

Resumo de Aprendizagem de Máquina 2014-2

Eduardo M. B. de A. Tenório

embat@cin.ufpe.br

CIn-UFPE

Resumo

Este documento tem por finalidade ser um resumo dos assuntos abordados na disciplina Aprendizagem de Máquina do período 2014-2 do CIn-UFPE, ministrada pelos professores Francisco Carvalho e Teresa Ludermit. A maioria do documento referencia o livro “Pattern Classification”, de Duda, Hart & Stork. Os códigos utilizados como exercício de fixação encontram-se em github.com/embatbr/resumo-aprendizagem.

1 Teoria da Decisão Bayesiana

1.1 Introdução

Teoria da Decisão Bayesiana é uma abordagem estatística para a classificação de padrões, baseada em quantificar os tradeoffs associados a tomar uma determinada decisão (classificar) utilizando probabilidade e considerando os custos associados.

O **estado natural** é denotado por ω , de modo que $\omega = \omega_i$, para $i = 1, 2, \dots, c$, significa que o exemplo foi classificado como pertencente à classe ω_i . Cada uma dessas classes possui uma **probabilidade a priori** $P(\omega_i)$, com

$$\sum_{i=1}^c P(\omega_i) = 1, \quad (1)$$

refletindo o conhecimento prévio da chance de

um elemento da classe ω_i aparecer. A **regra de decisão** fica:

$$\text{Decida } \omega_i \text{ se } i = \arg \max_j P(\omega_j). \quad (2)$$

Neste caso a classe ω_i sempre é escolhida e a probabilidade de erro é dada por:

$$P_{err}(\omega_i) = 1 - P(\omega_i). \quad (3)$$

Utilizando uma característica x que seja contínua e aleatória, sua **densidade de probabilidade estado-condicional** é dada por $p(x|\omega)$. Logo, a diferença entre $p(x|\omega_i)$ e $p(x|\omega_j)$ descreve a diferença da característica x entre as populações das classes ω_i e ω_j .

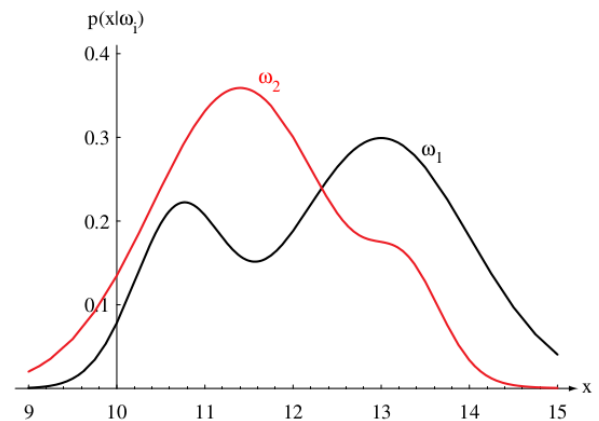


Figura 1: Para $\omega_i = \omega_2$, é mais frequente observar x entre 11 e 12 que $x = 13$ (valor mais provável se $\omega_i = \omega_1$).

Sabendo $P(\omega_i)$ e $p(x|\omega_i)$, e medindo um valor x , a probabilidade conjunta de achar

um padrão na classe ω_i e com x é dado por:
 $p(\omega_i, x) = P(\omega_i|x)p(x) = p(x|\omega_i)P(\omega_i)$, que
 pela **fórmula de Bayes** fica:

$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)}, \quad (4)$$

com a evidência para c classes

$$p(x) = \sum_{j=1}^c p(x|\omega_j)P(\omega_j). \quad (5)$$

A probabilidade a posteriori das classes ω_1 e ω_2 para um conjunto de valores de x é mostrada em Fig. (2). A regra de decisão fica:

$$\text{Decida } \omega_i \text{ se } i = \arg \min_j P(\text{erro}_j|x), \quad (6)$$

onde

$$P(\text{erro}_i|x) = \sum_{j \neq i} P(\omega_j|x), \quad (7)$$

ou simplesmente $P(\text{erro}_i|x) = 1 - P(\omega_i|x)$.
 Então a regra torna-se:

$$\text{Decida } \omega_i \text{ se } i = \arg \max_j P(\omega_j|x), \quad (8)$$

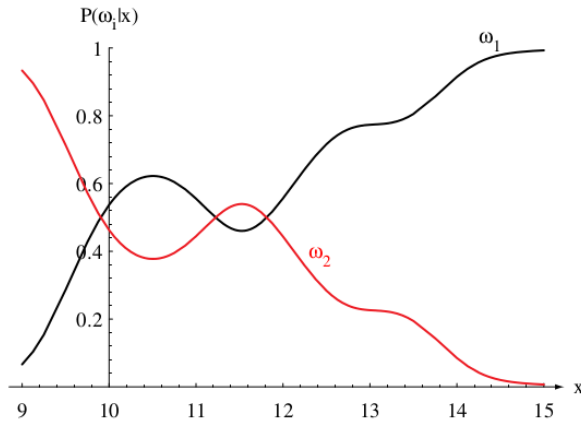


Figura 2: Probabilidades a posteriori para $P(\omega_1) = \frac{2}{3}$ e $P(\omega_2) = \frac{1}{3}$, e para as densidades de probabilidade estado-condicional mostradas em Fig. (1).

Esta regra minimiza a probabilidade média de erro, dada por

$$P(\text{erro}_i) = \int_{-\infty}^{\infty} P(\text{erro}_i|x)p(x)dx. \quad (9)$$

1.2 Características Contínuas

É de fácil compreensão que a característica x pode ser trocada por um vetor de características $\mathbf{x} = (x_1, x_2, \dots, x_d)$, onde \mathbf{x} pertence ao espaço \mathbf{R}^d (espaço de características). A região que decide ω_i é denotada por \mathcal{R}_i .

Outras ações além de apenas classificar um elemento podem ser tomadas, como por exemplo a **rejeição**: recusar-se a tomar uma decisão; uma opção válida quando o custo de ser indeciso é aceitável. Para isso **funções de custo** são inseridas, permitindo tratar de situações onde alguns erros de classificação são mais importantes que outros.

Seja $\{\omega_1, \dots, \omega_c\}$ o conjunto finito de c classes e seja $\{\alpha_1, \dots, \alpha_a\}$ o conjunto finito de possíveis ações. A função de custo $\lambda(\alpha_i|\omega_j)$ descreve o custo de tomar a ação α_i quando a classe é ω_j . Logo, observado um \mathbf{x} em particular, tomar a ação α_i quando a classe é ω_j leva a um custo esperado (**risco**)

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|\mathbf{x}). \quad (10)$$

$R(\alpha_i|\mathbf{x})$ é chamado de **risco condicional**. Qualquer que seja o \mathbf{x} observado, o risco pode ser minimizado selecionando a ação que minimiza $R(\alpha_i|\mathbf{x})$.

A regra de decisão geral é uma função $\alpha(\mathbf{x})$ que diz qual ação tomar para cada possível observação, ou seja, para cada \mathbf{x} a **função de decisão** $\alpha(\mathbf{x})$ assume um dos a valores $\alpha_1, \dots, \alpha_a$.

$$\alpha(\mathbf{x}) = \alpha_i \text{ se } i = \arg \min_j R(\alpha_j|\mathbf{x}), \quad (11)$$

Logo, o **risco global** é dado por

$$R = \int R(\alpha(\mathbf{x}))p(\mathbf{x})d\mathbf{x}. \quad (12)$$

O risco global mínimo é chamado de **risco de Bayes**, denotado por R^* , sendo a melhor performance alcançável.

1.3 Taxa de Erro Mínima

Para evitar erros, a regra de decisão procurada é aquela que minimiza a probabilidade de erro, i.e. minimiza a **taxa de erro**. A função de custo de interesse para este caso é chamada de **simétrica** ou **zero-um**,

$$\lambda(\alpha_i|\omega_j) = \begin{cases} 0 & \text{se } i = j \\ 1 & \text{se } i \neq j \end{cases} \quad (13)$$

para $i, j = 1, \dots, c$. Como todos os erros tem custo igual, o risco condicional é dado por

$$R(\alpha_i|\mathbf{x}) = 1 - P(\omega_i|\mathbf{x}) \quad (14)$$

com $P(\omega_i|\mathbf{x})$ sendo a probabilidade condicional da ação α_i estar correta. A regra de decisão neste caso continua:

$$\text{Decida } \omega_i \text{ se } i = \arg \max_j P(\omega_j|\mathbf{x}). \quad (15)$$

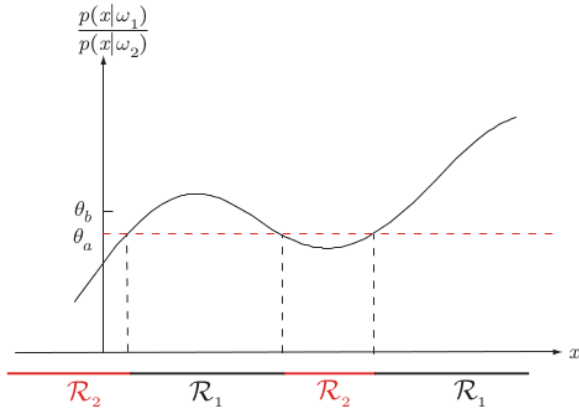


Figura 3: Se a penalização de classificar ω_1 como ω_2 for maior que o oposto, então a razão tende ao threshold θ_b .

1.4 Funções Discriminantes

A maneira mais usual de representar classificadores de padrões é através de um conjunto de **funções discriminantes** $g_i(\mathbf{x})$, $i = 1, \dots, c$. O classificador atribui um vetor de características \mathbf{x} à classe ω_i se

$$g_i(\mathbf{x}) = \arg \max_j g_j(\mathbf{x}) \quad (16)$$

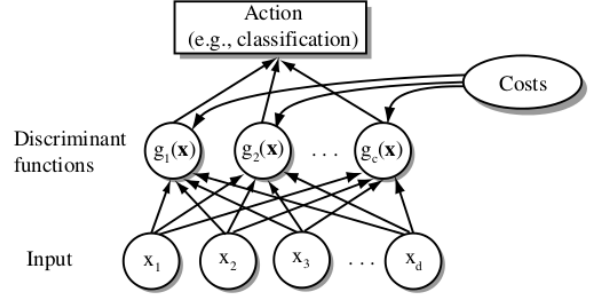


Figura 4: Classificador com c funções discriminantes e entradas d -dimensional. A ação geralmente é “escolher o maior $g_i(\mathbf{x})$ ”.

Para o caso geral com riscos, pode-se fazer $g_i(\mathbf{x}) = -R(\alpha_i|\mathbf{x})$, enquanto para o caso “taxa de erro mínima”, $g_i(\mathbf{x}) = P(\omega_i|\mathbf{x})$. A função discriminante $g_i(\mathbf{x})$ pode ser substituída por $f(g_i(\mathbf{x}))$, com $f(\cdot)$ sendo uma função monotonicamente crescente (e.g. logaritmo), com o resultado da classificação ficando inalterado. Como resultado, \mathbf{R}^d é dividido em regiões de decisão \mathcal{R}_i (não necessariamente conectadas) para cada classe ω_i .

Para o caso em que $c = 2$, o classificador é chamado **dicotomizador**, e apenas uma função discriminante $g(\mathbf{x}) \equiv g_1(\mathbf{x}) - g_2(\mathbf{x})$ é necessária. Logo a regra de decisão torna-se:

$$\text{Decida } \omega_1 \text{ se } g(\mathbf{x}) > 0; \text{ senão, } \omega_2. \quad (17)$$

1.5 Características Discretas

Em muitas aplicações práticas as componentes de \mathbf{x} são valores inteiros binários, ternários ou outro “ário”, de modo que \mathbf{x} pode assumir apenas um dos m valores discretos $\mathbf{v}_1, \dots, \mathbf{v}_m$. Nestes casos, a função de densidade de probabilidade $p(\mathbf{x}|\omega_j)$ torna-se uma função de massa de probabilidade $P(\mathbf{x}|\omega_j)$ e

$$\int p(\mathbf{x}|\omega_j) d\mathbf{x} \quad (18)$$

é substituída por

$$\sum_{\mathbf{x}} P(\mathbf{x}|\omega_j). \quad (19)$$

Na fórmula de Bayes as densidades de probabilidade são trocadas por probabilidades

$$P(\omega_j|\mathbf{x}) = \frac{P(\mathbf{x}|\omega_j)P(\omega_j)}{P(\mathbf{x})} \quad (20)$$

onde

$$P(\mathbf{x}) = \sum_{j=1}^c P(\mathbf{x}|\omega_j)P(\omega_j). \quad (21)$$

A definição do risco condicional $R(\alpha|\mathbf{x})$ mantém-se inalterada.

2 Estimação Paramétrica

2.1 Introdução

Sabendo $P(\omega_i)$ e $p(\mathbf{x}|\omega_i)$ é possível projetar um classificador ótimo. Contudo, em aplicações reais de classificação de padrões raramente tem-se este conhecimento completo a respeito da estrutura probabilística do problema. Geralmente o que está disponível é um conhecimento superficial da situação e um conjunto de **dados de treinamento**.

Como exemplo, sejam 300 objetos divididos em 2 classes de modo que a classe 1 tenha 200 objetos e a classe 2 tenha 100. Cada objeto é gerado a partir de uma das 3 gaussianas bi-variadas:

1a: $\mu_1 = 60, \mu_2 = 30, \sigma_1^2 = 9$ e $\sigma_2^2 = 144$

1b: $\mu_1 = 52, \mu_2 = 30, \sigma_1^2 = 9$ e $\sigma_2^2 = 9$

2: $\mu_1 = 45, \mu_2 = 22, \sigma_1^2 = 100$ e $\sigma_2^2 = 9$

Nota-se que Σ_{1a} , Σ_{1b} e Σ_2 são matrizes diagonais, reduzindo-as a vetores de variância, o que significa independência entre as variáveis.

As duas imagens em Fig. (5) mostram os dados de treinamento. Na de cima os dados são divididos em clusters gerados pelas gaussianas definidas anteriormente. Na segunda são divididos em duas classes.

Pode-se utilizar estas amostras para estimar as probabilidades e densidades de probabilidade desconhecidas. Para problemas

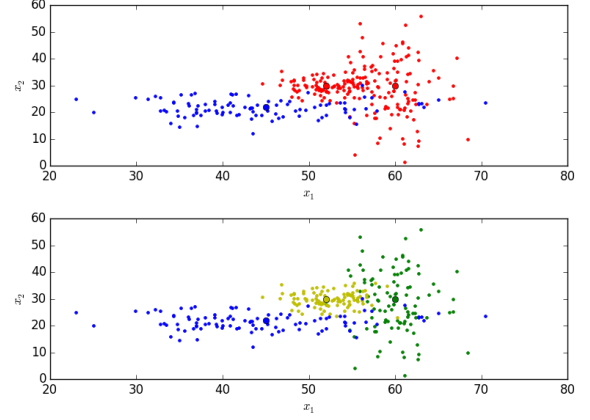


Figura 5: Classes 1a (verde), 1b (amarelo) e 2 (azul).

de aprendizagem supervisionada, a estimação das $P(\omega_i)$ é bastante simples. Sabendo a frequência relativa de aparição de cada classe ω_i no conjunto de treinamento, pode-se estimar suas $P(\omega_i)$ (desde que o conjunto de treinamento seja uma boa representação do “mundo real”). Contudo, estimar as $p(\mathbf{x}|\omega_i)$ é complicado quando se tem poucos dados, especialmente quando a dimensionalidade de \mathbf{x} aumenta. Uma maneira mais prática é assumir cada $p(\mathbf{x}|\omega_i)$ como uma gaussiana, cujas média e matriz de covariância são dadas por μ_i e Σ_i . Logo, o problema passa de estimar uma função desconhecida $p(\mathbf{x}|\omega_i)$ para estimar os **parâmetros** μ_i e Σ_i .

Para a estimação dos parâmetros, será utilizada a técnica da **máxima-verossimilhança**.

2.2 Máxima-Verossimilhança

Estimação por Máxima-Verossimilhança (Maximum-Likelihood, ML) considera os parâmetros como quantidades cujos valores são fixos, mas desconhecidos. A melhor estimativa de seus valores é aquela que maximiza a probabilidade de obter as amostras observadas. Este método possui uma gama de atributos atrativos. Primeiramente, possui quase sempre boas propriedades de convergência quando o número de amostras de treinamento aumenta. Além disso, ML

geralmente é mais simples que métodos alternativos, tais como técnicas bayesianas.

Princípio Geral

Um conjunto \mathcal{D} de amostras é dividido em conjuntos $\mathcal{D}_1, \dots, \mathcal{D}_c$ cujas amostras são desenhadas independentemente, de acordo com $p(\mathbf{x}|\omega_j)$, resultado de variáveis aleatórias independentes e identicamente distribuídas. Assume-se que $p(\mathbf{x}|\omega_j)$ tem uma forma paramétrica, como por exemplo $N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, de modo que pode ser escrita como $p(\mathbf{x}|\omega_j, \boldsymbol{\theta}_j)$, onde $\boldsymbol{\theta}_j$ consiste nos componentes de $\boldsymbol{\mu}_j$ e $\boldsymbol{\Sigma}_j$. O problema torna-se utilizar as informações dadas pelas amostras para obter boas estimativas dos vetores desconhecidos $\boldsymbol{\theta}_j$, para cada ω_j . Uma simplificação é assumir que as amostras de cada \mathcal{D}_j provêm informações relevantes apenas aos seus respectivos $\boldsymbol{\theta}_j$.

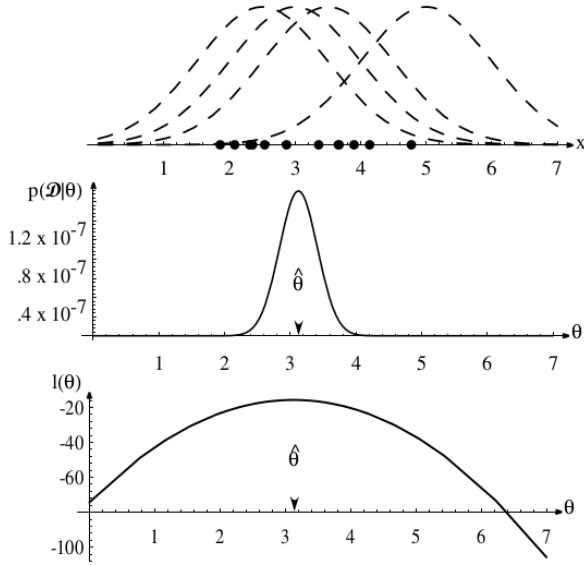


Figura 6: Estimação do parâmetro $\hat{\theta}$, o valor $\hat{\theta}$ achado e $l(\hat{\theta})$.

Usa-se o conjunto \mathcal{D} , composto pelas amostras de treinamento $\mathbf{x}_1, \dots, \mathbf{x}_n$ desenhadas independentemente de $p(\mathbf{x}|\boldsymbol{\theta})$, para estimar os $\boldsymbol{\theta}$ desconhecidos, tal que

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x}_k|\boldsymbol{\theta}). \quad (22)$$

A densidade $p(\mathcal{D}|\boldsymbol{\theta})$ é conhecida como a **verossimilhança** de $\boldsymbol{\theta}$ em relação às amostras, e sua estimativa máxima é, por definição, o valor $\hat{\boldsymbol{\theta}}$ que maximiza $p(\mathcal{D}|\boldsymbol{\theta})$ (Fig. (6)).

Devido a ser monotonicamente crescente e mais fácil de trabalhar que a verossimilhança, o uso da log-verossimilhança

$$l(\boldsymbol{\theta}) \equiv \ln p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{k=1}^n \ln p(\mathbf{x}_k|\boldsymbol{\theta}) \quad (23)$$

é preferível. Logo, a solução para $\hat{\boldsymbol{\theta}}$ é

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}), \quad (24)$$

onde a dependencia do conjunto de dados \mathcal{D} fica implícito. Se o número de parametros a serem estimados é p , então $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^t$ e $\nabla_{\boldsymbol{\theta}} = (\frac{\partial}{\partial \theta_1}, \dots, \frac{\partial}{\partial \theta_p})^t$, e, pela Eq. (23)

$$\nabla_{\boldsymbol{\theta}} l = \sum_{k=1}^n \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}_k|\boldsymbol{\theta}). \quad (25)$$

Então, um conjunto de condições necessárias para a estimação da máxima verossimilhança para $\boldsymbol{\theta}$ pode ser obtido do conjunto de p equações

$$\nabla_{\boldsymbol{\theta}} l = 0. \quad (26)$$

Uma solução $\hat{\boldsymbol{\theta}}$ para Eq. (26) pode ser uma máximo global, um máximo/mínimo local, ou (raramente) um ponto de inflexão de $l(\boldsymbol{\theta})$. Se todas as soluções forem achadas, é garantido que representa o máximo verdadeiro, senão deve-se checar todas as soluções individualmente (ou calcular derivadas de segunda ordem) para identificar qual é o ótimo global.

Gaussiana com μ desconhecido

Neste caso o vetor $\theta = \mu$, levando a $l(\theta) = l(\mu) \equiv \ln p(\mathcal{D}|\mu)$. Logo pela Eq. (23),

$$l(\mu) = \sum_{k=1}^n \ln p(x_k|\mu), \quad (27)$$

e pela Eq. (25),

$$\nabla_{\mu} l = \sum_{k=1}^n \nabla_{\mu} \ln p(x_k|\mu). \quad (28)$$

Como $\nabla_{\mu} \ln p(x_k|\mu) = \Sigma^{-1}(x_k - \mu)$ e $\nabla_{\mu} l = 0$ (Eq. (26)), então

$$\sum_{k=1}^n \Sigma^{-1}(x_k - \hat{\mu}) = 0. \quad (29)$$

Multiplicando por Σ e rearranjando os termos, tem-se

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k. \quad (30)$$

Fica evidente que para o caso em que apenas μ é desconhecido, a estimativa para a máxima verossimilhança é apenas a média das amostras, às vezes escrita como $\hat{\mu}_n$ para clarificar sua dependência do número de amostras.

Gaussiana com μ e Σ desconhecidos

Este é o caso mais típico, quando μ e Σ são desconhecidos. Desta forma, estes parâmetros constituem as componentes do vetor paramétrico θ , ou seja $\theta = (\mu, \Sigma)^t$. Considerando o caso univariado, $\theta_1 = \mu$ e $\theta_2 = \sigma^2$. Então

$$\ln p(x_k|\theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2}(x_k - \theta_1)^2 \quad (31)$$

e seu gradiente é

$$\nabla_{\theta} l = \nabla_{\theta} \ln p(x_k|\theta) = \begin{bmatrix} \frac{1}{\theta_2}(x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix} \quad (32)$$

e pela Eq. (29) tem-se

$$\sum_{k=1}^n \frac{1}{\hat{\theta}_2}(x_k - \hat{\theta}_1) = 0 \quad (33)$$

e

$$-\sum_{k=1}^n \frac{1}{2\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{2\hat{\theta}_2^2} = 0, \quad (34)$$

Rearrmando tudo, tem-se

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k \quad (35)$$

e

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2. \quad (36)$$

Dando um salto de fé (embora esta seja demonstrável, diferente da cristã) e trocando os $\hat{\mu}$ e $\hat{\sigma}$ por $\hat{\mu}$ e $\hat{\Sigma}$, respectivamente, tem-se

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k \quad (37)$$

e

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^t. \quad (38)$$

3 Misturas

3.1 Introdução

Esta seção trata sobre aprendizado **não-supervisionado**, no qual as amostras não são rotuladas com suas classes. Existem pelo menos cinco razões básicas para utilizá-lo: (1) coletar e rotular um grande conjunto de dados

pode ser extremamente custoso; (2) o procedimento pode ser efetuado na ordem inversa, treinando com um grande número de amostras não rotuladas e então utilizar supervisão para rotular os grupos criados; (3) em muitas aplicações as características dos padrões podem mudar com o tempo; (4) métodos não-supervisionados podem ser utilizados para encontrar características que serão úteis para a categorização; (5) em estágio iniciais de uma investigação, pode ser interessante efetuar análise exploratória de dados para ganhar algum *insight* sobre a natureza ou estrutura dos dados.