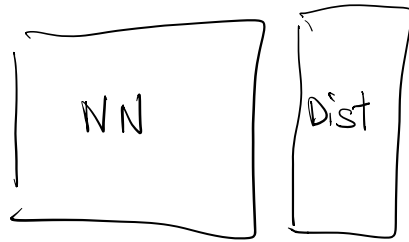


## Maximum likelihood



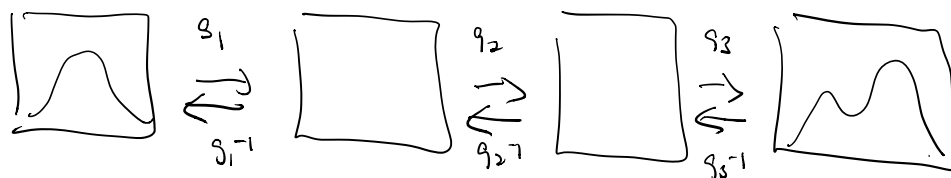
NN outputs parameter of chosen distribution

Loss function is negative log likelihood of distribution

E.g. for a classification model,

the distribution could be a multinomial.

## Normalizing flows



↓  
bijections  $\rightarrow$  i.e. reversible functions

$g = g_1 \circ g_2 \circ g_3$   
is also a bijection

↓  
these functions have  
parameters  $\theta$

↓  
 $\theta$  estimated using  
maximum likelihood

## Math Review

Transformations of RV

$$x = g(y)$$

$$f_y(y) = f_x(g(y)) \underbrace{g'(y)}$$

Replace  $\mathbb{Z}$  Jacobian  
determinant in  $\mathbb{R}^1$

## Variational Bayes

$$D_{KL}(P||Q) = \sum P(x) \log \frac{P(x)}{Q(x)} \approx \int P(x) \log \frac{P(x)}{Q(x)} dx$$

In VB, we have a family of distributions parameterized by  $\lambda$ ,  $Q_\lambda(\theta)$  that we want to use to approximate  $P(\theta|D)$

$$D_{KL}[Q_\lambda(\theta) || P(\theta|D)]$$

note flip!

$$= \int Q_\lambda(\theta) \log \frac{Q_\lambda(\theta)}{P(\theta|D)} d\theta$$

$$= \int Q_\lambda(\theta) \log \frac{Q_\lambda(\theta) P(D)}{P(D|\theta) P(\theta)} d\theta$$

$$= \int Q_\lambda(\theta) \log P(D) d\theta - \int Q_\lambda(\theta) \log \frac{P(D|\theta) P(\theta)}{Q_\lambda(\theta)} d\theta$$

$$= \log P(D) \int Q_\lambda(\theta) d\theta + \int Q_\lambda(\theta) \log \frac{Q_\lambda(\theta)}{P(\theta)} d\theta - \int Q_\lambda(\theta) \log P(D|\theta) d\theta$$

$$= \underbrace{\log P(D)}_{\text{no role in minimization}} + \underbrace{D_{KL}[Q_\lambda(\theta) || P(\theta)]}_{\text{regularizer}} - \underbrace{\mathbb{E}_{\theta \sim Q_\lambda} [\log P(D|\theta)]}_{\text{mean NLL when } \theta \text{ is sampled from } Q_\lambda}$$

custom loss  
func!

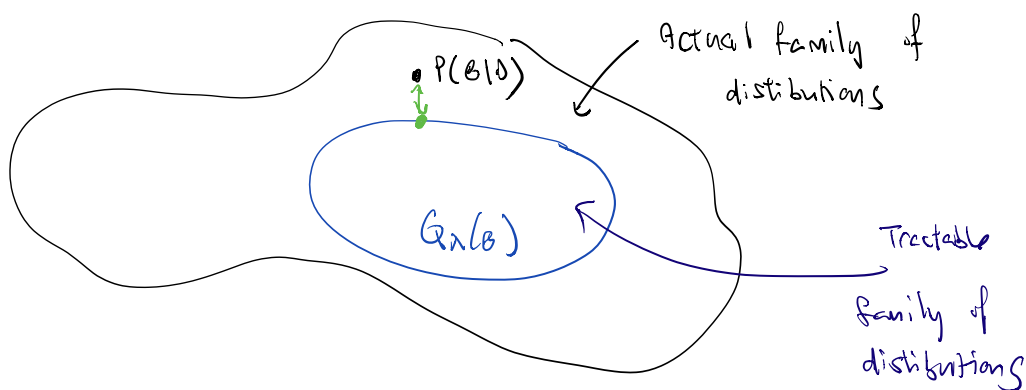
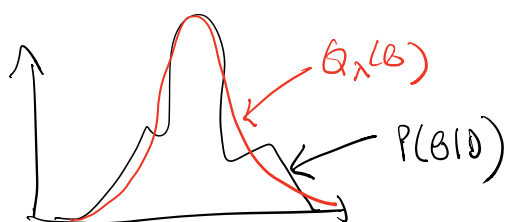
no role in  
minimization

regularizer

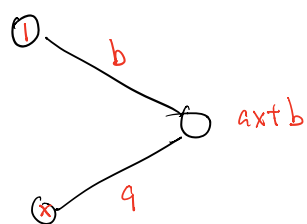
mean NLL when

$\theta$  is sampled from  $Q_\lambda$

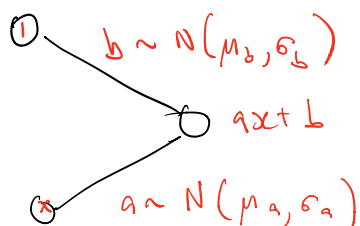
$P(\theta)$  chosen to be  
small & centered  
around 0



## Bayesian NN



Regular NN



Bayesian NN

## Variational parameters

note:  $\sigma_a, \sigma_b > 0$

use softmax

$$f(x) = \log(1 + e^x)$$

$$\lambda = \{\mu_a, \sigma_a, \mu_b, \sigma_b\}$$

$Q(\lambda) \sim$  Product of independent Gaussian distributions

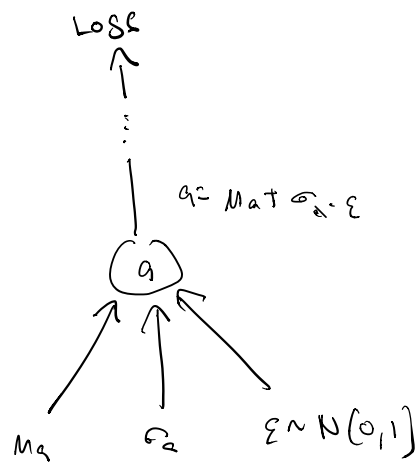
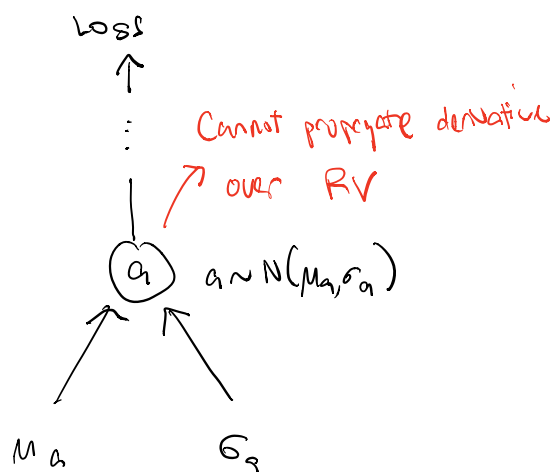
$$\text{Loss} = \underbrace{D_{KL} [Q_{\lambda}(\theta) \parallel P(\theta)]}_{\theta \sim Q_{\lambda}} - \underbrace{E \left[ \log P(D|\theta) \right]}_{\theta \sim Q_{\lambda}}$$

$$P(\theta) \sim N(0, 1)$$

$$Q_{\lambda}(\theta) \sim N(\mu, \sigma)$$

$$-\frac{1}{2} (1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

single sample of  $\theta$   
is used in practice!



MC dropout

[illegible]

In MC depend, turn ON during training AND test

During test  $\rightarrow$  evaluate each sample multiple times

Averaging the outputs gives the Bayesian predictive distribution.



VB use  $\mu, \sigma$  to  
model distribution  
of  $w$

