



Análise de custos com cuidados médicos

Vinícius H. P. Cavalcanti



PERGUNTAS

- Gênero, índice de massa corporal (imc), quantidade de filhos, se é fumante ou não explicam os gastos?
- Se explicam, quanto explicam ? Todas as variáveis explicam 100% dos gastos?
- Existe diferença dos custos entre regiões para a variável mais explicativa ?



SOBRE O CONJUNTO DE DADOS

- Foi obtido na plataforma Kaggle
- É utilizado no livro `Machine Learning with R` por Brett Lantz
- Segundo o autor do livro, foi desenvolvido com base no censo demográfico dos EUA
- Possui 1338 observações
- Possui 7 colunas, sendo essas:
 - idade
 - gênero
 - índice de massa corporal
 - quantidade de filhos
 - fumante?
 - região
 - custos



METODOLOGIAS



Estatística descritiva & Análise exploratória

- Quais os tipos de dados presentes (numérico ou categórico) ?
- Para os valores numéricos, qual a média? moda ? mediana? os dados variam?
- Para os valores categóricos, quantos valores distintos existem?
- Existem valores fora da curva?
- Visualmente, qual a relação entre custo e gênero, se é fumante ou não, região e quantidade de filho?



Inferência estatística

- Regressão linear múltipla
 - as variáveis explicam o custo?
 - qual a variável mais explicativa?
 - as variáveis explicam 100% dos custos?
- Análise de variância
 - entre a variável mais explicativa, existe diferença entre região ou por quantidade de filhos?

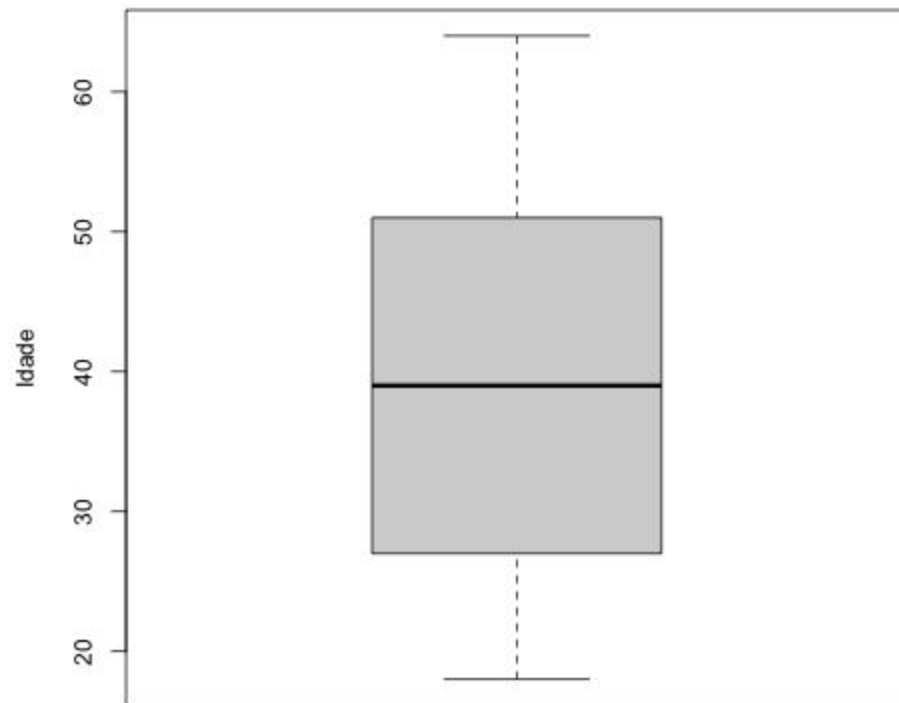


RESULTADOS E ANÁLISE

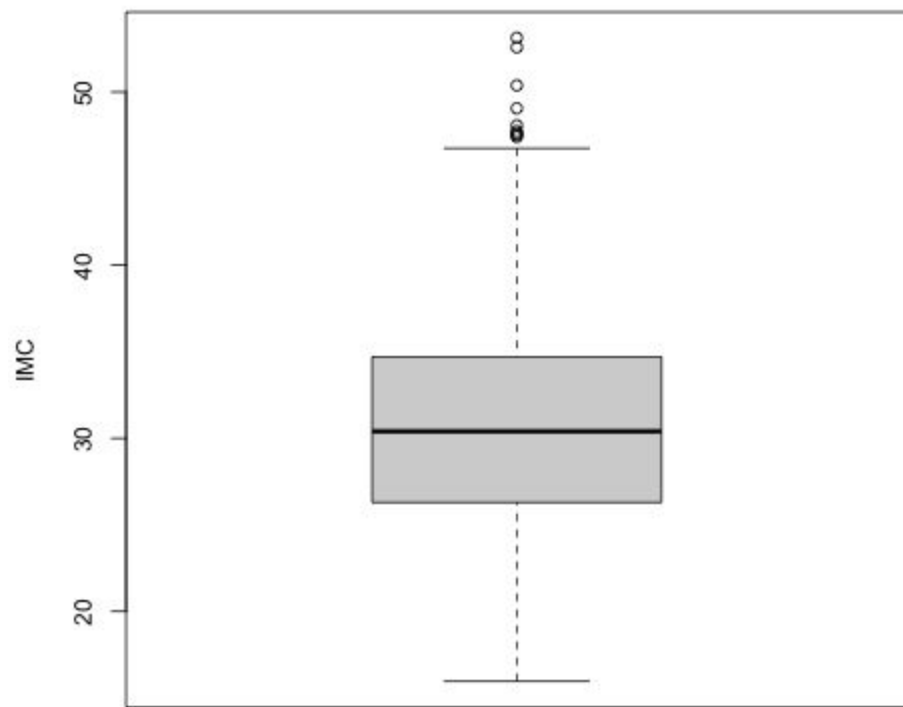


Estatística descritiva

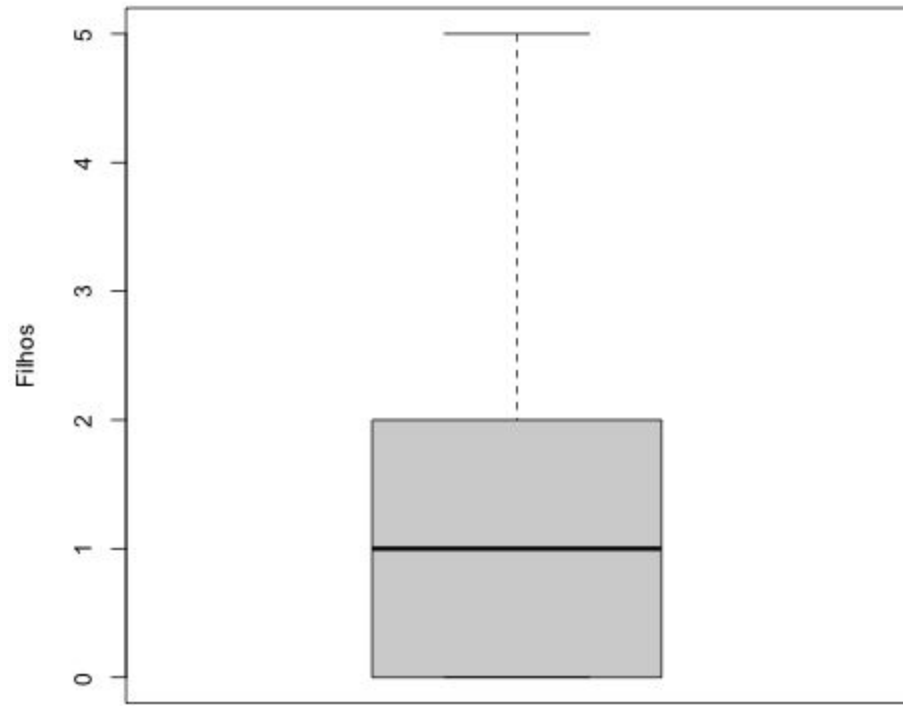
Idade



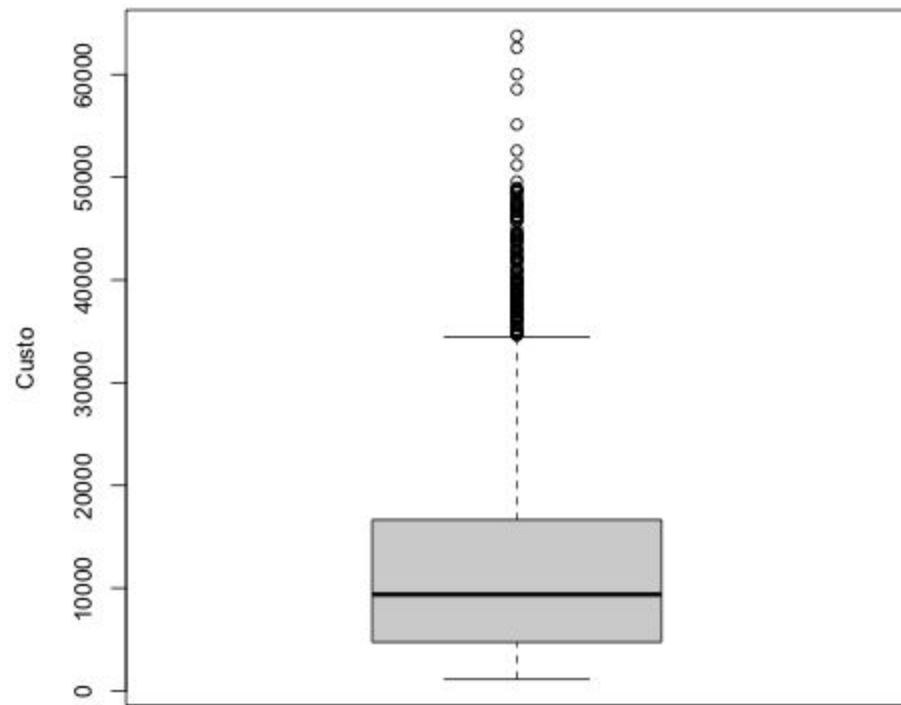
IMC

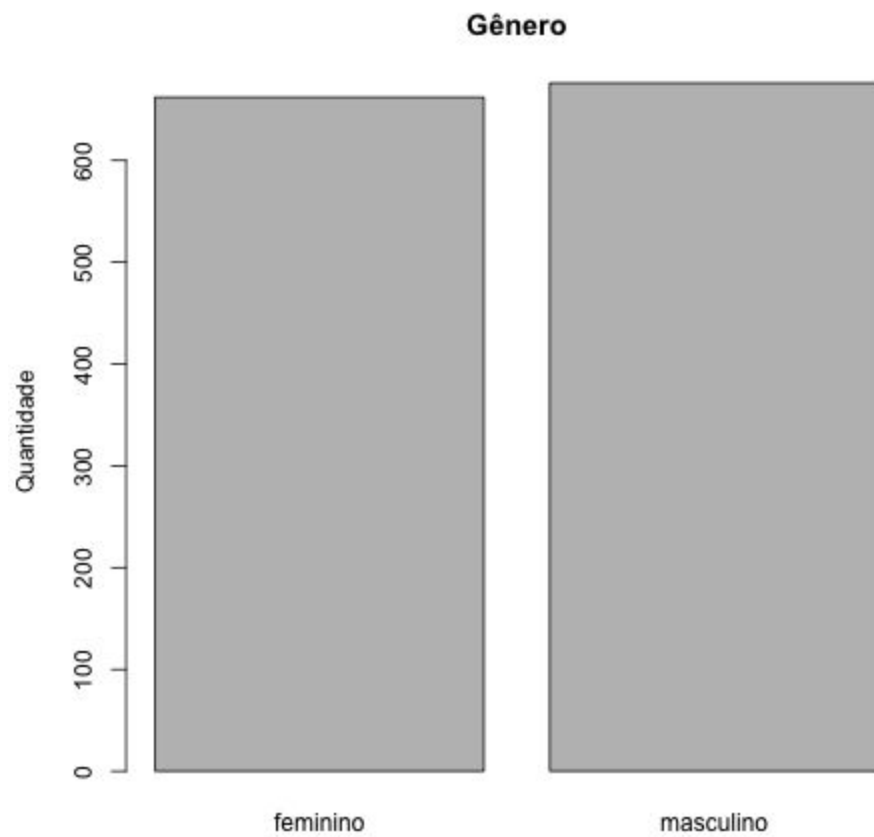


Filhos

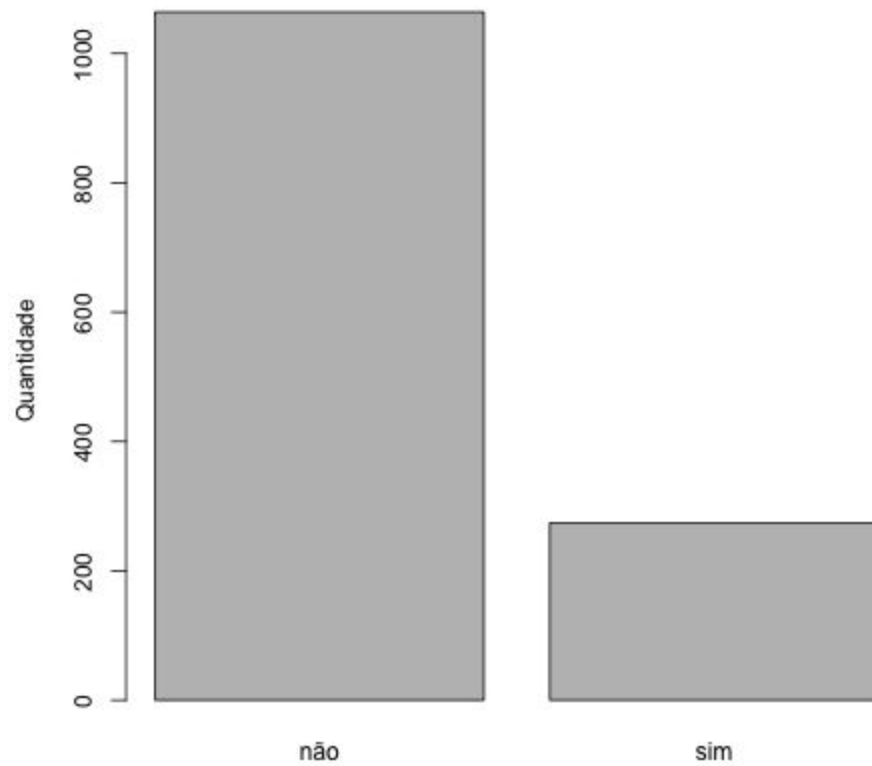


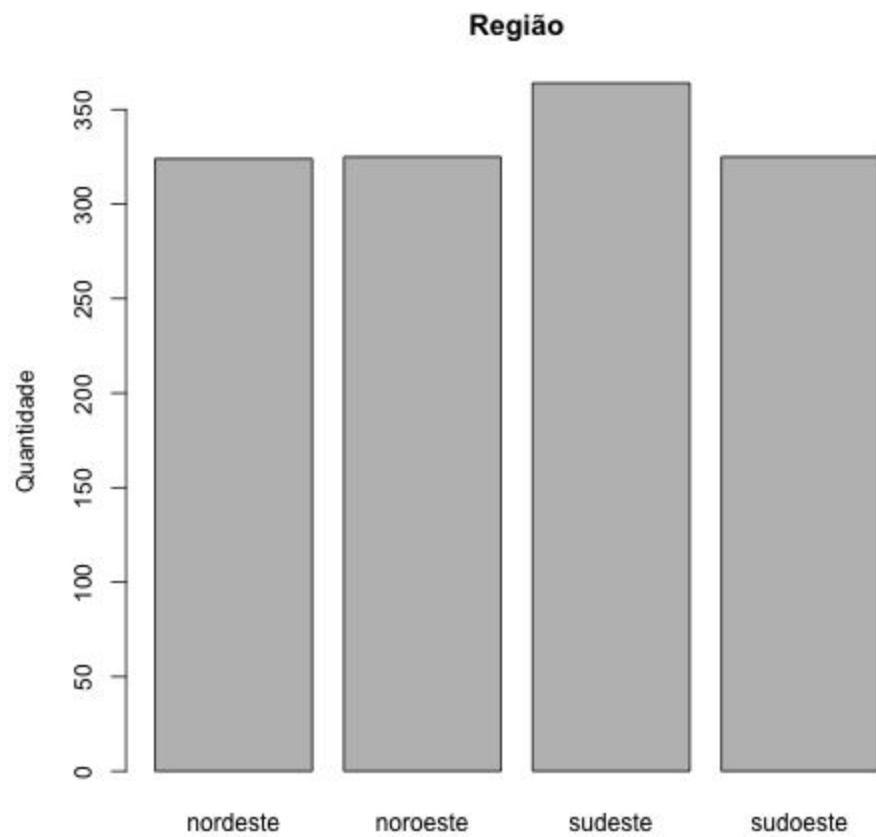
Custo





Fumante







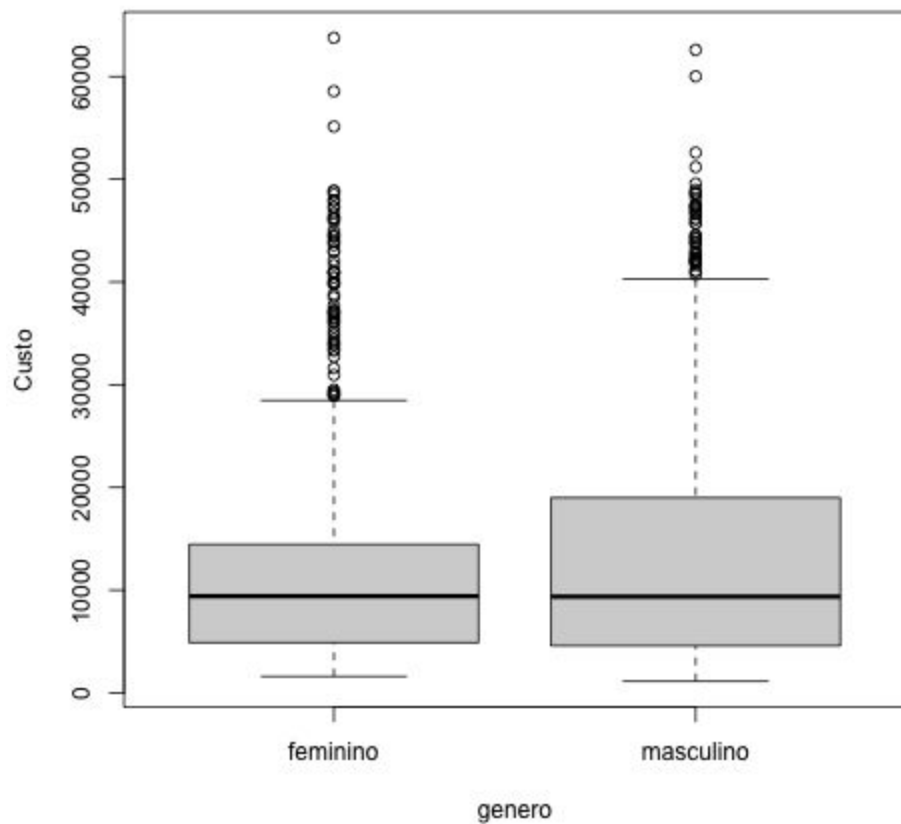
Observações vazias

idade	0
gênero	0
imc	0
filhos	0
fumante	0
região	0
custo	0

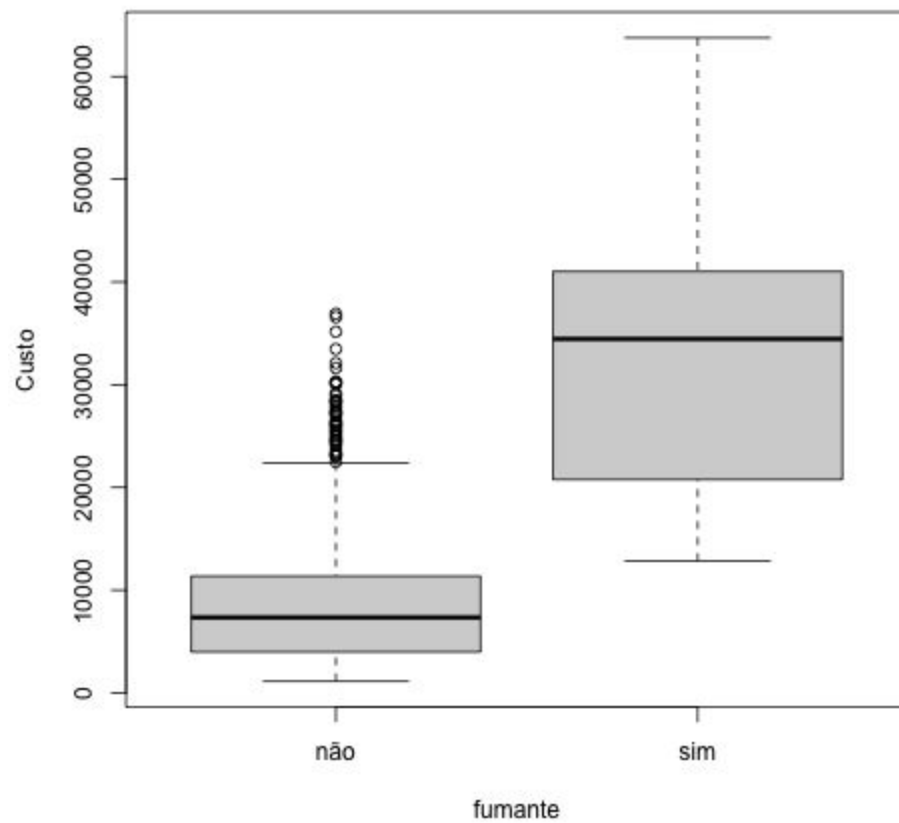


Análise exploratória

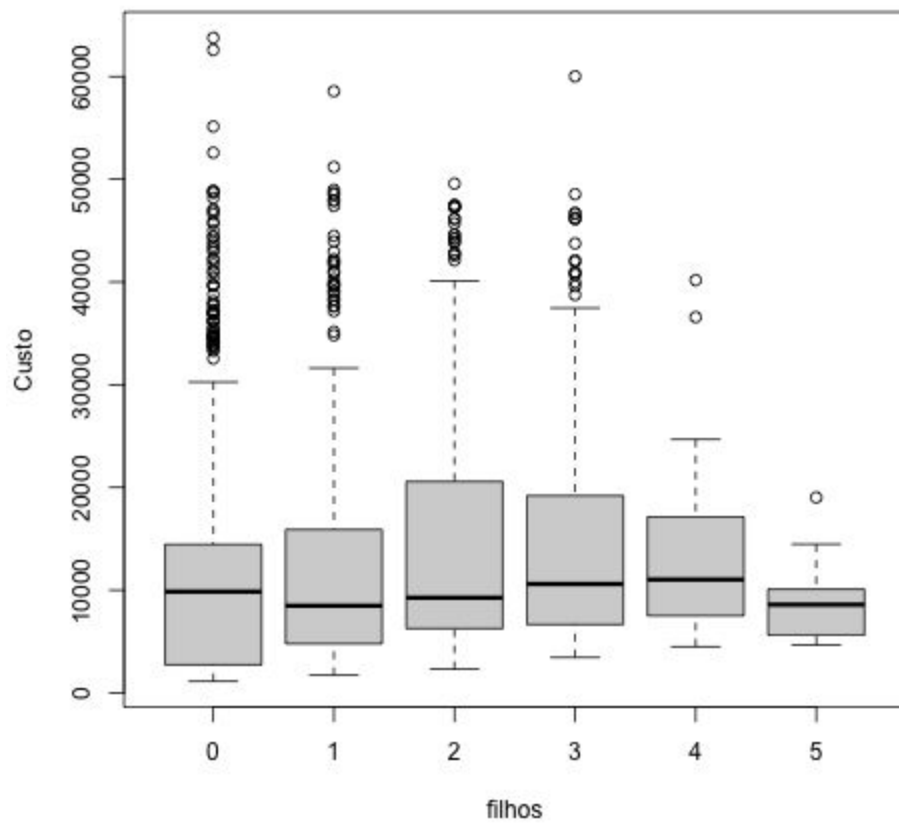
Custo por gênero



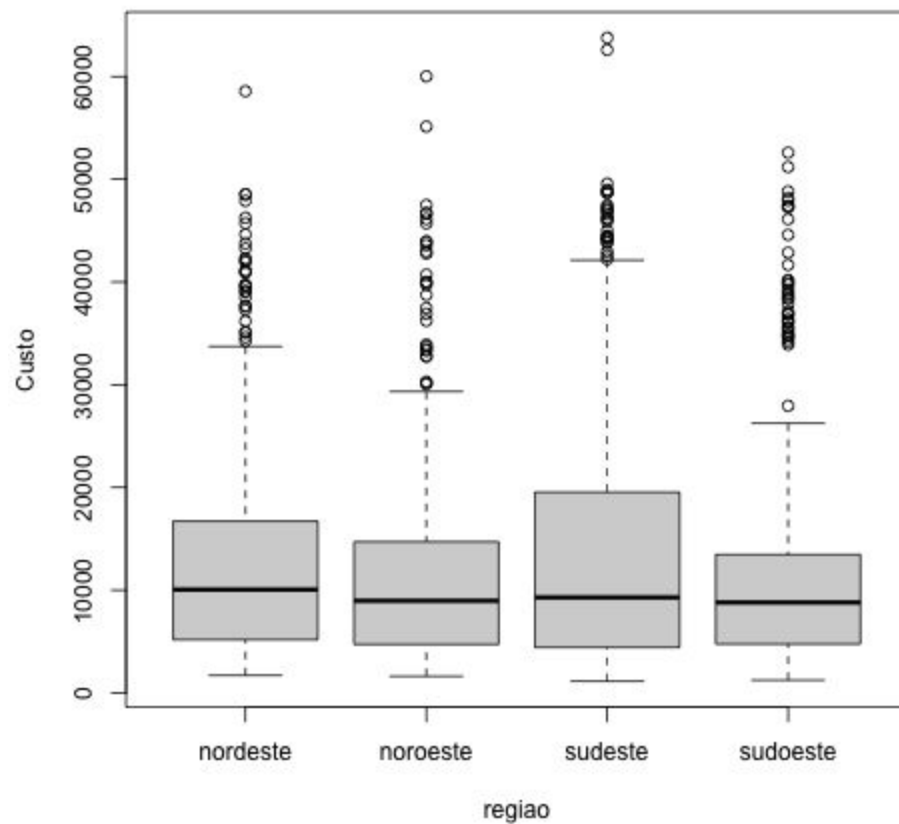
Custo por fumante



Custo por filhos



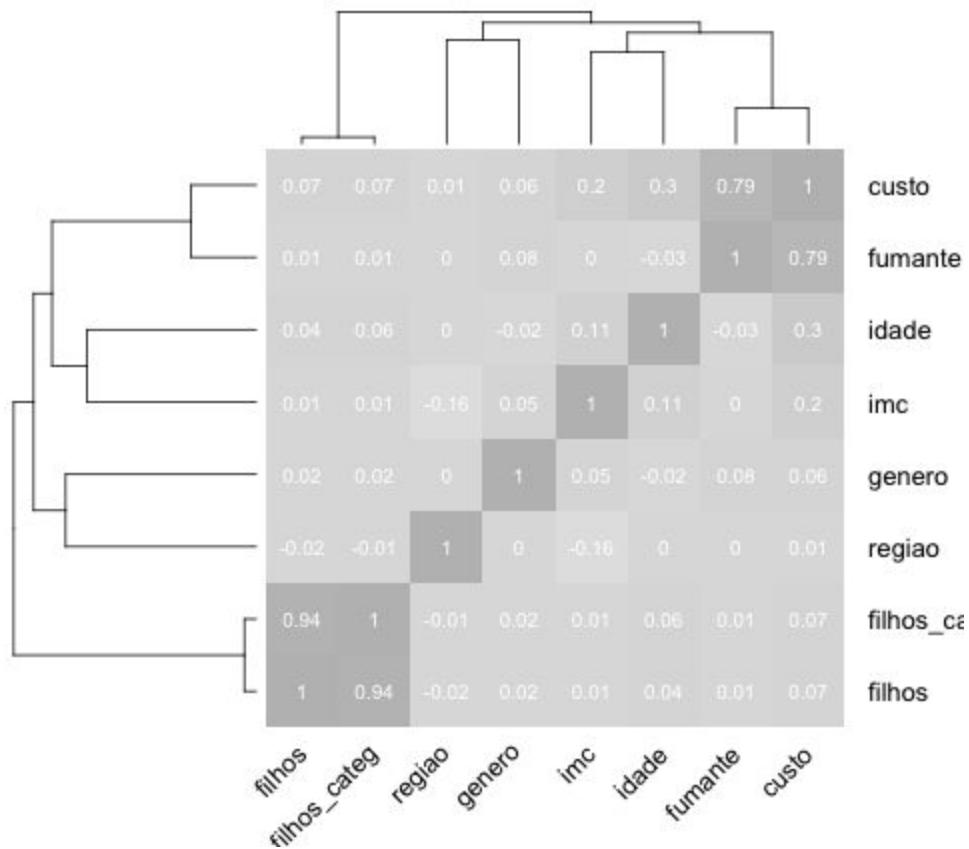
Custo por região





Inferência Estatística

Correlation Heatmap





lm(formula = custo ~ idade + imc + filhos + fumante, data = numeric_df)

	Estimativa	Erro	T-Valor	P-Valor
Intercepto	-12102.77	941.98	-12.848	<2e-16
Idade	257.85	11.90	21.675	<2e-16
IMC	321.85	27.38	11.756	<2e-16
Filhos	473.50	137.79	3.436	6,08E-04
Fumante	23811.40	411.22	57.904	2,00E-16

Erro Padrão Residual	6068	
R² Múltiplo	0,75	
R² Ajustado	0,75	
Estatística F	998,1	p < 2.2e-16



lm(formula = custo ~ fumante, data = numeric_df)

	Estimativa	Erro	T-Valor	P-Valor
Intercepto	8434.3	229.0	36.83	<2e-16
Fumante	23616.0	506.1	46.66	<2e-16

Erro Padrão Residual	7470	
R² Múltiplo	0,62	
R² Ajustado	0,62	
Estatística F	998,1	p < 2.2e-16



Shapiro-Wilk fumante\$custo | nao_fumante\$custo

	W	P-Valor
Fumantes	0,94	3,63E-06
Não fumantes	0,87286	2,20E-16



ANOVA custo para fumantes e não fumantes por região

	P-Valor
Fumantes	0,0205
Não fumantes	0,0917



Teste de Tukey

	diff	lwr	upr	p adj
noroeste-nordeste	518	-4.766	5.803	0,99
sudeste-nordeste	5.171	428	9.914	0,03
sudoeste-nordeste	2.596	-2.689	7.880	0,58
sudeste-noroeste	4.653	-297	9.603	0,07
sudoeste-noroeste	2.077	-3.394	7.548	0,76
sudoeste-sudeste	-2.576	-7.526	2.374	0,54



CONSIDERAÇÕES FINAIS

- A variável que mais explica o custo é a fumante
 - Somente a variável fumante explica 62% do custo
 - Onde se a pessoa é fumante existe um aumento de \$23811.40
- Para fumantes existe uma diferença significativa somente para a região sudeste-nordeste