

Mixture density models

Victor Kitov

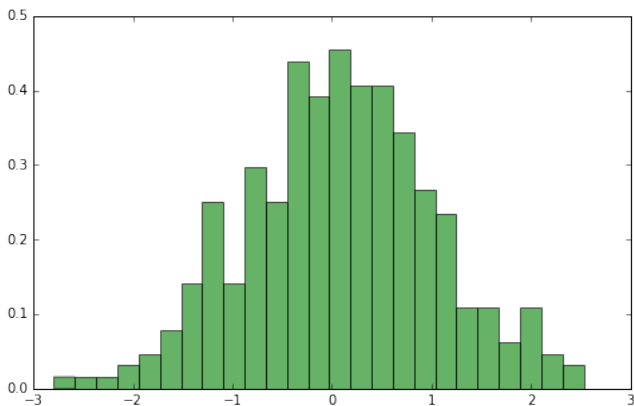
v.v.kitov@yandex.ru

Table of Contents

- 1 Mixture models
- 2 K-means
- 3 Simplifications of Gaussian mixtures

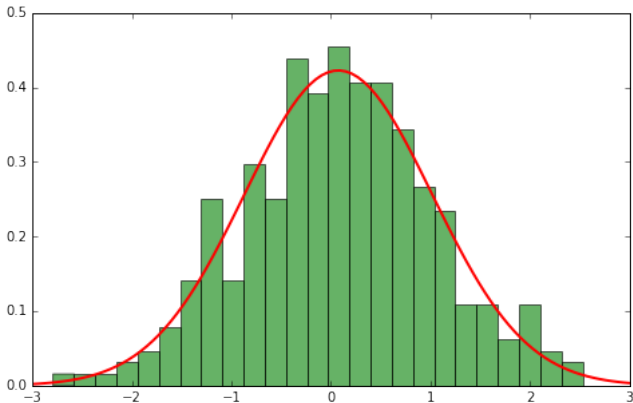
Sample density

Consider sample density:



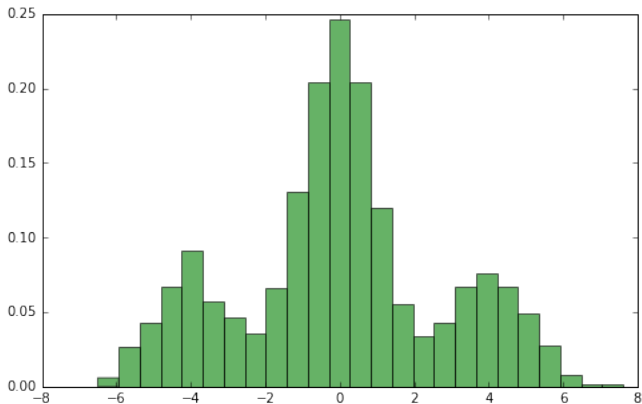
Parametric density approximation

It can be accurately modelled with existing parametric family -
Normal



Non-standard sample distribution

What to do if no parametric model fits well?



Mixture models

$$p(x) = \sum_{z=1}^Z \phi_z p(x; \theta_z)$$

- Z - number of components
- $\phi_z, z = 1, 2, \dots, Z$ - mixture component probabilities,
 $\phi_z \geq 0, \sum_{z=1}^Z \phi_z = 1$
- $p(x; \theta_z)$ - component density functions
- Parameters of mixture model $\Theta = \{\phi_z, \theta_z, z = 1, 2, \dots, Z\}$

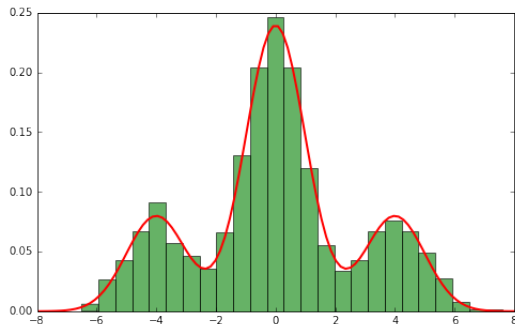
$p(x, \theta_z)$ may be of single or different parametric families.

Mixture of Gaussians

Gaussians model continuous r.v. on $(-\infty, +\infty)$.

$$p(x, \theta_z) = N(x, \mu_z, \Sigma_z), \theta_z = \{\mu_z, \Sigma_z\}.$$

$$p(x) = \sum_{z=1}^Z \phi_z N(x, \mu_z, \Sigma_z) \quad (1)$$



Mixtures of other distributions

Mixture of random variables:

- continuous, distributed on $(-\infty, +\infty)$
- continuous, distributed on $[a, \infty)$
- continuous, distributed on $[a, b]$
- discrete, distributed on $[a, \infty)$
- discrete, distributed on $[a, b]$

Mixtures of other distributions

Mixture of random variables:

- continuous, distributed on $(-\infty, +\infty)$
 - Normal, Laplace, Student
- continuous, distributed on $[a, \infty)$
 - Gamma
- continuous, distributed on $[a, b]$
 - Beta
- discrete, distributed on $[a, \infty)$
 - Poisson
- discrete, distributed on $[a, b]$
 - Binomial

Sampling from mixture

- 1 Sample mixture component z with random probabilities $\phi_1, \phi_2, \dots, \phi_Z$

Sampling from mixture

- ① Sample mixture component z with random probabilities $\phi_1, \phi_2, \dots, \phi_Z$
 - to do that we sample $u \sim \text{Uniform}[0, 1]$ and select component z if $\sum_{k=1}^{z-1} \phi_k < u \leq \sum_{k=1}^z \phi_k$
- ② Sample observation $x \sim p(x|\theta_z)$

Classification using mixtures

Model within class probability with mixtures:

$$p(x|y) = \sum_{z=1}^{Z_y} \phi_{y,z} p(x; \theta_{y,z})$$

where Z_y , $\pi_{y,z}$ and $p(x; \theta_{y,z})$ are specific for each class y .

Bayes decision rule:

$$\hat{y} = \arg \max_y \lambda_y p(y) p(x|y)$$

λ_y - cost for misclassifying class y

$p(y)$ - prior for class y

$p(x|y)$ - within class probability

EM-algorithm for normal mixtures

Initialize ϕ_j, μ_j and Σ_j , $j = 1, 2, \dots, g$.

REPEAT UNTIL convergence:

E-STEP. Calculate correspondences of x_n
to component z :

FOR $n = 1, 2, \dots, N$:

for $z = 1, 2, \dots, Z$:

$$w_{nz} = \frac{\phi_z N(x_n; \mu_z, \Sigma_z)}{\sum_k \phi_k N(x_n; \mu_k, \Sigma_k)} \quad \# = p(z|x(n))$$

M-STEP. Update component parameters:

FOR $z = 1, 2, \dots, Z$:

$$\hat{\phi}_z = \frac{1}{N} \sum_{n=1}^N w_{nz}$$

$$\hat{\mu}_z = \frac{\sum_{n=1}^N w_{nz} x_n}{\sum_{n=1}^N w_{nz}}$$

$$\hat{\Sigma}_z = \frac{1}{\sum_{n=1}^N w_{nz}} \sum_{n=1}^N w_{nz} (x_n - \hat{\mu}_z)(x_n - \hat{\mu}_z)^T$$

Interpretation

$$\begin{aligned}
 w_{nz} &= P(z|x_n) = \frac{P(z, x_n)}{P(x_n)} = \frac{P(z, x_n)}{\sum_k P(k, x_n)} = \\
 &= \frac{P(z)P(x_n|z)}{\sum_k P(k)P(x_n|k)} = \frac{\hat{\phi}_z N(x_n; \hat{\mu}_z, \hat{\Sigma}_z)}{\sum_k \hat{\phi}_k N(x_n; \hat{\mu}_k, \hat{\Sigma}_k)}
 \end{aligned}$$

$\hat{\phi}_z, \hat{\mu}_z, \hat{\Sigma}_z$ are weighted averages, weighted with $w_{nz} = P(z|x_n)$:

$$\begin{aligned}
 \hat{\phi}_z &= \frac{1}{N} \sum_{n=1}^N w_{nz} & \hat{\mu}_z &= \frac{\sum_{n=1}^N w_{nz} x_n}{\sum_{n=1}^N w_{nz}} \\
 \hat{\Sigma}_z &= \frac{1}{\sum_{n=1}^N w_{nz}} \sum_{n=1}^N w_{nz} (x_n - \hat{\mu}_z)(x_n - \hat{\mu}_z)^T
 \end{aligned} \tag{2}$$

Table of Contents

- 1 Mixture models
- 2 K-means
- 3 Simplifications of Gaussian mixtures

K-means algorithm

- Suppose we want to cluster our data into g clusters.
- Cluster i has a center μ_i , $i=1,2,\dots,g$.
- Consider the task of minimizing

$$\sum_{n=1}^N \rho(x_n, \mu_{z_n})^2 \rightarrow \min_{z_1, \dots, z_N, \mu_1, \dots, \mu_g} \quad (3)$$

where $z_i \in \{1, 2, \dots, g\}$ is cluster assignment for x_i and μ_1, \dots, μ_g are cluster centers.

- Direct optimization requires full search and is impractical.
- K-means is a suboptimal algorithm for optimizing (3).

K-means algorithm

Initialize μ_j , $j = 1, 2, \dots, g$ # usually by setting them
to randomly chosen $x(n)$

REPEAT UNTIL convergence:

FOR $i = 1, 2, \dots, N$: # cluster assignments

$$z_i = \arg \min_{j \in \{1, 2, \dots, g\}} \|x_i - \mu_j\|$$

FOR $j = 1, 2, \dots, g$: # means recalculation

$$\mu_j = \frac{1}{\sum_{n=1}^N \mathbb{I}[z_n = j]} \sum_{n=1}^N \mathbb{I}[z_n = j] x_i$$

Possible stop conditions:

- cluster assignments z_1, \dots, z_N stop to change (typical)
- maximum number of iterations reached
- cluster means $\{\mu_i, i = 1, 2, \dots, g\}$ stop changing significantly

K-means versus EM clustering

K-means versus EM clustering

- For each x_n EM algorithm gives $w_{nz} = p(z|x_n)$.
- This is soft or probabilistic clustering into Z clusters, having priors ϕ_1, \dots, ϕ_Z and probability distributions $p(x; \theta_1), \dots, p(x; \theta_Z)$.
- We can make it hard clustering using $z_n = \arg \max_z w_{nz}$.
- EM clustering becomes K-means clustering when:
 - applied to Gaussians
 - with equal priors
 - with unity covariance matrices
 - with hard clustering

Initialization for Gaussian mixture EM

- 1 Fit K-means to x_1, x_2, \dots, x_N , obtain cluster centers $\mu_z, z = 1, 2, \dots, Z$ and cluster assignments z_1, z_2, \dots, z_N .
- 2 Initialize mixture probabilities

$$\hat{\phi}_z \propto \sum_{n=1}^N \mathbb{I}[z_n = z]$$

- 3 Initialize Gaussian means with cluster centers $\mu_z, z = 1, 2, \dots, Z$.
- 4 Initialize Gaussian covariance matrices with

$$\hat{\Sigma}_z = \frac{1}{\sum_{n=1}^N \mathbb{I}[z_n = z]} \sum_{n=1}^N \mathbb{I}[z_n = z] (x_n - \mu_z)(x_n - \mu_z)^T$$

Properties of EM

- Many local optima exist
 - in particular likelihood $\rightarrow \infty$ when $\mu_z = x_i$ and $\sigma_z \rightarrow 0$
- Only local optimum is found with EM
- Results depends on initialization
 - It is common to run algorithm multiple times with different initializations and then select the result maximizing the likelihood function.
- Number of components may be selected with:

Properties of EM

- Many local optima exist
 - in particular likelihood $\rightarrow \infty$ when $\mu_z = x_i$ and $\sigma_z \rightarrow 0$
- Only local optimum is found with EM
- Results depends on initialization
 - It is common to run algorithm multiple times with different initializations and then select the result maximizing the likelihood function.
- Number of components may be selected with:
 - cross-validation on the final task
 - out-of-sample maximum likelihood
 - statistical tests, heuristics, such as AIC/BIC information criteria

Table of Contents

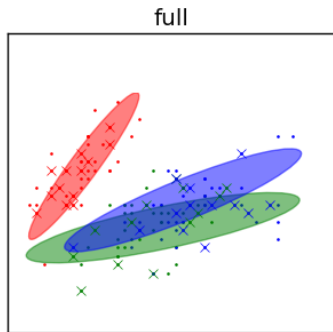
- 1 Mixture models
- 2 K-means
- 3 Simplifications of Gaussian mixtures**

Simplifications of Gaussian mixtures

- $\Sigma_z \in \mathbb{R}^{D \times D}$ requires $\frac{D(D+1)}{2}$ parameters.
- Covariance matrices for Z components require $Z \frac{D(D+1)}{2}$ parameters.
- Components can be poorly identified when
 - $Z \frac{D(D+1)}{2}$ is large compared to N
 - when components are not well separated
- In these cases we can impose restrictions on covariance matrices.

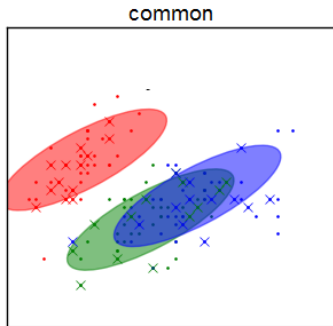
Unrestricted covariance matrices

- full covariance matrices Σ_z , $z = 1, 2, \dots, Z$.



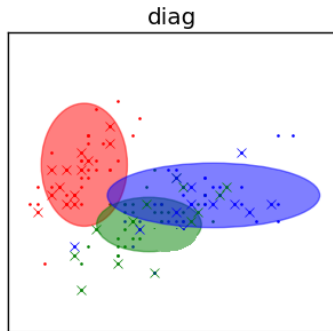
Common covariance matrix

- $\Sigma_1 = \Sigma_2 = \dots = \Sigma_Z$



Diagonal covariance matrices

- $\Sigma_z = \text{diag}\{\sigma_{z,1}^2, \sigma_{z,2}^2 \dots \sigma_{z,D}^2\}$



Spherical matrices

- $\Sigma_z = \sigma_z^2 I$, $I \in \mathbb{R}^{D \times D}$ - identity matrix

