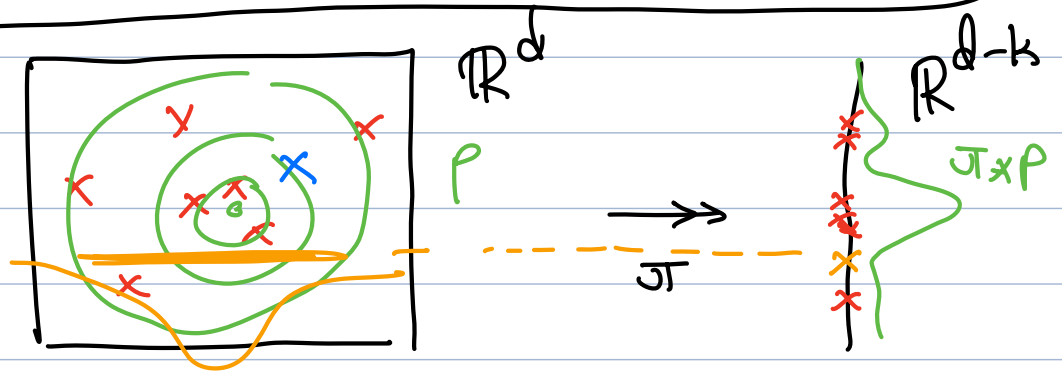


# HOW TO FAKE ANYTHING

TALK ON VARIATIONAL AUTOENCODERS  
FRIDAY, APRIL 19

## I. INTRO + GOAL STATEMENT

MAIN GOAL: GIVEN DATA  $\{x_i\}_{i=1}^n \subseteq \mathbb{R}^d$ ,  
WANT TO FIND + BE ABLE TO USE PDF  $x_i \sim p$ .

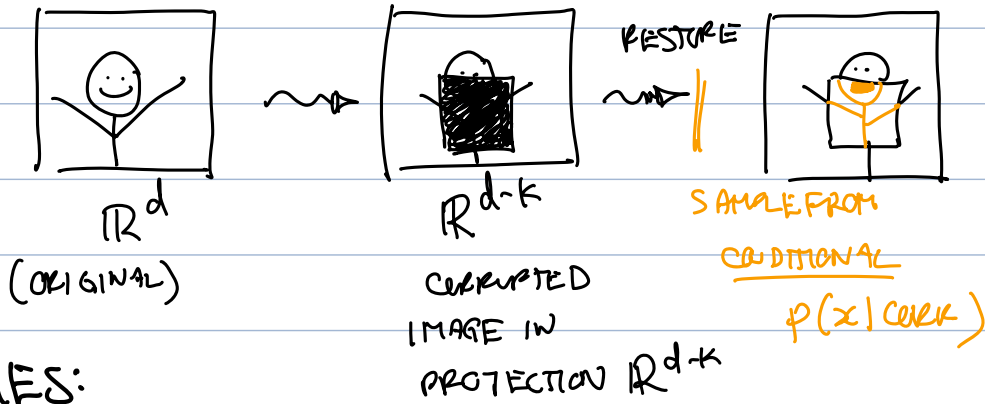


USE? TO SAMPLE, FROM  $p$  (AND RELATED DISTS)!

① SAMPLE FROM  $p =$  GENERATE NEW DATA

② SAMPLE FROM CONDITIONAL  $p(-|c) =$

# INPAINT / RESTORE CORRUPTED DATA



## APPROACHES:

- ① NON-PARAMETRIC:
- (a) HISTOGRAM
  - (b) DENSITY ESTIMATION USING
    - KERNELS:  $K_2 * P_{empirical}$
    - OR THE  $\delta$  FUNCTIONS
    - ...

PROS

CONS

EASY!

GENERAL!

ASYMPTOTICALLY

PRECISE

HARD TO

SAMPLE

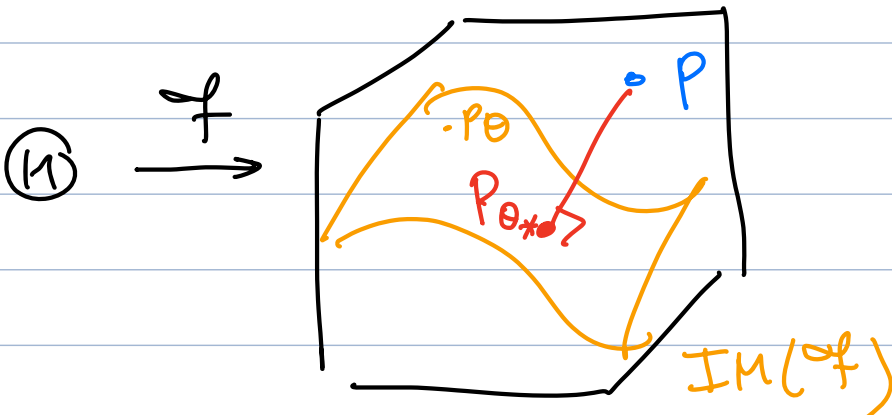
(NEED TO LEARN A FLOW TO A KNOWN DIST)

② PARAMETRIC:

PICK FAMILY  $\mathcal{H}$  OF NICE DISTRIBUTIONS AND FIND APPROX  $P_{0*}$ :

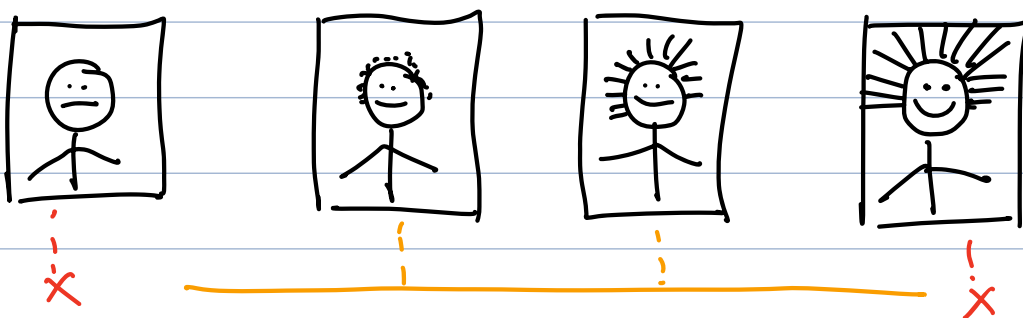
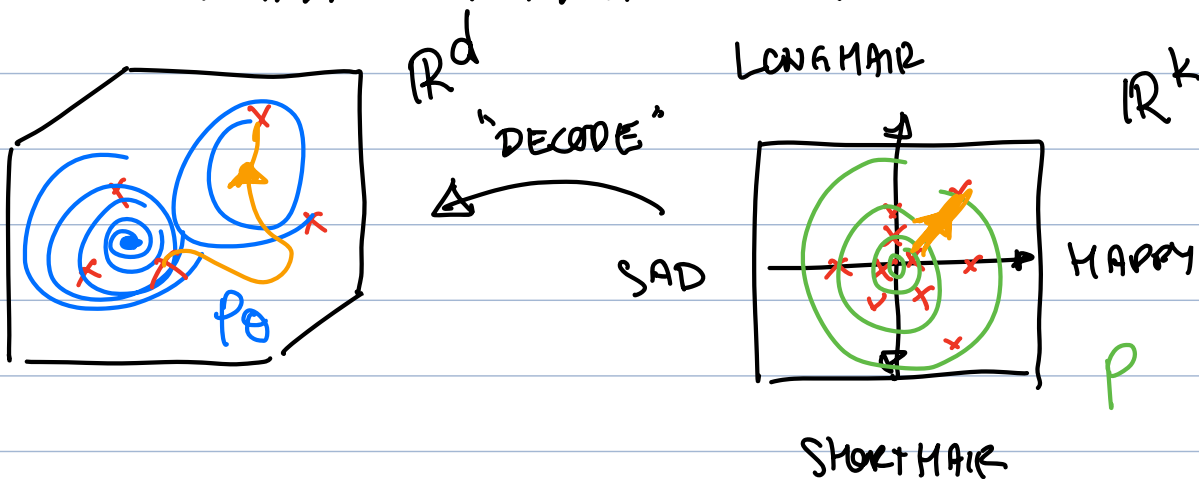
EASIER TO SAMPLE

FUNDAMENTALLY IMPRECISE



# ADDING LATENT VARIABLES

USE (3): INTERPOLATION BIT DATA  
IN A LEARNED FEATURE SPACE



I.E. WE POSIT THE EXISTENCE OF A LATENT SPACE  $\mathbb{R}^k$   
LATENT VAR -  $z$

$$\text{ST } P_0(x) = \int_{\mathbb{R}^k} \underbrace{P_0(x|z)}_{\text{LIKELIHOOD}} \underbrace{p(z)}_{\text{PRIOR}} dz$$

CHOOSE THIS TO  
BE GAUSSIAN

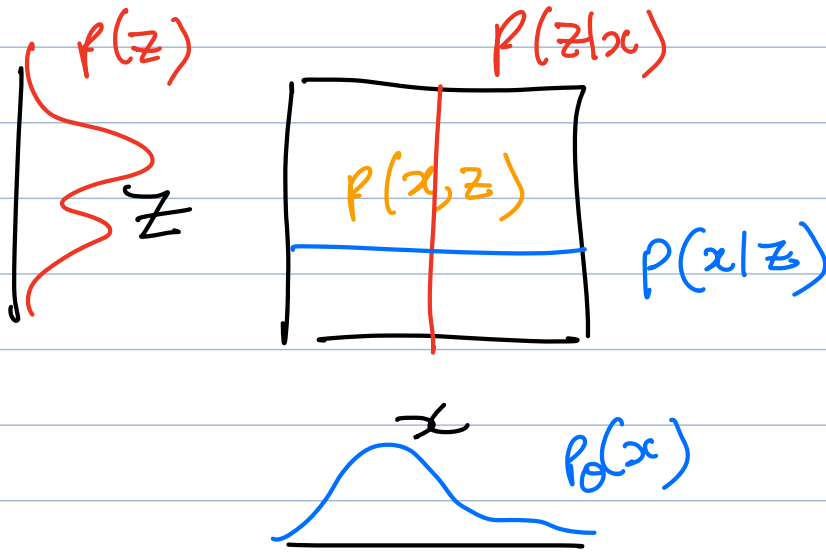
CHOOSE / FIX THIS TO  
BE EASY DIST:  
 $\sim N(0, I)$

$$N(f(z), cI)$$

PRIOR

$$\text{I.E. } \theta = (f, c) \in \mathbb{F} \times \mathbb{R}_{>0}$$

= (1)



$$p(x|z)p(z) = p(x, z) = p(z|x)p(x)$$

$$\therefore p(x) := \frac{p(x|z)p(z)}{p(z|x)} dz$$

POSTERIOR:

APPROXIMATE IT BY ... GAUSSIAN:

$$q_\phi(z|x) \sim N(g(x), h(x) \mathbb{I})$$

## II. WHICH OBJECTIVE FUNCTION?

WE :: ARE AFTER MLE ESTIMATE

$$\theta^* = \underset{\theta, \phi}{\text{ARGMAX}} \prod_{i=1}^n p_{\theta}(x_i)$$

$$= \underset{\theta}{\text{ARGMAX}} \sum_{i=1}^n \log p_{\theta}(x_i)$$

$$= \underset{\theta}{\text{ARGMAX}} \sum_{i=1}^n \log \int_{\mathbb{R}^k} p_{\theta}(x_i | z) p(z) dz$$

REPLACE BY:

$$(\theta^*, \phi^*) = \underset{\theta, \phi}{\text{ARGMAX}} \sum_{i=1}^n \log \int \frac{p_{\theta}(x_i | z)}{q_{\phi}(z | x_i)} p(z) \underline{q_{\phi}(z | x_i)} dz$$

$$= \underset{\theta, \phi}{\text{ARGMAX}} \sum_{i=1}^n \log \mathbb{E}_{z \sim q_{\phi}(z | x_i)} \left[ \frac{p_{\theta}(x_i | z)}{q_{\phi}(z | x_i)} p(z) \right]$$

LET'S LOOK AT:

$$\log \mathbb{E}_{z \sim q_{\phi}(z | x_i)} \left[ \frac{p_{\theta}(x_i | z)}{q_{\phi}(z | x_i)} p(z) \right]$$

JENSEN

$$\geq \mathbb{E}_{z \sim q_{\phi}(z | x_i)} \left[ \log \frac{p_{\theta}(x_i | z)}{q_{\phi}(z | x_i)} p(z) \right]$$

$$\textcircled{1} \rightarrow \mathbb{E}_{z \sim q_{\phi}(-|x_i)} [\log p_{\theta}(x_i, z)] + \underbrace{\mathbb{H}(q_{\phi}(-|x_i))}_{\geq 0}$$

$$\boxed{\text{ELBO}} =: \mathcal{L}_i(\theta, \phi)$$

$$\textcircled{2} \rightarrow = \mathbb{E}_{z \sim q_{\phi}(-|x_i)} \left[ \log \frac{p_{\theta}(z|x_i)}{q_{\phi}(z|x_i)} + \log p_{\theta}(x_i) \right]$$

$$= \underbrace{-D_{\text{KL}}(q_{\phi}(-|x_i) \| p_{\theta}(-|x_i))}_{\geq 0} + \boxed{\log p_{\theta}(x_i)}$$

THE START OF OUR WHOLE JOURNEY!

$$\therefore \log p_{\theta}(x_i) = \underbrace{D_{\text{KL}}(q_{\phi}(-|x_i) \| p_{\theta}(-|x_i))}_{\geq 0} + \mathcal{L}_i(\theta, \phi)$$

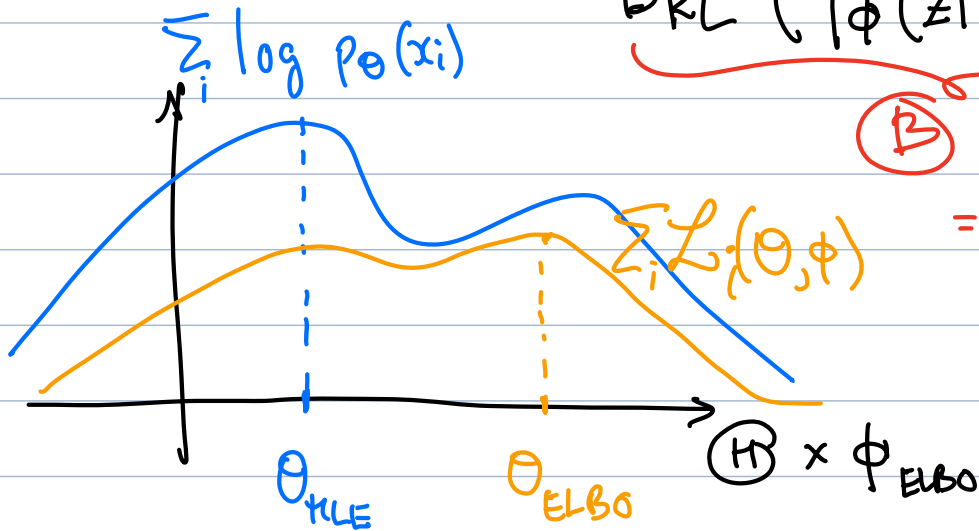
RECALL OUR GOAL: MAXIMIZE  $\sum_{i=1}^n \log p_{\theta}(x_i)$

WE NOW REPLACE IT BY THE DIFFERENT BUT RELATED

NEW GOAL: MAXIMIZE

$$\mathcal{L}(\theta, \phi) = \sum_{i=1}^n \mathbb{E}_{z \sim q_{\phi}(-|x_i)} \left[ \log \frac{p_{\theta}(x_i|z)}{q_{\phi}(z|x_i)} p(z) \right]$$

$$= \sum_{i=1}^n \left\{ \underbrace{\mathbb{E}_{z \sim q_{\phi}(\cdot | x_i)} [\log p_{\theta}(x_i | z)]}_{\text{(A)}} - \underbrace{\text{D}_{\text{KL}}(q_{\phi}(z | x_i) \| p(z))}_{\text{(B)}} \right\}$$



Q: WHAT DOES THIS NEW GOAL MEAN?

A: (1) LOOK AT ELBO ABOVE: TAKE DATA PT  $x_i$  AND

(i) SAMPLE  $z \sim q_{\phi}(\cdot | x_i)$ .

$z$  IS CODE DESCRIBING  $x_i$

(SO CALL  $q_{\phi}$  THE ENCODER)

(ii) TERM (A) IS THE LOG-LIKELIHOOD OF THE OBSERVED  $x_i$ , GIVEN THE CODE

$z$  THAT WE SAMPLED.

THIS TERM IS BIG WHEN  $p(x_i | z)$  IS BIG.

$p_\theta(x_i | z)$  TRIES TO RECONSTRUCT  $x_i$  FROM CODE

(SO CALL  $p_\theta(x | z)$  THE DECODER)

$\therefore$  = RECONSTRUCTION ERROR !

(iii) **TERM (B)**: PRESSURE ON DKL TO BE SMALL ... I.E FOR  $q_\phi$  TO LOOK LIKE  $p(z)$  ... I.E FOR CODES  $z$  FOR  $x_i$  TO LOOK

STD GAUSSIAN. I.E

PREVENTS CODE  $q_\phi$  FROM SPOKING

TO REMEMBER JUST THE IDENTITY MAP,

AND  $\therefore$  LEARN A MORE INTERESTING MAP

$\therefore$  REGULARIZATION TERM !



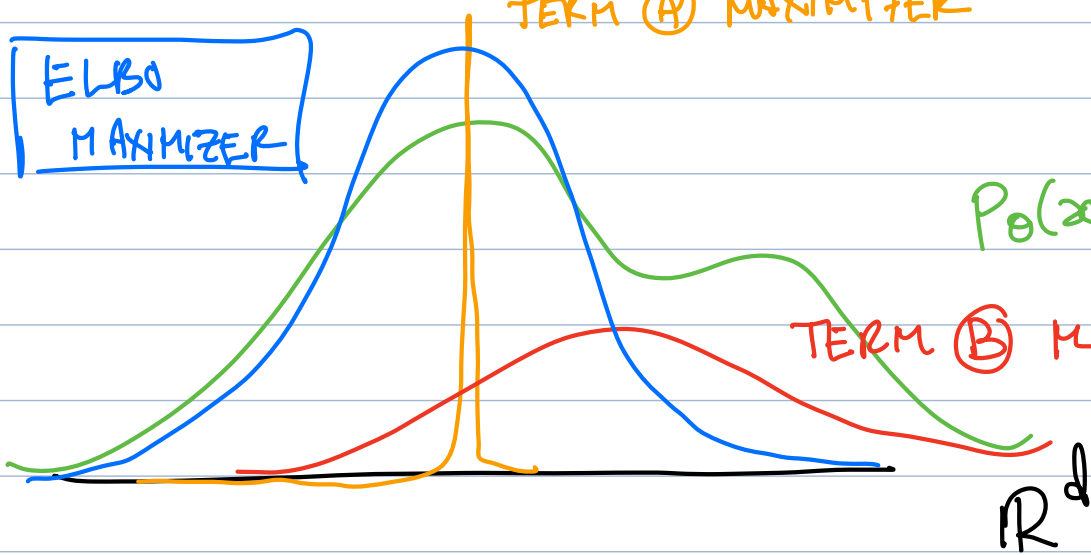
ELBO  
MAXIMIZER

TERM (A) MAXIMIZER

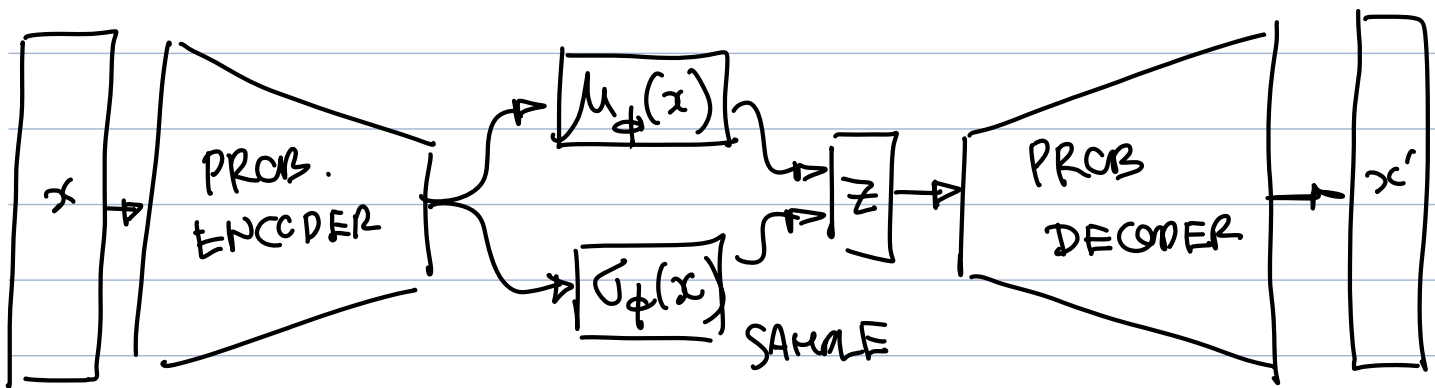
$P_{\theta}(x_i|z)$

TERM (B) MAXIMIZER

$\mathbb{R}^d$



### III. TRAINING



$$q_\phi(z|x) \sim N(\mu_\phi(x), \sigma_\phi(x))$$

$$p_\theta(x|z) \sim N(\mu_\theta(z), \sigma_\theta)$$

RECALL OUR OBJECTIVE:

$$\mathcal{L}_i(\theta, \phi)$$

MAXIMIZE

$$\sum_{i=1}^N \left\{ \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \log(p_\theta(x_i|z) p(z)) \right] + \mathcal{D}(q_\phi(z|x_i)) \right\}$$

NAIVE

(ONE ITERATION)

ALGORITHM FOR EACH  $x_i$  (OR MINI-BATCH)

① CALCULATE

$$\nabla_\theta \mathcal{L}_i(\theta, \phi)$$

② CALCULATE

$$\nabla_{\phi} \mathcal{L}_i(\theta, \phi)$$

③ TAKE STEP  $\theta \mapsto \theta + \alpha \nabla_{\theta} \mathcal{L}_i$

$$\phi \mapsto \phi + \alpha \nabla_{\phi} \mathcal{L}_i$$

HOW TO CALCULATE ① AND ②?

IDEA #1: JUST DO IT. FOLLOW THE APPROACH:

$$\nabla_{\theta} \mathbb{E}_{q_{\phi}} [r(\theta, \phi)] = \mathbb{E}_{q_{\phi}} [\nabla_{\theta} r(\theta, \phi)]$$

$$\begin{aligned} \nabla_{\phi} \mathbb{E}_{q_{\phi}} [r(\theta, \phi)] &= \nabla_{\phi} \int r(\theta, \phi) q_{\phi} dz \\ &= \int (\nabla_{\phi} r(\theta, \phi)) q_{\phi} dz \\ &\quad + \int r(\theta, \phi) \nabla_{\phi} q_{\phi} dz \end{aligned}$$

$$\nabla_{\theta} \mathcal{L}_i = \mathbb{E}_{q_{\phi}} \left[ \nabla_{\theta} r(\theta, \phi) + r(\theta, \phi) \nabla_{\theta} \log q_{\phi} \right]$$

$$r(\phi) = -\log q_{\phi}$$

For us, this is:

$$\begin{cases} \nabla_{\theta} \mathcal{L}_i & \stackrel{(*)}{=} \mathbb{E}_{q_{\phi}} \left[ \nabla_{\theta} \log p_{\theta}(x_i | z) \right] \\ \nabla_{\phi} \mathcal{L}_i & \stackrel{(*)}{=} \mathbb{E}_{q_{\phi}} \left[ \log p_{\theta}(x_i | z) \nabla_{\phi} \log q_{\phi}(z) \right] \\ & + \mathbb{E}_{q_{\phi}} \left[ \frac{-\nabla_{\phi} q_{\phi}}{q_{\phi}} - \log q_{\phi} \nabla_{\phi} \log q_{\phi} \right] \end{cases}$$

IDEA #2: WAIT... SO EVERYTHING IS AN EXPECTATION OVER AN EASY DIST.  $q_{\phi, x_i}$ !

So: INSTEAD OF CALCULATING, APPROXIMATE BY SAMPLING!

$$z_1, \dots, z_k \sim q_{\phi, x_i}$$

$$\nabla_{\theta} \mathcal{L}_i \stackrel{(*)}{\approx} \frac{1}{k} \sum_{j=1}^k \nabla_{\theta} \log p_{\theta}(x_i | z_j)$$

$\nabla_{\phi} \mathcal{L}_i \stackrel{(*)}{\approx}$  SIMILAR

IDEA #3: WE'RE DOING THIS OVER AND OVER...

MAYBE WE CAN SAMPLE JUST ONE  $z$  ( $k=1$ ) ?  
(FOR EACH OF  $\nabla_{\theta} \mathcal{L}_i, \nabla_{\phi} \mathcal{L}_i$ )

IDEA #4: ACTUALLY... JUST SAMPLE

ONE  $\epsilon \sim N(0, 1)$  :

THAT'S BECAUSE

$$z \sim q_{\phi, x_i}(z) = N(\mu_{\phi}(x_i), \sigma_{\phi}(x_i))$$

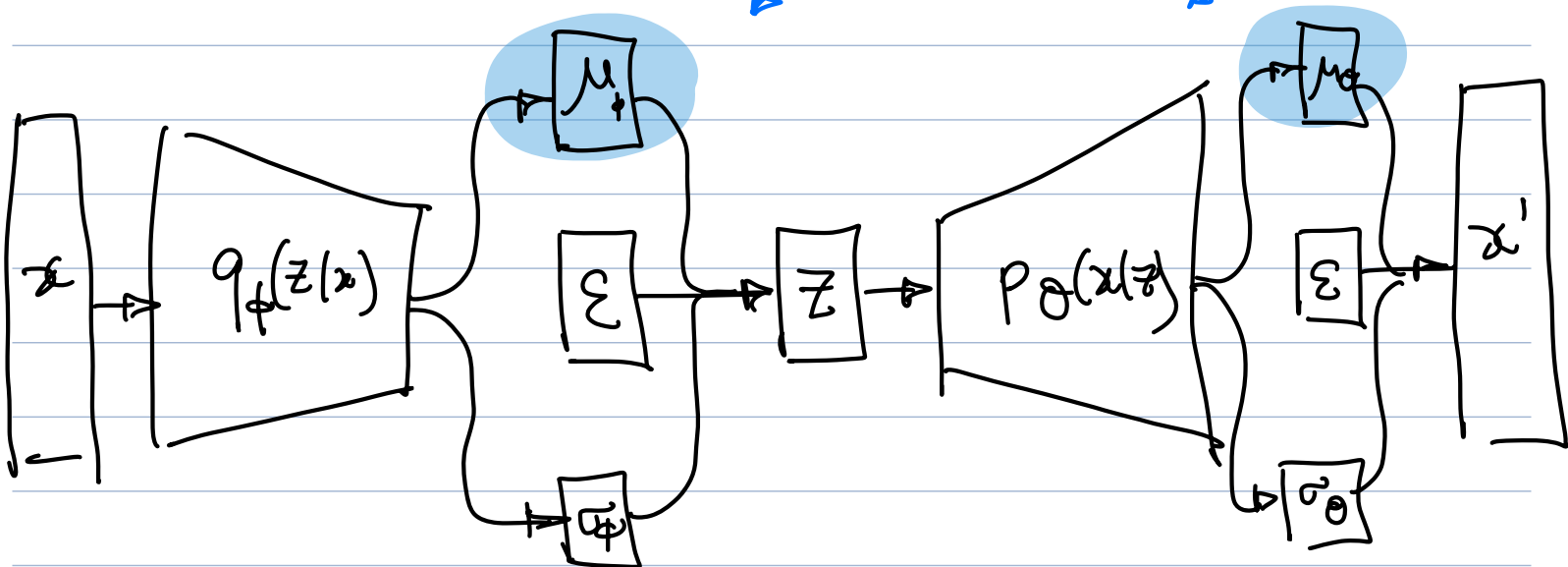
IE.

$$z = \underbrace{\mu_{\phi}(x_i)}_{\text{DETERMINISTIC (BUILT BY NN'S)}} + \underbrace{\sigma_{\phi}(x_i)}_{\sim N(0, 1)} \epsilon$$

... I'LL SOAK YOU THE FINAL FORMULA.

## FINAL SCHEMATIC:

SHADOWS OF A  
VANILLA AUTOENCODER



## FINAL TRAINING ALGO:

FOR EACH  $x_i$ : (IE ONE ITERATION)

① SAMPLE ONE  $\epsilon_1 \sim N(0, 1)$  TO

APPROXIMATE  $\nabla_{\theta} \mathcal{L}_i(\theta, \phi)$

② SAMPLE A SECOND  $\epsilon_2 \sim N(0, 1)$  TO

APPROXIMATE  $\nabla_{\phi} \mathcal{L}_i(\theta, \phi)$

③ STEP  $\theta \mapsto \theta + \alpha \widetilde{\nabla_{\theta} \mathcal{L}_i}$

$\phi \mapsto \phi + \alpha \widetilde{\nabla_{\phi} \mathcal{L}_i}$

## IV: GENERATION

TRAINING PRODUCES A DISTRIBUTION THAT

APPROXIMATES  $P_{\text{EMPIRICAL}}(x)$ :

$$x_i \sim \int_{\mathbb{R}^d} \underbrace{P_{\theta^*}(x|z)}_{\text{DECODER}} p(z) dz =: P_{\theta^*}(x)$$

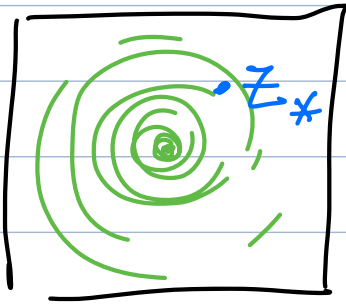
HERE ARE FUN THINGS YOU CAN NOW DO:

① GENERATE NEW DATA:

TO GENERATE "NEW DATA," YOU "SAMPLE FROM  $P_{\theta^*}(x)$ " BY:

① SAMPLE  $z_* \sim p(z) = N(0, I)$

② SAMPLE  $x_* \sim P_{\theta^*}(x|z_*)$   
 $= N(\mu_{\theta^*}(z_*), \Sigma_{\theta^*}^{-1})$



$z \in \mathbb{R}^d$



$x \in \mathbb{R}^n$

THIS  $x_*$  SHOULD LOOK LIKE  $\{x_i\}_{i=1}^N$  !

## ② INTERPOLATE USING LATENT SPACE

GIVEN TWO DATA POINTS  $x_0, x_{k+1}$

$$\left( \rightarrow z_i := \mu_\phi(x_i) \right)$$

OR <sup>JUST</sup> TWO LATENT SAMPLES  $z_0, z_{k+1} \sim N(\mu, \sigma^2)$

CAN

①

DRAW A STRAIGHT LINE

$z_0 \rightsquigarrow z_{k+1}$  IN LATENT SPACE,

②


PICK EVENLY-SPACED  $z_i$  ON

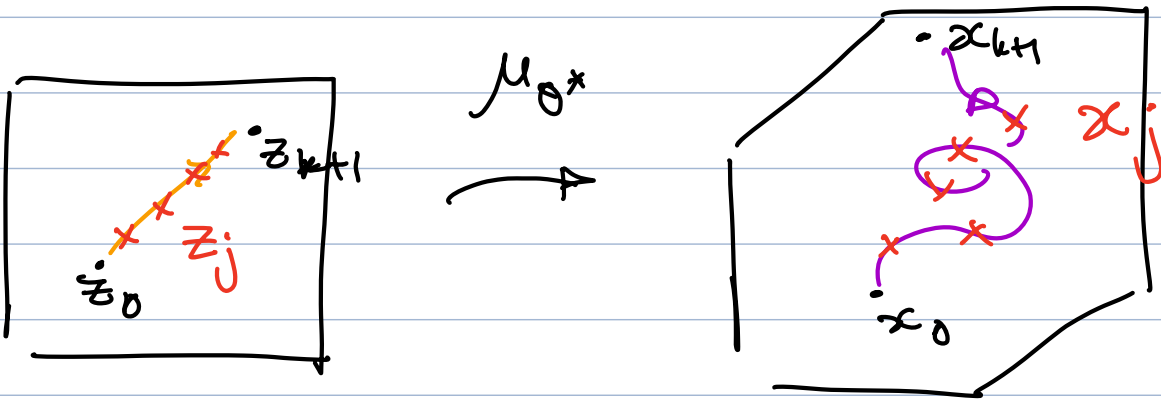
THE LINE



③ YOUR INTERPOLATION WILL BE

$$x_j := \mu_{\theta^*}(z_j)$$

$z_0$    $z_{k+1}$



YOU CAN HOPE THAT THE NN'S LEARNED

TO IDENTIFY CERTAIN FEATURES OF THE TRAINING

DATA WITH DIRECTIONS IN LATENT SPACE!

... HENCE THE SMILE - FROWN

LONG - SHORT HAIR

CARTOON FROM EARLIER.

THANKS FOR LISTENING!

