

Foundation of Generative Models

We will discuss neural networks
and machine learning in terms
of probabilities and distributions

instead of input and output.

Terminologies :

$p(x) : X \rightarrow \mathbb{R}^+$ such that

$$\left\{ \begin{array}{l} \sum_{x \in X} p(x) = 1 \\ \int_X p(x) dx = 1 \end{array} \right.$$

$p_\theta(x)$ is a learned model of $p(x)$

Joint distribution $p(x, y)$

Conditional distribution $p(x|y)$

From distribution $p(\cdot)$ to instance x :

$x \in$ Discrete Space:

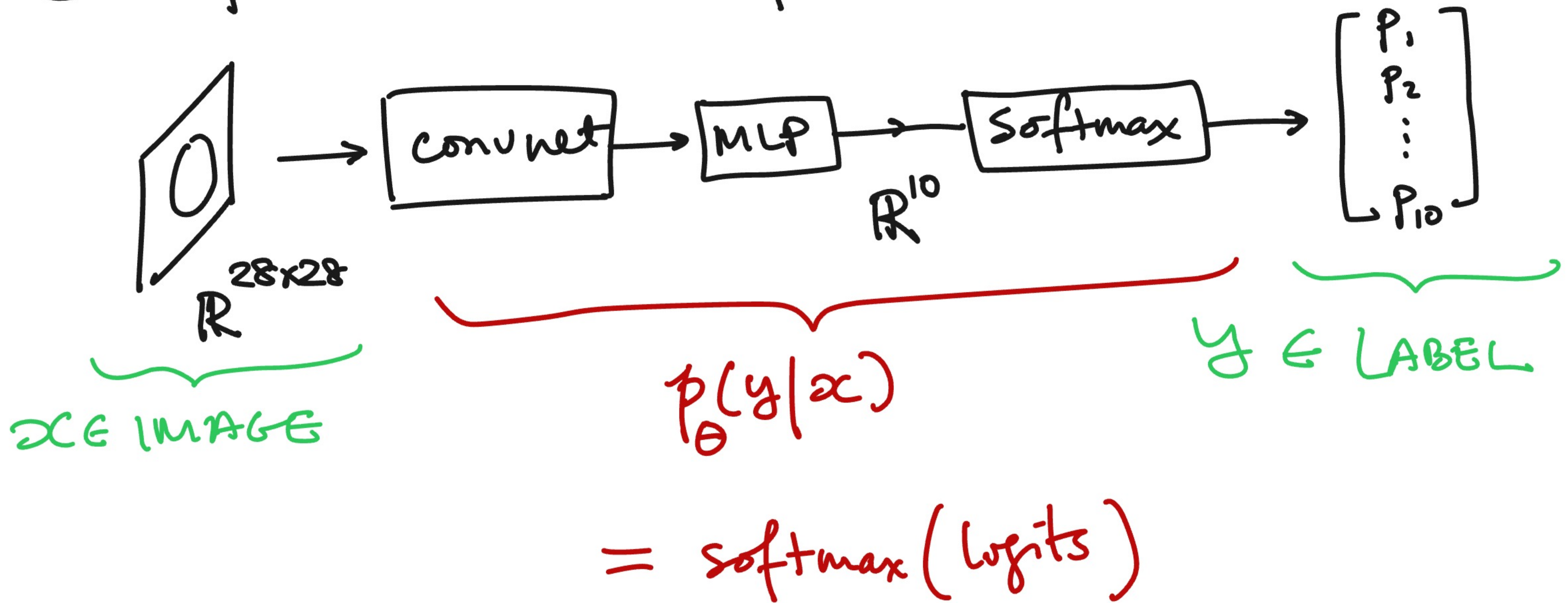
Sample $(P) \rightarrow x$

Pick from discrete space
and favour samples x
w/ larger $p(x)$.

$x \in \mathbb{R}^d$:

$$x = E[X|P] = \int_{\mathbb{R}^d} x \cdot p(x) dx$$

Classification as a probabilistic model



CLASSIFICATION \equiv DISCRIMINATIVE MODEL

Training data = pairs of random variables
 (x, y)

Learning the conditional distribution

$$P_{\theta}(y|x) \approx P_{\text{data}}(y|x)$$

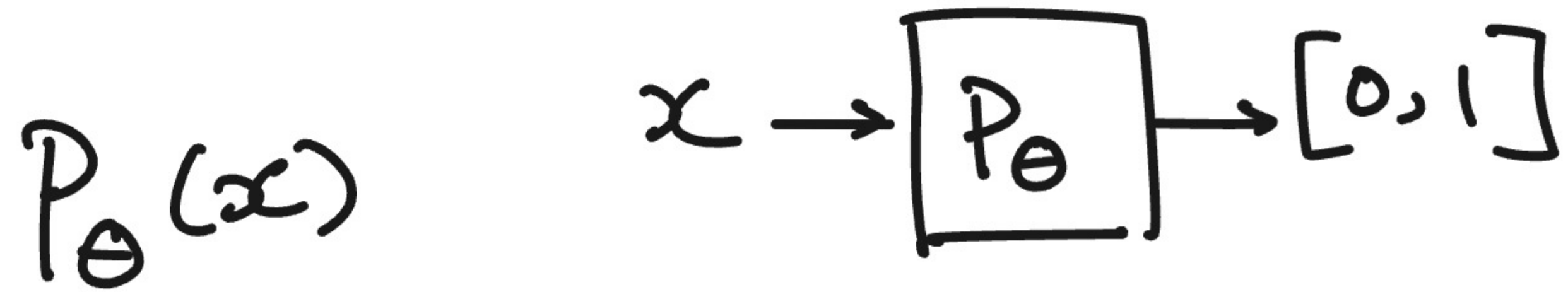
GENERATIVE MODEL

Training data $D = \{x_1, x_2, \dots, x_N\}$
are samples of a random var $x \in \mathbb{R}^d$.

Want to learn $P_\theta(x) \approx P_{\text{data}}(x)$

Given x , we know its likelihood $p(x)$!
(Explicit model)

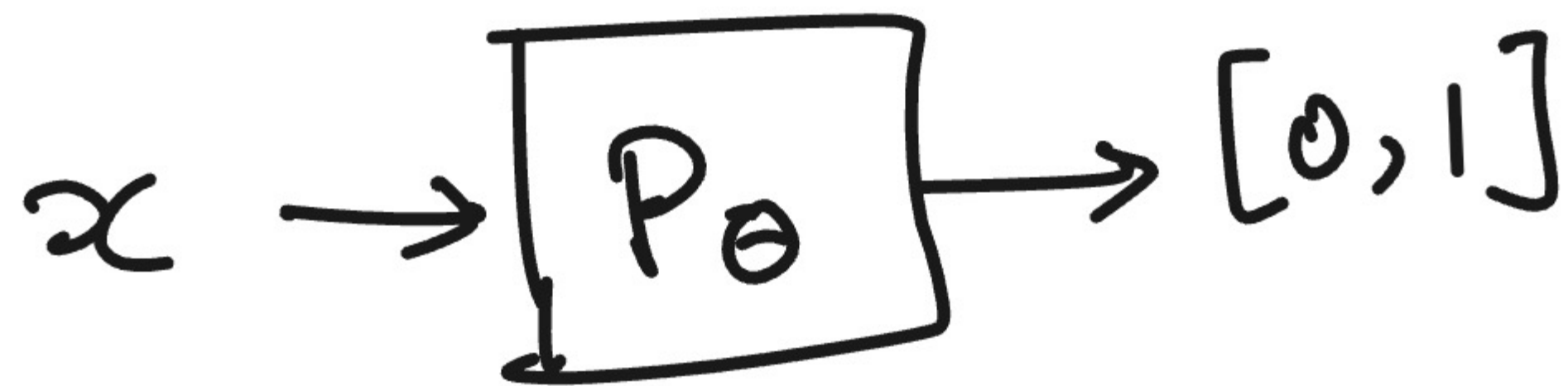
Why is generative model useful?



Under conditions, we can make new instances of x ! (i.e. generate new instances of samples not in the observed data.)

Explicit Density Model vs Implicit Density Model

- We can get the likelihood from model given instance x



- Cannot measure likelihood,
- Generate new "likely" instances directly



Autoregressive LM is an explicit density model

$x = [x_1 \ x_2 \ \dots \ x_L]$ is a sentence,
 x_i are tokens.

What is $P_{\theta}(x) = ?$

$$P(x_1, x_2, \dots, x_L) = P(x_L | x_1, \dots, x_{L-1}) \cdot P(x_1, \dots, x_{L-1})$$



This is cond prob

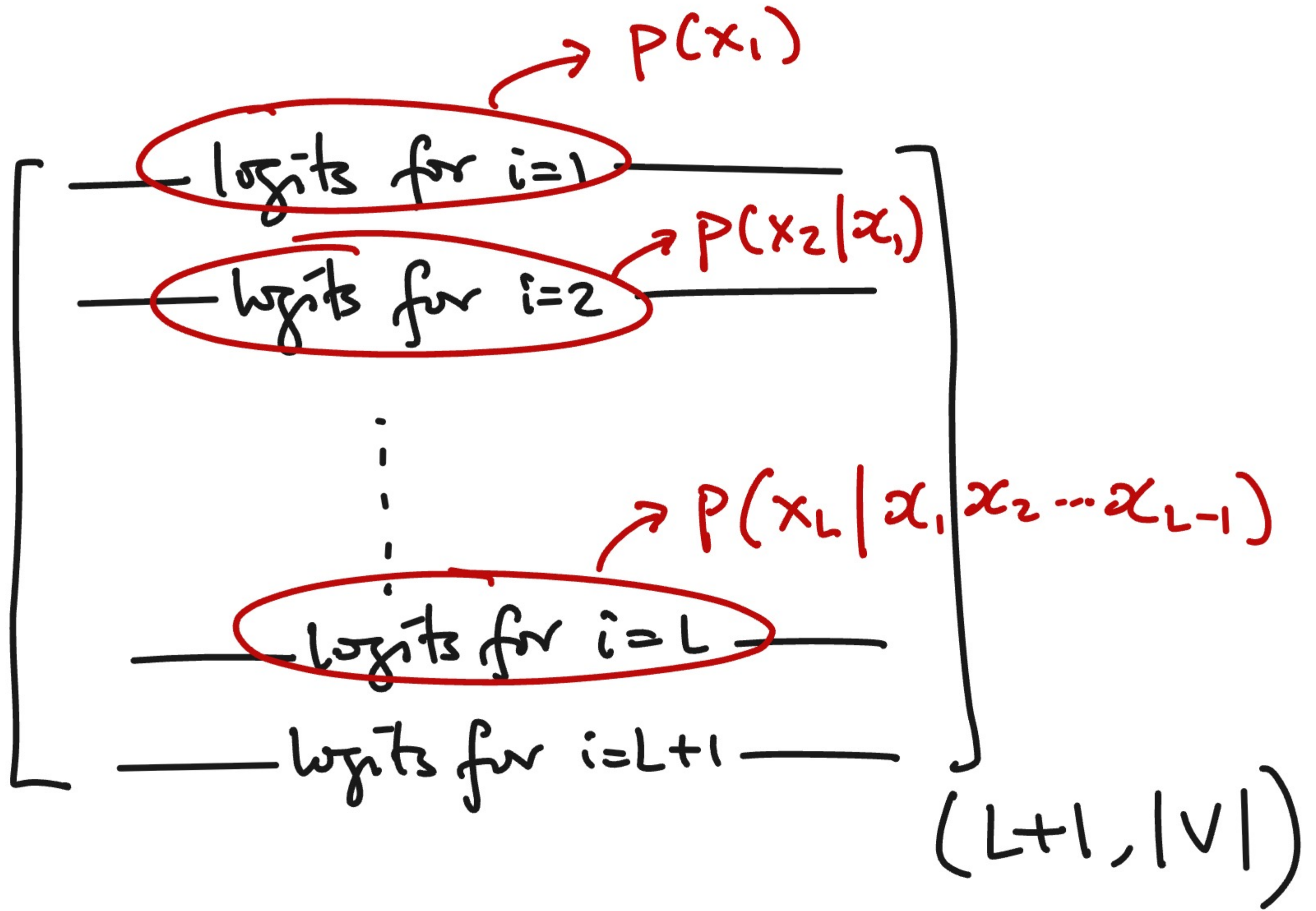
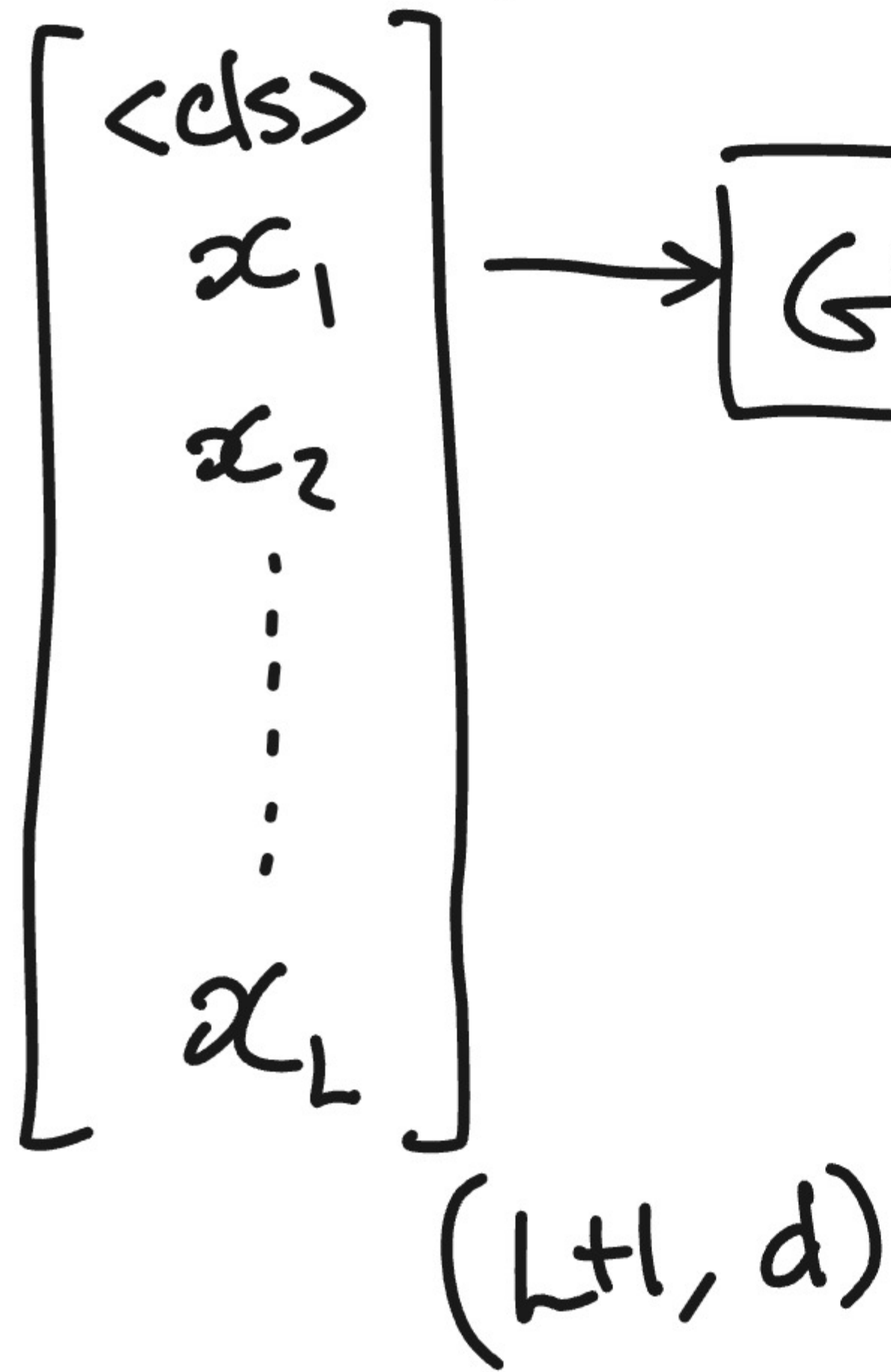
$$P(A, B) = P(A|B) P(B)$$

$$p(x_1, x_2 \dots x_L) = \underbrace{p(x_1, x_2 \dots x_{L-1})}_{\text{next token}} \cdot p(x_L | x_1 \dots x_{L-1})$$

$$= p(x_1 \dots x_{L-2}) \cdot p(x_{L-2} | x_1 \dots x_{L-3}) \cdot p(x_L | x_1 \dots x_{L-1})$$

$$= \prod_{i=1}^L \underbrace{p(x_i | \underbrace{x[1:i-1]}_{\text{already generated tokens}})}_{\text{next token}}$$

Autoregressive



How to generate from explicit density models?

⇒ SAMPLING

Relies on finite discrete vocabulary for each x_i

Since $p(x_{i+1} | x_1, \dots, x_i)$ is known for all $x_{i+1} \in V$,

Sample x_{i+1} by picking most likely next tok.

What about continuous space sampling ⇒ Variational Autoencoder...

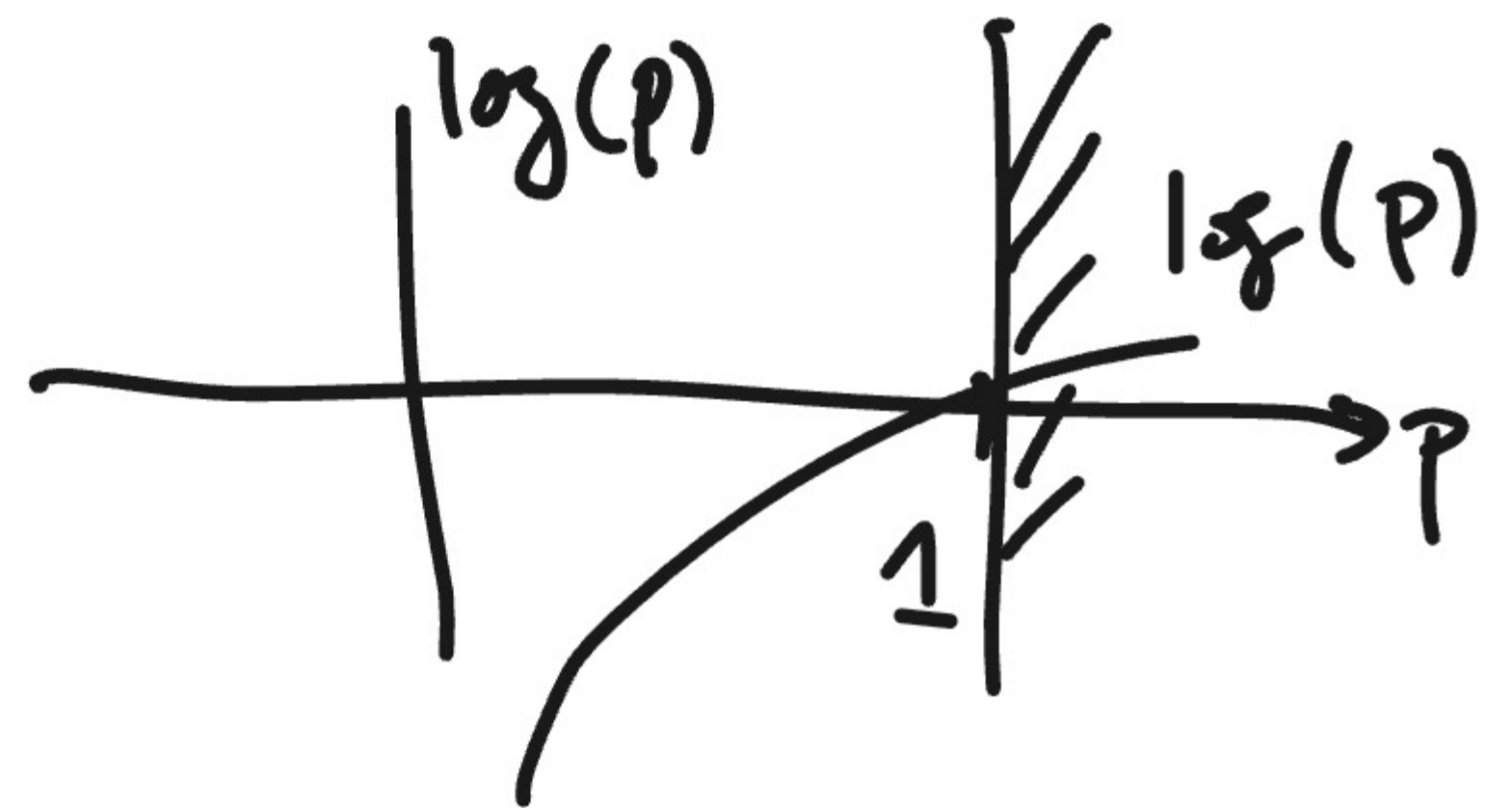
Maximal Likelihood Training of Explicit Model

$\mathcal{D} = \{x^{(1)}, x^{(2)} \dots x^{(N)}\}$ is the training data.

for LM, each $x^{(i)}$ is a sequence of tokens!

P_{θ} is trained so \mathcal{D} are likely samples of P_{θ} .

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P_{\theta}(x^{(i)})$$

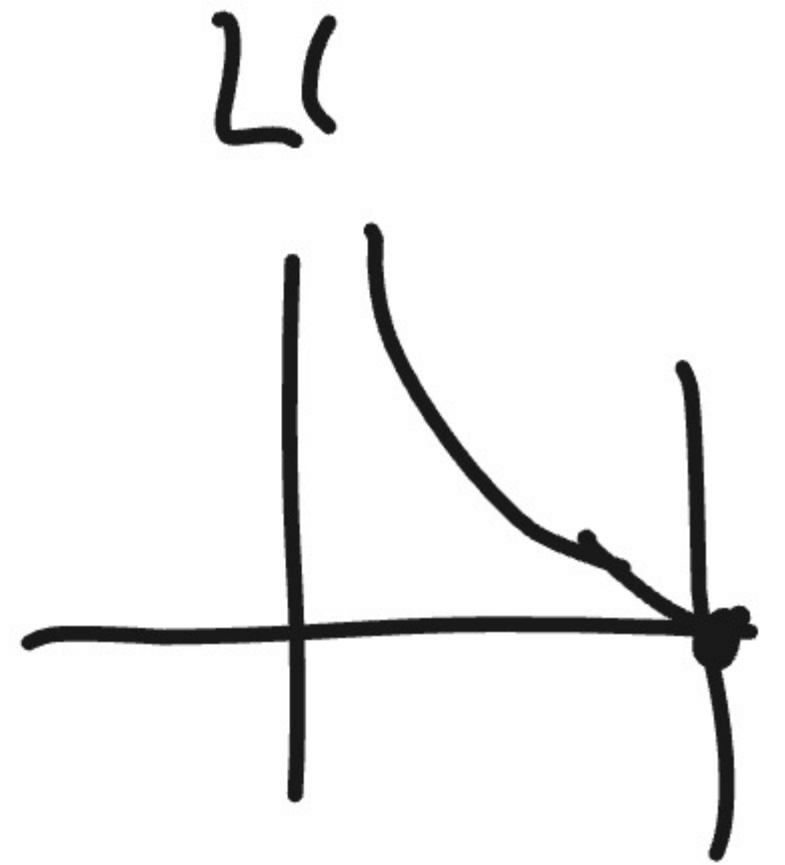


Max. Likelihood Loss Function

$$L_{ML} = -\frac{1}{N} \sum_{i=1}^N \log P_{\theta}(x^{(i)})$$

distance measure
between two
distributions

Connection to P_{data} via KL-divergence



$$L_{ML} = -E \left[\underbrace{\log P_{\theta}(x)}_{\text{observing } \log P_{\theta}(x)} \mid \underbrace{P_{data}(x)}_{\text{sample } x \text{ from } P_{data}} \right]$$

About KL Divergence (a quick look)

$$KL(P_{\text{data}} \parallel P_{\theta}) = E_{x \sim P_{\text{data}}} \left[\log \left(\frac{P_{\text{data}}(x)}{P_{\theta}(x)} \right) \right]$$

$$= E_{x \sim P_{\text{data}}} \left[\log P_{\text{data}}(x) \right] - E_{x \sim P_{\text{data}}} \left[\log P_{\theta}(x) \right]$$

Not dependent

on $\theta \Rightarrow$ constant

$L_{ML}(\theta)$

$$\theta^* = \arg \min_{\theta} KL(P_{\text{data}} \parallel P_{\theta})$$

Latent Variable Models as Explicit Models

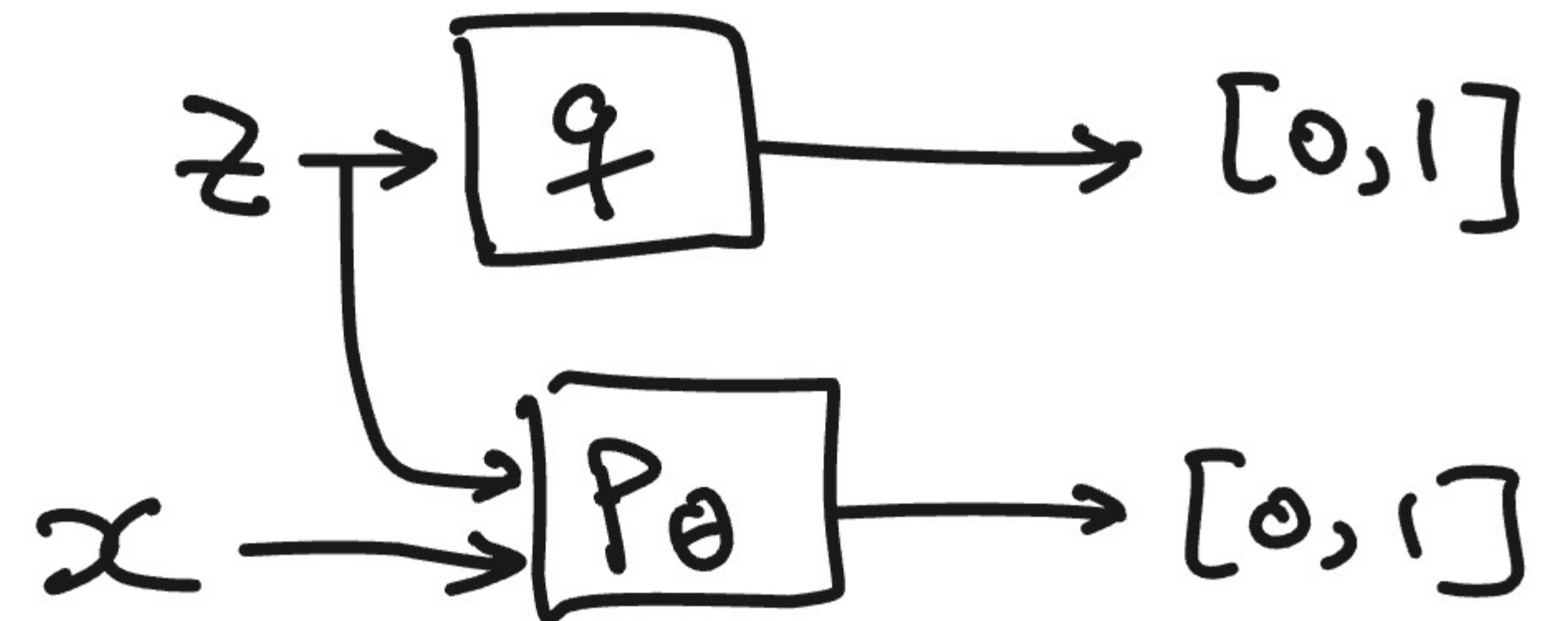
2 stage models

Stage #1 : explicit model on z

$$z \sim q(z)$$

Stage #2 : conditional model

$$x \sim p_{\theta}(x|z)$$



Latent Var. Model

Given x , what is $p(x)$?

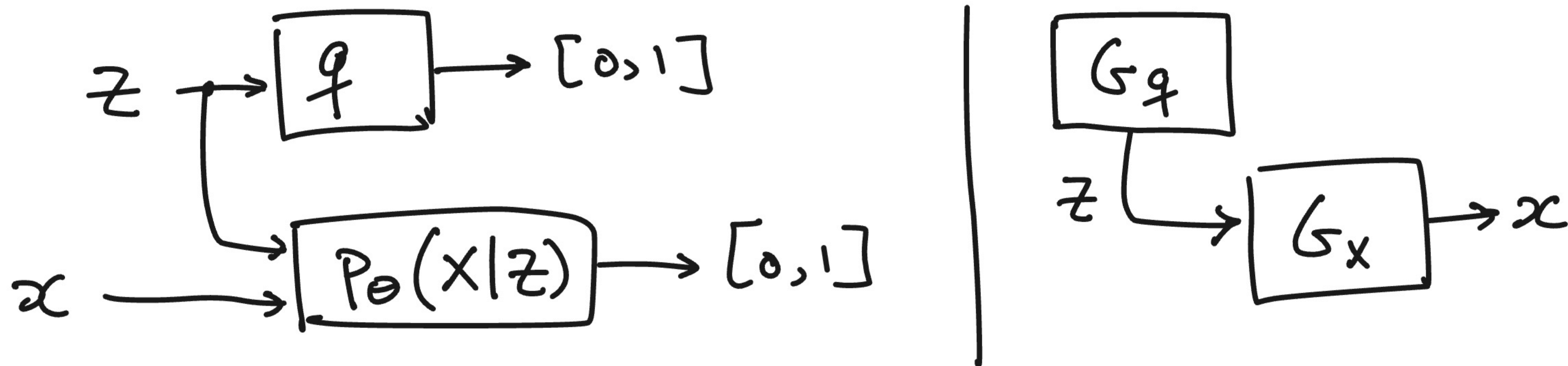
$$p_{\theta}(x) = \int p_{\theta}(x|z) \cdot q(z) dz$$

continuous

$$= \sum_{z \in \mathcal{Z}} p_{\theta}(x|z) q(z)$$

discrete

More on latent variable explicit models



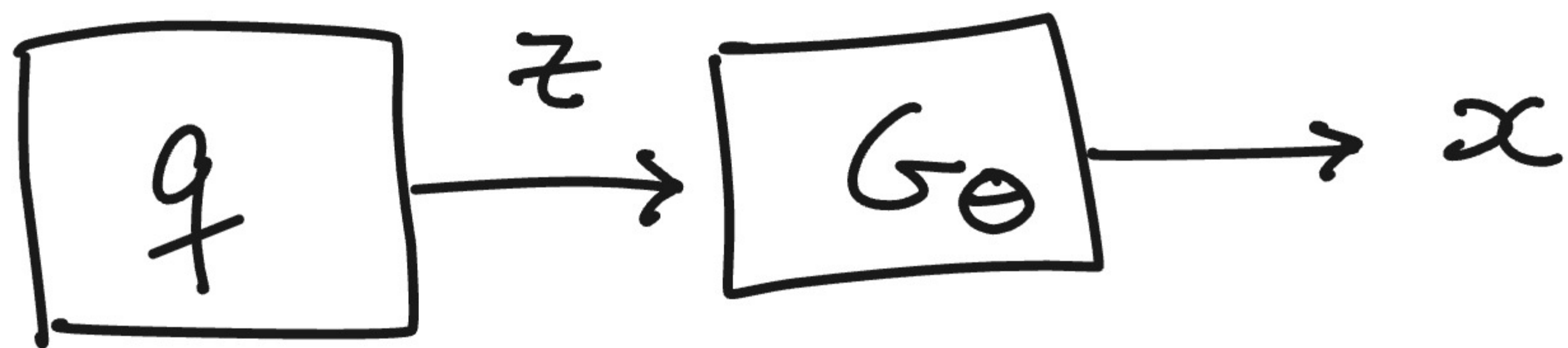
SAMPLING:

$$x \sim P_{\theta}(x)$$

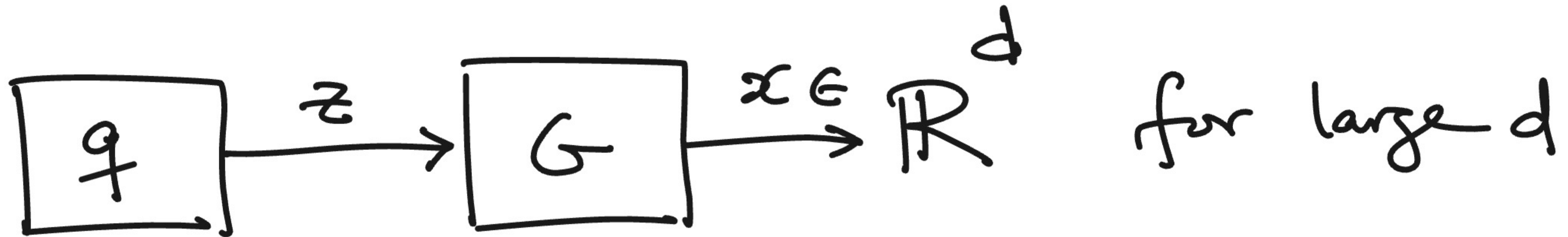
INFERENCE

$$P(z|x)$$

Latent Variable Implicit Models



LATENT VARIABLE z is a control



Simple

Gaussian

$\mathcal{N}(0, 1)$ in \mathbb{R}^k for small k

NEXT WEEK:

Variational Autoencoder ← latent variable
explicit model

Diffusion Model ← latent variable
implicit model

IMAGE GENERATION

