

Decision theory. Fisher's LDA, QDA. Mixtures. EM algorithm

V. Kitov



Yandex School of Data Analysis

Imperial College London
Department of Physics

January 2015

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions
- 4 Discriminant functions
- 5 Gaussian classifiers
- 6 Mixture models

Bayes Decision Rule for Minimum Error

- C classes: $\omega_1, \omega_2, \dots, \omega_C$.
- Need to predict class for target object.

Initial solution (features unknown)

$$\hat{c} = \arg \max_c p(\omega_c)$$

- Now features x of object are observed.

Solution after observing x

$$\hat{c} = \arg \max_c p(\omega_c | x)$$

Reformulation

Rewrite class posterior probability:

$$p(\omega_c|x) = \frac{p(\omega_c, x)}{p(x)} = \frac{p(\omega_c)p(x|\omega_c)}{p(x)}$$

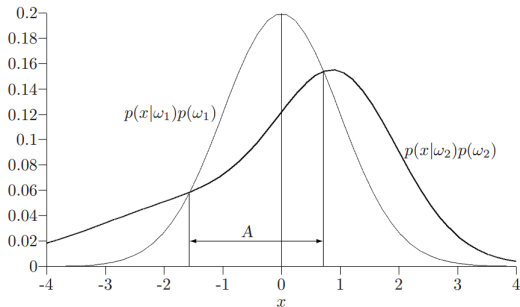
Reformulated solution after observing x

$$\hat{c} = \arg \max_c p(\omega_c)p(x|\omega_c)$$

Two class case:

$$\text{Assign to class } \omega_1 \text{ if } \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{p(\omega_2)}{p(\omega_1)}$$

Illustration of decision rule



Class ω_1 is selected in region A , class ω_2 selected outside of A

Minimum error property

Let $\Omega_1, \Omega_2, \dots, \Omega_C$ be decision regions for classes $\omega_1, \omega_2, \dots, \omega_C$ and Ω denotes the whole features space.

Then probability of correct classification equals:

$$\begin{aligned} p(\text{correct}) &= \sum_{c=1}^C p(\text{correct}, \omega_c) = \sum_{c=1}^C p(x \in \Omega_c, \omega_c) = \\ &= \sum_{c=1}^C \int_{x \in \Omega_c} p(x, \omega_c) dx = \sum_{c=1}^C \int_{x \in \Omega_c} p(x) p(\omega_c | x) dx \end{aligned}$$

Probability of correct classification is maximized when for each c
 $\Omega_c = \{x : c = \arg \max_i p(\omega_i | x)\}$ - Bayes rule!

$$p(\text{error}) = 1 - p(\text{correct})$$

Probability of error is minimized for Bayes decision rule.

Probabilities for Bayes minimum error decision rule

By definition of Bayes decision rule for each c :

$$p(\text{correct}, \omega_c) = \int_{\Omega_c} p(\omega_c) p(x|\omega_c) dx = \int_{\Omega_c} \max_i p(\omega_i) p(x|\omega_i) dx$$

Probability of correct classification equals (using property $\bigcap_c \Omega_c = \emptyset, \bigcup_c \Omega_c = \Omega$):

$$p(\text{correct}) = \sum_{c=1}^C p(\text{correct}, \omega_c) = \int \max_i p(\omega_i) p(x|\omega_i) dx$$

Probability of erroneous classification:

$$p(\text{error}) = 1 - \int \max_i p(\omega_i) p(x|\omega_i) dx$$

Reject option

- Most of the misclassification errors occurs when algorithm is unsure:

$$p(\omega_{\hat{c}}|x) \text{ is small, where } \hat{c} = \arg \max_c p(\omega_c)$$

- Introduce **reject option**: if $\max_c p(\omega_c)$ is small, then reject classification and
 - leave as third class
 - classify is later
 - when new information about object becomes available
 - classify using different algorithm

Rejection region

Rejection region

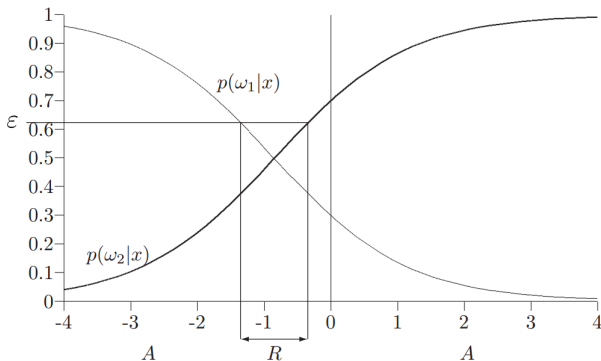
$$R = \{x : \max_c p(\omega_c|x) < \varepsilon\}$$

Acceptance region

$$A = \{x : \max_c p(\omega_c|x) \geq \varepsilon\}$$

Trade-off: larger reject region will decrease error-rate, but more otherwise correctly classified objects will be rejected.

Illustration of rejection region



Acceptance and rejection regions.

Probabilities

Correct classification with reject option implies:

- object was not rejected
- object was correctly classified

Probability of correct classification:

$$p(\text{correct}) = \int_A \max_i p(\omega_i) p(x|\omega_i) dx$$

Probability of rejection:

$$p(\text{reject}) = \int_R p(x) dx$$

Probability of error:

$$p(\text{error}) = \int_A (1 - \max_i p(\omega_i) p(x|\omega_i)) dx = 1 - p(\text{reject}) - p(\text{correct})$$

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions
- 4 Discriminant functions
- 5 Gaussian classifiers
- 6 Mixture models

Motivation

- In most practical situations loss from misclassifying different classes is different.
- Examples:
 - e-mail filtering: spam/not spam
 - medical diagnostics: ill/healthy
 - bank crediting: reliable borrower/unreliable borrower
 - network intrusion detection: intrusion/usual use
- Possible solution:

Assign different costs to different misclassifications.

Derivation

- Define $C \times C$ loss matrix $\mathbf{\Lambda}$, where λ_{ij} equals cost when object belongs to class ω_i but was assigned to class ω_j .
- Average loss equals

$$\begin{aligned}\mathbf{E}[L] &= \sum_i \sum_j \lambda_{ij} p(\omega_i, \hat{\omega}_j) = \\ &= \sum_i \sum_j \lambda_{ij} \int_{\Omega_j} p(x, \omega_i) dx = \sum_j \int_{\Omega_j} \sum_i \lambda_{ij} p(x, \omega_i) dx\end{aligned}$$

- It follows that minimum average solution satisfies:

$$\begin{aligned}\Omega_c &= \{x : c = \arg \min_j \sum_i \lambda_{ij} p(x, \omega_i) \\ &= \arg \min_j \sum_i \lambda_{ij} p(\omega_i | x) p(x) = \arg \min_j \sum_i \lambda_{ij} p(\omega_i | x)\}\end{aligned}$$

Bayes minimum risk classification

Bayes minimum risk classification

$$\hat{c} = \arg \min_c \sum_i \lambda_{ic} p(\omega_i | x)$$

Average risk attained by this decision rule equals:

$$\begin{aligned} \mathbf{E}[L] &= \sum_i \sum_j \lambda_{ij} p(\omega_i, \hat{\omega}_j) = \\ &= \sum_i \sum_j \lambda_{ij} \int_{\Omega_j} p(x, \omega_i) dx = \sum_j \int_{\Omega_j} \sum_i \lambda_{ij} p(x, \omega_i) dx \\ &= \int \min_j \sum_i \lambda_{ij} p(x, \omega_i) dx \end{aligned}$$

Bayes minimum risk classification

Two class case for $\lambda_{11} = \lambda_{22} = 0$:

$$\text{Assign to class } \omega_1 \text{ if } \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{\lambda_{21}p(\omega_2)}{\lambda_{12}p(\omega_1)}$$

For zero-one loss matrix

$$\lambda_{ij} = \begin{cases} 0, & i = j \\ 1, & i \neq j \end{cases}$$

Bayes minimum risk solution reduces to Bayes minimum error solution.

Reject option

Define conditional risk of assigning object x to class ω_j :

$$loss^j(x) = \sum_i \lambda_{ij} p(\omega_i | x)$$

For Bayes minimum risk decision rule expected loss for object x is

$$loss(x) = \min_j loss^j(x)$$

We can reject classification if expected cost of misclassification object is too high.

$$\begin{array}{ll} \text{rejection region:} & R = \{x : loss(x) > \varepsilon\} \\ \text{acceptance region:} & A = \{x : loss(x) \leq \varepsilon\} \end{array}$$

Reject option

Expected risk with rejection option equals:

$$\mathbf{E}[L|\text{no rejection}] = \int_A \min_j \sum_i \lambda_{ij} p(x, \omega_i) dx$$

Probability of rejection:

$$p(\text{reject}) = \int_R p(x) dx$$

Bayes minimum risk for equal within class costs

Earlier we have obtained Bayes minimum risk decision rule:

$$\hat{c} = \arg \min_c \sum_i \lambda_{ic} p(\omega_i | x)$$

Suppose the costs matrix Λ has simpler structure:

$$\lambda_{ic} = \begin{cases} 0, & i = c \\ \gamma_i, & i \neq c \end{cases}$$

Then Bayes minimum risk becomes:

$$\begin{aligned} \hat{c} &= \arg \min_c \sum_{i \neq c} \gamma_i p(\omega_i | x) = \arg \min_c \left\{ \sum_i \gamma_i p(\omega_i | x) - \gamma_c p(\omega_c | x) \right\} \\ &= \arg \max_c \gamma_c p(\omega_c | x) \end{aligned}$$

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions**
- 4 Discriminant functions
- 5 Gaussian classifiers
- 6 Mixture models

Neyman-Pearson rule

- Applied for two class cases. Class ω_1 will be the target class
- Two types of errors incurred:
 - False alarm: $p(\hat{\omega}_1|\omega_2)$
 - Missed detection: $p(\hat{\omega}_2|\omega_1)$
- Instead of maximizing general performance we minimize false alarm rate, given tolerable missed detection rate.

$$\begin{cases} \int_{\Omega_2} p(x|\omega_1) \rightarrow \min_{\Omega_1} \\ \int_{\Omega_1} p(x|\omega_2) = \varepsilon \end{cases}$$

- Neyman-Pearson solution:

Assign to class ω_1 if $\frac{p(x|\omega_1)}{p(x|\omega_2)} > \mu$, where μ is such that $\int_{\Omega_1} p(x|\omega_2) = \varepsilon$

Minimax criterion

- Bayes minimum error rule:

$$\begin{aligned} p(\text{error}) &= p(\omega_1) \int_{\Omega_2} p(x|\omega_1) dx + p(\omega_2) \int_{\Omega_1} p(x|\omega_2) dx \\ &= p(\omega_1) \int_{\Omega_2} p(x|\omega_1) dx + (1 - p(\omega_1)) \int_{\Omega_1} p(x|\omega_2) dx \end{aligned}$$

- Suppose, prior class probabilities may change in future.
- Optimize Bayes rule for worst case of prior probabilities:

$$\max_{p(\omega_1)} p(\text{error}) \rightarrow \min_{\Omega_1}$$

Minimax criterion

- $p(\text{error})$ is a linear function of $p(\omega_1)$, so maximum is achieved on the edge and is equal to

$$\max \left\{ \int_{\Omega_2} p(x|\omega_1) dx, \int_{\Omega_1} p(x|\omega_2) dx \right\}$$

- Solution:

$$\int_{\Omega_2} p(x|\omega_1) dx = \int_{\Omega_1} p(x|\omega_2) dx$$

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions
- 4 Discriminant functions**
- 5 Gaussian classifiers
- 6 Mixture models

Definition

Discriminant functions approach

- have C discriminant functions $g_i(x)$, $i = 1, 2 \dots C$.
- assign x to class having maximum discriminant function value:

$$\hat{c} = \arg \max_c g_c(x)$$

Discriminant functions are not unique:

$g_i(x)$ and $g'(x) = f(g(x))$ lead to equivalent classification for any monotonically increasing function $f(x)$.

Two class case

- For two class case we may define a single function $g(x)$ such that

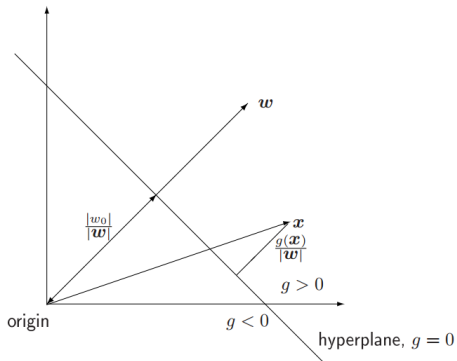
$$\hat{c} = \begin{cases} 1, & g(x) > k, \\ 2 & g(x) \leq k. \end{cases}$$

- Possible reductions from multiclass case:
 - $g(x) = g_1(x) - g_2(x)$, $k = 0$
 - $g(x) = g_1(x)/g_2(x)$, $k = 1$ for positive $g_1(x), g_2(x)$.

Linear discriminant function

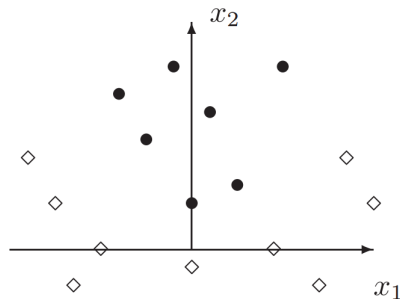
Simplest case - linear discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$



Limitation of linear decision boundary

The objects below can't be separated with linear boundary.



However, objects may be linearly separated in transformed space:

$$\phi_1(\mathbf{x}) = x_1^2, \phi_2(\mathbf{x}) = x_2.$$

Non-linear discriminant functions

Natural way to make non-linear decision boundaries is to apply standard linear discriminant functions with transformed features.

Most well-known examples:

- linear: $\phi_i(\mathbf{x}) = x_i$
- polynomial: $\phi_i(\mathbf{x}) = x_{k_1}^{s_1} x_{k_2}^{s_2} \dots x_{k_q}^{s_q}$
- radial basis functions: $\phi_i(\mathbf{x}) = \phi(|\mathbf{x} - \boldsymbol{\nu}_i|)$, where $\phi(\cdot)$ is non-increasing function, meaning proximity.
- multi-layer perceptron: $\phi_i(\mathbf{x}) = f(\mathbf{x}^T \boldsymbol{\nu}_i + \nu_{i0})$, where $f(z) = 1/(1 + e^{-z})$ - logistic or any other “step” function.

Probability calibration

- Assume we have an arbitrary discrimination function $g(x)$ and two classes.
- Classification is based on the sign: $y = \text{sign } g(x)$
- Need to estimate posterior class probabilities

$$p(\omega_1|x) = F(R_\theta(g(x)))$$

where $R_\theta(z)$ is monotone transform and F maps \mathbb{R} to $[0, 1]$.

- Can assume $F(z) = \sigma(z)$, $R_\theta(z) = \theta_0 + \theta_1 z$ (Platt's calibration)
- Then, using the property $1 - \sigma(z) = \sigma(-z)$:
$$p(y = 1|x) = \sigma(\theta_0 + \theta_1 g(x)), \quad p(y = -1|x) = \sigma(-\theta_0 - \theta_1 g(x))$$
- Estimate using ML:

$$\prod_{i=1}^n \sigma(y_i(\theta_0 + \theta_1 g(x))) \rightarrow \max_{\theta}$$

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions
- 4 Discriminant functions
- 5 Gaussian classifiers**
- 6 Mixture models

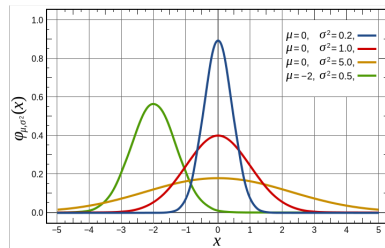
Gaussian (normal) distribution

Univariate case:

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

Multivariate case (D -dimensionality of data):

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



Gaussian (normal) distribution - sample estimates

Univariate case:

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})^2$$

Multivariate case (D -dimensionality of data):

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \hat{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})(x_n - \hat{\mu})^T$$

Since $\hat{\sigma}_{ML}^2$ and $\hat{\Sigma}_{ML}$ are biased estimates, it is common to use unbiased estimates:

$$\hat{\sigma}_{ML}^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu})^2, \quad \hat{\Sigma}_{ML} = \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu})(x_n - \hat{\mu})^T$$

Gaussian classifier

- In Gaussian classifier

$$p(x|\omega_j) = \frac{1}{(2\pi)^{D/2} |\Sigma_j|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma_j^{-1} (x - \mu) \right\}$$

- It follows that

$$\begin{aligned} \log p(x|\omega_j) &= \log p(x|\omega_j) + \log p(\omega_j) - \log p(x) \\ &= -\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) - \frac{1}{2} \log |\Sigma_j| \\ &\quad - \frac{d}{2} \log(2\pi) + \log p(\omega_j) - \log p(x) \end{aligned}$$

- Removing common additive terms, we obtain discriminant functions:

$$g_j(x) = \log p(\omega_j) - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) \quad (1)$$

Practical application

- In practice we replace theoretical terms μ_j , Σ_j with their sample estimates $\hat{\mu}_j$, $\hat{\Sigma}_j$.
- $p(\omega_j)$ may be estimated using sample frequency (assuming that experiment conditions don't change): $p(\omega_j) = \frac{n_j}{n}$.

$$g_j(x) = \log p(\omega_j) - \frac{1}{2} \log |\hat{\Sigma}_j| - \frac{1}{2} (x - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (x - \hat{\mu}_j)$$

- Analysis:
 - depends on normality assumptions (in particular - on unimodality)
 - needs to specify:
 - CD parameters to estimate $\hat{\mu}_j$, $j = 1, 2, \dots, C$.
 - $CD(D+1)/2$ parameters to estimate $\hat{\Sigma}_j$, $j = 1, 2, \dots, C$.

Simplifying assumptions

- $CD(D + 3)/2$ may be too large for multidimensional tasks with small training sets.
- Simplifying assumptions:
 - **Naive Bayes**: assume that $\Sigma_1, \Sigma_2, \dots, \Sigma_C$ are diagonal.
 - **Project onto a subspace**: for example on first few principal components.
 - **Proportional covariance matrices**: assume that $\Sigma_1 = \alpha_1 \Sigma, \Sigma_2 = \alpha_2 \Sigma, \dots, \Sigma_C = \alpha_C \Sigma$.
 - **Fisher's linear discriminant analysis**: assume that $\Sigma_1 = \Sigma_2 = \dots = \Sigma_C$.

Fisher's LDA

- Quadratic discriminant analysis classification was obtained above in (1):

$$g_j(x) = \log p(\omega_j) - \frac{1}{2} \log |\hat{\Sigma}_j| - \frac{1}{2} (x - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (x - \hat{\mu}_j)$$

- Under assumption $\Sigma_1 = \Sigma_2 = \dots = \Sigma_C$ we obtain:

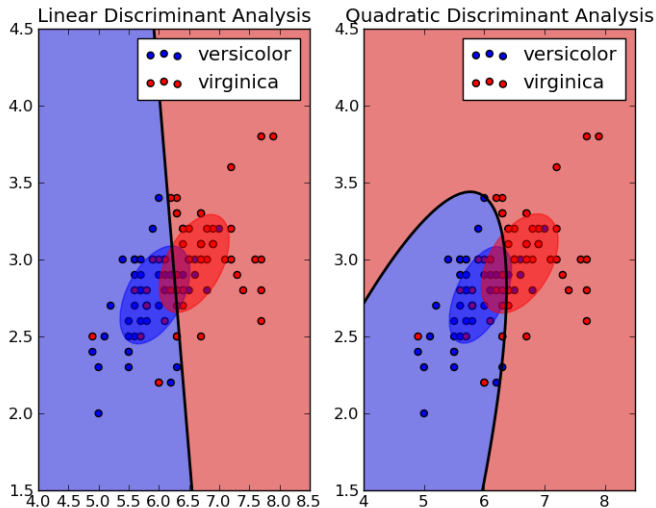
$$g_j(x) = \log p(\omega_j) - \frac{1}{2} \hat{\mu}_j^T S_W^{-1} \hat{\mu}_j + x^T S_W^{-1} \hat{\mu}_j$$

- where estimator of common covariance matrix is given by pooled within class covariance matrix

$$S_W = \sum_{j=1}^C \frac{N_j}{N} \hat{\Sigma}_j$$

- Unbiased estimate of common covariance matrix is $\frac{N}{N-C} S_W$.

LDA vs. QDA



Addition

- For LDA: if $p(\omega_1) = p(\omega_2) = \dots = p(\omega_C)$ and $\Sigma_1 = \Sigma_2 = \dots = \Sigma_C = I$, then LDA reduces to **nearest mean classifier**.
- Regularized discriminant analysis - intermediate case between different, common and identity covariance matrices:

$$\tilde{\Sigma}_j = \alpha \Sigma_j + \beta \Sigma + (1 - \alpha - \beta)I, \quad \alpha \geq 0, \beta \geq 0, \alpha + \beta \leq 1.$$

- α and β for RDA are selected using cross-validation.

Table of Contents

- 1 Bayes Decision Rule for Minimum Error
- 2 Bayes Decision Rule for Minimum Risk
- 3 Class prior independent decisions
- 4 Discriminant functions
- 5 Gaussian classifiers
- 6 Mixture models**

Mixture models

- Mixture models have the form:

$$p(x) = \sum_{j=1}^g \pi_j p(x; \theta_j)$$

- May be viewed as a two-step random process:

- 1 generate cluster number j from discrete distribution

$$j \sim \text{Disc}(\pi_1, \pi_2, \dots, \pi_g)$$

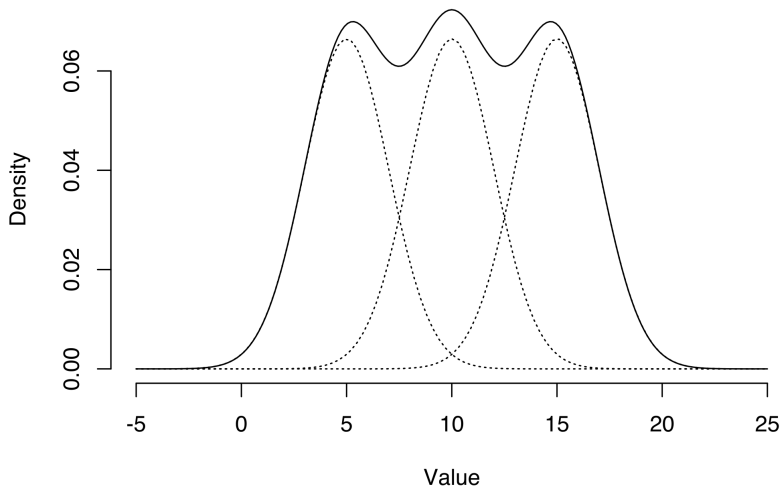
- 2 generate observation x using density, representing given cluster
 $p(x; \theta_j)$

$$x \sim p(x; \theta_j)$$

- Normal mixtures:

$$p(x) = \sum_{j=1}^g \pi_j N(x; \mu_j, \Sigma_j)$$

Normal mixture example



Likelihood maximization

- Direct optimization of likelihood is not possible:

$$L(\pi_1, \dots, \pi_g, \mu_1, \dots, \mu_g, \Sigma_1, \dots, \Sigma_g) = \prod_{i=1}^N \left\{ \sum_{j=1}^g \pi_j N(x_i; \mu_j, \Sigma_j) \right\}$$

- EM-algorithm is an iterative algorithm for likelihood maximization:

EM algorithm

- while parameters not converged:
 - for each sample $i = 1, 2, \dots, N$; for each cluster $j = 1, 2, \dots, g$:
 - recalculate cluster correspondences:

$$w_{ij} = \frac{\pi_j N(x_i; \mu_j, \Sigma_j)}{\sum_k \pi_k N(x_i; \mu_k, \Sigma_k)}$$

- recalculate parameters of each cluster:

$$\hat{\pi}_j = \frac{1}{N} \sum_{i=1}^N w_{ij}, \quad \hat{\mu}_j = \frac{1}{\sum_{i=1}^N w_{ij}} \sum_{i=1}^N w_{ij} x_i$$

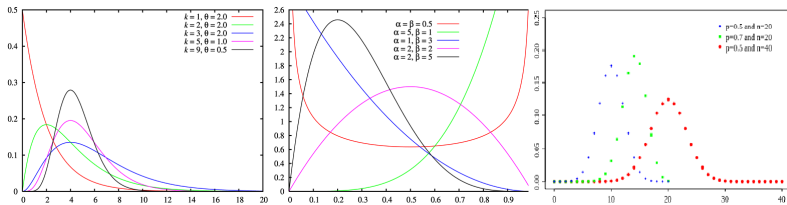
$$\hat{\Sigma}_j = \frac{1}{\sum_{i=1}^N w_{ij}} \sum_{i=1}^N w_{ij} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$$

Application

- May converge to local optimum, depending on initial condition
 - starting EM from different initial conditions is advised
- Initial conditions may be obtained from k-means algorithm
- Initial covariance should be wide enough
- Degenerate optimums: $\hat{\mu}_j = x_k, \hat{\Sigma}_j \rightarrow \Theta \in \mathbb{R}^{D \times D}$ (Θ is the zero matrix)
- $\hat{\Sigma}_j$ may become singular, so $N(x_i; \mu_j, \Sigma_j)$ will not be computable
 - $\tilde{\Sigma}_j \leftarrow \Sigma_j + \alpha I$
 - consider only diagonal Σ_j
 - consider only spherical $\Sigma_j = \alpha_j I$
 - assume common covariance matrix $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g$

Different distributions

- Other distributions may be taken instead of Gaussian:
 - **Gamma distribution** for non-negative random variables
 - **Beta distribution** for non-negative bounded random variables
 - **Polynomial distribution** for non-negative discrete random variables
 - etc.



Gamma, beta and binomial distributions