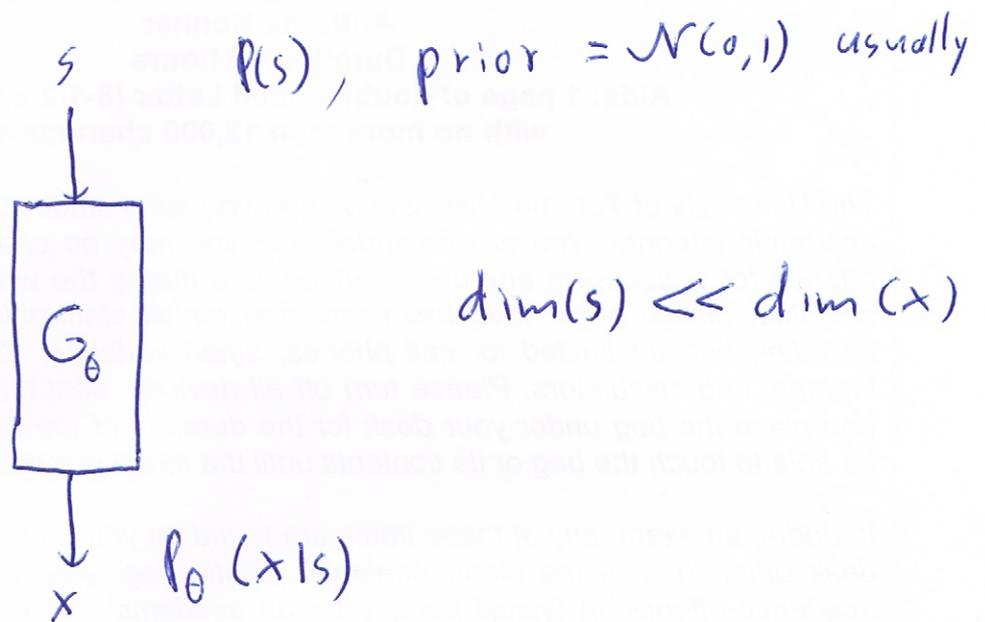


# Generative Models

1

Learn a generative model for a class of images (eg. MNIST).



Goal 1:

find  $\theta$  so that  $P_\theta(x)$  is the distribution of our class of images. (eg, MNIST, or room interiors, or...).

i.e, find  $\theta$  to maximize the probability of the training data,  $x^1 \dots x^N$ .

i.e, maximize  $\prod_n P_\theta(x^n) = P_\theta(x^1 \dots x^n)$

or equivalently  $\sum_n \log P_\theta(x^n) = \log P_\theta(x^1 \dots x^n)$

1.1

assume we can compute  $\frac{\partial P_\theta(x)}{\partial \theta}$

eg,  $G_\theta$  is a NN with weights  $\theta$   
or a "differentiable renderer"  
with parameters  $\theta$ .

Goal 2: compute  $P(s|x)$ , so we can  
do inverse graphics.  
ie, infer the scene from an image.

Lets work on Goal 1 first

(2)

$P_{\theta}(x)$  is unknown, but

$$P_{\theta}(x) = \sum_s P_{\theta}(x, s) = \sum_s P_{\theta}(x|s) \cdot P(s)$$

$P(s)$  is known, &  $P_{\theta}(x|s)$  can be learned (as we shall see),  $\sum_s$  is very large & intractable, but we will see how to deal with it later.

$$\begin{aligned} \therefore \log P_{\theta}(x^1, \dots, x^N) &= \sum_n \log P_{\theta}(x^n) \\ &= \sum_n \log \left[ \sum_s P_{\theta}(x^n, s) \right] \end{aligned}$$

Difficult to maximize because  $\sum_s$  is inside log.

(3)

# Inference

## Variational Approximation

Let  $q_n(s)$  be any distribution (but preferably one that approximates  $P_\theta(s|X^n)$ , which we want to compute.)

Then,

$$\log P_\theta(x^1 \dots x^n)$$

$$= \sum_n \log \left[ \sum_s \frac{q_n(s) P_\theta(x^n, s)}{q_n(s)} \right]$$

$$\geq \sum_n \sum_s^{q_n(s)} \log \frac{P_\theta(x^n, s)}{q_n(s)} \stackrel{\text{def}}{=} \mathcal{L}(\theta, q)$$

$q = (q_1, \dots, q_n)$

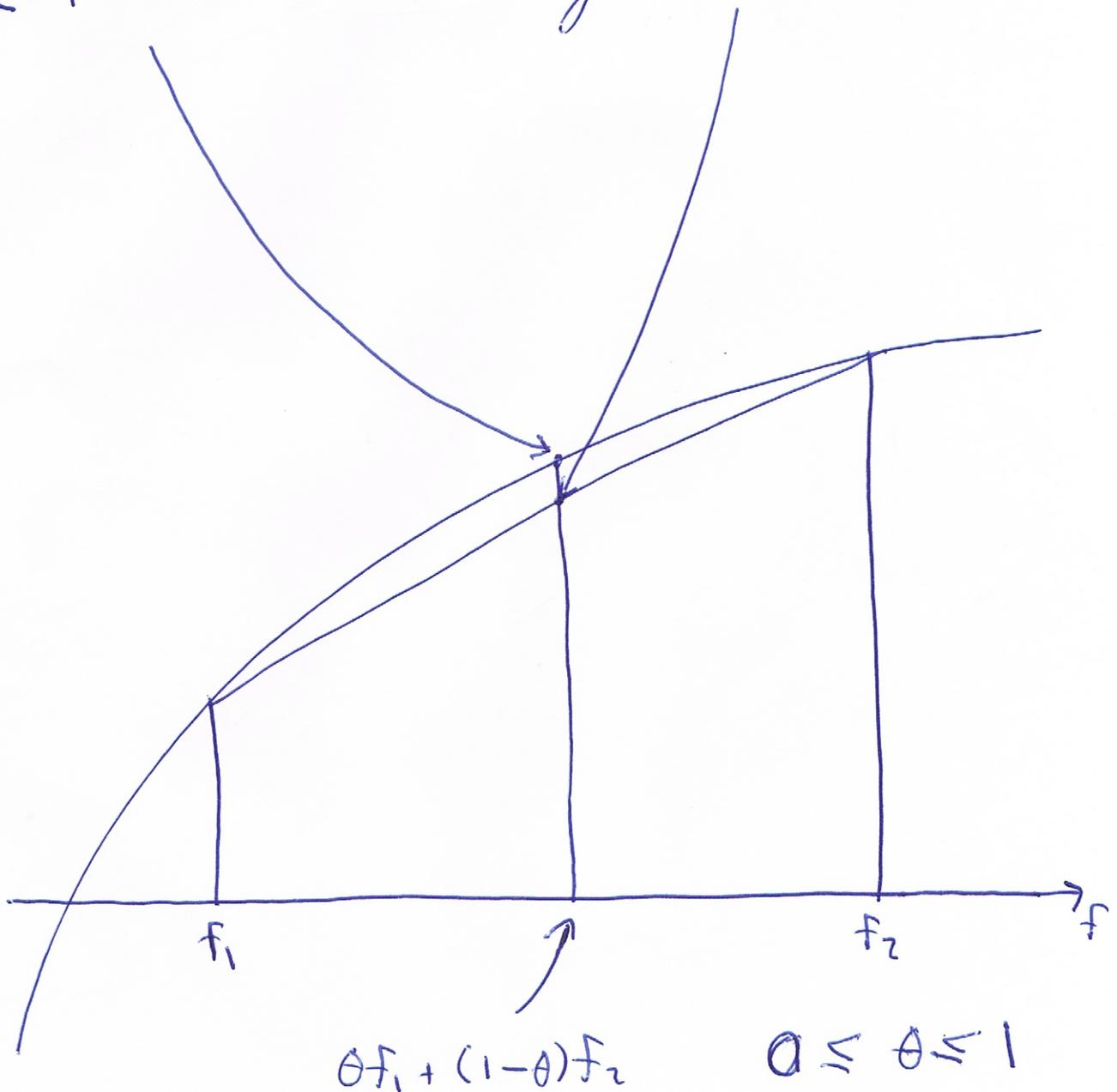
sum is now outside the log!

Why can we do this?

(4)

Because log is concave.

$$\log[\theta f_1 + (1-\theta)f_2] \geq \theta \log f_1 + (1-\theta) \log f_2$$



(5)

$$\text{ie } \log(\theta_1 f_1 + \theta_2 f_2) \geq \theta_1 \log f_1 + \theta_2 \log f_2$$

$$\text{if } \theta_1, \theta_2 \geq 0 \text{ + } \theta_1 + \theta_2 = 1,$$

more generally,

$$\log \sum_s \theta_s f_s \geq \sum_s \theta_s \log f_s$$

$$\text{if } \theta_s \geq 0 \text{ + } \sum_s \theta_s = 1$$

~~with~~ ~~if~~  ~~$\theta_1, \theta_2$~~

note:  $\sum$  comes outside of  $\log$ .

$$\text{let } \theta_s = q_n(s) \text{ + } f_s = P(X^n, s) / q_n(s)$$

then

$$\log \sum_s q_n(s) \frac{P(X^n, s)}{q_n(s)} \geq \sum_s q_n(s) \log \frac{P(X^n, s)}{q_n(s)}$$

note:  $q_n(s) \geq 0$  +  $\sum_s q_n(s) = 1$  since it is a prob. distribution.

so  $\log P_\theta(x^1 \dots x^n) \geq$

~~$\sum_{n,s} q_n(s) \log \frac{P_\theta(x^n, s)}{q_n(s)}$~~

when does this become equality?

when  $q_n(s) = P_\theta(s|x^n)$  = posterior prob. of  $s$ , which is what we also want to compute approximately.  
Proof.

~~$\sum_{n,s} q_n(s) P_\theta(s|x^n) \cdot \log \frac{P_\theta(x^n, s)}{P_\theta(s|x^n)}$~~

$= \sum_{n,s} P_\theta(s|x^n) \cdot \log \frac{P_\theta(x^n, s)}{P_\theta(x^n, s) / P(x^n)}$

$= \sum_{n,s} P_\theta(s|x^n) \cdot \log P_\theta(x^n)$

$= \sum_n \log P_\theta(x^n) \cdot \underbrace{\sum_s P_\theta(s|x^n)}_1$

$= \sum_n \log P_\theta(x^n) = \log P_\theta(x^1 \dots x^n)$

9

Instead of maximizing  $\log P_{\theta}(x^1 \dots x^n)$ , which is intractable, we shall maximize

$$L = \sum_{n, s} q_n(s) \log \frac{P_{\theta}(x^n, s)}{q_n(s)} \quad \text{over } \theta.$$

This is much easier, if we choose  $q_n(s)$  wisely.

~~Let's simplify~~

Note that making  $L$  big forces  $\log P_{\theta}(x^1 \dots x^n)$  to be big too.

~~However, we must choose an appropriate  $q_n(s)$ .~~

# Simplifying

8

$$L = \sum_{n,s} q_n(s) \log \frac{P_\theta(x^n, s)}{q_n(s)}$$

$$= \sum_{n,s} q_n(s) \log \frac{P_\theta(x^n | s) \cdot P(s)}{q_n(s)}$$

$$= \sum_{n,s} q_n(s) \log P_\theta(x^n | s)$$

$$- \underbrace{\sum_n \sum_s q_n(s) \log \frac{q_n(s)}{P(s)}}_{KL(q_n, P)}$$

note:  $KL(q_n, P)$  does not depend on  $\theta$ .

# KL divergence

9

Given two distributions  $q, p$ ,

$$KL(q, p) = \sum_s q(s) \log \frac{q(s)}{p(s)}$$

Exercise:

$$KL(q, p) \geq 0$$

$$KL(q, p) = 0 \quad \text{iff} \quad q = p$$

interpretation:

~~$KL(q, p)$~~  measures the "distance" between two distributions.

but note:  $KL(q, p) \neq KL(p, q)$

∴ When maximizing  $\mathcal{L}$ , we put downward pressure on  $KL(q_n, p)$ .

i.e.,  $KL(q_n, p)$  acts like a regularization term, pressuring  $q_n(s) \approx p(s)$ , the prior on  $s$ .

As we shall see later, there is also pressure to make  $q_n(s) \approx P(s|X^n)$ , the posterior of  $s$  given  $X^n$ .

$L$  is a function of  $\theta$  +  $q = (q_1, \dots, q_n)$ .  
ie,  $L = L(\theta, q)$ .

To maximize  $L$ , we will maximize over both  $\theta$  +  $q$ . In fact, we will alternate ~~we~~ between  $\theta$  +  $q$ .

iii

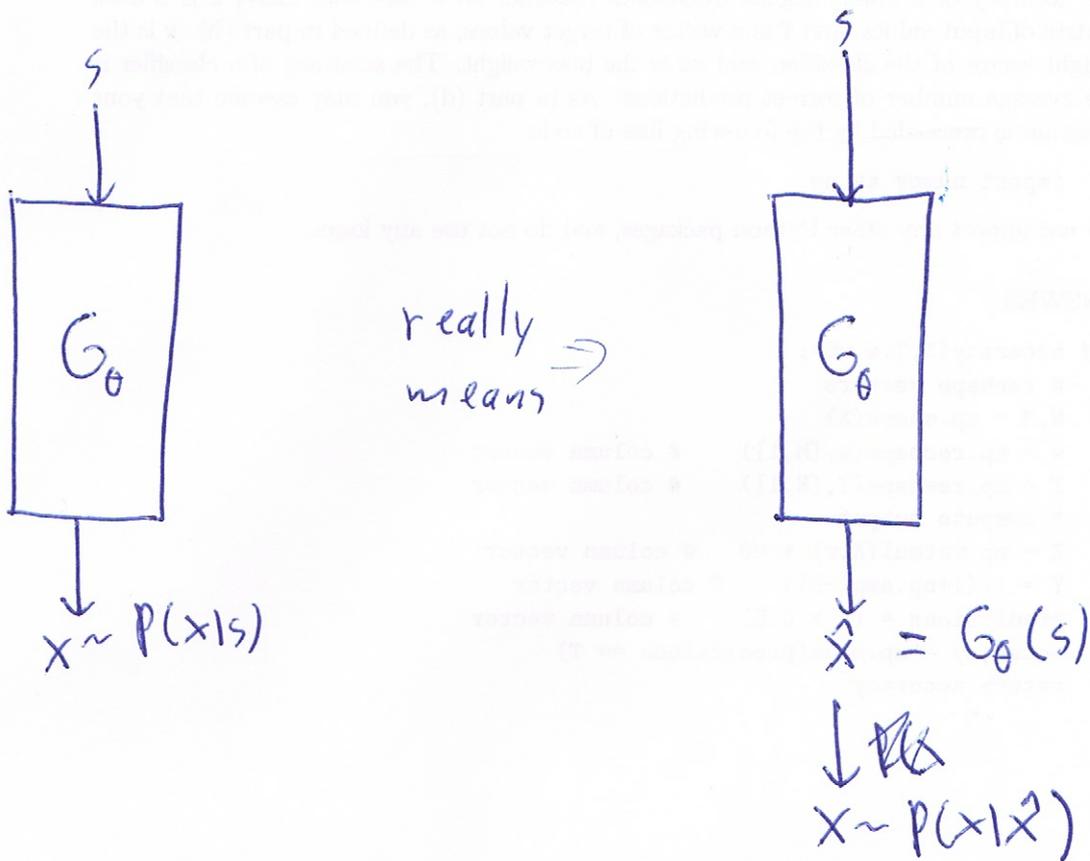
~~do until convergence~~  
~~maximize  $L$  over~~

① we can ~~also~~ maximize over  $\theta$  by doing back-propagation through

$$\sum_{n, s} q_n(s) \log p_\theta(x^n | s)$$

# Our Generative Model

11



usually, so,  $P(x|s) = P(x|\hat{x} = G_\theta(s))$ .

$P$  is chosen so that  $P(x|\hat{x})$  is high iff  $x \approx \hat{x}$

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

$$k = \dim(x).$$

$$\therefore \log P_{\theta}(x|s)$$

$$= -\|x - \hat{x}\|^2 / 2\sigma^2 - k \log \sqrt{2\pi}\sigma$$

$$= -\|x - G_{\theta}(s)\|^2 / 2\sigma^2 - k \log \sqrt{2\pi}\sigma$$

~~$$\frac{\partial \log P_{\theta}(x|s)}{\partial \theta} = \frac{1}{\sigma^2} \|x - G_{\theta}(s)\| \frac{\partial G_{\theta}(s)}{\partial \theta}$$~~

so, maximizing  $\log P_{\theta}(x|s)$  over  $\theta$   
 means minimizing  $\|x - G_{\theta}(s)\|^2$  over  $\theta$ .

we can easily do this by backpropagation  
 if  $G_{\theta}$  is a NN.

So, we can maximize  $L(\theta, q)$   
over  $\theta$  by backprop.

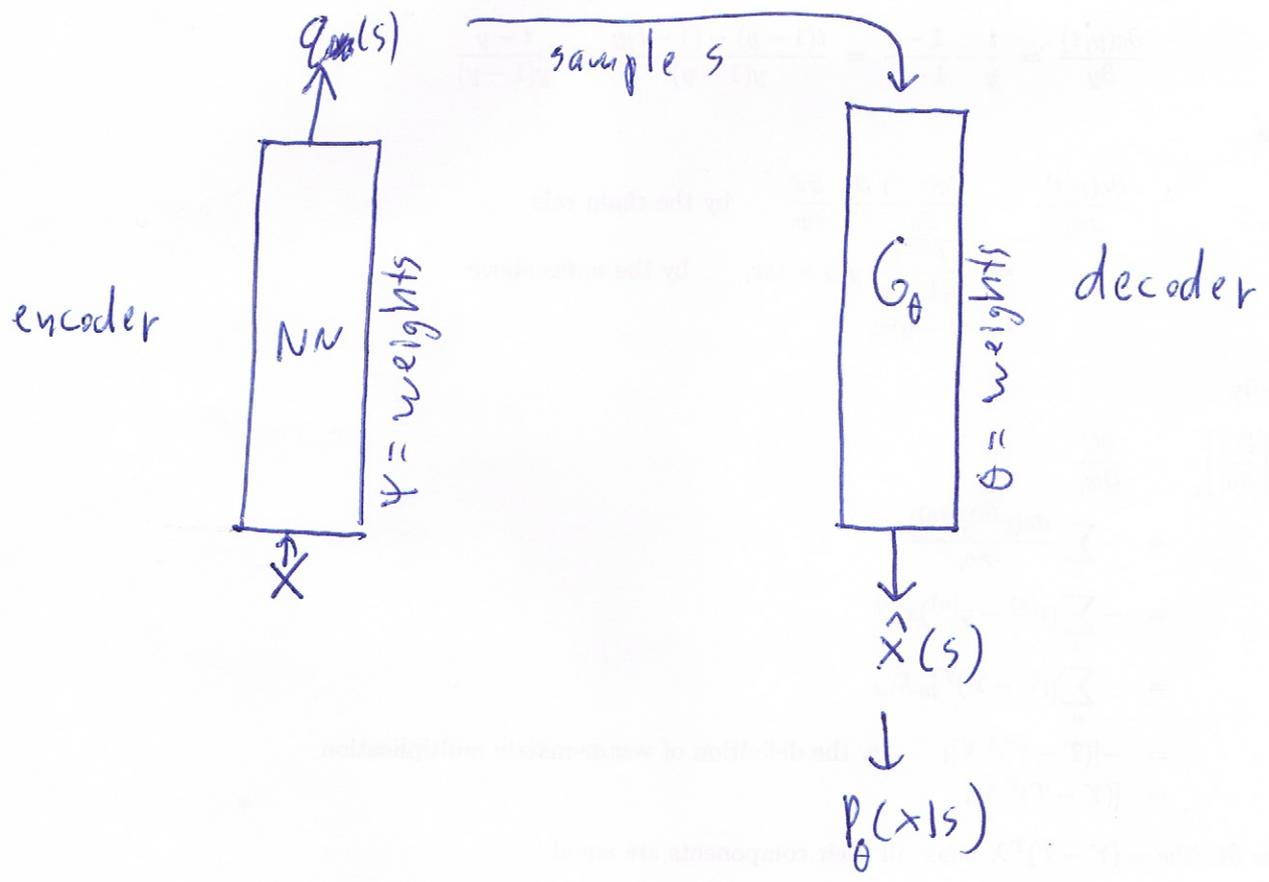
How can we maximize over  $q$ ?

~~This~~

The maximum occurs ~~at~~ when  $q_n(s) = P(S|X^n)$   
as we shall see later. We would love to  
~~can~~ compute this, but it is intractable  
in general. Instead, we ~~assume~~  ~~$q_n(s)$~~  has  
~~a particular~~ limit ourselves to a  
tractable ~~set~~ class of  $q_n$  functions,  
& maximize over that set. This  
is called variational approximation

# Variational Autoencoders

if we assume that  $q_n(s)$  is defined by a neural net with ~~parameters  $\psi$~~  input  $x^n$ , then we have a variational autoencoder.



$$q_n(s) = q_\psi(s|x^n)$$

so,  $L(\theta, q) = L(\theta, \psi)$       it, need to maximize  $L$  over  $\psi$ .

We can also use backprop to maximize ~~over~~  $\mathcal{L}$  over  $\Psi$ . (how?)

But how do we backpropagate through the sampling operation?

(See Kingma & Welling).

In the end, when  $\mathcal{L}$  is maximized,  $q_n(s)$  will be an approximation of  $p(s|X^n)$ , which we wanted.

ie, the ~~de~~ prob. dist. of the scene given an image.

To see this, ~~note~~ <sup>recall</sup> that

$$\mathcal{L} = \sum_{n,s} q_n(s) \log \frac{p(X^n, s)}{q_n(s)} \quad \text{by def.}$$

$$= \sum_{n,s} q_n(s) \log \frac{p(s|X^n) \cdot p(X^n)}{q_n(s)}$$

$$= - \sum_{n,s} q_n(s) \log \frac{q_n(s)}{p(s|X^n)} + \sum_{n,s} q_n(s) \log p(X^n)$$

$$= -\text{KL}[q_n(s), p(s|X^n)] + \sum_n \log p(X^n) \cdot \underbrace{\sum_s q_n(s)}_1$$

$$\therefore \mathcal{L} = -\text{KL}[q_n(s), P(s|X^n)] + \sum_n \log P(X^n)$$

So, maximizing  $\mathcal{L}$  over  $q_n$  is equivalent to minimizing the KL divergence between  $q_n(s)$  &  $P(s|X^n)$ .  
i.e., making  $q_n(s)$  similar to  $P(s|X^n)$

---

Note: by maximizing  $\mathcal{L}$  we are both

- making  $q_n(s)$  similar to  $P(s|X^n)$ , &
- maximizing the likelihood of the data,  $P(x^1, \dots, x^n)$ .