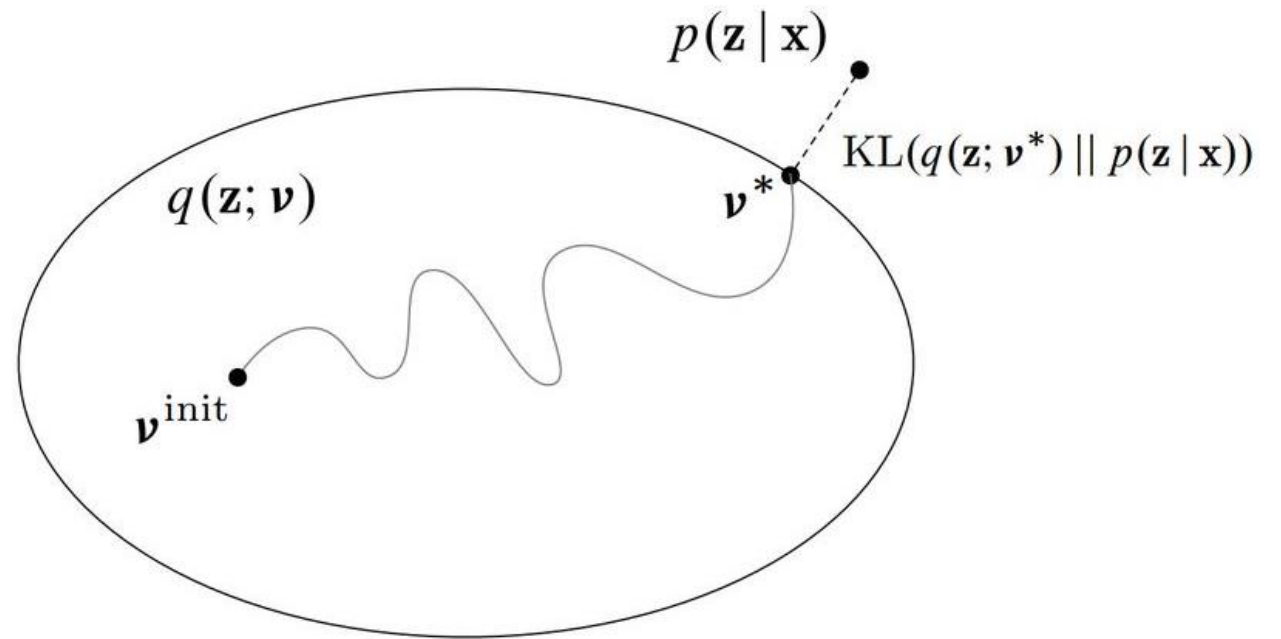


Variational inference



Marginal likelihood

- Assume you have statistics $p(x, z)$ of eye colors x per nationality z

	Dutch	Greek	Chinese	Indian	Italian	German	US	Spanish
Brown	0.02	0.03	0.02	0.01	0.07	0.03	0.01	0.00
Blue	0.09	0.01	0.03	0.03	0.08	0.04	0.03	0.06
Green	0.01	0.01	0.08	0.06	0.07	0.08	0.06	0.06

- If we want to know the distribution of one of our variables (*e.g.*, eye colors)
 - → we sum up over all possible outcomes (marginalize) of the other variable
 - *E.g.*, nationalities for the marginal likelihood $p(x) = \sum_z p(x, z)$

Color	
Brown	0.19
Blue	0.37
Green	0.44
Total	1.00

Marginal likelihood

- Or assume that our bottom half pixels are visible (\mathbf{x}) and the upper half not (\mathbf{z})



- Assume we somehow know a good model $p(\mathbf{x}, \mathbf{z})$ of how the visible and hidden/latent pixels interact
- Let's say we want to know how likely bottom half the image is to be observed
- We must marginalize out all possible latents \mathbf{z} to fill the rest of the image
 - For instance, some pictures might contain one, two, or more elephants

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$$

Marginal likelihood in latent variable models

- When “learning to represent” an input \mathbf{x} we assume a latent variable \mathbf{z}
 - and try to explain \mathbf{x} using all possible \mathbf{z}

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

- Hence, a latent variable model can be viewed as a generation process
 - First, we generate a new \mathbf{z} from $p(\mathbf{z})$ by sampling
 - Then, we generate a new \mathbf{x} by sampling from the $p(\mathbf{x}|\mathbf{z})$ given the sampled \mathbf{z}



Intractable $p_{\theta}(\mathbf{x}, \mathbf{z})$

- Question: how to find the optimal parameters θ ?
- Maximizing log-likelihood, again

$$\log \prod_{\mathbf{x} \in D} p(\mathbf{x}) = \sum_{\mathbf{x}} \log p(\mathbf{x}) = \sum_{\mathbf{x}} \log \sum_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z})$$

- Like in Boltzmann machines, the $\sum_{\mathbf{z}} \dots$ is a nasty one
 - E.g., for a 3-dimensional binary \mathbf{z} iterate over $[0,0,0], [0,0,1], [0,1,1], \dots$
 - For 20 dimensions $2^{20} \approx 1M$ latents and generations. Per image \mathbf{x} !
 - For continuous \mathbf{z} even harder, we cannot even enumerate

Making $p_{\theta}(\mathbf{x}, \mathbf{z})$ tractable with naive Monte Carlo

- We want to optimize per data point $\mathbf{x} : \log \sum_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z})$
- The sum contains a bunch of probabilities
 - equivalent to **expected value** times the **number of summands**

$$\log \sum_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) = \log |Z| \mathbb{E}[p_{\theta}(\mathbf{x}, \mathbf{z})]$$

- Do we need all the summands to compute the expected value (average)
 - No, if we sample randomly \mathbf{z} (uniformly) and average, it gives us an estimate
 - Basically replace whole sum with a weighted smaller sum

$$\log |Z| \mathbb{E}[p_{\theta}(\mathbf{x}, \mathbf{z})] \approx \log \frac{|Z|}{K} \mathbb{E}_{\mathbf{z} \sim \text{Uniform}}^{(K)} [p_{\theta}(\mathbf{x}, \mathbf{z})] = \log \frac{|Z|}{K} \sum_{k=1}^K p_{\theta}(\mathbf{x}, \mathbf{z}_k)$$

- Doesn't scale, too many samples for the expectation estimate to be accurate
 - Most \mathbf{z}_k would be in 'very low density regions' \rightarrow Unimportant $p_{\theta}(\mathbf{x}, \mathbf{z}_k)$
 - In technical terms, this is a 'high variance' estimator

Making $p_{\theta}(\mathbf{x}, \mathbf{z})$ tractable with importance sampling MC

- Better if select few good summands in the sum $\sum_{k=1}^K p_{\theta}(\mathbf{x}, \mathbf{z}_k)$
- If, theoretically, we had a nice distribution around the mass of relevant \mathbf{z}_k
 - we could **use that distribution** to sample \mathbf{z}_k and get a better sample average with fewer k

$$\begin{aligned}\log \sum_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) &= \log \sum_{\mathbf{z}} \mathbf{q}_{\varphi}(\mathbf{z}) \frac{1}{q_{\varphi}(\mathbf{z})} p_{\theta}(\mathbf{x}, \mathbf{z}) \\ &= \log \mathbb{E}_{\mathbf{z} \sim \mathbf{q}_{\varphi}(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z})} \right] \approx \log \frac{1}{K} \sum_{k=1}^K \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_k)}{q_{\varphi}(\mathbf{z}_k)}, \text{ where } \mathbf{z}_k \text{ are sampled from } \mathbf{q}_{\varphi}(\mathbf{z})\end{aligned}$$

- Note the dual use of $q_{\varphi}(\mathbf{z})$
 - In the nominator $\mathbf{q}_{\varphi}(\mathbf{z})$ is the density function we use as sampling mechanism. By sampling from it (e.g., Gaussian samples if it is Gaussian) this quantity is used and disappears by the sum
 - In the denominator $q_{\varphi}(\mathbf{z}_k)$ is simply a function. We feed it \mathbf{z}_k and returns how important \mathbf{z}_k is for our probability space
- Scales much better and with much lower variance, but we don't know what is a good $q_{\varphi}(\mathbf{z}_k)$

Learning the importance sampling distribution

- Importance sampling is promising but how to determine $q_\varphi(\mathbf{z}_k)$?
- Learn $q_\varphi(\mathbf{z}_k)$ from data!
- Our learning objective is to maximize the log probability

$$\log \mathbb{E}_{\mathbf{z} \sim q_\varphi(\mathbf{z})} \left[\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z})} \right] \approx \log \frac{1}{K} \sum_{k=1}^K \frac{p_\theta(\mathbf{x}, \mathbf{z}_k)}{q_\varphi(\mathbf{z}_k)}$$

- The $\log \mathbb{E}$ stands for logarithm of an unknown integral
 - not very convenient for derivations and computations
- Would be much nicer if we could swap the $\log \mathbb{E}$ to $\mathbb{E} \log$
 - Then we would simply need the expectation of the logarithm of a function
 - Especially convenient if $p_\theta(\mathbf{x}, \mathbf{z})$ belongs to the exponential family

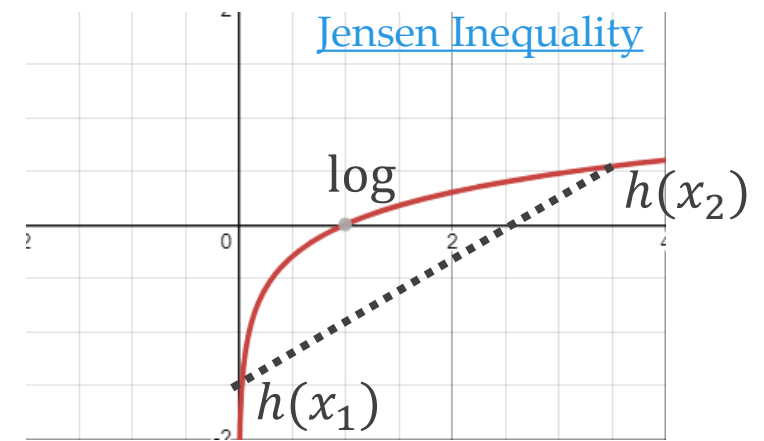
Jensen's inequality

- A concave function h (like a logarithm) on a sum will always be larger than the sum of h on individual summands
 - Basically, a line connecting two points of a function will be always below the function

$$h(tx_1 + (1 - t)x_2) \geq th(x_1) + (1 - t)h(x_2)$$

- With probabilities and random variables this translates to

$$h(\mathbb{E}[\mathbf{x}]) \geq \mathbb{E}[h(\mathbf{x})]$$



A lower bound on the maximum likelihood

- By applying Jensen's inequality

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

- We replaced the original ML objective with a quantity that is always smaller
 - (1) By improving $\mathbb{E}_{\mathbf{z} \sim q_{\varphi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z})} \right]$ we always improve $\log \mathbb{E}_{\mathbf{z} \sim q_{\varphi}(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z})} \right]$
 - 'Lower bound'
- (2) $\mathbb{E}_{\mathbf{z} \sim q_{\varphi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z})} \right]$ is a tractable & comfortable quantity → easy optimization
 - An expectation → Monte Carlo sampling is possible
 - The log can couple nicely with p_{θ} if chosen properly

Making $p(\mathbf{z}|\mathbf{x})$ tractable with variational inference

- We can also view variational inference from the lens of intractability
- The problematic quantity in our latent model is the posterior
 - The reason is the **intractable normalization**

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} = \frac{p(\mathbf{x}, \mathbf{z})}{\int \mathbf{p}(\mathbf{x}, \mathbf{z}) d\mathbf{z}}$$

- Variational inference approximates the true posterior $p(\mathbf{z}|\mathbf{x})$ with $q_\phi(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} \text{KL}(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) &= \int q_\phi(\mathbf{z}) \log \frac{q_\phi(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= - \int q_\phi(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})q_\phi(\mathbf{z})} d\mathbf{z} = - \int q_\phi(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z})} + \int q_\phi(\mathbf{z}) \log p(\mathbf{x}) d\mathbf{z} \\ &= -\mathbb{E}_{q_\phi(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z})} \right] + \log p(\mathbf{x}) \end{aligned}$$

Evidence lower bound

- $\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right]$ also known as ‘evidence lower bound’

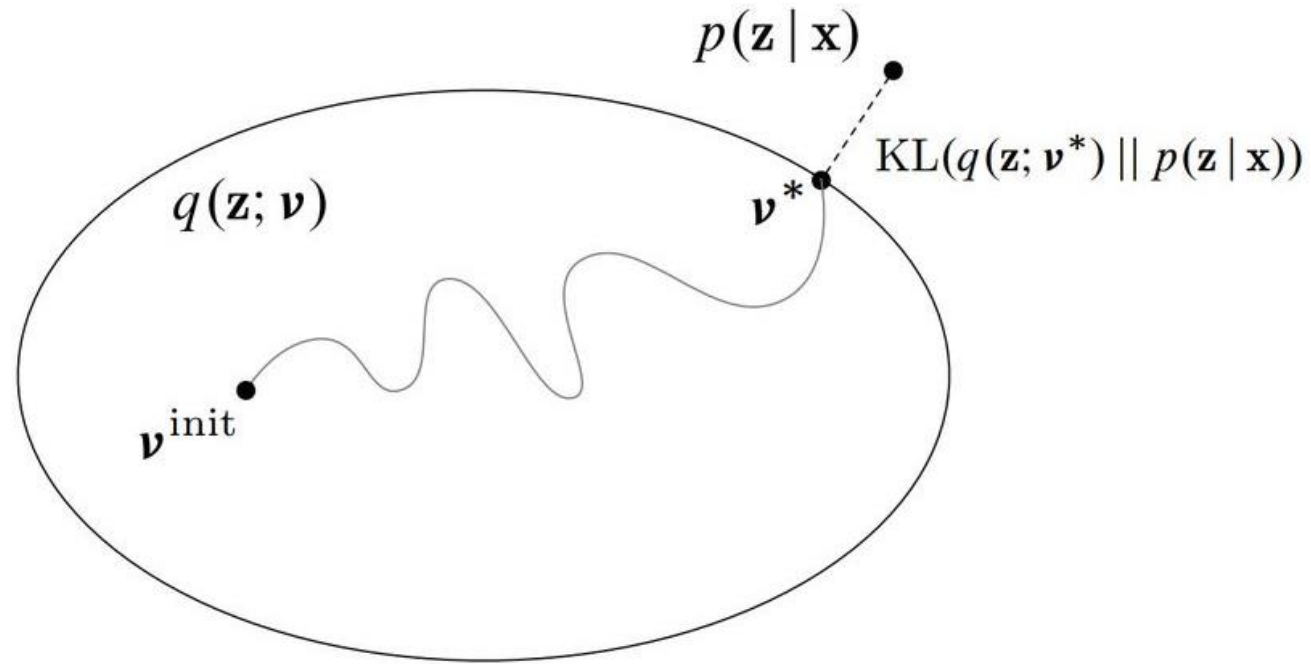
$$\begin{aligned} \log p(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right] + \text{KL}(q_{\phi}(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \\ &= \text{ELBO} + \text{KL}(q_{\phi}(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \end{aligned}$$

- Why ‘evidence’?
 - The KL term is always positive
 - If we drop it, we bound the log evidence $\log p(\mathbf{x})$ from below

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right]$$

- Higher ELBO \rightarrow smaller difference to true $p_{\theta}(\mathbf{z}|\mathbf{x}) \rightarrow$ better latent representation
- Higher ELBO \rightarrow gap to log-likelihood tightens \rightarrow better density model

Variational inference, graphically



ELBO balancing reconstruction and the prior

- We can expand the ELBO as

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

- The **first term** encourages the reconstructions that the maximize likelihood
- The **second term** minimizes the distance of the variational distribution from the prior

ELBO and entropy regularization

- We can also expand the ELBO as

$$\begin{aligned}\mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right] &= \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z})} \right] \\ &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z})} [\log q_{\phi}(\mathbf{z})] \\ &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}))\end{aligned}$$

where $H(\cdot)$ is the entropy

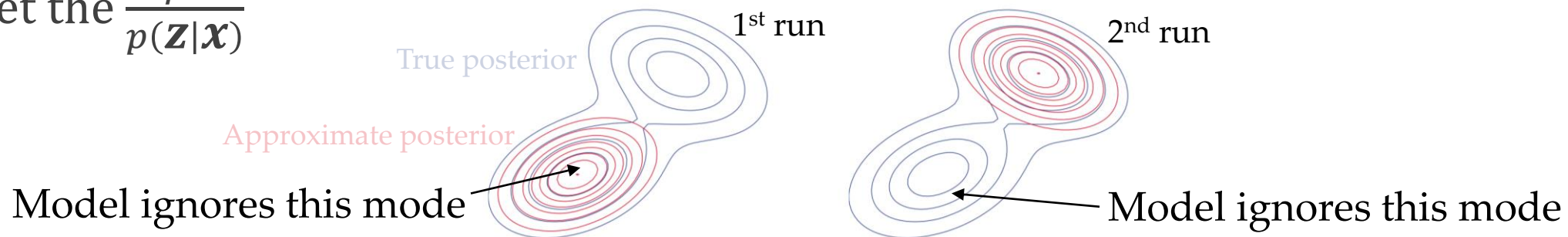
- Maximizing the **joint likelihood** → Something like the Boltzmann energy
- While **maintaining enough entropy** (‘uncertainty’) in the distribution of latents
 - Avoiding latents to collapse to pathological, point estimates (\mathbf{z} as single values)

Variational inference underestimates variance

- If you noticed, for the second way to derive the ELBO we minimized

$$\text{KL}(q_{\phi}(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) = \int q_{\phi}(\mathbf{z}) \log \frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

- We want to sample from $q_{\phi}(\mathbf{z})$ in expectations later on, as $p(\mathbf{z}|\mathbf{x})$ is intractable
- The model wants to approximate $p(\mathbf{z}|\mathbf{x}) \rightarrow$ can't really know where $p(\mathbf{z}|\mathbf{x})$ is low
- The model prefers to hedge and 'bias' $q_{\phi}(\mathbf{z})$ towards 0 for regions it can't be certain
 - Better pick one mode (randomly) than miss a 'zero' density region of $p(\mathbf{z}|\mathbf{x})$ and skyrocket the $\frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$



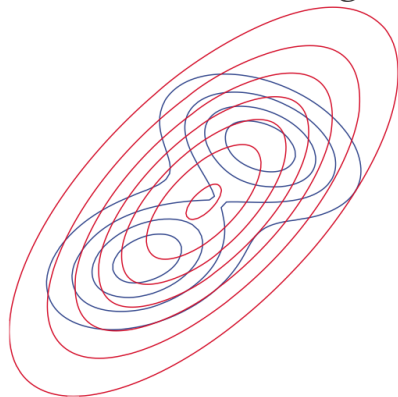
How to overestimate variance?

- You would need to use the forward KL

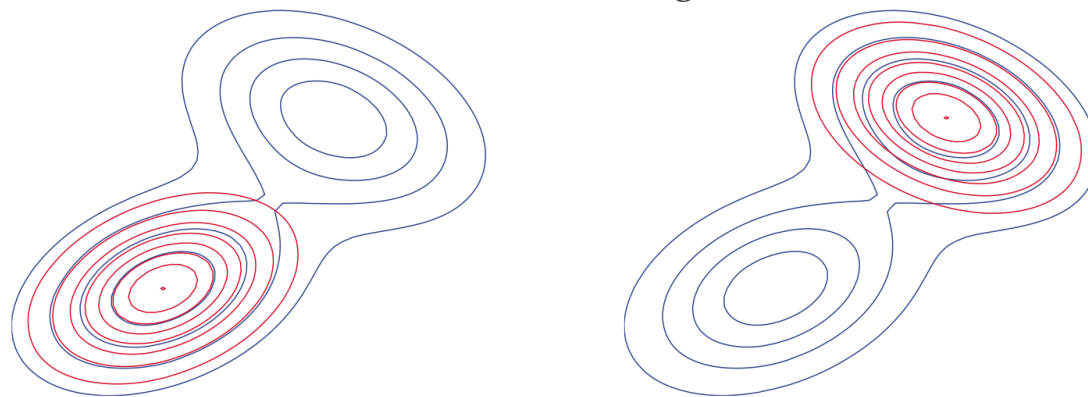
$$\text{KL}(p(\mathbf{z}|\mathbf{x}) \parallel q_{\boldsymbol{\varphi}}(\mathbf{z})) = \int p(\mathbf{z}|\mathbf{x}) \log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})}{q_{\boldsymbol{\varphi}}(\mathbf{z})} d\mathbf{z}$$

- The model would prefer placing some density everywhere
 - That way it avoids $\frac{p(\mathbf{z}|\mathbf{x})}{q_{\boldsymbol{\varphi}}(\mathbf{z})}$ skyrocketing if it misses areas where $p(\mathbf{z}|\mathbf{x})$

Overestimating variance



Underestimating variance



Variational Inference and Variational Autoencoders

- Variational Inference is a general machine learning methodology
 - Not specific to Variational Autoencoders
- In the original Variational Inference, we have $q_{\phi}(\mathbf{z}) = q(\mathbf{z}; \boldsymbol{\phi})$ where $\boldsymbol{\phi} = \{\boldsymbol{\phi}_i\}$ per data point
 - Not $q_{\phi}(\mathbf{z}|\mathbf{x})$, this is specific in Variational Autoencoders
 - That is, we optimize a different neural network $\boldsymbol{\phi}_i$ per data point
 - This is inefficient and cannot generalize ($\boldsymbol{\phi}_i$ are only for the training)
- As we will see, Variational Autoencoders share the parameters $\boldsymbol{\phi}$ among all data points (single neural network in the end)
 - This is called amortization
 - And introduce specialization by conditioning on the data point $q_{\phi}(\mathbf{z}|\mathbf{x})$