

Machine Learning Lecture 8

Guassian Mixture Models

Qian Ma

Sun Yat-sen University

Spring Semester, 2020



Acknowledgement

- A large part of slides in this lecture are originally from
 - Prof. Andrew Ng (Stanford University)
 - Prof. Shuai Li (Shanghai Jiao Tong University)



Prof. Andrew Ng
Stanford University



Prof. Shuai Li
Shanghai Jiao Tong University

Generative Models (review)

Discriminative / Generative Models

- Discriminative models
 - Modeling the **dependence** of unobserved variables on observed ones
 - also called conditional models
 - Deterministic: $y = f_{\theta}(x)$
 - Probabilistic: $p_{\theta}(y|x)$
- Generative models
 - Modeling the **joint** probabilistic distribution of data
 - Given some hidden parameters or variables

$$p_{\theta}(x, y)$$

- Then do the conditional inference

$$p_{\theta}(y|x) = \frac{p_{\theta}(x, y)}{p_{\theta}(x)} = \frac{p_{\theta}(x, y)}{\sum_{y'} p_{\theta}(x, y')}$$

Discriminative Models

- Discriminative models
 - Modeling the **dependence** of unobserved variables on observed ones
 - also called conditional models
 - Deterministic: $y = f_{\theta}(x)$
 - Linear regression
 - Probabilistic: $p_{\theta}(y|x)$
 - Logistic regression
- Directly model the dependence for label prediction
- Easy to define dependence on specific features and models
- Practically yielding higher prediction performance
- E.g. linear regression, logistic regression, k nearest neighbor, SVMs, (multi-layer) perceptrons, decision trees, random forest

Generative Models

- Generative models
 - Modeling the **joint** probabilistic distribution of data
 - Given some hidden parameters or variables

$$p_{\theta}(x, y)$$

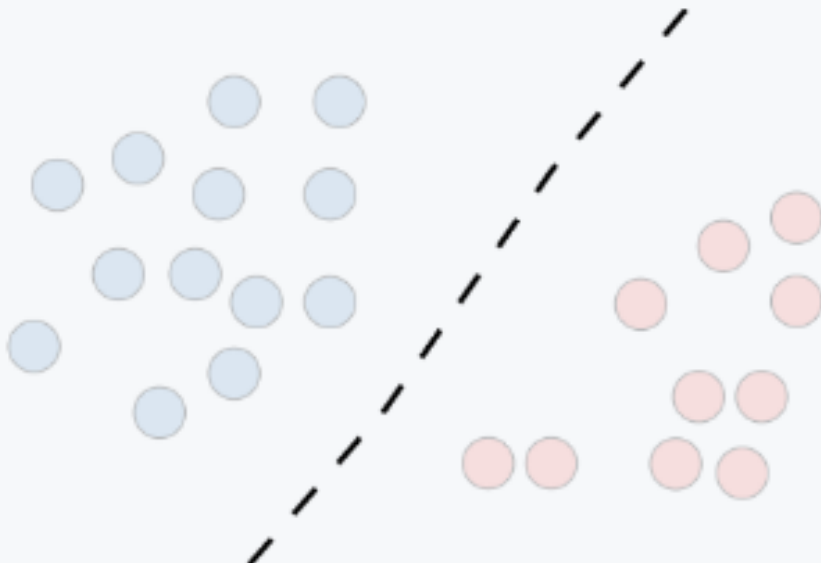
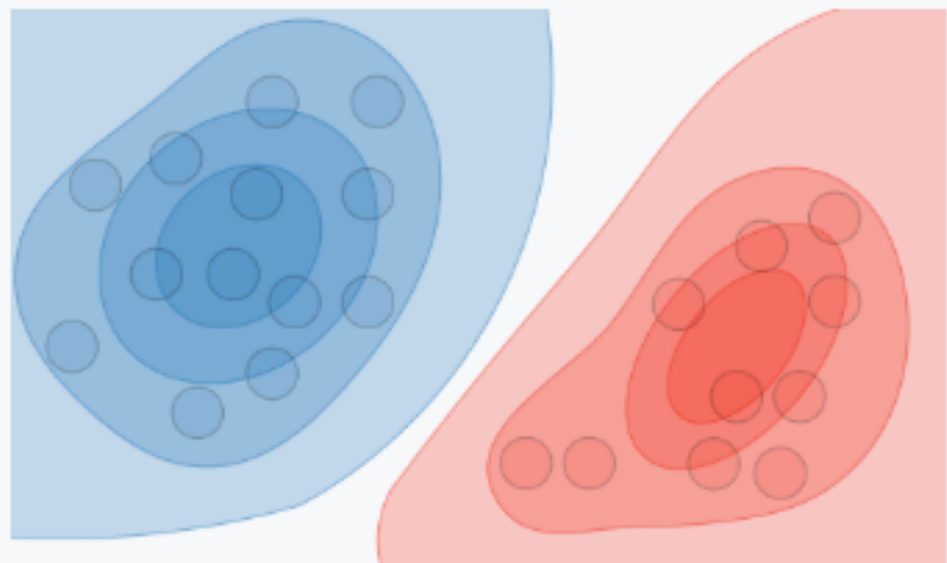
- Then do the conditional inference

$$p_{\theta}(y|x) = \frac{p_{\theta}(x, y)}{p_{\theta}(x)} = \frac{p_{\theta}(x, y)}{\sum_{y'} p_{\theta}(x, y')}$$

- Recover the data distribution [essence of data science]
- Benefit from hidden variables modeling
- E.g. Naive Bayes, Hidden Markov Model, Mixture Gaussian, Markov Random Fields, Latent Dirichlet Allocation

Discriminative Models vs Generative Models

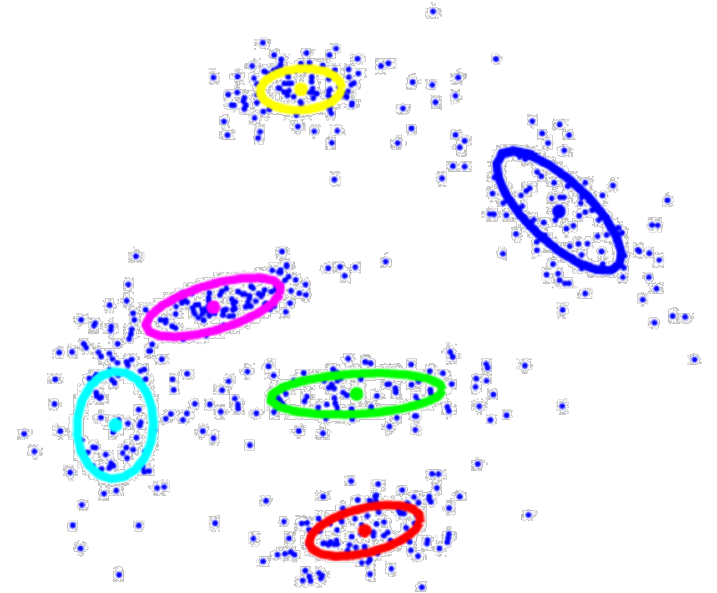
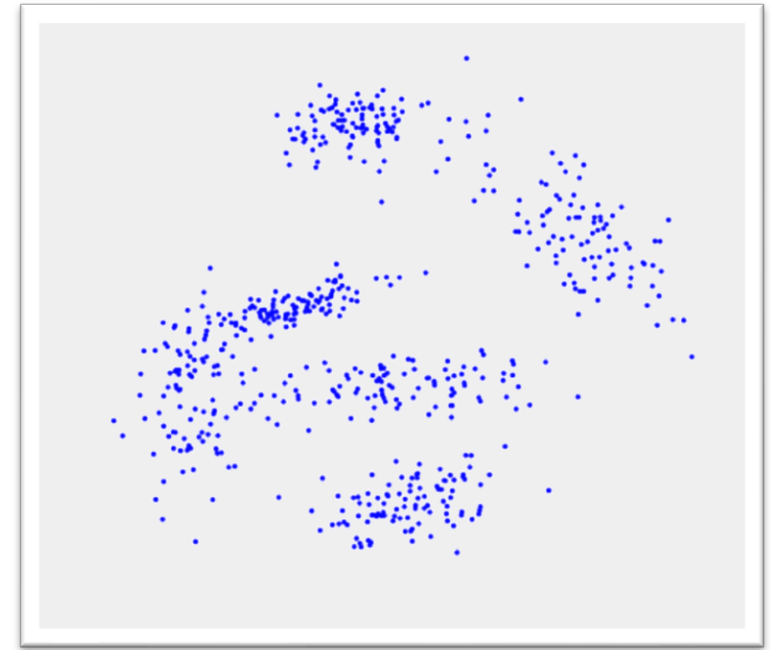
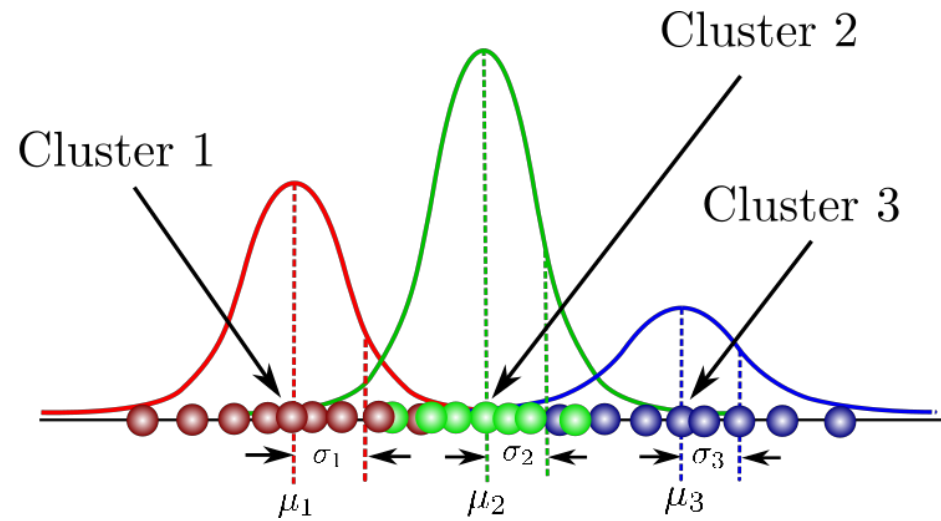
- In General
 - A Discriminative model models the **decision boundary between the classes**
 - A Generative Model explicitly models the **actual distribution of each class**
- Example: Our training set is a bag of fruits. Only **apples** and **oranges** are labeled. Imagine a post-it note stuck to the fruit
 - A generative model will model various attributes of fruits such as color, weight, shape, etc
 - A discriminative model might model color alone, **should that suffice** to distinguish apples from oranges

	Discriminative model	Generative model
Goal	Directly estimate $P(y x)$	Estimate $P(x y)$ to then deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data
Illustration		
Examples	Regressions, SVMs	GDA, Naive Bayes

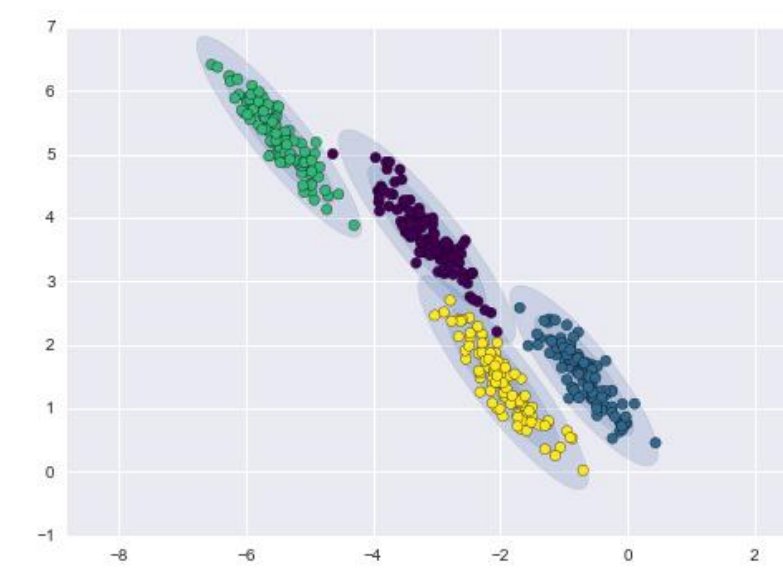
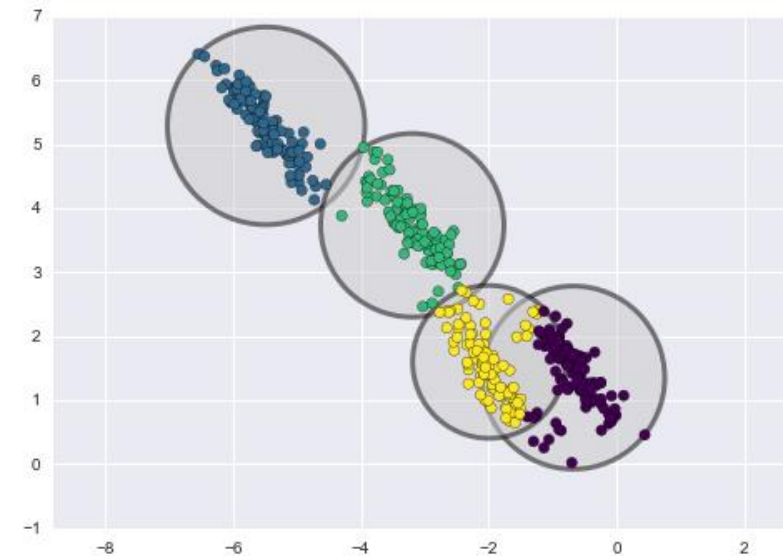
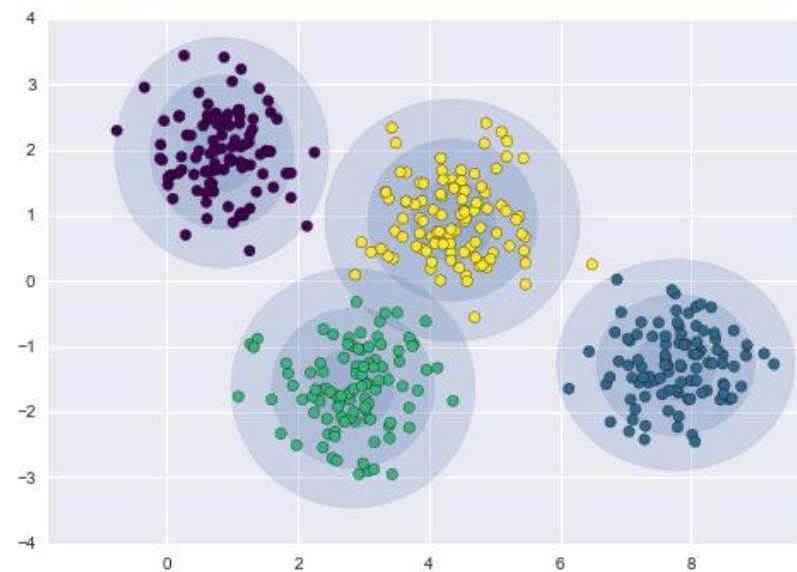
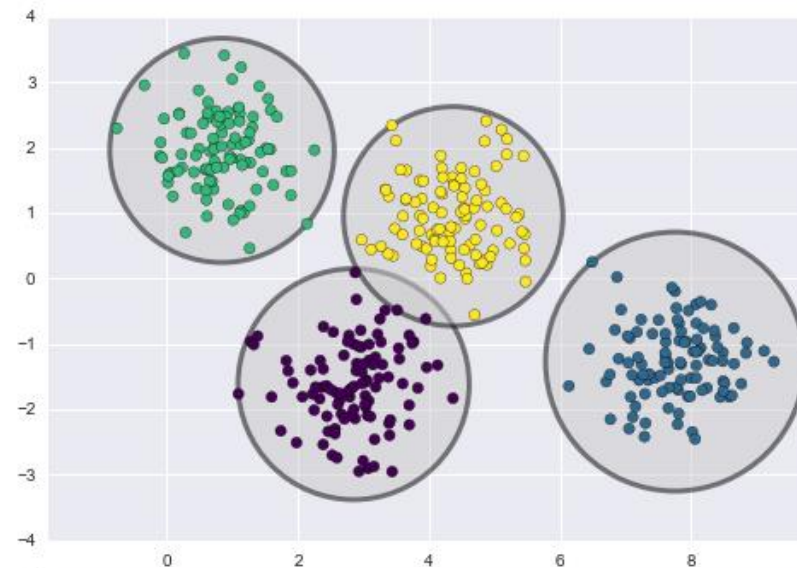
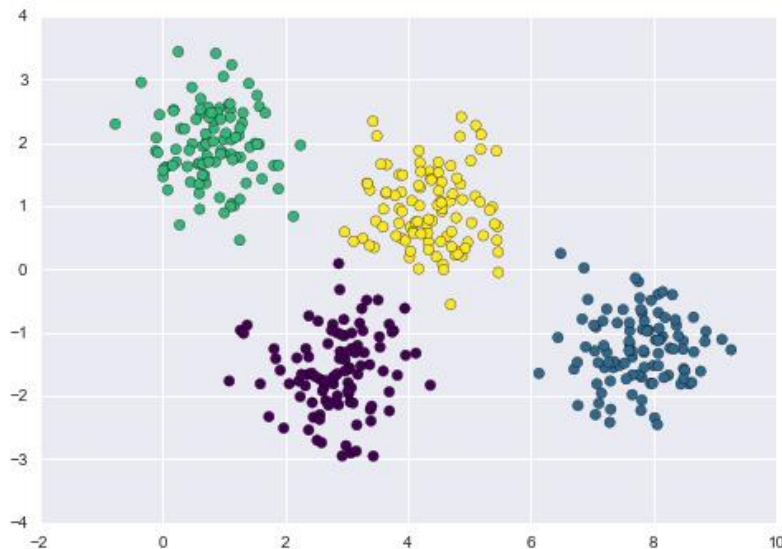
Gaussian Mixture Models

Gaussian Mixture Models

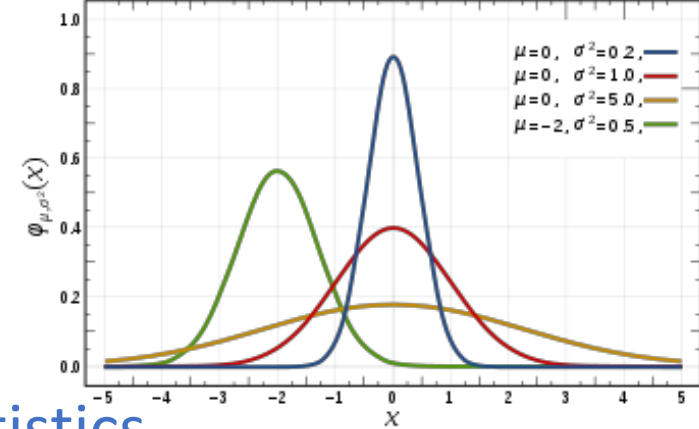
- Is a clustering algorithms
- Difference with K-means
 - K-means outputs the label of a sample
 - GMM outputs the probability that a sample belongs to a certain class
 - GMM can also be used to **generate** new samples!



K-means vs GMM



Gaussian distribution



- Very common in **probability theory** and important in **statistics**
- often used in the natural and social sciences to represent real-valued random variables whose distributions are **not known**
- is useful because of the **central limit theorem**
 - averages of samples independently drawn from the same distribution converge in distribution to the normal with the true mean and variance, that is, they become normally distributed when the number of observations is sufficiently large
- Physical quantities that are expected to be the sum of many independent processes often have distributions that are nearly normal
- The probability density of the Gaussian distribution is

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

High-dimensional Gaussian distribution

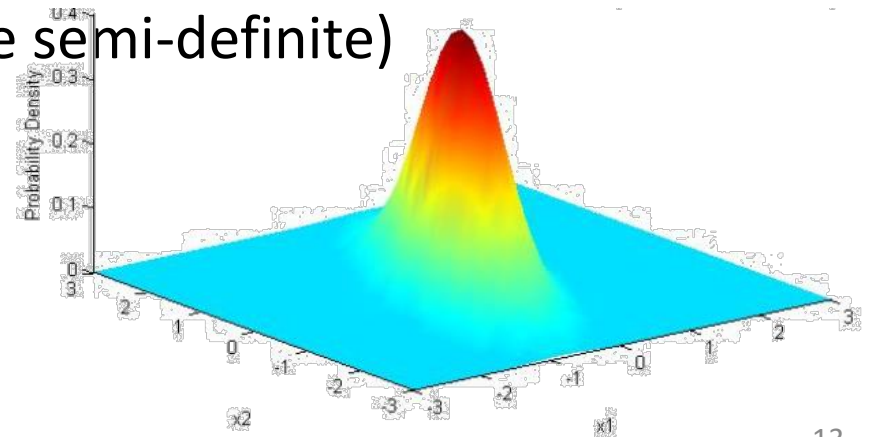
$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- The probability density of Gaussian distribution on $x = (x_1, \dots, x_d)^\top$ is

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

- where μ is the mean vector
- Σ is the symmetric covariance matrix (positive semi-definite)
- E.g. the Gaussian distribution with

$$\mu = (0,0)^T \quad \Sigma = \begin{pmatrix} 0.25 & 0.30 \\ 0.30 & 1.00 \end{pmatrix}$$



Mixture of Gaussian

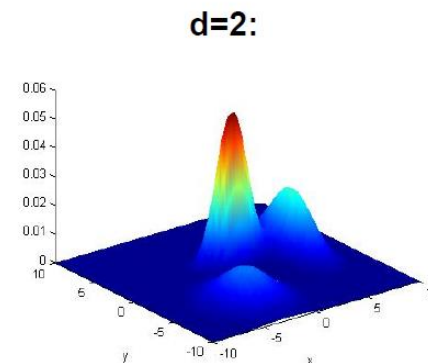
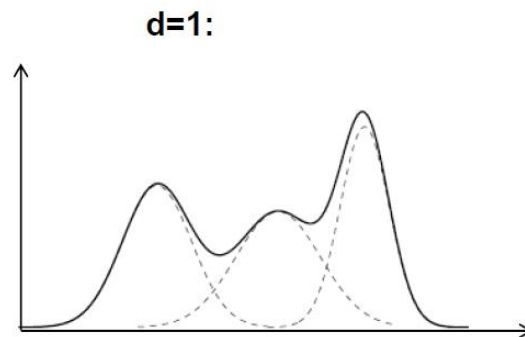
- The probability given in a mixture of K Gaussians is:

$$p(x) = \sum_{j=1}^K w_j \cdot \mathcal{N}(x|\mu_j, \Sigma_j)$$

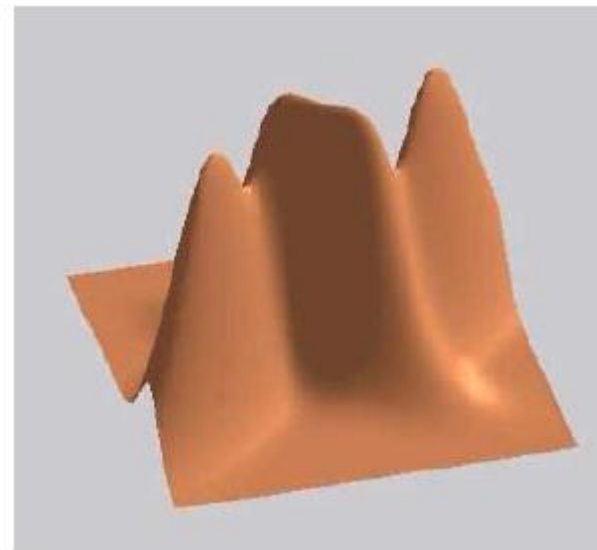
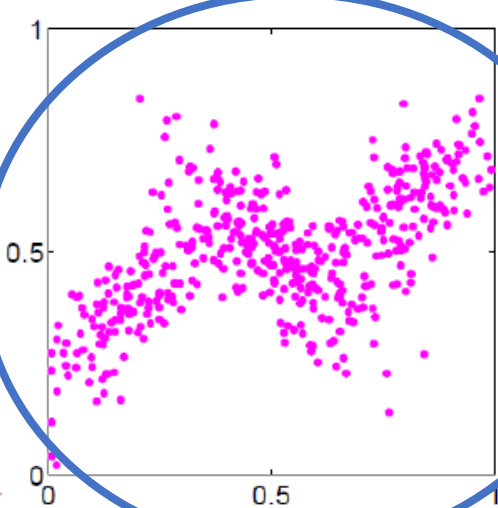
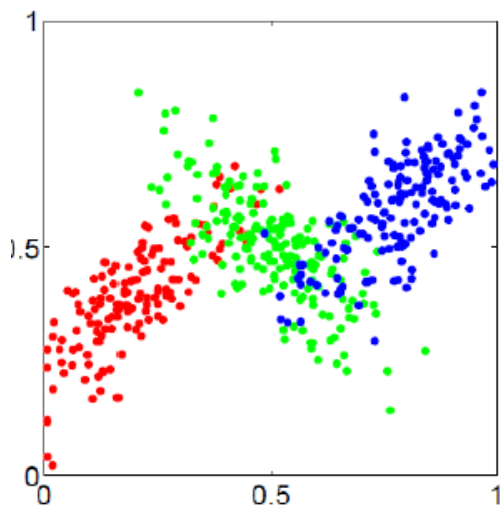
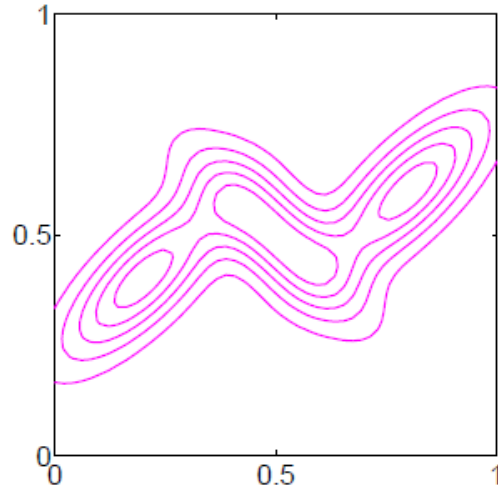
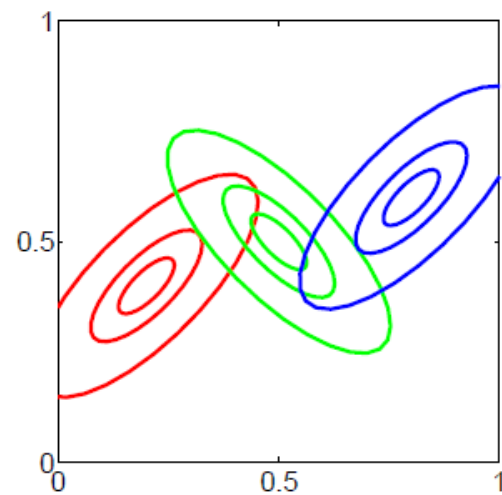
where w_j is the prior probability of the j -th Gaussian

$$\sum_{j=1}^K w_j = 1 \quad \text{and} \quad 0 \leq w_j \leq 1$$

- Example



Examples



Observation

Data generation

- Let the parameter set $\theta = \{w_j, \mu_j, \Sigma_j : j\}$, then the probability density of mixture Gaussian can be written as

$$p(x|\theta) = \sum_{j=1}^K w_j \cdot \mathcal{N}(x|\mu_j, \Sigma_j)$$

- Equivalent to generate data points in two steps
 - Select which component j the data point belongs to according to the categorical (or multinoulli) distribution of (w_1, \dots, w_K)
 - Generate the data point according to the probability of j -th component

Learning task

- Given a dataset $X = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$ to train the GMM model
- Find the best θ that maximizes the probability $\mathbb{P}(X|\theta)$
- Maximal likelihood estimator (MLE)

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

Introduce latent variable

- For data points $x^{(i)}, i = 1, \dots, N$, let's write the probability as

$$\mathbb{P}(x^{(i)}|\theta) = \sum_{j=1}^K w_j \cdot \mathcal{N}(x^{(i)}|\mu_j, \Sigma_j)$$

where $\sum_{j=1}^K w_j = 1$

- Introduce latent variable
 - $z^{(i)}$ is the Gaussian cluster ID indicates which Gaussian $x^{(i)}$ comes from
 - $\mathbb{P}(z^{(i)} = j) = w_j$ “Prior”
 - $\mathbb{P}(x^{(i)}|z^{(i)} = j; \theta) = \mathcal{N}(x^{(i)}|\mu_j, \Sigma_j)$
 - $\mathbb{P}(x^{(i)}|\theta) = \sum_{j=1}^K \mathbb{P}(z^{(i)} = j) \cdot \mathbb{P}(x^{(i)}|z^{(i)} = j; \theta)$

Maximal likelihood

- Let $l(\theta) = \sum_{i=1}^N \log \mathbb{P}(x^{(i)}; \theta) = \sum_{i=1}^N \log \sum_j \mathbb{P}(x^{(i)}, z^{(i)} = j; \theta)$ be the log-likelihood

- We want to solve

$$\operatorname{argmax} l(\theta) = \operatorname{argmax} \sum_{i=1}^N \log \sum_{j=1}^K \mathbb{P}(z^{(i)} = j) \cdot \mathbb{P}(x^{(i)} | z^{(i)} = j; \theta)$$

$$= \operatorname{argmax} \sum_{i=1}^N \log \sum_{j=1}^K w_j \cdot \mathcal{N}(x^{(i)} | \mu_j, \Sigma_j)$$

- No closed solution by solving

$$\frac{\partial \mathbb{P}(X|\theta)}{\partial w} = 0, \quad \frac{\partial \mathbb{P}(X|\theta)}{\partial \mu} = 0, \quad \frac{\partial \mathbb{P}(X|\theta)}{\partial \Sigma} = 0$$

Likelihood maximization

- If we know $z^{(i)}$ for all i , the problem becomes

$$\operatorname{argmax} l(\theta) = \operatorname{argmax} \sum_{i=1}^N \log \mathbb{P}(x^{(i)} | \theta)$$

$$= \operatorname{argmax} \sum_{i=1}^N \log \mathbb{P}(x^{(i)}, z^{(i)} | \theta)$$

$$= \operatorname{argmax} \sum_{i=1}^N \log \mathbb{P}(x^{(i)} | \theta, z^{(i)}) + \log \mathbb{P}(z^{(i)} | \theta)$$

$$= \operatorname{argmax} \sum_{i=1}^N \log \mathcal{N}(x^{(i)} | \mu_{z^{(i)}}, \Sigma_{z^{(i)}}) + \log w_{z^{(i)}}$$

The solution is

- $w_j = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}$
- $\mu_j = \frac{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}}$
- $\Sigma_j = \frac{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}}$

Average over
each cluster

Likelihood maximization (cont.)

- Given the parameter $\theta = \{w_j, \mu_j, \Sigma_j : j\}$, the **posterior** distribution of each latent variable $z^{(i)}$ can be inferred

- $$\mathbb{P}(z^{(i)} = j | x^{(i)}; \theta) = \frac{\mathbb{P}(x^{(i)}, z^{(i)} = j | \theta)}{\mathbb{P}(x^{(i)} | \theta)}$$

- $$= \frac{\mathbb{P}(x^{(i)} | z^{(i)} = j; \mu_j, \Sigma_j) \mathbb{P}(z^{(i)} = j | w)}{\sum_{j'=1}^K \mathbb{P}(x^{(i)} | z^{(i)} = j'; \mu_{j'}, \Sigma_{j'}) \mathbb{P}(z^{(i)} = j' | w)}$$

- Or $w_j^{(i)} = \mathbb{P}(z^{(i)} = j | x^{(i)}; \theta) \propto \mathbb{P}(x^{(i)} | z^{(i)} = j; \mu_j, \Sigma_j) \mathbb{P}(z^{(i)} = j | w)$

Likelihood maximization (cont.)

$$w_j^{(i)} = \mathbb{P}(z^{(i)} = j | x^{(i)}; \theta)$$

- For every possible values of $z^{(i)}$'s

The solution is

- $w_j = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}$
- $\mu_j = \frac{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}}$
- $\Sigma_j = \frac{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^\top}{\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}}$



Which is equivalent to

- $w_j = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}$
- $(\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}) \mu_j = \sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} x^{(i)}$
- $(\sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\}) \Sigma_j = \sum_{i=1}^N \mathbf{1}\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^\top$



Take
expectation

- Take the expectation on **probability** of $z^{(i)}$

The solution is

- $w_j = \frac{1}{N} \sum_{i=1}^N w_j^{(i)}$
- $\mu_j = \frac{\sum_{i=1}^N w_j^{(i)} x^{(i)}}{\sum_{i=1}^N w_j^{(i)}}$
- $\Sigma_j = \frac{\sum_{i=1}^N w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^\top}{\sum_{i=1}^N w_j^{(i)}}$



Take expectation on two sides

- $w_j = \frac{1}{N} \sum_{i=1}^N \mathbb{P}(z^{(i)} = j)$
- $(\sum_{i=1}^N \mathbb{P}(z^{(i)} = j)) \mu_j = \sum_{i=1}^N \mathbb{P}(z^{(i)} = j) x^{(i)}$
- $(\sum_{i=1}^N \mathbb{P}(z^{(i)} = j)) \Sigma_j = \sum_{i=1}^N \mathbb{P}(z^{(i)} = j) (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^\top$

Expectation maximization methods

- E-step:
 - Infer the **posterior** distribution of the latent variables given the model parameters
- M-step:
 - Tune parameters to maximize the data likelihood given the latent variable distribution
- EM methods
 - Iteratively execute E-step and M-step until convergence

EM for GMM

- Repeat until convergence: {

(E-step) For each i, j , set

$$w_j^{(i)} = \mathbb{P}(z^{(i)} = j | x^{(i)}, w, \mu, \Sigma)$$

(M-step) Update the parameters

$$w_j = \frac{1}{N} \sum_{i=1}^N w_j^{(i)}, \quad \mu_j = \frac{\sum_{i=1}^N w_j^{(i)} x^{(i)}}{\sum_{i=1}^N w_j^{(i)}}$$

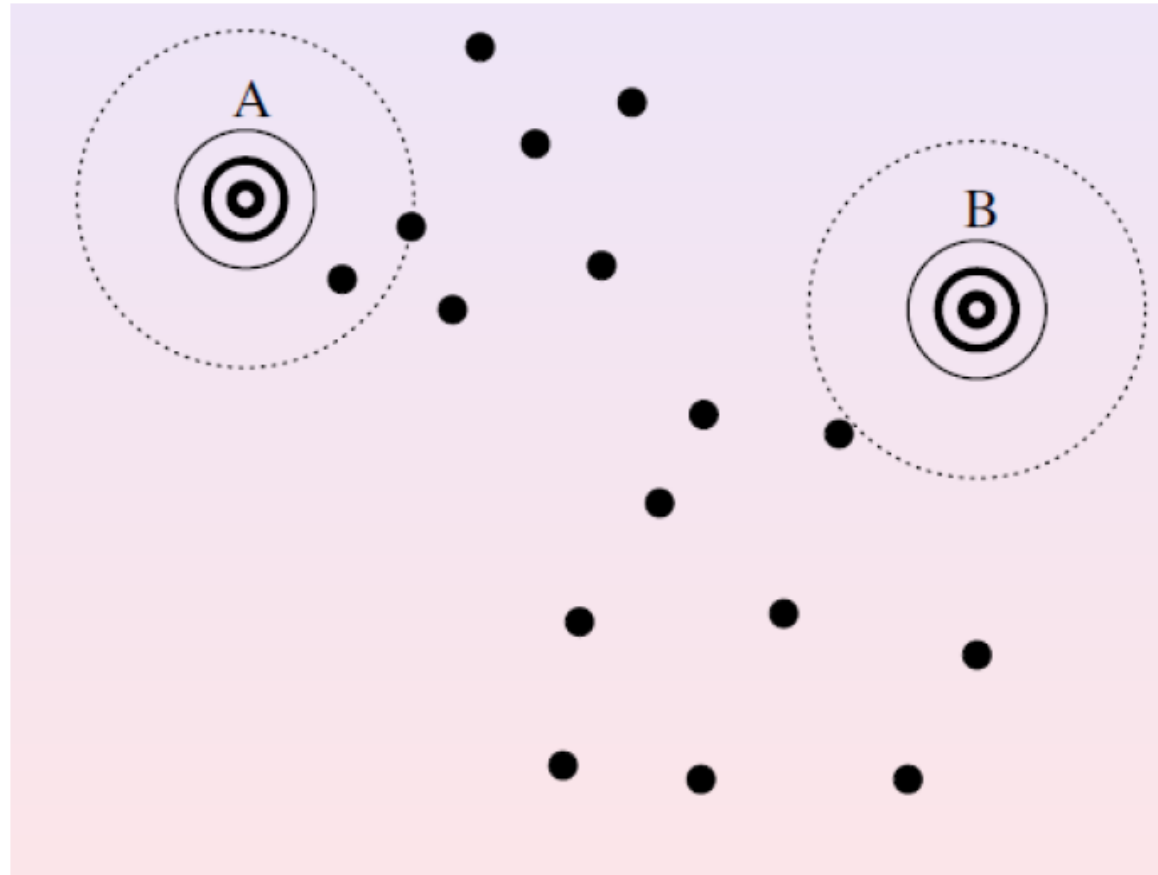
$$\Sigma_j = \frac{\sum_{i=1}^N w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^{\top}}{\sum_{i=1}^N w_j^{(i)}}$$

}

Example

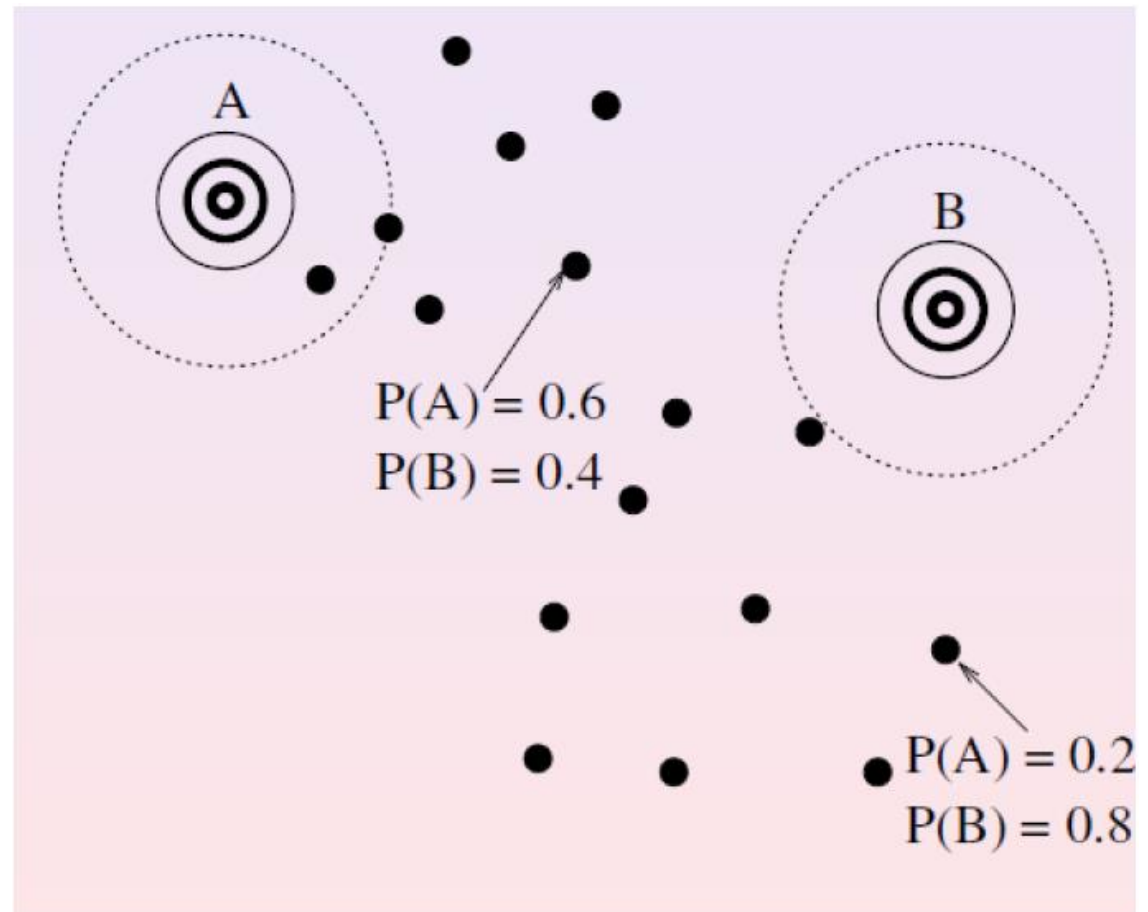
$$p(x|\theta) = \sum_{j=1}^K w_j \cdot \mathcal{N}(x|\mu_j, \Sigma_j)$$

- Hidden variable: for each point, which Gaussian generates it?



Example (cont.)

- E-step: for each point, estimate the probability that each Gaussian component generated it



Example (cont.)

- M-Step: modify the parameters according to the hidden variable to maximize the likelihood of the data

