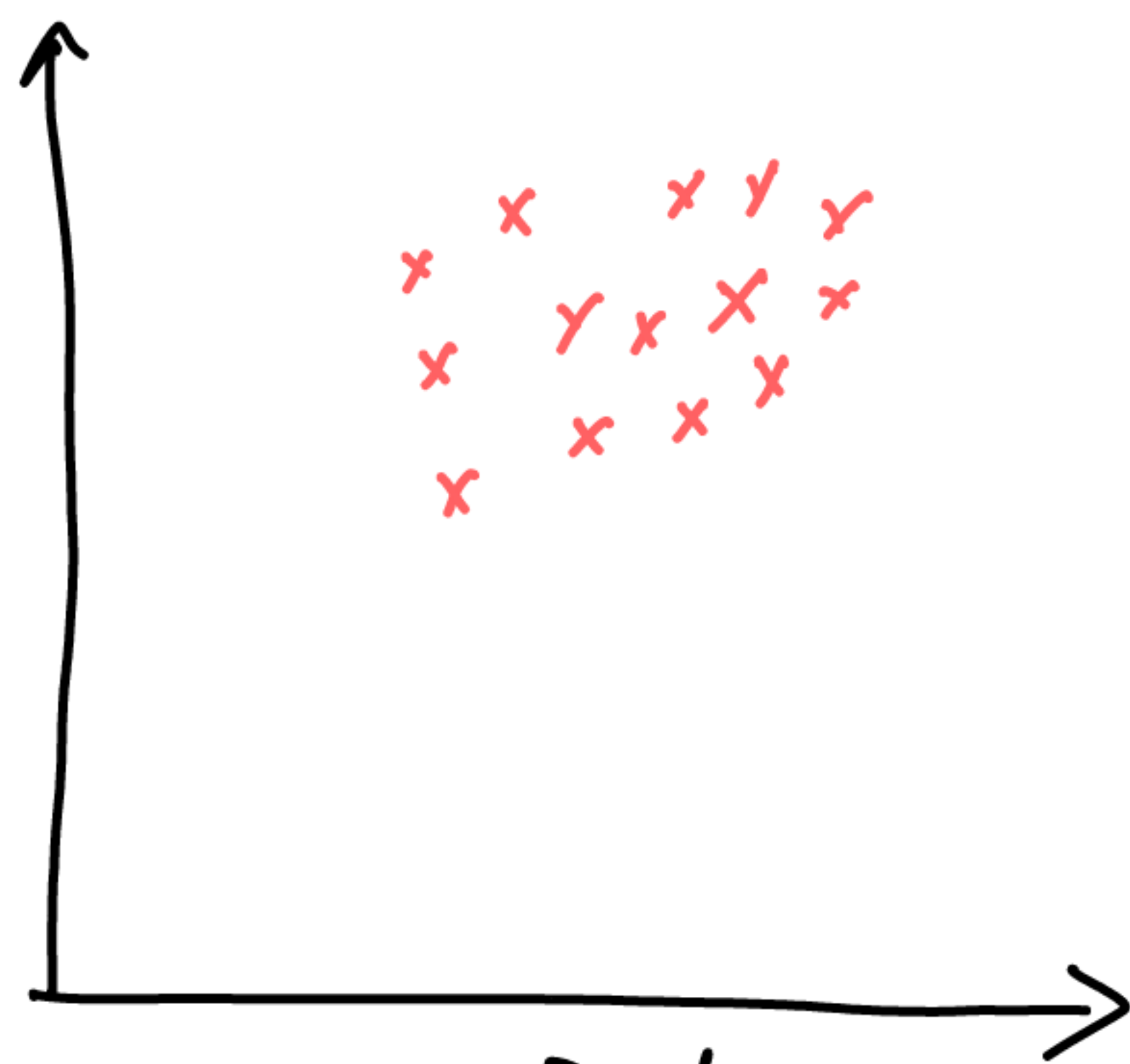$$P(\theta|x) = \frac{P(x|\theta) \, P(\theta)}{P(x)}$$

where

$$P(x) = \int_\theta P(x|\theta) P(\theta) d\theta$$

$$= \int_{\theta_1} \int_{\theta_2} \cdots \int_{\theta_m} P(x|\theta_1, \theta_2 \cdots \theta_m) d\theta_1, \cdots d\theta_m$$

$X =$ Data
$= \{x_1, x_2 \cdots x_n\}$

multiple integral intractable to compute.

To approximate the posterior $P(\theta|x)$ we have to find another posterior from familly of dist$^n$ $Q$.

$p(\theta|x) \longleftrightarrow q(\theta)$

$Q$

$$KL\left(q(\theta) \| P(\theta|x)\right) = \int q(\theta) \log \frac{q(\theta)}{P(\theta|x)} d\theta$$

$$= \int q(\theta) \left( \log q(\theta) - \log P(\theta|x) \right) d\theta$$

$$= \int q(\theta) \log q(\theta) d\theta - \int q(\theta) \log P(\theta|x) d\theta$$

$$= \mathbb{E}_q \left[ \log q(\theta) \right] - \mathbb{E}_q \left[ \log P(\theta|x) \right]$$

$$\therefore P(\theta|x) = \frac{P(\theta,x)}{P(x)}$$

$$= \mathbb{E}_q\left[\log q(\theta)\right] - \mathbb{E}_q\left[\log p(x,\theta) - \log p(x)\right]$$

$$KL\left(q(\theta) \| p(\theta|x)\right) = \mathbb{E}_q\left[\log q(\theta) - \log p(x,\theta)\right] + \log p(x)$$

$$\mathbb{E}_q\left[\log p(x,\theta) - \log q(\theta)\right] = \log p(x) - KL\left(q(\theta) \| p(\theta|x)\right)$$

Evidence lower bound

$$ELBO(q) = \mathbb{E}_q\left[\log \frac{p(x,\theta)}{q(\theta)}\right]$$

we need to maximize ELBO to find parameter of Approximate posterior

Let's us unpack ELBO

$$ELBO(q) = \mathbb{E}_q\left[\log p(x,\theta) - \log q(\theta)\right]$$

$$= \mathbb{E}_q\left[\log p(x,\theta)\right] - \mathbb{E}_q\left[\log q(\theta)\right]$$

$$= \mathbb{E}_q\left[\log\left(p(x|\theta)p(\theta)\right)\right] - \mathbb{E}_q\left[\log q(\theta)\right]$$

$$ELBO(q) = \mathbb{E}_q\left[\log p(x|\theta)\right] + \mathbb{E}_q\left[\log p(\theta) - \log q(\theta)\right]$$

$$= \mathbb{E}_q\left[\log p(x|\theta)\right] + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta$$

$$ELBO(q) = \mathbb{E}_q\left[\log p(x|\theta)\right] - KL\left(q(\theta) \| p(\theta)\right)$$

log liklihood          KL-div

In our case of 2-dim data $\left(\begin{array}{c}\text{dim}\\ K=2\end{array}\right)$

$$\theta = (\mu, \Sigma)$$

$q(\theta) \longrightarrow$ parameterize by variational parameter $\mu_q, \Sigma_q$ we have to find this

$p(\theta) \longrightarrow$ parametrize by prior parameter obtain from data

$$\mu_p = \left(\mu_x = \frac{1}{n}\sum_{i=1}^{n} x_i, \ \mu_y = \frac{1}{n}\sum_{i=1}^{n} y_i\right)$$

$$\Sigma_p = \frac{1}{n}\sum_{i=1}^{n} \underbrace{(x_i - \mu_x)}_{2\times 1}\underbrace{(y_i - \mu_y)^T}_{1\times 2}$$

**<u>Now</u> log <u>liklihood</u>:**

$$\text{log-liklihood} = \log p(x|\theta) \qquad \theta = (\mu_q, \Sigma_q)$$

$$= \log\left(\prod_{i=1}^{n} p(x_i|\theta)\right) \qquad \text{data } X \text{ in iid}$$

$$= \log\left(\prod_{i=1}^{n} p(x_i | \mu_q, \Sigma_q)\right)$$

$$= \log\left(\prod_{i=1}^{n} \frac{1}{(2\pi)^{\frac{K}{2}} |\Sigma_q|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-\mu_q)^T \Sigma_q^{-1}(x-\mu_q)\right\}\right)$$

$$\log p(x|\theta) = \sum_{i=1}^{n}\left[\log\left(\frac{1}{(2\pi)^{K/2}|\Sigma_q|^{1/2}}\right) - \frac{1}{2}(x_i - \mu_q)^T \Sigma_q^{-1}(x_i - \mu_q)\right]$$

$$\log p(X|\theta) = n \log \left( \frac{1}{(2\pi)^{K/2} |\Sigma_q|^{1/2}} \right)$$

$$- \frac{1}{2} \sum_{i=1}^{n} (X_i - \mu_q)^T \Sigma_q^{-1} (X_i - \mu_q)$$

$$\log p(X|\theta) = n \left[ -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_q| \right]$$

$$- \frac{1}{2} \sum_{i=1}^{n} (X_i - \mu_q)^T \Sigma_q^{-1} (X_i - \mu_q) \quad \text{(A)}$$

$$\text{where } X_i \in X = \underbrace{\{x_1, x_2, \ldots x_n\}}_{Data}$$

**KL Divergence :**

$$KL \left( q(\theta) \| p(\theta) \right)$$

$$q(\theta) = \mathcal{N}(\mu_q, \Sigma_q)$$

$$p(\theta) = \mathcal{N}(\mu_p, \Sigma_p)$$

$$KL(q \| p) = \frac{1}{2} \left[ -\log \frac{|\Sigma_q|}{|\Sigma_p|} - K + (\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q)^T \right.$$

$$\left. + tr \left\{ \Sigma_p^{-1} \Sigma_q \right\} \right] \quad \text{(B)}$$

To update variational distribution $q(\theta)$ parameter we need to find gradient of ELBO w.r.t parameters

$$\mu_{q\_new} = \mu_{q_{old}} + \eta \frac{\partial ELBO}{\partial \mu_q}$$

$$\Sigma_{q_{new}} = \Sigma_{q_{old}} + \eta \frac{\partial ELBO}{\partial \Sigma_q}$$

$\left.\begin{array}{r}\end{array}\right\}$ Applying Gradient Ascent

Now

$$\frac{\partial ELBO}{\partial \mu_q} = \frac{\partial \log P(X|\theta)}{\partial \mu_q} - \frac{\partial KL(q(\theta) \| P(\theta))}{\partial \mu_q}$$

from Ⓐ

let

$$\frac{\partial \log P(X|\theta)}{\partial \mu_q} = -\frac{1}{2} \sum_{i=1}^{n} (-1) \, 2 \, (X_i - \mu_q) \, \Sigma_q^{-1}$$

$$= \sum_{i=1}^{n} (X_i - \mu_q) \Sigma_q^{-1}$$

from Ⓑ

$$\frac{\partial KL(q(\theta) \| P(\theta))}{\partial \mu_q} = -\Sigma_p^{-1} (\mu_p - \mu_q)$$

$$\boxed{\frac{\partial ELBO}{\partial \mu_q} = \sum_{i=1}^{n} (X_i - \mu_q) \Sigma_q^{-1} + \Sigma_p^{-1} (\mu_p - \mu_q)}$$

<u>Now</u> $\dfrac{\partial ELBO}{\partial \Sigma_q} = \dfrac{\partial \log p(X|\theta)}{\partial \Sigma_q} - \dfrac{\partial KL\left(q(\theta)\|p(\theta)\right)}{\partial \Sigma_q}$

let

$$\frac{\partial \log p(X|\theta)}{\partial \Sigma_q} = \frac{\partial}{\partial \Sigma_q}\left[ n\left\{ -\frac{K}{2}\log(2\pi) - \frac{1}{2}\log|\Sigma_q| \right\} - \frac{1}{2}\sum_{i=1}^{n}(X_i - \mu_q)^T \Sigma_q^{-1}(X_i - \mu_q) \right]$$

$$= -\frac{n}{2}\frac{1}{|\Sigma_q|}\frac{\partial|\Sigma_q|}{\partial \Sigma_q}$$

$$-\frac{1}{2}\sum_{i=1}^{w}\left\{ -\Sigma_q^{-1}(X_i-\mu_q)(X_i-\mu_q)^T \Sigma_q^{-1} \right\}$$

$$\boxed{\begin{array}{l}\dfrac{\partial |X|}{\partial X} = |X|(X^{-1})^T \\[2mm] \dfrac{\partial a^T X^{-1} b}{\partial X} = -X^{-1}ab^T X^{-1} \\[2mm] \text{if } X \text{ is symmetrical} \\ \text{matrix}\end{array}}$$

$$= -\frac{n}{2}\frac{|\Sigma_q|}{|\Sigma_q|}\left(\Sigma_q^{-1}\right)^T + \frac{1}{2}\sum_{i=1}^{w}\left[ \Sigma_q^{-1}(X_i-\mu_q)^2 \Sigma_q^{-1} \right]$$

$$\frac{\partial \log p(X|\theta)}{\partial \Sigma_q} = -0.5\, n\left(\Sigma_q^{-1}\right)^T + 0.5\sum_{i=1}^{n}\left[ \Sigma_q^{-1}(X_i-\mu_q)^2 \Sigma_q^{-1} \right]$$

$$\boxed{\frac{\partial \log p(X|\theta)}{\partial \Sigma_q} = -0.5\left\{ n\left(\Sigma_q^{-1}\right)^T - \sum_{i=1}^{n}\left[ \Sigma_q^{-1}(X_i-\mu_q)^2 \Sigma_q^{-1} \right] \right\}}$$

$$\frac{\partial KL(q(\theta) \| p(\theta))}{\partial \Sigma_q} = \frac{\partial}{\partial \Sigma_q} \left\{ \frac{1}{2} \left\{ -\log \frac{|\Sigma_q|}{|\Sigma_p|} - k \right.\right.$$

$$\left.\left. + (\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q) + tr\left(\Sigma_p^{-1} \Sigma_q\right) \right]\right.$$

$$= \frac{1}{2} \left[ -\frac{\partial}{\partial \Sigma_q} \left\{ \log \frac{|\Sigma_q|}{|\Sigma_p|} \right\} + \frac{\partial}{\partial \Sigma_q} \left\{ tr\left(\Sigma_p^{-1} \Sigma_q\right) \right\} \right]$$

$$\Rightarrow = \frac{\partial}{\partial \Sigma_q} \left[ \log |\Sigma_q| - \log |\Sigma_p| \right]$$

$$= \frac{1}{|\Sigma_q|} \frac{\partial}{\partial \Sigma_q} |\Sigma_q|$$

$$= \frac{1}{|\Sigma_q|} |\Sigma_q| \left(\Sigma_q^{-1}\right)$$

$$= \left(\Sigma_q^{-1}\right)^T \qquad \boxed{\because \frac{\partial |X|}{\partial X} = |X|\left(X^{-1}\right)^T}$$

$$\Rightarrow \frac{\partial}{\partial \Sigma_q} \left[ tr\left\{ \Sigma_p^{-1} \Sigma_q \right\} \right]$$

$$= \left(\Sigma_p^{-1}\right)^T \qquad \boxed{\because \frac{\partial Tr(AX)}{\partial X} = A^T}$$

$$\frac{\partial ELBO}{\partial \Sigma_q} = \frac{\partial \log p(X|\theta)}{\partial \Sigma_q} - \frac{\partial KL(q(\theta) \| p(\theta))}{\partial \Sigma_q}$$

$$\frac{\partial ELBO}{\partial \Sigma_q} = -0.5 \left[ n \left( \Sigma_q^{-1} \right)^T - \sum_{i=1}^{n} \left( \Sigma_q^{-1} (x_i - \mu_q)^2 \Sigma_q^{-1} \right) \right]$$

$$- \frac{1}{2} \left[ -\Sigma_q^{-1} + \left( \Sigma_p^{-1} \right)^T \right]$$