



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 5212
Machine Learning
Lecture 11

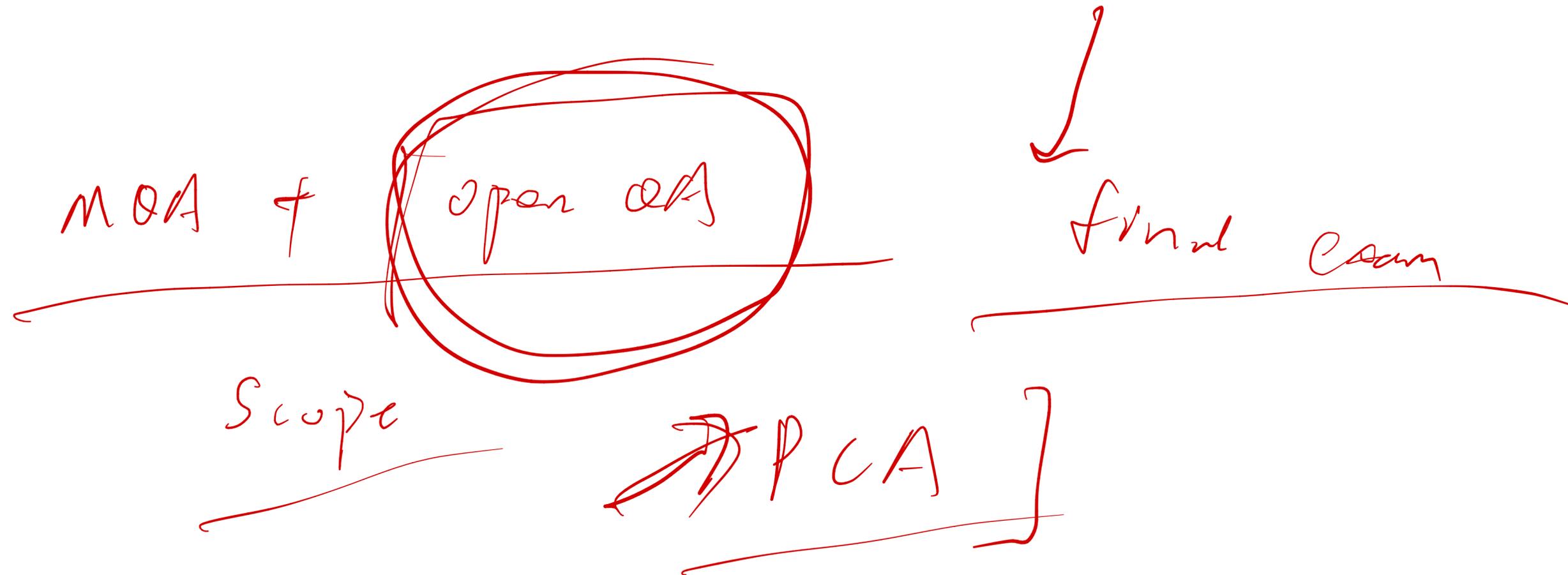
Expectation Maximization

Junxian He
Mar 19, 2026

Midterm Exam

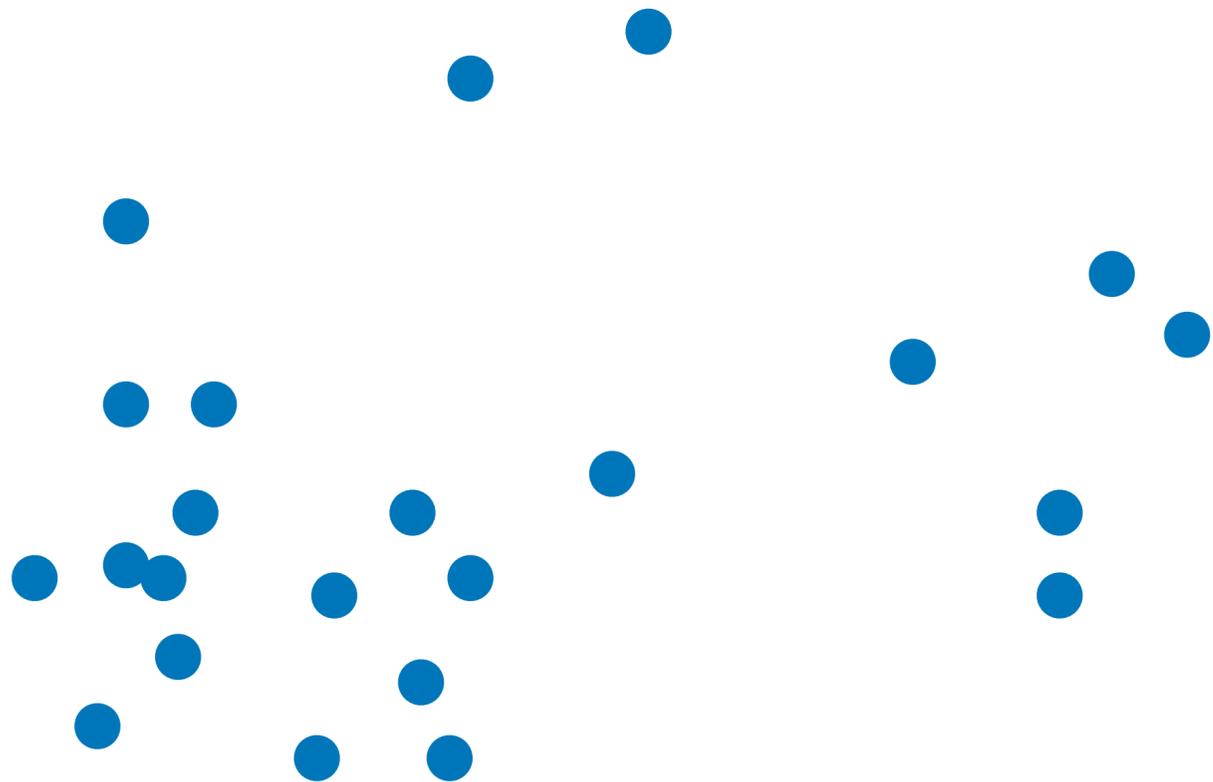
80 min

31/03, in-class, same classroom, 430pm-550pm, one A4-size double-sided cheetsheet is allowed (either printing or handwriting is fine).



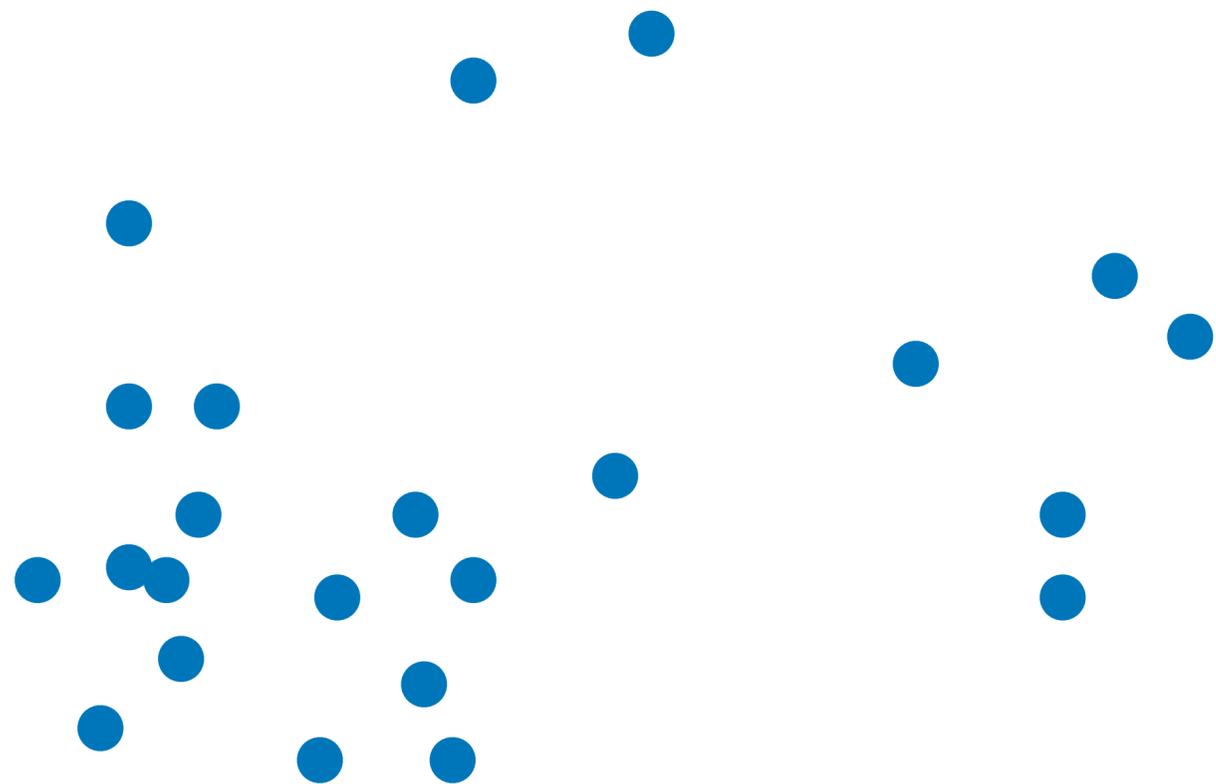
Recap: Generative Models

Recap: Generative Models



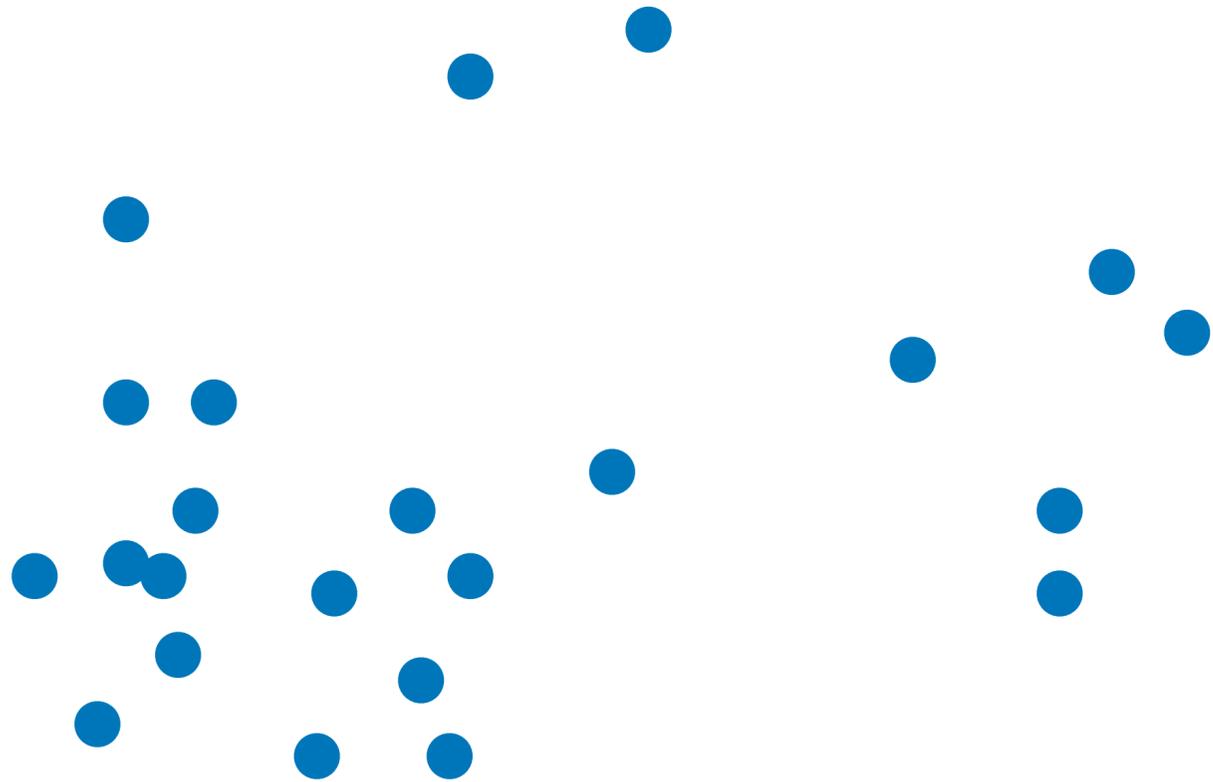
Recap: Generative Models

We want to model $p(x)$



Recap: Generative Models

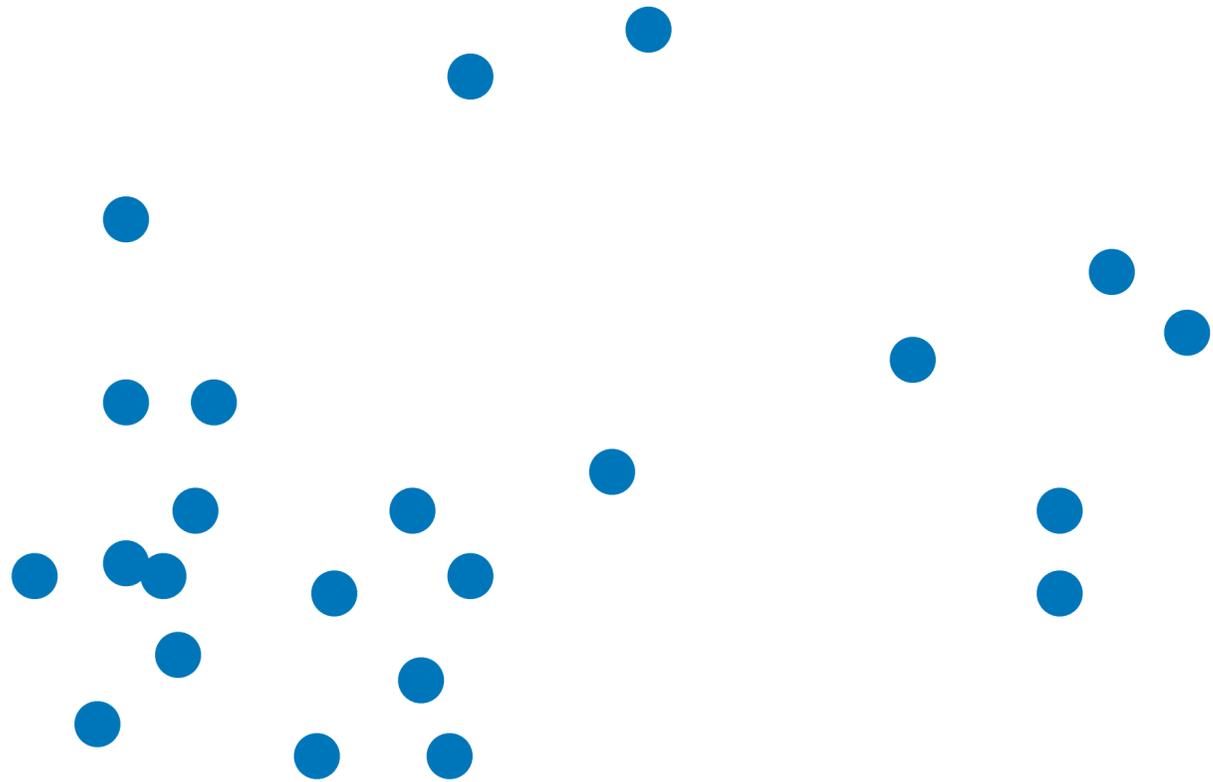
We want to model $p(x)$



In discriminative models, we need to “design” model to make assumption about the function: linear regression, logistic regression, kernel methods

Recap: Generative Models

We want to model $p(x)$



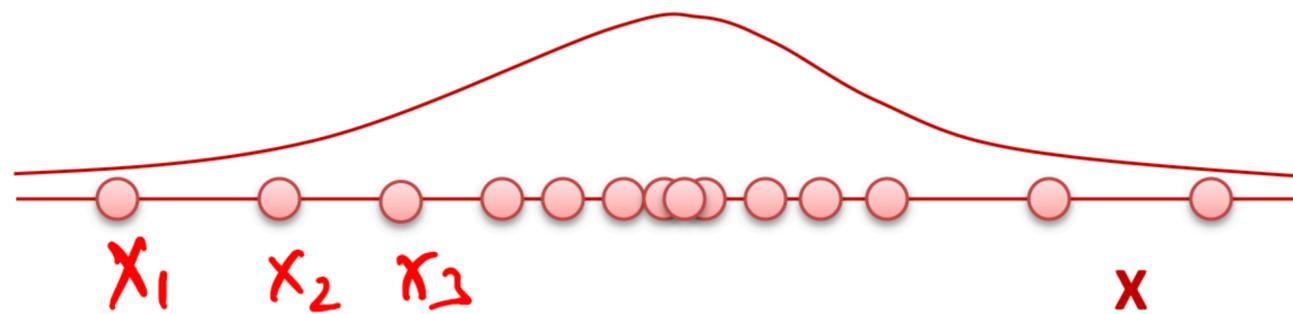
In discriminative models, we need to “design” model to make assumption about the function: linear regression, logistic regression, kernel methods

In generative models, we “design” the model and make assumptions about the data, through defining a distribution family

Recap: Generative Models

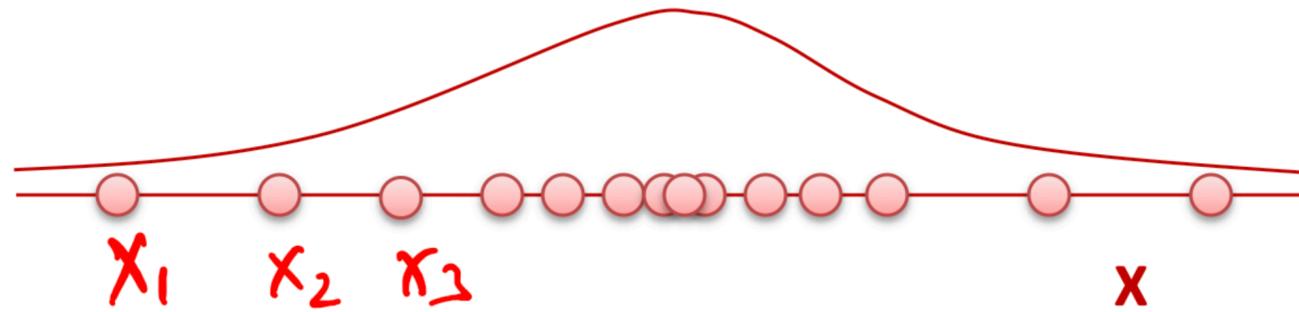
Recap: Generative Models

Data, $D =$



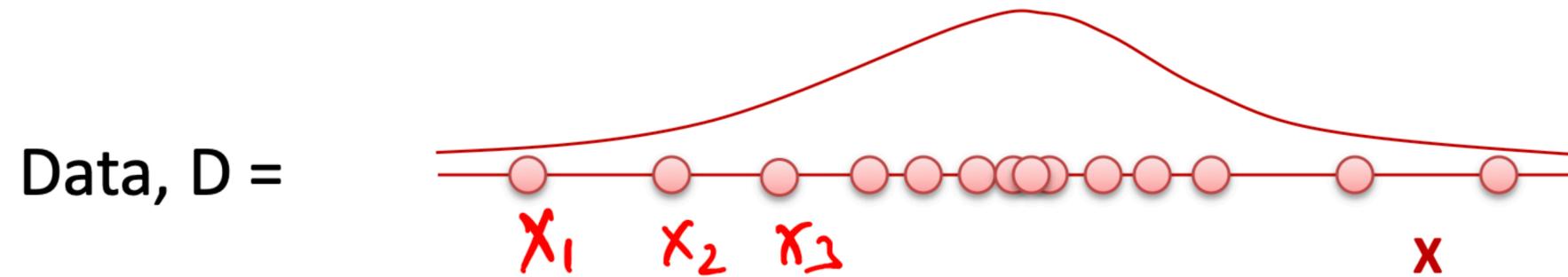
Recap: Generative Models

Data, $D =$



As a simplest case, we directly assume $x \sim N(\mu, \Sigma)$

Recap: Generative Models

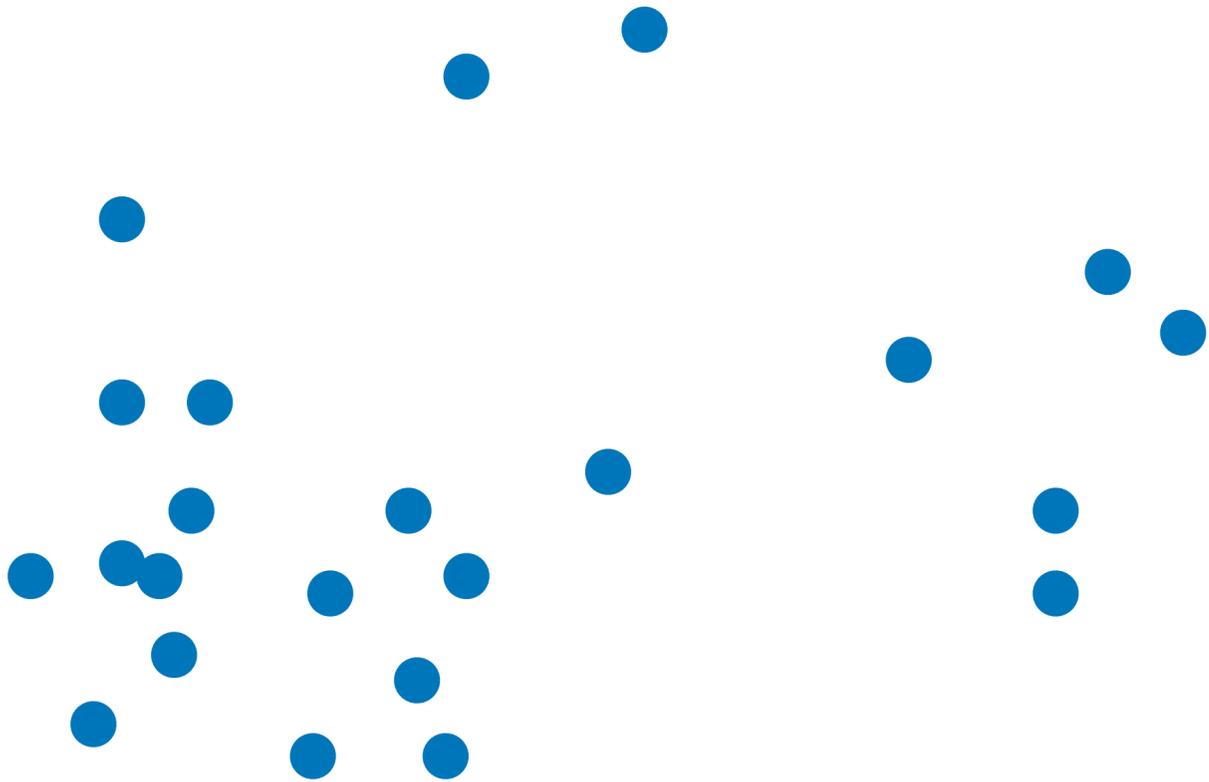


As a simplest case, we directly assume $x \sim N(\mu, \Sigma)$

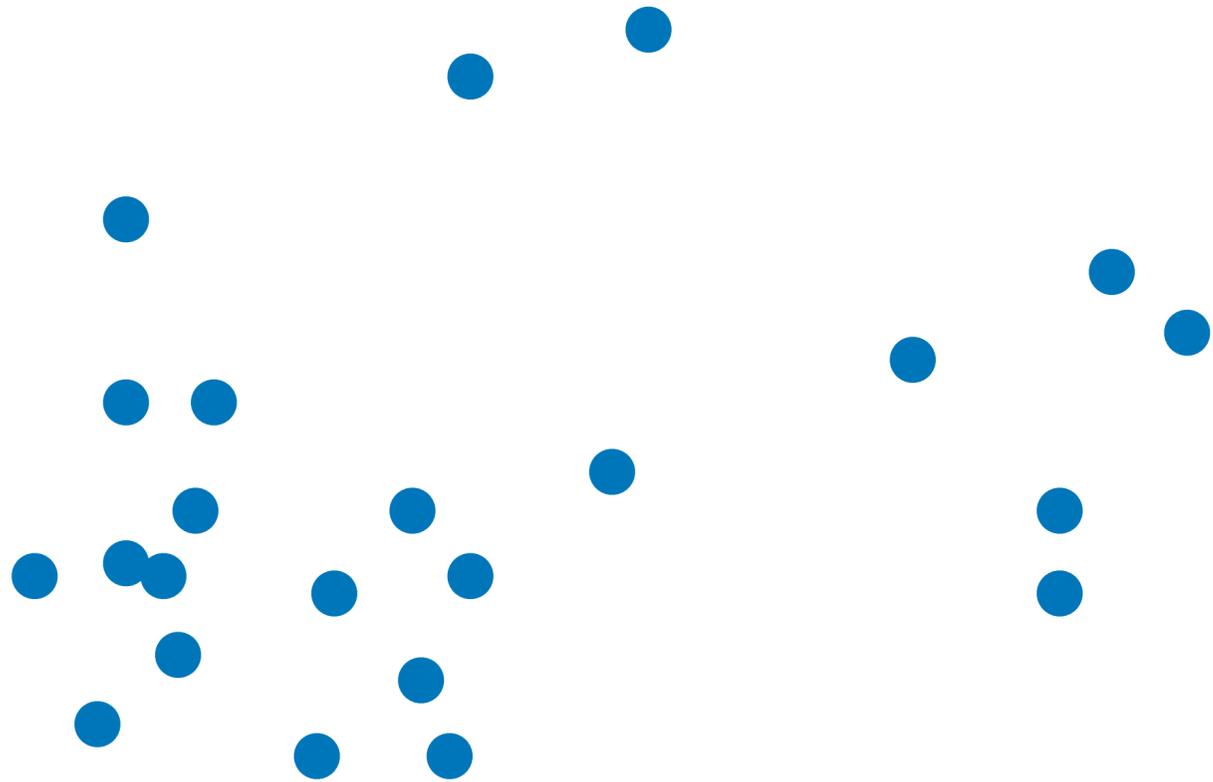
By varying the parameters (μ, Σ) , the model represents different distributions that belong to the Gaussian family

Recap: Generative Models

How to construct more complex distribution family?



Recap: Generative Models

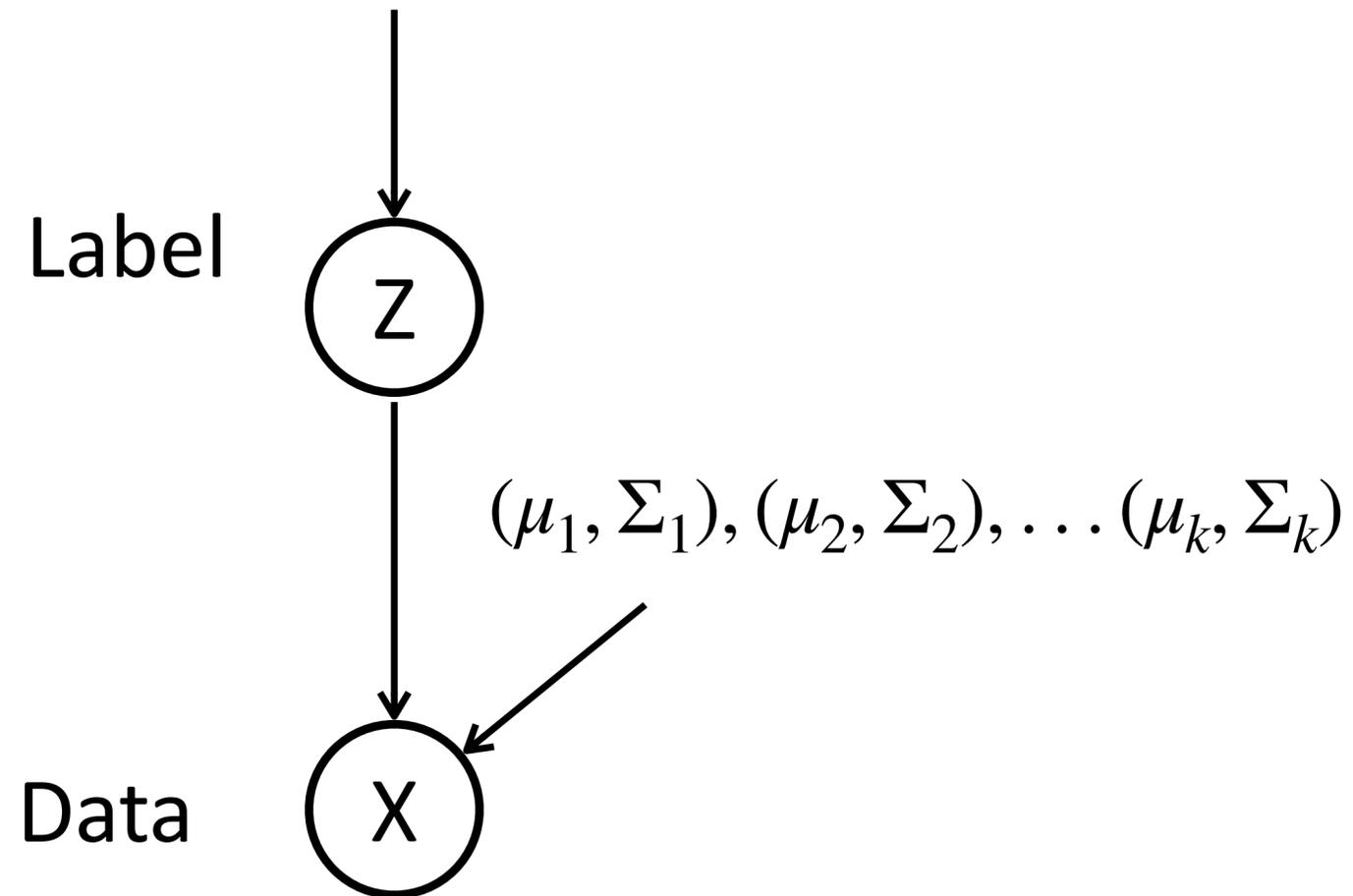


How to construct more complex distribution family?

Introducing more latent variables

Recap: Gaussian Mixture Model

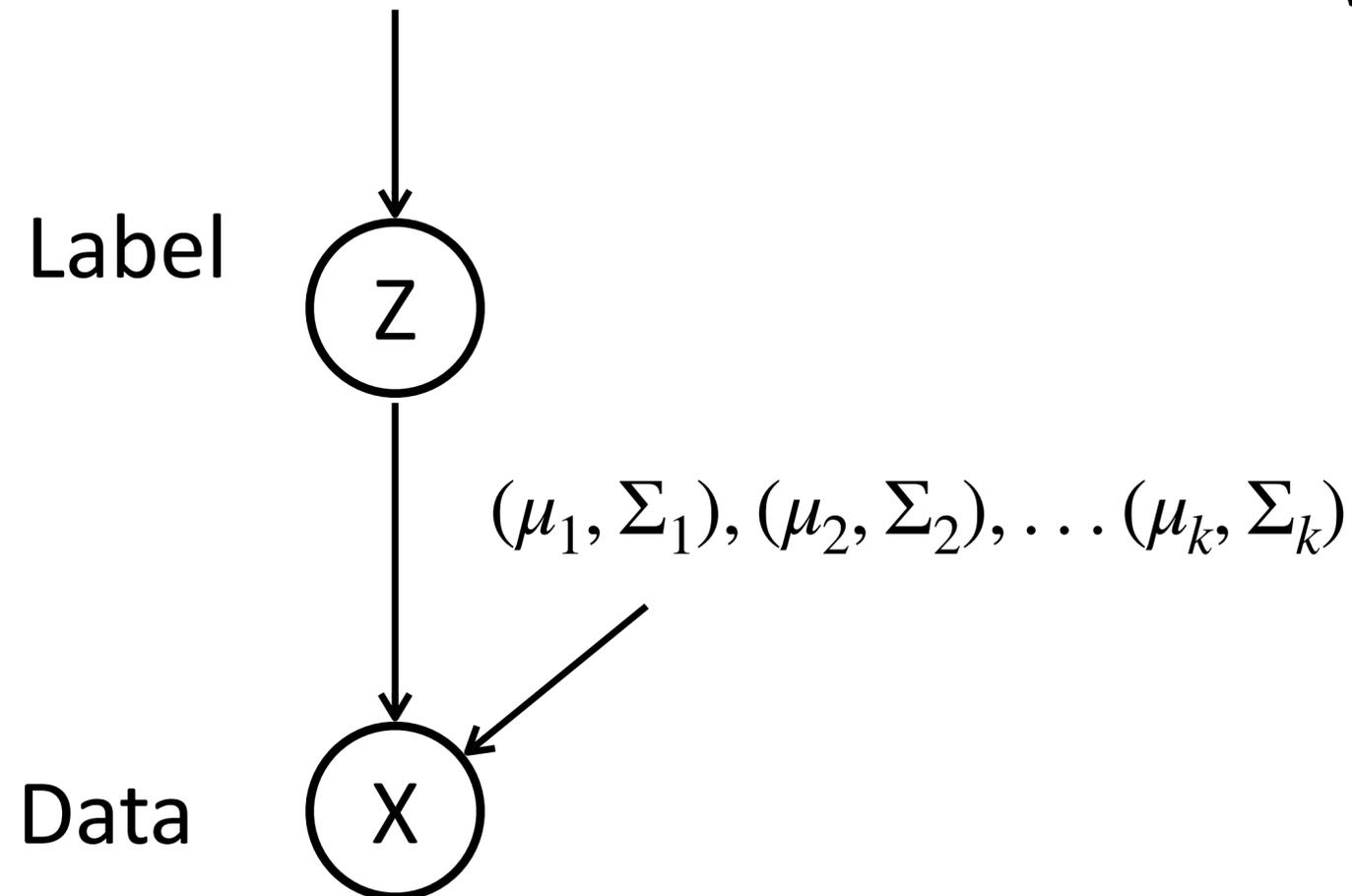
$p(z)$: multinomial, k
classes (e.g. uniform)



Recap: Gaussian Mixture Model

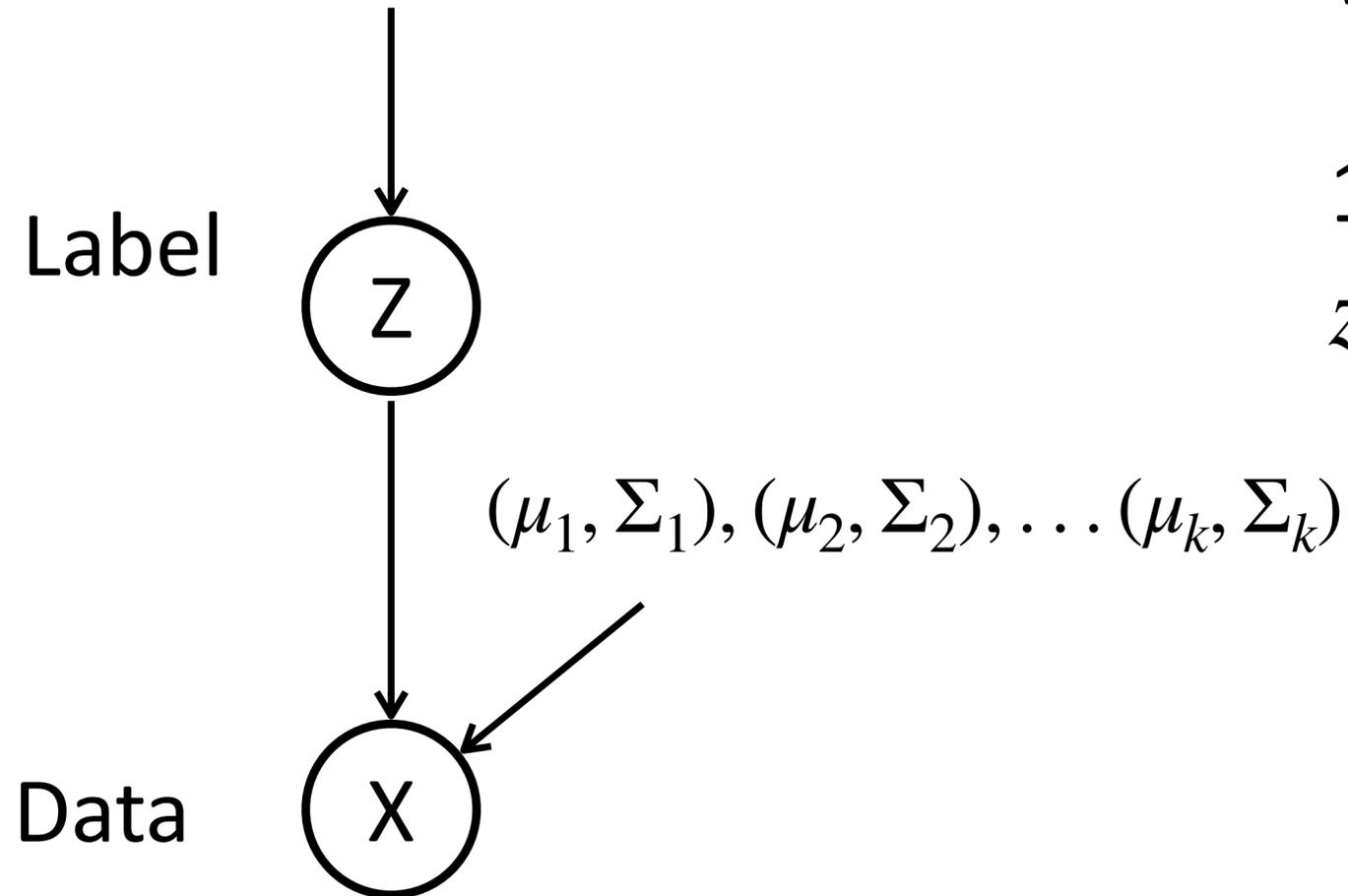
$p(z)$: multinomial, k
classes (e.g. uniform)

We assume the generative process as:



Recap: Gaussian Mixture Model

$p(z)$: multinomial, k
classes (e.g. uniform)

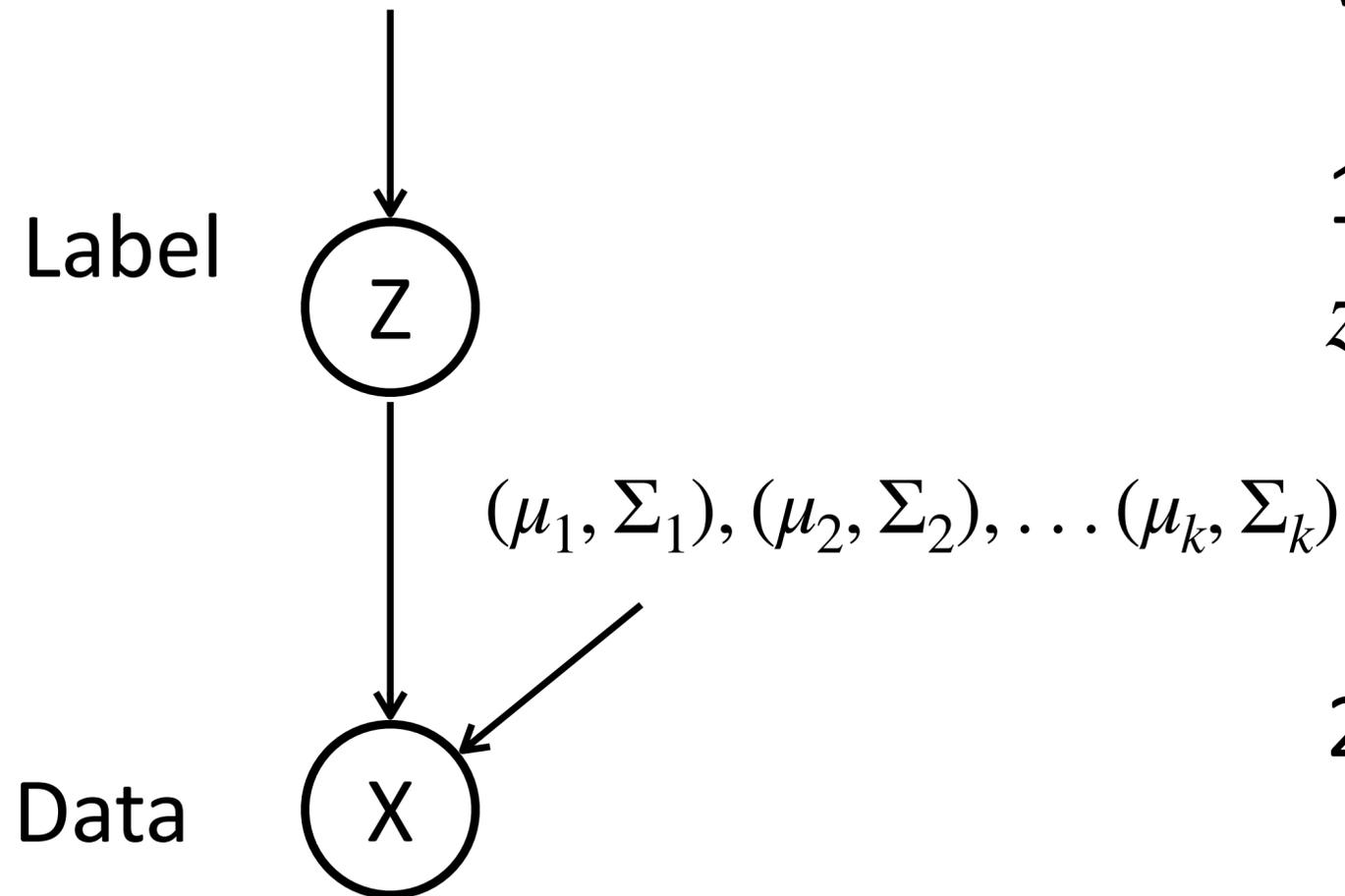


We assume the generative process as:

1. For each data point, sample its label z_i from $p(z)$

Recap: Gaussian Mixture Model

$p(z)$: multinomial, k classes (e.g. uniform)



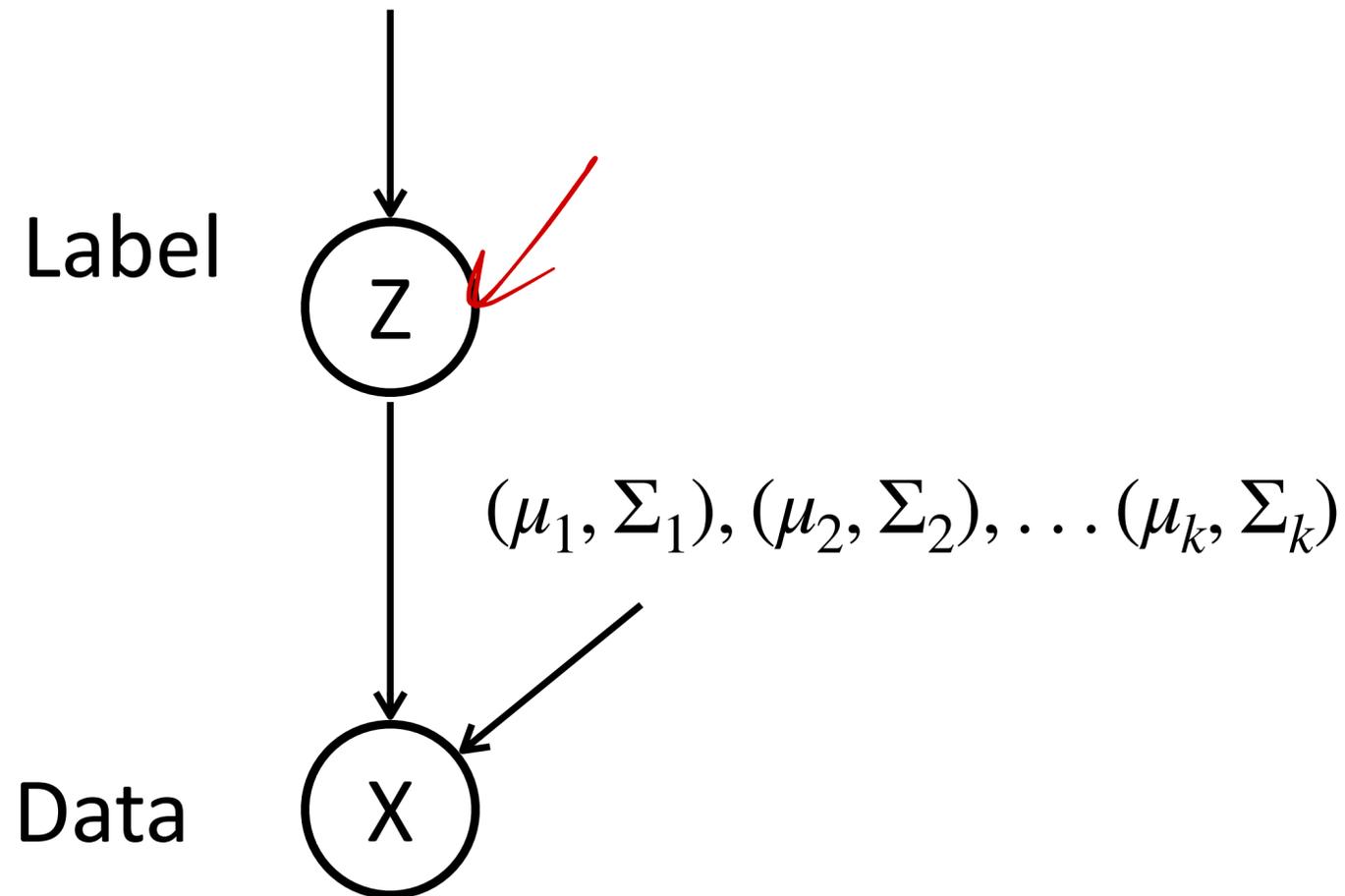
We assume the generative process as:

1. For each data point, sample its label z_i from $p(z)$

2. Sample $x_i \sim N(\mu_{z_i}, \Sigma_{z_i})$

Recap: MLE for GMM

$p(z)$: multinomial, k classes (e.g. uniform)



Unsupervised:

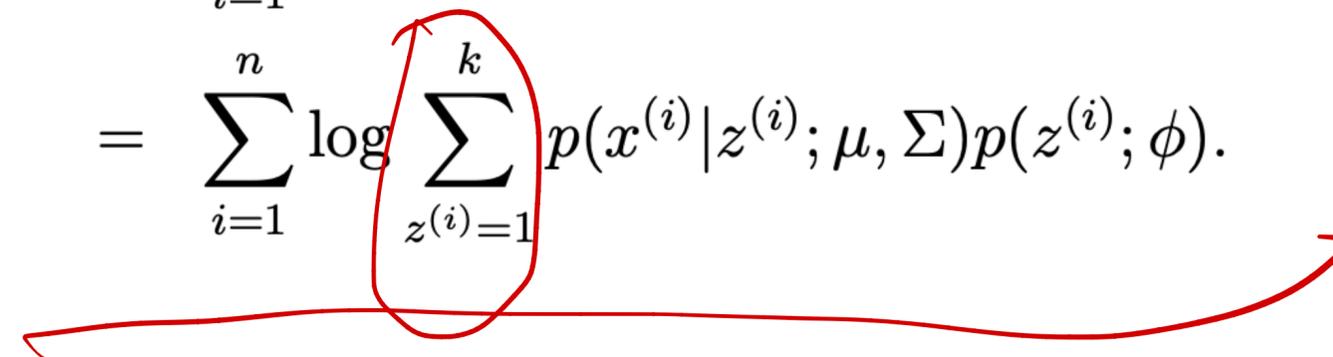
$$\operatorname{argmax}_{\phi, \mu, \Sigma} \log p(x)$$

How to compute this?

$\log p(x, \tilde{z})$ supervised learning

$\log p(x)$

Recap: MLE for GMM

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$


Recap: MLE for GMM

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

if z is continuous,

1. Intractable (no closed-form for the solution)

Recap: MLE for GMM

$$\begin{aligned} \ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi). \end{aligned}$$

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

$z^{(i)} \sim p(z; \phi)$
 $p(x | z; \mu, \Sigma)$
 $\log \sum_{z=1}^k p(x | z; \mu, \Sigma) p(z; \phi)$
 $p(x | z; \mu, \Sigma)$

Recap: MLE for GMM

$$\begin{aligned}\ell(\phi, \mu, \Sigma) &= \sum_{i=1}^n \log p(x^{(i)}; \phi, \mu, \Sigma) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi).\end{aligned}$$

1. Intractable (no closed-form for the solution)
2. Large variance in gradient descent

Expectation Maximization is to address the MLE optimization problem

Things are easy when we know z .

In case we know z

Things are easy when we know z ..

In case we know z

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^n \log p(x^{(i)} | z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi).$$

$\log p(x, z)$



Things are easy when we know z .

In case we know z

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^n \log p(x^{(i)} | z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi).$$

$$\phi_j = \frac{1}{n} \sum_{i=1}^n 1\{z^{(i)} = j\},$$

$$\mu_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^n 1\{z^{(i)} = j\}},$$

$$\Sigma_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n 1\{z^{(i)} = j\}}.$$

Counting GDA

Things are easy when we know z .

In case we know z

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^n \log p(x^{(i)} | z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi).$$

$$\phi_j = \frac{1}{n} \sum_{i=1}^n 1\{z^{(i)} = j\},$$

$$\mu_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^n 1\{z^{(i)} = j\}},$$

$$\Sigma_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n 1\{z^{(i)} = j\}}.$$

Expectation maximization is to infer the latent variables first (z here), and maximize the likelihood given the inferred z

Expectation Maximization for GMM

Repeat until convergence:

{

}

Expectation Maximization for GMM

Repeat until convergence:

{

(E-step) For each i, j , set

$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

infer z

posterior distribution

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

}

Expectation Maximization for GMM

Repeat until convergence:

{

(E-step) For each i, j , set

$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

Compute the posterior distribution,
given current parameters

}

Expectation Maximization for GMM

Repeat until convergence:

{

No parameter change in E-step

(E-step) For each i, j , set

$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

Compute the posterior distribution,
given current parameters

}

Expectation Maximization for GMM

Repeat until convergence:

No parameter change in E-step

(E-step) For each i, j , set

$$w_j^{(i)} = p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

Compute the posterior distribution, given current parameters

(M-step) Update the parameters:

$$\phi_j := \frac{1}{n} \sum_{i=1}^n w_j^{(i)},$$

$$\mu_j := \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}},$$

$$\Sigma_j := \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$

$\log p(x, z)$

Expectation Maximization

- Why does it work?
- What is its relation to MLE estimation? *log P(x)*
- How is convergence guaranteed?
- When we perform EM, what is the real objective that we are optimizing?

General EM Algorithm

General EM Algorithm

$$p(x; \theta) = \sum_z p(x, z; \theta)$$

General EM Algorithm

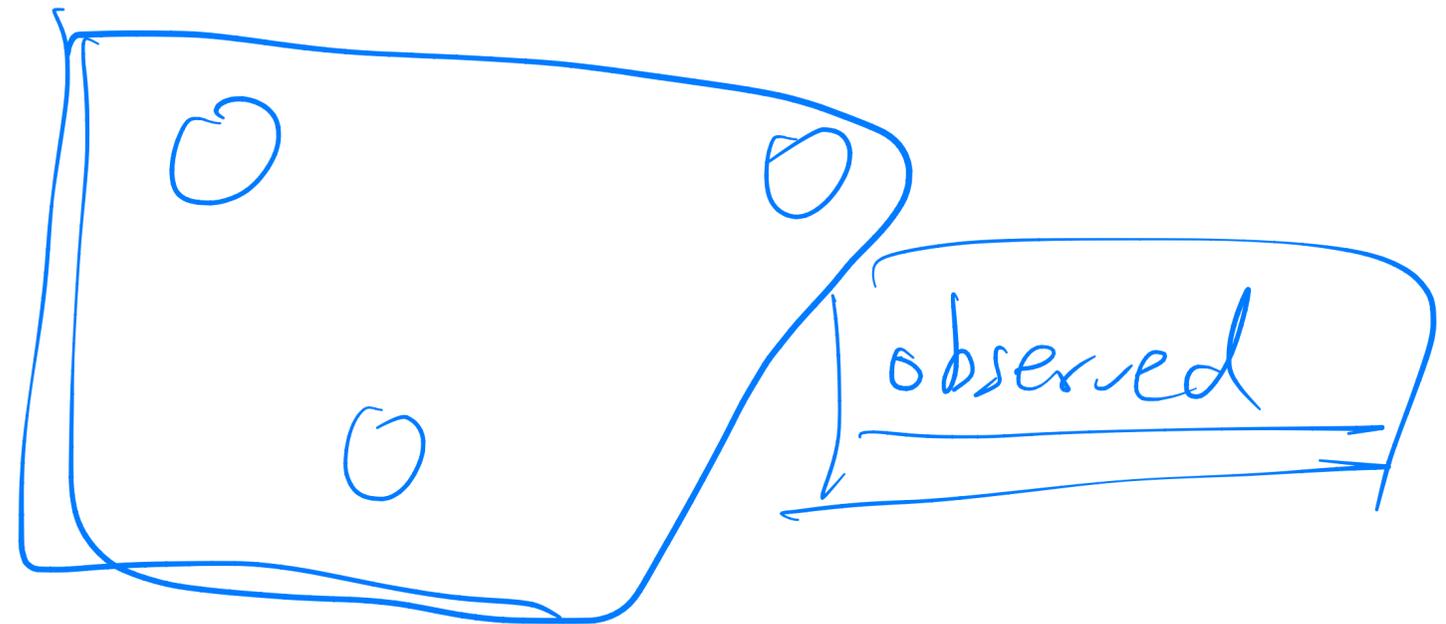
$$p(x; \theta) = \sum_z p(x, z; \theta)$$

z

$$\vec{z} = \{z_1, z_2, \dots, z_n\}$$

$$\ell(\theta) = \sum_{i=1}^n \log p(x^{(i)}; \theta)$$

$$= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta).$$



General EM Algorithm

$$p(x; \theta) = \sum_z p(x, z; \theta)$$

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^n \log p(x^{(i)}; \theta) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta). \end{aligned}$$

Let Q to be a distribution over z

General EM Algorithm

$$p(x; \theta) = \sum_z p(x, z; \theta)$$

$\max \log p(x)$

$$\ell(\theta) = \sum_{i=1}^n \log p(x^{(i)}; \theta)$$

lower bound $\leq \log p(x)$

$$= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta).$$

$\log p(x)$

Let Q to be a distribution over z

$Q(z)$

$\log p(x; \theta)$

$$= \log \sum_z p(x, z; \theta)$$

$$= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)}$$

$$\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

Jensen inequality

KL 7/0

General EM Algorithm

$$p(x; \theta) = \sum_z p(x, z; \theta)$$

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^n \log p(x^{(i)}; \theta) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta). \end{aligned}$$

Let Q to be a distribution over z

This lower bound holds for any $Q(z)$

$$\begin{aligned} \log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} \end{aligned}$$

General EM Algorithm

$$p(x; \theta) = \sum_z p(x, z; \theta)$$

$$\ell(\theta) = \sum_{i=1}^n \log p(x^{(i)}; \theta)$$

$$= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta).$$

Let Q to be a distribution over z

$$\log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

$$= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)}$$

$$\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

Jensen inequality

This lower bound holds for any $Q(z)$

$$\mathbb{E}_{z \sim Q(z)} \left[\frac{p(x, z; \theta)}{Q(z)} \right]$$

log concave

Jensen Inequality

For a convex function f , and $t \in [0,1]$

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

Jensen Inequality

For a convex function f , and $t \in [0,1]$

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

$$P(x_1) + P(x_2) \dots = 1$$

In probability:

$$f(E[x]) = f(P(x_1)x_1 + P(x_2)x_2 + \dots)$$

$$f(E[X]) \leq \cancel{E[f(X)]} \quad E[f(x)]$$

Jensen Inequality

For a convex function f , and $t \in [0,1]$

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

In probability:

x is constant $f(\mathbb{E}[x]) = \mathbb{E}[f(x)]$

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$$

If f is strictly convex, then equality holds only when X is a constant

Evidence Lower Bound (ELBO)

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}\end{aligned}$$

↓ ELBO

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} && \text{ELBO} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}\end{aligned}$$

ELBO

$Q(z)$
uniform

$Q(z)$
 $N(0, 1)$

Because the log likelihood is intractable, people often optimize its lower bound instead

Evidence Lower Bound (ELBO)

$$\begin{aligned} \log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} && \text{ELBO} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} \end{aligned}$$

tight lower bound

Because the log likelihood is intractable, people often optimize its lower bound instead

Why optimizing lower bound works? How to choose $Q(z)$, why we computed posterior in the E step, what is the benefit?

log P(x)

$Q_2(z)$

ELBO

$f(x) = -\infty$

gap?

$Q_1(z)$ ELBO

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

$$= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)}$$

Jensen

Inequality

$$\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

When is the lower bound tight?

$$\log p(x) = \text{ELBO}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}\end{aligned}$$

$E[]$

When is the lower bound tight?

$$\frac{p(x, z; \theta)}{Q(z)} = c$$

Evidence Lower Bound (ELBO)

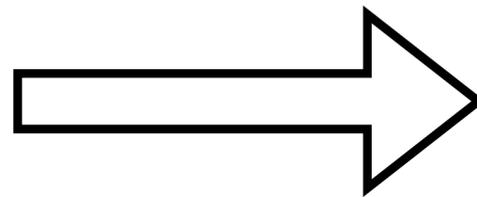
$$\begin{aligned} \log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z Q(z) \frac{p(x, z; \theta)}{Q(z)} \\ &\geq \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} \end{aligned}$$

$\int_z Q(z) = 1$

$Q(z)$

When is the lower bound tight?

$$\frac{p(x, z; \theta)}{Q(z)} = c$$



$$\begin{aligned} Q(z) &= \frac{p(x, z; \theta)}{\sum_z p(x, z; \theta)} \\ &= \frac{p(x, z; \theta)}{p(x; \theta)} \\ &= p(z|x; \theta) \end{aligned}$$

$Q(z) = \frac{p(x, z; \theta)}{c}$

Evidence Lower Bound (ELBO)

Evidence Lower Bound (ELBO)

Verify $\sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$ when $Q(z) = p(z|x)$?

$\approx \log P(x)$

$$\sum_z P(z|x) \log \left[\frac{P(x, z; \theta)}{P(z|x)} \right]$$

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

$$= \sum_z P(z|x) \log P(x)$$

$$= \log P(x) \cdot \left(\sum_z P(z|x) \right) = \log P(x) \geq \text{ELBO}$$

Evidence Lower Bound (ELBO)

Verify $\sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$ when $Q(z) = p(z|x)$?

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

Evidence Lower Bound (ELBO)

Verify $\sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$ when $Q(z) = p(z|x)$?

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

$$\forall Q, \theta, x, \quad \log p(x; \theta) \geq \text{ELBO}(x; Q, \theta)$$

Evidence Lower Bound (ELBO)

Verify $\sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$ when $Q(z) = p(z|x)$?

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

$$\forall Q, \theta, x, \quad \log p(x; \theta) \geq \text{ELBO}(x; Q, \theta)$$

For a dataset of many data samples

$$\begin{aligned} \ell(\theta) &\geq \sum_i \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$?

$P(z|x)$

$Q^*(z) = P(z|x)$

$\text{ELBO} \leq \log P(x)$

The General EM Algorithm

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \theta).$$

(M-step) Set

$$\theta := \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

}

no parameter change

parameter

$$\mathbb{E}_{z \sim Q(z)} \left[\log \frac{p(x, z)}{Q(z)} \right]$$

$$z \sim Q(z) \left[\log \frac{p(x, z)}{Q(z)} \right]$$

The General EM Algorithm

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

Based on current θ , model parameters does not change in E-step

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

}

The General EM Algorithm

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

Based on current θ , model parameters does not change in E-step

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

$Q(z)$

$Q(z)$ is not relevant to θ , and $Q(z)$ does not change in the M-step

$Q(z)$

}

The General EM Algorithm

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

Based on current θ , model parameters does not change in E-step

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

$Q(z)$ is not relevant to θ , and $Q(z)$ does not change in the M-step

}

E-step is maximizing ELBO over $Q(z)$, M-step is maximizing ELBO over θ

$\tilde{E}M$:

repeat:

\tilde{E} step:

$\max \tilde{E}LBO$

$Q(z)$

M step:

$\max \tilde{E}LBO$

θ

The General EM Algorithm

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

Based on current θ , model parameters does not change in E-step

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

$Q(z)$ is not relevant to θ , and $Q(z)$ does not change in the M-step

}

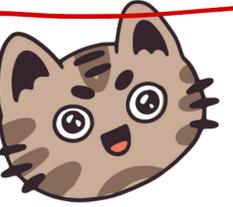
E-step is maximizing ELBO over $Q(z)$, M-step is maximizing ELBO over θ

Why is maximizing lower-bound sufficient?

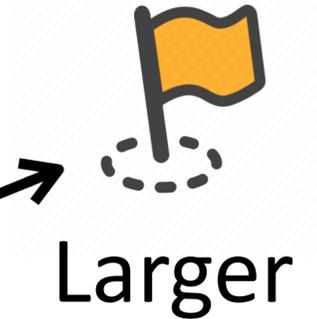
EM is Hill Climbing



$\log p(x; \theta)$



ELBO



EM is Hill Climbing



$\log p(x; \theta)$

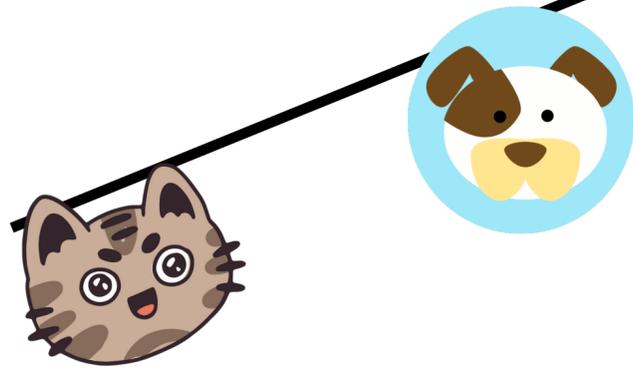
Only related to θ , no z



Larger



ELBO



EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

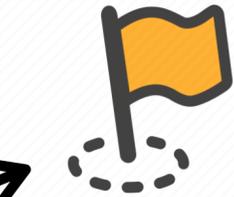
EM is Hill Climbing



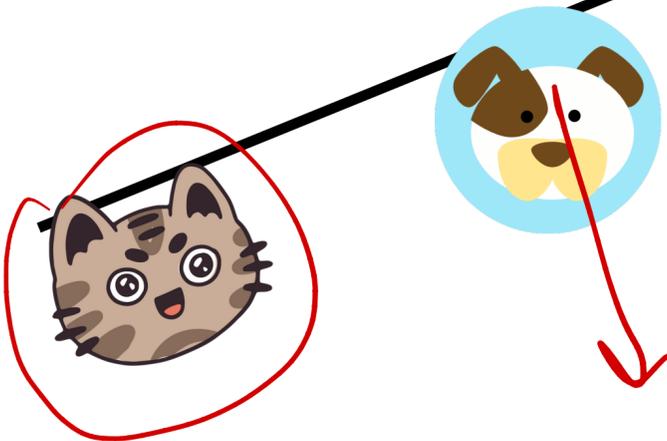
$\log p(x; \theta)$



ELBO



Larger



E-step: $Q(z) = p(z | x; \theta)$, making ELBO tight

not change

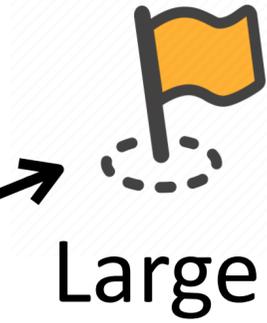
EM is Hill Climbing



$\log p(x; \theta)$



ELBO



E-step: $Q(z) = p(z | x; \theta)$, making ELBO tight
“dog” doesn’t change, because θ does not change

EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

E-step: $Q(z) = p(z | x; \theta)$, making ELBO tight
“dog” doesn’t change, because θ does not change

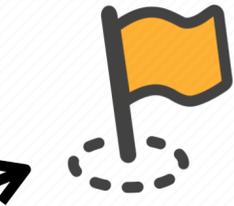
EM is Hill Climbing



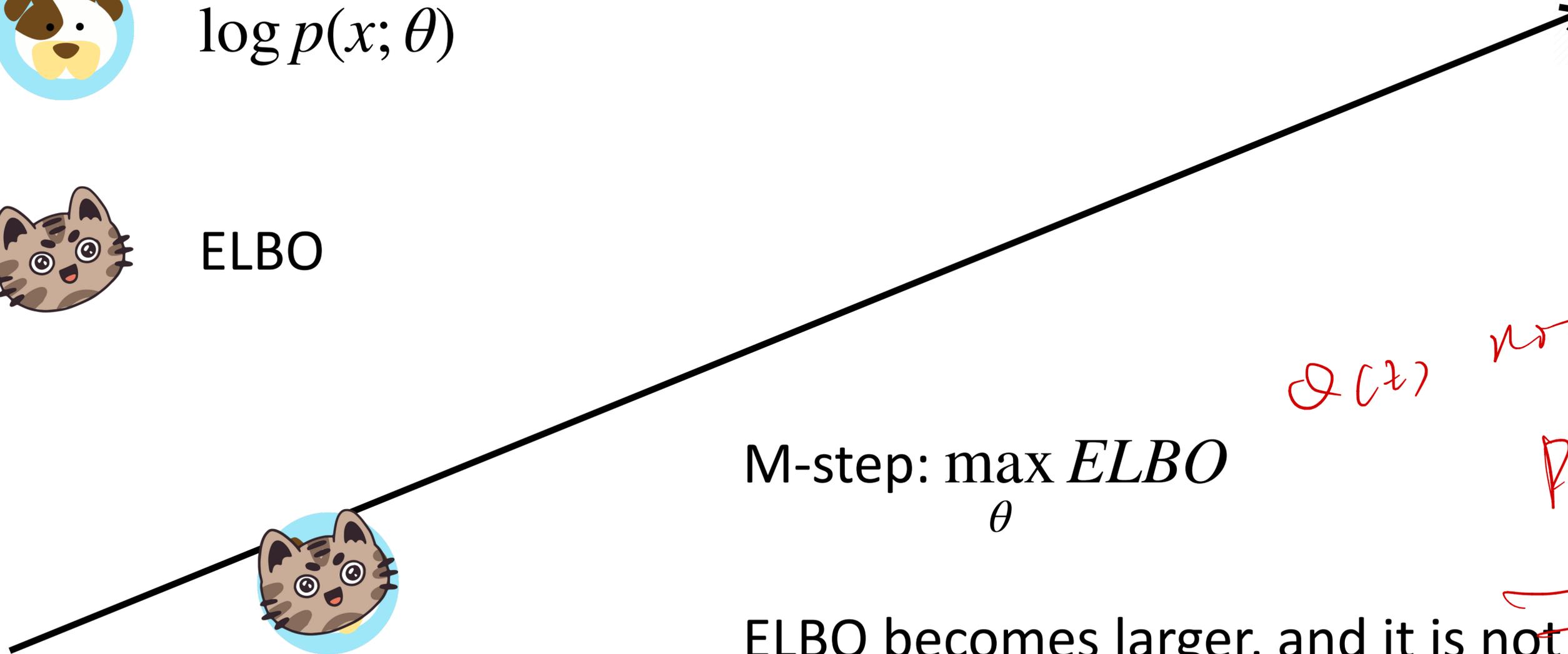
$\log p(x; \theta)$



ELBO



Larger



M-step: $\max_{\theta} ELBO$

Q(z) not change

P(z|x; theta)

ELBO becomes larger, and it is not tight anymore because posterior changes

EM is Hill Climbing



$\log p(x; \theta)$



ELBO



M-step: $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes



Larger

EM is Hill Climbing



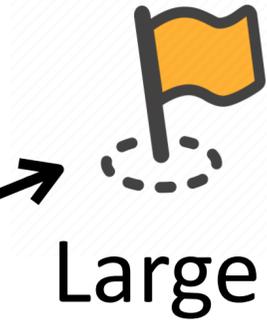
$\log p(x; \theta)$



ELBO



M-step: $\max_{\theta} ELBO$



Larger

ELBO becomes larger, and it is not tight anymore because posterior changes

EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

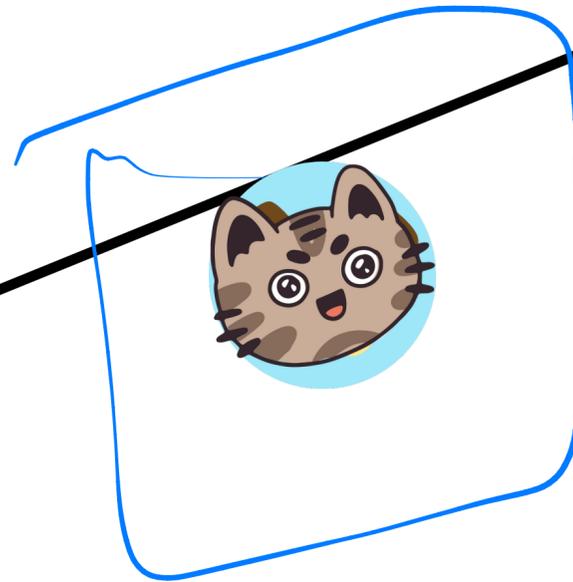
EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

E-step: $Q(z) = p(z | x; \theta)$, making ELBO tight
“dog” doesn’t change, because θ does not change

EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

M-step: $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes

EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

M-step: $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes

$\log p(x; \theta)$ is monotonically increasing!

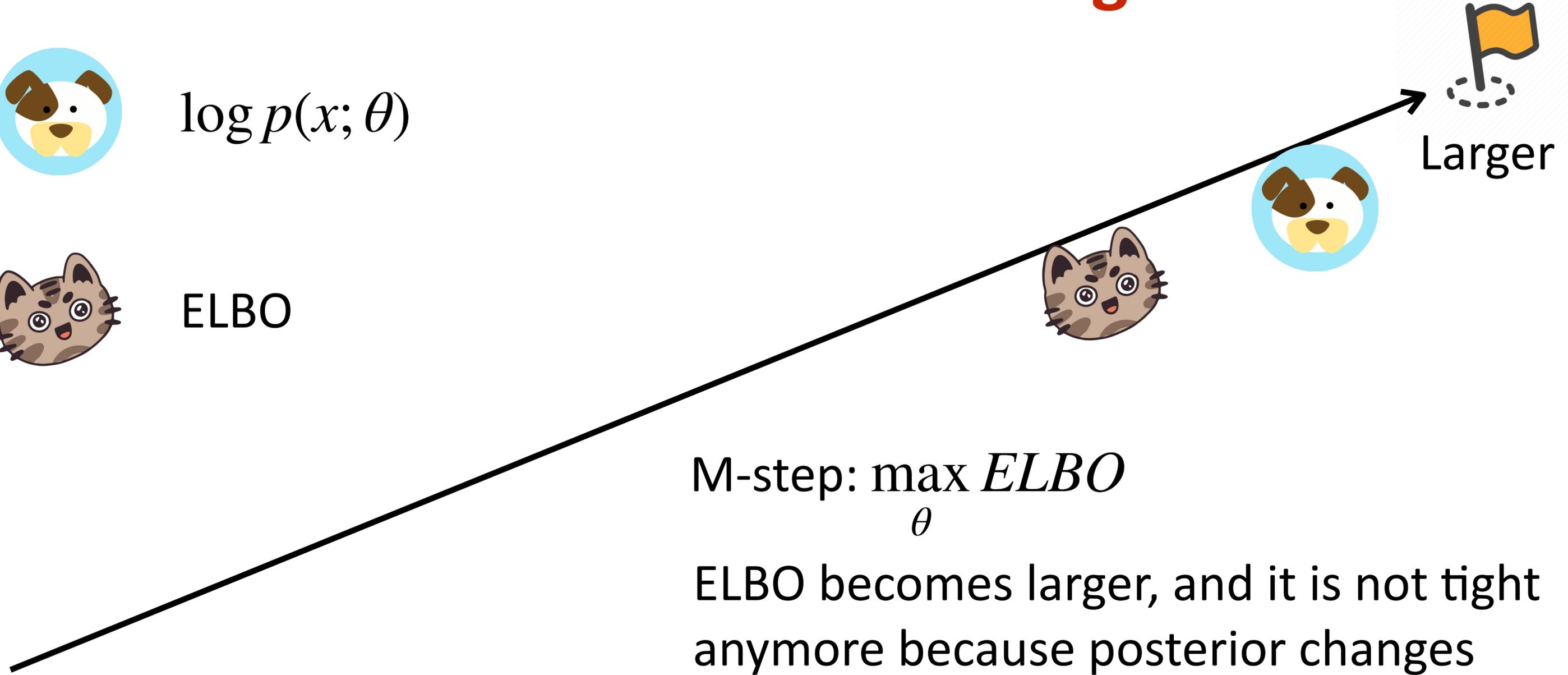
EM is Hill Climbing



$\log p(x; \theta)$



ELBO



M-step: $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes

$\log p(x; \theta)$ is monotonically increasing!
We are doing MLE implicitly!

EM is Hill Climbing



$\log p(x; \theta)$



ELBO

$$\log p(x) \leq \theta$$



Larger

M-step: $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes

$\log p(x; \theta)$ is monotonically increasing!

We are doing MLE implicitly!

Convergence is guaranteed

Revisit the E-Step

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

}

Revisit the E-Step

Computable posterior is important. If $Q(z)$ is not the posterior, then there is no guarantee that $\log p(x)$ is improved at every iteration

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

}

Revisit the E-Step

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

(M-step) Set

$$\begin{aligned} \theta &:= \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta) \\ &= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}. \end{aligned}$$

}

Computable posterior is important. If $Q(z)$ is not the posterior, then there is no guarantee that $\log p(x)$ is improved at every iteration

~~Still remember conjugate prior?~~ Which is for easy-to-compute posterior

Revisit the M-Step

Revisit the M-Step

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} = \operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta)$$

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) \quad \leftarrow \quad \cancel{\sum_z Q(z) \log Q(z)}$$

constant

Revisit the M-Step

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} = \operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta)$$

Sometimes the sum is computable, but sometimes not

Revisit the M-Step

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} = \operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta)$$

Sometimes the sum is computable, but sometimes not

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

Revisit the M-Step

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)} = \operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta)$$

Sometimes the sum is computable, but sometimes not

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

$$\frac{\sum_{i=1}^M \log p(x, z^{(i)}; \theta)}{M}$$

$$z^{(i)} \sim Q(z)$$

We can use Monte-Carlo sampling to approximate the expectation

Comparing Direct Maximization and EM

Direct maximization:

$$\operatorname{argmax}_{\theta} \log \sum_z p(x | z; \theta) p(z) = \operatorname{argmax}_{\theta} \log \mathbb{E}_{z \sim p(z)} p(x | z; \theta)$$

Comparing Direct Maximization and EM

Direct maximization:

$$\operatorname{argmax}_{\theta} \log \sum_z p(x|z; \theta) p(z) = \operatorname{argmax}_{\theta} \log \mathbb{E}_{z \sim p(z)} p(x|z; \theta)$$

Handwritten notes: $z \sim p(z)$ above the expectation term, and a blue circle around the entire equation.

M-Step in EM:

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

Handwritten note: $z \sim Q(z)$

Comparing Direct Maximization and EM

Direct maximization:

$$\operatorname{argmax}_{\theta} \log \sum_z p(x|z; \theta) p(z) = \operatorname{argmax}_{\theta} \log \mathbb{E}_{z \sim p(z)} p(x|z; \theta)$$

M-Step in EM:

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

Why don't we use MC sampling to approximate expectation in direct maximization?

Comparing Direct Maximization and EM

$p(z)$

Direct maximization:

$$\operatorname{argmax}_{\theta} \log \sum_z p(x|z; \theta) p(z) = \operatorname{argmax}_{\theta} \log \mathbb{E}_{z \sim p(z)} p(x|z; \theta)$$

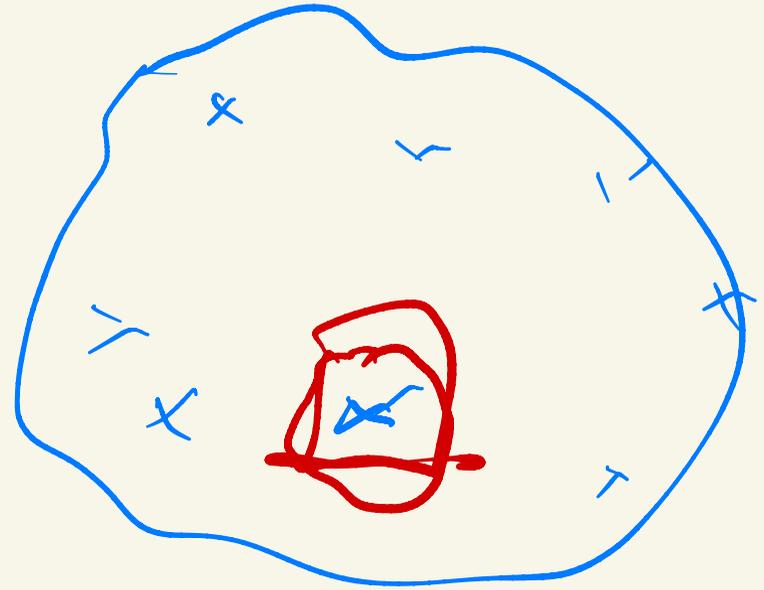
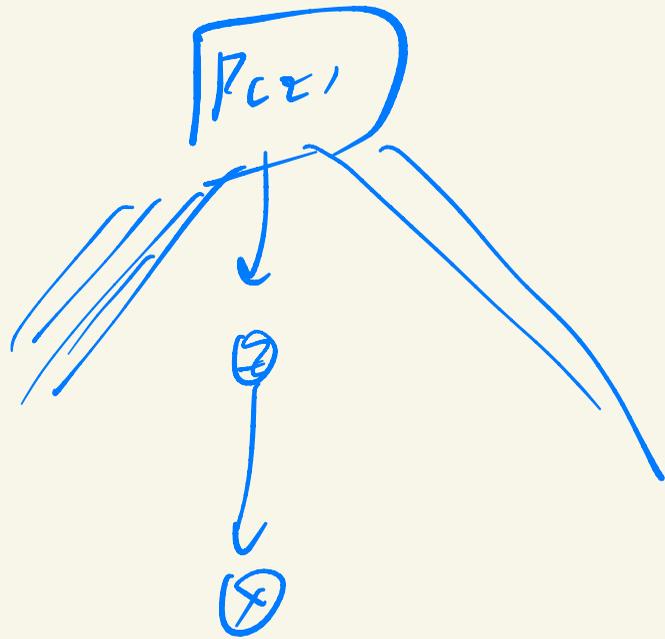
log p(x) (handwritten above the sum)
large variance (handwritten above the expectation)

M-Step in EM:

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

Why don't we use MC sampling to approximate expectation in direct maximization?

It may need a large number of samples to have a good approximation



$$\sum_{z \in P(z)} P(x)$$

$$P(z) = \underline{\underline{P(z|x)}}$$

Other Interpretations of ELBO

$$\begin{aligned} \text{ELBO}(x; Q, \theta) &= \mathbb{E}_{z \sim Q}[\log p(x, z; \theta)] - \mathbb{E}_{z \sim Q}[\log Q(z)] \\ &= \mathbb{E}_{z \sim Q}[\log p(x|z; \theta)] - D_{KL}(Q || p_z) \end{aligned}$$

Regularize $Q(z)$ towards the prior $p(z)$

$$\text{ELBO} = \mathbb{E}_{z \sim Q} \left[\frac{\log p(x, z)}{Q(z)} \right]$$

$$\text{ELBO}(x; Q, \theta) = \log p(x) - D_{KL}(Q || p_{z|x})$$

Maximizing ELBO over $Q(z)$ is essentially solving the posterior distribution $p(z|x)$

$$\text{gap} = \underline{D_{KL}(Q(z) || p(z|x))}$$

$$\max_{\theta} \text{ELBO} = \mathbb{E}_{Q(z)} \left[\log \frac{P(x, z)}{Q(z)} \right]$$

$$= \mathbb{E}_{Q(z)} [\log P(x, z)] - \mathbb{E}_{Q(z)} [\log Q(z)]$$

$$= \mathbb{E}_{Q(z)} [\log P(x|z) + \log P(z)] - \mathbb{E}_{Q(z)} [\log Q(z)]$$

$$= \underbrace{\mathbb{E}_{Q(z)} [\log P(x|z)]}_{\text{reconstruction}} - \underbrace{\mathbb{E}_{Q(z)} \log \frac{Q(z)}{P(z)}}_{\text{KL}(Q(z) \parallel P(z))}$$

Further Questions

Further Questions

- What if we do not have closed-form model posterior?

Further Questions

- What if we do not have closed-form model posterior? → Variational EM

Further Questions

- What if we do not have closed-form model posterior? —> Variational EM

The process of approximating the model posterior is called variational inference

Further Questions

- What if we do not have closed-form model posterior? → Variational EM

The process of approximating the model posterior is called variational inference

We will learn variational autoencoder later

Thank You!
Q & A