

# Variational inference and VAEs

## CPSC 440/550: Advanced Machine Learning

`cs.ubc.ca/~dsuth/440/24w2`

University of British Columbia, on unceded Musqueam land

2024-25 Winter Term 2 (Jan–Apr 2025)

## Last time: approximate inference

- Bayesian inference requires computing **expectations with respect to posterior**,

$$\mathbb{E}[f(\theta)] = \int_{\theta} f(\theta) p(\theta | x) d\theta$$

- If  $f(\theta) = \theta$ , we get **posterior mean** of  $\theta$
- If  $f(\theta) = p(\tilde{x} | \theta)$ , we get **posterior predictive**
- If  $f(\theta) = \mathbb{1}(\theta \in S)$  we get probability of  $S$  (e.g., **marginals**)
  
- But posterior often **doesn't have a closed-form** expression
  - Bayesian linear regression –  $w \sim \mathcal{N}(m, V)$ ;  $y | x, w \sim \mathcal{N}(w^T x, \sigma^2)$  – does
  - **Bayesian logistic regression** –  $p(y | x, w) = 1/(1 + \exp(-y w^T x))$  – doesn't
  - More complex models almost never do
  
- Our two main tools for **approximate inference**:
  - 1 Monte Carlo methods
  - 2 Variational methods

# Approximate Inference

Two main strategies for **approximate inference**:

## ① Monte Carlo methods:

- Approximate expectations based on **samples**,

$$\mathbb{E}_{X \sim p} f(X) \approx \frac{1}{n} \sum_{i=1}^n f(x^{(i)})$$

- Turns **inference into sampling**
- **Simple Monte Carlo**: exactly as above, if we can take iid samples from  $p$
- **Rejection sampling**: get  $p$  samples from  $q$  samples and  $M \geq \max_x \tilde{p}(x)/q(x)$
- **Importance sampling**: estimate  $p$  expectations based on reweighting  $q$  samples
- **Markov chain Monte Carlo**: a little later in the course

## ② Variational methods:

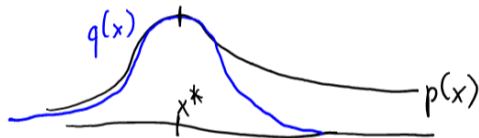
- Approximate  $p$  with “closest” **distribution  $q$  from a tractable family**,

$$\mathbb{E}_{X \sim p} f(X) \approx \mathbb{E}_{X \sim q} f(X)$$

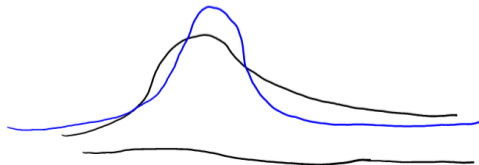
- $q$  could be Gaussian, product of Bernoulli, any other model with easy inference. . .
- Turns **inference into optimization**

## Variational Inference Illustration

- Example: approximate a non-Gaussian  $p$  by a Gaussian  $q$ 
  - Theoretical justification that most posteriors are “eventually” Gaussian
- Laplace approximation: find the mode  $x^*$ , then match first two derivatives of log-likelihood at the mode:  $\mathcal{N}(x^*, [\nabla^2 \log p(x^*)]^{-1})$ 
  - Still works with an unnormalized  $\tilde{p}(x)/Z = p(x)$ :  
 $\log p(x) = \log \tilde{p}(x) - \log Z$  has same mode and derivatives as  $\log p$



- Is this the “best” Gaussian approximation? What if we want non-Gaussian approx?



## Kullback-Leibler (KL) Divergence

- We'd like to find the “closest”  $q$  to our target  $p$
- How do we define “closeness” between a distribution  $p$  and  $q$ ?
- A common measure is **Kullback-Leibler (KL) divergence** between  $p$  and  $q$ :

$$\text{KL}(p \parallel q) = \mathbb{E}_{X \sim p} \log \frac{p(X)}{q(X)}$$

- Also called **information gain**: “information lost when  $p$  is approximated by  $q$ ”
- If  $p = q$ , we have  $\text{KL}(p \parallel q) = 0$  (no information lost)
- Otherwise,  $\text{KL}(p \parallel q)$  **grows as it becomes hard to predict  $p$  from  $q$**
- KL is **not symmetric**: in general,  $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$
- Maximizing likelihood = minimizing  $\text{KL}(p_{\text{true}} \parallel p_{\theta})$  (**bonus slide**)
- Unfortunately, computing this **requires integrating over or sampling from  $p$** 
  - ... exactly the problem we're trying to avoid

## Minimizing Reverse KL Divergence

- Most variational methods minimize “reverse KL”, as showed up with the ELBO:

$$\text{KL}(q \parallel p) = \mathbb{E}_{X \sim q} \log \frac{q(X)}{p(X)} = \mathbb{E}_{X \sim q} \log \left( \frac{q(x)}{\tilde{p}(x)} Z \right)$$

- Not very intuitive: “how much information is lost when we approximate  $q$  by  $p$ ”
- Does give some guarantee on approximating bounded functions (**bonus**)
- “Reverse” KL **only needs unnormalized distribution  $\tilde{p}/Z = p$  and expectations in  $q$**

$$\text{KL}(q \parallel p) = \mathbb{E}_{X \sim q} [\log q(X)] - \mathbb{E}_{X \sim q} [\log \tilde{p}(X)] + \underbrace{\mathbb{E}_{X \sim q} [\log(Z)]}_{\text{const. in } q}$$

- $-\mathbb{E}_{x \sim q} \log q(x) = \text{Entropy}[q]$  is the (differential) **entropy** of  $q$
- Value is known for many common choices of  $q$

$$\arg \min_q \text{KL}(q \parallel p) = \arg \max_q \text{Entropy}[q] + \mathbb{E}_{x \sim q} \log \tilde{p}(x)$$

## Example: Best Multivariate Gaussian

- We want to find  $\max_q \text{Entropy}[q] + \mathbb{E}_{x \sim q}[\log \tilde{p}(x)]$
- For multivariate Gaussians, we have  $\text{Entropy}[q] = \frac{1}{2} \log |\boldsymbol{\Sigma}| + \frac{d}{2} \log(2\pi e)$
- So to find the **best multivariate Gaussian approximation**, we need to find

$$\arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \frac{1}{2} \log |\boldsymbol{\Sigma}| + \mathbb{E}_{x \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} \log \tilde{p}(x) = \arg \max_{\boldsymbol{\mu}, \mathbf{L}} \log |\mathbf{L}| + \mathbb{E}_{z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \log \tilde{p}(\boldsymbol{\mu} + \mathbf{L}z)$$

- How to optimize this? Can't autodiff through expectation. . .
- **Reparameterization trick**: take relevant variable out of the expectation
  - Use Leibniz rule  $\frac{\partial}{\partial a} \mathbb{E}_{x \sim p} f(a, x) = \mathbb{E}_{x \sim p} \frac{\partial}{\partial a} f(a, x)$  **when  $p$  doesn't depend on  $a$**
  - Change variables to  $q = \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^\top)$ ; use  $|\mathbf{L}\mathbf{L}^\top| = |\mathbf{L}||\mathbf{L}^\top| = |\mathbf{L}|^2$
  - If  $L$  is **lower-triangular** with  $L_{jj} > 0$  (Cholesky factor), then  $|L| = \prod_j L_{jj}$  is **easy**
- Can take samples for  $z$  and run SGD to optimize (but note it's **non-convex**)

## Reparameterization trick

- Another view on why we can't autodiff through the expectation:

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} f(x) = \nabla_{\theta} \int f(x) p_{\theta}(x) dx = \int f(x) \nabla_{\theta} p_{\theta}(x) dx$$

and what do we do with that? (well, see **bonus slide**)

- But if we write  $x = g(\varepsilon, \theta)$  for  $\varepsilon \sim r$  (standard normal, uniform, ...),

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} f(x) &= \nabla_{\theta} \mathbb{E}_{\varepsilon} f(g(\varepsilon, \theta)) = \nabla_{\theta} \int f(g(\varepsilon, \theta)) r(\varepsilon) d\varepsilon \\ &= \int \nabla_{\theta} [f(g(\varepsilon, \theta))] r(\varepsilon) d\varepsilon = \mathbb{E}_{\varepsilon} [\nabla_{\theta} f(g(\varepsilon, \theta))] \end{aligned}$$

which is autodiff-friendly if we take Monte Carlo samples for  $\varepsilon$

- Need  $g$  to be differentiable (i.e.  $x$  should be continuous)
  - Tricks to avoid this; “Gumbel-Softmax” = “Concrete” distribution



- Another common scheme is **coordinate optimization** with an appropriate  $q$
- Consider choosing  $q$  as **a product of independent  $q_j$**

$$q(x) = \prod_{j=1}^d q_j(x_j)$$

- If we fix  $q_{\neg j}$  and optimize  $q_j$  among all distributions, we get (see PML2 10.2)

$$q_j(x_j) \propto \exp \left( \mathbb{E}_{q_{\neg j}} [\log \tilde{p}(x)] \right)$$

- Iterative algorithm: pick  $j$ , choose (discrete or conjugate)  $q_j$  to match above
  - Each iteration improves the (non-convex) reverse KL

# Outline

- 1 Variational inference
- 2 Variational Auto-Encoders**

## Deep latent variable model

- So far, we've built generative models out of relatively simple parts
  - Gaussian mixture is  $Z \sim \text{Cat}(\pi)$ ,  $X | (Z = z) \sim \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$
  - Can maximize  $p(x) = \sum_{z=1}^k \pi_z \mathcal{N}(x; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$  with EM/GD/...
- Discriminative models allow using arbitrary functions (e.g. **deep nets**) inside

$$Y | (X = x) \sim \mathcal{N}(g_\theta(z), \sigma^2)$$

- Can maximize  $p(y | x) = \mathcal{N}(x; f_\theta(x), \sigma^2)$  with GD/...
- Can we get a generative model with an arbitrary (deep) function in it?
- An example **deep latent variable model** ( $X$  is  $d$ -dimensional,  $Z$  is  $k$ -dimensional):

$$Z \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k) \quad X | (Z = z) \sim \mathcal{N}(g_\theta(z), \sigma^2 \mathbf{I}_d)$$

- Want to maximize  $p(x) = \int \mathcal{N}(z; \mathbf{0}_k, \mathbf{I}_k) \mathcal{N}(x; g_\theta(z), \sigma^2 \mathbf{I}_d) dz$
- How?

## Deep latent variable model

- We'd like to do MLE/similar for

$$Z \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k) \quad X | (Z = z) \sim \mathcal{N}(g_\theta(z), \sigma^2 \mathbf{I}_d)$$

which is

$$\max_{\theta} \sum_i \log \mathbb{E}_{z^{(i)} \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k)} \left[ \mathcal{N}(x^{(i)}; g_\theta(z^{(i)}), \sigma^2 \mathbf{I}) \right]$$

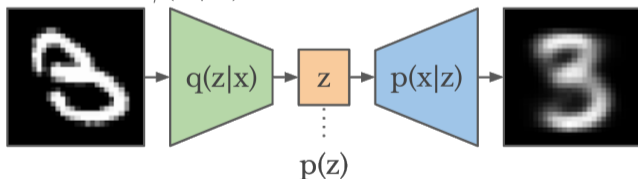
- Could potentially approximate this integral with Monte Carlo + reparam trick:

$$\max_{\theta} \sum_i \log \left( \frac{1}{M} \sum_{j=1}^M \mathcal{N}(x^{(i)}; g_\theta(z^{(i,j)}), \sigma^2 \mathbf{I}) \right) \quad \text{for } z^{(i,j)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- But this converges **really slowly** when  $k$  is large /  $g_\theta$  is complicated: need a huge  $M$

## Amortized inference

- EM would alternate between
  - Expectation:  $q(Z | X, \Theta) = p(Z | X, \Theta)$ , compute  $\mathbb{E}_{Z \sim q} \log p(X, Z | \Theta)$
  - Maximization of  $\mathbb{E}_{Z \sim q} \log p(X, Z | \Theta)$  in  $\Theta$
- The E step (inferring  $Z$  given  $X$ ) is **hard** here!
  - (So is the M step, in the normal deep network way)
- Have some complicated function from  $X$  to  $Z$
- Idea: instead of **exactly** solving the inference problem, let's **approximate** it
- ... with a neural network  $q_\phi(z | x)$



<https://danijar.com/building-variational-auto-encoders-in-tensorflow/>

- Named **variational autoencoder (VAE)**:  
encode image into latent code  $z$ , decode back to approximation of original image

# ELBO

- We'd like to maximize  $p_\theta(x) = \int p_\theta(x | z)p_\theta(z)dz$

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

- Since  $\text{KL} \geq 0$ ,  $\text{ELBO}_{\theta, \phi}(x) = \log p_\theta(x) - \text{KL}(q_\phi(z | x) \parallel p_\theta(z | x)) \leq \log p_\theta(x)$ 
  - ELBO is the **Evidence Lower Bound**
  - Same as we used in EM, except that we've separated it per-sample here

## Maximizing the ELBO

- Once we know how to evaluate it, we can use as our loss

$$\sum_{i=1}^n \text{ELBO}_{\theta, \phi}(x^{(i)}) = \sum_{i=1}^n \log p_{\theta}(x^{(i)}) - \text{KL}(q_{\phi}(z^{(i)} | x^{(i)}) || p_{\theta}(z^{(i)} | x^{(i)}))$$

- Because  $\text{KL} \geq 0$ , this is a **lower bound** on the log-likelihood
- Maximizing over the encoder/recognition parameters  $\phi$  is

$$\arg \max_{\phi} \sum_{i=1}^n \text{ELBO}_{\theta, \phi}(x^{(i)}) = \arg \min_{\phi} \sum_{i=1}^n \text{KL}(q_{\phi}(z^{(i)} | x^{(i)}) || p_{\theta}(z^{(i)} | x^{(i)}))$$

- Finds a network that gives you a **low reverse KL**, for any training input  $x^{(i)}$
- Making the inference network better **makes the likelihood bound tighter**
- If  $q_{\phi}(z | x) \approx p_{\theta}(z | x)$  (on the training set), maximizing over the probability parameters  $\theta$  **(approximately) maximizes likelihood**

## Evaluating the ELBO

- To efficiently evaluate the ELBO here:

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

- First term:  $q_{\phi}(z|x)$  should give a latent distribution where decoding to  $x$  is likely
- Second term:  $q_{\phi}(z|x)$  should be “near”  $p_{\theta}(z)$  (**regularization**)



## Computing the ELBO and its gradient: the reparameterization trick

- We want to maximize the average of

$$\text{ELBO}_{\theta, \phi}(x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x | z)] - \text{KL}(q_{\phi}(z | x) \parallel p(z))$$

- KL term for a given  $x$  is available **in closed form** if  $p(z)$ ,  $q_{\phi}(z | x)$  are Gaussian (if  $p(z)$  is  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $q_{\phi}(z | x)$  is  $\mathcal{N}(\boldsymbol{\mu}_{\phi}(x), \boldsymbol{\Sigma}_{\phi}(x))$ ; regularizes  $\|\boldsymbol{\mu}_{\phi}(x)\|^2$  and  $\boldsymbol{\Sigma}_{\phi}(x)$  to be near  $\mathbf{I}$  – **bonus**)

- For the other term, we need Monte Carlo

- Usually  $p_{\theta}(x | z)$  is  $\mathcal{N}(f_{\theta}(z), \sigma^2 \mathbf{I})$ , so  $\log p_{\theta}(x | z) = -\frac{1}{2\sigma^2} \|x - f_{\theta}(z)\|^2 + \text{const}$

- We need  $\mathbb{E}_{z \sim q_{\phi}(z|x)} \log p_{\theta}(x | z)$

- Estimate with Monte Carlo

- Can usually use just a single step for simplicity: if  $q_{\phi}$  is “good,” should be okay

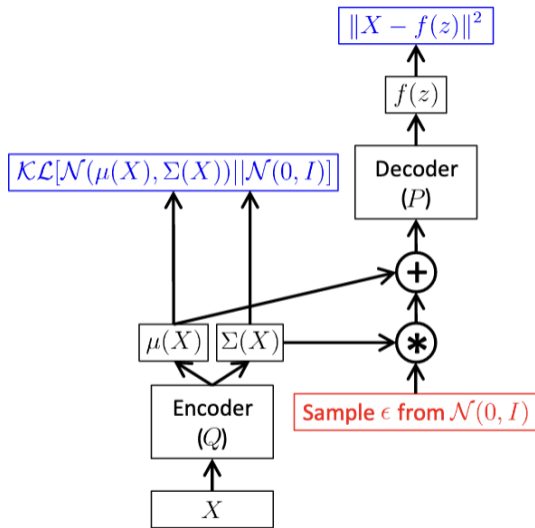
- But how do we take  $\nabla_{\phi}$  of this expectation? Use **reparameterization trick** again:

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x | z)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \log p_{\theta}(x | z = \boldsymbol{\mu}_{\phi}(x) + \boldsymbol{\Sigma}_{\phi}(x)^{\frac{1}{2}} \epsilon)$$

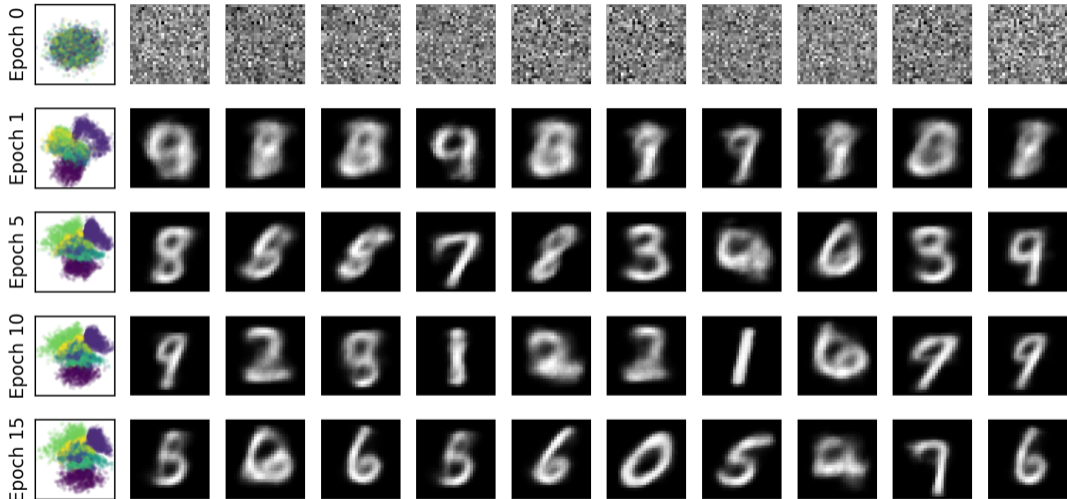
- Take a Monte Carlo sample for  $\epsilon$ ; now have something we can autodiff

- Now just do SGD to maximize  $\frac{1}{n} \sum_{i=1}^n \widehat{\text{ELBO}}_{\theta, \phi}(x^{(i)})$

# A VAE



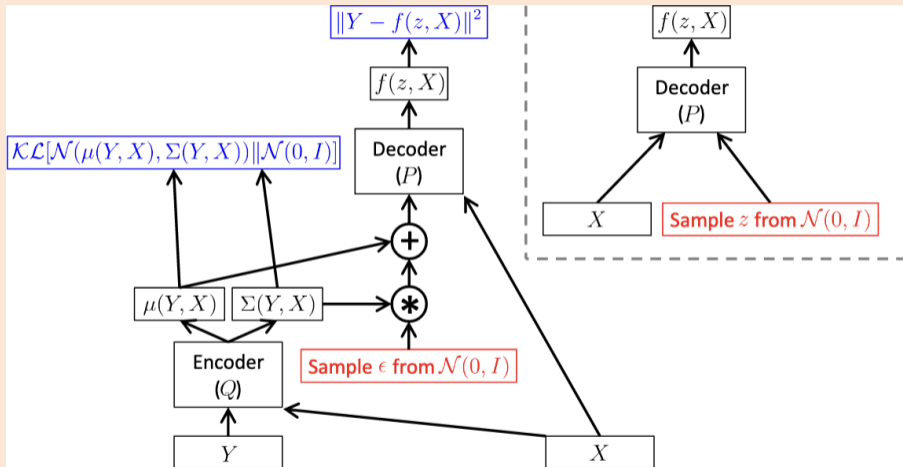
# A VAE on MNIST



<https://danijar.com/building-variational-auto-encoders-in-tensorflow/>

# Conditional VAE

bonus!



<https://arxiv.org/pdf/1606.05908.pdf>

# Conditional VAE to “in-paint” on MNIST

bonus!

ground-truth	7	3	6	2	3	5	0	0	5	6	2	6	2	3	5	4	1	0	4
NN	7	3	6	2	3	3	0	0	5	2	2	6	2	3	5	9	1	0	4
CVAE	7	3	6	2	3	3	0	0	5	2	2	6	2	3	5	4	1	0	4
	7	3	6	2	3	3	0	0	5	4	2	6	2	3	5	4	1	0	4
	7	3	6	2	5	3	0	0	5	2	2	6	2	3	5	4	1	0	4
	7	3	6	2	3	3	0	0	5	4	2	6	2	3	5	4	1	0	4
	7	3	6	2	5	3	0	0	5	4	2	6	2	3	5	4	1	0	4
	7	3	6	2	5	3	0	0	5	6	2	6	2	3	5	4	1	0	4
	7	5	6	2	5	3	0	0	3	4	2	6	2	3	5	4	1	0	4

[https://papers.nips.cc/paper\\_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf](https://papers.nips.cc/paper_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf)

# Summary

- Variational inference: choose  $q \approx p$  and estimate  $\mathbb{E}_{X \sim p} f(X) \approx \mathbb{E}_{X \sim q} f(X)$ 
  - Reparameterization trick often helps find best  $q$  with SGD
  - Not consistent ( $q$  doesn't converge to  $p$  if we run the algorithm forever)
  - Can be complicated and hard to parallelize
  - Only needs unnormalized density
  - Often better approximation for a given amount of computation than MCMC
- Variational auto-encoders (VAEs) are a deep latent variable model
  - $p(x) = \int p(z)p(x | z)dz$
  - Learn an “inference” / “recognition” network  $q(z | x)$  to approximate inference
  - Minimizing the ELBO maximizes a lower bound on the likelihood
- Next up: how to design a VAE and how to use it

## Maximum likelihood minimizes KL

bonus!

$$\begin{aligned}\arg \min_{\theta} \text{KL}(p_{\text{true}} \parallel p_{\theta}) &= \arg \min_{\theta} \int p_{\text{true}}(x) \log \frac{p_{\text{true}}(x)}{p_{\theta}(x)} dx \\ &= \arg \min_{\theta} \underbrace{\int p_{\text{true}}(x) \log p_{\text{true}}(x) dx}_{\text{doesn't depend on } \theta} - \int p_{\text{true}}(x) \log p_{\theta}(x) dx \\ &= \arg \min_{\theta} - \int p_{\text{true}}(x) \log p_{\theta}(x) dx \\ &= \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{true}}} \log p_{\theta}(x) \\ &\approx \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(x^{(i)})\end{aligned}$$

- Alternative gradient estimator to the reparameterization trick:

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} f(x) = \nabla_{\theta} \int f(x) p_{\theta}(x) dx = \int f(x) \nabla_{\theta} p_{\theta}(x) dx$$

- Notice that  $\nabla_{\theta} \log p_{\theta}(x) = \frac{1}{p_{\theta}(x)} \nabla_{\theta} p_{\theta}(x)$ , so

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} f(x) = \int f(x) \nabla_{\theta} \log p_{\theta}(x) p_{\theta}(x) dx = \mathbb{E}_{x \sim p_{\theta}} [f(x) \nabla_{\theta} \log p_{\theta}(x)]$$

- So if we can evaluate  $\nabla_{\theta} \log p_{\theta}(x)$  for any  $x$ , we can estimate the gradient!
- Called REINFORCE estimator:
  - Doesn't require  $x$  to be continuous, e.g. it can be categorical with parameters a differentiable function of  $\theta$
  - Unbiased estimator
  - Tends to have large variance



- Does  $\text{KL}(q \parallel p)$  being small tell us anything about  $|\mathbb{E}_{X \sim p} f(X) - \mathbb{E}_{X \sim q} f(X)|$ ?
- For a **bounded**  $f$  with  $|f(x)| \leq F$  for all  $x$ , we have for any  $p$  and  $q$  that

$$\begin{aligned} \left| \mathbb{E}_{X \sim p} f(X) - \mathbb{E}_{X \sim q} f(X) \right| &= 2F \frac{1}{2} \left| \mathbb{E}_{X \sim p} \frac{f(X)}{F} - \mathbb{E}_{X \sim q} \frac{f(X)}{F} \right| \\ &\leq 2F \text{TV}(p, q) \\ &\leq F \sqrt{2 \text{KL}(q \parallel p)} \end{aligned}$$

- Here we used the **total variation** distance, which is (supremum is a fancy version of max)

$$\text{TV}(p, q) = \frac{1}{2} \sup_{f: \forall x, |f(x)| \leq 1} \left| \mathbb{E}_{X \sim p} f(X) - \mathbb{E}_{Y \sim q} f(Y) \right|$$

- **Pinsker's inequality** says that  $\text{TV}(p, q) \leq \sqrt{\frac{1}{2} \text{KL}(p \parallel q)}$ ; note TV is symmetric
  - [arxiv.org/abs/2202.07198](https://arxiv.org/abs/2202.07198) gives a nice overview of when this is tight
- Not a super tight bound, but does give some reassurance
- Can min. another “integral probability metric” for better bound, but usually harder

- We want to maximize the average of

$$\text{ELBO}_{\theta, \phi}(x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x | z)] - \text{KL}(q_{\phi}(z | x) \parallel p(z))$$

- KL term for a given  $x$  is often available **in closed form**
- Typically we choose  $p_{\theta}(z)$  to be  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $q_{\phi}(z | x)$  to be  $\mathcal{N}(\boldsymbol{\mu}_{\phi}(x), \boldsymbol{\Sigma}_{\phi}(x))$
- Then the KL is just (see PML2 eq 5.80)

$$\frac{1}{2} (\|\boldsymbol{\mu}_{\phi}(x)\|^2 + \text{Tr} \boldsymbol{\Sigma}_{\phi}(x) - \log |\boldsymbol{\Sigma}_{\phi}(x)| - d)$$

- Most of the time we also choose  $\boldsymbol{\Sigma}_{\phi}(x)$  to be diagonal; determinant is easy
- This is just an expression in terms of  $\phi$ ; we can use autodiff