

Sequential Data Modeling

Tomoki Toda
Graham Neubig
Sakriani Sakti

Augmented Human Communication Laboratory
Graduate School of Information Science

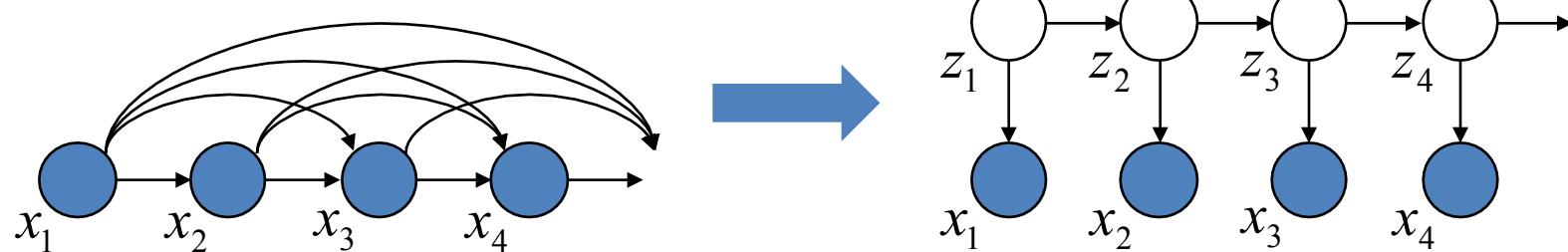
Syllabus

Date	Course description	Lecturer
6/03	Basics of sequential data modeling 1	Graham Neubig
6/10	Basics of sequential data modeling 2	Graham Neubig
6/17	Discrete latent variable models 1	Tomoki Toda
6/24	Discrete latent variable models 2	Tomoki Toda
7/1	Continuous latent variable models 1	Tomoki Toda
7/15	Discriminative models for sequential labeling 1	Sakriani Sakti
7/29	Continuous latent variable models 2	Tomoki Toda
8/1	Discriminative models for sequential labeling 2	Sakriani Sakti

Questions to: **sdm2016@is.naist.jp**

1st and 2nd Classes (6/3 and 6/10)

- Lecturer: Graham Neubig
- Contents: Basics of sequential data modeling
 - Markov process
 - Latent variables
 - Mixture models
 - Expectation-maximization (EM) algorithm



Sequential Data Modeling

2nd class

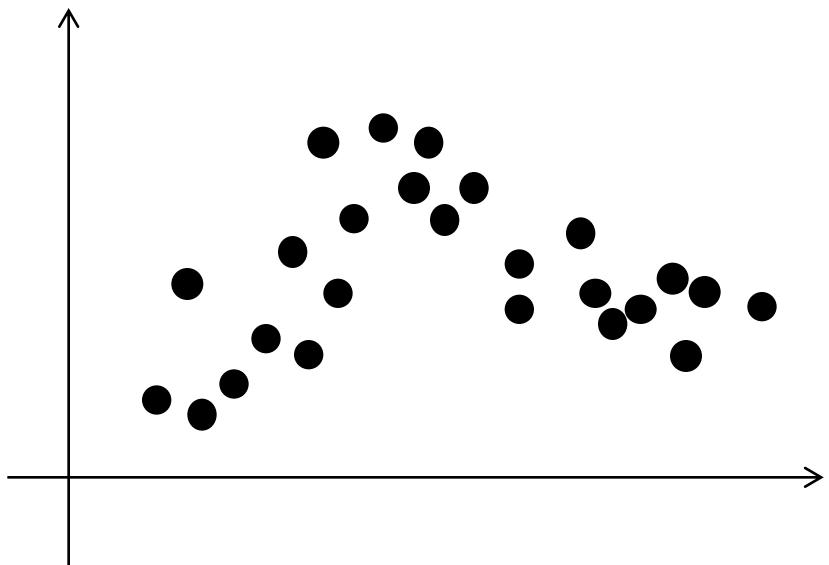
“Basics of sequential data modeling 2”

Graham Neubig

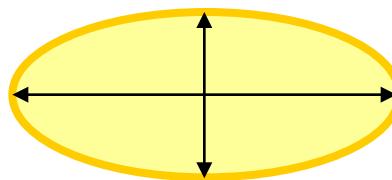
Augmented Human Communication Laboratory
Graduate School of Information Science

Question

2-dimensional data samples are observed as follows:



We can set the following types of distributions to any place.



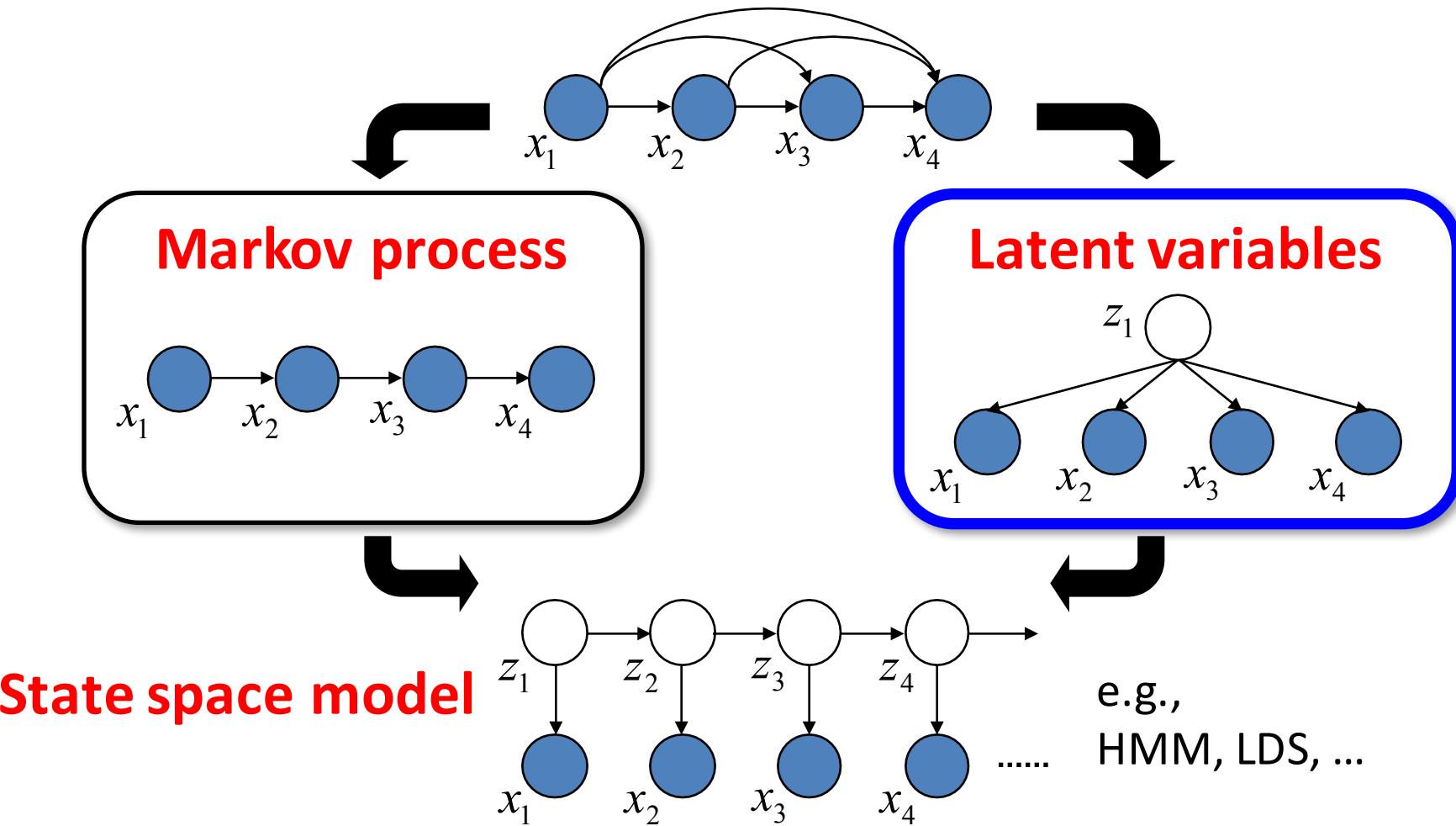
We can control each variance on x-axis and y-axis but we can not rotate this distribution.

Q. How can we model the distribution of these data samples using the above distributions?

After this class, you can answer these questions!

Basic Approaches

How to efficiently model joint probability of high-dimensional data
(i.e., sequential data)



Towards Understanding Latent Variables

- In this class, we focus on a **mixture model** using a **discrete latent variable**.
- We assume that the length of sequential data is constant.

e.g., if the length of each data is 5,

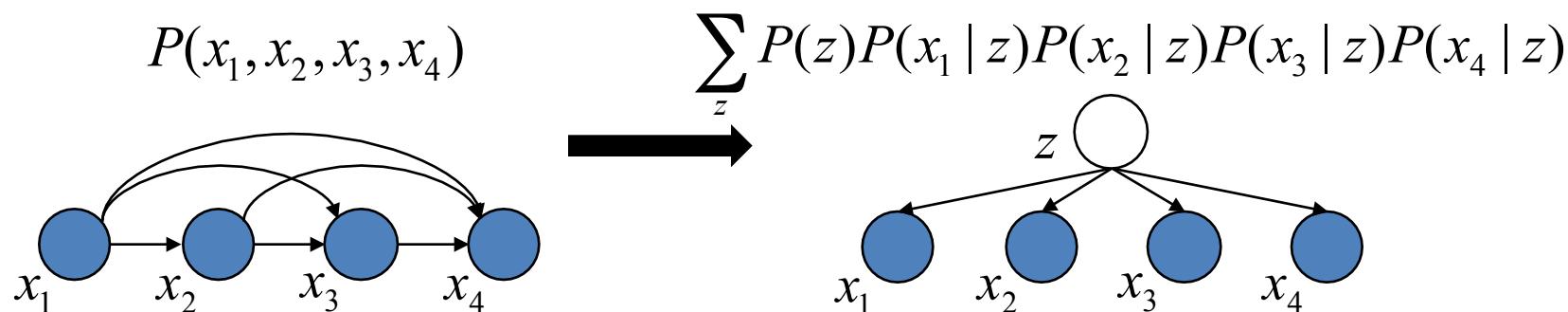
Data 1: { 0.1, 0.2, -0.3, 1.3, -5.8 }

Data 2: { 3.2, -4.8, -0.5, 0.6, -1.3 }

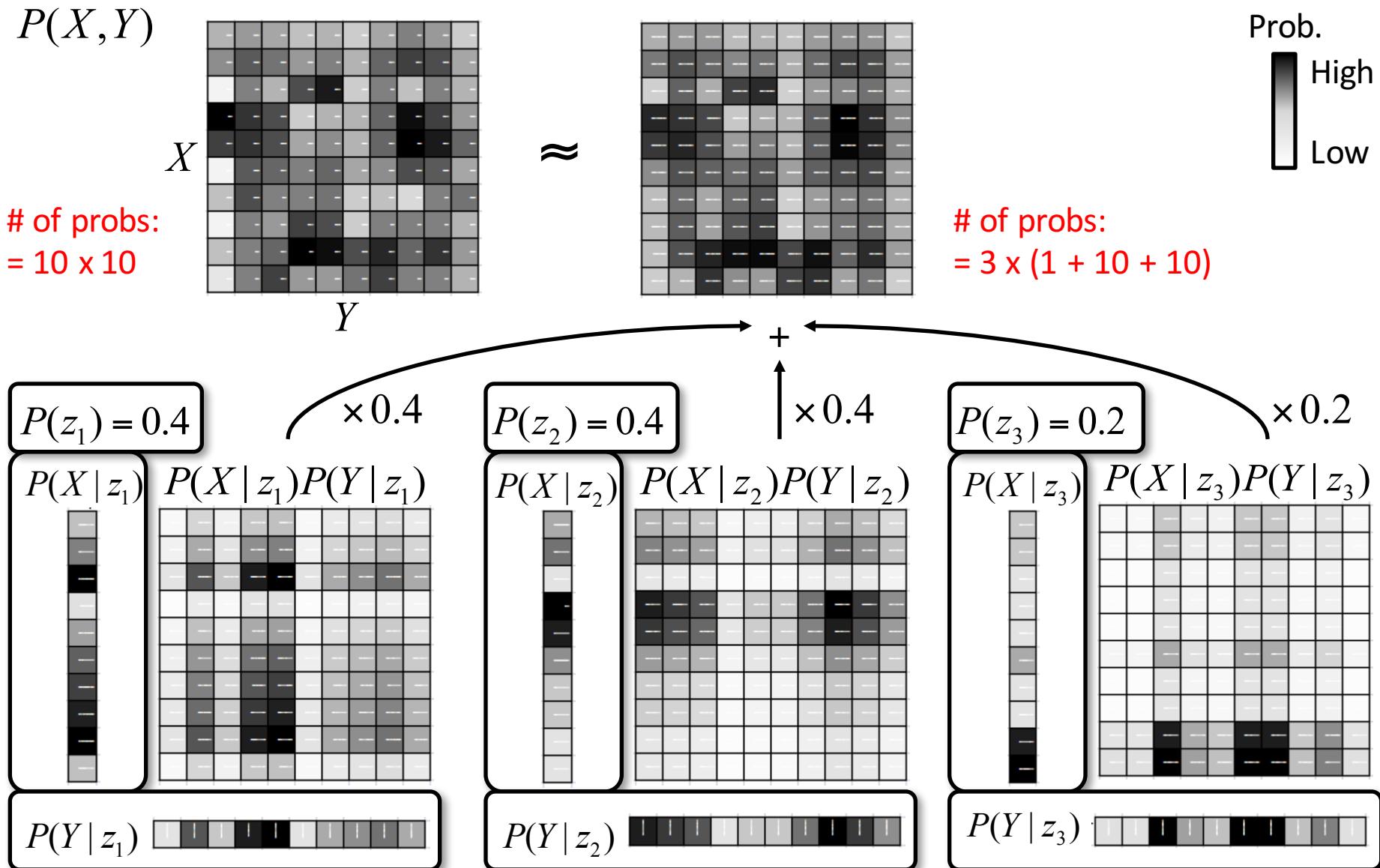
⋮

each data sample is represented as a 5 dimensional vector.

- We learn how to model **joint probability density** in such a high-dimensional space using the **mixture model**.



Effectiveness of Mixture Model

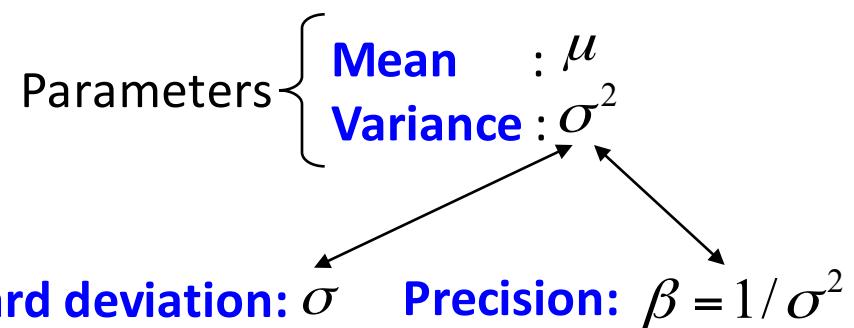
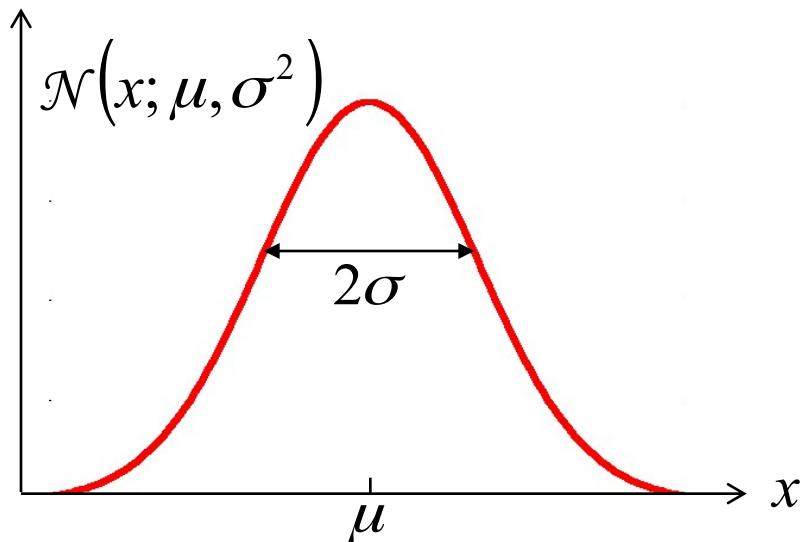


Mixture Model

“Gaussian Mixture Model”

The Gaussian Distribution

- Normal distribution or Gaussian distribution



$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

Conditions to be satisfied:

$$\mathcal{N}(x; \mu, \sigma^2) > 0$$

$$\int_{-\infty}^{\infty} \mathcal{N}(x; \mu, \sigma^2) dx = 1$$

Multivariate Gaussian Distribution

- Gaussian distribution over a D -dimensional vector \mathbf{x} of continuous variables

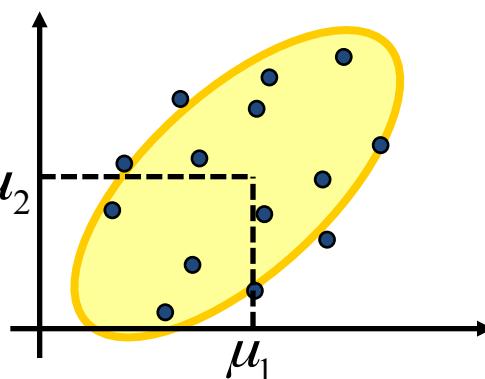
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Parameters $\begin{cases} \text{Mean vector} & : \boldsymbol{\mu} \\ \text{Covariance matrix} & : \boldsymbol{\Sigma} \end{cases}$

Example of 2-dimensional case: $\begin{cases} \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} & \boldsymbol{\Sigma} = \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix} \end{cases}$

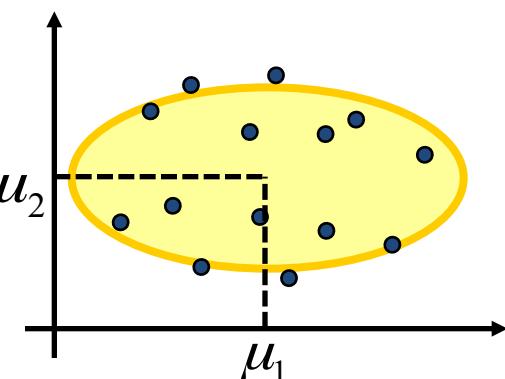
Use of full covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix}$$



Use of diagonal covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} s_{11} & 0 \\ 0 & s_{22} \end{bmatrix}$$

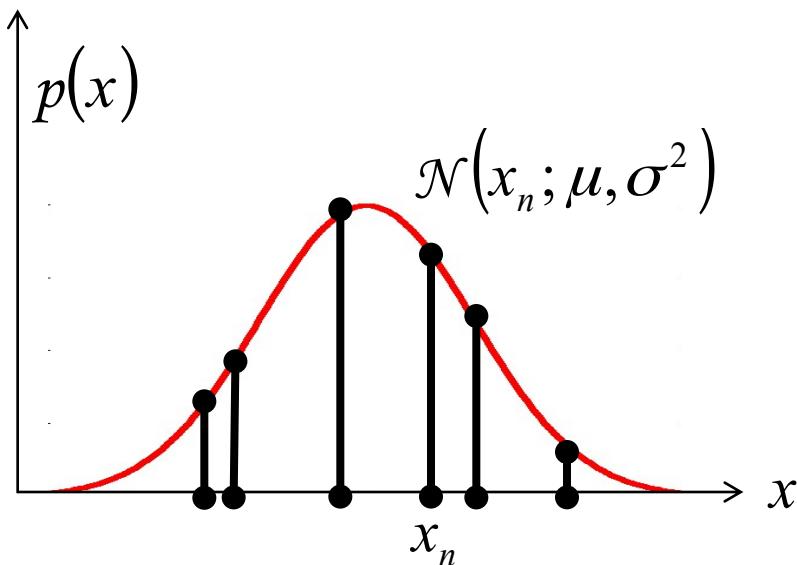


Maximum Likelihood Estimation (MLE)

- Observation data set (N observations of variable x) $X = \{x_1, \dots, x_N\}$
 - Independent and identically distributed: i.i.d.
 - Data points drawn independently from the same distribution
- Likelihood function

$$p(X | \lambda) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \Sigma)$$

Because the data set X is i.i.d.,
 $p(X | \lambda)$ is given by $\prod_{n=1}^N p(x_n | \lambda)$.



Determine parameters $\lambda = \{\mu, \Sigma\}$ using an observation data set by maximizing likelihood function!

ML Estimates of Mean and Covariance

- Determine model parameters by maximizing likelihood function

$$\hat{\lambda} = \arg \max_{\lambda} \ln p(X | \lambda)$$

ML estimate of mean vector:

$$\hat{\mu} = \frac{1}{N} \langle \mathbf{x} \rangle \quad \text{Sample mean vector}$$

The mean vector is the center of the data points

ML estimate of covariance matrix:

$$\hat{\Sigma} = \frac{1}{N} \langle \mathbf{x} \mathbf{x}^T \rangle - \hat{\mu} \hat{\mu}^T \quad \text{Sample covariance}$$

The covariance matrix is the average square of the samples, minus the square of the mean

Sufficient Statistics

Log-scaled likelihood function:

$$\begin{aligned}
 \ln p(X | \lambda) &= \sum_{n=1}^N \ln p(x_n | \lambda) \\
 &= \sum_{n=1}^N -\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} \underline{(x_n - \mu)^\top \Sigma^{-1} (x_n - \mu)} \\
 &\propto \frac{N}{2} \left(\ln |\Sigma^{-1}| - \underline{\mu^\top \Sigma^{-1} \mu} \right) + \frac{1}{2} \underbrace{\left(\sum_{n=1}^N x_n \right)^\top \Sigma^{-1} \mu}_{\text{blue}} + \frac{1}{2} \underline{\mu^\top \Sigma^{-1} \left(\sum_{n=1}^N x_n \right)} - \frac{1}{2} \underbrace{\text{tr} \left[\Sigma^{-1} \sum_{n=1}^N x_n x_n^\top \right]}_{\text{red}} \\
 &= \frac{N}{2} \left(\ln |\Sigma^{-1}| - \underline{\mu^\top \Sigma^{-1} \mu} \right) + \frac{1}{2} \underbrace{\langle x \rangle^\top \Sigma^{-1} \mu}_{\text{blue}} + \frac{1}{2} \underline{\mu^\top \Sigma^{-1} \langle x \rangle} - \frac{1}{2} \text{tr} \left[\Sigma^{-1} \underbrace{\langle x x^\top \rangle}_{\text{blue}} \right]
 \end{aligned}$$

$$\text{tr}[ABC] = \text{tr}[BCA] = \text{tr}[CAB]$$



# of samples	:		N
Sum of samples	:		$\langle x \rangle = \sum_{n=1}^N x_n$ ← $\sum_{n=1}^N \underline{n}$
Sum of squared samples:		$\langle x x^\top \rangle = \sum_{n=1}^N x_n x_n^\top$	← $\sum_{n=1}^N \underline{\underline{n}}$

Deriving ML Estimates

- Determine model parameters by maximizing likelihood function

$$\hat{\lambda} = \arg \max_{\lambda} \ln p(X | \lambda)$$

$$\left\{ \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}} \right\} = \arg \max_{\{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}} \frac{N}{2} \left(\ln |\boldsymbol{\Sigma}^{-1}| - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) + \frac{1}{2} \langle \mathbf{x} \rangle^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \langle \mathbf{x} \rangle - \frac{1}{2} \text{tr} [\boldsymbol{\Sigma}^{-1} \langle \mathbf{x} \mathbf{x}^T \rangle]$$

ML estimate of mean vector:

$$\left. \frac{\partial \ln p(X | \lambda)}{\partial \boldsymbol{\mu}} \right|_{\lambda=\hat{\lambda}} = -N \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}} + \hat{\boldsymbol{\Sigma}}^{-1} \langle \mathbf{x} \rangle = \mathbf{0}$$

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \langle \mathbf{x} \rangle \quad \text{Sample mean vector}$$

ML estimate of covariance matrix:

$$\left. \frac{\partial \ln p(X | \lambda)}{\partial \boldsymbol{\Sigma}^{-1}} \right|_{\lambda=\hat{\lambda}} = \frac{N}{2} \hat{\boldsymbol{\Sigma}} - \frac{N}{2} \hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^T + \frac{1}{2} \langle \mathbf{x} \rangle \hat{\boldsymbol{\mu}}^T + \frac{1}{2} \hat{\boldsymbol{\mu}} \langle \mathbf{x} \rangle^T - \frac{1}{2} \langle \mathbf{x} \mathbf{x}^T \rangle = \mathbf{0}$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N} \langle \mathbf{x} \mathbf{x}^T \rangle - \hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^T \quad \text{Sample covariance}$$

Mixture Model

- Probability density function (*p.d.f.*) given by marginal *p.d.f.* of joint *p.d.f.* over latent variables

$$p(\mathbf{x} | \lambda) = \sum_z p(z | \lambda) p(\mathbf{x} | z, \lambda)$$

Latent variables (1-of- M):

$$\mathbf{z} = \{z_1, \dots, z_M\} \quad z_m \in \{0,1\} \quad \sum_{m=1}^M z_m = 1$$

$\mathbf{z} = \{1,0,0\}$ use 1st component

$$p(\mathbf{x} | z_1 = 1, \lambda) \text{ (blue bell)} \times p(z_1 = 1 | \lambda) \longrightarrow p(\mathbf{x} | \lambda)$$

$\mathbf{z} = \{0,1,0\}$ use 2nd component

$$p(\mathbf{x} | z_2 = 1, \lambda) \text{ (blue bell)} \times p(z_2 = 1 | \lambda) \longrightarrow + = p(\mathbf{x} | \lambda)$$

$\mathbf{z} = \{0,0,1\}$ use 3rd component

$$p(\mathbf{x} | z_3 = 1, \lambda) \text{ (blue bell)} \times p(z_3 = 1 | \lambda) \longrightarrow + = p(\mathbf{x} | \lambda)$$

Gaussian Mixture Model (GMM)

- Mixture of multiple Gaussian distributions

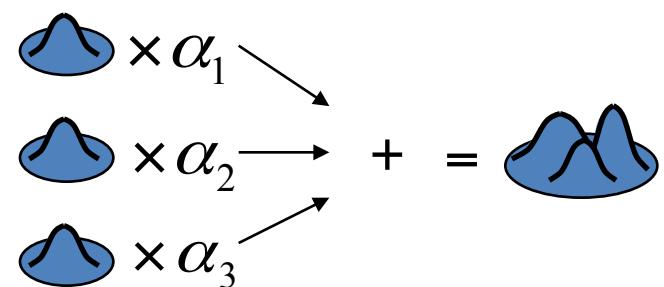
$$\begin{aligned} p(\mathbf{x} | \lambda) &= \sum_z p(z | \lambda) p(\mathbf{x} | z, \lambda) \\ &= \sum_{m=1}^M \alpha_m \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \end{aligned}$$

Prior probability of m^{th} mixture component

$$p(z_m = 1 | \lambda) = \alpha_m \quad 0 \leq \alpha_m \leq 1 \quad \sum_{m=1}^M \alpha_m = 1$$

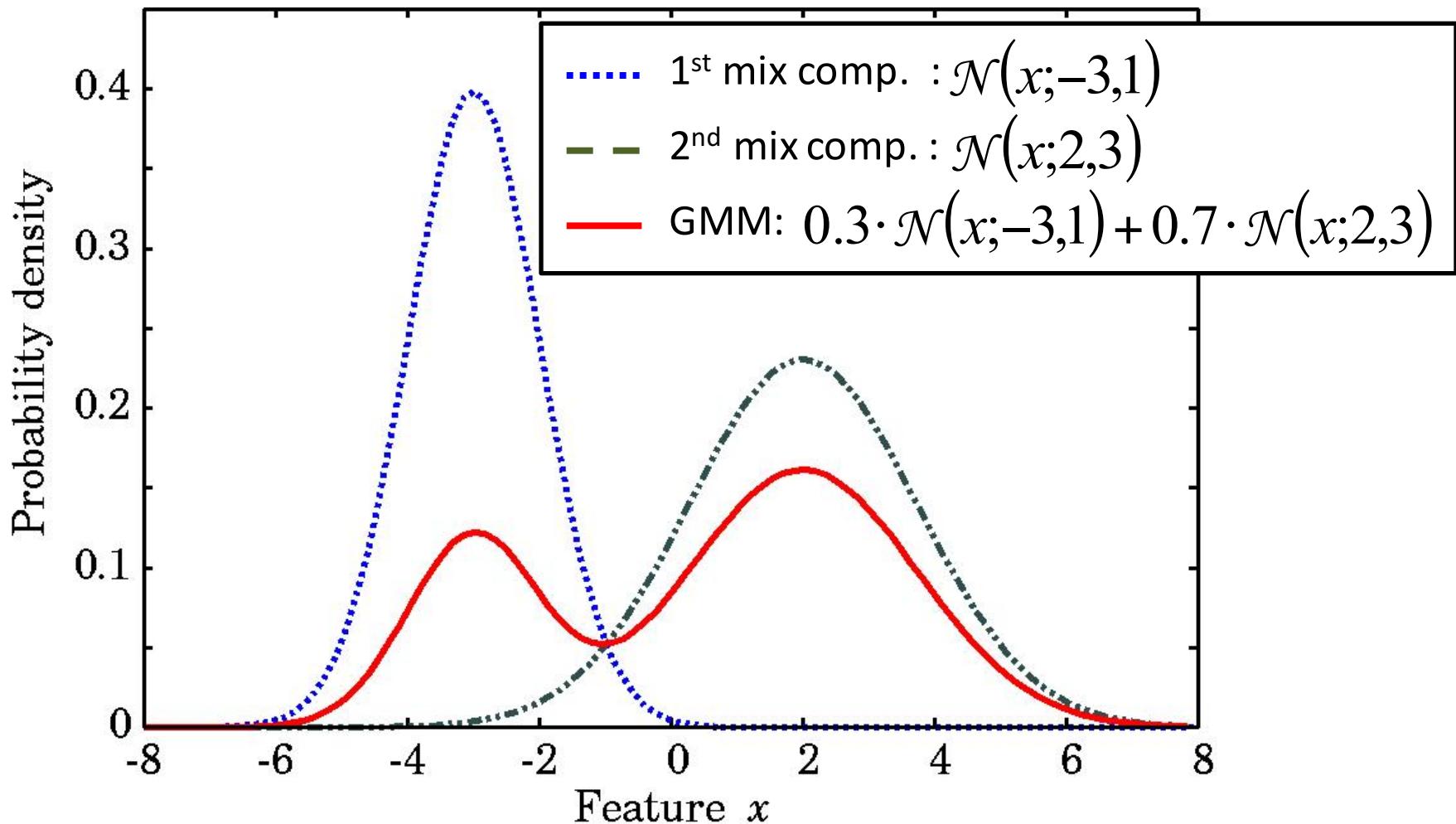
p.d.f. of m^{th} mixture component

$$p(\mathbf{x} | z_m = 1, \lambda) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$



Example of *p.d.f.* of GMM

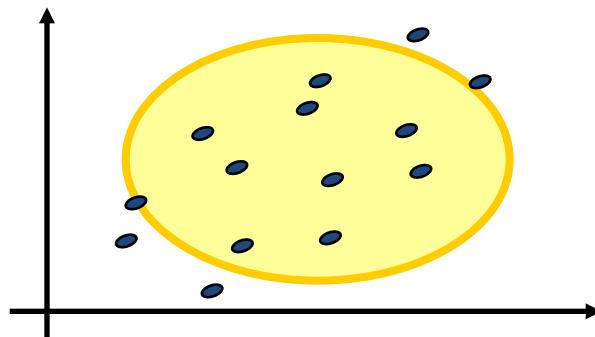
- GMM with 2 mixture components on 1-dimensional space



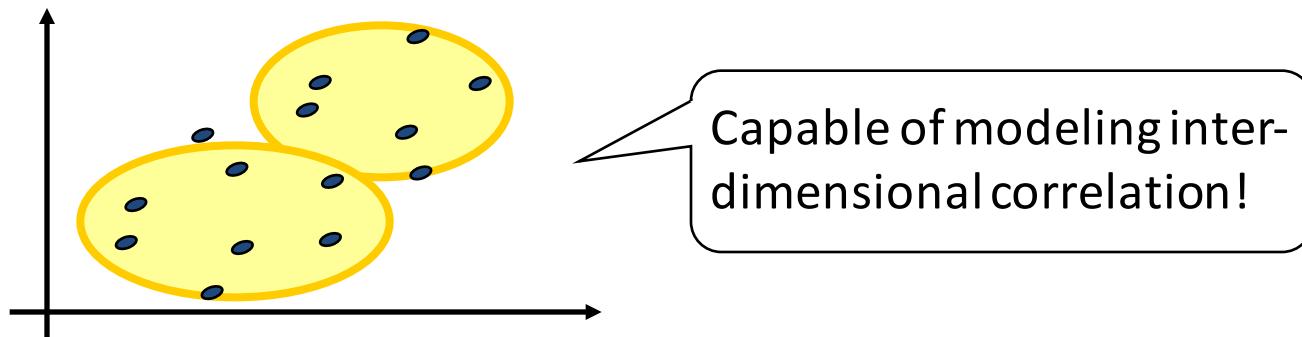
Effectiveness of GMM

- Capable of modeling inter-dimensional correlation even if using diagonal covariance matrices

If using a single Gaussian distribution w/ diagonal covariance matrix...



If using two Gaussian distributions w/ diagonal covariance matrices...



How to Implement MLE?

Log-scaled likelihood of a mixture model:

$$\begin{aligned}\ln p(\mathbf{X} | \lambda) &= \sum_{n=1}^N \ln p(\mathbf{x}_n | \lambda) \\ &= \sum_{n=1}^N \ln \sum_z p(z | \lambda) p(\mathbf{x}_n | z, \lambda) \\ &= \sum_{n=1}^N \ln \sum_{m=1}^M \alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)\end{aligned}$$

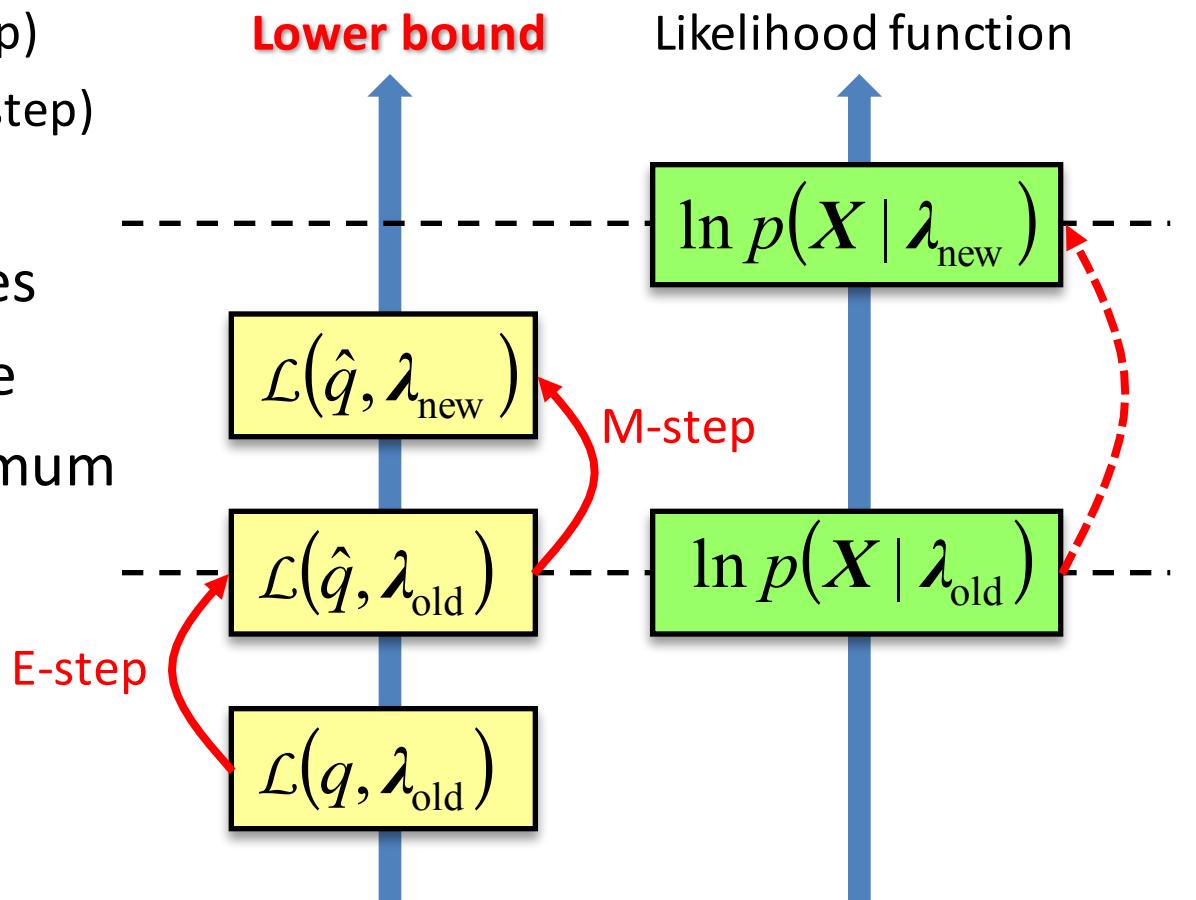
Because of the logarithm of the summation,
its derivatives with respect to model parameters
are not represented as linear equations...
Gradient methods will be necessary...

Can we push the logarithm inside the summation?

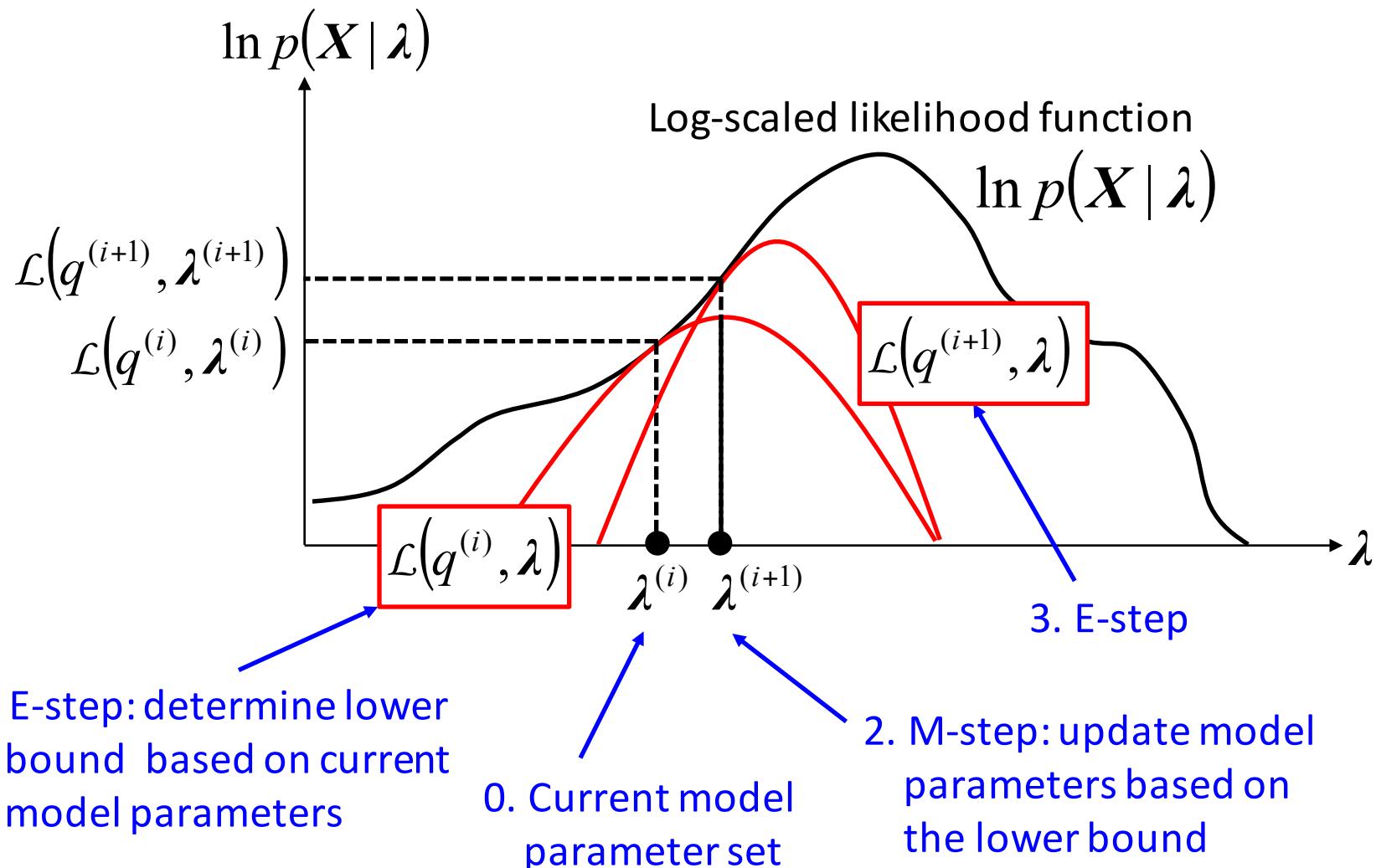
Expectation- Maximization (EM) Algorithm

EM Algorithm

- Iteratively update lower bound of likelihood function through two steps:
 - Expectation step (E-step)
 - Maximization step (M-step)
- Enable closed-form solution of ML estimates
- Guarantee convergence
- Converge to local maximum



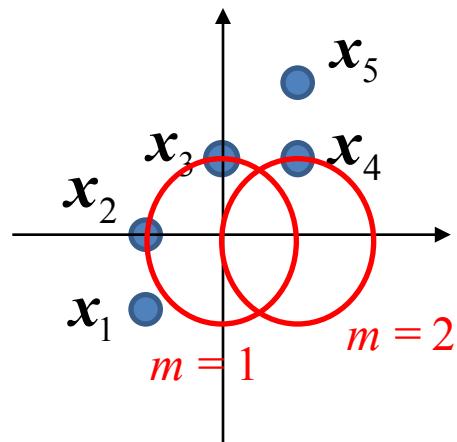
Schematic Image of EM Algorithm



Example in Gaussian Mixture Models

EM Algorithm for GMM

- Initialize: Prepare some training data \mathbf{x}
- Iterate:
 - E-step: Guess which distribution each data point in \mathbf{x} belongs to
 - M-step: Update the parameters based on these guesses



Example Data

- Five 2-Dimensional samples given as training data

$$\mathbf{x}_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix} \quad \mathbf{x}_2 = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \quad \mathbf{x}_3 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \mathbf{x}_4 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{x}_5 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

- Initial model

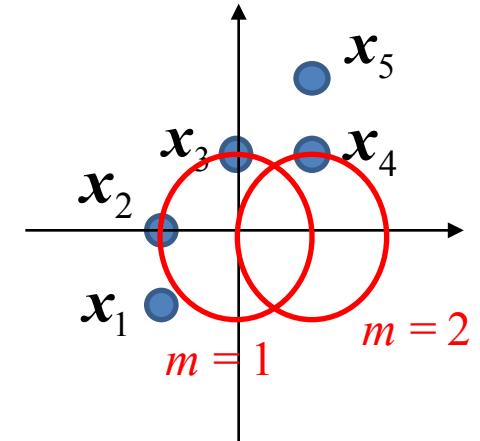
$$p(\mathbf{x}_n | \lambda) = \sum_{m=1}^2 \alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

1st mixture component w/ a diagonal covariance matrix:

$$\alpha_1 = 0.5 \quad \boldsymbol{\mu}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

2nd mixture component w/ a diagonal covariance matrix:

$$\alpha_2 = 0.5 \quad \boldsymbol{\mu}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



E-Step

1. Calculate joint probabilities for individual mixture components

$$p(\mathbf{x}_n, z_n = m | \boldsymbol{\lambda}) = \alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \quad (\text{as } \boldsymbol{\Sigma}_m = \mathbf{I})$$

$$= \alpha_m \frac{1}{\sqrt{(2\pi)^2}} \exp\left(-\frac{1}{2} \sum_{d=1}^2 (\mathbf{x}_{n,d} - \boldsymbol{\mu}_{m,d})^2\right)$$



	$m = 1$	$m = 2$
\mathbf{x}_1	0.029	0.007
\mathbf{x}_2	0.048	0.011
\mathbf{x}_3	0.048	0.029
\mathbf{x}_4	0.029	0.048
\mathbf{x}_5	0.007	0.011

2. Calculate posterior probabilities (i.e., **expected # of samples**)

$$\gamma_{n,m} = p(z_n = m | \mathbf{x}_n, \boldsymbol{\lambda}) = \frac{\alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)}{\sum_{k=1}^M \alpha_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$



	$m = 1$	$m = 2$
\mathbf{x}_1	0.82	0.18
\mathbf{x}_2	0.82	0.18
\mathbf{x}_3	0.62	0.38
\mathbf{x}_4	0.38	0.62
\mathbf{x}_5	0.38	0.62

Calculation of Sufficient Statistics

3. Calculate sufficient statistics

			
	$m = 1$	$m = 2$	
# of samples	$\gamma_m = \sum_{n=1}^N \gamma_{n,m}$	3.01	1.99
Sum of samples	$\langle \mathbf{x} \rangle_m = \sum_{n=1}^N \gamma_{n,m} \mathbf{x}_n$	$\begin{bmatrix} -0.88 \\ 0.94 \end{bmatrix}$	$\begin{bmatrix} 0.88 \\ 2.06 \end{bmatrix}$
Sum of squared samples	$\langle \mathbf{x} \mathbf{x}^\top \rangle_m = \sum_{n=1}^N \gamma_{n,m} \mathbf{x}_n \mathbf{x}_n^\top$	$\begin{bmatrix} 2.39 & * \\ * & 3.33 \end{bmatrix}$	$\begin{bmatrix} 1.61 & * \\ * & 3.67 \end{bmatrix}$

M-Step

4. Update model parameters

Mixture weights

$$\hat{\alpha}_m = \frac{\gamma_m}{\sum_{n=1}^M \gamma_n}$$

Mean vectors

$$\hat{\mu}_m = \frac{1}{\gamma_m} \langle \mathbf{x} \rangle_m$$

Covariance matrices

$$\hat{\Sigma}_m = \frac{1}{\gamma_m} \langle \mathbf{x} \mathbf{x}^\top \rangle_m - \hat{\mu}_m \hat{\mu}_m^\top$$



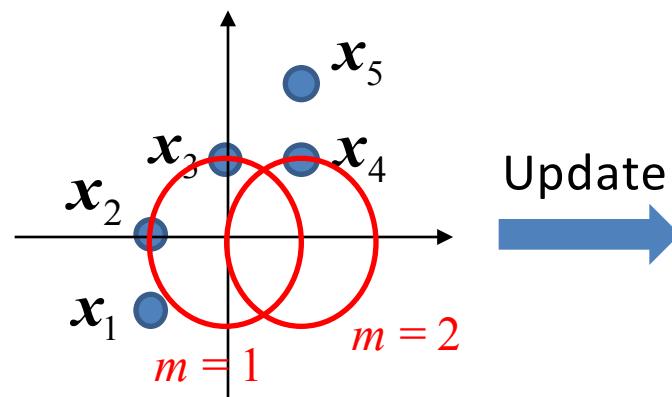
$m = 1$

0.6

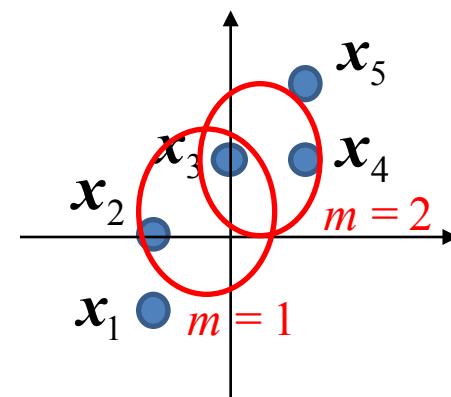


$m = 2$

0.4



Update



Increase of Likelihood

5. Calculate likelihood and compare it before and after update

$$\prod_{n=1}^N p(\mathbf{x}_n | \lambda) = \prod_{n=1}^N \sum_{m=1}^M \alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

Before update:



	$m = 1$	$m = 2$	$P(\mathbf{x}_n \lambda)$
\mathbf{x}_1	0.029	0.007	0.036
\mathbf{x}_2	0.048	0.011	0.059
\mathbf{x}_3	0.048	0.029	0.078
\mathbf{x}_4	0.029	0.048	0.078
\mathbf{x}_5	0.007	0.011	0.017

After update:



	$m = 1$	$m = 2$	$P(\mathbf{x}_n \lambda)$
\mathbf{x}_1	0.034	0.001	0.035
\mathbf{x}_2	0.076	0.008	0.084
\mathbf{x}_3	0.084	0.078	0.163
\mathbf{x}_4	0.028	0.071	0.099
\mathbf{x}_5	0.008	0.039	0.048

$$\frac{1}{N} \sum_{n=1}^N \ln p(\mathbf{x}_n | \lambda) = \underline{-3.07}$$



$$\frac{1}{N} \sum_{n=1}^N \ln p(\mathbf{x}_n | \lambda) = \underline{-2.60}$$

Derivation

Lower Bound of Likelihood Function

- Derivation of lower bound of log-scaled likelihood function

Log-scaled likelihood function:

$$\begin{aligned}\ln p(X | \lambda) &= \sum_{n=1}^N \ln \sum_{z_n} p(x_n, z_n | \lambda) \\ &= \sum_{n=1}^N \ln \sum_{z_n} q(z_n) \frac{p(x_n, z_n | \lambda)}{q(z_n)} \\ &\geq \sum_{n=1}^N \sum_{z_n} q(z_n) \ln \frac{p(x_n, z_n | \lambda)}{q(z_n)} \\ &= \mathcal{L}(q, \lambda)\end{aligned}$$

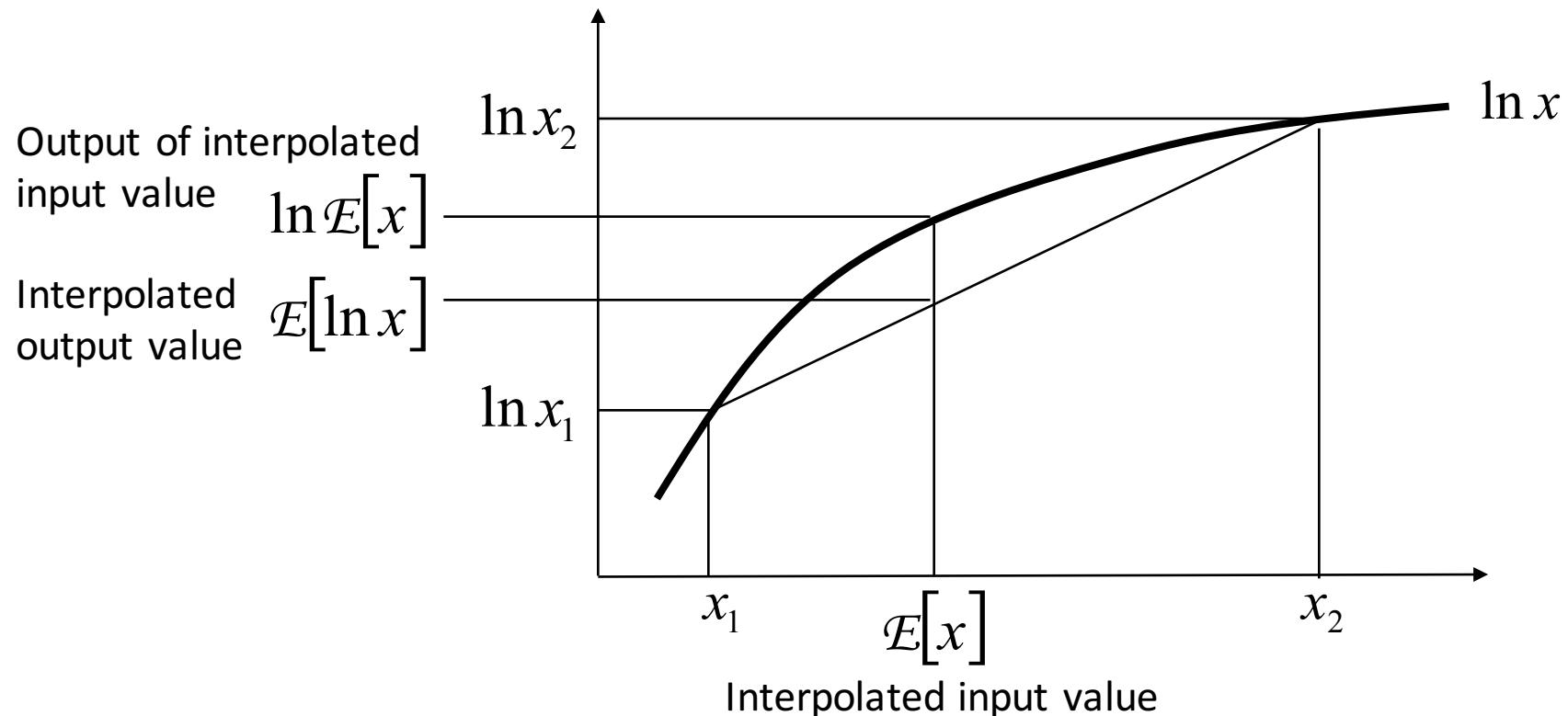
$q(z)$ Probability distribution function of latent variables

Jensen's inequality

Lower bound: $\mathcal{L}(q, \lambda) = \sum_{n=1}^N \sum_{z_n} q(z_n) \ln \frac{p(x_n, z_n | \lambda)}{q(z_n)}$

Jensen's Inequality

$$\ln \sum_n p_n x_n \geq \sum_n p_n \ln x_n \quad \text{subject to } p_n \geq 0, \sum_n p_n = 1$$



Lower Bound as Functional of q

Lower bound:

$$\begin{aligned}\mathcal{L}(q, \lambda) &= \sum_{n=1}^N \sum_{z_n} q(z_n) \ln \frac{p(x_n, z_n | \lambda)}{q(z_n)} \\ &= \sum_{n=1}^N \left\{ \sum_{z_n} q(z_n) \ln p(x_n | \lambda) + \sum_{z_n} q(z_n) \ln \frac{p(z_n | x_n, \lambda)}{q(z_n)} \right\} \\ &= \sum_{n=1}^N \left\{ \underbrace{\ln p(x_n | \lambda)}_{\text{Log-scaled likelihood}} - \underbrace{\sum_{z_n} q(z_n) \ln \frac{q(z_n)}{p(z_n | x_n, \lambda)}}_{\text{KL divergence}} \right\} \\ &= \underbrace{\ln p(X | \lambda)}_{\text{Log-scaled likelihood}} - \underbrace{\sum_{n=1}^N \text{KL}(q(z_n) \| p(z_n | x_n, \lambda))}_{\text{KL divergence}}\end{aligned}$$

“Lower bound” = “log-scaled likelihood” – “KL divergence”

KL Divergence

Discrete variable case:

$$\text{KL}(p \parallel q) = \sum_z p(z) \ln \frac{p(z)}{q(z)}$$

Continuous variable case:

$$\text{KL}(p \parallel q) = \int p(x) \ln \frac{p(x)}{q(x)} dx$$

$$\text{KL}(p \parallel q) = \int p(x) (\ln p(x) - \ln q(x)) dx$$

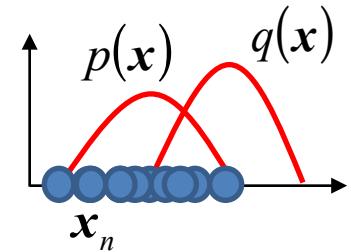
Approximated using samples x_n randomly generated from $p(x)$

$$\approx \frac{1}{N} \sum_{n=1}^N \ln p(x_n) - \ln q(x_n)$$

Log-scaled likelihood
of true distribution

} Difference of
log-scaled likelihoods

Log-scaled likelihood of
another distribution



$$\text{KL}(p \parallel q) \geq 0 \quad \text{If } q = p, \text{KL}(p \parallel q) = 0$$

Lower Bound as Function of λ

Lower bound:

$$\begin{aligned}\mathcal{L}(q, \lambda) &= \sum_{n=1}^N \sum_{z_n} q(z_n) \ln \frac{p(x_n, z_n | \lambda)}{q(z_n)} \\ &= \underbrace{\sum_{n=1}^N \sum_{z_n} q(z_n) \ln p(x_n, z_n | \lambda)}_{\text{Auxiliary function}} - \underbrace{\sum_{z_n} q(z_n) \ln q(z_n)}_{\text{Constant}}\end{aligned}$$

“Lower bound” \propto “Auxiliary function”

EM: Maximization of Lower Bound

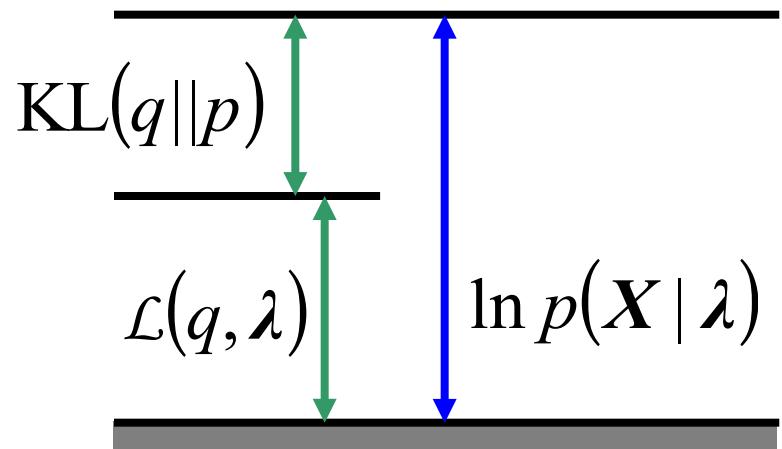
- Lower bound: functional of q and function of λ

E-step: Maximize lower bound with respect to q while fixing λ :

$$\mathcal{L}(q, \lambda) = \ln p(\mathbf{X} | \lambda) - \sum_{n=1}^N \text{KL}(q(z_n) || p(z_n | \mathbf{x}_n, \lambda))$$

M-step: Maximize lower bound
with respect to λ while fixing q

$$\mathcal{L}(q, \lambda) \propto \sum_{n=1}^N \sum_{z_n} q(z_n) \ln p(\mathbf{x}_n, z_n | \lambda)$$



* Need to set initial model parameters λ

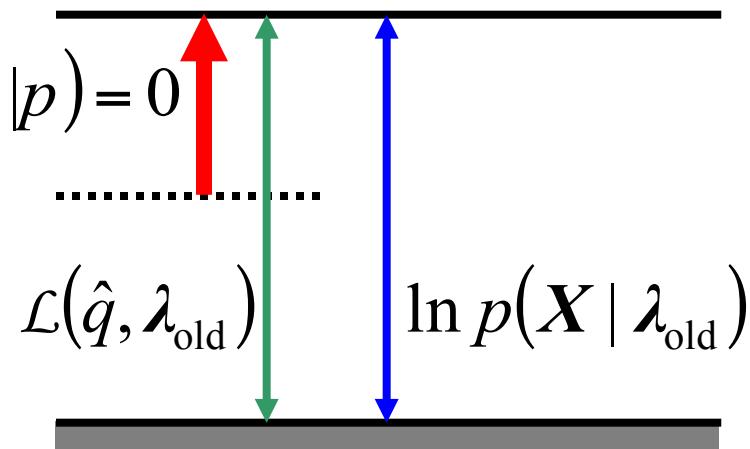
E-Step: Update q

- Set KL divergence to 0 under the fixed model parameters λ_{old}

$$\sum_{n=1}^N \text{KL}(q(z_n) \| p(z_n | \mathbf{x}_n, \lambda_{\text{old}})) = 0 \quad \longrightarrow \quad \hat{q}(z_n) = p(z_n | \mathbf{x}_n, \lambda_{\text{old}})$$

Calculate posterior probabilities of latent variables for each sample

$$\begin{aligned}
 & p(z_{n,m} = 1 | \mathbf{x}_n, \lambda_{\text{old}}) \\
 &= \frac{p(z_{n,m} = 1 | \lambda_{\text{old}}) p(\mathbf{x}_n | z_{n,m} = 1, \lambda_{\text{old}})}{\sum_{m=1}^M p(z_{n,m} = 1 | \lambda_{\text{old}}) p(\mathbf{x}_n | z_{n,m} = 1, \lambda_{\text{old}})} \\
 &= \frac{\alpha_m \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)}{\sum_{k=1}^M \alpha_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \\
 &= \gamma_{n,m}
 \end{aligned}$$



M-Step: Update λ

- Maximize auxiliary function with respect to model parameters λ_{new}

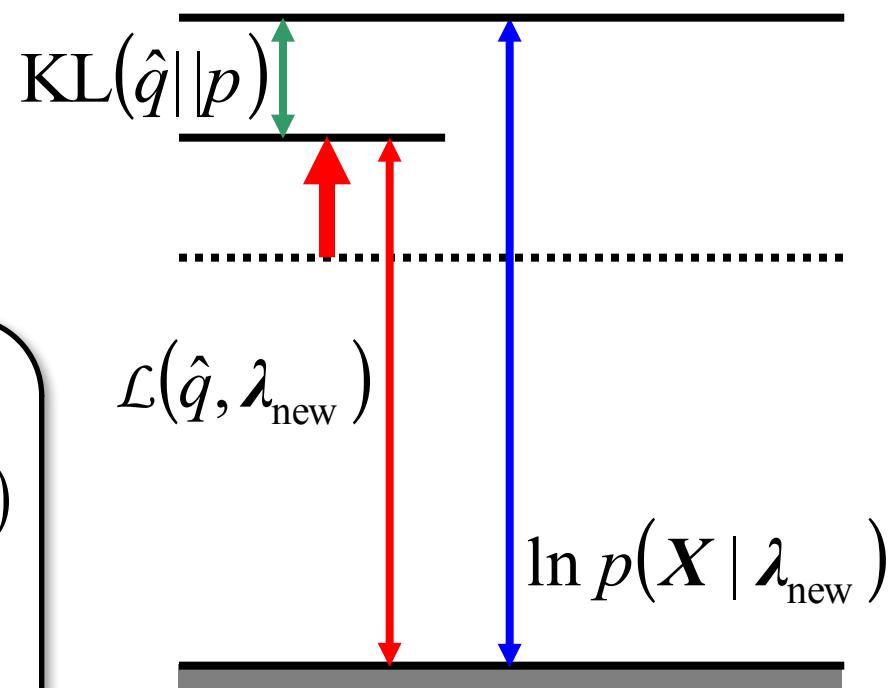
$$\mathcal{L}(\hat{q}, \lambda) \propto \sum_{n=1}^N \sum_{z_n} \hat{q}(z_n) \ln p(x_n, z_n | \lambda)$$

Auxiliary function

$$Q(\lambda_{\text{new}}, \lambda_{\text{old}})$$

$$= \sum_{n=1}^N \sum_{z_n} p(z_n | x_n, \lambda_{\text{old}}) \ln p(x_n, z_n | \lambda_{\text{new}})$$

$$= \sum_{n=1}^N \sum_{m=1}^M \gamma_{n,m} \left\{ \ln \alpha_m + \ln \mathcal{N}(x_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \right\}$$



Sufficient Statistics

Auxiliary function:

$$\begin{aligned}
 Q(\lambda_{\text{new}}, \lambda_{\text{old}}) &= \sum_{n=1}^N \sum_{m=1}^M \gamma_{n,m} \left\{ \ln \alpha_m + \underbrace{\ln \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)} \right\} \\
 &= \sum_{m=1}^M \underbrace{\gamma_m \left\{ \ln \alpha_m + \frac{1}{2} \ln |\boldsymbol{\Sigma}_m^{-1}| - \frac{1}{2} \boldsymbol{\mu}_m^\top \boldsymbol{\Sigma}_m^{-1} \boldsymbol{\mu}_m \right\}} \\
 &\quad + \sum_{m=1}^M \underbrace{\left\{ \frac{1}{2} \boldsymbol{\mu}_m^\top \boldsymbol{\Sigma}_m^{-1} \langle \mathbf{x} \rangle_m + \frac{1}{2} \langle \mathbf{x} \rangle_m^\top \boldsymbol{\Sigma}_m^{-1} \boldsymbol{\mu}_m - \frac{1}{2} \text{tr} \left[\boldsymbol{\Sigma}_m^{-1} \langle \mathbf{x} \mathbf{x}^\top \rangle_m \right] \right\}}
 \end{aligned}$$

For each mixture component,

of samples : \blacksquare_m $\gamma_m = \sum_{n=1}^N \gamma_{n,m}$ $\xleftarrow{\quad} \sum_{n=1}^N \blacksquare_m$

Sum of samples : \blacksquare_m $\langle \mathbf{x} \rangle_m = \sum_{n=1}^N \gamma_{n,m} \mathbf{x}_n$ $\xleftarrow{\quad} \sum_{n=1}^N \blacksquare_m \blacksquare_n$

Sum of squared samples: \blacksquare_m $\langle \mathbf{x} \mathbf{x}^\top \rangle_m = \sum_{n=1}^N \gamma_{n,m} \mathbf{x}_n \mathbf{x}_n^\top$ $\xleftarrow{\quad} \sum_{n=1}^N \blacksquare_m \blacksquare_n \blacksquare_n$

ML Estimates

Auxiliary function:

$$Q(\lambda_{\text{new}}, \lambda_{\text{old}}) = \sum_{m=1}^M \gamma_m \left\{ \ln \alpha_m + \frac{1}{2} \ln |\Sigma_m^{-1}| - \frac{1}{2} \boldsymbol{\mu}_m^\top \Sigma_m^{-1} \boldsymbol{\mu}_m \right\} \\ + \sum_{m=1}^M \left\{ \frac{1}{2} \boldsymbol{\mu}_m^\top \Sigma_m^{-1} \langle \mathbf{x} \rangle_m + \frac{1}{2} \langle \mathbf{x} \rangle_m^\top \Sigma_m^{-1} \boldsymbol{\mu}_m - \frac{1}{2} \text{tr} [\Sigma_m^{-1} \langle \mathbf{x} \mathbf{x}^\top \rangle_m] \right\}$$

$$\frac{\partial Q(\lambda_{\text{new}}, \lambda_{\text{old}}) + \beta \left(\sum_{n=1}^M \alpha_n - 1 \right)}{\partial \alpha_m} \Bigg|_{\lambda=\hat{\lambda}} = 0$$

$$\frac{\partial Q(\lambda_{\text{new}}, \lambda_{\text{old}})}{\partial \boldsymbol{\mu}_m} \Bigg|_{\lambda=\hat{\lambda}} = \boldsymbol{0}$$

$$\frac{\partial Q(\lambda_{\text{new}}, \lambda_{\text{old}})}{\partial \Sigma_m^{-1}} \Bigg|_{\lambda=\hat{\lambda}} = \boldsymbol{0}$$

For each mixture component,

Mixture weights $\hat{\alpha}_m = \frac{\gamma_m}{\sum_{n=1}^M \gamma_n}$

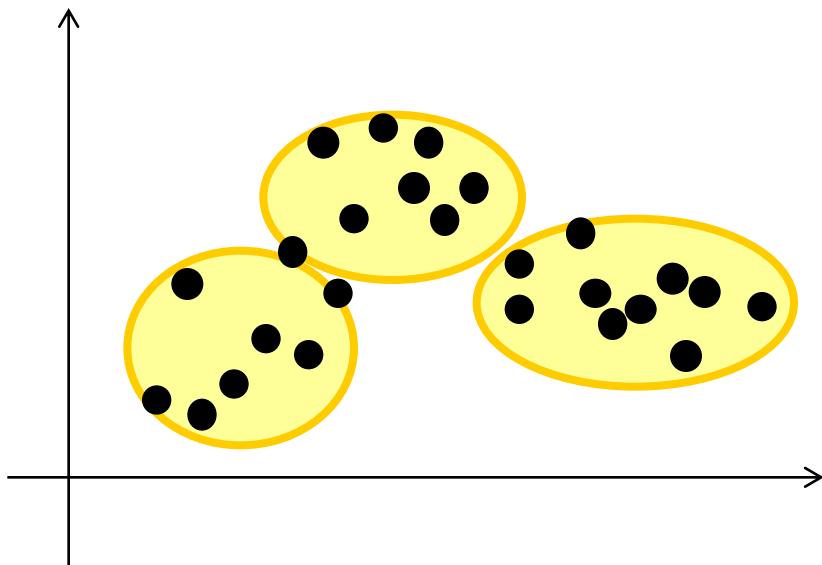
Mean vectors $\hat{\boldsymbol{\mu}}_m = \frac{1}{\gamma_m} \langle \mathbf{x} \rangle_m$

Covariance matrices $\hat{\Sigma}_m = \frac{1}{\gamma_m} \langle \mathbf{x} \mathbf{x}^\top \rangle_m - \hat{\boldsymbol{\mu}}_m \hat{\boldsymbol{\mu}}_m^\top$

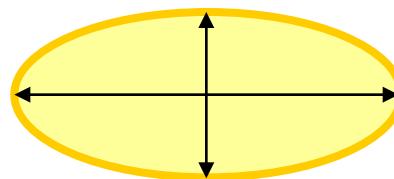
Conclusion

You Can Answer the Question!

2-dimensional data samples are observed as follows:



We can set the following types of distributions to any place.



We can control each variance on x-axis and y-axis but we can not rotate this distribution.

Q. How can we model the distribution of these data samples using the above distributions?

Let's use GMM and optimize it with EM algorithm!