

Introdução a machine learning com aplicações em R

Aula 2 - Pré-processamento dos dados

Felipe Barletta

Departamento de Estatística

03 novembro, 2020



Sumário

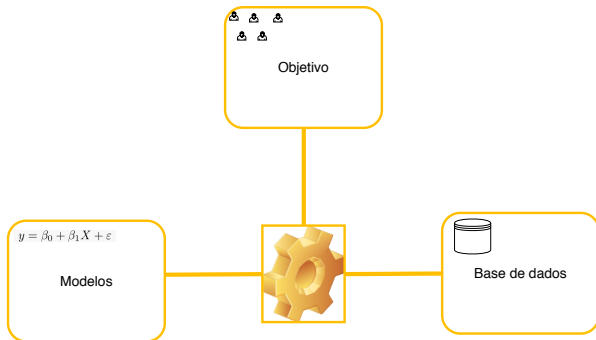
- 1 Revisão

Revisão - *Machine Learning*

- Intersecção de métodos estatísticos, Ciência da Computação (por meio de dados)
- É a tarefa de aprender com os dados e fazer generalizações



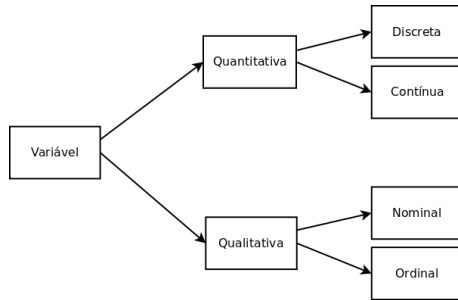
Revisão - Os três pilares



Revisão - Termos comuns

| <i>Machine Learning</i> | Estatística | O que é? |
|----------------------------------|-------------------------|----------------------------------|
| <i>Label</i> | Variável resposta | Variável a ser predita - (Y) |
| <i>Feature ou característica</i> | Variável preditora | Variável que pode prever (X) |
| <i>output</i> | Resposta | Predição do modelo |
| <i>input</i> | Variável preditora | Variável que está no modelo |
| Variável binária | Variável dicotômica | Variável com duas classes |
| Treinamento | Ajuste do modelo | Estrutura matemática |
| <i>Missing values</i> | <i>Missing values</i> | Valores faltantes |
| <i>Cross validation</i> | <i>Cross validation</i> | Validação cruzada |

Tipos de características ou variáveis



- **Nominal:** São categorias que não podem ser ordenadas
- **Ordinal:** São categorias que seguem uma ordem natural
- **Discreta:** Seu suporte é um conjunto finito ou infinito enumerável
- **Continua:** Seu suporte é um conjunto infinito não-enumerável

Pré-processamento

- **Formatar os dados.**
- Retirar amostras
- Remover ou corrigir dados faltantes
- Agregar ou criar novas características(variáveis)
- Verificar associação entre características(variáveis)

Pré-processamento

- Formatar os dados.
- Retirar amostras
- Remover ou corrigir dados faltantes
- Agregar ou criar novas características(variáveis)
- Verificar associação entre características(variáveis)

Pré-processamento

- Formatar os dados.
- Retirar amostras
- Remover ou corrigir dados faltantes
- Agregar ou criar novas características(variáveis)
- Verificar associação entre características(variáveis)

Pré-processamento

- Formatar os dados.
- Retirar amostras
- Remover ou corrigir dados faltantes
- Agregar ou criar novas características(variáveis)
- Verificar associação entre características(variáveis)

Pré-processamento

- Formatar os dados.
- Retirar amostras
- Remover ou corrigir dados faltantes
- Agregar ou criar novas características(variáveis)
- Verificar associação entre características(variáveis)

| Desfecho | Quantidade |
|--------------------------|------------|
| Ainda Internado | 246800 |
| Alta Médica | 4000 |
| Alta Melhorado | 502 |
| Alta | 70 |
| Atendendimento Cancelado | 203 |
| Evadiu-se | 5 |
| Internação Cancelada | 6 |
| Internamento Cancelado | 14 |
| Óbito | 400 |
| Óbito - Pós-Operatório | 29 |

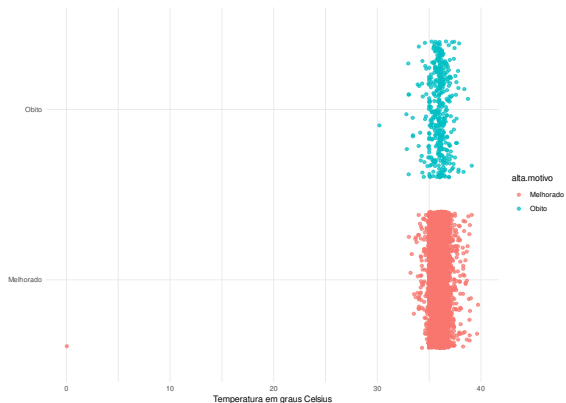
| Desfecho | Quantidade |
|---------------|------------|
| Internado | 246800 |
| Alta | 4000 |
| Alta | 502 |
| Alta | 70 |
| Não Internado | 203 |
| Não Internado | 5 |
| Não Internado | 6 |
| Não Internado | 14 |
| Óbito | 400 |
| Óbito | 29 |

| Desfecho | Quantidade |
|---------------|------------|
| Internado | 246800 |
| Alta | 4572 |
| Não Internado | 228 |
| Óbito | 429 |

| Desfecho | Quantidade |
|-----------|------------|
| Melhorado | 4572 |
| Óbito | 429 |

○○○○○

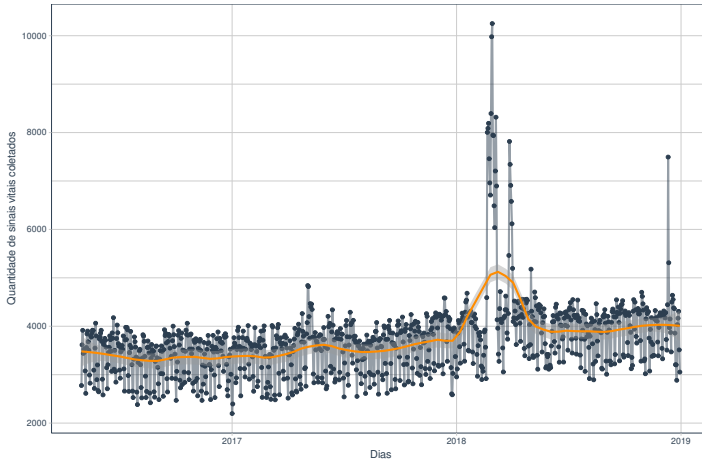




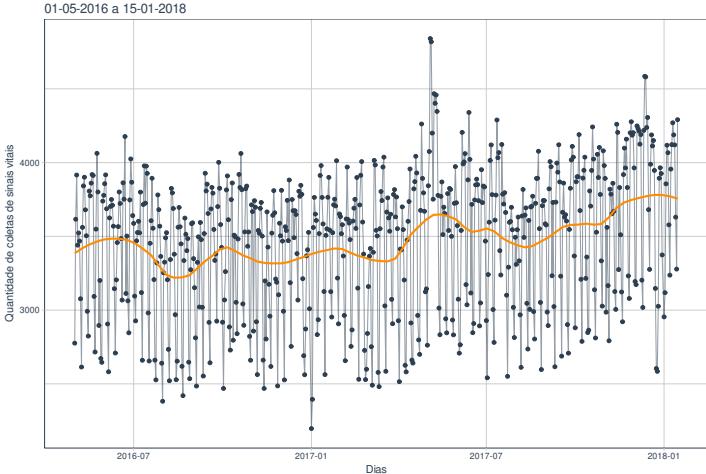


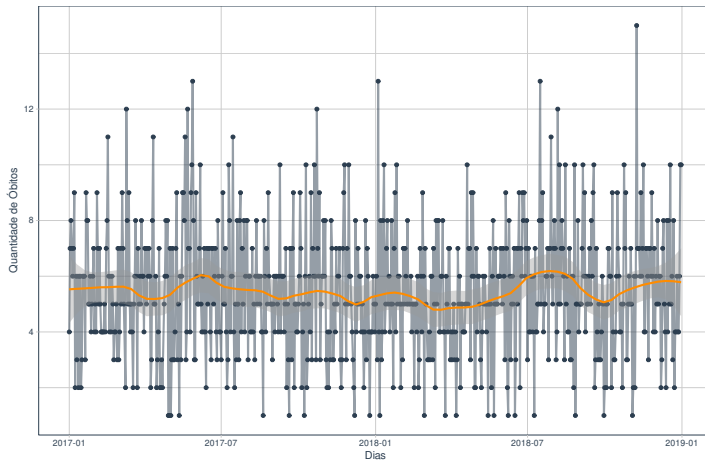
| | |
|-----------|-----------|
| Melhorado | Obito |
| 0.3526037 | 0.8087152 |

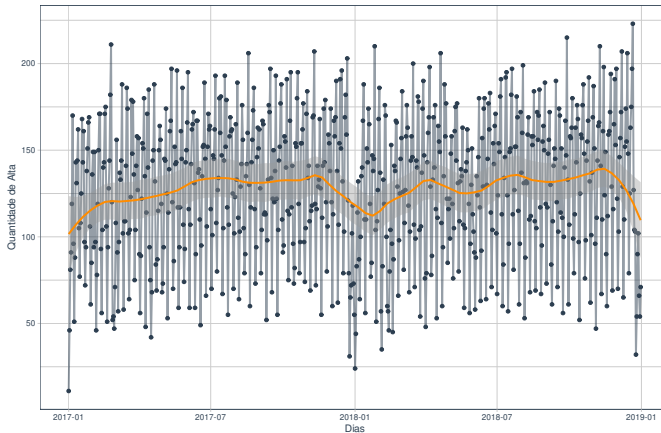
hospital



hospital







domingo
61

quarta
144

quinta
148

sábado
107

segunda
97

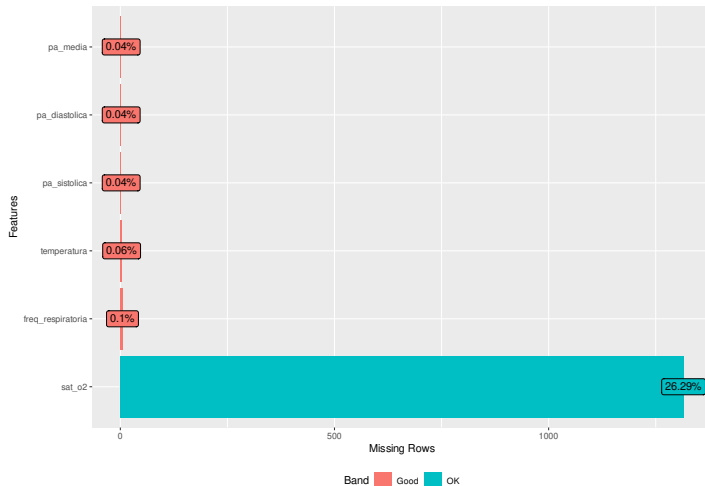
sexta
174

terça
135

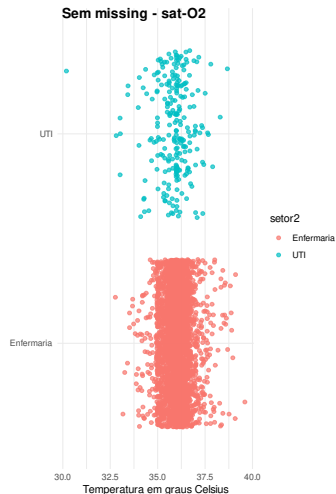
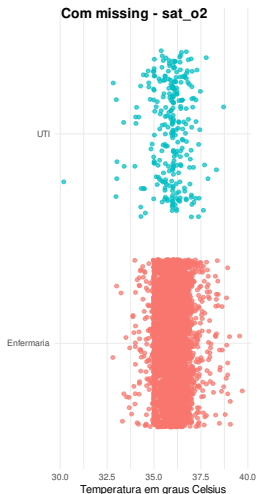
Dados faltantes - missing data

- **Missing completely at random (MCAR):** quando a probabilidade dos dados faltantes é a mesma entre as observações. Por exemplo, dados perdidos por um backup incorreto;
- **Missing at random (MAR):** quando os dados faltantes variam de acordo com outras variáveis. Por exemplo, missing em uma variável idade pode ser diferente entre mulheres e homens;
- **Missing not at random (MNAR):** quando a probabilidade de missing está relacionada ao missing. Por exemplo, dependendo da renda do cliente, é mais provável que ele não responda sobre a renda.

Exemplo - dados faltantes



Exemplo - dados faltantes



Exemplo - dados faltantes

- Distribuição do desfecho na **UTI**

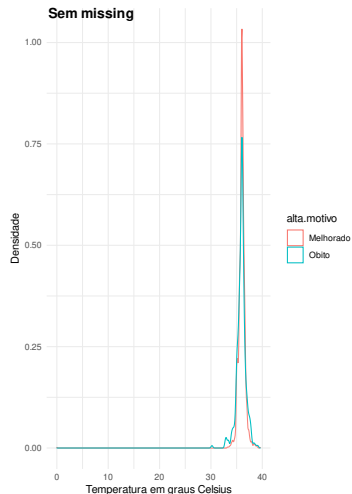
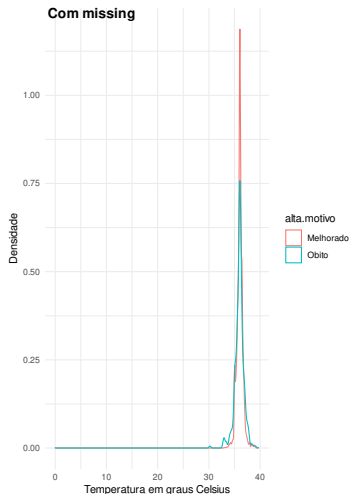
Com missing - sat_o2

| | |
|-----------|--------|
| Melhorado | Obito |
| 0.0438 | 0.9562 |

Sem missing - sat_o2

| | |
|-----------|--------|
| Melhorado | Obito |
| 0.0412 | 0.9588 |

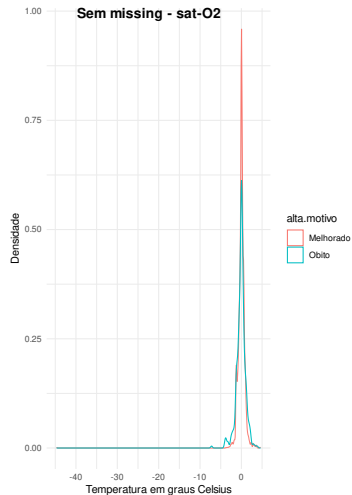
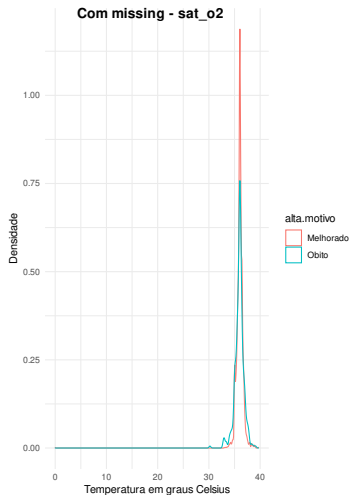
Exemplo - dados faltantes



Imputação

- Algoritmo **KNN** (*k-nearest neighbors*)
 - 1 Recebe um dado não classificado;
 - 2 Define o tamanho da vizinhança (K);
 - 3 Mede a distância (Euclidiana) do novo; dado com todos os outros dados que já estão classificados;
 - 4 Obtém a menor ou menores distâncias;
 - 5 Verifica o valor de cada um dos dados que tiveram a menor distância;
 - 6 Toma como resultado o valor que mais apareceu dentre os dados que tiveram as menores distâncias.

Imputação



Distribuições de probabilidade

- Há diversos modelos probabilísticos
- Distribuições de probabilidade de V.A. (discretas ou contínuas)
- Descrição das probabilidades associadas com os possíveis valores de X
 - **Variáveis discretas** \Rightarrow suporte em um conjunto de valores enumeráveis (finitos ou infinitos)
 - **Variáveis contínuas** \Rightarrow suporte em um conjunto não enumerável de valores

Distribuição de probabilidade - variável aleatória contínua

Definição: função densidade de probabilidade

A função densidade de probabilidade $f(x)$ de uma V.A. contínua, é uma função definida em um intervalo de valores na reta dos reais, ou seja,

$$P[a \leq X \leq b] = \int_b^a f_x(x).$$

- i. A área total da abaixo da curva de $f(x)$ é 1.

$$\int_{-\infty}^{\infty} f_x(x) = 1$$

- ii. A função $f(x)$ é positiva ou zero.

$$f_x(x) \geq 0$$

Padronização

Gráfico Função de Densidade de Probabilidade

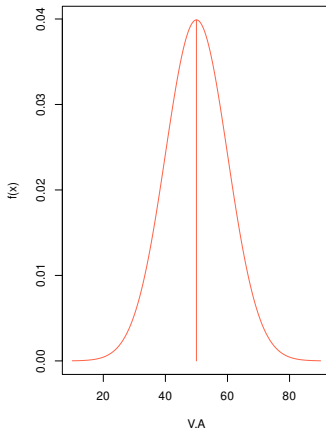
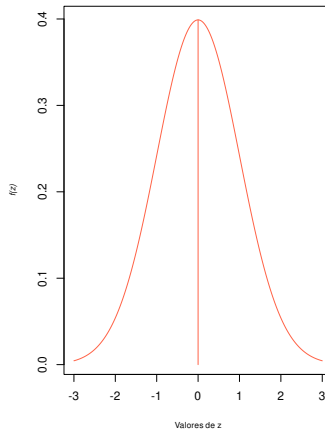


Gráfico Função de Densidade de Probabilidade



Modelo Gaussiano(ou Normal)

Definição:

Uma V.A. contínua é dita ter distribuição gaussiana com parâmetros μ e σ^2 se sua f.d.p. for:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right], \quad -\infty < x < \infty$$

Notação: $X \sim N(\mu, \sigma^2)$

Normal padrão

Definição:

Uma V.A. contínua que segue uma distribuição normal com parâmetros $\mu = 0$ e $\sigma^2 = 1$, é conhecida com uma distribuição normal padrão.

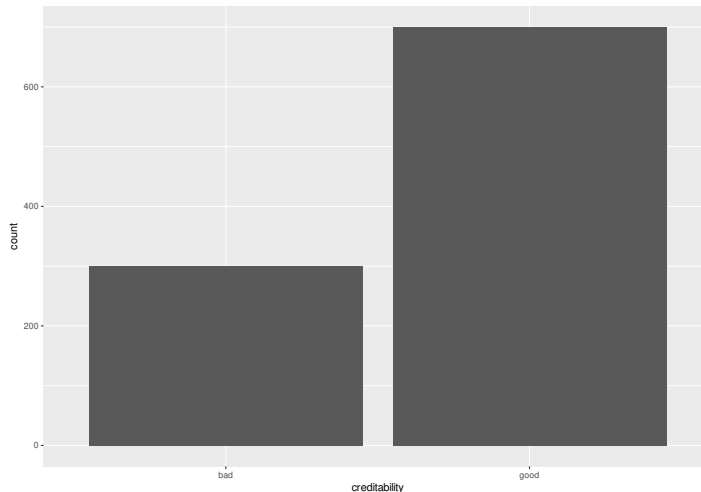
Notação: $Z \sim N(\mu = 0, \sigma^2 = 1)$

e sua f.d.p. é:

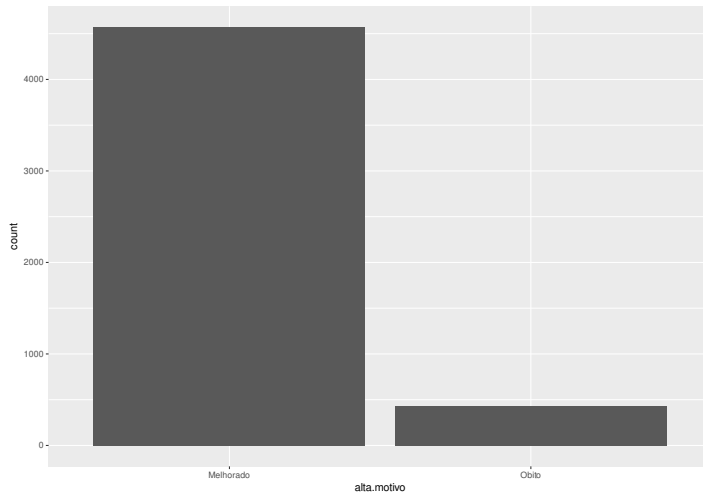
$$f(z) = \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{1}{2} (z)^2 \right], \quad -\infty < z < \infty$$

$$Z = \frac{X - \mu}{\sigma}$$

Dados desbalanceados - Exemplo 1



Dados desbalanceados - Exemplo 2



Probabilidade condicional

Intuitivamente, a probabilidade de ocorrer um evento pode muito bem ser influenciada pelo ocorrência ou não de outros eventos

Exemplo:

- Probabilidade de que um aluno não comparecer a uma aula será influenciado:
 - esteja chovendo ou não,
 - trote dos calouros foi ou não realizado na noite anterior.
 - Internet falhou ou não
- Esta é a probabilidade condicional

Probabilidade condicional

Intuitivamente, a probabilidade de ocorrer um evento pode muito bem ser influenciada pelo ocorrência ou não de outros eventos

Exemplo:

- Probabilidade de que um aluno não comparecer a uma aula será influenciado:
 - esteja chovendo ou não,
 - trote dos calouros foi ou não realizado na noite anterior.
 - Internet falhou ou não
- Esta é a probabilidade condicional

Probabilidade condicional

Definição

- Dados dois eventos A e B, a probabilidade condicional de A ocorrer, dado que ocorreu B é representado por $P(A|B)$ e dada por

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad \text{para } P(B) > 0.$$

Regra do produto

- Da definição de probabilidade condicional

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

temos que

$$P(A \cap B) = P(A|B) \cdot P(B).$$

- Caso $P(B) = 0$, definimos $P(A|B) = P(A)$.

Independência de eventos

Os eventos A e B são **eventos independentes** se a ocorrência de B não altera a probabilidade de ocorrência de A, ou seja, eventos A e B são independentes se

$$P(A|B) = P(A) \quad \text{e também que} \quad P(B|A) = P(B).$$

Com isso, e a regra do produto, temos que

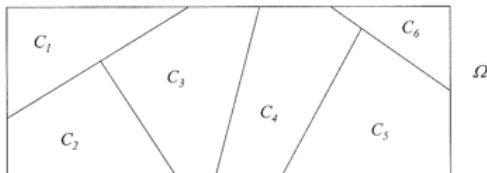
$$P(A \cap B) = P(B) \cdot P(A|B) = P(B) \cdot P(A).$$

$$P(A \cap B) = P(A) \cdot P(B|A) = P(A) \cdot P(B).$$

Partição do espaço amostral

Dizemos que os eventos C_1, C_2, \dots, C_k formam uma **partição** do espaço amostral, se eles não tem interseção entre si, e se sua união é igual ao espaço amostral. Isto é,

$$C_i \cap C_j = \emptyset \quad \text{para} \quad i \neq j \quad \text{e} \quad \bigcup_{i=1}^k C_i = \Omega.$$



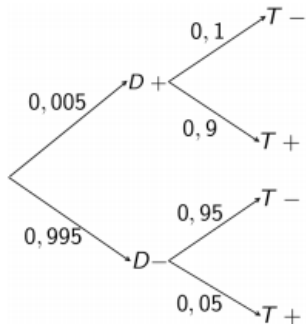
Teorema de Bayes

Suponha que os eventos C_1, C_2, \dots, C_k formem uma partição de Ω e que suas probabilidades sejam conhecidas. Suponha, ainda, que para um evento A , se conheçam as probabilidades $P(A|C_i)$ para todo $i = 1, 2, \dots, k$. Então, para qualquer j ,

$$P(C_j|A) = \frac{P(C_j)P(A|C_j)}{\sum_{i=1}^k P(C_i)P(A|C_i)}, \quad j = 1, 2, \dots, k.$$

Exemplo - Regra de Bayes

- Considere um teste de triagem de rotina para uma doença



- Teste altamente preciso
 - $P(FP) = P(T_+|D_-) = 0.05$
 - $P(FN) = P(T_-|D_+) = 0.10$

Exemplo - Regra de Bayes

- Qual a probabilidade de ter a doença, dado que o teste foi positivo?

- $$P(D_+|T_+) = \frac{P(D_+) * P(T_+|D_+)}{P(T_+)}$$

- Calculamos o denominador utilizando a lei da probabilidade total:

- $$P(T_+) = P(D_+) * P(T_+|D_+) + P(D_-) * P(T_+|D_-) =$$

$$0.005 * .9 + 0.995 * .05 = 0.05425$$

Assim,

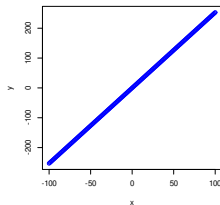
- $$P(D_+|T_+) = \frac{P(D_+) * P(T_+|D_+)}{P(T_+)} = \frac{0.005 * 0.9}{0.05425} = .083$$

Dados desbalanceados

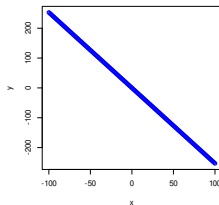
- **SMOTE** (*Synthetic Minority Over-sampling Technique*): reamostrar o conjunto de dados original por superamostragem da classe minoritária.
- Utilizar uma amostra balanceada excluindo registros da classe majoritária.
- Definir um ponto de corte diferente do *default*.

Correlação

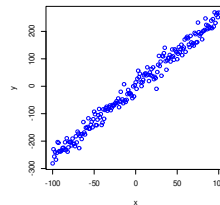
Correlação perfeita positiva



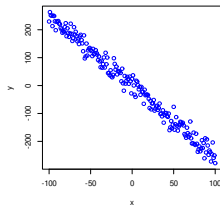
Correlação perfeita negativa



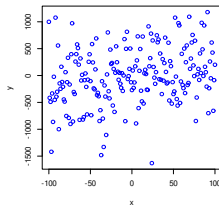
Alta correlação positiva



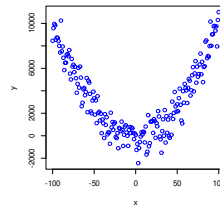
Alta correlação negativa



Baixa correlação



Correlação não linear



Correlação

Coeficientes de correlação linear

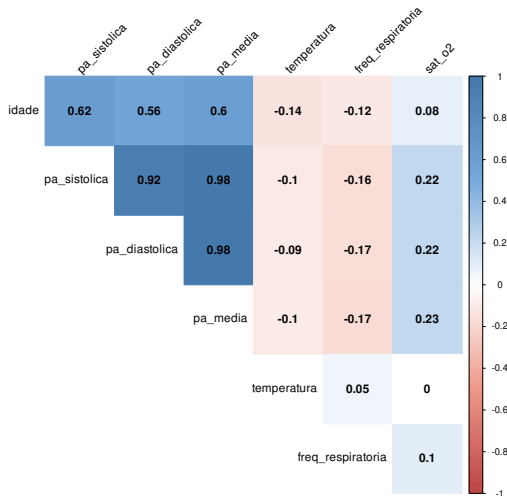
- ***Spearman***

$$r = 1 - \frac{6 \sum_i d_i^2}{n(n^2 - 1)}$$

- **Pearson**

$$\rho(X_1, X_2) = \frac{\text{cov}(X_1, X_2)}{\sigma_{x1}\sigma_{x2}}$$

Matriz de correlação



Desafio2

Com os dados coletados no desafio anterior, faça:

- Engenharia de características
 - 1 Faça uma análise exploratória de seus dados. Com gráficos e medidas descritivas;
 - 2 Crie novas características com base em outras já presente no conjunto de dados;
 - 3 Se houver dados faltantes, identifique-os e use dois tipos de imputação e compare os resultados;
 - 4 Se houver dados desbalanceados em sua label(Y), utilize a técnica *SMOTE* ou redução da classe majoritária;
 - 5 Faça uma matriz de correlação linear, interprete e decida se irá excluir ou não alguma característica.

Desafio2

Com os dados coletados no desafio anterior, faça:

- Engenharia de características
 - 1 Faça uma análise exploratória de seus dados. Com gráficos e medidas descritivas;
 - 2 Crie novas características com base em outras já presente no conjunto de dados;
 - 3 Se houver dados faltantes, identifique-os e use dois tipos de imputação e compare os resultados;
 - 4 Se houver dados desbalanceados em sua label(Y), utilize a técnica *SMOTE* ou redução da classe majoritária;
 - 5 Faça uma matriz de correlação linear, interprete e decida se irá excluir ou não alguma característica.

Desafio2

Com os dados coletados no desafio anterior, faça:

- Engenharia de características
 - 1 Faça uma análise exploratória de seus dados. Com gráficos e medidas descritivas;
 - 2 Crie novas características com base em outras já presente no conjunto de dados;
 - 3 Se houver dados faltantes, identifique-os e use dois tipos de imputação e compare os resultados;
 - 4 Se houver dados desbalanceados em sua label(Y), utilize a técnica *SMOTE* ou redução da classe majoritária;
 - 5 Faça uma matriz de correlação linear, interprete e decida se irá excluir ou não alguma característica.

Desafio2

Com os dados coletados no desafio anterior, faça:

- Engenharia de características
 - 1 Faça uma análise exploratória de seus dados. Com gráficos e medidas descritivas;
 - 2 Crie novas características com base em outras já presente no conjunto de dados;
 - 3 Se houver dados faltantes, identifique-os e use dois tipos de imputação e compare os resultados;
 - 4 Se houver dados desbalanceados em sua label(Y), utilize a técnica *SMOTE* ou redução da classe majoritária;
 - 5 Faça uma matriz de correlação linear, interprete e decida se irá excluir ou não alguma característica.

