

# Latent variable models and EM-algorithm

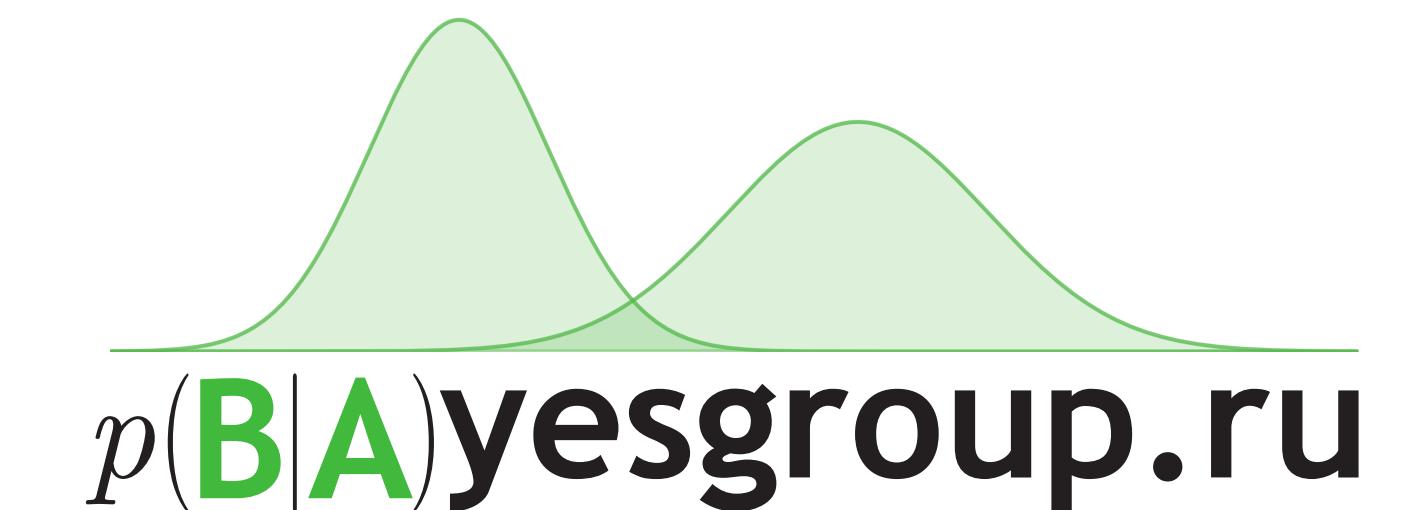
Nadia Chirkova

Higher School of Economics, Samsung-HSE Laboratory  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



# Machine learning overview

- Supervised learning: regression and classification
  - Linear models
  - Decision trees and ensembles
  - K Nearest Neighbors
  - Neural networks
- Unsupervised learning: understanding data structure
  - Clustering
  - Dimensionality reduction, e. g. PCA or autoencoders

# (Bayesian) Machine learning overview

- Supervised learning: regression and classification
  - Linear models — **Bayesian linear regression**
  - Decision trees and ensembles
  - K Nearest Neighbors — **Gaussian processes**
  - Neural networks — **Bayesian neural networks**
- Unsupervised learning: understanding data structure
  - Clustering
  - Dimensionality reduction, e. g. PCA or autoencoders

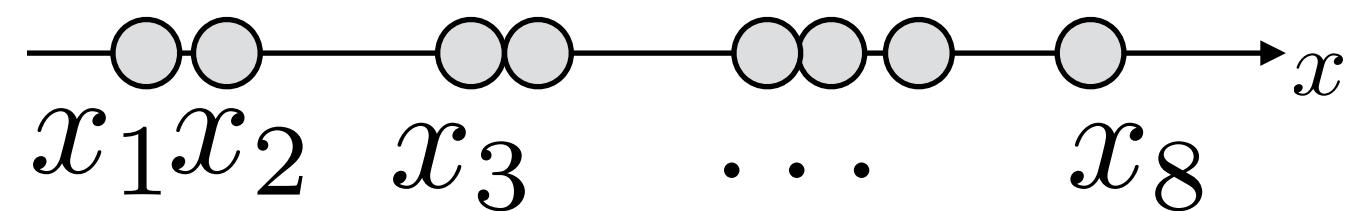
# (Bayesian) Machine learning overview

- Supervised learning: regression and classification
  - Linear models — **Bayesian linear regression**
  - Decision trees and ensembles
  - K Nearest Neighbors — **Gaussian processes**
  - Neural networks — **Bayesian neural networks**
- Unsupervised learning: **probabilistic models of data**
  - Clustering — **Gaussian Mixture Models**
  - Dimensionality reduction, e. g. PCA — **Probabilistic PCA**  
or autoencoders — **Variational autoencoders**

# Probabilistic models

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

1-dim data



images



$x_1$



$x_2$



$x_3$

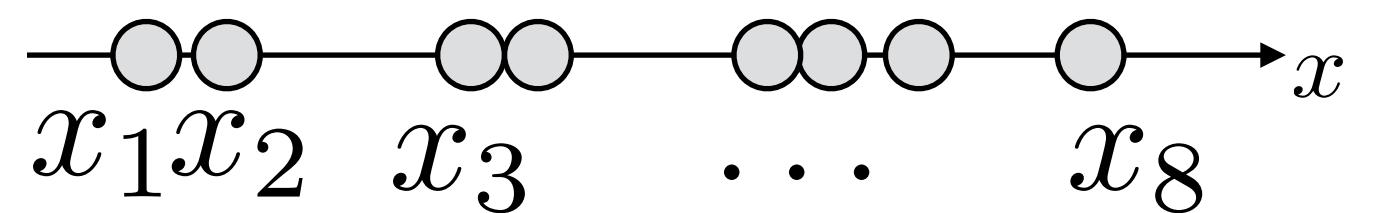
tabular data

	Age	Attrition	BusinessTravel	DailyRate	Department
$x_1$	41	Yes	Travel_Rarely	1102	Sales
$x_2$	49	No	Travel_Frequently	279	Research & Development
$x_3$	37	Yes	Travel_Rarely	1373	Research & Development
$x_4$	33	No	Travel_Frequently	1392	Research & Development

# Probabilistic models

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

1-dim data



images



$x_1$



$x_2$



$x_3$

$p(x|\theta)$  — probabilistic model

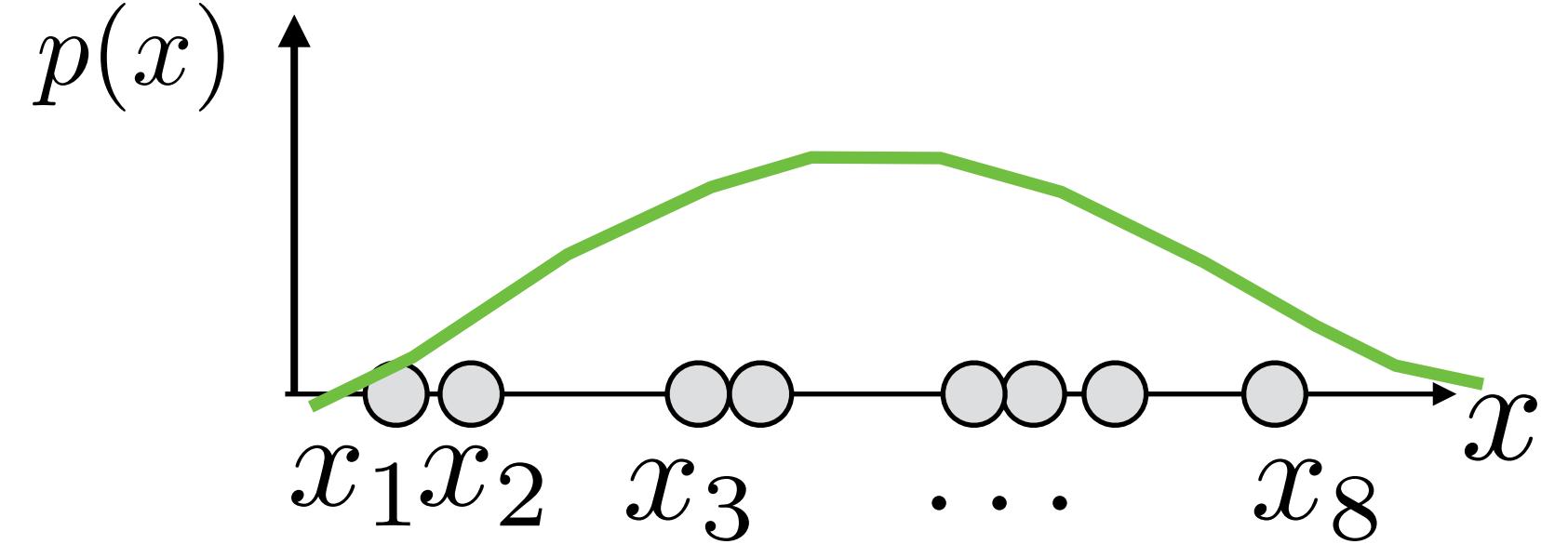
$$\prod_{i=1}^N p(x_i|\theta) \rightarrow \max_{\theta}$$

tabular data

	Age	Attrition	BusinessTravel	DailyRate	Department
$x_1$	41	Yes	Travel_Rarely	1102	Sales
$x_2$	49	No	Travel_Frequently	279	Research & Development
$x_3$	37	Yes	Travel_Rarely	1373	Research & Development
$x_4$	33	No	Travel_Frequently	1392	Research & Development

# Probabilistic models

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$
$$i = 1, \dots, N$$



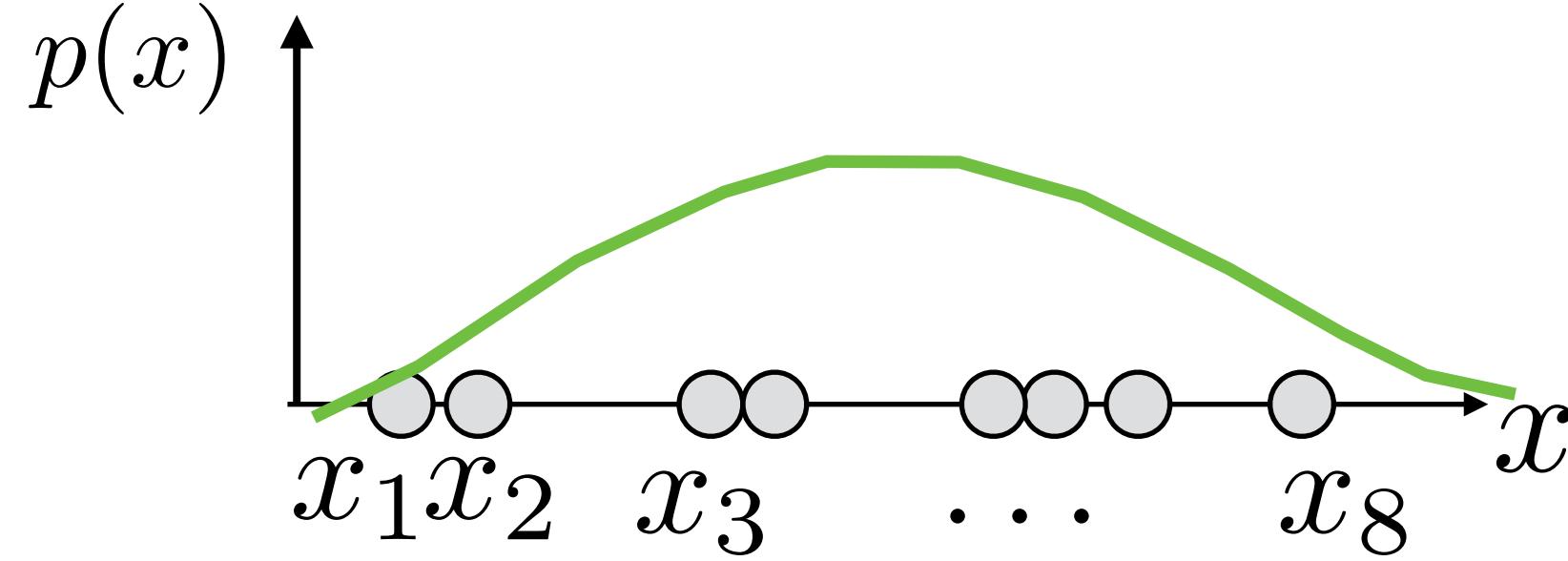
$$p(x|\mu, \sigma) = \mathcal{N}(x|\mu, \sigma^2) =$$
$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

$$\theta = \{\mu, \sigma\}$$

# Probabilistic models

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$i = 1, \dots, N$



Only simple models :(

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

$$\prod_{i=1}^N p(x_i|\mu, \sigma) \rightarrow \max_{\mu, \sigma}$$

Solution:

$\mu$  — sample mean

$\sigma^2$  — sample variance

# Latent variable models

$$X = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$Z = \{z_1, \dots, z_N\}$  — latent variables  
(one per object)

$$p(x, z|\theta) = p(x|z, \theta_1)p(z|\theta_2)$$

$$p(x|\theta) = \int p(x, z|\theta) dz$$

(or sum in case of discrete  $z$ )



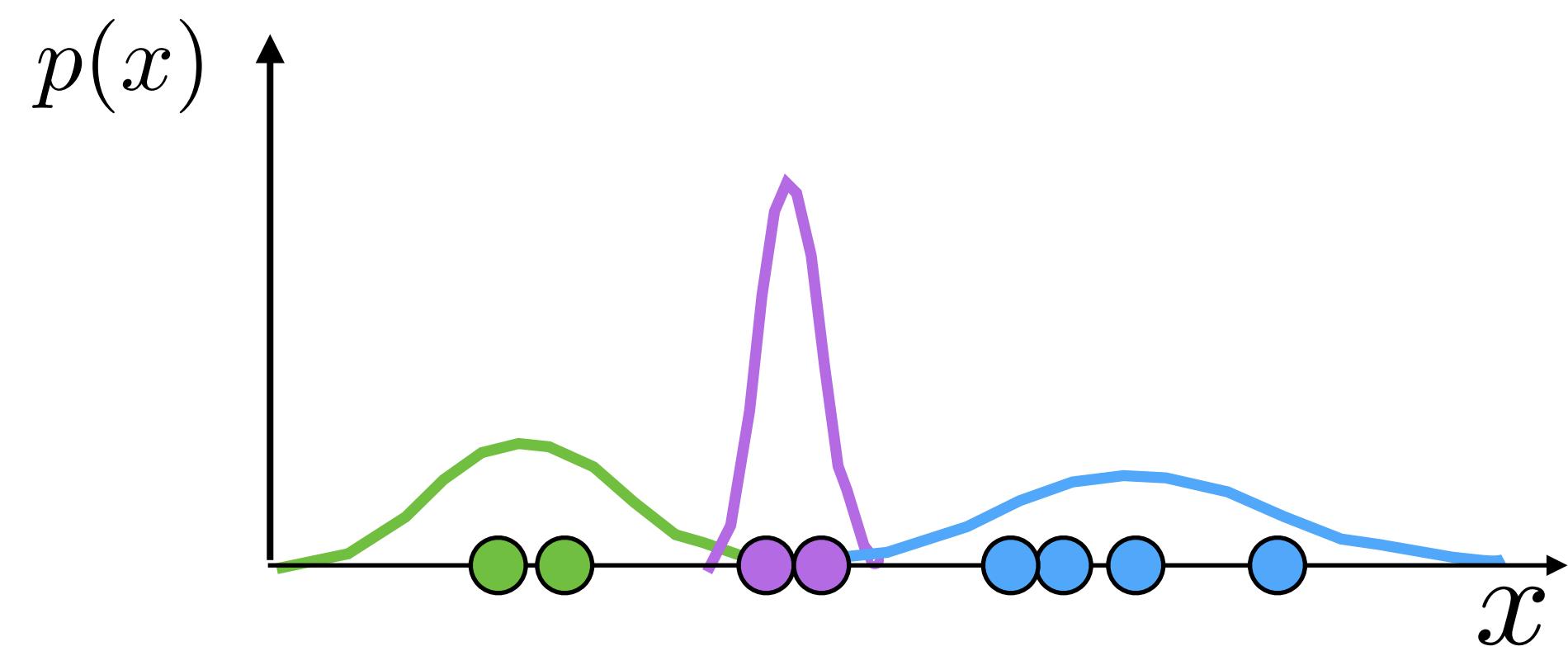
# Latent variable models: example 1

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, 2, 3\}$$

$$i = 1, \dots, N$$



Example 1: assume each object comes from one of three Gaussians

Latent variables  $z$ :  
what Gaussian does an object come from?

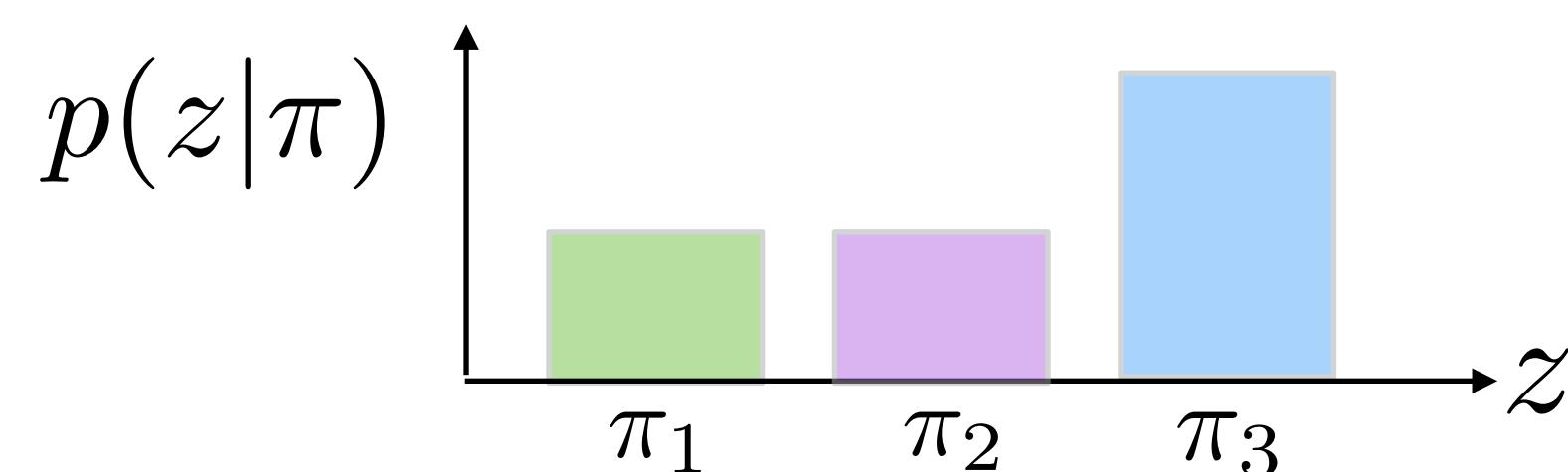
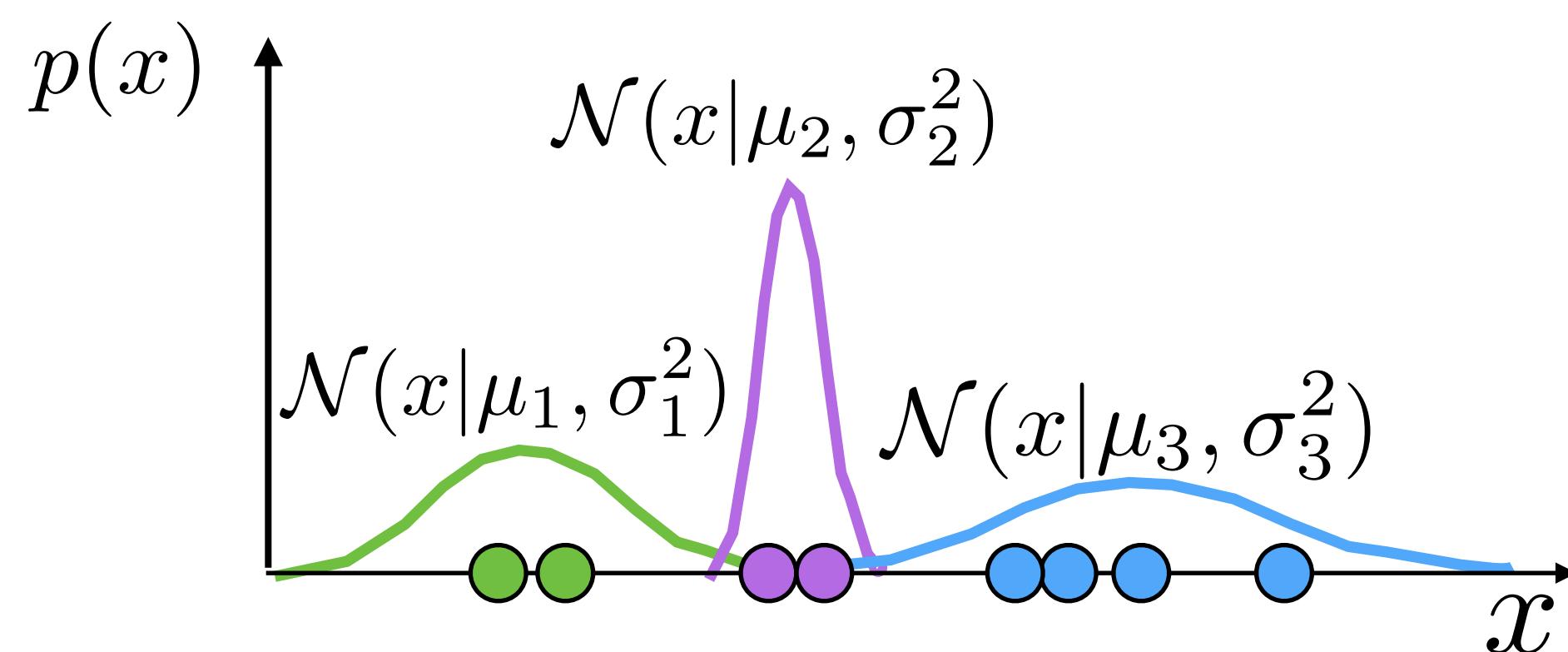
If we knew the values of latent variables,  
estimating  $\theta$  would be trivial

clustering!

# Latent variable models: example 1

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}_{i=1, \dots, N}$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, 2, 3\}_{i=1, \dots, N}$$



$$\begin{aligned} p(x|\underbrace{\pi, \mu, \sigma}_\theta) &= \sum_{z=1}^3 p(x|z, \underbrace{\mu, \sigma}_\theta) p(z|\underbrace{\pi}_\theta) = \\ &= \sum_{z=1}^3 \mathcal{N}(x|\mu_z, \sigma_z^2) \pi_z \end{aligned}$$

Example 1: assume each object comes from one of three Gaussians

Latent variables  $z$ : what Gaussian does an object come from?

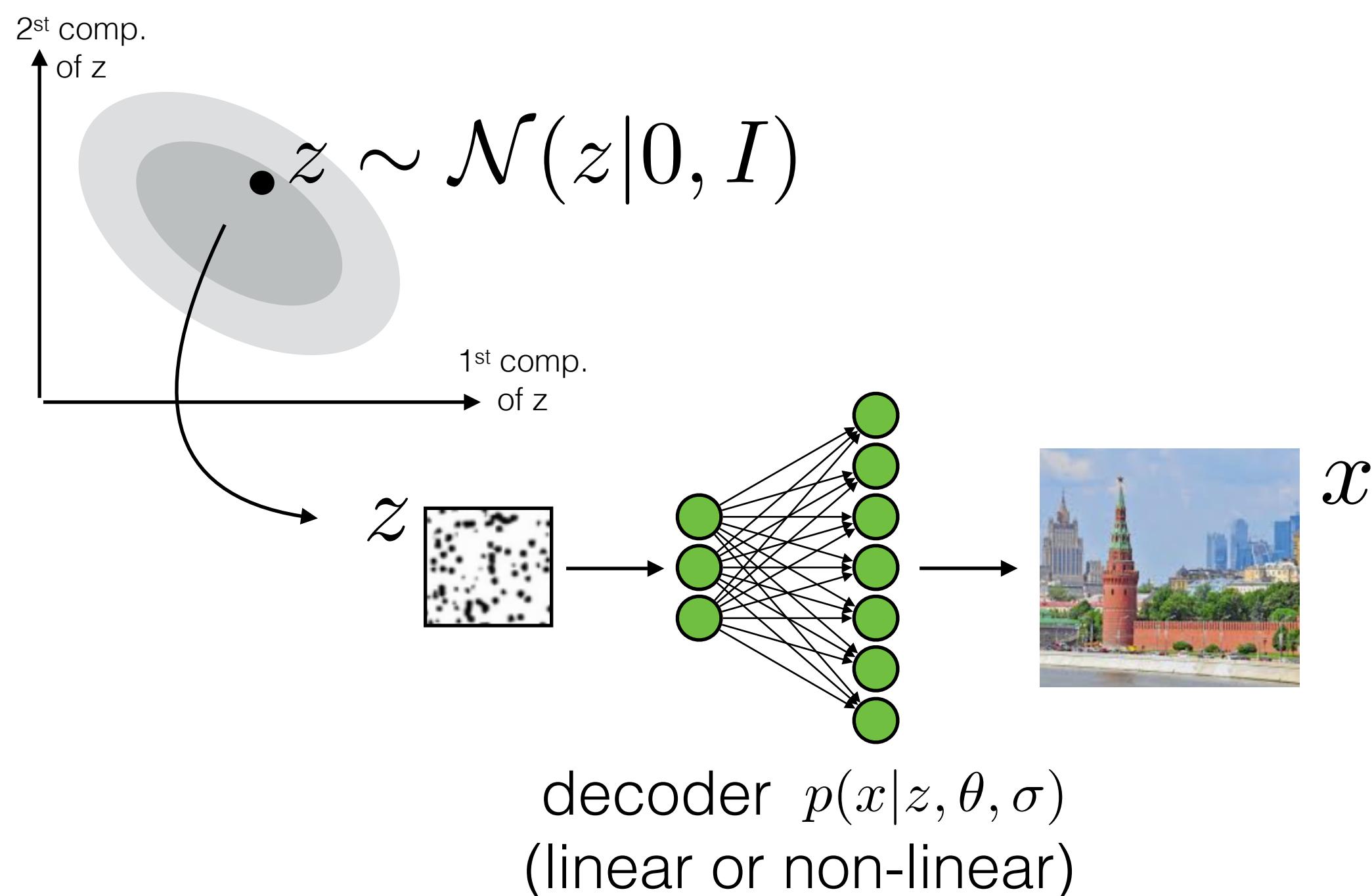
If we knew the values of latent variables, estimating  $\theta$  would be trivial

clustering!

# Latent variable models: example 2

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \mathbb{R}^{d_{\text{small}}} \\ i = 1, \dots, N$$



Example 2: assume each object  
is generated by the decoder  
from an embedding

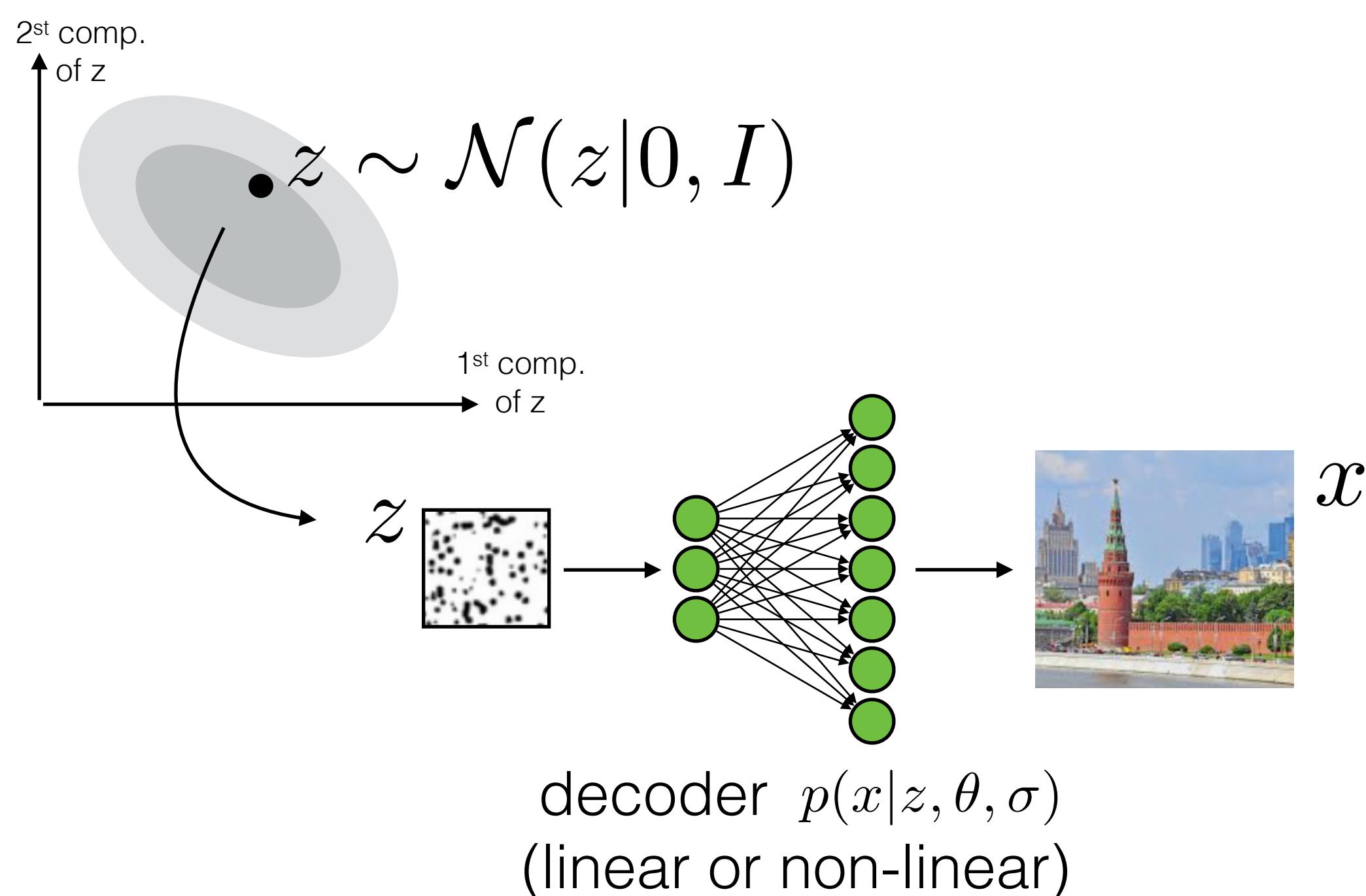
Latent variables  $z$ :  
what is an object's embedding?

If we knew the values of latent variables,  
estimating parameters  $\theta$  would be trivial.

# Latent variable models: example 2

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \mathbb{R}^{d_{\text{small}}} \\ i = 1, \dots, N$$



$$p(x|\theta, \sigma) = \int p(x|z, \theta, \sigma)p(z)dz \\ = \int \mathcal{N}(x|\mu_\theta(z), \sigma^2 I)\mathcal{N}(z|0, I)dz$$

Example 2: assume each object  
is generated by the decoder  
from an embedding

Latent variables  $z$ :  
what is an object's embedding?

If we knew the values of latent variables,  
estimating parameters  $\theta$  would be trivial.

# Latent variable models

$$X = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$Z = \{z_1, \dots, z_N\}$  — latent variables

$$p(x, z|\theta) = p(x|z, \theta_1)p(z|\theta_2)$$

$$p(x|\theta) = \int p(x, z|\theta) dz$$

(or sum in case of discrete  $z$ )

How can l. v. models be useful?

- Understanding data structure:  $p(z|x, \theta) = \frac{p(x|z, \theta_1)p(z|\theta_2)}{\int(x|\tilde{z}, \theta_1)p(\tilde{z}|\theta_2)d\tilde{z}}$  e. g.  
what cluster  
does an object  
come from?
- Generating new objects:  $z \sim p(z|\theta_2), \quad x \sim p(x|z, \theta_1)$  e. g.  
generating  
new images  
using decoder

# Latent variable models

$$X = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$Z = \{z_1, \dots, z_N\}$  — latent variables

$$p(x, z|\theta) = \underbrace{p(x|z, \theta_1)}_{\text{likelihood}} \underbrace{p(z|\theta_2)}_{\text{prior}}$$

$$p(x|\theta) = \int p(x, z|\theta) dz \\ (\text{or sum in case of discrete } z)$$

How can l. v. models be useful?

- Understanding data structure:  $p(z|x, \theta) = \frac{\underbrace{p(x|z, \theta_1)}_{\text{posterior}} \underbrace{p(z|\theta_2)}_{\text{over l. v.}}}{\int(x|\tilde{z}, \theta_1)p(\tilde{z}|\theta_2)d\tilde{z}}$
- Generating new objects:  $z \sim p(z|\theta_2), \quad x \sim p(x|z, \theta_1)$

# Training latent variable models

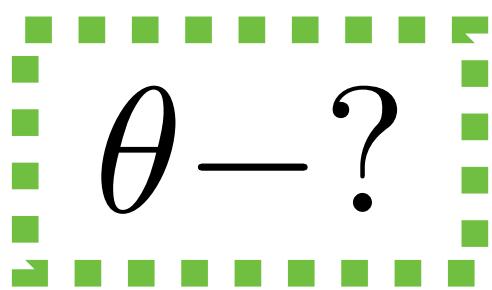
$$X = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$Z = \{z_1, \dots, z_N\}$  — latent variables

$$p(x, z|\theta) = p(x|z, \theta_1)p(z|\theta_2)$$

$$p(x|\theta) = \int p(x, z|\theta) dz$$

(or sum in case of discrete  $z$ )



$$\log p(X|\theta) = \sum_{i=1}^N \log p(x_i|\theta) \rightarrow \max_{\theta}$$

# Variational lower bound

$$\begin{aligned}\underbrace{\log p(X|\theta)}_{\text{lower bound}} &= \int q(Z) \log p(X|\theta) dz = \int q(Z) \log \frac{p(X, Z|\theta)}{p(Z|X, \theta)} dZ = \\ &= \int q(Z) \log \frac{p(X, Z|\theta)q(z)}{p(Z|X, \theta)q(z)} dZ = \\ &= \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ + \int q(Z) \log \frac{q(Z)}{p(Z|X, \theta)} dZ = \\ &= \mathcal{L}(q(Z), \theta) + KL(q(Z)||p(Z|X, \theta))\end{aligned}$$

(holds for any  $q(Z)$ )

# Training latent variable models

$$\begin{aligned}\log p(X|\theta) &= \int q(Z) \log p(X|\theta) dz = \int q(Z) \log \frac{p(X, Z|\theta)}{p(Z|X, \theta)} dZ = \\ &= \int q(Z) \log \frac{p(X, Z|\theta)q(z)}{p(Z|X, \theta)q(z)} dZ = \\ &= \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ + \int q(Z) \log \frac{q(Z)}{p(Z|X, \theta)} dZ = \\ &= \mathcal{L}(q(Z), \theta) + \underbrace{KL(q(Z)||p(Z|X, \theta))}_{\geq 0} \end{aligned}$$

.....  
 $\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta) \rightarrow \max_{q, \theta}$  .....

**Optimization problem  
for training l. v. models**

# EM-algorithm

$$\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta) \rightarrow \max_{q, \theta}$$

Block-coordinate ascent: alternate optimizing w. r. t.  $q(Z)$  and  $\theta$

E-step:  $\mathcal{L}(q(Z), \theta) \rightarrow \max_q$   
(expectation)

M-step:  $\mathcal{L}(q(Z), \theta) \rightarrow \max_\theta$   
(maximization)

# EM-algorithm

$$\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta) \rightarrow \max_{q, \theta}$$

Block-coordinate ascent: alternate optimizing w. r. t.  $q(Z)$  and  $\theta$

E-step:  
(expectation)  $\mathcal{L}(q(Z), \theta) \rightarrow \max_q$

$$\log p(X|\theta) = \mathcal{L}(q(Z), \theta) + KL(q(Z)||p(Z|X, \theta))$$

constant w.r.t.q  $\rightarrow \max_q$

M-step:  
(maximization)  $\mathcal{L}(q(Z), \theta) \rightarrow \max_\theta$

# EM-algorithm

$$\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta) \rightarrow \max_{q, \theta}$$

Block-coordinate ascent: alternate optimizing w. r. t.  $q(Z)$  and  $\theta$

E-step:  $\mathcal{L}(q(Z), \theta) \rightarrow \max_q \quad \Leftrightarrow \quad KL(q(Z)||p(Z|X, \theta)) \rightarrow \min_q$   
(expectation)

$$q(Z) = p(Z|X, \theta)$$

M-step:  $\mathcal{L}(q(Z), \theta) \rightarrow \max_\theta$   
(maximization)

# EM-algorithm

$$\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta) \rightarrow \max_{q, \theta}$$

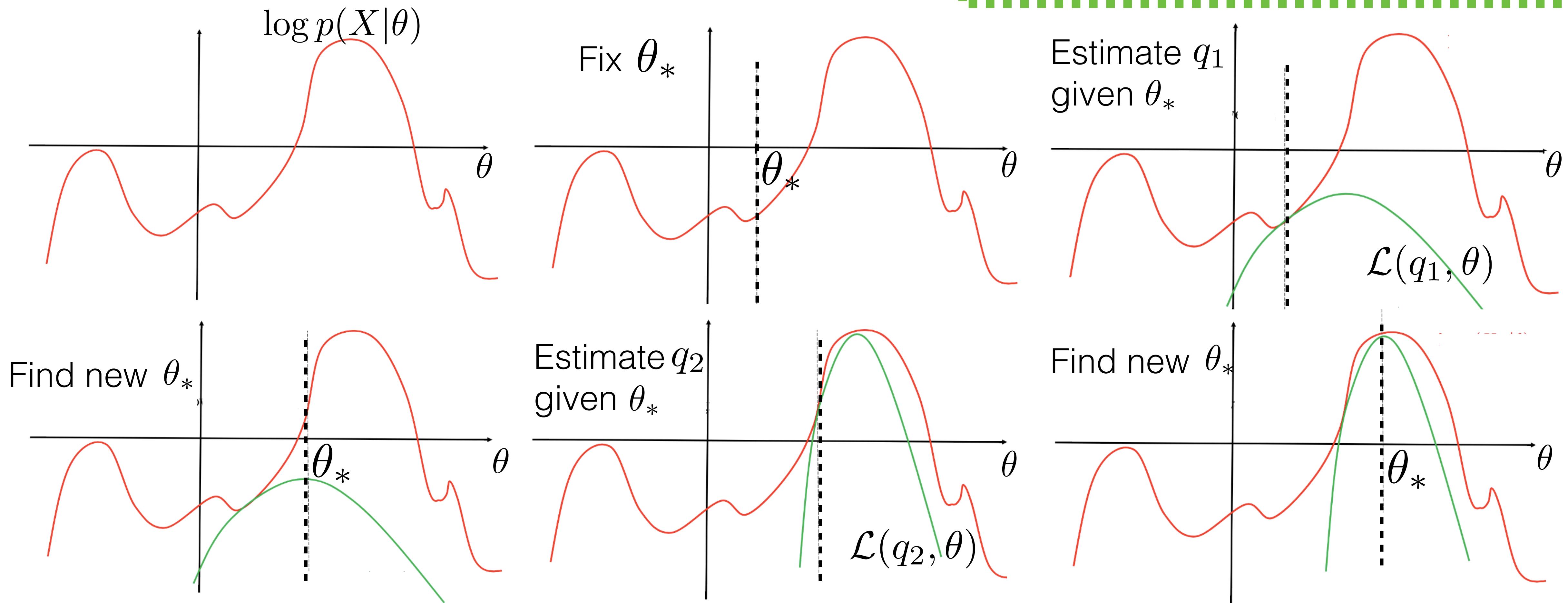
Initialize  $\theta_*$  and repeat

E-step:  $q(Z) = p(Z|X, \theta_*)$

M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

# EM-algorithm: visualization

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

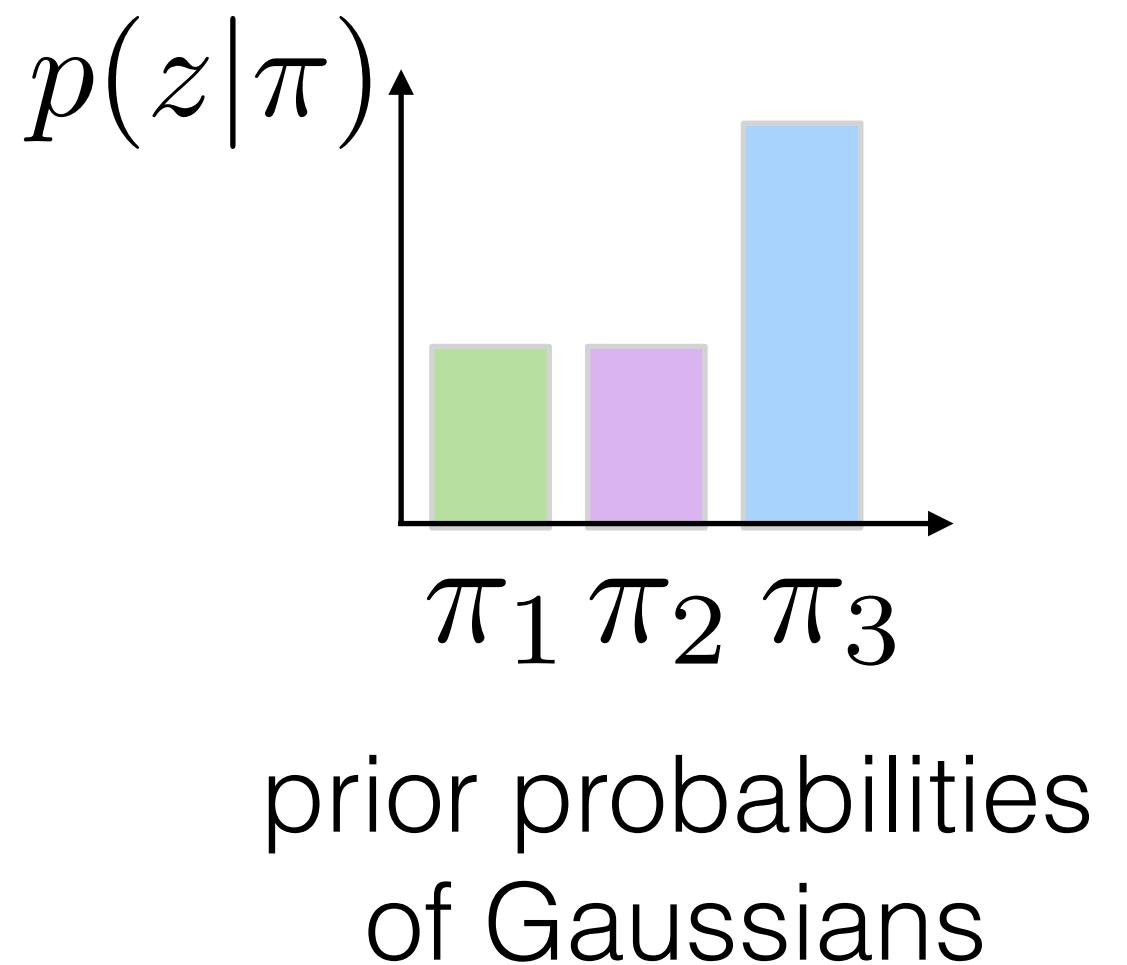
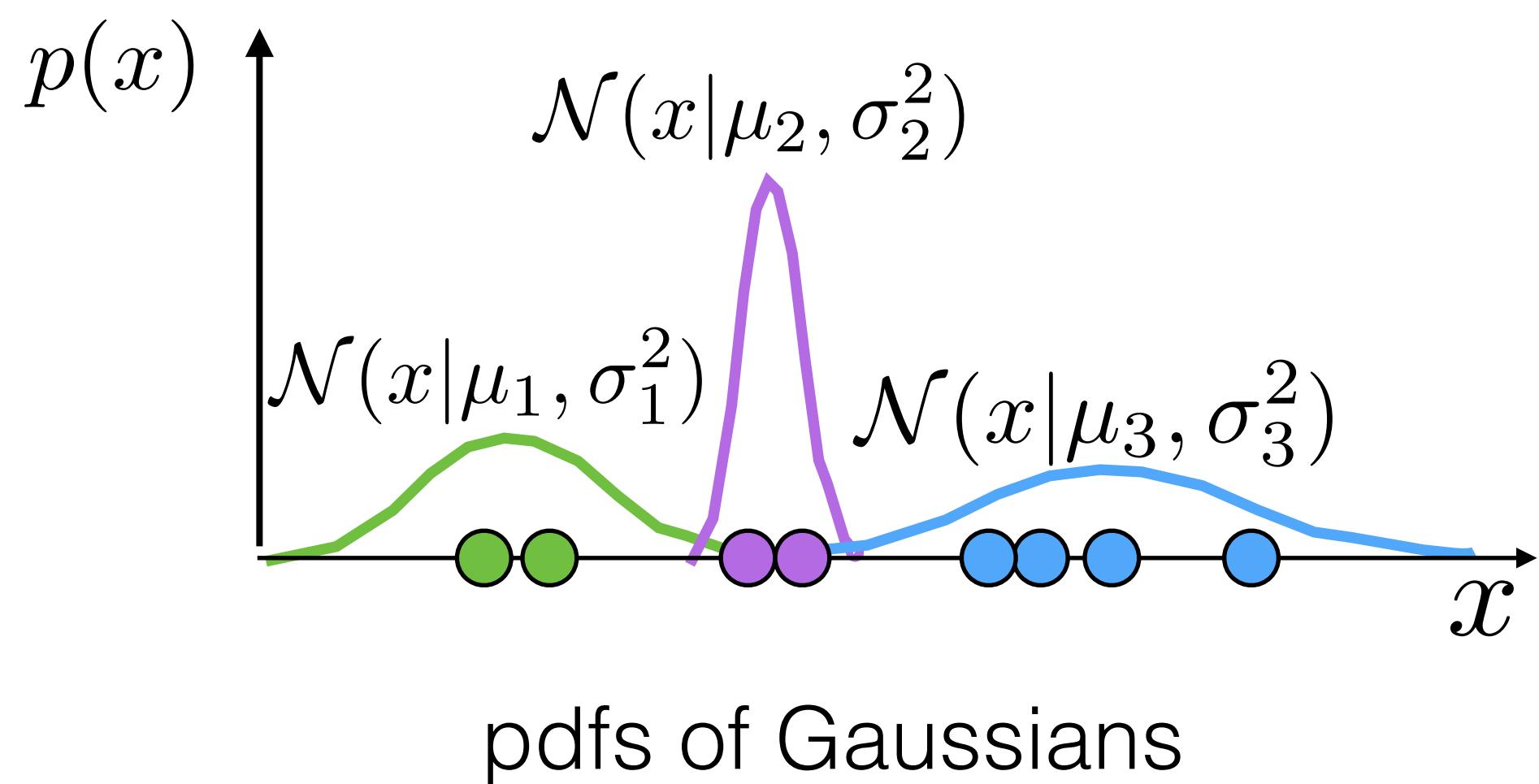


# Gaussian mixture model (1-dim case)

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\} \\ i = 1, \dots, N$$

$$p(x|\pi, \mu, \sigma) = \sum_{z=1}^K p(x|z, \mu, \sigma)p(z|\pi) = \\ = \sum_{z=1}^K \mathcal{N}(x|\mu_z, \sigma_z^2)\pi_z$$



$$\pi = \{\pi_1, \dots, \pi_K\}, \\ \pi_k \geq 0, \sum_{k=1}^K \pi_k = 1$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

E-step: ?

Initialize  $\theta_*$  and repeat

E-step:  $q(Z) = p(Z|X, \theta_*)$

M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

$$\text{E-step: } p(Z|X, \theta_*) = \frac{p(X|Z, \theta_{*1})p(Z|\theta_{*2})}{\int p(X|\tilde{Z}, \theta_{*1})p(\tilde{Z}|\theta_{*2})d\tilde{Z}} =$$

$$= \prod_{i=1}^N \frac{p(x_i|z_i, \theta_{*1})p(z_i|\theta_{*2})}{\int p(x_i|\tilde{z}_i, \theta_{*1})p(\tilde{z}_i|\theta_{*2})d\tilde{z}_i} = \prod_{i=1}^N p(z_i|x_i, \theta_*)$$

(or sum in the denominator case of discrete  $z$ )

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_{\theta} \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

E-step:

Matrix  $q$ :

	$K$			
$N$	$q(z_i)$			
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				

- Initialize  $\theta_*$  and repeat
- E-step:  $q(Z) = p(Z|X, \theta_*)$
- M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

$$\begin{aligned} q_{ik} &= q(z_i = k) = \\ &= p(z_i = k | x_i, \mu, \sigma, \pi) - ? \end{aligned}$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

E-step:

Matrix  $q$ :

	$K$
$N$	$q(z_i)$
1	0.1 0.2 0.3 0.4
2	0.2 0.3 0.4 0.3
3	0.3 0.4 0.3 0.2
4	0.4 0.3 0.2 0.1
5	0.5 0.4 0.3 0.2
6	0.6 0.5 0.4 0.3
7	0.7 0.6 0.5 0.4
8	0.8 0.7 0.6 0.5
9	0.9 0.8 0.7 0.6
10	1.0 0.9 0.8 0.7

Initialize  $\theta_*$  and repeat

E-step:  $q(Z) = p(Z|X, \theta_*)$

M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

$$q_{ik} = p(z_i = k|x_i, \mu, \sigma, \pi) = \frac{p(x_i|z_i = k, \mu, \sigma)p(z_i = k|\pi)}{\sum_{k'=1}^K p(x_i|z_i = k', \mu, \sigma)p(z_i = k'|\pi)}$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

E-step:

Matrix  $q$ :

	$K$			
$N$	$q(z_i)$			
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				

Initialize  $\theta_*$  and repeat

E-step:  $q(Z) = p(Z|X, \theta_*)$

M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

$$q_{ik} = \frac{\mathcal{N}(x|\mu_k, \sigma_k)\pi_k}{\sum_{k'=1}^K \mathcal{N}(x|\mu_{k'}, \sigma_{k'})\pi_{k'}}$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

M-step - ?

$$\mathcal{L}(q(Z), \theta) = \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ = \int q(Z) \log \frac{p(X|Z, \theta_1)p(Z|\theta_2)}{q(Z)} dZ$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

M-step - ?

$$\mathcal{L}(q(Z), \theta) = \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ = \int q(Z) \log \frac{p(X|Z, \theta_1)p(Z|\theta_2)}{q(Z)} dZ$$

$$= \sum_{i=1}^N \int q(z_i) \log \frac{p(x_i|z_i, \theta_1)p(z_i|\theta_2)}{q(z_i)} dz_i = \sum_{i=1}^N \mathcal{L}(q(z_i), \theta)$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

M-step - ?

$$\mathcal{L}(q(z_i), \theta) = ?$$

$$\sum_{i=1}^N \int q(z_i) \log \frac{p(x_i|z_i, \theta_1)p(z_i|\theta_2)}{q(z_i)} dz_i = \sum_{i=1}^N \mathcal{L}(q(z_i), \theta)$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

M-step - ?

$$\mathcal{L}(q(z_i), \mu, \sigma, \pi) = \sum_{k=1}^K q_{ik} \log \frac{\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k}{q_{ik}}$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

# Gaussian mixture model

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}$$

$$i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\}$$

$$i = 1, \dots, N$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$p(x|z, \mu, \sigma) = \mathcal{N}(x|\mu_z, \sigma_z)$$

$$p(z|\pi) = \pi_z$$

M-step:

$$\mathcal{L}(q(z_i), \mu, \sigma, \pi) = \sum_{k=1}^K q_{ik} \log \frac{\mathcal{N}(x_i|\mu_k, \sigma_k)\pi_k}{q_{ik}}$$

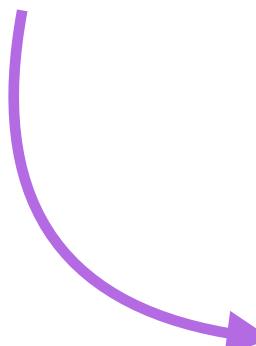
$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i|\mu_k, \sigma_k)\pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \mu_j?$$

# Gaussian mixture model: M-step

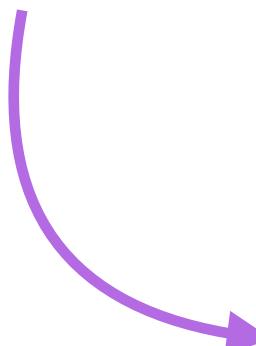
$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \mu_j?$$


$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \left( -\frac{1}{2\sigma_k^2} (x_i - \mu_k)^2 \right) + \text{const}$$

$$\frac{\partial}{\partial \mu_j} :$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \mu_j?$$


$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \left( -\frac{1}{2\sigma_k^2} (x_i - \mu_k)^2 \right) + \text{const}$$

$$\frac{\partial}{\partial \mu_j} : - \sum_{i=1}^N \frac{q_{ij}}{\sigma_j^2} (x_i - \mu_j) = 0$$

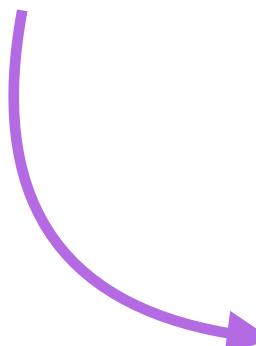
$$\boxed{\mu_j = \frac{\sum_{i=1}^N q_{ij} x_i}{\sum_{i=1}^N q_{ij}}}$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \sigma_j?$$

# Gaussian mixture model: M-step

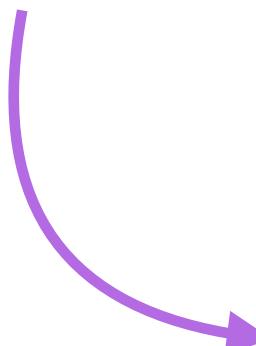
$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \sigma_j?$$


$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \left( -\log \sigma_k - \frac{1}{2\sigma_k^2} (x_i - \mu_k)^2 \right) + \text{const}$$

$$\frac{\partial}{\partial \sigma_j} :$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi} \quad \text{Maximum w.r.t. } \sigma_j?$$


$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \left( -\log \sigma_k - \frac{1}{2\sigma_k^2} (x_i - \mu_k)^2 \right) + \text{const}$$

$$\frac{\partial}{\partial \sigma_j} : \sum_{i=1}^N q_{ij} \left( -\frac{1}{\sigma_j} + \frac{2}{2\sigma_j^3} (x_i - \mu_j)^2 \right) = 0$$

$$\boxed{\sigma_j^2 = \frac{\sum_{i=1}^N q_{ij} (x_i - \mu_j)^2}{\sum_{i=1}^N q_{ij}}}$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\mathcal{N}(x_i | \mu_k, \sigma_k) \pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

Maximum w.r.t.  $\pi$  ?

$$\sum_{k=1}^K \pi_k = 1$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\cancel{\mathcal{N}(x_i | \mu_k, \sigma_k)} \pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

Maximum w.r.t.  $\pi$  ?

$$\sum_{k=1}^K \pi_k = 1$$

Lagrangian!

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\cancel{\mathcal{N}(x_i | \mu_k, \sigma_k)} \pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

Maximum w.r.t.  $\pi$  ?

$$\sum_{k=1}^K \pi_k = 1$$

$$L(\pi, \lambda) = \sum_{i=1}^N \sum_{k=1}^K q_{ik} \log \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial L}{\partial \pi_j} =$$

$$\frac{\partial L}{\partial \lambda} =$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\cancel{\mathcal{N}(x_i | \mu_k, \sigma_k)} \pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

Maximum w.r.t.  $\pi$  ?

$$\sum_{k=1}^K \pi_k = 1$$

$$L(\pi, \lambda) = \sum_{i=1}^N \sum_{k=1}^K q_{ik} \log \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial L}{\partial \pi_j} = \sum_{i=1}^N \frac{q_{ij}}{\pi_j} + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1 = 0$$

$$\boxed{\pi_j = \frac{\sum_{i=1}^N q_{ij}}{\sum_{i=1}^N \sum_{k=1}^K q_{ij}}}$$

# Gaussian mixture model: M-step

$$\sum_{i=1}^N \sum_{k=1}^K q_{ik} \log(\cancel{\mathcal{N}(x_i | \mu_k, \sigma_k)} \pi_k) \rightarrow \max_{\mu, \sigma, \pi}$$

Maximum w.r.t.  $\pi$  ?

$$\sum_{k=1}^K \pi_k = 1$$

$$L(\pi, \lambda) = \sum_{i=1}^N \sum_{k=1}^K q_{ik} \log \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial L}{\partial \pi_j} = \sum_{i=1}^N \frac{q_{ij}}{\pi_j} + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1 = 0$$

$$\boxed{\pi_j = \frac{1}{N} \sum_{i=1}^N q_{ij}}$$

# Gaussian mixture model: EM-algorithm

Initialize  $\mu, \sigma, \pi$  and repeat

$$\text{E-step: } q_{ik} = \frac{\mathcal{N}(x|\mu_k, \sigma_k) \pi_k}{\sum_{k'=1}^K \mathcal{N}(x|\mu_{k'}, \sigma_{k'}) \pi_{k'}}$$

$$\text{M-step: } \mu_k = \frac{\sum_{i=1}^N q_{ik} x_i}{\sum_{i=1}^N q_{ik}} \quad \pi_k = \frac{1}{N} \sum_{i=1}^N q_{ik}$$

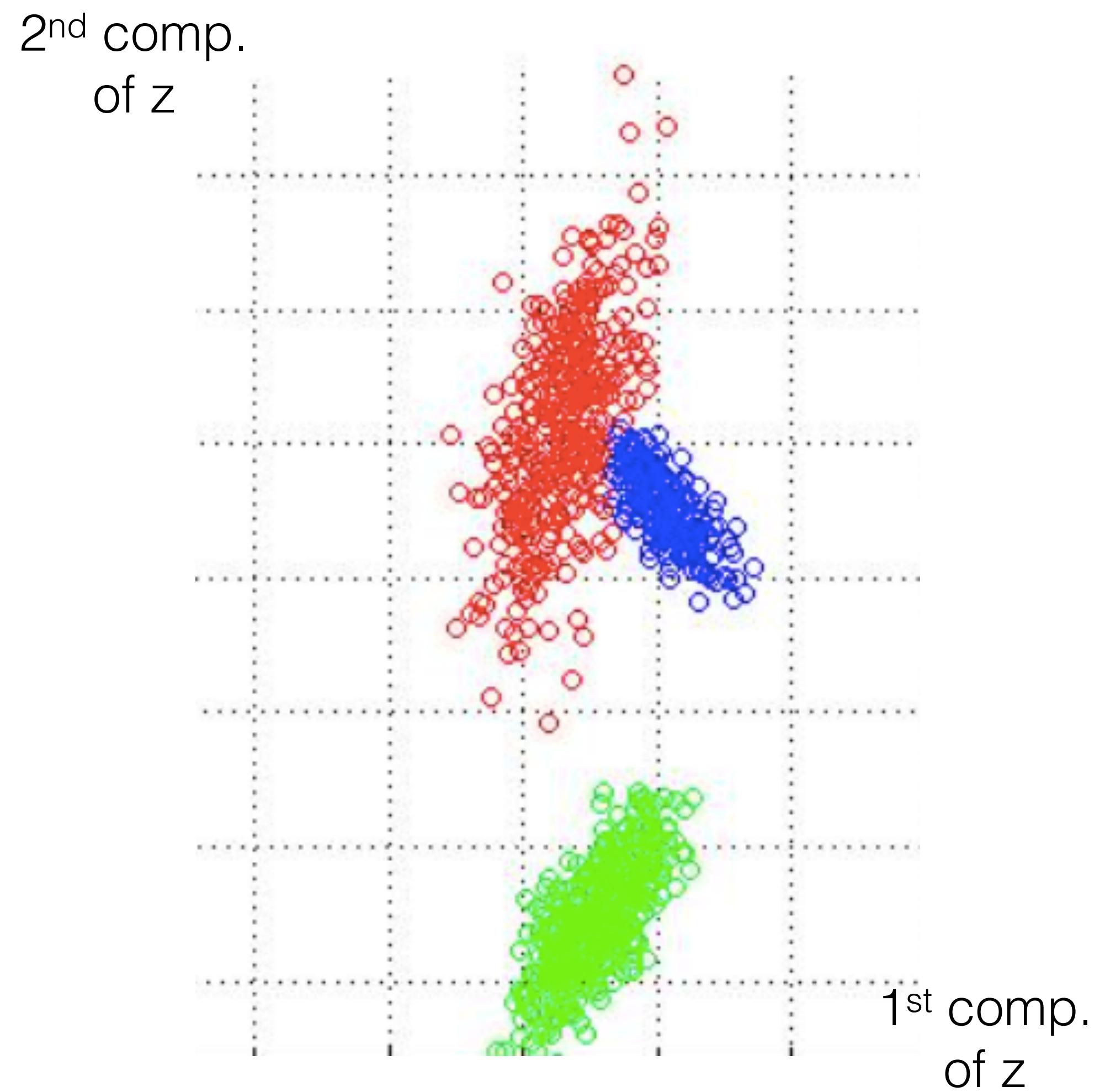
$$\sigma_k^2 = \frac{\sum_{i=1}^N q_{ik} (x_i - \mu_k)^2}{\sum_{i=1}^N q_{ik}}$$

# Gaussian mixture model: multi-dim case

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \{1, \dots, K\} \\ i = 1, \dots, N$$

$$\begin{aligned} p(x|\pi, \mu, \Sigma) &= \sum_{z=1}^K p(x|z, \mu, \Sigma)p(z|\pi) \\ &= \sum_{z=1}^K \mathcal{N}(x|\mu_z, \Sigma_z)\pi_z \end{aligned}$$



# Gaussian mixture model: EM-algorithm

Initialize  $\mu, \Sigma, \pi$  and repeat

$$\text{E-step: } q_{ik} = \frac{\mathcal{N}(x|\mu_k, \Sigma_k)\pi_k}{\sum_{k'=1}^K \mathcal{N}(x|\mu_{k'}, \Sigma_{k'})\pi_{k'}}$$

$$\text{M-step: } \mu_k = \frac{\sum_{i=1}^N q_{ik} x_i}{\sum_{i=1}^N q_{ik}} \quad \pi_k = \frac{1}{N} \sum_{i=1}^N q_{ik}$$

$$\sigma_k^2 = \frac{\sum_{i=1}^N q_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N q_{ik}}$$

# Gaussian mixture model: summary

- GMM provides “soft” clustering of points
- All formulas are easily interpretable
- There are more complicated models that involve priors on GMM parameters

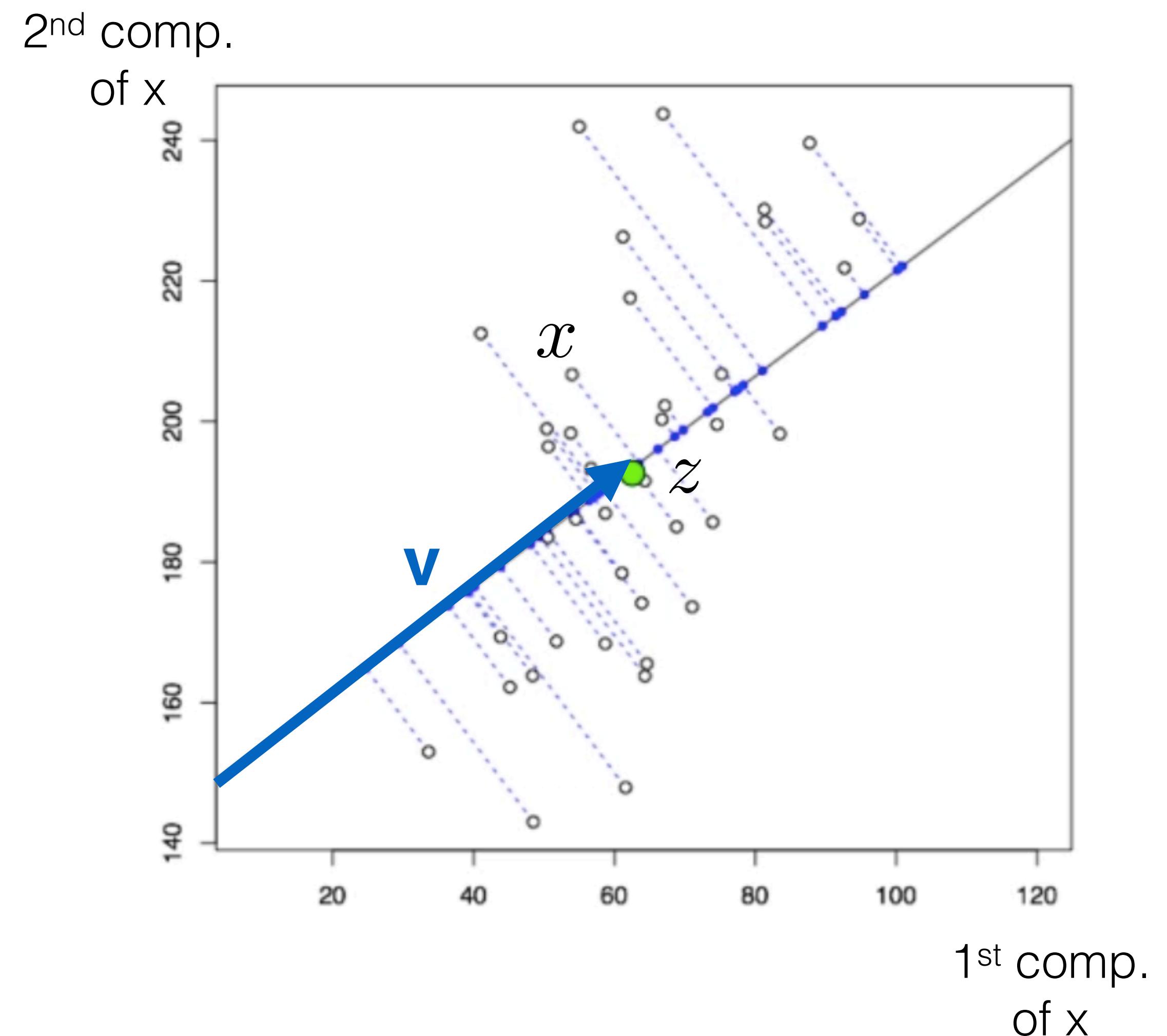
# Probabilistic PCA

$$X = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$$Z = \{z_1, \dots, z_N\}, z_i \in \mathbb{R}^{d_{\text{small}}} \\ i = 1, \dots, N$$

$$p(X, Z | V, \sigma^2) = \prod_{i=1}^n p(x_i | z_i, V, \sigma^2) p(z_i) =$$

$$= \prod_{i=1}^N \mathcal{N}(x_i | V z_i, \sigma^2 I) \mathcal{N}(z_i | 0, I)$$



# Probabilistic PCA

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

$$\text{E-step: } q(Z) = \prod_{i=1}^N q(z_i) = \prod_{i=1}^N p(z_i|x_i, V, \sigma^2)$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

# Probabilistic PCA

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

$$\text{E-step: } q(Z) = \prod_{i=1}^N q(z_i) = \prod_{i=1}^N p(z_i|x_i, V, \sigma^2)$$

$$q(z_i) \propto \mathcal{N}(x_i|Vz_i, \sigma^2 I) \mathcal{N}(z_i|0, I) =$$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

**conjugance!**

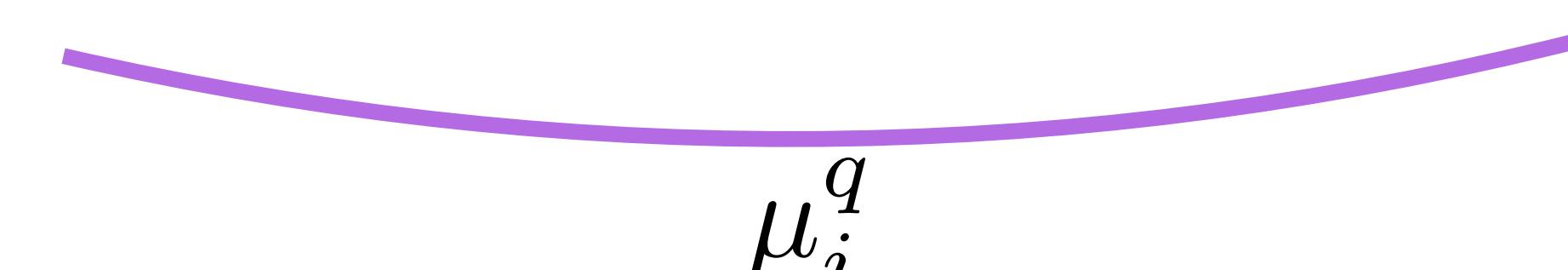
# Probabilistic PCA

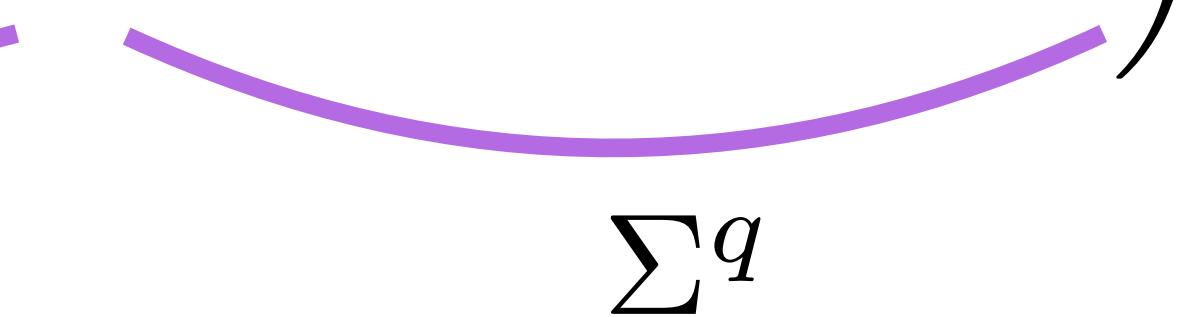
$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

$$\text{E-step: } q(Z) = \prod_{i=1}^N q(z_i) = \prod_{i=1}^N p(z_i|x_i, V, \sigma^2)$$

$$\begin{aligned} q(z_i) &\propto \mathcal{N}(x_i|Vz_i, \sigma^2 I)\mathcal{N}(z_i|0, I) = \\ &= \mathcal{N}\left(z_i \mid \left(\frac{1}{\sigma^2}V^T V + I\right)^{-1} \left(\sum_{i=1}^N V^T x_i\right), \left(\frac{1}{\sigma^2}V^T V + I\right)^{-1}\right) \end{aligned}$$

  
 $\mu_i^q$

  
 $\Sigma^q$

formulas of  $V$ ,  $\sigma$  and  $X$

Initialize  $\theta_*$  and repeat  
E-step:  $q(Z) = p(Z|X, \theta_*)$   
M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

# Probabilistic PCA

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

- Initialize  $\theta_*$  and repeat
- E-step:  $q(Z) = p(Z|X, \theta_*)$
- M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$\mathcal{L}(q(Z), \theta) = \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ$$

M-step:

$$\sum_{i=1}^N \int q(z_i) \log p(x_i, z_i | V, \sigma^2) dz_i = \mathbb{E}_{q(z_i)} \log p(x_i, z_i | V, \sigma^2) \rightarrow \max_{V, \sigma^2}$$

# Probabilistic PCA

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

- Initialize  $\theta_*$  and repeat
- E-step:  $q(Z) = p(Z|X, \theta_*)$
- M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$\mathcal{L}(q(Z), \theta) = \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ$$

M-step:

$$\sum_{i=1}^N \int q(z_i) \log p(x_i, z_i | V, \sigma^2) dz_i = \mathbb{E}_{q(z_i)} \log p(x_i, z_i | V, \sigma^2) \rightarrow \max_{V, \sigma^2}$$

1-dim z case:

$$\mathbb{E}_{q(z_i)} \log (\mathcal{N}(x_i | vz_i, \sigma^2) \mathcal{N}(z_i | 0, 1)) = \mathbb{E}_{q(z_i)} (-2 \log \sigma + \exp(0.5\sigma^{-2} ||x_i - vz_i||^2))$$

$$\mathbb{E}_{q(z_i)} z_i = \mu_i^q \quad \mathbb{E}_{q(z_i)} z_i^2 = \sigma_i^q + \mu_i^q$$

# Probabilistic PCA

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

$$p(z) = \mathcal{N}(z|0, I)$$

- Initialize  $\theta_*$  and repeat
- E-step:  $q(Z) = p(Z|X, \theta_*)$
- M-step:  $\theta_* = \operatorname{argmax}_\theta \mathcal{L}(q(Z), \theta)$

$$\mathcal{L}(q(Z), \theta) = \int q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} dZ$$

M-step:

$$V = \left( \sum_{i=1}^N x_i (\mathbb{E}_{q(z_i)} z_i)^T \right) \left( \sum_{i=1}^N \mathbb{E}_{q(z_i)} z_i z_i^T \right)^{-1}$$

$$\sigma^2 = \frac{1}{Nd} \sum_{i=1}^N \mathbb{E}_{q(z_i)} \|x_i - Vz_i\|^2$$

formulas  
of q params  
and X

# Probabilistic PCA: EM-algorithm

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I) \quad p(z) = \mathcal{N}(z|0, I)$$

E-step:  $q(z_i) = p(z_i|x_i, V, \sigma^2) = \mathcal{N}(x_i|\mu_i^q, \Sigma^q)$  **formula of V, σ and X**

M-step:  $V = \left( \sum_{i=1}^N x_i (\mathbb{E}_{q(z_i)} z_i)^T \right) \left( \sum_{i=1}^N \mathbb{E}_{q(z_i)} z_i z_i^T \right)^{-1}$  **formulas of q params and X**

$$\sigma^2 = \frac{1}{Nd} \sum_{i=1}^N \mathbb{E}_{q(z_i)} \|x_i - Vz_i\|^2$$

# Probabilistic PCA: EM-algorithm

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I) \quad p(z) = \mathcal{N}(z|0, I)$$

E-step:  $q(z_i) = p(z_i|x_i, V, \sigma^2) = \mathcal{N}(x_i|\mu_i^q, \Sigma^q)$

formula of  $V$ ,  $\sigma$  and  $X$

M-step:  $V = \left( \sum_{i=1}^N x_i (\mathbb{E}_{q(z_i)} z_i)^T \right) \left( \sum_{i=1}^N \mathbb{E}_{q(z_i)} z_i z_i^T \right)^{-1}$

formulas of  $q$  params  
and  $X$

$$\sigma^2 = \frac{1}{Nd} \sum_{i=1}^N \mathbb{E}_{q(z_i)} \|x_i - Vz_i\|^2$$

# Probabilistic PCA: pros

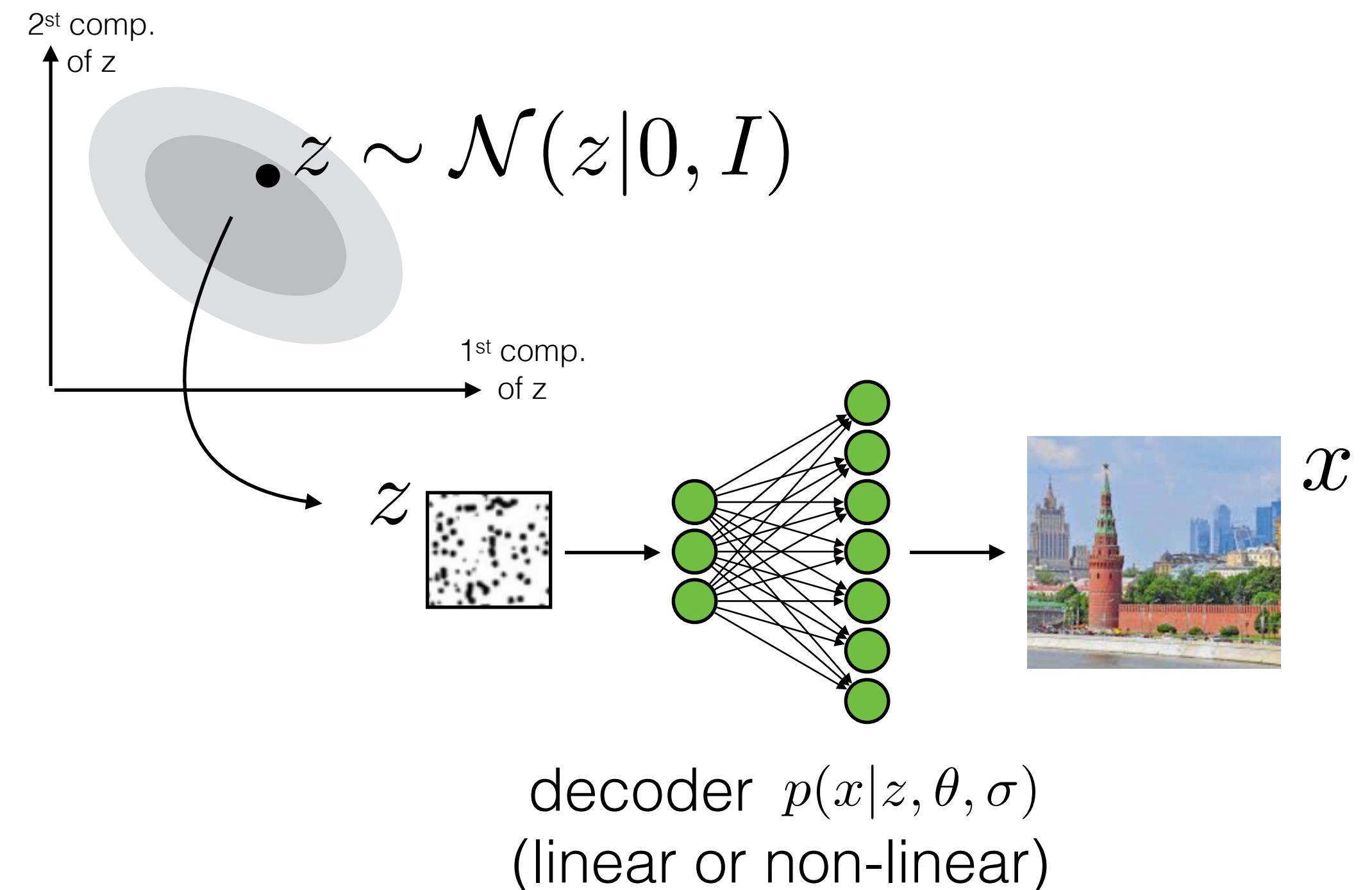
- Probabilistic PCA provides iterative EM-algorithm for finding transformation matrix  $V$  that may be more efficient than SVD-decomposition
- EM-algorithm allows combining models, e. g. Mixture of PPCA
- PPCA can determine  $d_{\text{small}}$  (more complex model)

# Next lecture: variational autoencoders

$$p(x|z, V, \sigma^2) = \mathcal{N}(x|Vz, \sigma^2 I)$$

next lecture:  
nonlinear transformation

$$p(z) = \mathcal{N}(z|0, I)$$



# Latent variable models

$$X = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^d \\ i = 1, \dots, N$$

$Z = \{z_1, \dots, z_N\}$  — latent variables

$$p(x, z|\theta) = p(x|z, \theta_1)p(z|\theta_2)$$

$$p(x|\theta) = \int p(x, z|\theta) dz$$

(or sum in case of discrete  $z$ )

How can l. v. models be useful?

- Understanding data structure:  $p(z|x, \theta) = \frac{p(x|z, \theta_1)p(z|\theta_2)}{\int(x|\tilde{z}, \theta_1)p(\tilde{z}|\theta_2)d\tilde{z}}$  e. g.  
what cluster  
does an object  
come from?
- Generating new objects:  $z \sim p(z|\theta_2), \quad x \sim p(x|z, \theta_1)$  e. g.  
generating  
new images  
using decoder