

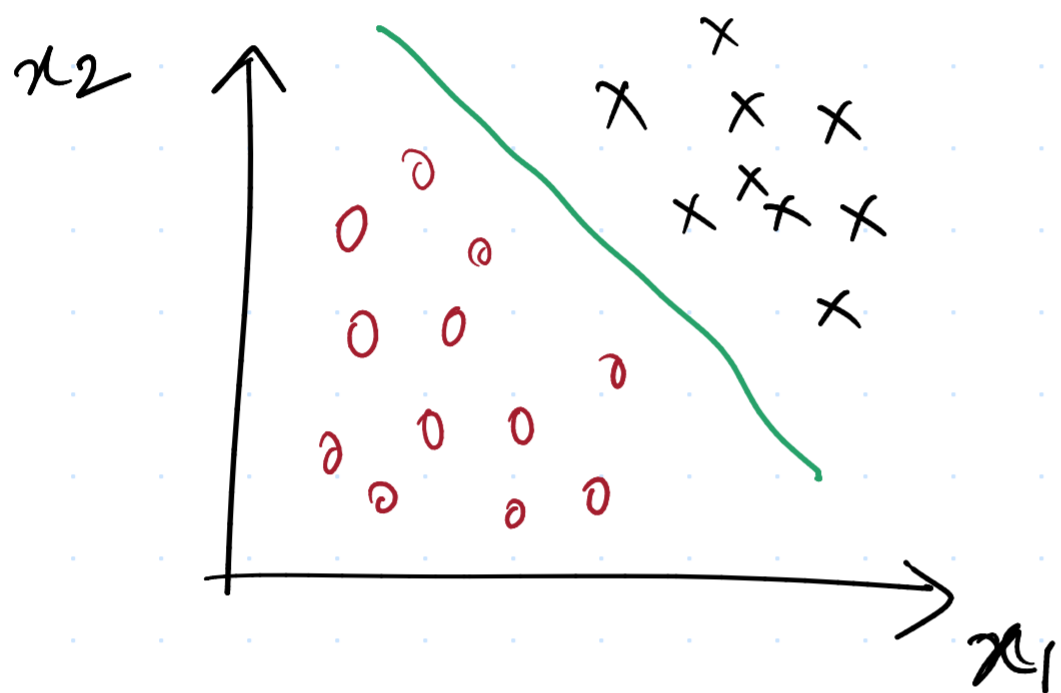
08/03/25

Unsupervised Learning, GMM, EM & VAEs using EMs.

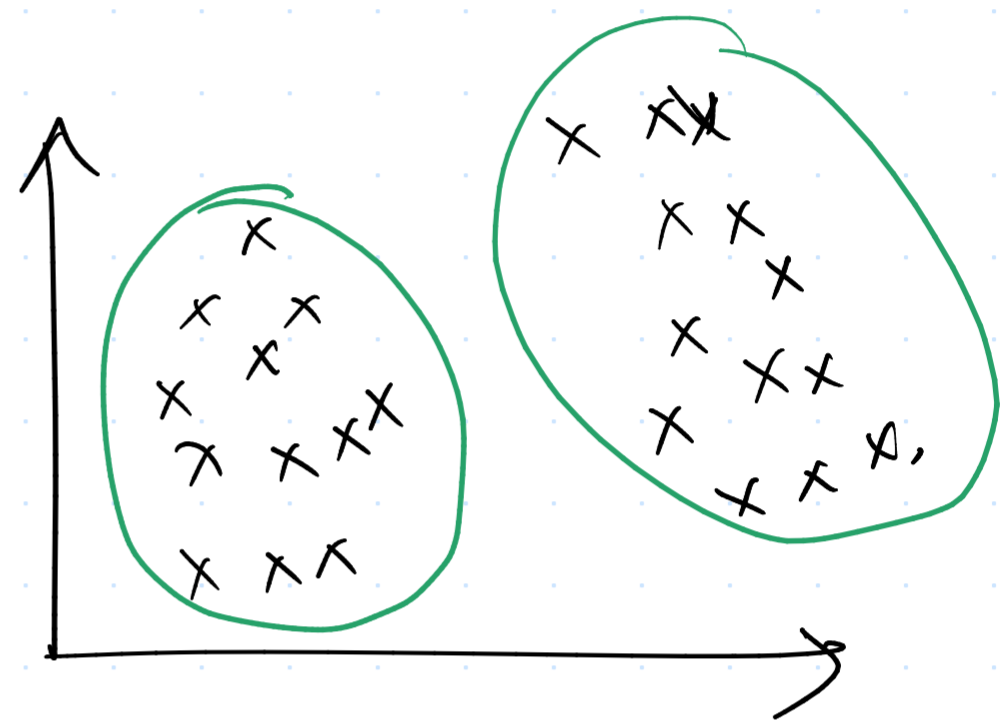
Unsupervised Learning

$X \rightarrow Y$ (Supervised Learning)

$$S = \{X^{(1)}, X^{(2)}, \dots, X^{(n)}\}$$



Supervised Learning.

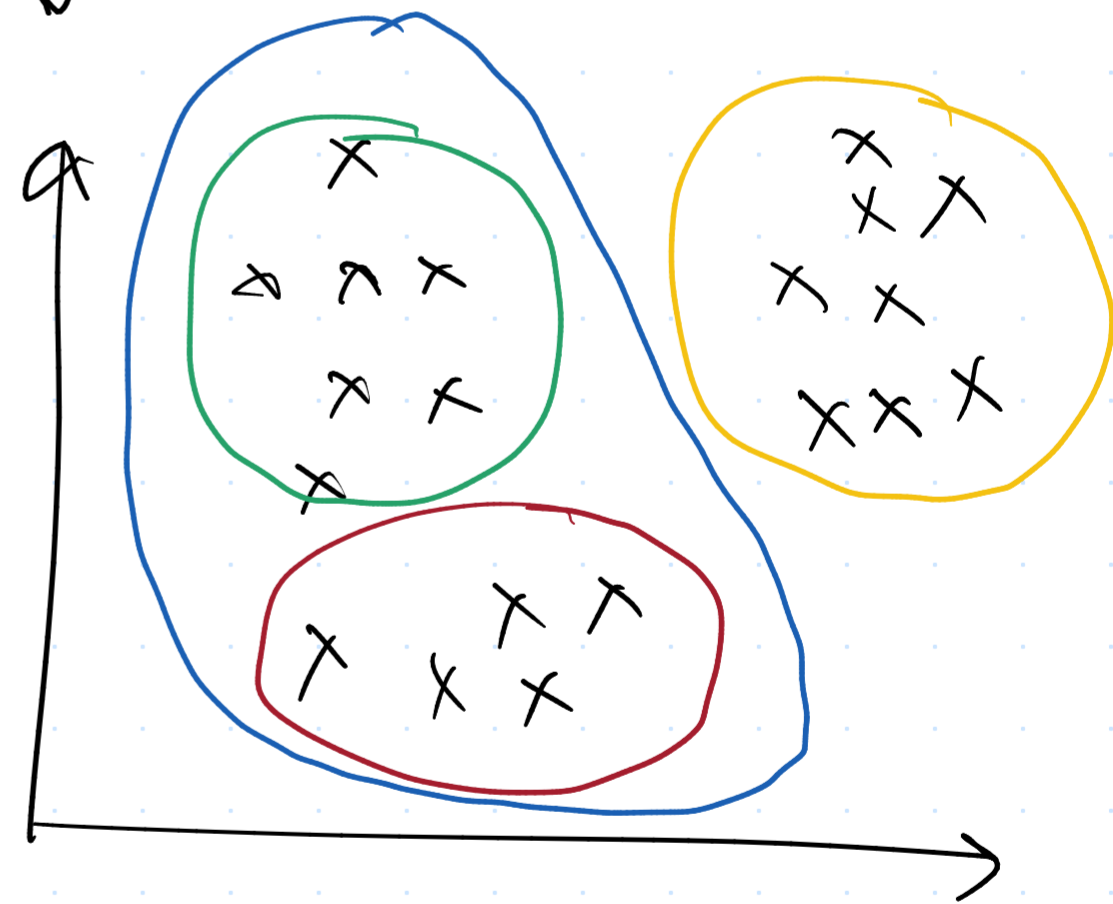


Unsupervised Learning

Unsupervised Learning \rightarrow 3

\rightarrow 2

No correct way to think that only a particular cluster assignment is possible. & there will always be "ambiguity" here



Classification problems in supervised settings can be related to clustering problems in unsupervised (settings) learning.

K-Means clustering Algorithm

$$S = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}, \quad x^{(i)} \in \mathbb{R}^d$$

K-clusters are given.

1. Initialize the cluster centroids :-

$$\mu_1, \mu_2, \mu_3, \dots, \mu_k \in \mathbb{R}^d \rightarrow \text{randomly.}$$

μ_i 's \rightarrow vector

2. Repeat until convergence

- For every i , set,

$$c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|_2^2$$

- For every j , set :-

$$\mu_j = \frac{\sum_{i=1}^n \mathbb{1}_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{\{c^{(i)}=j\}}}$$

Mean of all
 x_i 's for which
the $c_i = j$

Distortion function

$$J(c, \mu) = \sum_{i=1}^n \|x^{(i)} - \mu_{c(i)}\|_2^2$$

J -means \rightarrow minimizing the loss of the distortion function in the form of (coordinate descent) \rightarrow variation of gradient descent.

Instead of minimizing the loss w.r.t. all the variables, we minimize the loss w.r.t. few variables by keeping the others fixed.

step $\rightarrow 1$: μ -fixed & optimize c

step $\rightarrow 2$: c -fixed & optimize w.r.t. μ .

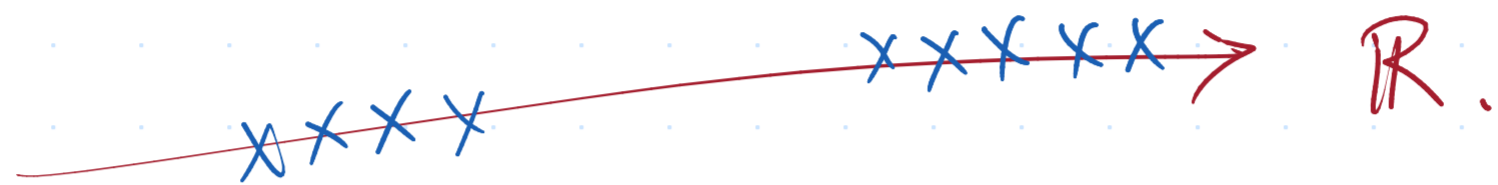
K -means converge eventually since we will reach some form of local minima of the J -function. We may toggle b/w two sets of μ 's & c 's alternatively once we reach a local minima.

J -non-convex, μ 's & c 's that we end up in can change from run to run. If we start with a different

initialization, we might end up with different μ 's & σ 's.

(Non-convex problem).

Density Estimation

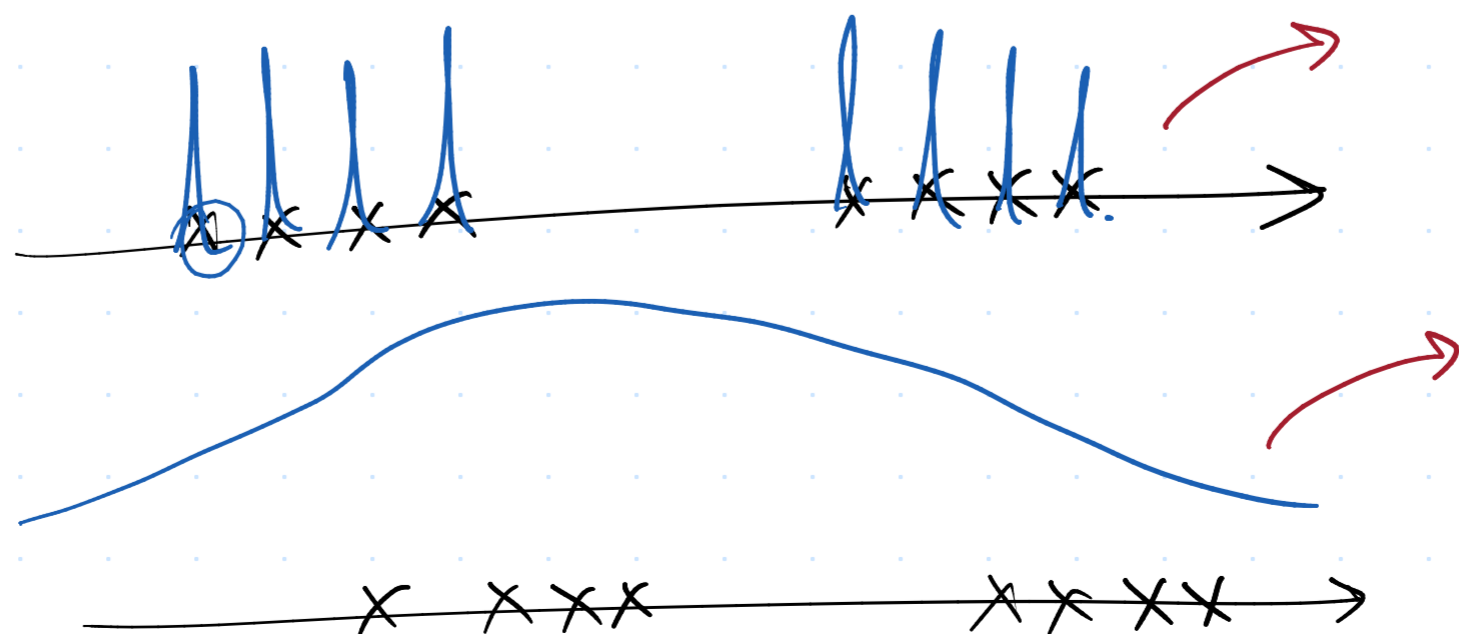


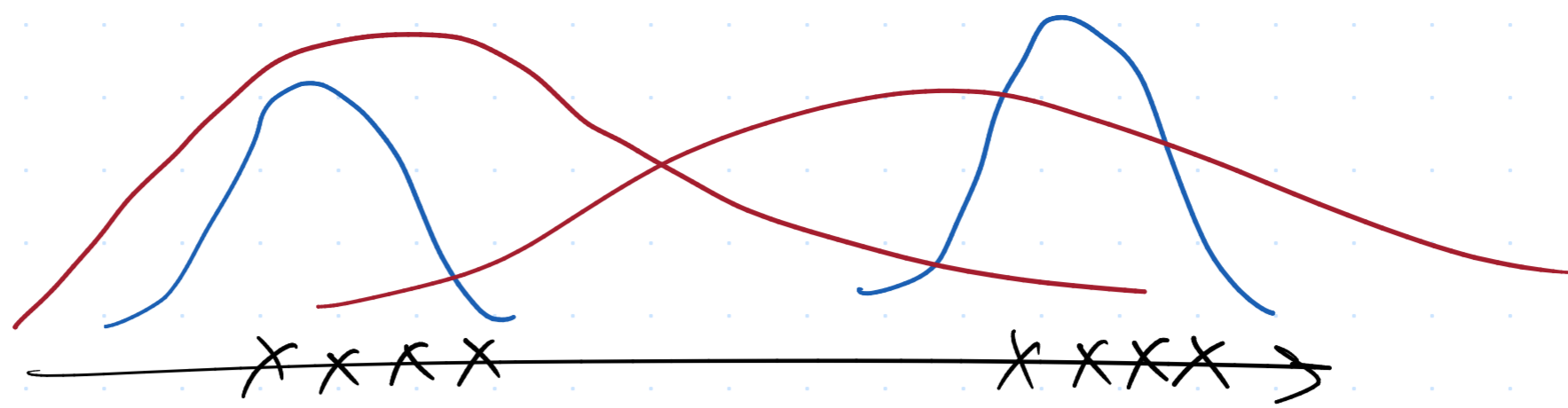
Points are residing in a continuous space. Sampled from some kind of probability distribution. The points come from a continuous distribution, the corresponding prob. dist has some kind of density.

PMF x PDF ✓

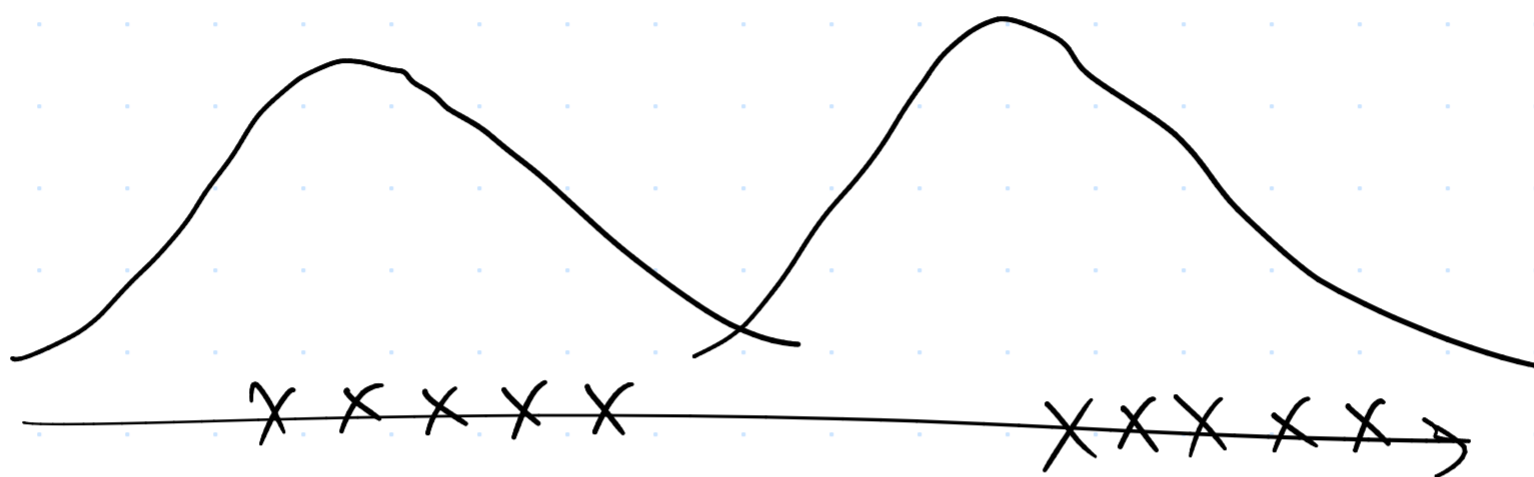
Not a probability mass function, but a probability density function.

Hard problem





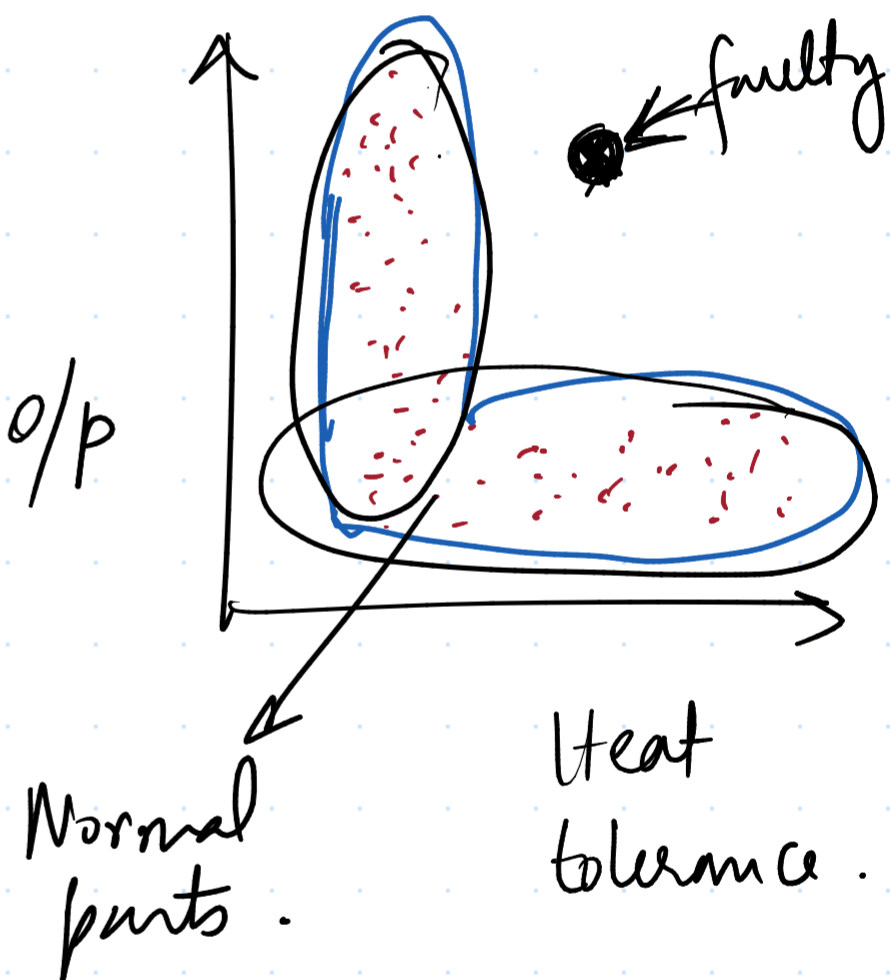
Gaussian Mixture Model :- Mixture of Gaussians.



$\phi(x)$ \rightarrow to be observed values.

mixture of Gaussians

Power \rightarrow o/p



(G.M.M. - Algorithm)

$$S = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$$

$z^{(i)} \sim$ Multinomial (ϕ)

$$\phi_j \geq 0 \quad \sum_{j=1}^K \phi_j = 1$$

$$\phi_j = p(z^{(i)} = j)$$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \Sigma_j)$$

$z^{(i)} \rightarrow$ latent variable, fancy name for R.V. that has not been observed.

GDA : Maximum Likelihood Estimation in GDA.

$$\log p(x, y; \mu, \Sigma, \phi) = \ell(\mu, \Sigma, \phi)$$

GMM

$$\log p(x; \mu, \Sigma, \phi) = \ell(\mu, \Sigma, \phi)$$

↑
maximize

$$= \log \sum_z p(x, z; \phi, \mu, \Sigma)$$

Writing out the full joint dist & marginalizing out the latent variable.

If we have observed z , then it could have contributed to the likelihood objective, but we don't know what z is, so it cannot be written - formulated like that.

$p(z)$ → class prior / if continuous.

$p(x, z)$ → model



z → latent (unobserved)

$p(z|x)$ = posterior

discriminating
process.

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

↑
generating process.

$p(x)$ → evidence.

Goal = maximize the likelihood given the evidence.

EM Algorithm - inspired by k-Means.

Repeat until convergence (randomly initialize μ, ϕ, Σ)

E-step

For each i, j , set:-

Given a point, we assign a prob. that this point belongs to a particular centroid.

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

weight / posterior dist.

M-step

Update the parameters

We don't need to know $z^{(i)}$, we are constructing the probability $z^{(i)} = j$ using the Bayes rule.

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

(Soft k-means)

$k=3,$

$$p(z^{(i)} = j | x^{(i)})$$

$$\Rightarrow \begin{bmatrix} 0.1 & k=1 \\ 0.7 & k=2 \\ 0.2 & k=3 \end{bmatrix} \rightarrow$$

We are doing soft assignment, where every point has a soft assignment in the form of probability distribution.

k-means \rightarrow hard assignment \rightarrow where each point belongs to one cluster only.

Gamma \quad Multinomial

$$p(x|z) = \frac{p(x|z) p(z)}{\sum_z p(x|z) p(z)}$$

\rightarrow this takes the form of a softmax.

EM algorithm gives a framework for performing

MLE when some variables are unobserved.

Used in cases where we have a functional form of

$$p(x, z)$$

x & $z \rightarrow$ could be anything, but

z - is unobserved.

Expectation Maximization

MLE in the presence of latent variables.

$$\text{True model} = p(x, z; \theta)$$

If x & $z \rightarrow$ observed, we perform simple MLE.

But when $z \rightarrow$ unobserved, we instead perform maximization.

$$l(\theta) = \log p(x; \theta)$$

EM gives us a framework for achieving this.

Jensen's inequality \rightarrow probabilistic tool.

$f \rightarrow$ to be convex function

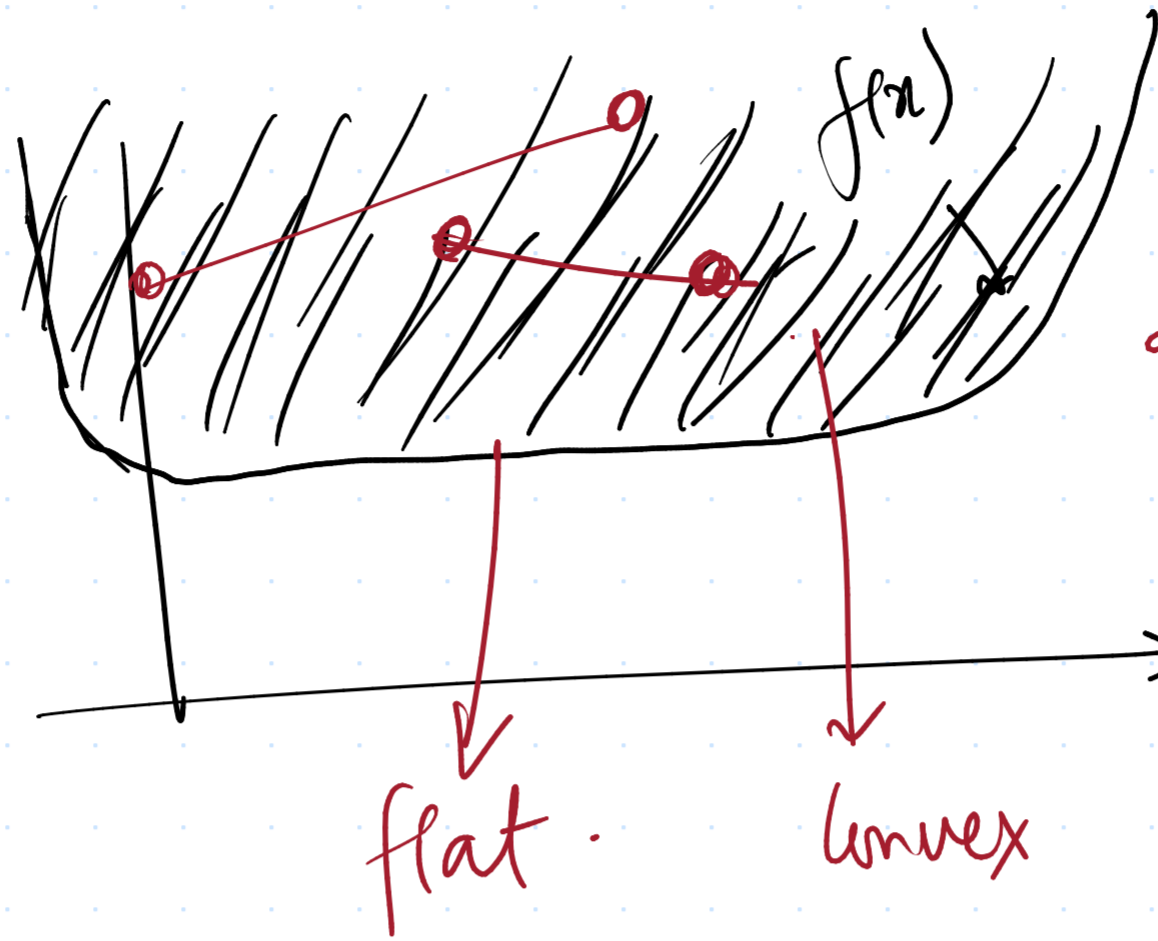
$$f''(x) \geq 0 \quad \forall x.$$

$f \rightarrow$ is strictly convex if, $f''(x) > 0 \quad \forall x.$

$$E[f(x)] \geq f(E[x]) \rightarrow \text{where } f \text{ is Convex.}$$

Convex function

R.V.



$$f(x) = 2x^2$$

$$f'(x) = 4x$$

$$f''(x) = 4$$

$$f''(x) > 0?$$

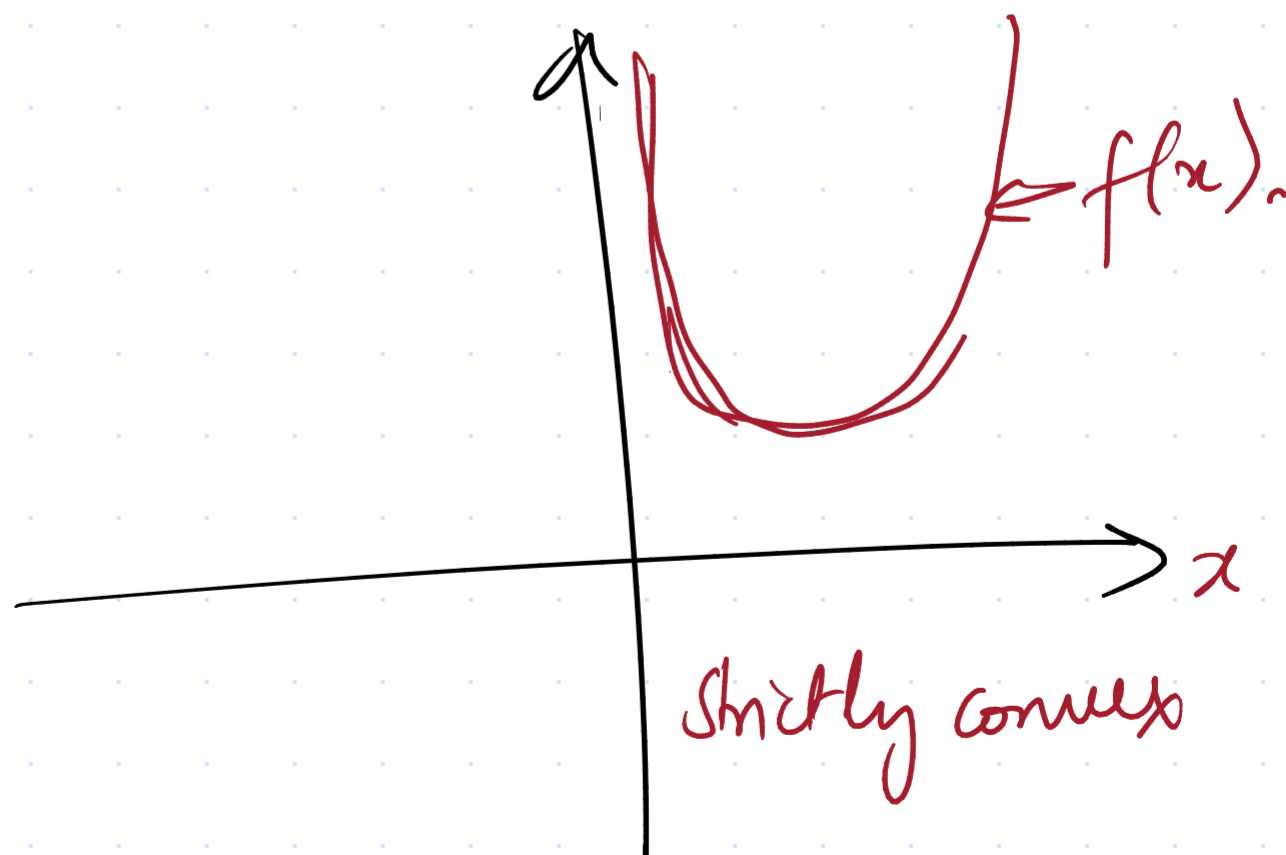
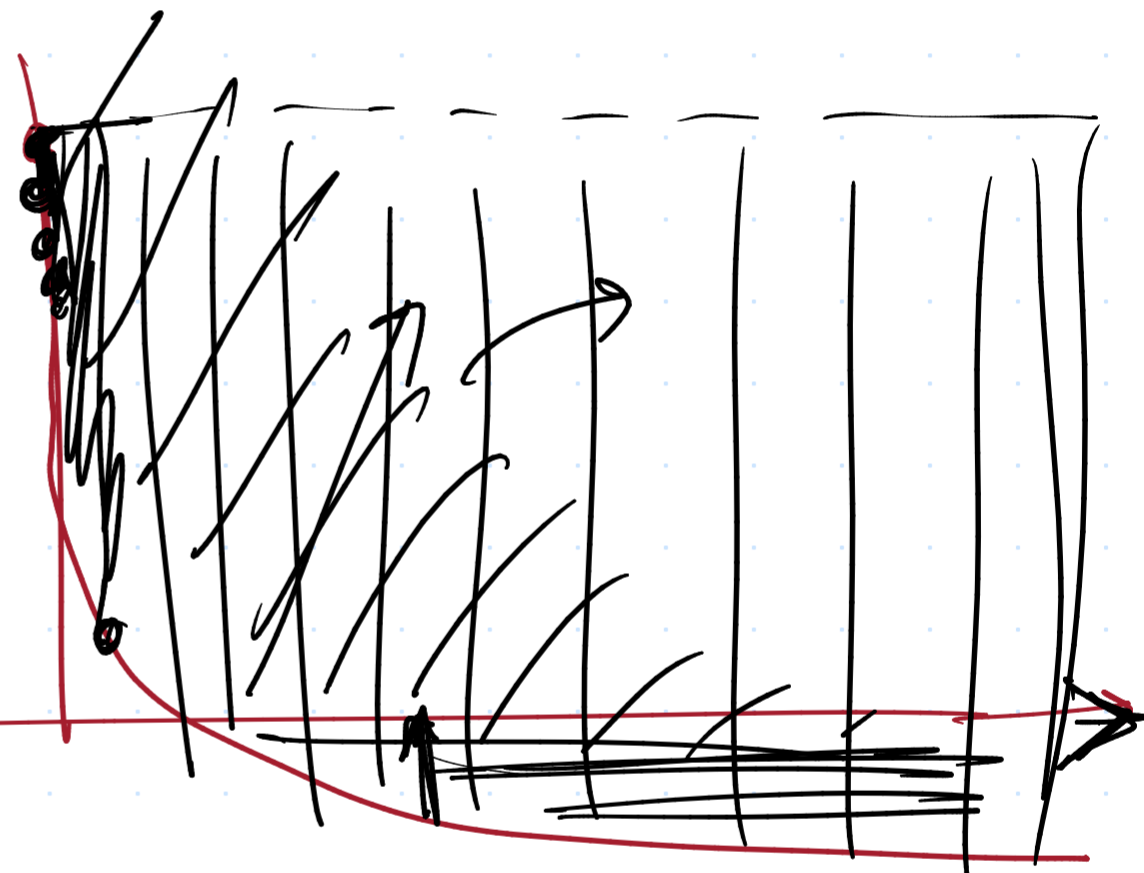
Convex.

How to define this region?

If f is strictly convex, then,

If $E[f(x)] = f(E[x])$ then,

$x = E[x]$, with prob = 1.



$$f''(x) \geq 0$$

convex

$$f''(x) \leq 0$$

concave

strict.

\rightarrow either \leftarrow

$$f(x) = x^2 \quad \checkmark$$

$$f(x) = -x^2 \quad \checkmark$$

$$f(x) = mx + c \quad \checkmark$$

\checkmark

\times

$$e^x \rightarrow \text{convex}$$

$$-e^x \rightarrow \text{concave}$$

\checkmark

$$e^{-x} \rightarrow$$

$$-e^{-x} \rightarrow$$

\checkmark

$$-\log(x) \rightarrow \text{convex}$$

$$\log(x) \rightarrow \text{concave}$$

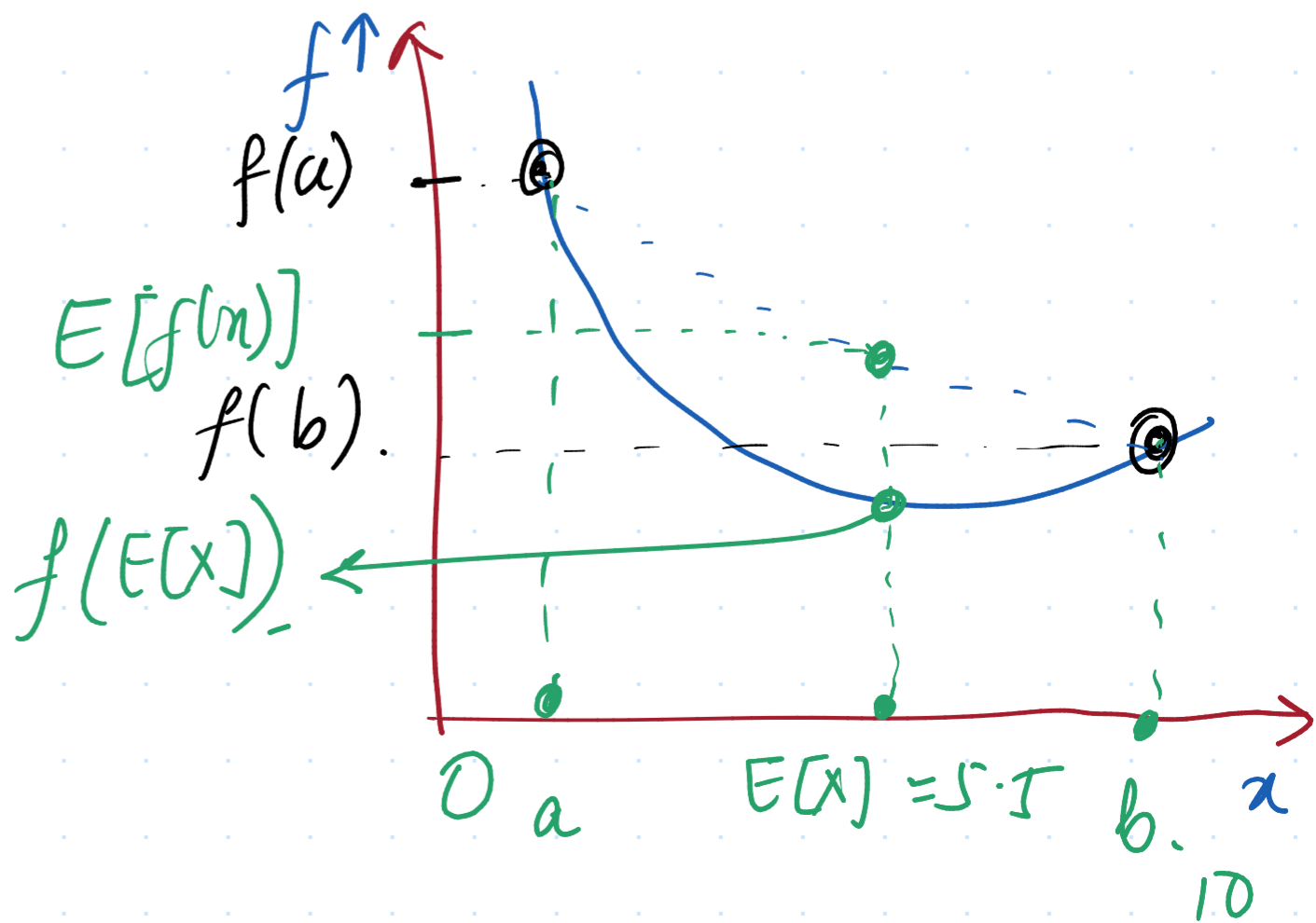
\checkmark

\checkmark If f is concave, e.g. $f(x) = \log(x)$

$$E[f(x)] \leq f(E[x])$$

Jensen's Equality

$$E[\log(x)] \leq \log(E[x])$$



EM Algorithm

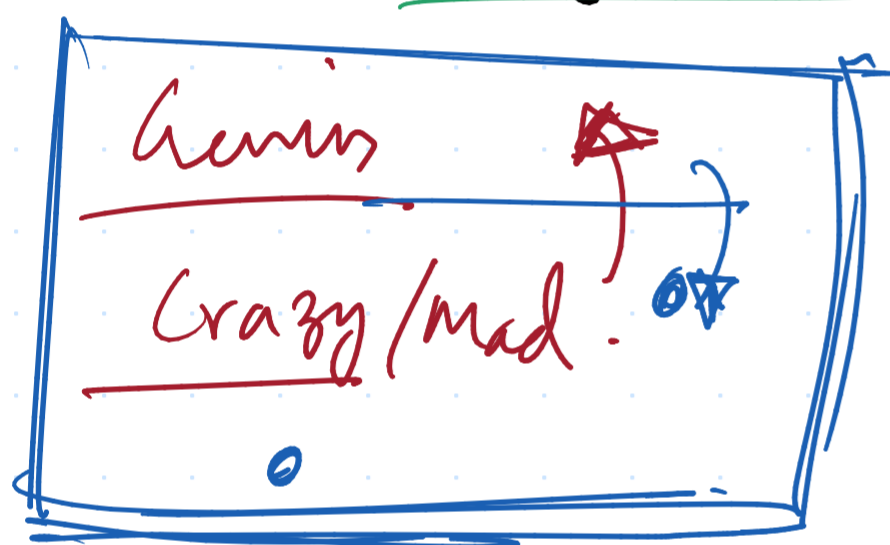
Goal: Maximize $\sum_{i=1}^n \log p(x^{(i)}; \theta)$.

z 's were unobserved, hence maximizing this is very hard.

$$\max \log p(x; \theta)$$

summation over z , here z will be continuous and hence the integration will be arbitrarily complex & hard.

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



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

$Q(z) > 0 \forall z$
↑
function.

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

Applying Jensen's Equality

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

→ concave function.

ELBO Evidence Lower Bound. $\therefore = \underline{\underline{ELBO(x; Q, \theta)}}$

If we are finding values of θ & Q , then the same values of θ , implicitly the $\log p(x; \theta)$ is going up.

Since, $ELBO(x; Q, \theta) \leq \log p(x; \theta)$ hence this is true.

There are cases when,

$$\boxed{\log p(x; \theta) = ELBO(x; Q, \theta)}$$

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

↓
if this is a constant.

Since, log is strictly convex :-

$$\frac{p(x, z; \theta)}{Q(z)} \rightarrow \text{constant, w.r.t. } z. = C.$$

$$Q(z) = \frac{1}{C} p(x, z; \theta) \rightarrow \text{proportionality constant.}$$

$$Q(z) \propto p(x, z; \theta)$$

distribution over z
that must sum to 1.

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

normalization
constant similar
to the
proportionality
constant.

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

$$\sum Q(z) = 1.$$

x & z \rightarrow doesn't have any specific form & it
can be mixture of Gaussians, it could be any algorithm etc.

This holds for any latent variable models, hence, we

will apply it to optimize/train VAE's in the next

lecture.

So, EM Algorithm :-

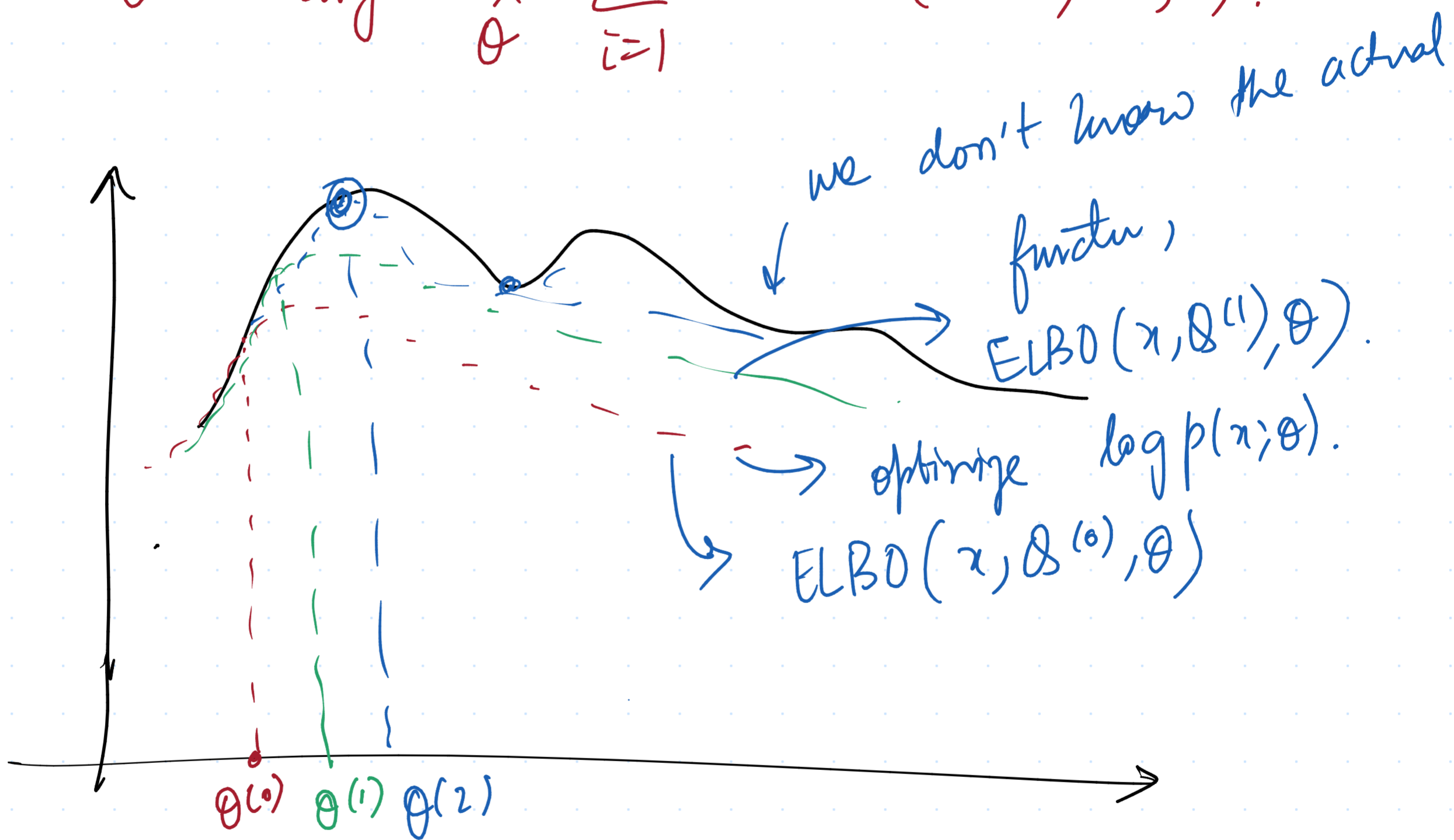
E-step

For each i , set :-

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

M-step

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



So, in every step, the ELBO that we get is tight & will be touching the $\log p(x)$.