

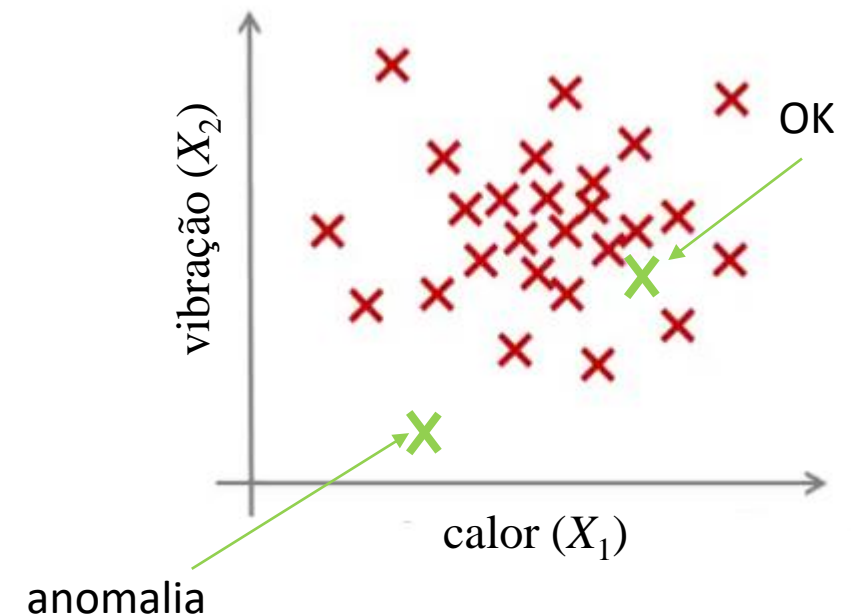
Avaliação de Desempenho

Aplicação da Distribuição Normal Multivariada

Detecção de Outliers

Detecção de outliers

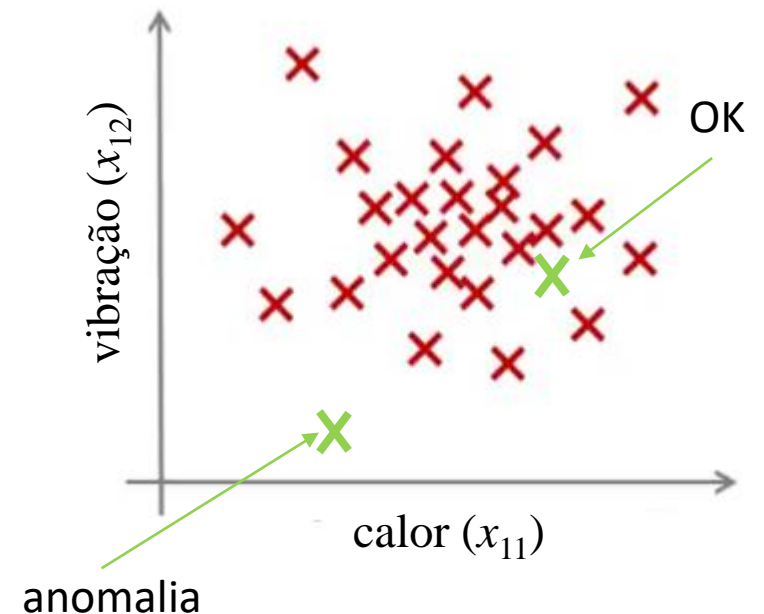
- Classificação ocorrências como comportamento fora do comum
- Funcionamento anormal de um motor
 - Características
 - aquecimento, intensidade de vibração
 - Observações $X = [X_1; X_2]$
 - Novo motor: $X_{\text{teste}} = [x_1; x_2]$



Detecção de outliers

- Dados de treinamento $X = [X_1, X_2, \dots, X_n]$
- Cada característica é X_i é uma coluna de X com m valores $X_j = [x_{1,j}; x_{2,j}; \dots; x_{m,j}]$
- Cada exemplo é a uma linha de X com n componentes $E_i = [x_{i,1}, x_{i,2}, \dots, x_{i,n}]$
- No exemplo $E_i = [\text{calor}, \text{vibração}]$

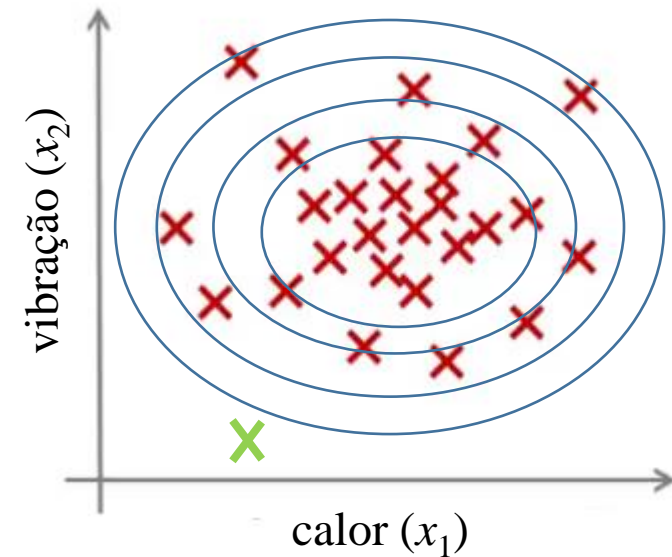
	X_1	X_2		X_n
E_1	x_{11}	x_{12}	...	x_{1n}
E_2	x_{21}	x_{22}	...	x_{2n}
E_m	x_{m1}	x_{m2}	...	x_{mn}



Detecção de outliers

- Formulação

- Dado um elemento E_{teste}
- E_{teste} é uma anomalia?
- Determinado pela probabilidade de ocorrência
- Se $f(E_{teste}) < \varepsilon$ então E_{teste} é uma anomalia



Detecção de outliers

- Outros exemplos
 - Detecção de fraude
 - X_1 quantidade de vezes que faz o login
 - X_2 quantas páginas visita
 - X_3 quantidade de posts em fóruns
 - ...
 - Monitoração de computadores em um data center
 - X_1 uso de memória
 - X_2 carga da CPU
 - X_3 tráfego de rede
 - ...

Detecção de outliers

Quando usar

Detecção de outliers (não supervisionado)

- Número pequeno de exemplos positivos
- Muitos tipos diferentes de anomalia
- Novas anomalias que não são semelhantes aos exemplos conhecidos

Aprendizado supervisionado

- Número grande de exemplos positivos e negativos
- Anomalias não vistas serão semelhantes àquelas encontradas no conjunto de treinamento
- Muitos exemplos rotulados

Detecção de outliers

Aplicações

Detecção de outliers (não supervisionado)

- Detecção de fraude
- Processos de fabricação
- Monitoração de máquinas em data center
- ...

Aprendizado supervisionado

- Classificação de e-mails (spam)
- Previsão do tempo
- Classificação de tumor
- ...

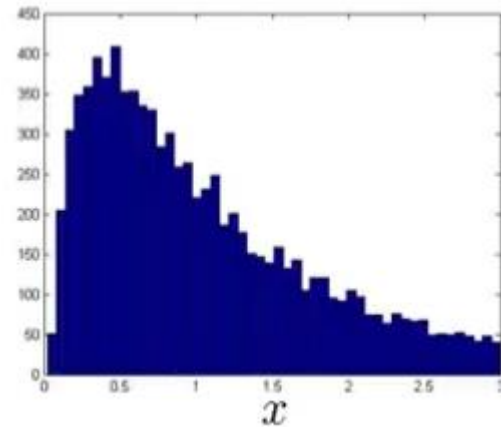
Dificuldades com o modelo

- Características não Gaussianas

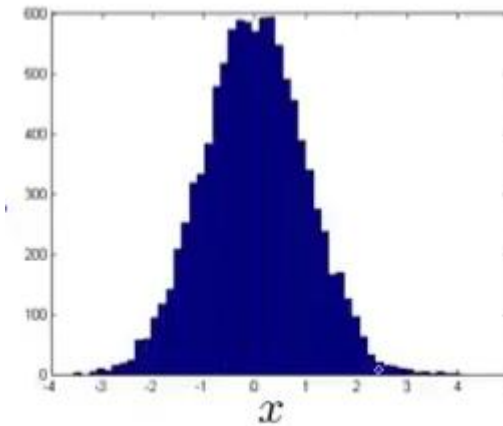
$$x_1 \leftarrow \log(x_1)$$

$$x_2 \leftarrow \log(x_2 + c)$$

$$x_3 \leftarrow x_3^{1/k}$$

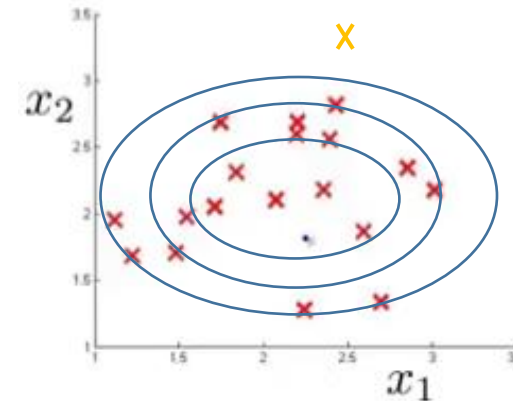
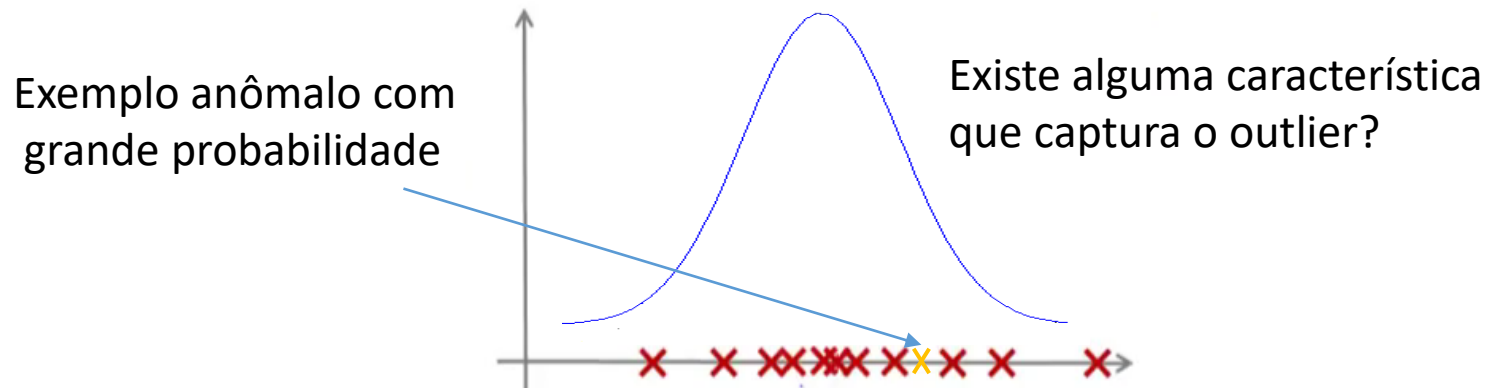


$\log(x)$

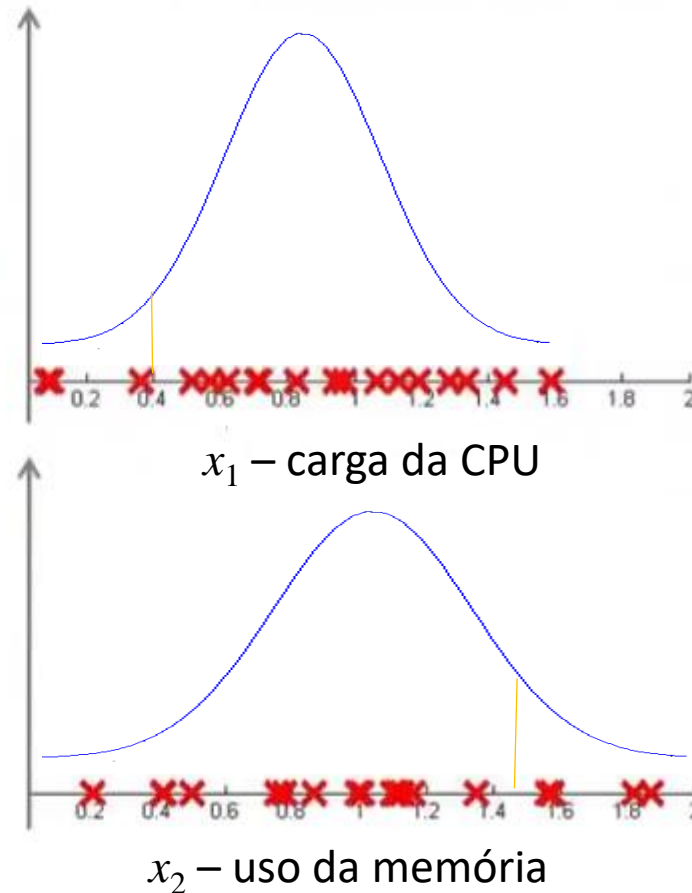
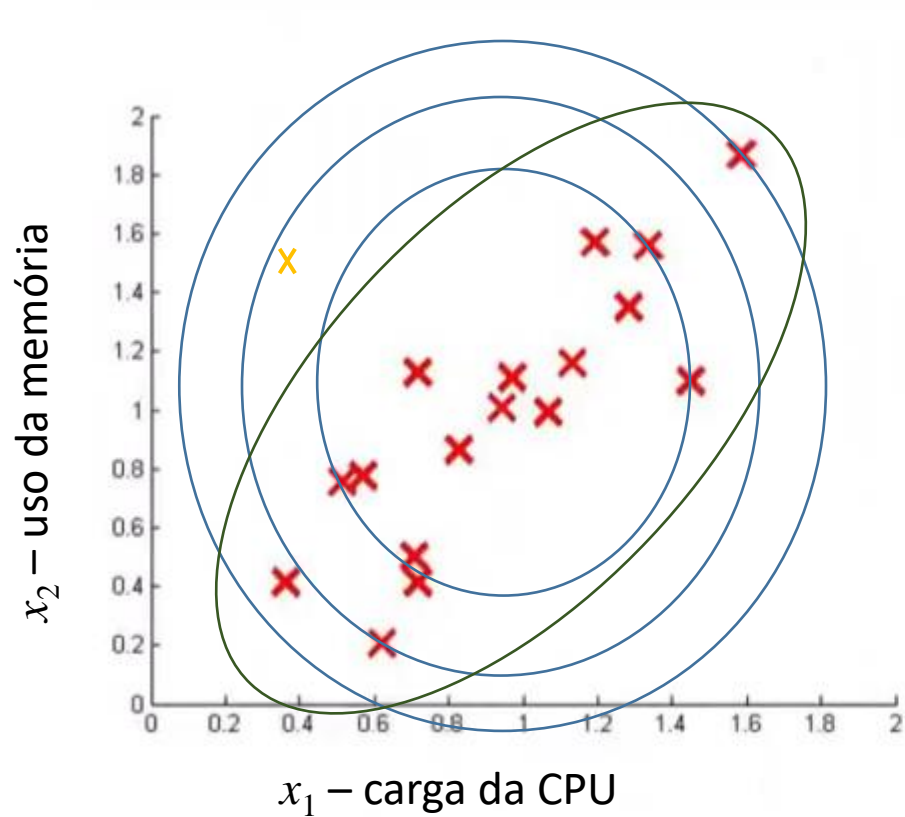


Dificuldades com o modelo

- Objetivo
 - $f(x)$ é grande para exemplos normais
 - $f(x)$ é pequeno para exemplos com anomalia
- Problema comum
 - $f(x)$ é semelhante para exemplos normais e anomalias



Distribuição Gaussiana Multivariada



- Usar a distribuição Gaussiana multivariada
- Captura possível correlação entre características

Distribuição Gaussiana Multivariada

- Um argumento (\mathbf{X}) e dois parâmetros ($\boldsymbol{\mu}$ e $\boldsymbol{\Sigma}$)

$$f(\mathbf{X}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu})\right)$$

$$\mathbf{X} = [\mathbf{X}_1 \quad \dots \quad \mathbf{X}_n] \quad \boldsymbol{\mu} = \begin{bmatrix} E[\mathbf{X}_1] \\ E[\mathbf{X}_2] \\ \dots \\ E[\mathbf{X}_n] \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} COV[\mathbf{X}_1, \mathbf{X}_1] & COV[\mathbf{X}_1, \mathbf{X}_2] & \dots & COV[\mathbf{X}_1, \mathbf{X}_n] \\ COV[\mathbf{X}_2, \mathbf{X}_1] & COV[\mathbf{X}_2, \mathbf{X}_2] & \dots & COV[\mathbf{X}_2, \mathbf{X}_n] \\ \vdots & \vdots & \ddots & \vdots \\ COV[\mathbf{X}_n, \mathbf{X}_1] & COV[\mathbf{X}_n, \mathbf{X}_2] & \dots & COV[\mathbf{X}_n, \mathbf{X}_n] \end{bmatrix}$$

- Estimativa dos parâmetros

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{X}^{(i)}$$

$$\boldsymbol{\Sigma} = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}^{(i)} - \boldsymbol{\mu}) \cdot (\mathbf{X}^{(i)} - \boldsymbol{\mu})^T$$

No MatLab

`mu = mean(X)`

`sigma = cov(X)`

Algoritmo

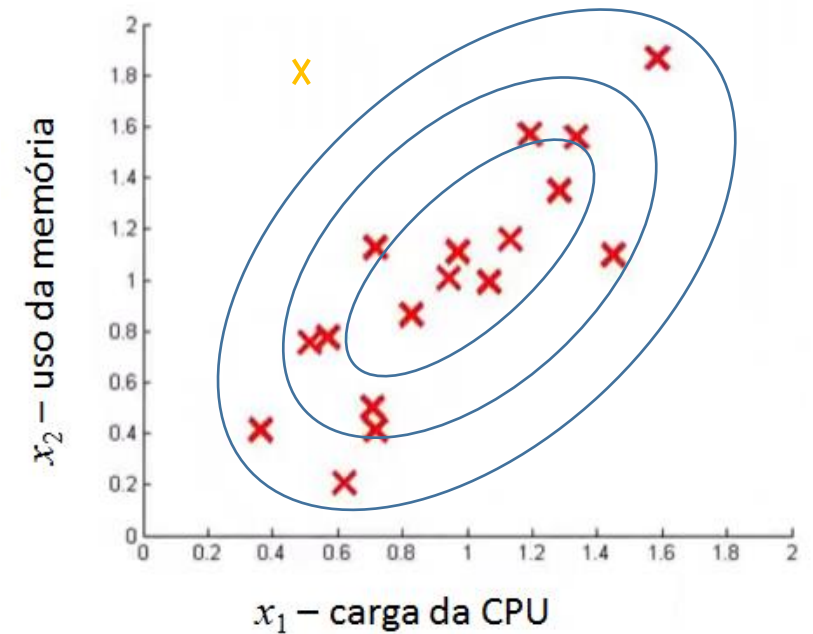
1. Ajustar o modelo $f(\mathbf{X})$ fazendo

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m X_i \quad \boldsymbol{\Sigma} = \frac{1}{m} \sum_{i=1}^m (\mathbf{X} - \boldsymbol{\mu}) \cdot (\mathbf{X} - \boldsymbol{\mu})^T$$

2. Dado um novo exemplo \mathbf{X} calcular

$$f(\mathbf{X}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu})\right)$$

3. Anomalia se $f(\mathbf{X}) < \varepsilon$ não anomalia caso contrário



`mvnpdf(X, mu, sigma)`

Métricas de avaliação

- O foco deve estar na capacidade preditiva do modelo e não no tempo que leva para classificar ou criar um modelo, na escalabilidade, etc.

		Categoria prevista	
		Sim	Não
Categoria real	Sim	Verdadeiro positivo	Falso negativo
	Não	Falso positivo	Verdadeiro negativo

Métricas de avaliação

- Precisão

$$prec = \frac{vp}{vp + fp}$$

- Revocação (recall)

$$rec = \frac{vp}{vp + fn}$$

- F1

$$F_1 = \frac{2 \cdot prec \cdot rec}{prec + rec}$$

Ajuste e avaliação do modelo

- Assume-se que temos exemplos que sabemos ser normais (maior parte) e exemplos que são anômalos (poucos)
- Conjunto de treinamento (não rotulado): precisa ter exemplos não anômalos (mais de 80%) e exemplos anômalos (menos de 80%)
- $X_{trei} = [X_1, X_2, \dots, X_n]$
- Conjunto de validação $X_{val} = [X_{val_1}, X_{val_2}, \dots, X_{val_n}, Y_{val}]$
- Conjunto teste $X_{test} = [X_{test_1}, X_{test_2}, \dots, X_{test_n}, Y_{test}]$

Ajuste e avaliação do modelo

- Ajustar o modelo usando o conjunto de validação
- Usar o conjunto de validação para selecionar ε (otimizar a métrica F1)

$$y = \begin{cases} 1 & \text{se } f(x) < \varepsilon \\ 2 & \text{se } f(x) \geq \varepsilon \end{cases}$$

- Realizar previsões no conjunto de teste

Avaliação de Desempenho

Aplicação da Distribuição Normal Multivariada

Detecção de Outliers