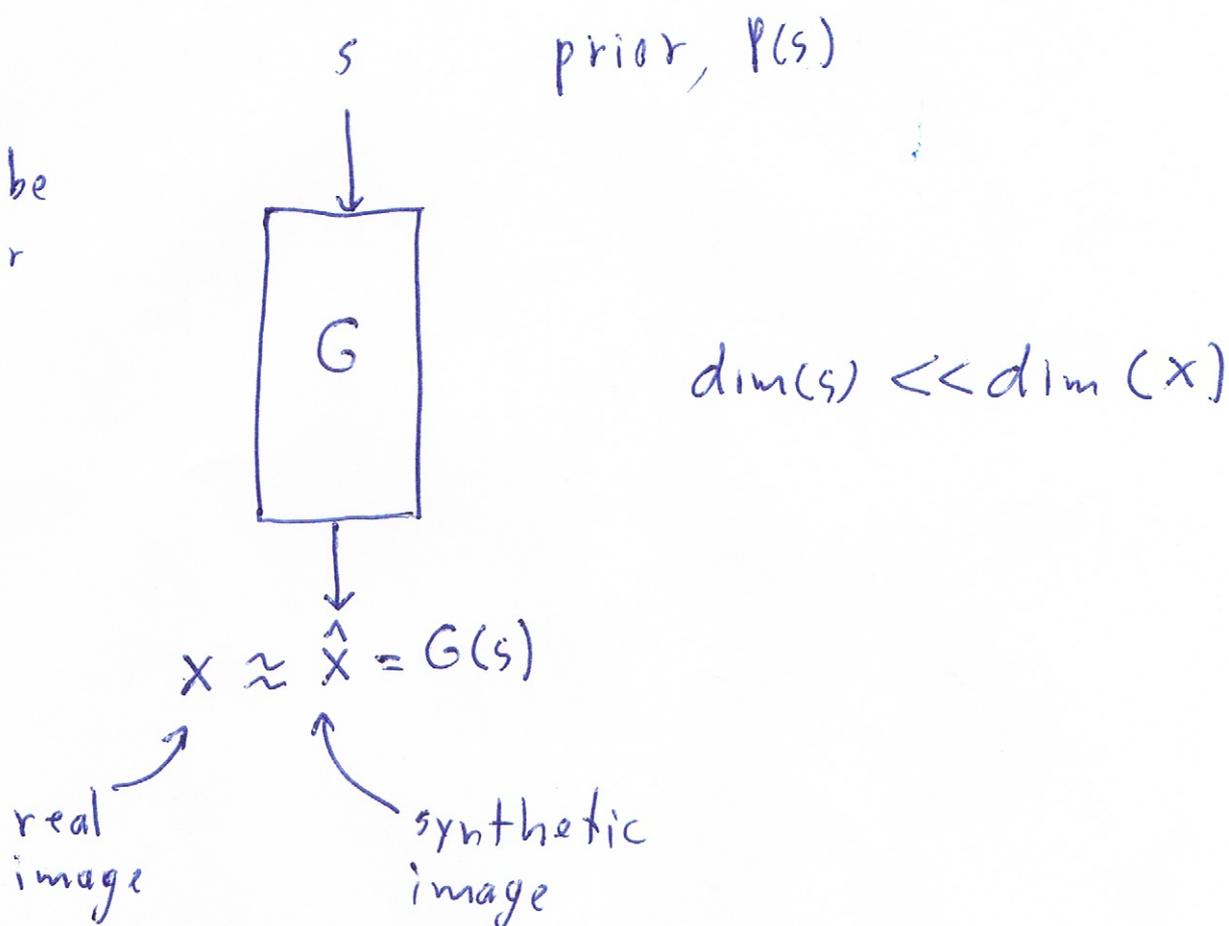


# Generative Models

①

G may be given or learned.



Pretend that we can randomly generate images "near"  $\hat{x}$  with probability  $P(x|\hat{x})$

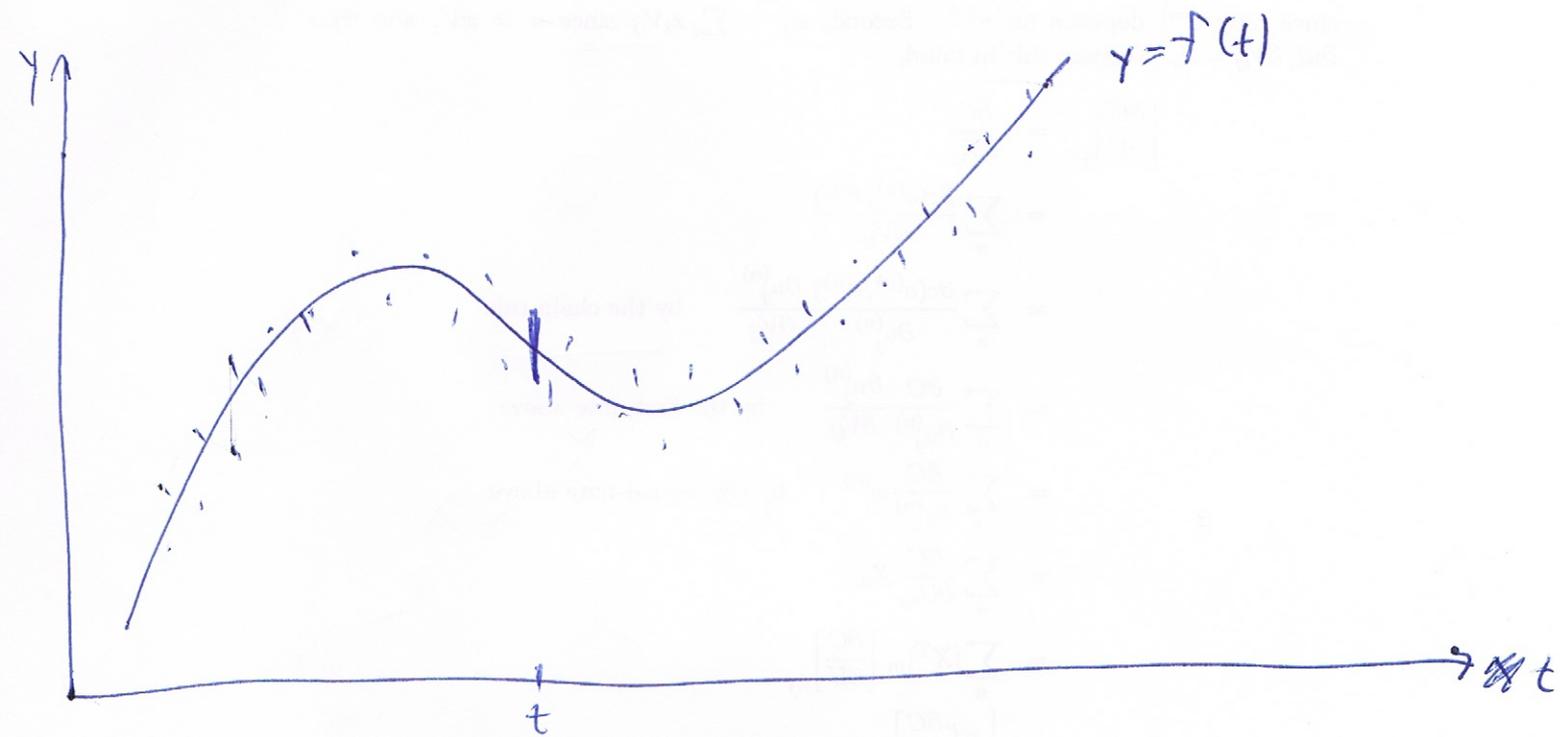
$$\text{eg. } P(x|\hat{x}) = \mathcal{N}(x|\hat{x}, \sigma) = \frac{e^{-\|x-\hat{x}\|^2/2\sigma^2}}{(\sqrt{2\pi}\sigma)^k}$$

$$k = \dim(x)$$

so, images similar to  $\hat{x}$  are most likely.

(2)

compare to regression



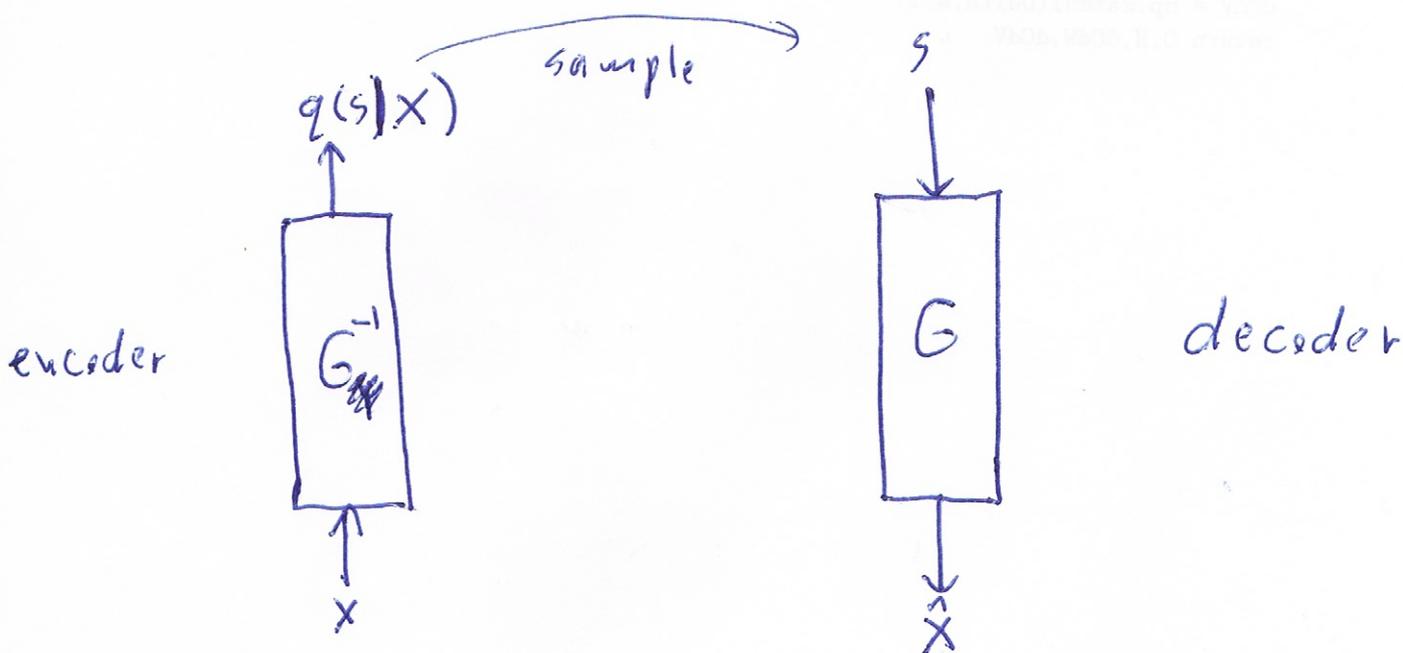
$$P(y|t) = \mathcal{N}(y|f(t), \sigma)$$

ie, given  $t$ , we predict a distribution for  $y$ . Typically,  $f(t)$  is the mean of the distribution. In simple models,  $\sigma$  is fixed, independent of  $t$ .

# Inverse Graphics

(3)

- compute (estimate)  $P(s|x)$
- intractable in general.
- approximate  $P(s|x)$  using a NN (encoder).



Train  $G^{-1}$  (+ maybe  $G$ )

so that  $x \approx \hat{x}$  for all training data.

Backprop through  $G$  +  $G^{-1}$ . How?

# Variational Autoencoders

Given training data  $x^1 \dots x^n$ ,  
train so that  $P(x^1, \dots, x^n)$  is maximal.

$$P(x) = \sum_s P(x, s) = \sum_s P(x|s) \cdot P(s)$$

$\uparrow$                      $\uparrow$                      $\uparrow$   
 huge                    known                    known  
 sum                    learnable

Assume iid data, so  $P(x^1 \dots x^n) = \prod_n P(x^n)$

$$\begin{aligned} \therefore \log P(x^1 \dots x^n) &= \sum_n \log P(x^n) \\ &= \sum_n \log \left[ \sum_s P(x^n, s) \right] \quad \left. \begin{array}{l} \text{sum inside} \\ \log \end{array} \right\} \\ &= \sum_n \log \left[ \sum_s q(s|x^n) \cdot \frac{P(x^n, s)}{q(s|x^n)} \right] \end{aligned}$$

(5)

$$\gg \sum_n \sum_s q(s|x^n) \log \frac{p(x^n, s)}{q(s|x^n)} \stackrel{\text{def}}{=} \mathcal{L}$$

- true for any distribution,  $q$ ,  
because  $\log$  is concave.

$\Delta$

exercise:

$$\mathcal{L} = \log P(x^1 \dots x^n) - \sum_n \text{KL}[q(s|x^n), p(s|x^n)]$$

$\therefore$  maximizing  $\mathcal{L}$  minimizes <sup>KL, i.e.,</sup> the  
difference between  $q(s|x^n)$  &  $p(s|x^n)$   
i.e., makes  $q(s|x^n)$  a better  
approximation of  $p(s|x^n)$ .

# Problem

(S.1)

The sum over  $S$  is huge  
& can't be evaluated. But we  
can approximate (to arbitrary accuracy):

$$\mathcal{L} = \sum_{n=1}^N \mathbb{E}_{s \sim q} \log \frac{P(X^n, s)}{q(s|X^n)} \quad *$$

$$\approx \sum_{n=1}^N \frac{1}{L} \sum_{\ell=1}^L \log \frac{P(X^n, s_\ell^n)}{q(s_\ell^n | X^n)}$$

where  $s_\ell^n \sim q(s|X^n)$

If  $N$  is large, can use  $L=1$ .

(6)

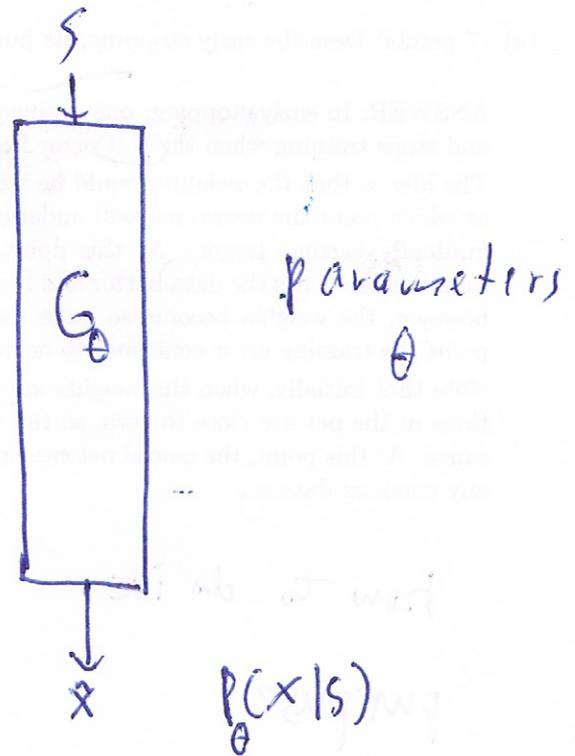
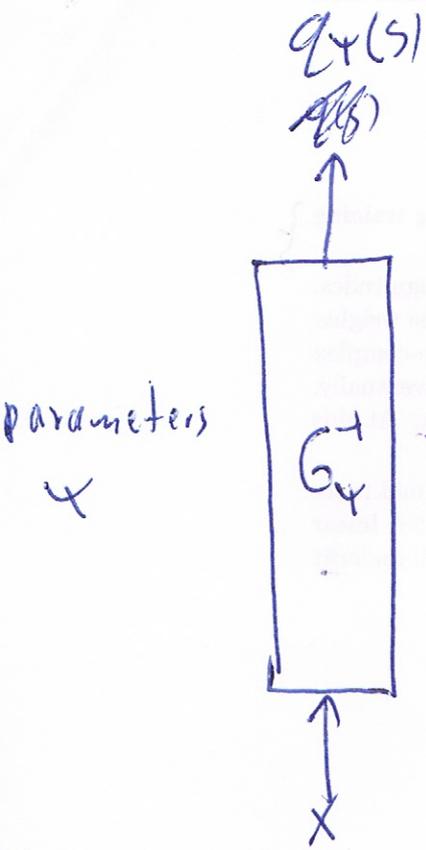
$$\therefore \mathcal{L} \approx \sum_n \log \frac{P(X^n, S^n)}{q(S^n | X^n)}$$

where  $S^n \sim q(S^n | X^n)$

$$= \sum_n \log \frac{P(X^n | S^n) \cdot P(S^n)}{q(S^n | X^n)}$$

$$= \sum_n \left[ \log P_\theta(X^n | S^n) + \log P(S^n) - \log q_\psi(S^n | X^n) \right]$$

Note:  $\nabla_\theta \mathcal{L} = \sum_n \nabla_\theta \log P_\theta(X^n | S^n)$



$$\begin{aligned}
 p_\theta(x|s) &= p(x|x^* = G_\theta(s)) \\
 &= \mathcal{N}(x | G_\theta(s), \alpha) \\
 &= \frac{e^{-(x - G_\theta(s))^2 / 2\alpha^2}}{(\sqrt{2\pi}\alpha)^k} \quad k = \dim(x)
 \end{aligned}$$

$$\begin{aligned}
 \log p_\theta(x|s) &= -\frac{(x - G_\theta(s))^2}{2\alpha^2} \\
 &\quad - k \log(\sqrt{2\pi}\alpha)
 \end{aligned}$$

$$\nabla_{\theta} \mathcal{L} = \sum_s \nabla_{\theta} \log P_{\theta}(x^n | \mathcal{S}^n)$$

(8)

$$= - \sum_s \nabla_{\theta} (x - G_{\theta}(s))^2 / 2\sigma^2$$

$$= \sum_s (x - G(s)) \cdot \nabla_{\theta} G(s)$$

so, back-propagate through NN G

easy

note: G must be differentiable.

$\nabla_{\psi} \mathcal{L}$  is harder

(9)

$\mathcal{L}$  depends on the samples,  ~~$s, s^n$~~   $s^n$ ,  
which ~~in turn depend on~~ we  
draw from  $q_{\psi}(s|x)$ , which  
depends on  $\psi$ .

Go back to \*.

$$\mathcal{L} = \sum_n E_{s \sim q_{\psi}} \log \frac{p_{\theta}(x^n, s)}{q_{\psi}(s|x^n)}$$

want  $\nabla_{\psi} \mathcal{L}$

$$\mathcal{L} = \sum_n E_{s \sim q_{\psi}} \log p_{\theta}(x^n, s) \left. \vphantom{\sum_n} \right\} \text{hard}$$
$$- E_{s \sim q_{\psi}} \log q_{\psi}(s|x^n) \left. \vphantom{\sum_n} \right\} \text{easy}$$

In general

(10)

$$\nabla_{\psi} \int_{\Sigma} E_{sq} f(s)$$

where  $q$  ~~is~~ depends on  $\psi$ .

Two approaches, depending on whether  $s$  is continuous or discrete.

① ~~continuous~~ continuous  $s$ , (reparameterization trick).

~~example: Suppose  $q = \mathcal{N}(\mu, \Sigma)$~~

~~where  $\mu$  &  $\Sigma$  are the outputs of a NN (typically assume  $\Sigma$  is diagonal)~~

~~example: Suppose  $q$  is 1D &  $q = \mathcal{N}(\mu, \sigma)$~~

~~(easily extended to any dimension)~~

# continuous case

(11)

~~IP~~

(reparameterization trick)

Example, Suppose  $s$  is 1D

$$q = \mathcal{N}(\mu, \sigma)$$

(easily extended to any dimension)

~~$q$~~  where  $\mu$  &  $\sigma$  depend on  $\psi$ .

(eg,  $\mu$  &  $\sigma$  are the output of a NN with weights  $\psi$ ).

Then,

$$E_{s \sim q} f(s) = E_{z \sim \mathcal{N}(0,1)} f(\mu + \sigma z)$$

~~$\mu + \sigma z$~~

$$\therefore \nabla_{\Psi} E_{\text{avg}} f(s)$$

$$= E_{\xi \sim \mathcal{N}(0,1)} \nabla_{\Psi} f(\mu + \sigma \xi)$$

~~$$\nearrow E_{\xi \sim \mathcal{N}(0,1)} \frac{\partial f(s)}{\partial s} \Big|_{s=\mu+\sigma\xi}$$~~

$$= E_{\xi \sim \mathcal{N}(0,1)} \cdot f'(\mu + \sigma \xi) \nabla_{\Psi} (\mu + \sigma \xi)$$

$$= E_{\xi \sim \mathcal{N}(0,1)} f'(\mu + \sigma \xi) (\underbrace{\nabla_{\Psi} \mu}_{\neq} + \underbrace{\xi \cdot \nabla_{\Psi} \sigma}_{\neq})$$

*vector-prop*

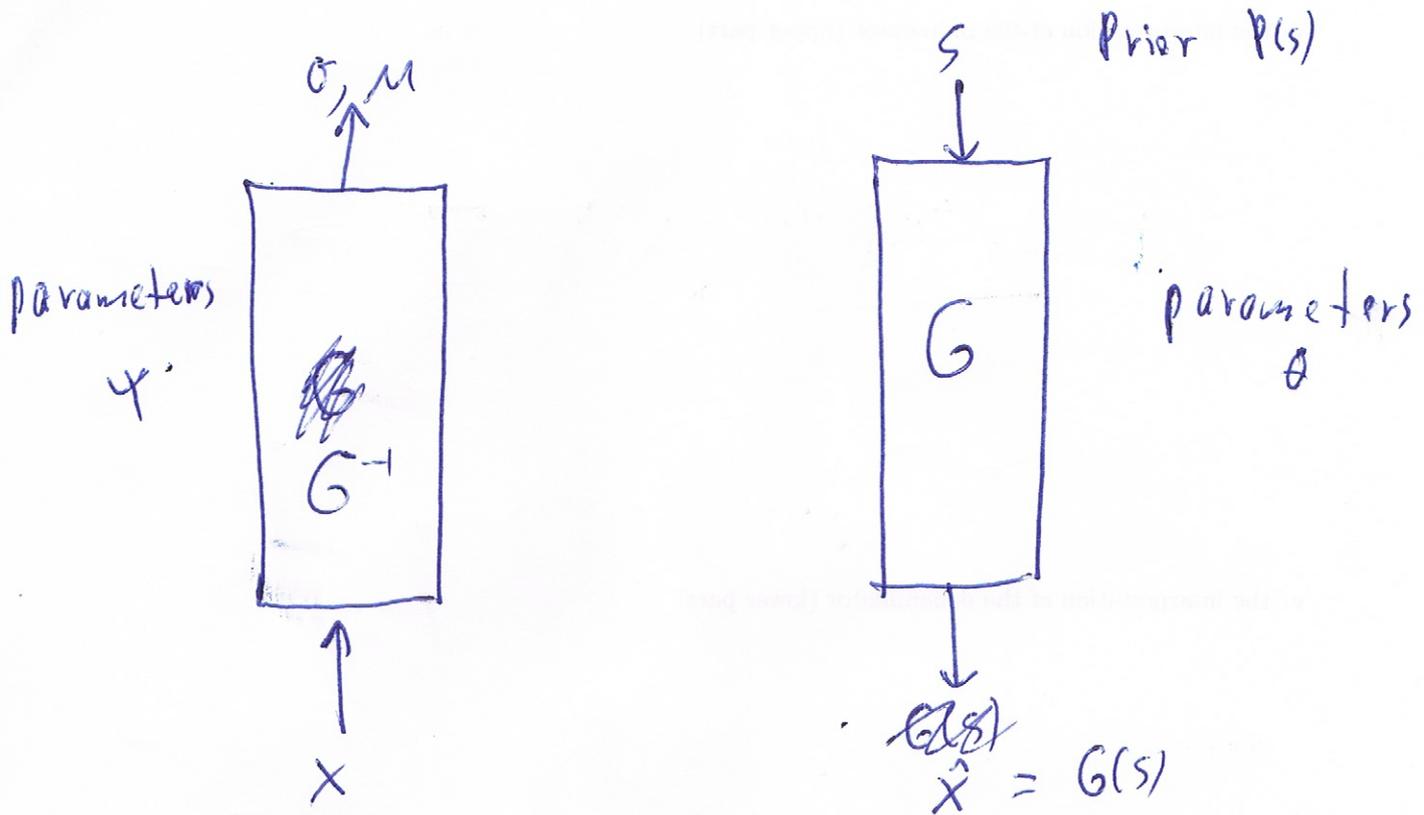
~~$$\approx \frac{1}{L} \sum_{l=1}^L f'(\mu + \sigma \xi_l) (\nabla_{\Psi} \mu + \xi_l \nabla_{\Psi} \sigma)$$~~

where  $\xi_l \sim \mathcal{N}(0,1)$

Note:  $f$  must be differentiable

# Variational Auto encoders

(13)



$$P(x|\hat{x}) = \mathcal{N}(x|\hat{x}, \alpha)$$

typically  $P(s) = \mathcal{N}(s|0, 1)$

$$\begin{aligned} \log P(x|s) &\approx \log P(x|\hat{x}) \\ &= -\frac{(x - \hat{x})^2}{2\alpha^2} + C \\ &= -\frac{(x - G(s))^2}{2\alpha^2} + C \end{aligned}$$

# Recall

(14)

$$\mathcal{L} = \sum_n E_{s \sim q(s|x^n)} \log \frac{p_\theta(x^n, s)}{q_\psi(s|x^n)}$$

$$= \sum_n E_{s \sim q(s|x^n)} \log p_\theta(x^n, s)$$

$$- \sum_n E_{s \sim q(s|x^n)} \log q_\psi(s|x^n)$$

$$= \sum_n E_{s \sim q(s|x^n)} \log p_\theta(x^n | s)$$

$$+ \sum_n H(q(s|x^n))$$

← Entropy of  $q(s|x^n)$   
(usually a simple  
formula)

eg. if  $q(s|x) = \mathcal{N}(s|\mu, \sigma)$

$$\text{then } H(s|x) = \log \sigma + c$$

$$\text{where } c = \frac{1}{2} + \frac{1}{2} \log(2\pi)$$

$$\begin{aligned} \therefore \nabla_{\Psi} \mathcal{L} &= \sum_n \nabla_{\Psi} E_{s \sim q(s|x^n)} \underbrace{\log p_{\theta}(x^n; s)}_{f_n(s)} \\ &\quad + \sum_n \nabla_{\Psi} H(q(s|x^n)) \end{aligned}$$

$$= \sum_n \nabla_{\Psi} E_{s \sim q(s|x^n)} f_n(s)$$

$$+ \sum_n \nabla_{\Psi} \log \sigma_n$$

where  $s|x^n \sim \mathcal{N}(s|\mu_n, \sigma_n)$

$$= \sum_n \nabla_{\Psi} E_{\epsilon \sim \mathcal{N}(0,1)} f_n(\mu_n + \sigma_n \epsilon) \} \text{hard.}$$

$$+ \sum_n \nabla_{\Psi} \log \sigma_n \} \text{easy (back prop through } G^{-1})$$

# The hard part

(15)

$$\sum_n \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} \nabla_{\Psi} f_n(\mu_n + \sigma_n \xi)$$

$$= \sum_n \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} \nabla_{\Psi} f_n(\mu_n + \sigma_n \xi)$$

$$= \sum_n \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} f_n'(\mu_n + \sigma_n \xi) \cdot \nabla_{\Psi}(\mu_n + \sigma_n \xi)$$

$$\approx \sum_n \frac{1}{L} \sum_{\ell=1}^L f_n'(\mu_n + \sigma_n \xi_n^{\ell}) \nabla_{\Psi}(\mu_n + \sigma_n \xi_n^{\ell})$$

where  $\xi_n^{\ell} \sim \mathcal{N}(0,1)$

(can use  $L=1$  when  $N$  is large)

$$\approx \sum_n f_n'(\mu_n + \sigma_n \xi_n) (\nabla_{\Psi} \mu_n + \xi_n \nabla_{\Psi} \sigma_n)$$

↑  
backprop  
through  $G$

↑  
backprop  
through  $G^{-1}$

where  $\xi_n \sim \mathcal{N}(0,1)$

Discrete Case

(s is discrete)

Recall the general problem:

evaluate  $\nabla_{\Psi} E_{s \sim q} f(s)$

where (the parameters of)  $q$  depends on  $\Psi$ .

Simple example:

$$s \in \{0, 1\}$$

$$q(s=1) = \alpha$$

$$q(s=0) = 1 - \alpha$$

where  $\alpha$  depends on  $\Psi$ .

eg,  $\alpha$  is the output of a neural net with weights  $\Psi$ .

(17)

Note that even though  $s$  is discrete, the expectation  $E$  over  $s$  is differentiable as long as  $\alpha$  is differentiable in  $\psi$ .

eg. in the simple example above,

$$\begin{aligned}
 E_{s \sim q} f(s) &= \alpha f(1) + \sum_s q(s) f(s) \\
 &= \alpha f(1) + q(0) f(0) \\
 &= \alpha \cdot f(1) + (1-\alpha) \cdot f(0) \\
 &= \alpha \cdot [f(1) - f(0)] + f(0)
 \end{aligned}$$

$$\therefore \nabla_{\psi} E_{s \sim q} f(s) = [f(1) - f(0)] \cdot \nabla_{\psi} \alpha$$

The problem is when  $s$  is a vector.  
Then the expectation (a sum) has a huge number of terms (eg,  $2^{100}$ ), far too many to enumerate.

But, because it is an expectation, we can try to approximate it using sampling.

~~Alas, for discrete distributions, this almost always returns 0 as the estimate of the gradient.~~

$$\nabla_{\psi} E_{s \sim q} f(s) \approx \nabla_{\psi} \frac{1}{L} \sum_{l=1}^L f(s_l) \quad \text{where } s_l \sim q$$

Since  $q$  depends on  $\psi$ , how do we differentiate through the sampling process?

# re Parameterization Trick doesn't work

example:  $S \in \{0, 1\}$

$$\left. \begin{aligned} q(s=1) &= \alpha \\ q(s=0) &= 1 - \alpha \end{aligned} \right\} \alpha \text{ depends on } \psi$$

To use reparam. trick, must make  $S$  a function of  $\alpha$  & some other random variable,  $\epsilon$ , ind. of  $\psi$

eg.  $S = [\epsilon < \alpha]$  where  $\epsilon \sim U(0, 1)$

$\therefore \Pr(\epsilon < \alpha) = \alpha$

note  $[True] = 1$   
 $[False] = 0$

$$\therefore \nabla_{\gamma} E_{\xi \sim q} f(\xi)$$

$$= \nabla_{\gamma} E_{\xi \sim u(0,1)} f([\xi < \alpha])$$

$$= E_{\xi \sim u(0,1)} \nabla_{\gamma} f([\xi < \alpha])$$

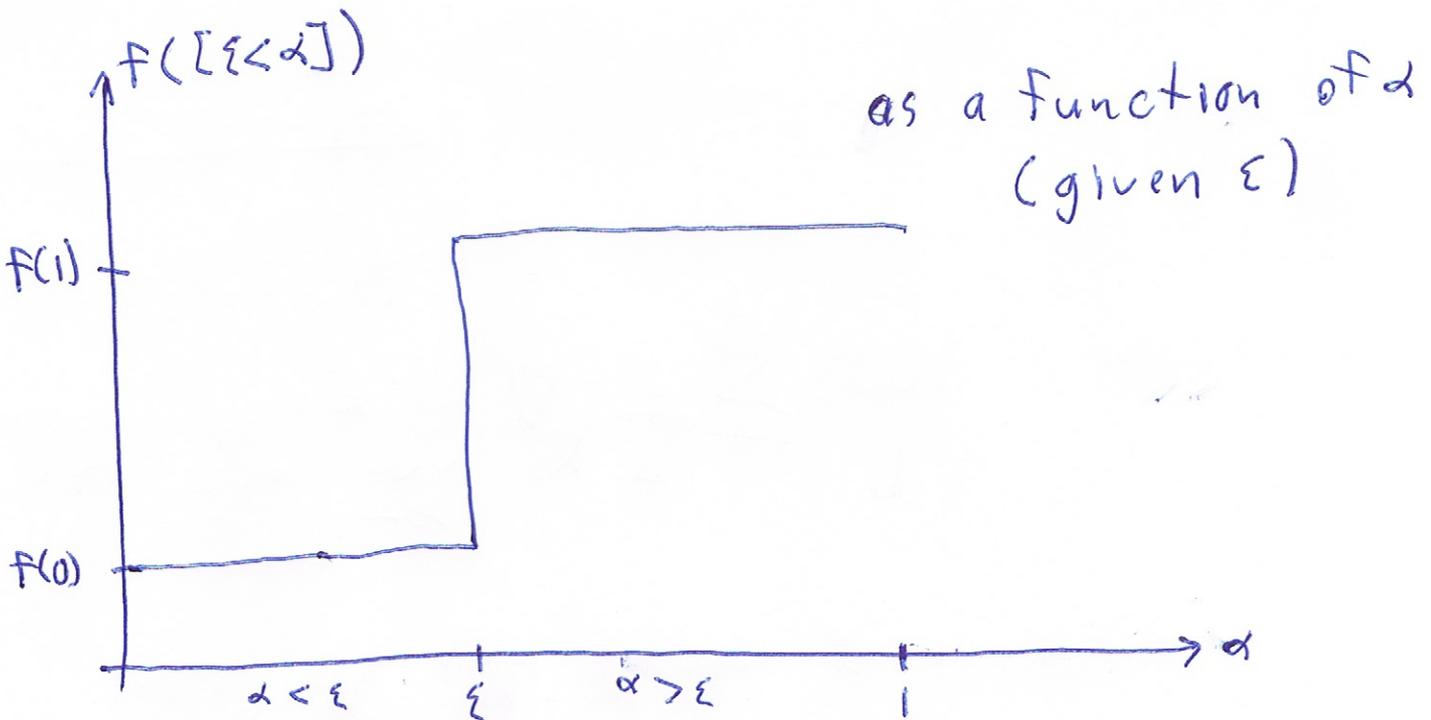
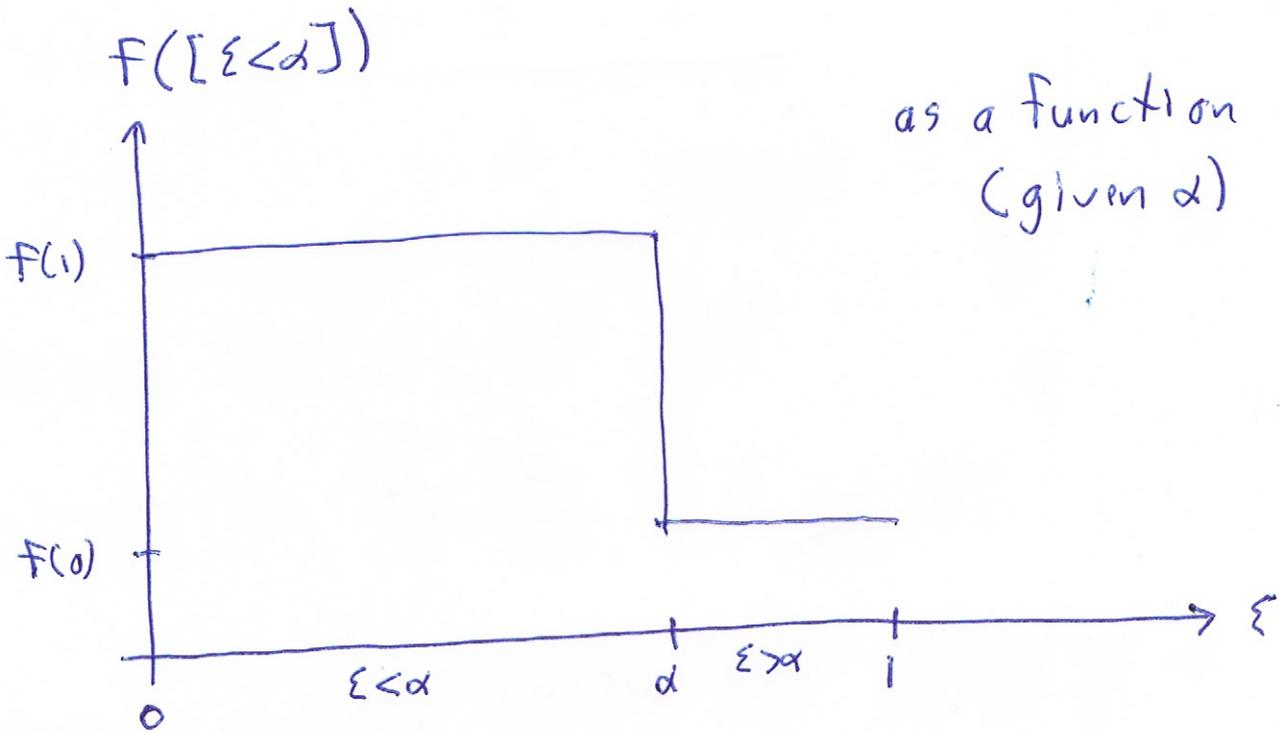
$$= E_{\xi \sim u(0,1)} \frac{\partial f([\xi < \alpha])}{\partial \alpha} \cdot \frac{\partial \alpha}{\partial \gamma}$$

} at try  
param.  
trick

but  $f([\xi < \alpha])$  is constant a.e.

$$f([\xi < \alpha]) = \begin{cases} f(1) & \text{if } \xi < \alpha \\ f(0) & \text{o.w.} \end{cases}$$

$$\therefore \frac{\partial f([\xi < \alpha])}{\partial \alpha} = \begin{cases} 0 & \text{if } \xi \neq \alpha \\ \infty & \text{a.s. if } \xi = \alpha \end{cases}$$



$\therefore \frac{\partial F([\varepsilon < \alpha])}{\partial \alpha} = 0$  almost everywhere

$$\therefore \nabla_{\psi} E_{\text{smg}} f(s)$$

$$= E_{\xi \sim U(0,1)} \frac{\partial f([\xi \leq \alpha])}{\partial \alpha}, \frac{\partial \alpha}{\partial \psi}$$

$$\approx \frac{1}{L} \sum_{e=1}^L \frac{\partial f([\xi_e \leq \alpha])}{\partial \alpha} \frac{\partial \alpha}{\partial \psi}$$

where  $\xi_e \sim U(0,1)$

$$= \begin{cases} 0 & \text{if all } \xi_e \neq \alpha \\ \infty & \text{if some } \xi_e = \alpha \end{cases}$$

ie, a very high variance estimate.

what can we do?