

# Why approximate inference?



# Analytical inference

$$p^*(z) = p(z|X) = \frac{p(X|z)p(z)}{p(X)}$$

- Easy for conjugate priors
  - Hard otherwise

Example:  $p(x|z) = \mathcal{N}(x|\mu(z), \sigma^2(z))$

sometimes, computing the posterior is intractable like this.

We want to APPROXIMATE the posterior so that it is easier for us, and that the approximation is acceptable.



Neural networks

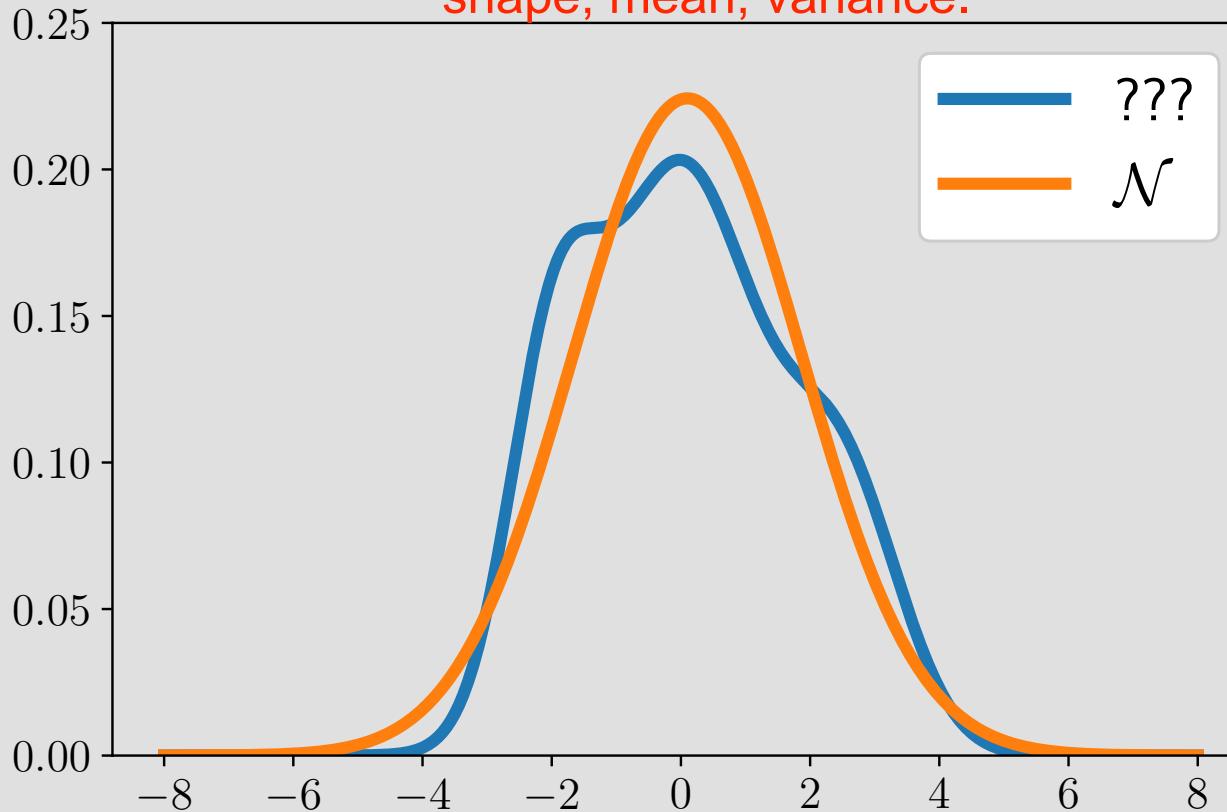
the likelihood here is a neural network. Here, there is obviously no conjugacy.



# Do we need exact posterior?

not really right? An approximation is good enough?

See the gaussian approximation below. It is a really good approximation. Matches shape, mean, variance.



This week we'll be studying how to approximate the posterior.



## Variational inference



1. Select a family of distributions  $Q$

Example:  $\mathcal{N}(\mu, \begin{pmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ 0 & & \ddots & \\ & & & \sigma_d^2 \end{pmatrix})$

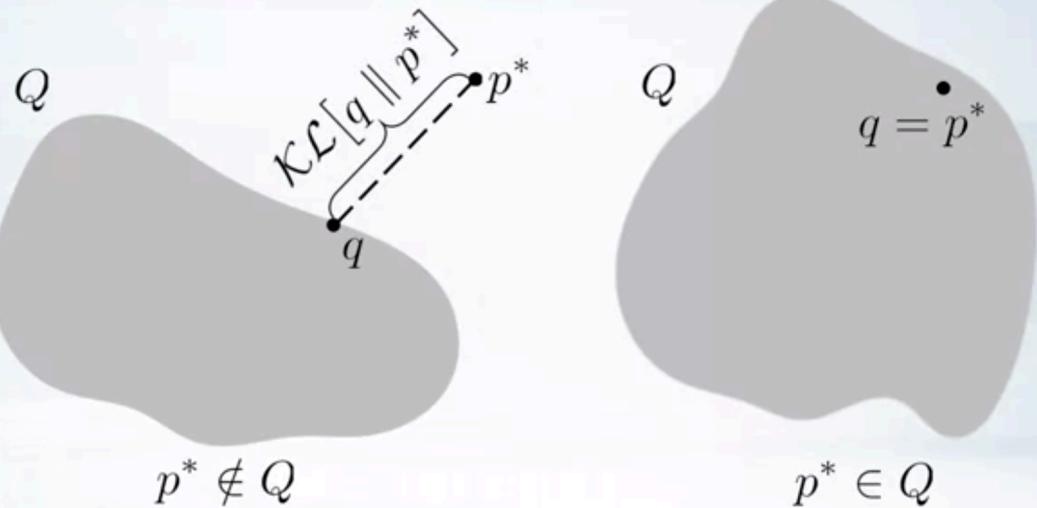
2. Find best approximation  $q(z)$  of  $p^*(z)$ :

$$\mathcal{KL}[q(z) \parallel p^*(z)] \rightarrow \min_{q \in Q}$$

Essentially, pick the best fitting distribution from a family  $Q$  (by KL)



## Choice of variational family



If  $Q$  is small, the true posterior will not lie in it.  
If  $Q$  is sufficiently large, we'll get a good posterior.  
If  $Q$  is larger, then it is harder to compute the variational inference. Search space too big, presumably?



## Unnormalized distribution

$$p^*(z) = p(z|X) = \frac{p(X|z)p(z)}{p(X)} = \frac{\hat{p}(z)}{Z}$$

This is hard at first because we need to compute  
the evidence  $p(X)$ .

This is really hard.

But we see a nice property of KL divergence:

$$\begin{aligned} \mathcal{KL}[q(z) \parallel \frac{\hat{p}(z)}{Z}] &= \int q(z) \log \frac{q(z)}{\hat{p}(z)/Z} dz \\ &= \int q(z) \log \frac{q(z)}{\hat{p}(z)} dz + \int q(z) \log Z dz \\ &= \mathcal{KL}[q(z) \parallel \hat{p}(z)] + \log Z \end{aligned}$$

$\mathcal{KL}[q(z) \parallel \hat{p}(z)] \rightarrow \min_z$

Let  $P(X|z) p(z)$  be  
 $\hat{p}$ .

We can set the evidence as  
a constant  $Z$ .

We take  $Z$  out of the KL  
divergence equation. Since  
we're computing KL w.r.t.  
parameters  $z$ , we can take  
this constant out of the  
equation.