

# Análise e Modelagem de Preços de Carros Usados

Gustavo Almeida Silva

Universidade Federal de Juiz de Fora

07/07/2023

# O Mercado de Carros Usados

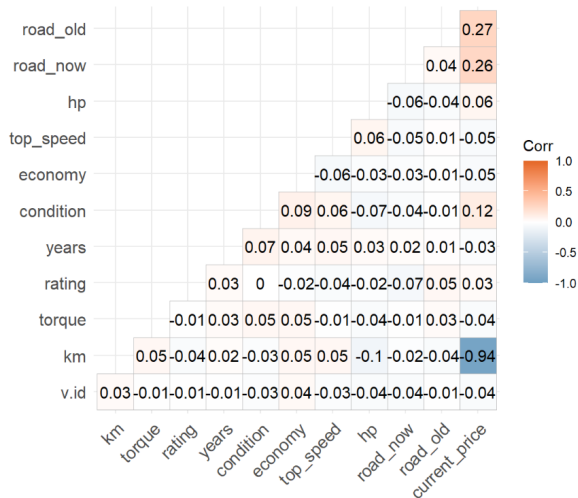
- Acessibilidade
- Sustentabilidade

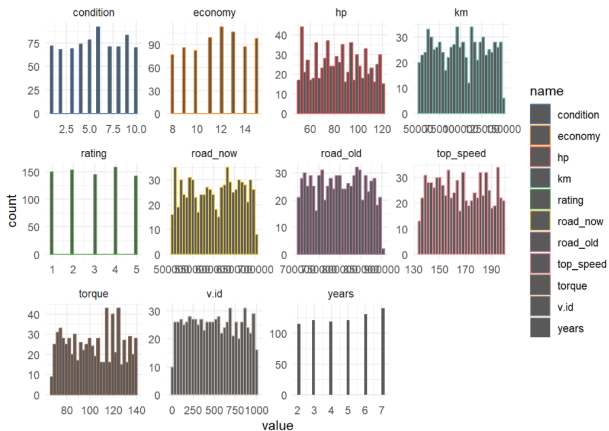
## 12 variáveis e 1000 observações

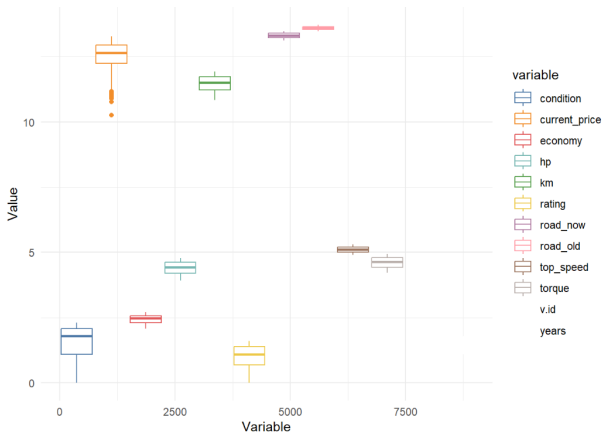
- **v.id**
- **road\_old**
- **road\_now**
- **years**
- **km**
- **rating**
- **condition**
- **economy**
- **top\_speed**
- **hp**
- **torque**
- **current\_price**

# Separação e Análise Exploratória dos Dados

- Divisão Treino e Teste (0.75, 0.25)
- Ausência de Valores Faltantes







# Modelos

- Saturado
- Correlacional

	Saturado	Correlacional
(Intercept)	-14651.9*	20605.4*
	p = <0.1	p = <0.1
road_now	0.5***	0.5***
	p = <0.1	p = <0.1
road_old	0.5***	0.5***
	p = <0.1	p = <0.1
years	-1595.4***	-1063.3**
	p = <0.1	p = <0.1
km	-4.0***	-4.0***
	p = <0.1	p = <0.1
rating	285.8	
	p = 0.2	
condition	4629.4***	
	p = <0.1	
economy	41.7	
	p = 0.8	
top_speed	-7.9	

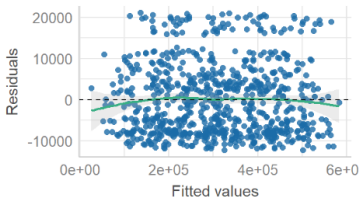
Num.Obs.	748	748
R2	0.995	0.985
R2 Adj.	0.995	0.984
AIC	15720.9	16585.0
BIC	15776.3	16612.7
Log.Lik.	-7848.427	-8286.511
RMSE	8722.06	15666.58



# Resíduos

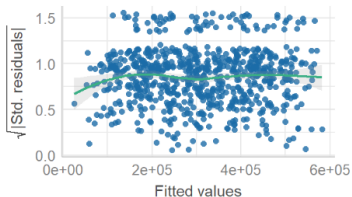
## Linearity

Reference line should be flat and horizontal



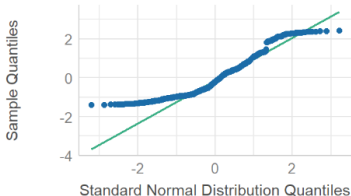
## Homogeneity of Variance

Reference line should be flat and horizontal



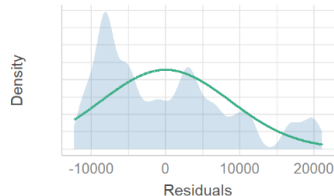
## Normality of Residuals

Dots should fall along the line



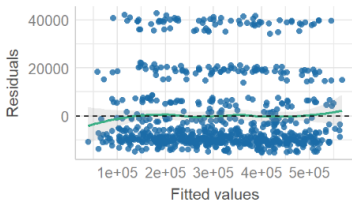
## Normality of Residuals

Distribution should be close to the normal curve



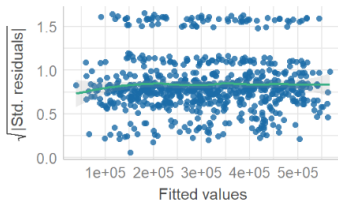
### Linearity

Reference line should be flat and horizontal



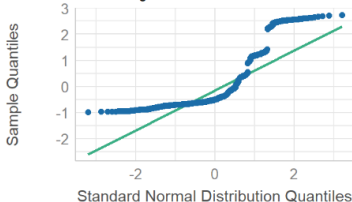
### Homogeneity of Variance

Reference line should be flat and horizontal



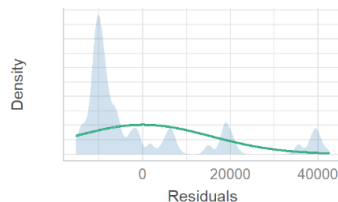
### Normality of Residuals

Dots should fall along the line



### Normality of Residuals

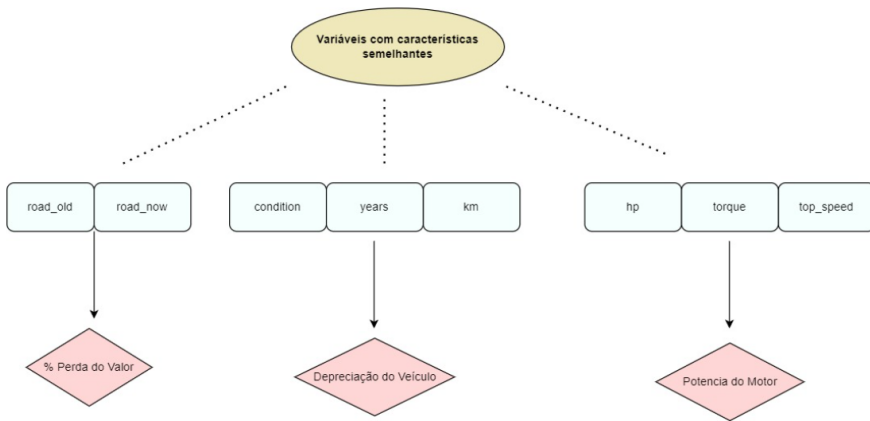
Distribution should be close to the normal curve



Teste	Modelo.Sat.P.valor	Modelo.Corr.P.valor
Shapiro	0	0
Cramer	0	0
Lilliefors	0	0

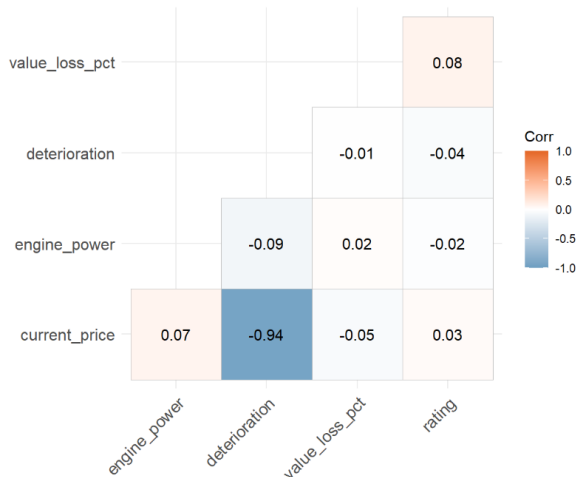
- Rejeição da Hipótese de Normalidade
- Heterogeneidade da Variância

- Criação de novas variáveis via transformação
- Tentativa de melhorar o modelo



## Transformações utilizadas

- **% Perda de Valor** =  $\frac{(ValorInicial - ValorAtual) \times 100}{ValorInicial}$
- **Deterioração do Veículo** =  $Media(condition, years, km)$
- **Potencia do Motor** =  $Media(top\ speed^2, hp^5, torque)$



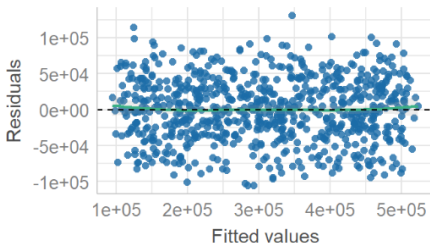
- Modelo:  $current\_price = \beta_0 + \beta_1 value\_loss\_pct + \beta_2 deterioration + \beta_3 engine\_power + E$



Features	
(Intercept)	739155.0***
	p = <0.1
value_loss_pct	-771.7***
	p = <0.1
deterioration	-12.3***
	p = <0.1
engine_power	0.0
	p = 0.1
Num.Obs.	748
R2	0.880
R2 Adj.	0.879
AIC	18115.6
BIC	18138.7
Log.Lik.	-9052.782
RMSE	43639.20

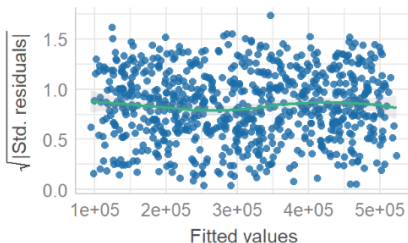
### Linearity

Reference line should be flat and horizontal



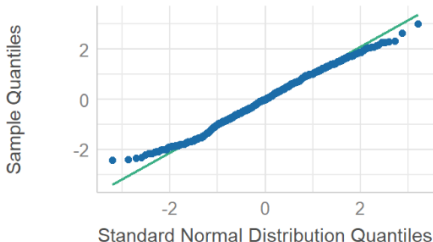
### Homogeneity of Variance

Reference line should be flat and horizontal



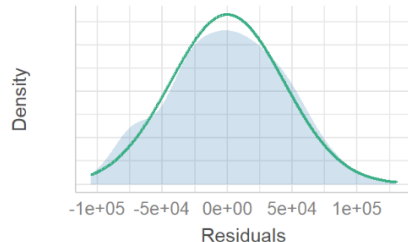
### Normality of Residuals

Dots should fall along the line

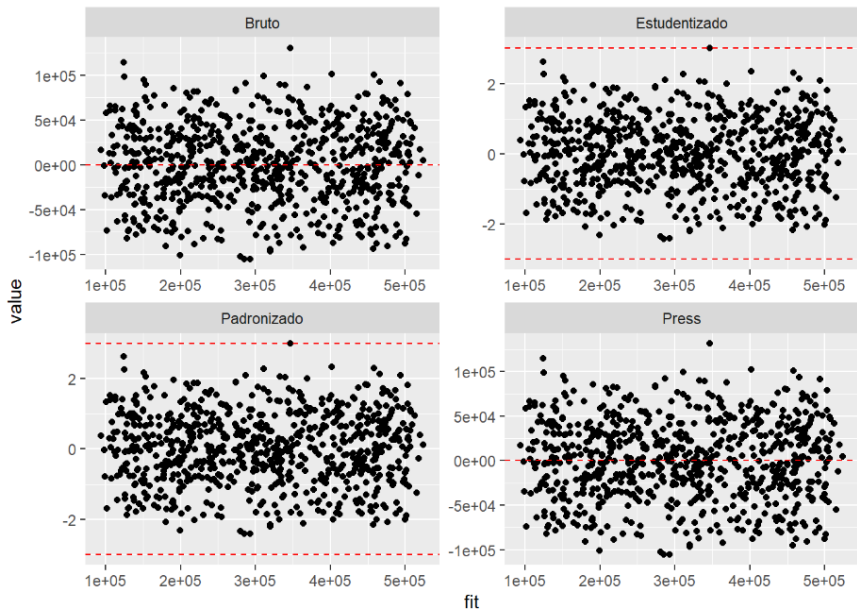


### Normality of Residuals

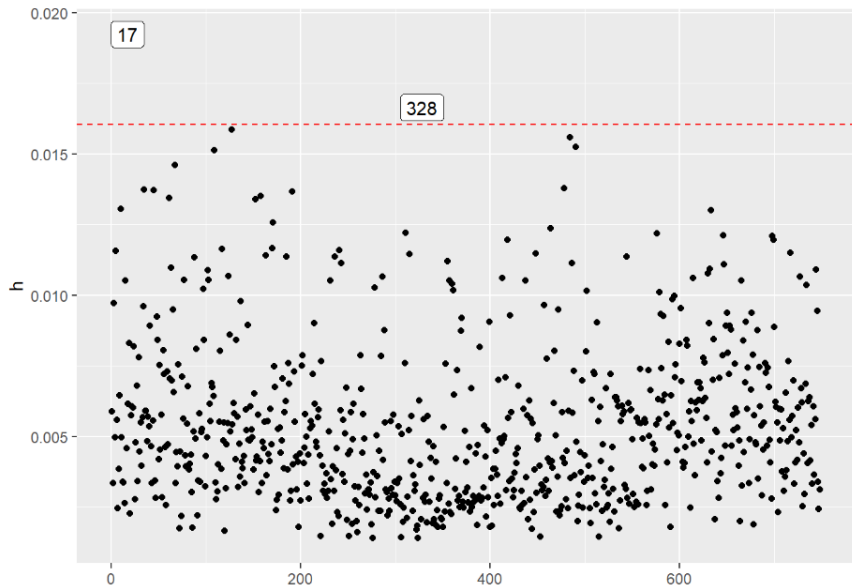
Distribution should be close to the normal curve

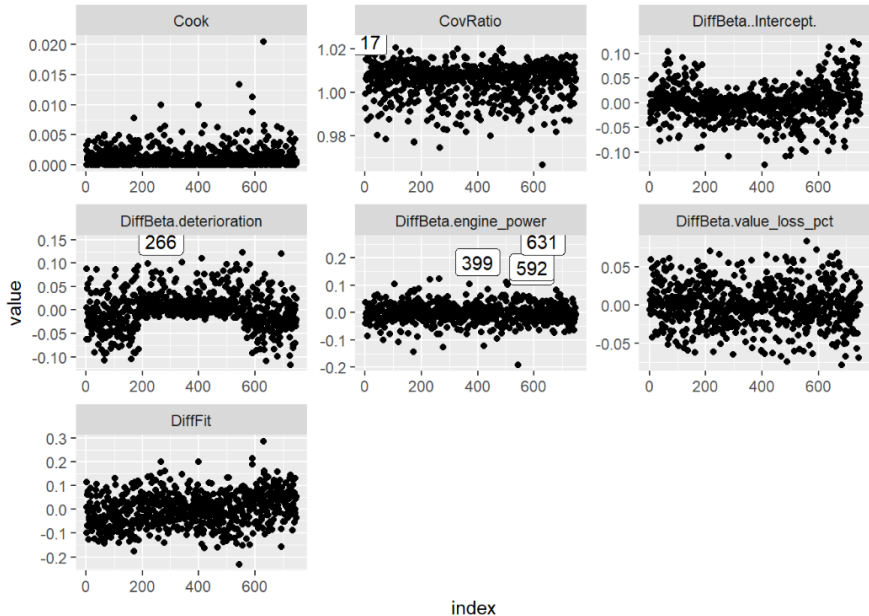


Teste	Modelo.Feat.P.valor
Shapiro	0.2048031
Cramer	0.3276374
Lilliefors	0.4188786



# Análise de Diagnósticos





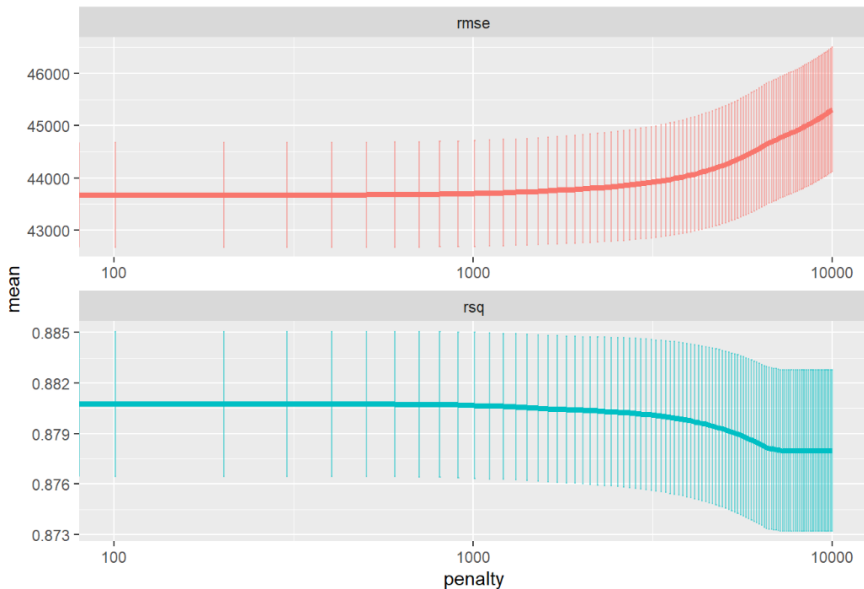
## Multicolinearidade

	VIF_Values
value_loss_pct	1.000291
deterioration	1.008633
engine_power	1.008797

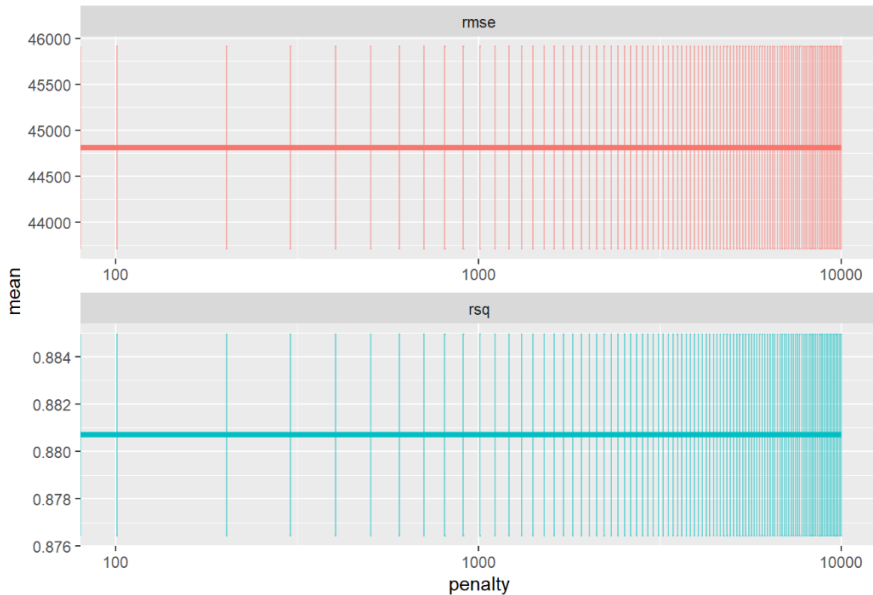
- Busca penalizar variáveis colineares
- Modelo não apresenta multicolinearidade, e portanto não espera-se melhora
- K-fold e Leave One Out para robustez de métricas
- **Lasso (L1)**
- **Ridge (L2)**
- **Elastic Net (Mistura 0.7 de penalty L1 e 0.3 de Penalty L2)**



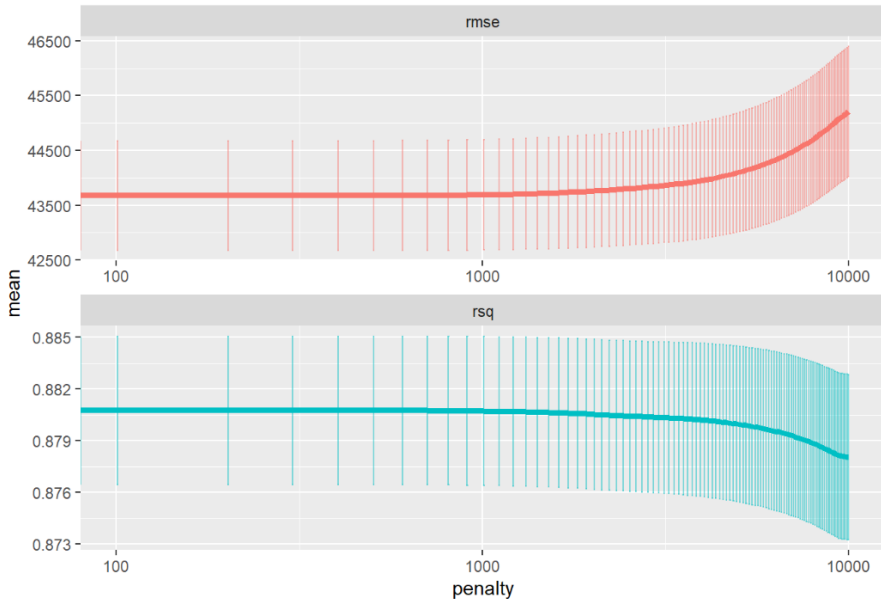
## Lasso Tuning



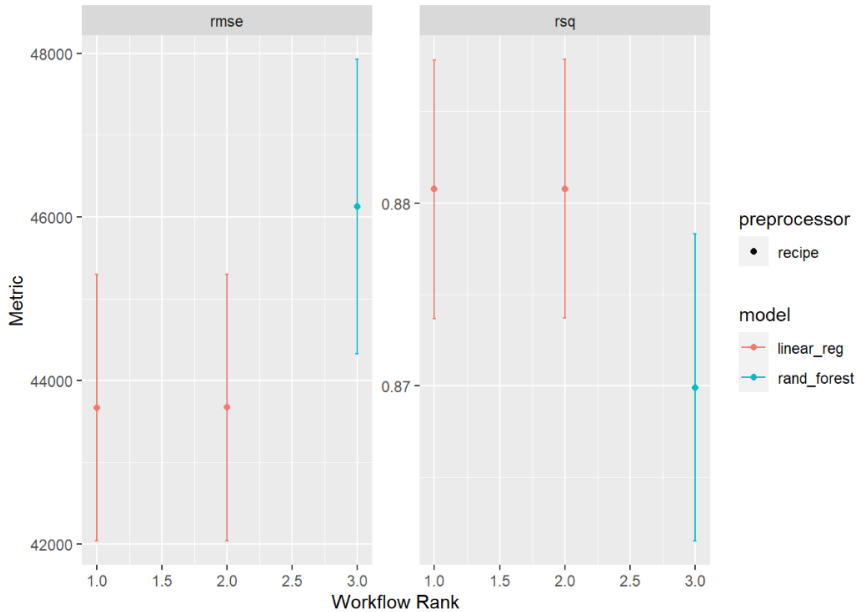
## Ridge Tuning



## Elastic Net Tuning



- Assuntos Comentados em Aula
- Modelo Bayesiano (família Gaussiana sem fixação de hiperparâmetros)
- Aprendizado de Máquina (Random Forest sem tuning de hiperparâmetros)
- K-fold e Leave One Out para robustez de métricas



- Métricas iguais e piores do que o modelo linear simples
- Complexidade de Interpretação

- Melhor Modelo:  $current\_price = \beta_0 + \beta_1 value\_loss\_pct + \beta_2 deterioration + \beta_3 engine\_power + E$

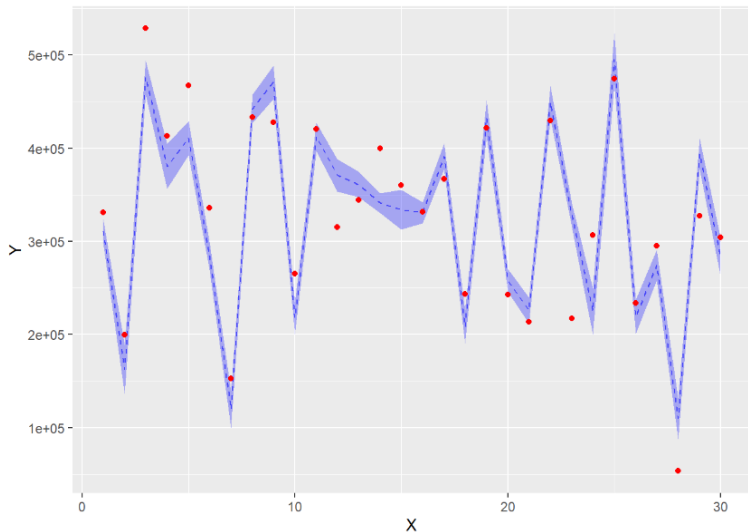
Features	
(Intercept)	739155.0***
	p = <0.1
value_loss_pct	-771.7***
	p = <0.1
deterioration	-12.3***
	p = <0.1
engine_power	0.0
	p = 0.1
Num.Obs.	748
R2	0.880
R2 Adj.	0.879
AIC	18115.6
BIC	18138.7
Log.Lik.	-9052.782
RMSE	43639.20



## Validação no Conjunto de Teste

Metric	Data_test	Data_Training
rmse	4.438387e+04	4.363920e+04
rsq	8.773829e-01	8.797584e-01

## Valores Preditos



- O trabalho utilizou técnicas de regressão linear, como modelos saturados e correlacionais, engenharia de características, análise de diagnósticos e análise de resíduos
- Os resíduos do modelo não rejeitaram a hipótese de normalidade, o que indica um bom ajuste.
- A engenharia de características permitiu criar um modelo com baixa multicolinearidade, conforme indicado pelos baixos valores de VIF
- A aplicação de regularização (L1, L2 e Elastic Net) não trouxe melhorias significativas ao modelo.
- Comparando com outros modelos, como os Bayesianos e Random Forest, o modelo de regressão linear simples obteve resultados melhores e uma boa capacidade de generalização.