

Script generated by TTT

Title: groh: profile1 (02.06.2015)

Date: Tue Jun 02 15:20:48 CEST 2015

Duration: 69:52 min

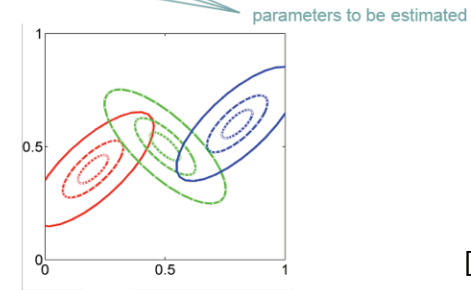
Pages: 61

- 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]

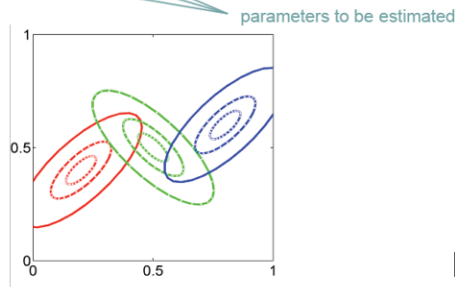


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- 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]

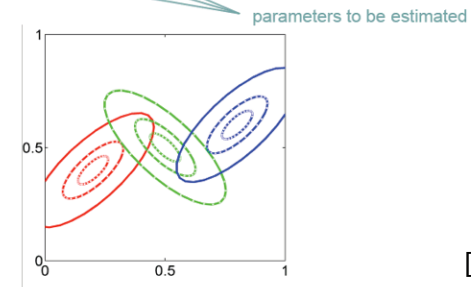


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- 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]



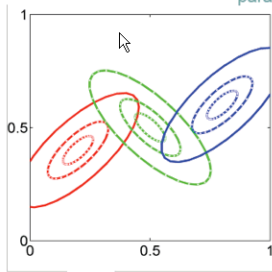
- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

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

- Linear combination of Gaussians

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

parameters to be estimated



[6]



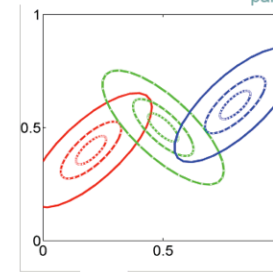
- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

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

- Linear combination of Gaussians

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

parameters to be estimated



[6]



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

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

- 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: $x \in \mathbb{R}^m$ 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: $x \in \mathbb{R}^m$ 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: $x \in \mathbb{R}^m$ 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: $x \in \mathbb{R}^m$ 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: $x \in \mathbb{R}^m$ 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:

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

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



Example: $x \in \mathbb{R}^m$ 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:

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

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



Example: $x \in \mathbb{R}^m$ 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)

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

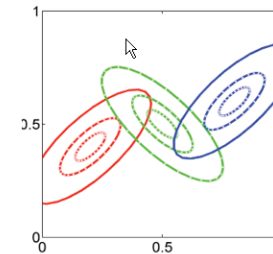
• Maximum log likelihood:

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

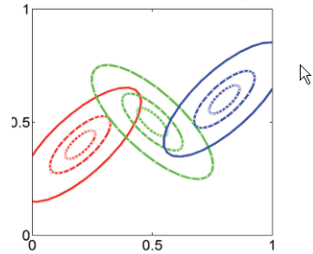
$$\left. \begin{aligned} \boldsymbol{\mu}_{ML} : \frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \\ \boldsymbol{\Sigma}_{ML} : \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \end{aligned}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



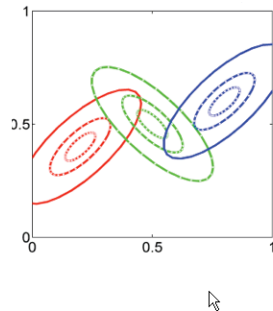
$$\text{GMM: } p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$



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

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

• 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

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

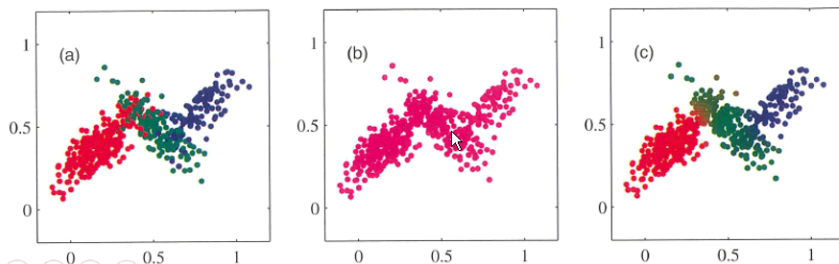
$p(\mathbf{x}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{p(\mathbf{x}, \mathbf{z})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$



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

• 1 of K representation

K -dimensional binary random variable \mathbf{z}

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

$p(z_k = 1) = \pi_k$

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

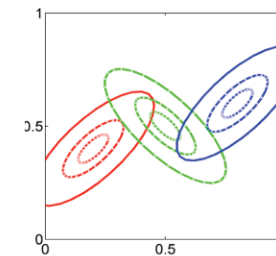
• conditional probability

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

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



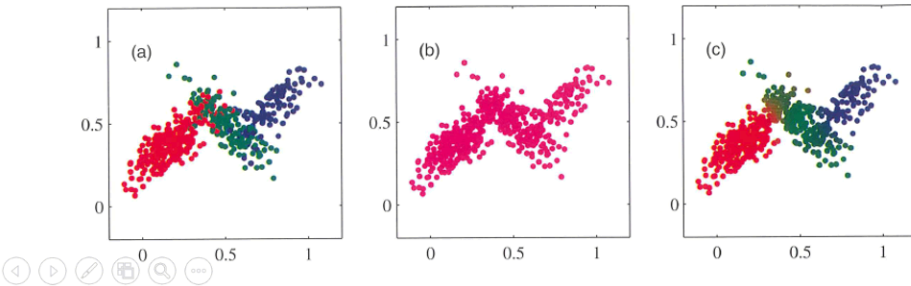
GMM: $p(\mathbf{x}|\theta) = p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$



Responsibilities

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

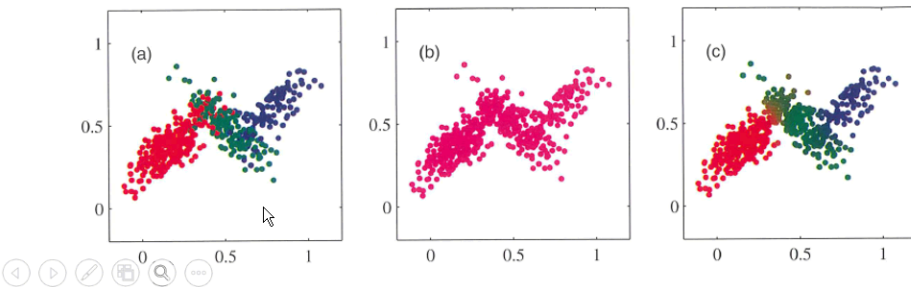
Example



Responsibilities

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

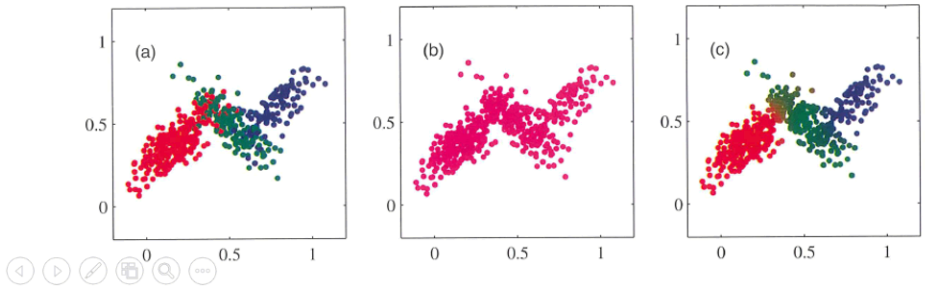
Example



Responsibilities

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

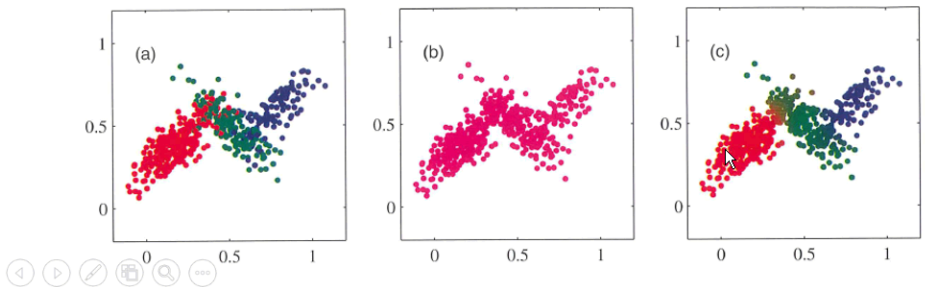
Example



Responsibilities

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

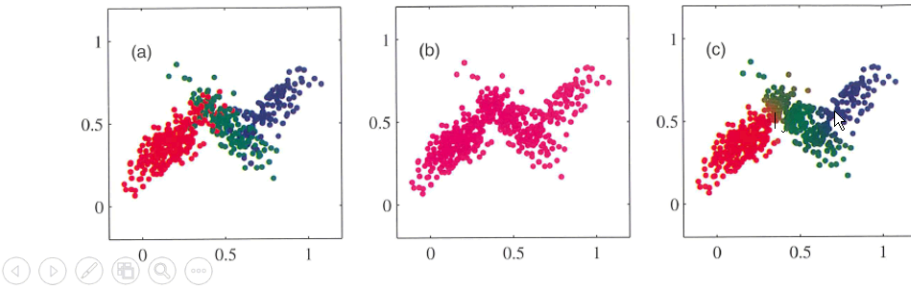
Example



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



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 $\boldsymbol{\Sigma} \rightarrow$

$$\begin{aligned} \boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n & \boldsymbol{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \\ \left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) & & \pi_k &= \frac{N_k}{N} \end{aligned}$$

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 $\boldsymbol{\Sigma} \rightarrow$

$$\begin{aligned} \boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n & \boldsymbol{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \\ \left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) & & \pi_k &= \frac{N_k}{N} \end{aligned}$$

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 $\boldsymbol{\Sigma} \rightarrow$

$$\begin{aligned} \boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n & \boldsymbol{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \\ \left(N_k = \sum_{n=1}^N \gamma(z_{nk}) \right) & & \pi_k &= \frac{N_k}{N} \end{aligned}$$

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 $\boldsymbol{\Sigma} \rightarrow$

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

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



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 $\boldsymbol{\Sigma} \rightarrow$

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

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

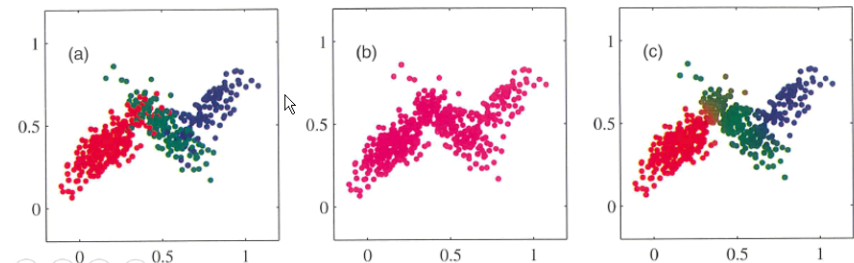


- Responsibilities

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

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

- Example



Maximum likelihood (GMM)

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

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

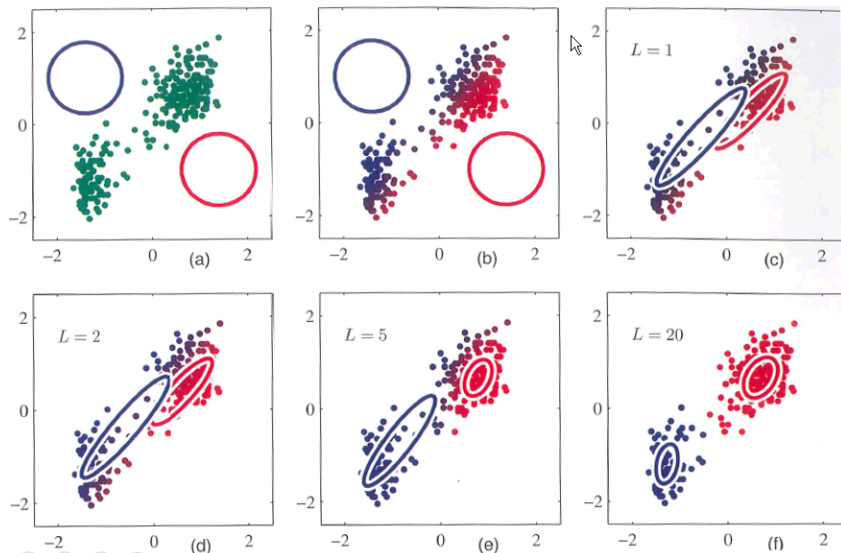
- maximizing w.r.t $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \rightarrow$

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

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



Maximum likelihood (GMM)



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?! \rightarrow **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

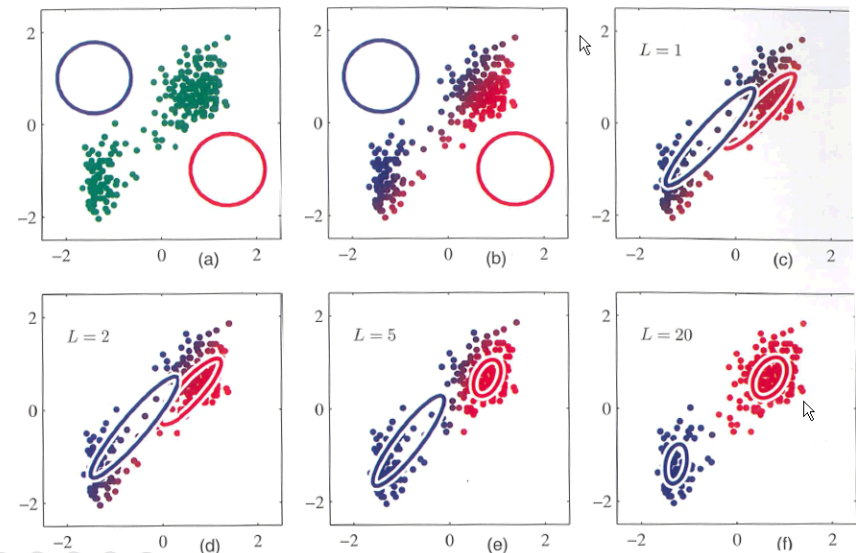
- Idea: Alternating approach (**EM-algorithm**):

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 $\gamma(z_{nk})_{(t)}$



Maximum likelihood (GMM)



- 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 $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$



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



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



- 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



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

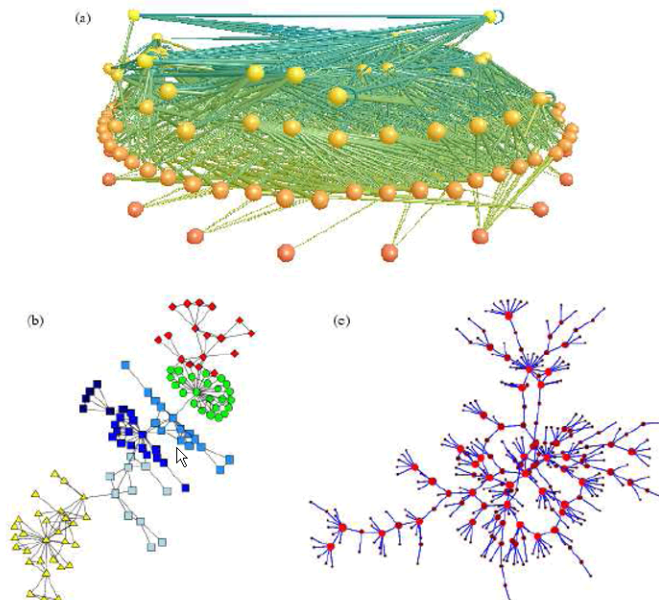
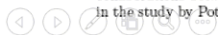


FIG. 2 Three examples of the kinds of networks that are the topic of this review. (a) A food web of predator-prey interactions between species in a freshwater lake [272]. Picture courtesy of Neo Martinez and Richard Williams. (b) The network of collaborations between scientists at a private research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



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



Studying Complex Networks

- Paradigm shift: small NW → large NW:
 - interest in **individual** elements (“centrality of node x ”) → interest in **global** statistical / topological properties (“degree distribution of NW”)
 - investigating **particular** NW instance → general **model** for types of NW with certain properties
 - 100 nodes → 10^8 nodes
 - visualization **possible** → **impossible** / pointless
- Typical sorts of NW investigated:
 - social, information, technological, biological**



Studying Complex Networks

- Paradigm shift: small NW → large NW:
 - interest in **individual** elements (“centrality of node x”) → interest in **global** statistical / topological properties (“degree distribution of NW”)
 - investigating **particular** NW instance → general **model** for types of NW with certain properties
 - 100 nodes → 10^8 nodes
 - visualization **possible** → **impossible** / pointless

- Typical sorts of NW investigated:
social, information, technological, biological



Social Networks

- Formalizations of **social context** (mostly long and medium term)

- **Until dawn of databases** : Collection via social science methods (questionnaires etc.) → „**laboratory effects**“ ;
- **Until dawn of Social Web**: „Indirect characterization of social relations“ (e.g. co-citation networks, collaboration NW (e.g. movie co-acting NW))
- **Today (Social Web / Mobile Social Web)** : self declared explication of social structures ; „collect data about“ / „observe“ Homo Sapiens in its „**natural habitat**“ (→ Twitter, Facebook etc.)



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



- 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



- 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



- 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



- 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



- 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

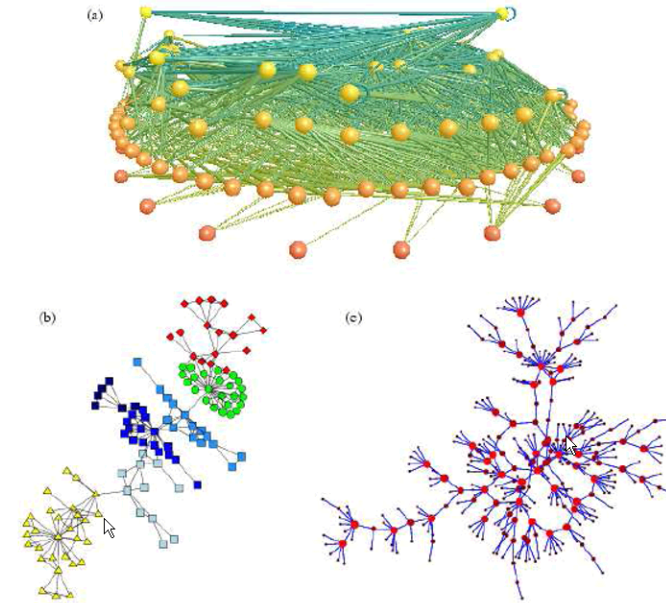


FIG. 2 Three examples of the kinds of networks that are the topic of this review. (a) A food web of predator-prey interactions between species in a freshwater lake [272]. Picture courtesy of Neo Martinez and Richard Williams. (b) The network of collaborations between scientists at a private research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

