_0)^2 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu_N)^2 \right]$$

$$\Rightarrow c_N = (k_0 + N) \frac{a_N}{b_N}$$

$$q_1(\mu) = N(\mu | \mu_N, \tau_N^{-1})$$

$$\text{and} \quad q_2(z) = P(z | a_N, b_N)$$

where

$$a_N = a_0 + \frac{N+1}{2}$$

$$\mu_N = \frac{k_0 \mu_0 + \sum_{n=1}^N x_n}{k_0 + N}$$

$$b_N = \left(1 - \frac{1}{2a_N}\right)^{-1} \left[b_0 + \frac{k_0}{2} (\mu_N - \mu_0)^2 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu_N)^2 \right]$$

$$\tau_N = (k_0 + N) \frac{a_N}{b_N}$$

$$P(\mu, \tau | \tilde{x}) \approx N(\mu | \mu_N, \tau_N) \cdot \Gamma(\tau | a_N, b_N)$$

$$\log P(x) \approx \text{ELBO}(\mathcal{N}(\mu|\mu_N, \tau_N^{-1}) \cdot P(z|a_N, b_N))$$

Machine Learning 2

Variational Inference
- Universal Variational Inference
as its own learning framework

Patrick Forré

Universal Variational Inference (UVI) - Setting

- Let $q(x)$ be an underlying ‘true’ distribution of interest and $\hat{X} = [x^{(1)}, \dots, x^{(N)}]^T$ data sampled from it.
- We aim at:
 - ‘learning the data distribution’,
 - prediction,
 - quantifying model uncertainty,
 - incorporating prior knowledge
 - regularization
 - computational efficiency, etc.

Universal Variational Inference - Ingredients

- In the general setting of UVI we need to specify:
 1. **Statistical model** $\mathcal{P} = \{p_\theta(x) = p(x | \theta) \mid \theta \in \Theta\}$ parameterized by θ .
 2. **Loss function**: $\mathcal{L}(\hat{X} | p_\theta)$ measuring the goodness of fit between (discrete) **data points** \hat{X} and (implicit/explicit) **model** p_θ .
 3. **Inference model**: $\mathcal{Q} = \{q_w(\theta) = q(\theta | w) \mid w \in \mathcal{W}\}$, which is usually given by comp. restrictions, hard constraints or simplifying assumptions.
 4. **Regularizing functional**: $\Omega(q)$ for $q \in \mathcal{Q}$ addressing uncertainty, prior knowledge, outcome expectations, soft constraints, complexity, etc.
 5. Optimizer \mathcal{O} and evalution scheme \mathcal{E} for the following:

UVI - Learning Objective + Predictive Distribution

- After p_θ , \mathcal{L} , \mathcal{Q} , Ω are specified we use an **optimizer** \mathcal{O} to solve:

$$q_{\hat{X}} \in \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} | p_\theta)] + \Omega(q) \right)$$

- Finally, for **predictions** do **model averaging** (via \mathcal{E}):

$$p(x | \hat{X}) := \int p(x | \theta) \cdot q_{\hat{X}}(\theta) d\theta$$

goal:
 $\approx q(x)$

Remark: Loss functions

- The loss function often comes in form:

$$\mathcal{L}(\hat{X} | p_\theta) = \sum_{n=1}^N \ell(x^{(n)} | p_\theta) \quad \text{with:}$$

- $\ell(\hat{x} | p_\theta) = -\log p_\theta(\hat{x})$ (log-loss)
- $\ell(\hat{x} | p_\theta) = \frac{1}{2} \left(\hat{x} - \mathbb{E}_{p_\theta(X)}[X] \right)^2$ (square-loss)
- $\ell(\hat{x} | p_\theta) = \alpha \cdot [\hat{x} - q_\alpha^{p_\theta}]_+ + (1 - \alpha) \cdot [q_\alpha^{p_\theta} - \hat{x}]_+$ (α -quantile loss)
- $\ell(\hat{x} | p_\theta) = 1 - p_\theta(\hat{x})$ (misclassification rate)
- $\ell(\hat{x} | p_\theta) = - (1 - p_\theta(\hat{x}))^\gamma \cdot \log p_\theta(\hat{x})$ (focal loss)
- Maximum-Mean-Discrepancies, Optimal Transport divergencies, Stein divergencies, etc.

Remark: Inference models

- Possible Inference models: $\mathcal{Q} = \{q_w(\theta) = q(\theta | w) \mid w \in \mathcal{W}\}$ are
 - Point estimators: $q_w(\theta) = \delta_w(\theta)$ for $w \in \Theta = \mathcal{W}$
 - Neural networks
 - Product distributions: $q_w(\theta) = q_{w_1}(\theta_1) \cdots q_{w_K}(\theta_K)$
 - All distributions

Remark: Regularizing functional

- Possible Regularizing functionals:
 - $\Omega(q) = 0$
 - $\Omega(q) = C(q\|\pi)$ with a prior $\pi(\theta)$
 - $\Omega(q) = \frac{1}{\beta} \text{KL}(q\|\pi)$ with a prior $\pi(\theta)$ and ‘inverse temperature’ β
 - $\Omega(q) = \pm H(q)$
 - $\Omega(q) = \mathbb{E}_{q(\theta)} [d(p_\theta\|p_0)]$ with reference model $p_0(x)$ and divergence d .

Remark: UVI for the supervised setting

- Assumption of true distribution $q(y, x)$ or at least $q(y | x)$
- Statistical model $p_\theta(y, x) = p_\theta(y | x) \cdot q(x)$,
where $q(x)$ is either assumed to be known, given or ignored or
coincides with empirical version $\hat{q}(x)$.
- Note we are here interested in $q(y | x)$ and not so much in $q(x)$.
- Then everything goes similar in the UVI framework.

Machine Learning 2

Variational Inference

- Universal Variational Inference**
- recovering special cases**

Patrick Forré

Standard Variational Inference via UVI

$$1. \mathcal{P} = \{p_{\theta}(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_{\theta}) = -\log p(\hat{x} \mid \theta)$$

$$3. \mathcal{Q} = \{\text{any/all}\}$$

$$4. \Omega(q) = KL(q \parallel \pi), \pi(\theta) \text{ prior}$$

$$5. \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_{\theta})] + \Omega(q) \right) = \mathbb{E}_{q(\theta)} \left[-\log p(\hat{x} \mid \theta) \right] + \mathbb{E}_{q(\theta)} \left[\log \frac{q(\theta)}{\pi(\theta)} \right]$$

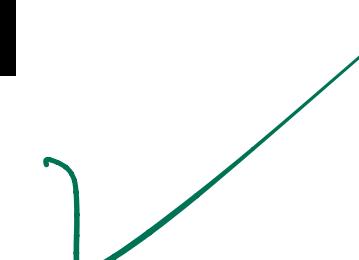
= $\mathbb{E}_{q(\theta)} \left[-\log p(\hat{x} \mid \theta) \right] + \mathbb{E}_{q(\theta)} \left[\log \frac{q(\theta)}{\pi(\theta)} \right]$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta \cancel{\neq}$$

$$\boxed{-ELBO(q)}$$

$$= -\mathbb{E}_{q(\theta)} \left[\log \frac{p(\hat{x} \mid \theta)}{q(\theta)} \right]$$

$$= \mathbb{E}_{q(\theta)} \log \frac{q(\theta)}{p(\hat{x} \mid \theta) \pi(\theta)}$$



Full Bayesian approach via UVI

$$1. \mathcal{P} = \{p_{\theta}(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_{\theta}) = -\log p(x \mid \theta)$$

$$3. \mathcal{Q} = \text{all } ?$$

$$4. \Omega(q) = \text{KL}(q \parallel \pi), \text{ prior } \pi(\theta)$$

$$5. \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_{\theta})] + \Omega(q) \right) = \arg \min_{q \in \mathcal{Q}} \mathbb{E}_{q(\theta)} \left[\log \frac{q(\theta)}{p(\theta \mid \hat{x})} \right]$$

= $p(\theta \mid \hat{x})$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta = \int p(x \mid \theta) p(\theta \mid \hat{x}) d\theta$$

Maximum Likelihood Estimation via UVI

$$1. \mathcal{P} = \{p_{\theta}(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_{\theta}) = -\log p_{\theta}(\hat{X})$$

$$3. \mathcal{Q} = \{ \delta_w(\theta) \mid w \in \tilde{\Theta} \}$$

$$4. \Omega(q) = 0$$

$$5. \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_{\theta})] + \Omega(q) \right) = \underset{\delta_w}{\operatorname{argmin}} -\log p_w(\hat{X})$$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta = p(x \mid \Theta_{MLE})$$

✓

Maximum A Posteriori Estimation via UVI

$$1. \mathcal{P} = \{p_{\theta}(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_{\theta}) = -\log p_{\theta}(\hat{x})$$

$$3. \mathcal{Q} = \{ \delta_w(\theta) \mid w \in \Omega \}$$

$$4. \Omega(q) = C(q \parallel \pi) + \text{prior } \pi(\theta)$$

$$5. \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_{\theta})] + \Omega(q) \right) = \arg \min_{\delta_w} -\log p_w(\hat{x}) - \log \pi(w)$$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta = p(x \mid \theta_{MAP})$$

Empirical Risk Minimization via UVI

$$1. \mathcal{P} = \{p_\theta(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_\theta) = \text{any}$$

$$3. Q = (\delta_w(\theta) \mid w \in \Theta)$$

$$4. \Omega(q) = 0 = \delta_{\Theta_{ERM}}(\theta)$$

$$5. \arg \min_{q \in Q} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_\theta)] + \Omega(q) \right) = \arg \min_{\delta_w} \mathcal{L}(\hat{X} \mid p_w)$$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta = p(x \mid \Theta_{ERM})$$

Variational Bayes - Mean Field Approximation via UVI

$$1. \mathcal{P} = \{p_{\theta}(x) \mid \theta \in \Theta\}$$

$$2. \mathcal{L}(\hat{X} \mid p_{\theta}) = -\log p_{\theta}(\hat{x})$$

$$3. \mathcal{Q} = \{q_w(\theta) = q_{w_1}(\theta_1) \cdot \dots \cdot q_{w_n}(\theta_n)\}$$

$$4. \Omega(q) = KL(q \parallel \pi), \quad \pi(\theta) \text{ prior.}$$

$$5. \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} \mid p_{\theta})] + \Omega(q) \right) = \operatorname{argmax}_{q \in \mathcal{Q}} \text{ELBO}(q)$$

$$6. p(x \mid \hat{X}) = \int p(x \mid \theta) \cdot q_{\hat{X}}(\theta) d\theta$$

Machine Learning 2

Variational Inference - Summary

Patrick Forré

Variational Inference

- True distribution $q(x)$, data \hat{X} , ~~latent variable model~~ or ~~Bayesian setting~~:
 $\{p_\theta(z, x) \mid \theta \in \Theta\}.$
- Introduce ~~inference model~~ $\{q_w(z \mid x) \mid w \in \mathcal{W}\}$ meant to approximate $p_\theta(z \mid \hat{X})$
- Key objective in VI is ~~Maximizing the Evidence Lower Bound~~ (ELBO):

$$L_{\hat{X}}(w, \theta) := \mathbb{E}_{q_w(Z \mid \hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z \mid \hat{X})} \right] \leq \log p_\theta(\hat{X})$$

w.r.t. parameters w and θ (if present).

$$\text{I = I} \quad \text{iff} \quad q_w(z \mid \hat{x}) = p_\theta(z \mid \hat{x})$$

EM-Algorithm and VAEs

- **Expectation-Maximization (EM) Algorithm:** **Alternate:**

- **E-step:** Evaluate: $Q(\theta, \theta^{(t)}) := \mathbb{E}_{p_{\theta^{(t)}}(Z|\hat{X})} \left[\log p_\theta(Z, \hat{X}) \right]$
- **M-step:** Compute: $\theta^{(t+1)} := \arg \max_{\theta} Q(\theta, \theta^{(t)}) + \text{constraints.}$

- **Variational AutoEncoders (VAE):**

- Parameterize everything with **deep neural networks.**
- Maximize ELBO with SGD and **reparameterization trick.**

Variational Bayes Mean Field Approximation

- **Variational Bayes Mean Field Approximation:**

$$p(z_1, \dots, z_K | \hat{X}) \approx \prod_{k=1}^K q_{w_k}(z_k | \hat{X}) =: q_w(z | \hat{X})$$

- lead to **VB-MF Fixpoint Equations:** for all $k = 1, \dots, K$:

$$q_k(z_k) \propto \exp \left(\mathbb{E}_{q(Z_{\neg k})} \left[\log p(\hat{X}, z_k, Z_{\neg k}) \right] \right)$$

- **solve system** for (the parameters of) the $q_k(z_k)$'s **explicitly** or via **iterative updating**.

Universal Variational Inference Learning Framework

- In the general setting of UVI we need to specify:
 1. **Statistical model**: $\mathcal{P} = \{p_\theta(x) = p(x | \theta) \mid \theta \in \Theta\}$
 2. **Loss function**: $\mathcal{L}(\hat{X} | p_\theta)$
 3. **Inference model**: $\mathcal{Q} = \{q_w(\theta) = q(\theta | w) \mid w \in \mathcal{W}\}$
 4. **Regularizing functional**: $\Omega(q)$ for $q \in \mathcal{Q}$.
 5. Solve **optimization** problem: $q_{\hat{X}} \in \arg \min_{q \in \mathcal{Q}} \left(\mathbb{E}_{q(\theta)} [\mathcal{L}(\hat{X} | p_\theta)] + \Omega(q) \right)$
 6. **Predictive distribution** via **model averaging**: $p(x | \hat{X}) := \int p(x | \theta) \cdot q_{\hat{X}}(\theta) d\theta$