

MTCARS 2.0

GABRIEL, MARLA,
MOISÉS

Tratamento de dados

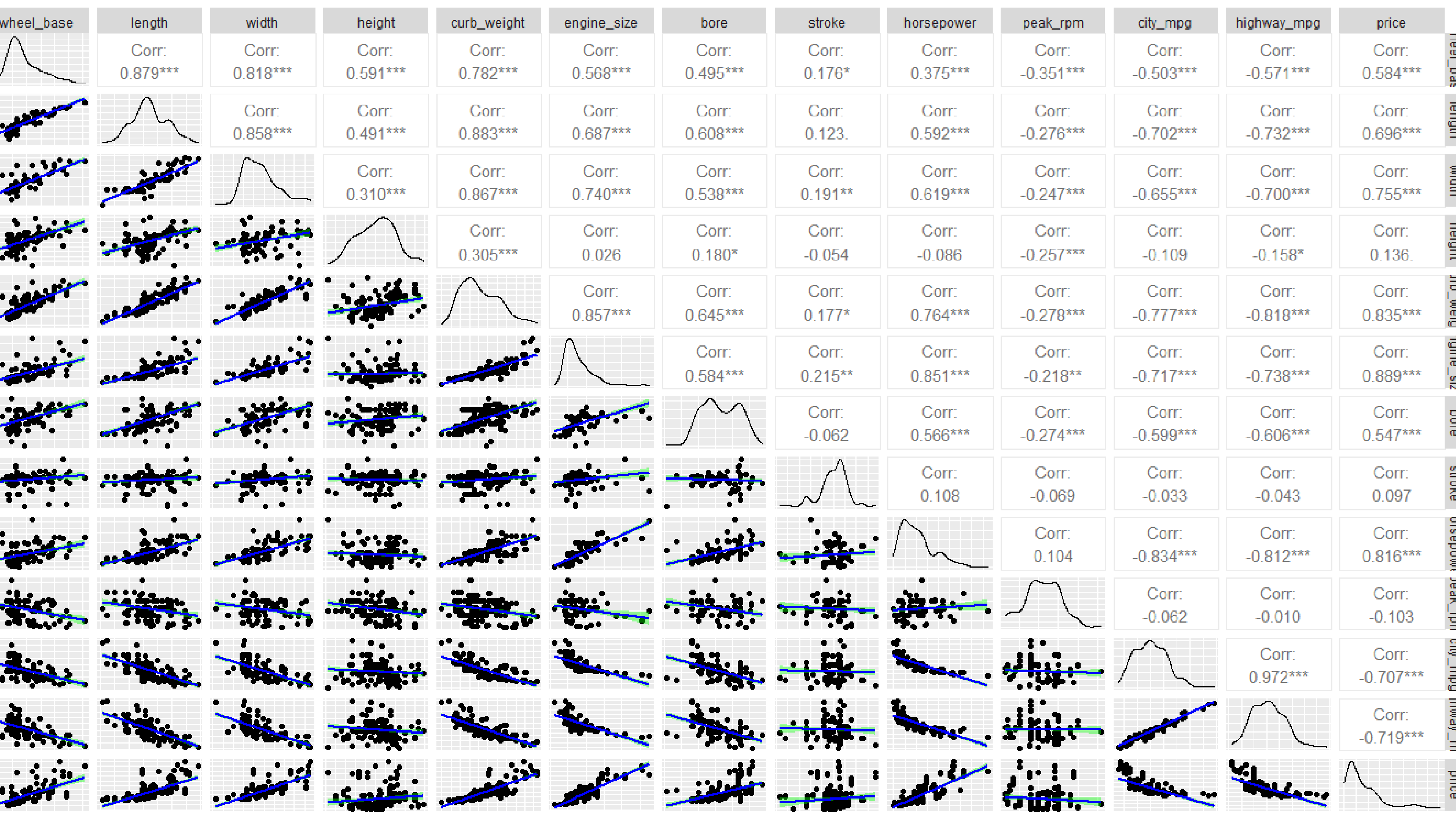
Resumo geral dos dados;
Variáveis sem nome;
Tipos/classes das variáveis;
Valores faltantes;
Categorização: `Compless_ratio`, `num_cylinders`.

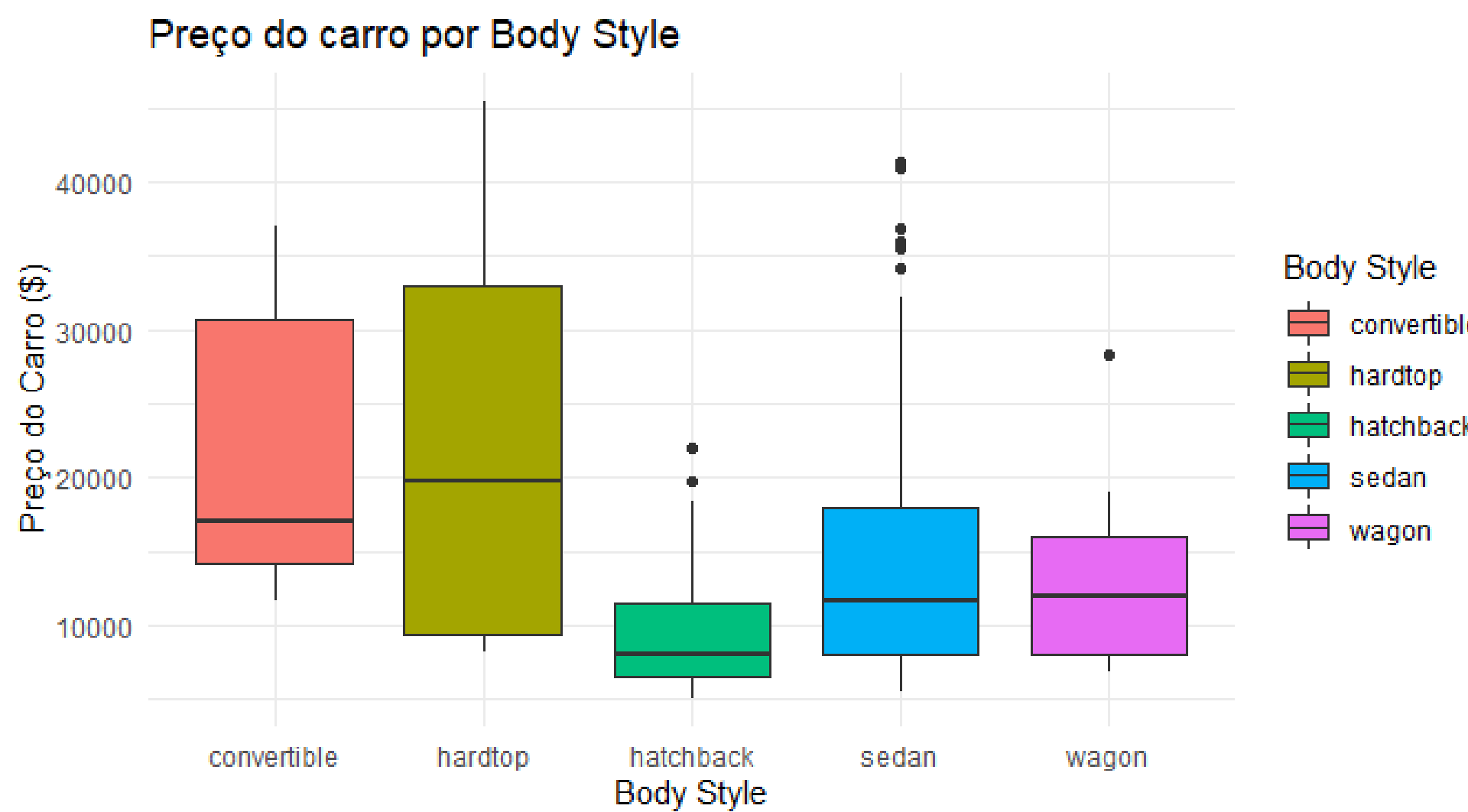
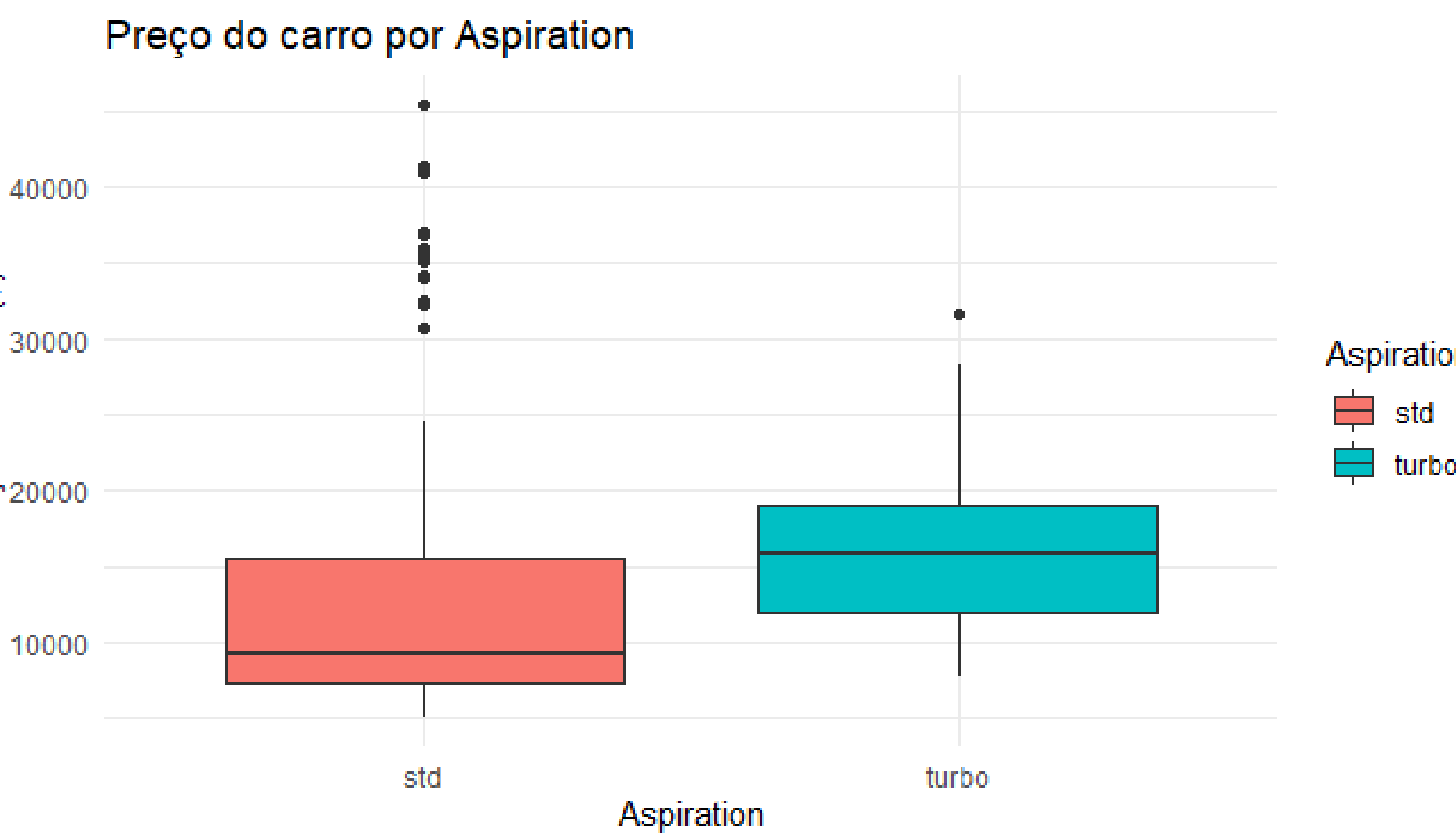
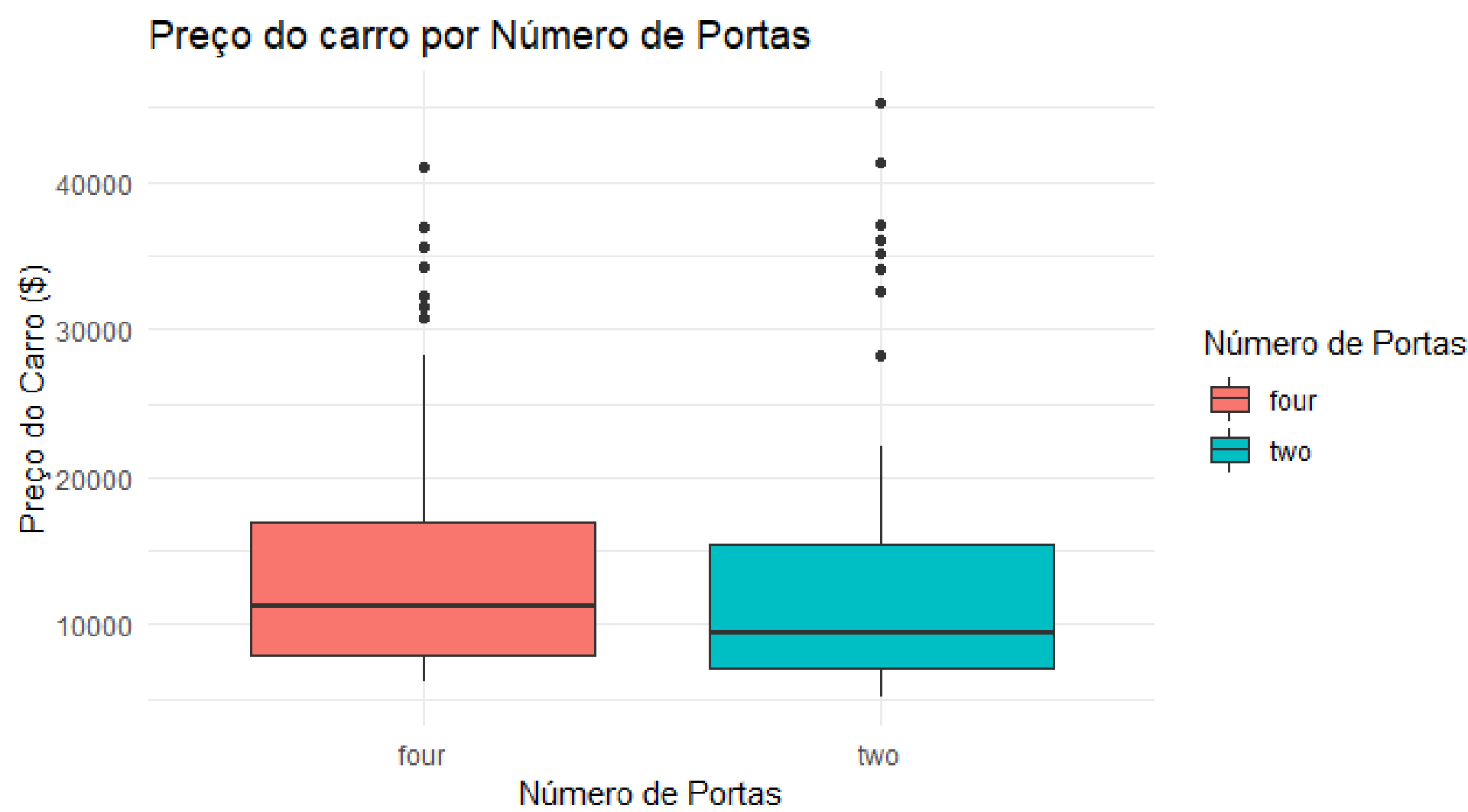
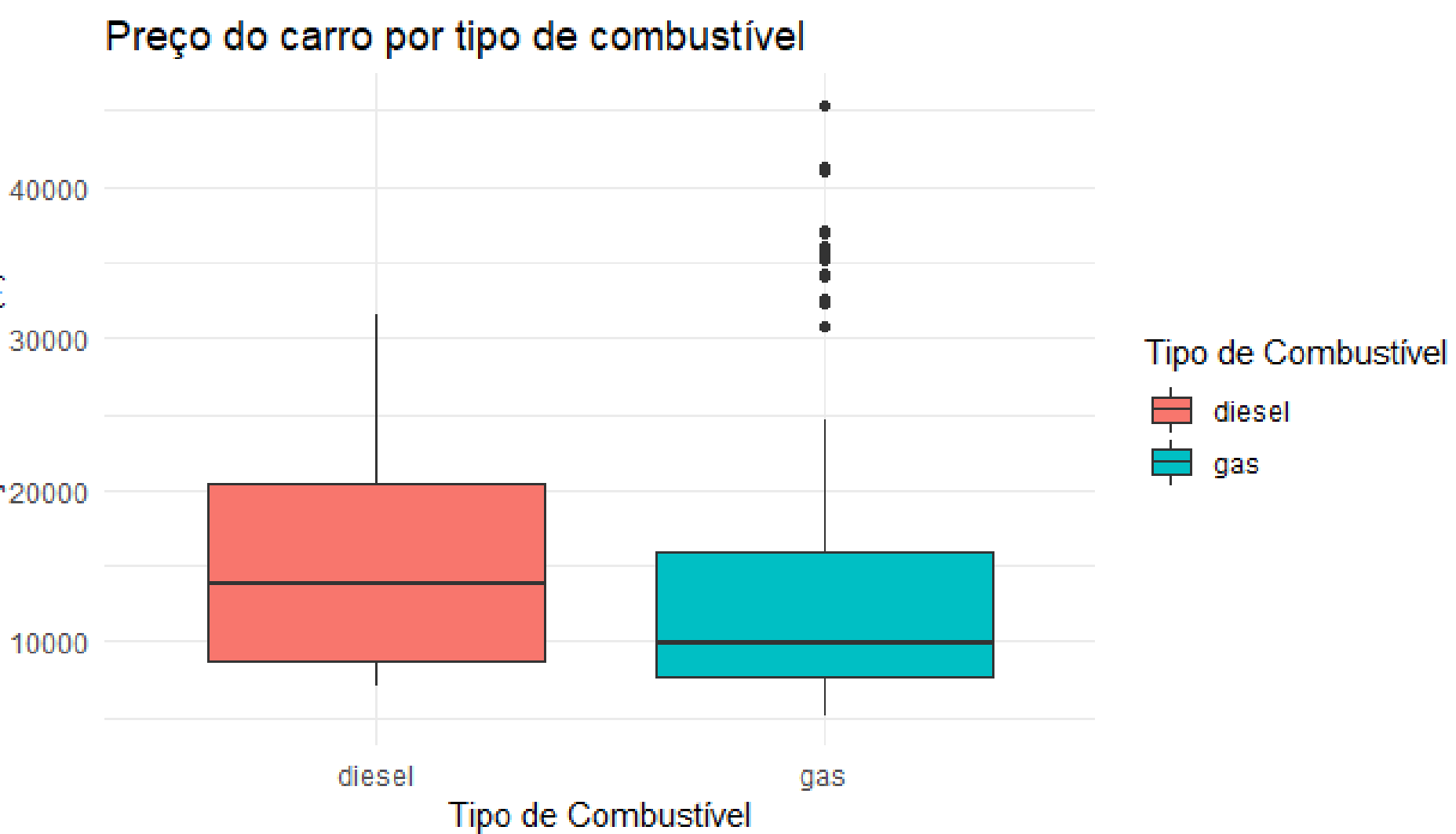


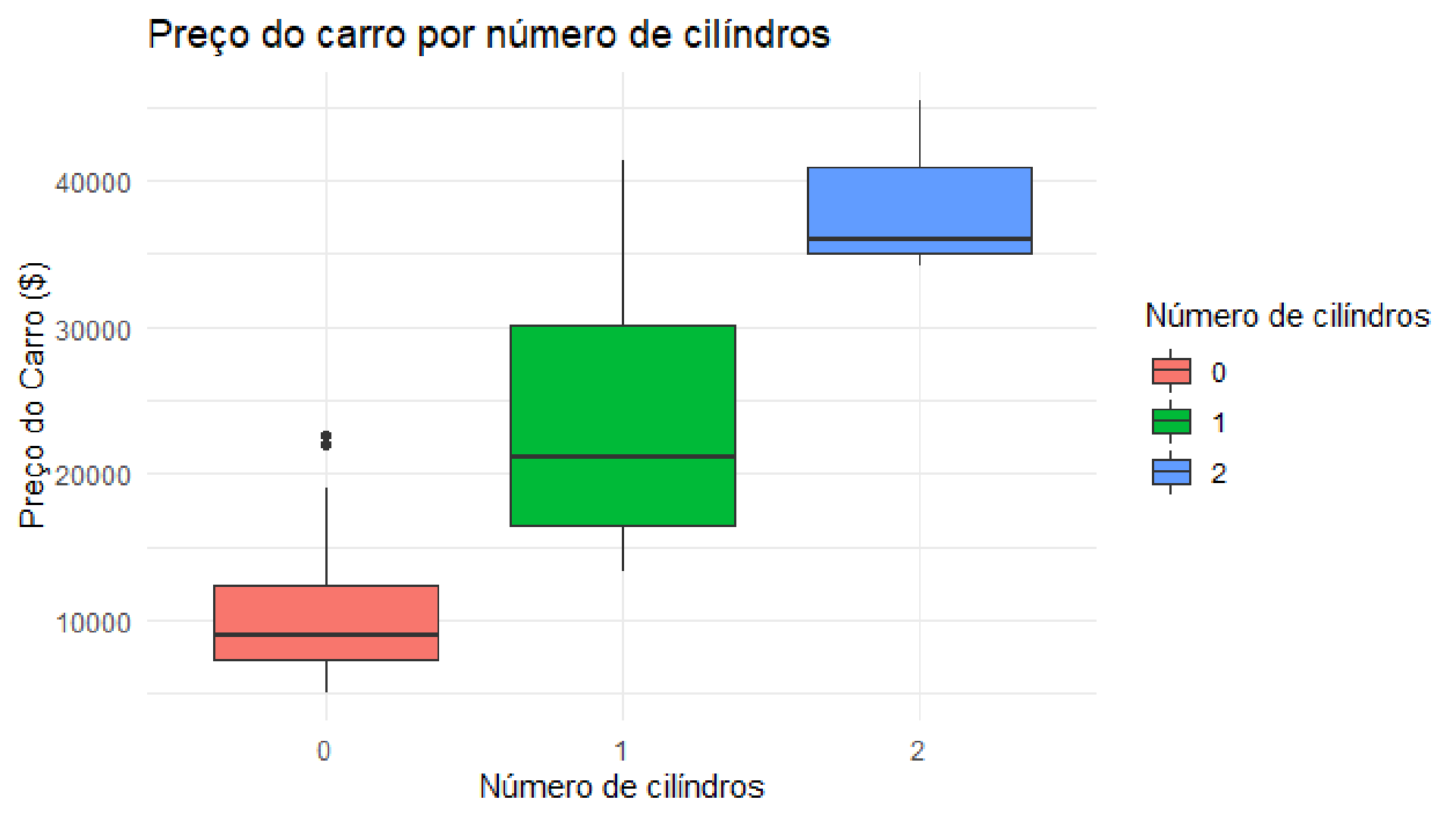
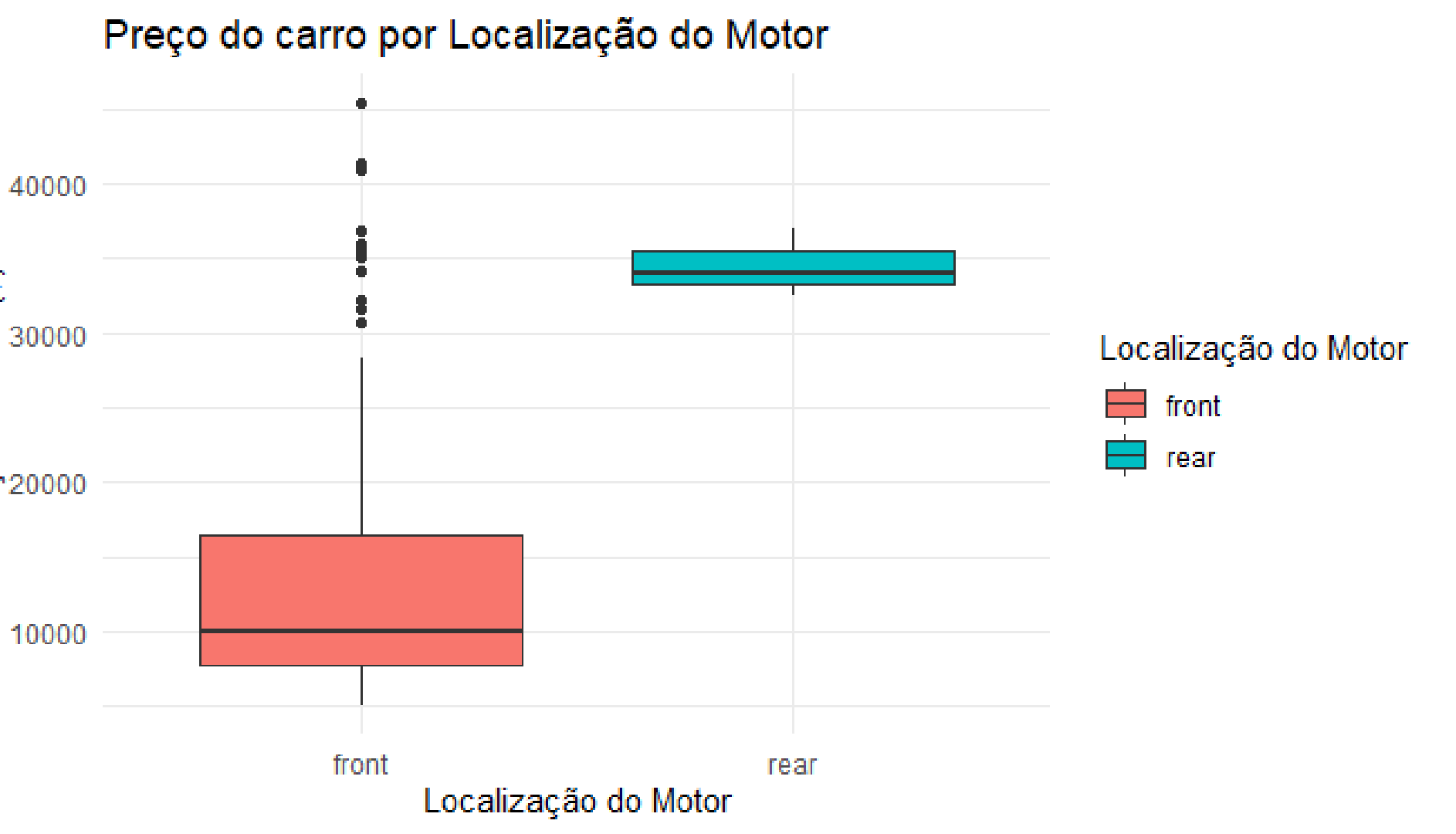
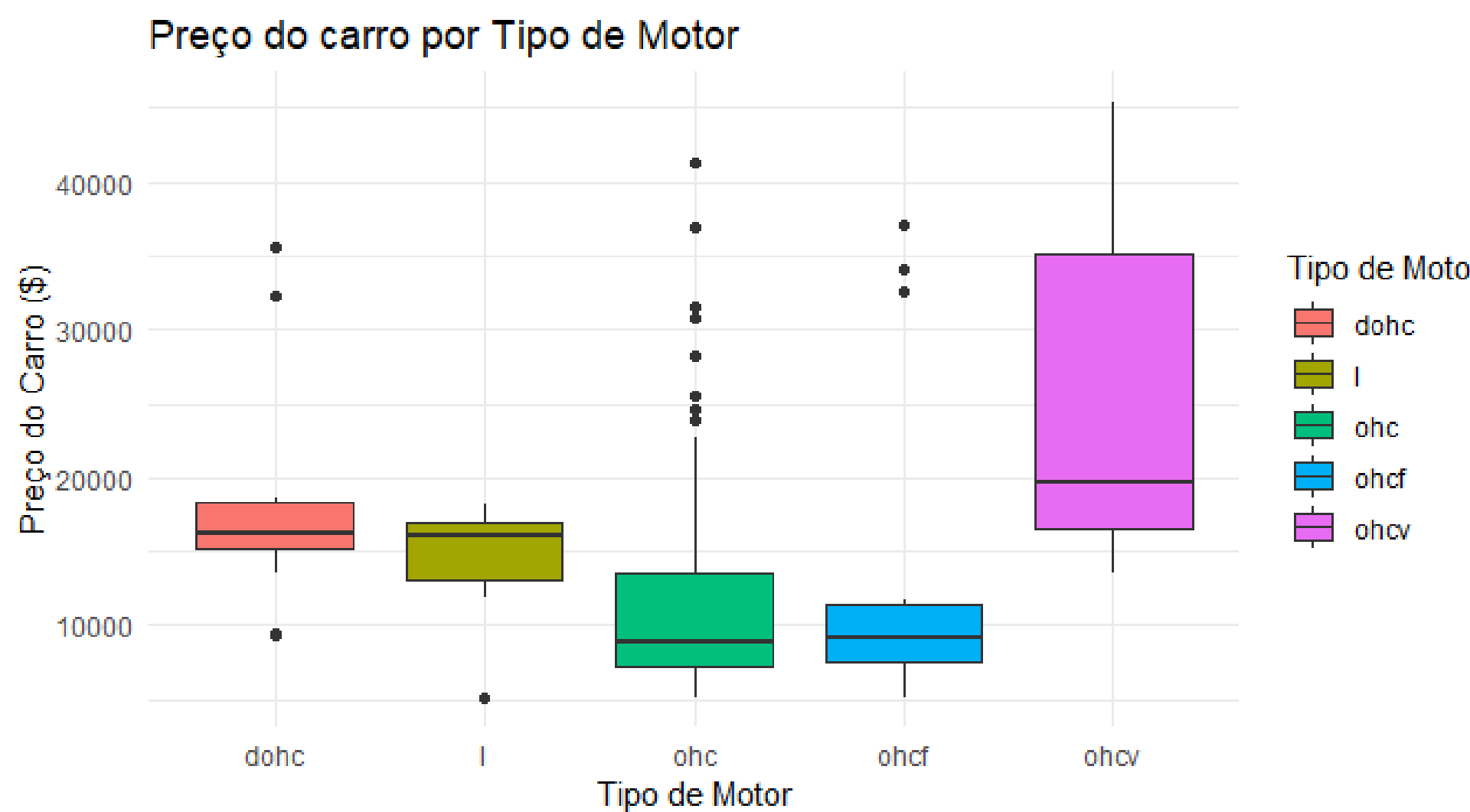
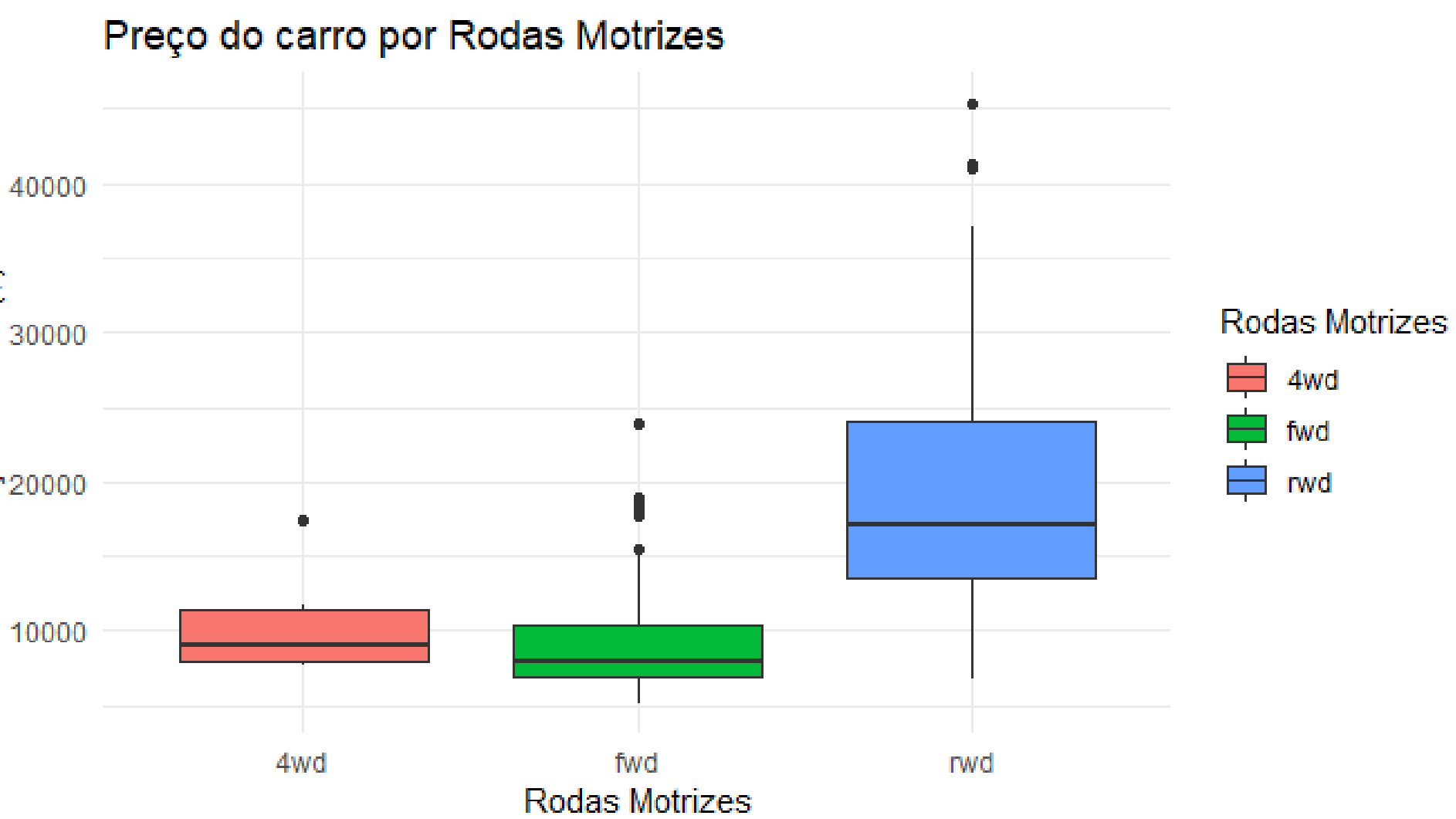


Exploratória

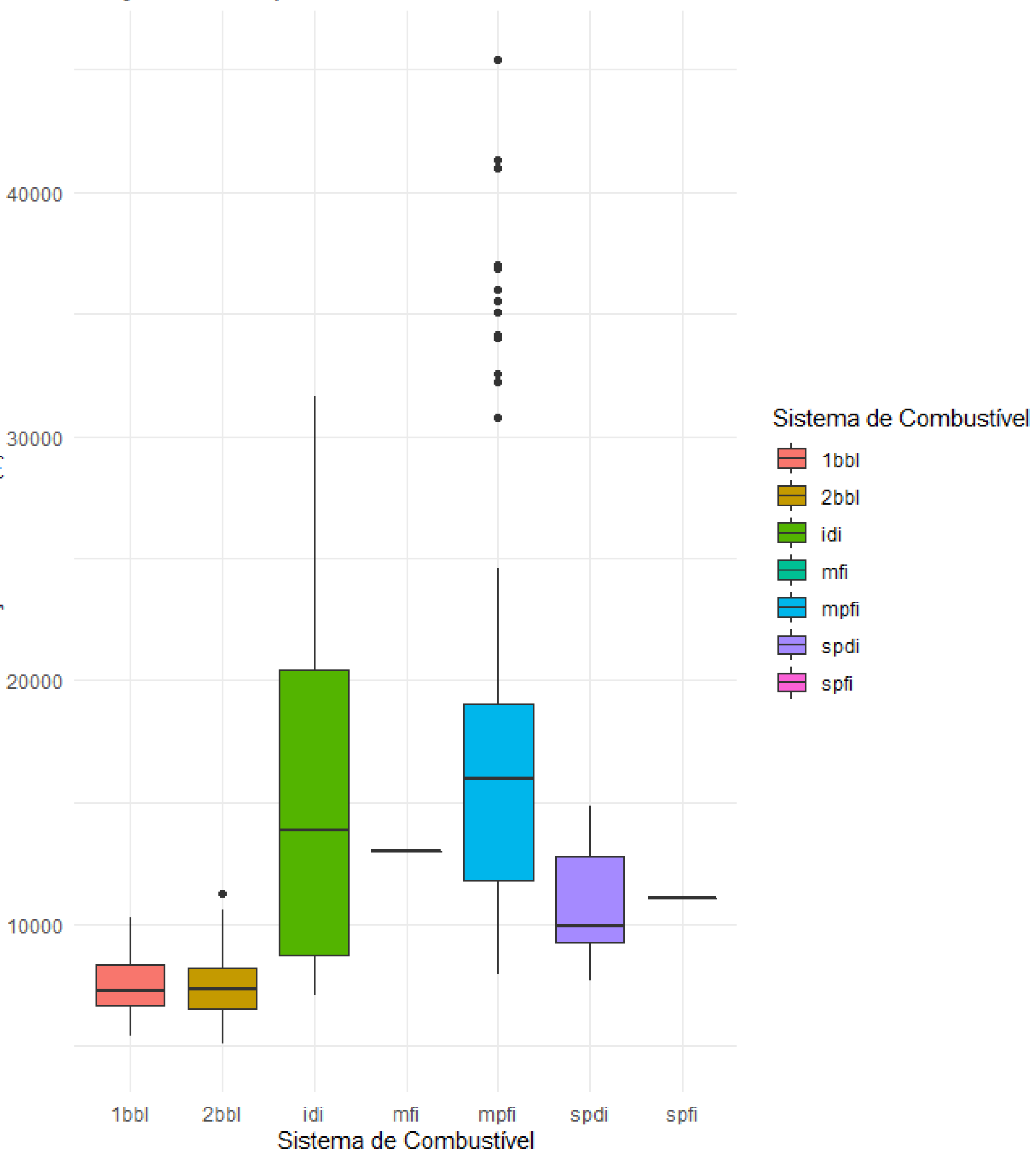
- Cálculo de descritivas;
- Gráficos bi-variados com o preço;
- Verificação de correlação entre variáveis explicativas;
- Testes de hipótese para correlação;
- Correlação de Crammer para variáveis categóricas.



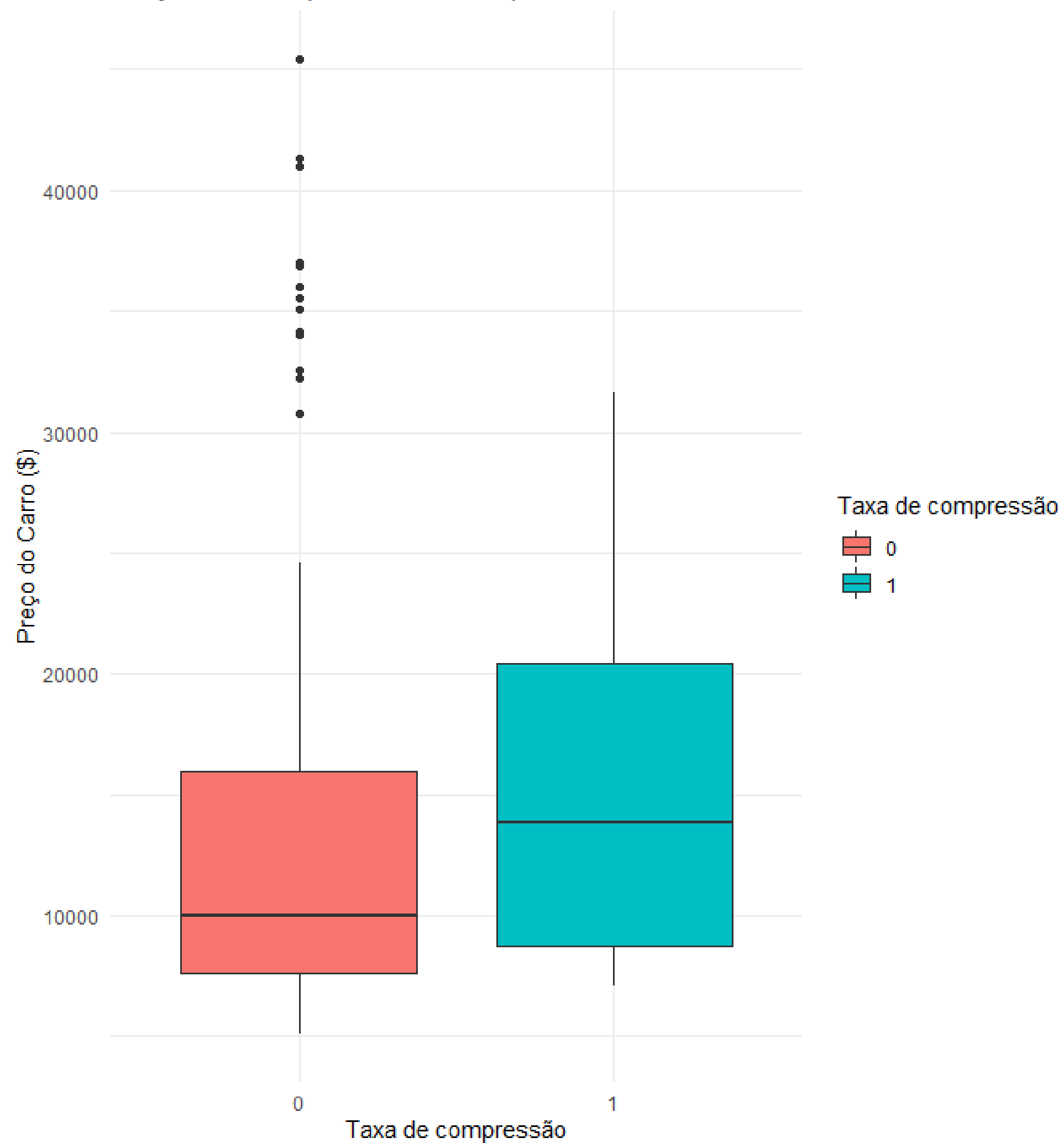




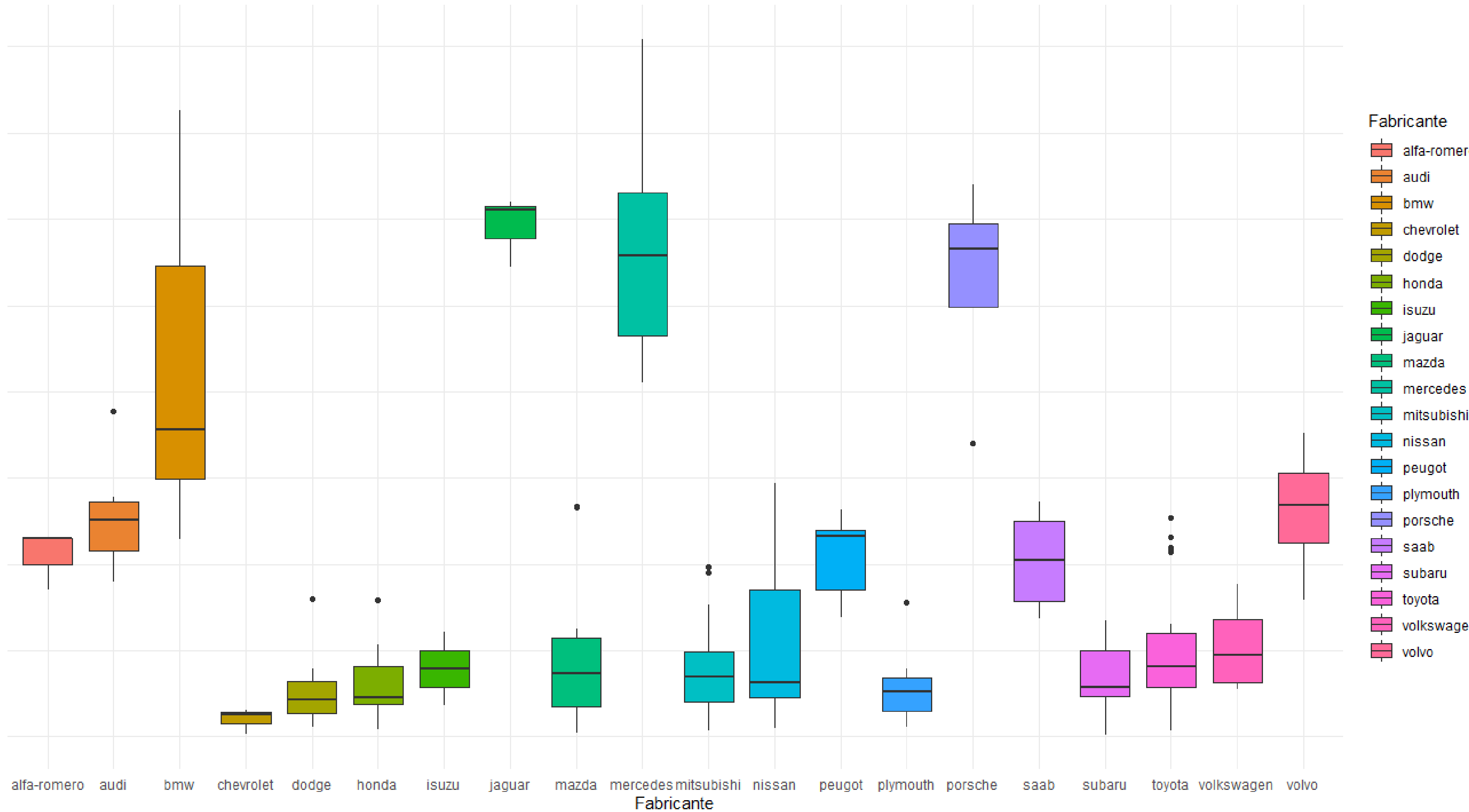
Preço do carro por Sistema de Combustível



Preço do carro por taxa de compressão



Preço do carro por fabricante





Modelos simples

$\text{price} \sim \text{length} + \text{length}^2;$

$\text{price} \sim \text{width} + \text{width}^2;$

$\text{price} \sim \text{curb_weight} + \text{curb_weight}^2;$

$\text{price} \sim \text{engine_size};$

$\text{price} \sim \text{horsepower};$

$\text{price} \sim \text{city_mpg} + \text{city_mpg}^2 + \text{city_mpg}^3 + \text{city_mpg}^4;$

$\text{price} \sim \text{highway_mpg} + \text{highway_mpg}^2 + \text{highway_mpg}^3.$

Modelos múltiplos

Ajuste de dois modelos múltiplos cheios para seleção de variáveis.

Consideraram-se: length, width, curb_weight, engine_size, horsepower, city_mpg, highway_mpg, compression_ratio, num_cylinders, engine_type, engine_location, drive_wheels, body_style, aspiration, make para o modelo 1.

Modelos múltiplos

MODELO 2

Igual ao modelo 1, sem a variável make (muito correlacionada com as outras segundo correlação de Crammer)

STEPWISE

Seleção de variáveis utilizando AIC

Melhores modelos



NORMALIDADE

Os dois modelos tiveram alguns problemas com normalidade, por conta de uns poucos pontos.



HOMOCEDASTICIDADE

Ambos pareciam homocedásticos e passaram no Goldfeld.Quandt. No entanto, falharam no Breush-Pagan.



MEDIDAS DE INFLUÊNCIA

Alguns pontos são extremamente influentes. Número razoável de pontos que talvez não tenham sido bem ajustados.

Melhores modelos



VARIAÇÕES

Os pontos influentes mudam a depender do modelo.



COMPARAÇÕES

Os Modelo 1 (após o STEPWISE) era o que tinha o menor AIC.



MULTICOLINEARIDADE

Uma imagem vale mais que mil palavras:

Modelos múltiplos

```
> aux2[5]
[[1]]
[[1]]$VIFs
```

	GVIF	Df
poly(width, degree = 2, raw = F)	3.693301e+01	2
poly(curb_weight, degree = 2, raw = F)	2.697714e+02	2
poly(engine_size, degree = 1, raw = F)	2.040108e+01	1
poly(city_mpg, degree = 4, raw = F)	6.152620e+01	4
engine_type	3.837155e+03	4
drive_wheels	9.817384e+00	2
body_style	6.002099e+00	4
aspiration	1.779449e+00	1
make	1.492038e+06	19

MODELO 1

Parece que make é realmente bem correlacionada com as outras;

```
[[1]]
[[1]]$VIFs
```

	GVIF	Df
poly(length, degree = 2, raw = F)	41.142117	2
poly(width, degree = 2, raw = F)	17.628598	2
poly(curb_weight, degree = 2, raw = F)	118.236216	2
poly(horsepower, degree = 1, raw = F)	14.379579	1
poly(city_mpg, degree = 4, raw = F)	14024.077264	4
poly(highway_mpg, degree = 3, raw = F)	9929.354454	3
compression_ratio	3.284065	1
num_cylinders	17.411708	2
engine_type	13.793931	4
engine_location	2.868308	1
drive_wheels	4.820219	2
body_style	4.886748	4

MODELO 2

O modelo 2 selecionou algumas variáveis muito correlacionadas entre si, sendo que algumas delas foram descartadas pelo Modelo 1.

Considerações

A FAZER

Verificar relações entre variáveis
categóricas e quantitativas
explicativas;
Possivelmente adicionar termos de
interação;

MAIS UMA MEDIDA

Utilizar outras medidas para
seleção de variáveis, como BIC e
o R^2 , verificar performance de
modelos sem algumas variáveis
que sabemos causar
multicolinearidade;

POR ENQUANTO É SÓ.

Obrigado pela atenção!
Dúvidas ou sugestões?