

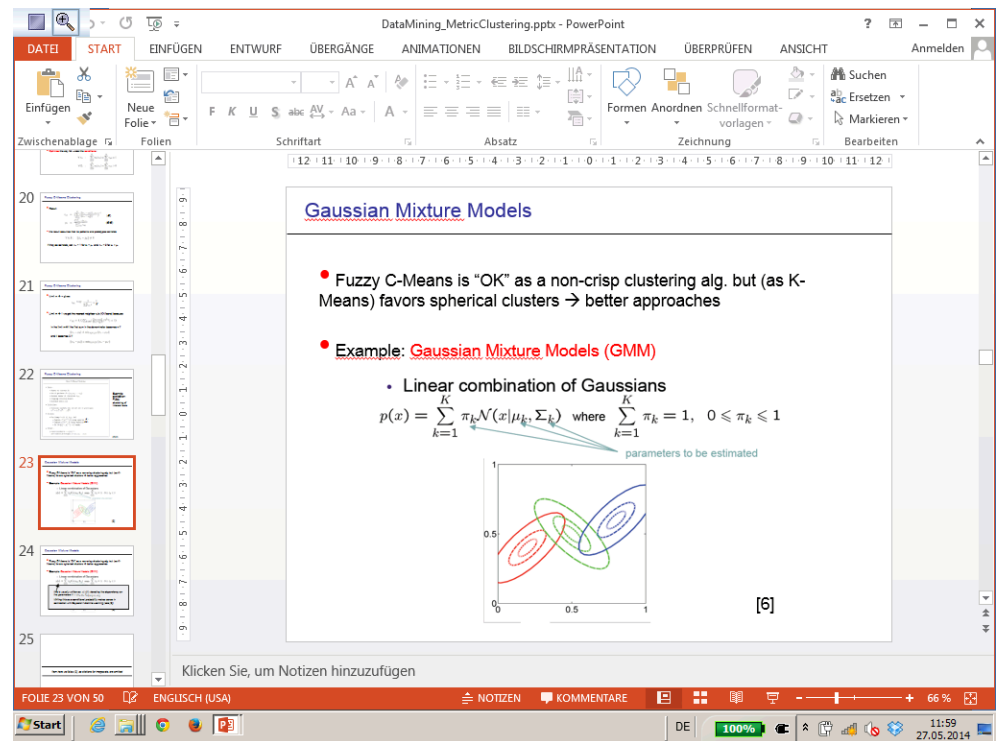
Script generated by TTT

Title: groh: profile1 (27.05.2014)

Date: Tue May 27 11:59:47 CEST 2014

Duration: 89:52 min

Pages: 91



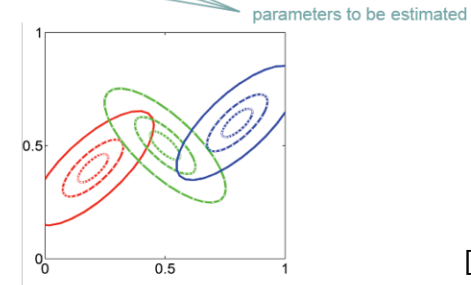
Gaussian Mixture Models

- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



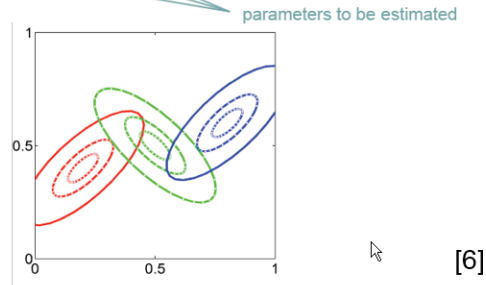
Gaussian Mixture Models

- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

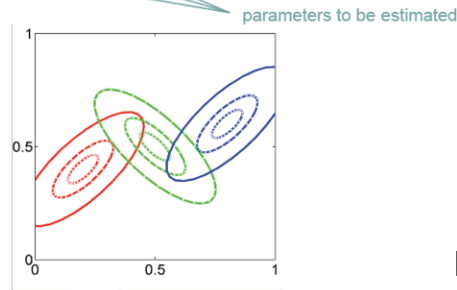


- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]

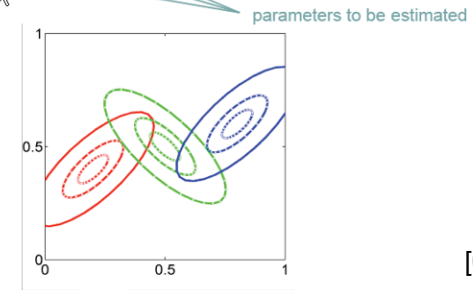


- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]

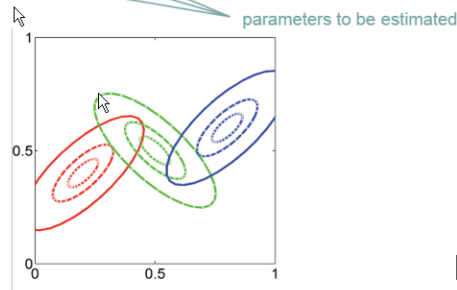


- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]

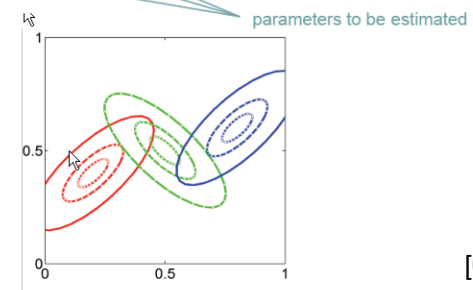


- Fuzzy C-Means is "OK" as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]

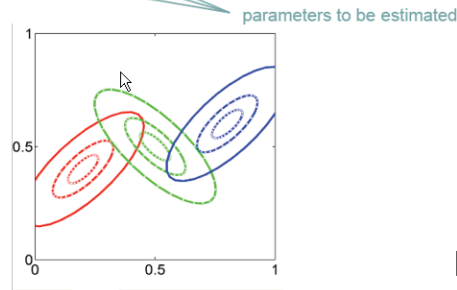


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]

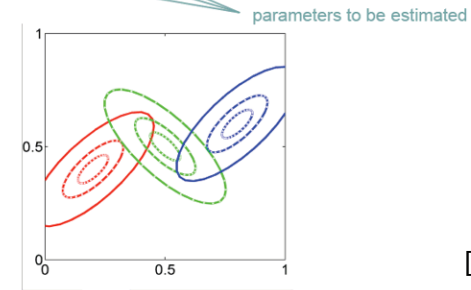


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$



[6]



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])

[6]



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])

[6]



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches
- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])

0 0.5 1 $[\theta]$



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches
- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])

0 0.5 1 $[\theta]$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid**: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_i|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid**: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_i|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

• iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

• $p(X|\theta)$ is called **likelihood**

• „finding the θ that best explains the data“:

Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

• convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$

$\Rightarrow \log p(X|\theta) = \sum_i \log p(x_i|\theta)$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian)

$$p(\mathbf{x}|\theta) = \mathcal{N}(\mathbf{x}|\underbrace{\mu, \Sigma}_{\theta}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

• **log likelihood:** (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1}(\mathbf{x}_n - \mu)$$

• **Maximum log likelihood:**

$$\theta_{ML} = \operatorname{argmax}_{\theta} \log p(X|\theta) \Rightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \mu_{ML} : \frac{\partial}{\partial \mu} \ln p(\mathbf{X}|\mu, \Sigma) = 0 \\ \Sigma_{ML} : \frac{\partial}{\partial \Sigma} \ln p(\mathbf{X}|\mu, \Sigma) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \mu_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \Sigma_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T \end{aligned} \right.$$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

• iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

• $p(X|\theta)$ is called **likelihood**

• „finding the θ that best explains the data“:

Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

• convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$

$\Rightarrow \log p(X|\theta) = \sum_i \log p(x_i|\theta)$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian)

$$p(\mathbf{x}|\theta) = \mathcal{N}(\mathbf{x}|\underbrace{\mu, \Sigma}_{\theta}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

• **log likelihood:** (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1}(\mathbf{x}_n - \mu)$$

• **Maximum log likelihood:**

$$\theta_{ML} = \operatorname{argmax}_{\theta} \log p(X|\theta) \Rightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \mu_{ML} : \frac{\partial}{\partial \mu} \ln p(\mathbf{X}|\mu, \Sigma) = 0 \\ \Sigma_{ML} : \frac{\partial}{\partial \Sigma} \ln p(\mathbf{X}|\mu, \Sigma) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \mu_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \Sigma_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T \end{aligned} \right.$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right.$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right.$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right.$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right.$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right\}$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right\}$$



Example: $x \in \mathbb{R}^n$ and $p(x|\theta)$ is one multivariate Gaussian

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

• log likelihood: (use base e)

$$\ln p(\mathbf{X}|\theta) = \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

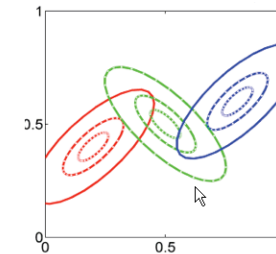
• Maximum log likelihood:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{X}|\theta) \Leftrightarrow \nabla_{\theta} (\sum_i \log p(x_i|\theta)) \stackrel{!}{=} 0$$

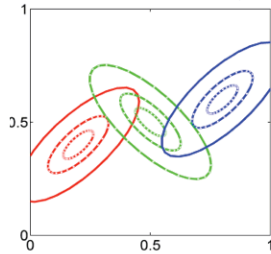
$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned} \right\} \Rightarrow \left. \begin{aligned} \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{aligned} \right\}$$



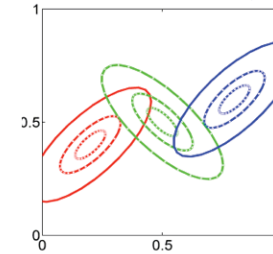
$$\text{GMM: } p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$



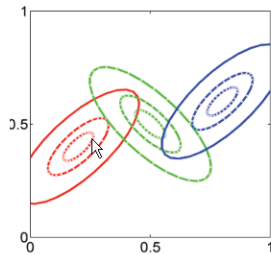
GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

• conditional probability

$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ $p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$

$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
 $p(\mathbf{x}, \mathbf{z})$

GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

• conditional probability

$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$

$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

• conditional probability

$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$

$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

• conditional probability

$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$

$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

• conditional probability

$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$

$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$

- 1 of K representation

K -dimensional binary random variable \mathbf{z}

$z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$p(z_k = 1) = \pi_k$

$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

remark: If we have several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, then, because we have represented the marginal distribution in the form $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$, it follows that for every observed data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n .

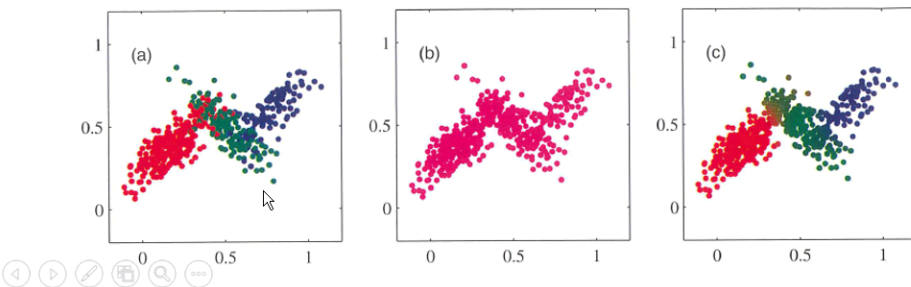
$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



- Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

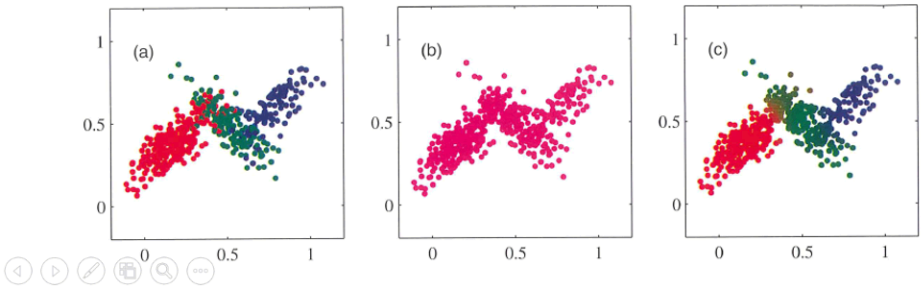
- Example



- Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

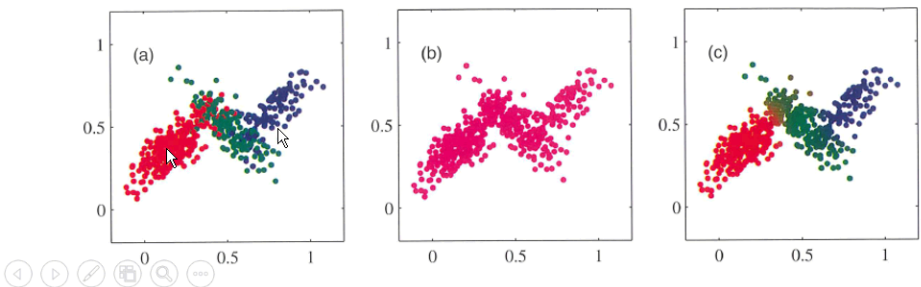
- Example



- Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

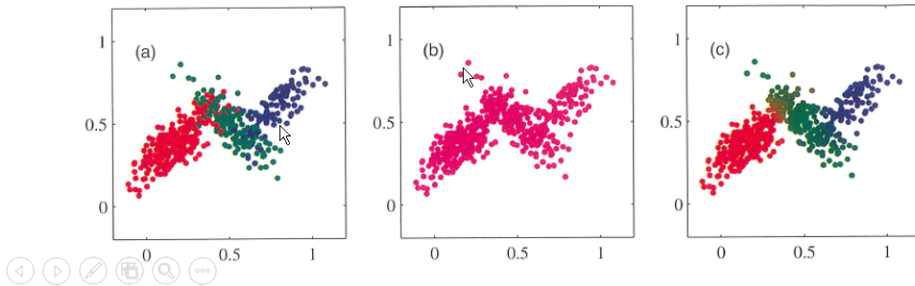
- Example



Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \end{aligned}$$

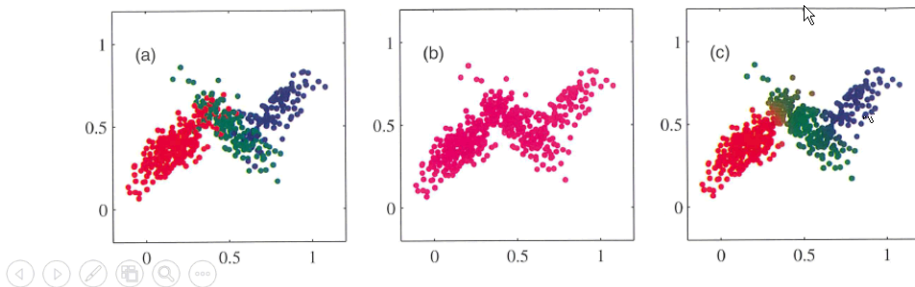
Example



Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \end{aligned}$$

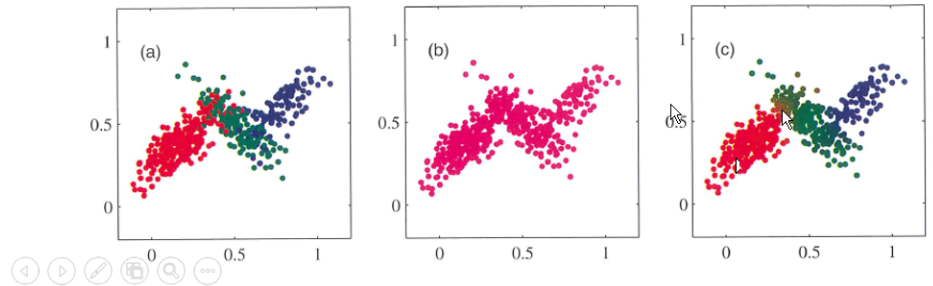
Example



Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \end{aligned}$$

Example



GMM: $p(\mathbf{x}|\boldsymbol{\theta}) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. $0 \leq \pi_k \leq 1$ $\sum_{k=1}^K \pi_k = 1$

1 of k representation K -dimensional binary random variable \mathbf{z}
 $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$$p(z_k = 1) = \pi_k$$

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

remark: If we have several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, then, because we have represented the marginal distribution in the form $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, it follows that for every observed data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n .

$$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

GMM:
$$p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$

- 1 of K representation $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$
 $p(z_k = 1) = \pi_k$
 $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

remark: If we have several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, then, because we have represented the marginal distribution in the form $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$, it follows that for every observed data point \mathbf{x}_n there is a corresponding latent variable z_n .

$$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Vector of K D -dim. means $\boldsymbol{\mu}_k$
 Vector of K $D \times D$ covariances $\boldsymbol{\Sigma}_k$

- maximizing w.r.t $\boldsymbol{\pi}, \boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Vector of K D -dim. means $\boldsymbol{\mu}_k$
 Vector of K $D \times D$ covariances $\boldsymbol{\Sigma}_k$

- maximizing w.r.t $\boldsymbol{\pi}, \boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

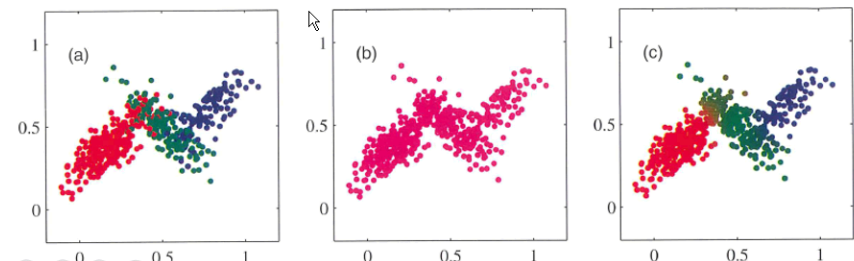


- Responsibilities

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- Example



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Vector of K D -dim. means $\boldsymbol{\mu}_k$
 Vector of K $D \times D$ covariances $\boldsymbol{\Sigma}_k$

- maximizing w.r.t $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Vector of K D -dim. means $\boldsymbol{\mu}_k$
 Vector of K $D \times D$ covariances $\boldsymbol{\Sigma}_k$

- maximizing w.r.t $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

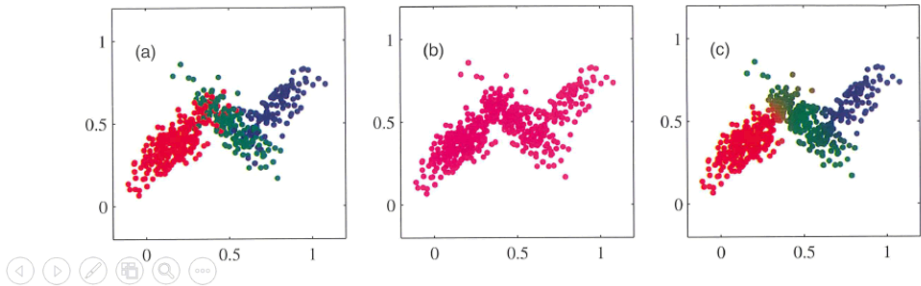


- Responsibilities

$$\gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- Example



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Vector of K D -dim. means $\boldsymbol{\mu}_k$
 Vector of K $D \times D$ covariances $\boldsymbol{\Sigma}_k$

- maximizing w.r.t $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

- so what?! → **Problem:** Expr. depend on $\gamma(z_{nk})$ which depends on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ which depends on $\gamma(z_{nk})$ which depends on

- Idea: Alternating approach (**EM-algorithm**):

Step t: Evaluate $\gamma(z_{nk})_{(t)}$ using $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t-1)}$

Evaluate $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t)}$ using $\gamma(z_{nk})_{(t-1)}$



The screenshot shows a PowerPoint presentation window titled 'DataMining_MetricClustering.pptx - PowerPoint'. The slide content is identical to the first slide, including the maximum likelihood equations and the EM algorithm steps. The slide number 33 is highlighted in the left sidebar. The status bar at the bottom indicates 'FOLIE 33 VON 50' and 'ENGLISCH (USA)'.

Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

- so what?! → **Problem:** Expr. depend on $\gamma(z_{nk})$ which depends on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ which depends on $\gamma(z_{nk})$ which depends on

- Idea: Alternating approach (**EM-algorithm**):

Step t: Evaluate $\gamma(z_{nk})_{(t)}$ using $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t-1)}$

Evaluate $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t)}$ using $\gamma(z_{nk})_{(t-1)}$



Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

- so what?! → **Problem:** Expr. depend on $\gamma(z_{nk})$ which depends on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ which depends on $\gamma(z_{nk})$ which depends on

- Idea: Alternating approach (**EM-algorithm**):

Step t: Evaluate $\gamma(z_{nk})_{(t)}$ using $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t-1)}$

Evaluate $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t)}$ using $\gamma(z_{nk})_{(t-1)}$



GMM-Basics

Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

• so what?! → **Problem:** Expr. depend on $\gamma(z_{nk})$ which depends on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ which depends on $\gamma(z_{nk})$ which depends on

• Idea: Alternating approach (**EM-algorithm**):

Step t: Evaluate $\gamma(z_{nk})_{(t)}$ using $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t-1)}$

Evaluate $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t)}$ using $\gamma(z_{nk})_{(t)}$

GMM-Basics

Maximum likelihood (GMM)

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) \quad \pi_k = \frac{N_k}{N}$$

• so what?! → **Problem:** Expr. depend on $\gamma(z_{nk})$ which depends on $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ which depends on $\gamma(z_{nk})$ which depends on

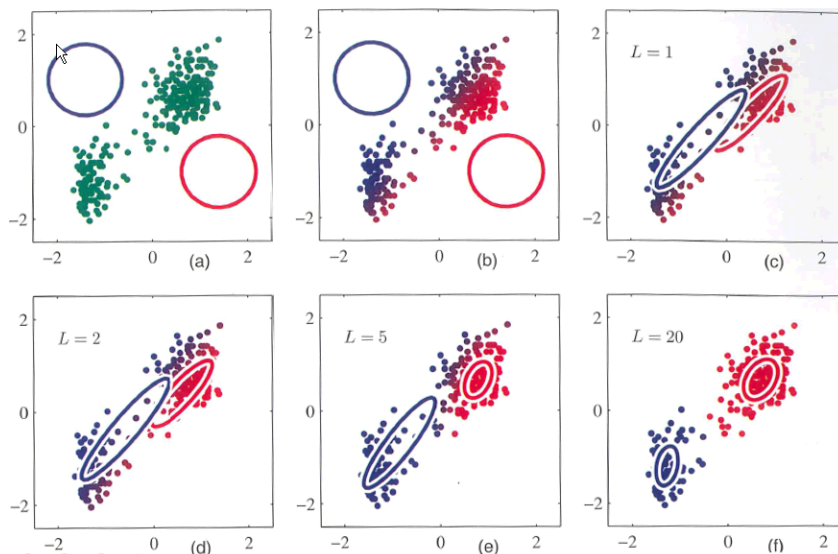
• Idea: Alternating approach (**EM-algorithm**):

Step t: Evaluate $\gamma(z_{nk})_{(t)}$ using $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t-1)}$

Evaluate $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})_{(t)}$ using $\gamma(z_{nk})_{(t)}$

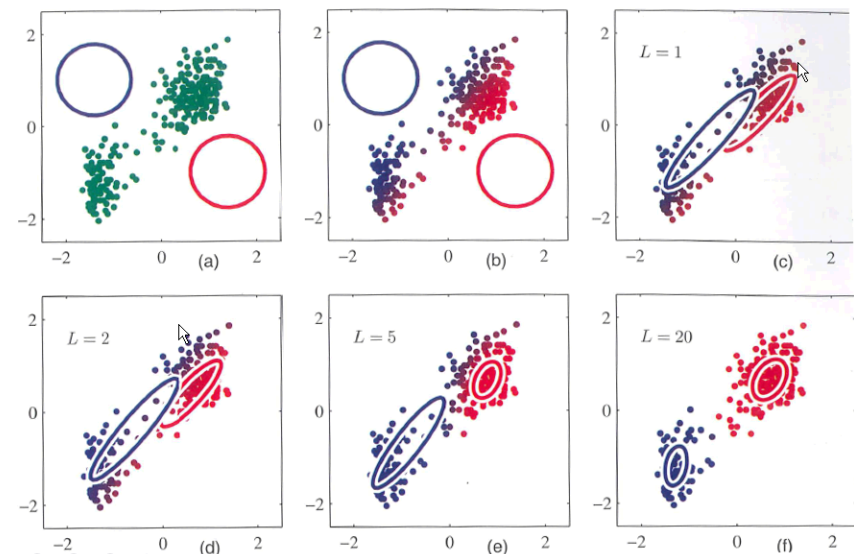
GMM-Basics

Maximum likelihood (GMM)



GMM-Basics

Maximum likelihood (GMM)



Maximum likelihood (GMM)

EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}. \quad (9.23)$$



- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !

- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))

- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



EM-algorithm: General View

- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !

- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))

- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



EM-algorithm: General View

- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !

- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))

- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !
- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))
- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !
- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))
- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



GMM-Basics

Maximum likelihood (GMM)

EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)} \quad (9.23)$$

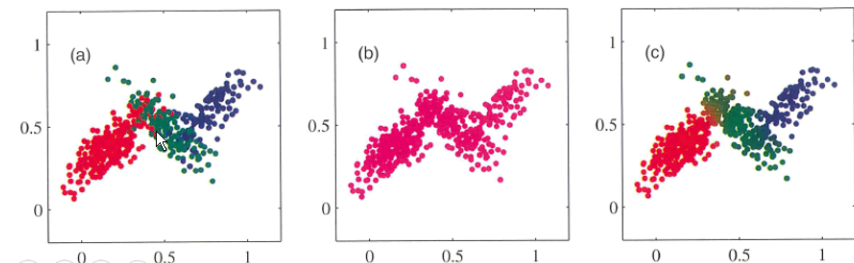


GMM-Basics

- Responsibilities

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \mu_j, \Sigma_j)}. \end{aligned}$$

- Example



EM-algorithm: General View

- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !

- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))

- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



EM-algorithm: General View

- alternating EM:

E-Step:
compute $Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta).$

M-Step:
compute $\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}).$



EM-algorithm: General View

- Having latent variables \mathbf{Z} , ML becomes

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Summation inside $\ln \rightarrow$ Problems !

- If we knew the complete dataset $\{\mathbf{X}, \mathbf{Z}\}$ (and thus the distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$), we could use ML to solve for θ with $p(\mathbf{X}, \mathbf{Z}|\theta)$ directly (which is easy, as we will see, because $p(\mathbf{X}, \mathbf{Z}|\theta)$ is of exponential family (the functional form is known!!))

- We only know $p(\mathbf{Z}|\mathbf{X}, \theta)$ (\rightarrow responsibilities, as we will see) \rightarrow compute expectation of (unknown) quantity $p(\mathbf{X}, \mathbf{Z}|\theta)$ or even better of the quantity $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



EM: Relation to K-Means

- If we use k Gaussians with $\Sigma = \epsilon \mathbf{I}$:

$$p(\mathbf{x}|\mu_k, \Sigma_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \mu_k\|^2 \right\}$$

that is why K-Means favors spherical clusters

- we get for the responsibilities:

$$\gamma(z_{nk}) = \frac{\pi_k \exp \{ -\|\mathbf{x}_n - \mu_k\|^2 / 2\epsilon \}}{\sum_j \pi_j \exp \{ -\|\mathbf{x}_n - \mu_j\|^2 / 2\epsilon \}}$$

- Letting $\epsilon \rightarrow 0$ and Taylor-Expansion:

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\mu, \Sigma, \pi)] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|^2 + \text{const}$$

\rightarrow same as on slide 18



EM: Relation to K-Means

that is why K-Means favors spherical clusters

- If we use k Gaussians with $\Sigma = \epsilon I$:

$$p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp\left\{-\frac{1}{2\epsilon}\|\mathbf{x} - \boldsymbol{\mu}_k\|^2\right\}$$

- we get for the responsibilities:

$$\gamma(z_{nk}) = \frac{\pi_k \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2/2\epsilon\}}{\sum_j \pi_j \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2/2\epsilon\}}$$

- Letting $\epsilon \rightarrow 0$ and Taylor-Expansion:

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + \text{const}$$

→ same as on slide 18



EM: Relation to K-Means

that is why K-Means favors spherical clusters

- If we use k Gaussians with $\Sigma = \epsilon I$:

$$p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp\left\{-\frac{1}{2\epsilon}\|\mathbf{x} - \boldsymbol{\mu}_k\|^2\right\}$$

- we get for the responsibilities:

$$\gamma(z_{nk}) = \frac{\pi_k \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2/2\epsilon\}}{\sum_j \pi_j \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2/2\epsilon\}}$$

- Letting $\epsilon \rightarrow 0$ and Taylor-Expansion:

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + \text{const}$$

→ same as on slide 18



EM: Relation to K-Means

that is why K-Means favors spherical clusters

- If we use k Gaussians with $\Sigma = \epsilon I$:

$$p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp\left\{-\frac{1}{2\epsilon}\|\mathbf{x} - \boldsymbol{\mu}_k\|^2\right\}$$

- we get for the responsibilities:

$$\gamma(z_{nk}) = \frac{\pi_k \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2/2\epsilon\}}{\sum_j \pi_j \exp\{-\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2/2\epsilon\}}$$

- Letting $\epsilon \rightarrow 0$ and Taylor-Expansion:

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + \text{const}$$

→ same as on slide 18



Real World Networks: Properties and Models



Lecture will mostly follow [1], thus corresponding citations are often omitted to increase readability



Information Networks / Knowledge NW

- **Most studied** examples: citation NW (tree), the WWW;
- **Example findings:**
 - $p(k)$ of author having k papers: $p(k) \sim k^{-\alpha}$: power law
 - distribution of in or out degrees of WWW pages (also for citation NW): $p(k) \sim k^{-\alpha}$

- Other **examples:**
 - **bipartite preference networks** :
→ recommender systems == link prediction on these NW;
example: collaborative filtering
 - **ontologies, semantic NW**
 - **word networks**
 - **tripartite tag/author/item networks**
→ Folksonomies

Information Networks / Knowledge NW

- **Most studied** examples: citation NW (tree), the WWW;
- **Example findings:**
 - $p(k)$ of author having k papers: $p(k) \sim k^{-\alpha}$: power law
 - distribution of in or out degrees of WWW pages (also for citation NW): $p(k) \sim k^{-\alpha}$

- Other **examples:**
 - **bipartite preference networks** :
→ recommender systems == link prediction on these NW;
example: collaborative filtering
 - **ontologies, semantic NW**
 - **word networks**
 - **tripartite tag/author/item networks**
→ Folksonomies

Information Networks / Knowledge NW

- **Most studied** examples: citation NW (tree), the WWW;
- **Example findings:**
 - $p(k)$ of author having k papers: $p(k) \sim k^{-\alpha}$: power law
 - distribution of in or out degrees of WWW pages (also for citation NW): $p(k) \sim k^{-\alpha}$

- Other **examples:**
 - **bipartite preference networks** :
→ recommender systems == link prediction on these NW;
example: collaborative filtering
 - **ontologies, semantic NW**
 - **word networks**
 - **tripartite tag/author/item networks**
→ Folksonomies

Technological Networks

- **Most studied examples:** **distribution NW:**
 - the **Internet**,
 - **electric power grids**,
 - **traffic NW** (roads, railway tracks etc.)

Biological / Chemistry / Physics Networks

- **Most studied examples:**
 - **biochemical pathways, gene-protein** and **protein-protein interaction NW**
 - **nervous systems, vascular systems** (also natural distribution NW),
 - **food NW, ecological dependency NW**

Technological Networks

- Most studied examples: distribution NW:
 - the Internet,
 - electric power grids,
 - traffic NW (roads, railway tracks etc.)

Biological / Chemistry / Physics Networks

- Most studied examples:
 - biochemical pathways, gene-protein and protein-protein interaction NW
 - nervous systems, vascular systems (also natural distribution NW),
 - food NW, ecological dependency NW



Mean Average Path Length

- “Small World Effect”: $l(n)$ “small” $\rightarrow l(n) \in O(\log(n))$
- undirected graph:

$$\ell = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}$$

formula also counts 0 distances from i to i : $\frac{1}{2}n(n+1) = \frac{1}{2}n(n-1) + n$

- Expression allowing for disconnected components (where $d_{ij} = \infty$ can occur): harmonic mean:

$$\ell^{-1} = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}^{-1}$$

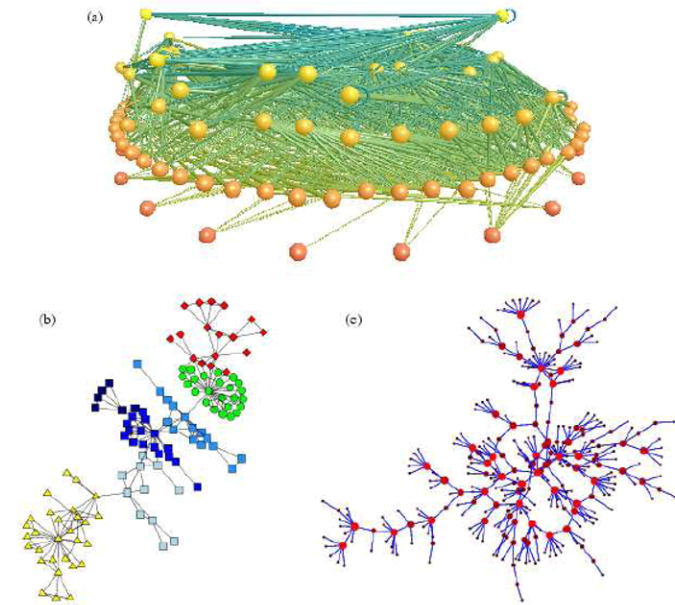


FIG. 2 Three examples of the kinds of networks that are the topic of this review. (a) A food web of predator-prey interactions between species in a freshwater lake [272]. Picture courtesy of Neo Martinez and Richard Williams. (b) The network of collaborations between scientists at a private research institution [171]. (c) A network of sexual contacts between individuals in the study by Potterat *et al.* [342].

[1]

Mean Average Path Length

- “Small World Effect”: $l(n)$ “small” $\rightarrow l(n) \in O(\log(n))$
- undirected graph:

$$\ell = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}$$

formula also counts 0 distances from i to i : $\frac{1}{2}n(n+1) = \frac{1}{2}n(n-1) + n$

- Expression allowing for disconnected components (where $d_{ij} = \infty$ can occur): harmonic mean:

$$\ell^{-1} = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}^{-1}$$



Mean Average Path Length

- “Small World Effect”: $l(n)$ “small” $\rightarrow l(n) \in O(\log(n))$
- undirected graph:

$$\ell = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}$$

formula also counts 0 distances from i to i : $\frac{1}{2}n(n+1) = \frac{1}{2}n(n-1) + n$

- Expression allowing for disconnected components (where $d_{ij} = \infty$ can occur): harmonic mean:

$$\ell^{-1} = \frac{1}{\frac{1}{2}n(n+1)} \sum_{i \geq j} d_{ij}^{-1}$$



Transitivity / Clustering Coefficient

- Clustering coefficient (whole graph):

$$C = C^{(1)} = \frac{3 \times \text{number of triangles in the network}}{\text{number of connected triples of vertices}} \quad p(\text{FOAF})$$

$$= \frac{6 \times \text{number of triangles in the network}}{\text{number of paths of length two}}$$

- Clustering coefficient (Watts-Strogatz-version, for node i):

$$C_i = \frac{\text{number of triangles connected to vertex } i}{\text{number of triples centered on vertex } i}$$

$$= \frac{|\{e_{\{kj\}} \mid v_k, v_j \in N_i\}|}{\frac{k_i(k_i-1)}{2}} \quad (\text{see Introduction, } k_i = \text{degree of node } i)$$

Clustering coefficient (Watts-Strogatz-version, for whole graph):

$$C = C^{(2)} = \frac{1}{n} \sum_i C_i$$



mean of ratio instead of ratio of means

Transitivity / Clustering Coefficient

- Clustering coefficient (whole graph):

$$C = C^{(1)} = \frac{3 \times \text{number of triangles in the network}}{\text{number of connected triples of vertices}} \quad p(\text{FOAF})$$

$$= \frac{6 \times \text{number of triangles in the network}}{\text{number of paths of length two}}$$

- Clustering coefficient (Watts-Strogatz-version, for node i):

$$C_i = \frac{\text{number of triangles connected to vertex } i}{\text{number of triples centered on vertex } i}$$

$$= \frac{|\{e_{\{kj\}} \mid v_k, v_j \in N_i\}|}{\frac{k_i(k_i-1)}{2}} \quad (\text{see Introduction, } k_i = \text{degree of node } i)$$

Clustering coefficient (Watts-Strogatz-version, for whole graph):

$$C = C^{(2)} = \frac{1}{n} \sum_i C_i$$



mean of ratio instead of ratio of means

Transitivity / Clustering Coefficient

- Clustering coefficient (whole graph):

$$C = C^{(1)} = \frac{3 \times \text{number of triangles in the network}}{\text{number of connected triples of vertices}} \quad p(\text{FOAF})$$

$$= \frac{6 \times \text{number of triangles in the network}}{\text{number of paths of length two}}$$

- Clustering coefficient (Watts-Strogatz-version, for node i):

$$C_i = \frac{\text{number of triangles connected to vertex } i}{\text{number of triples centered on vertex } i}$$

$$= \frac{|\{e_{\{kj\}} \mid v_k, v_j \in N_i\}|}{\frac{k_i(k_i-1)}{2}} \quad (\text{see Introduction, } k_i = \text{degree of node } i)$$

Clustering coefficient (Watts-Strogatz-version, for whole graph):

$$C = C^{(2)} = \frac{1}{n} \sum_i C_i$$



mean of ratio instead of ratio of means