

# APLICAÇÃO DE TÉCNICAS DE MINERAÇÃO DE DADOS PARA IDENTIFICAR ASPECTOS QUE INFLUENCIAM NO PROCESSO DE VENDA E ENTREGA DOS PRODUTOS DE UM E-COMMERCE

Acadêmico: Herikc Brecher

Orientador: Rafael Ballottin Martins





# Sumário

1. Contextualização
2. Objetivos
3. Projeto
4. Análise Exploratória
5. Processamento dos dados
6. Clusterização
7. Séries Temporais
8. Conclusão
9. Trabalhos Futuros



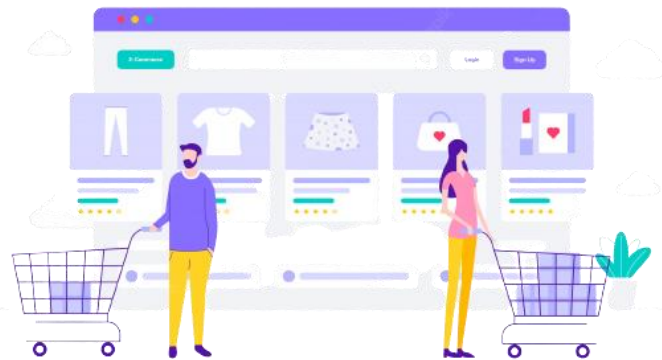
# 1. Contextualização

- ◎ Processo de KDD:
  - Analise Exploratória;
  - Pré-Processamento;
  - Mineração de Dados;
  - Avaliação dos resultados.



## 1.1 Problema de Pesquisa

- ◎ ↑ Geração de dados;
- ◎ ↑ Complexidade dos dados;
- ◎ Dados analisados superficialmente;
- ◎ Necessidade de analisar os dados.





## 1.2 Solução

- ◎ Base de dados publica;
- ◎ Processo de KDD:
  - Analise de dados;
  - Pré-Processamento;
  - Mineração de dados;
  - Interpretação e avaliação dos resultados.





## 2. Objetivo Geral

---

Identificar aspectos que influenciam no processo de avaliação dos pedidos, desempenho das vendas e entrega.



## 2.1 Objetivos Específicos

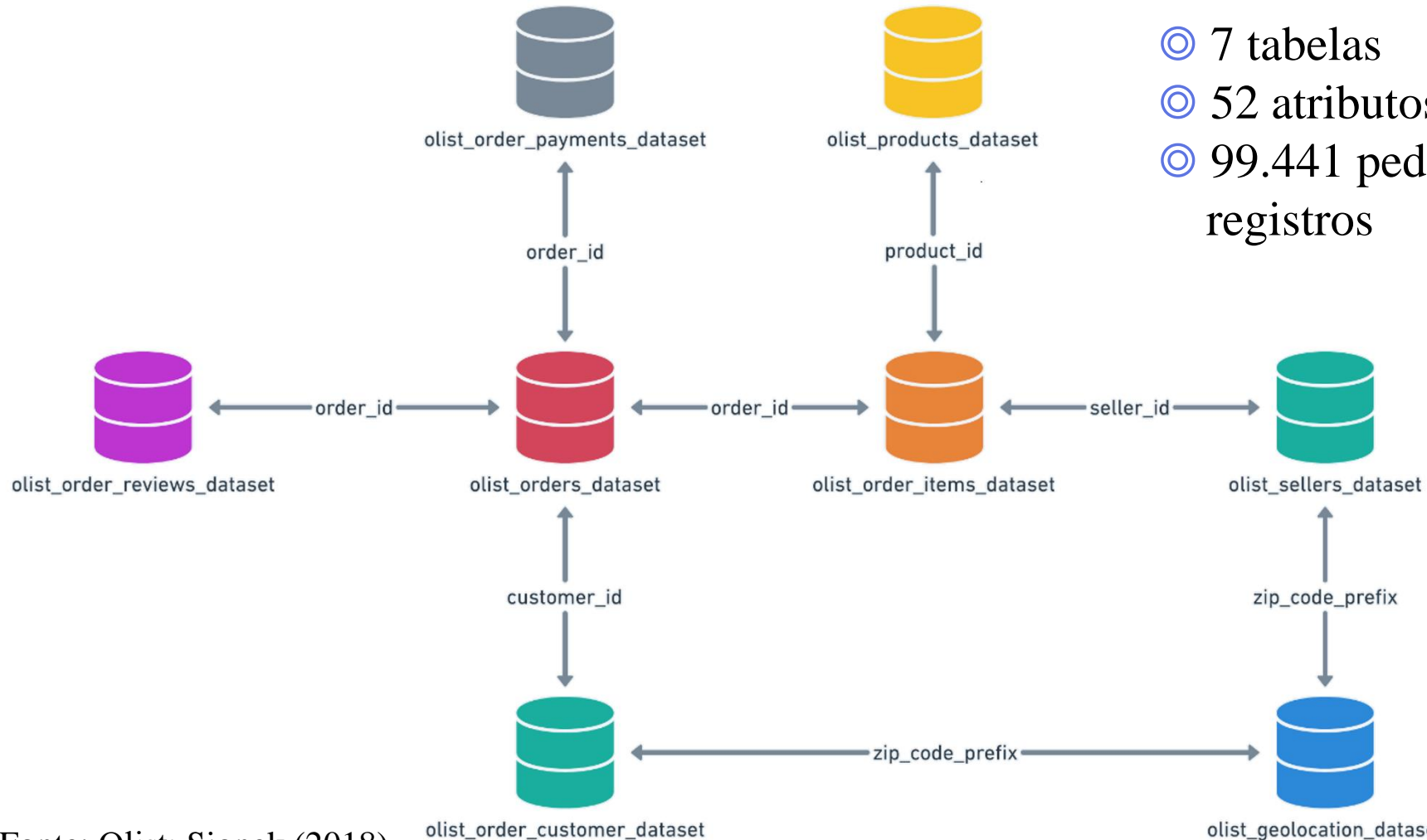
- ⦿ Definir os atributos mais relevantes através de uma análise exploratória dos dados;
- ⦿ Identificar os principais aspectos que influenciam a avaliação de pedidos, desempenho das vendas e entrega da empresa;
- ⦿ Identificar e definir algoritmos para regressão e clusterização;
- ⦿ Avaliar os resultados a fim de descobrir o nível de confiança das informações obtidas.

# PROJETO





- 7 tabelas
- 52 atributos
- 99.441 pedidos / registros





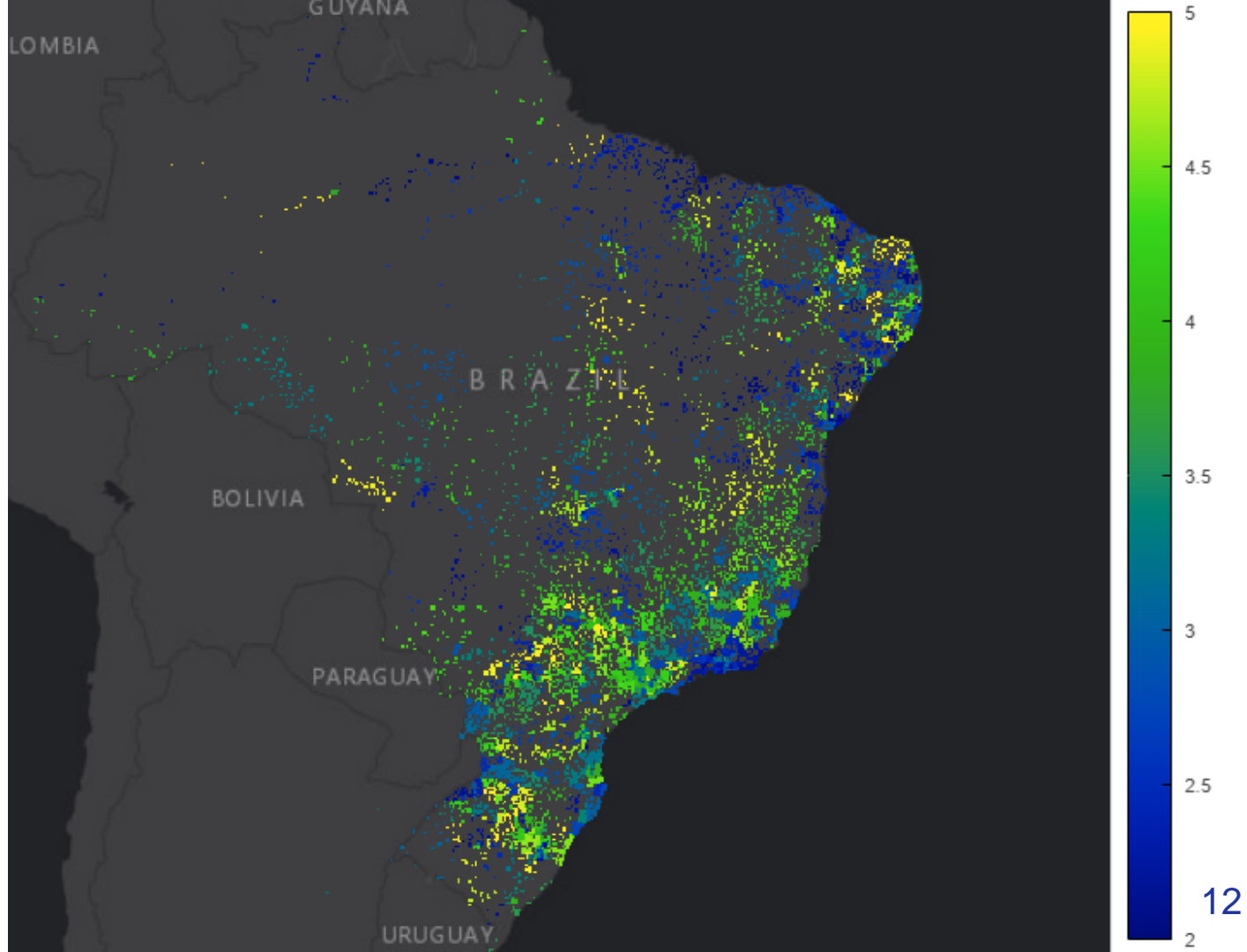
## 3. Projeto

Atributo	Valores Únicos
<b>Clientes</b>	96.096
<b>Pedidos</b>	99.441
<b>Avaliações</b>	98.410
<b>Lojas</b>	3.095
<b>Produtos</b>	71

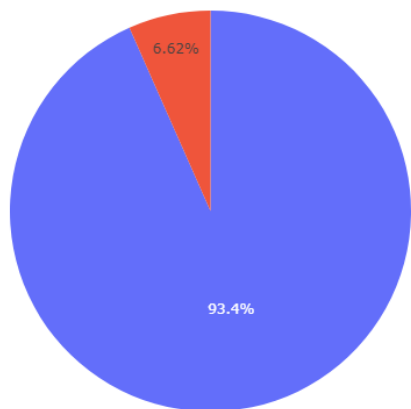
# ANALISE EXPLORATÓRIA



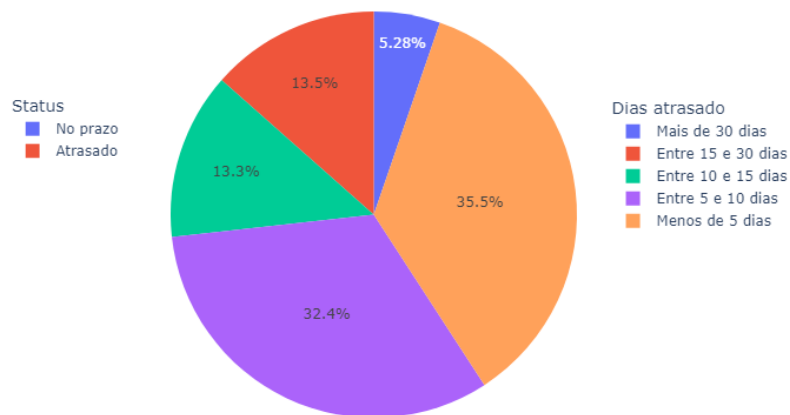
# Review médio por região no Brasil



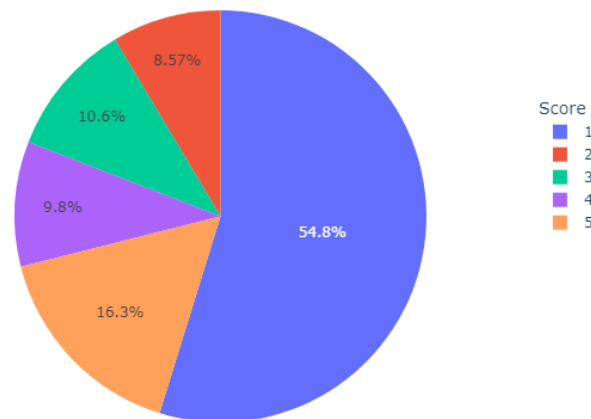
% de Pedidos Entregues no Período



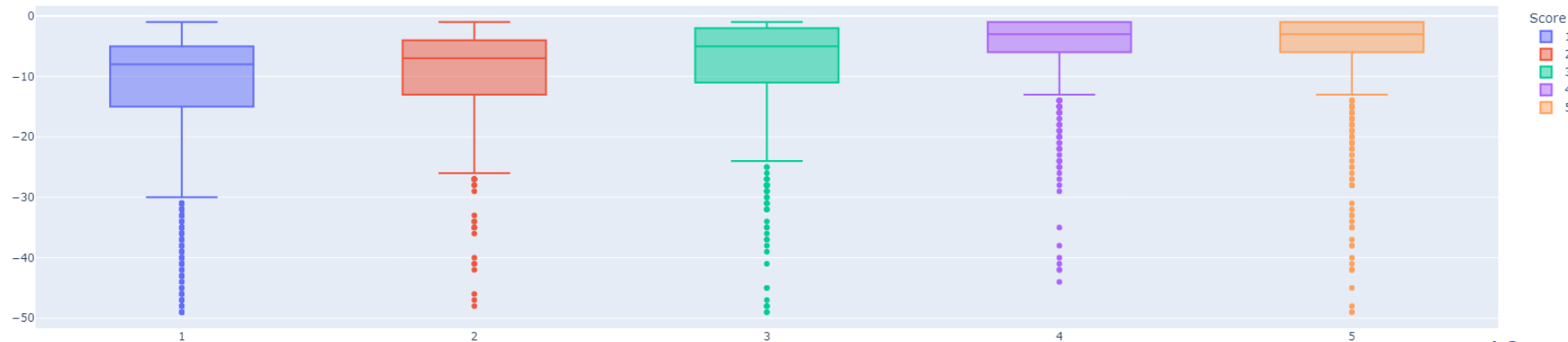
% Por Categoria de Dias de Atraso



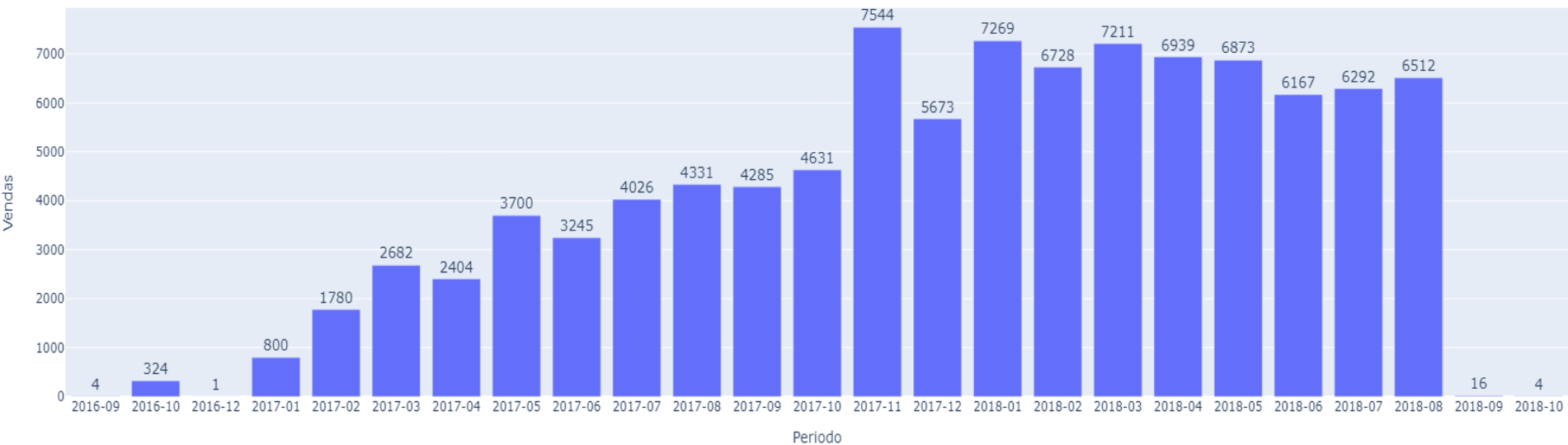
% de Pedidos Atrasados por Score



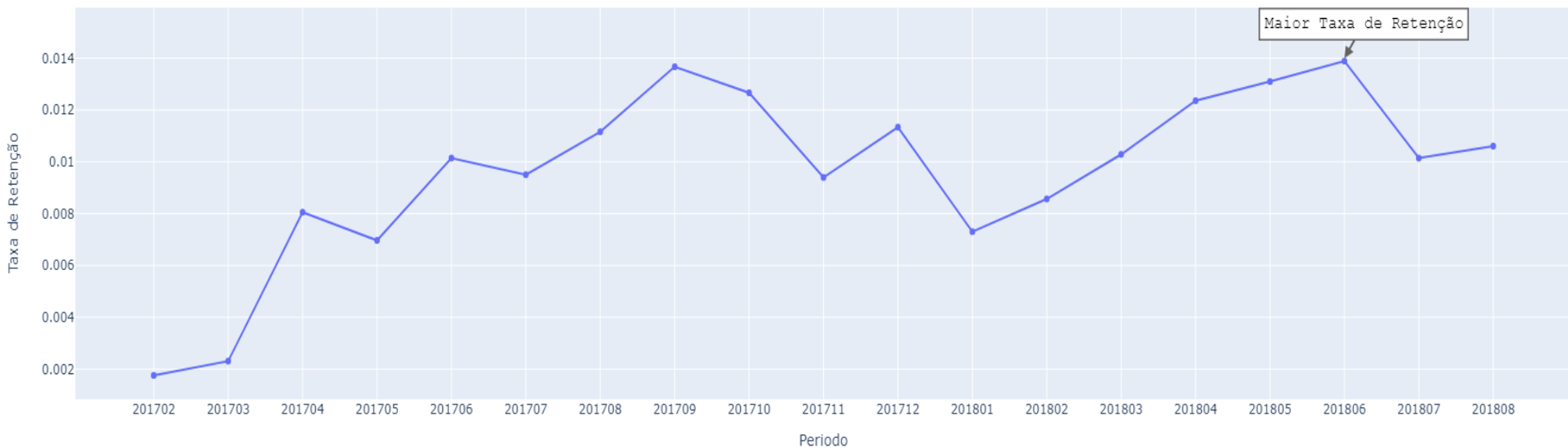
Dias de Atraso por Score



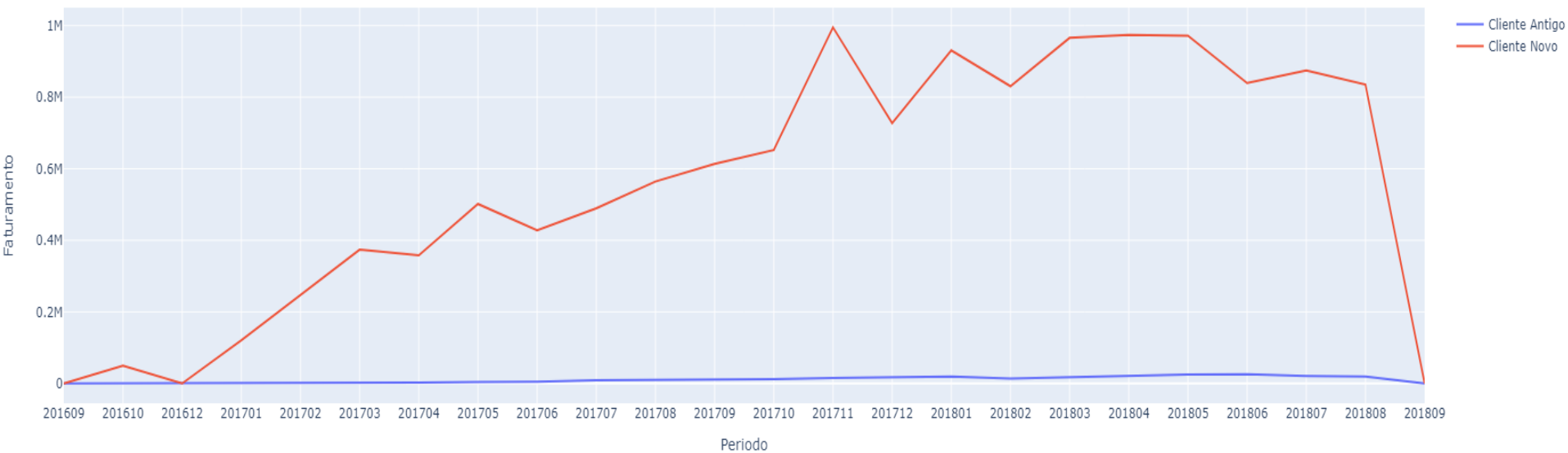
## Clientes Ativos Por Mês



## Taxa Mensal de Retenção de Clientes



Diferença de Faturamento ao Longo do Tempo Entre Clientes Novos e Antigos





# PROCESSAMENTO DOS DADOS



Teste de hipótese por meio de Grubs para detecção de outliers.

Atributo	Possui Outliers
freight_value	Sim
payment_value	Sim
quantity	Sim
diff_delivery_and_estimate	Sim
diff_delivery_and_purchase	Sim

# Aplicação de *Label Encoder*.

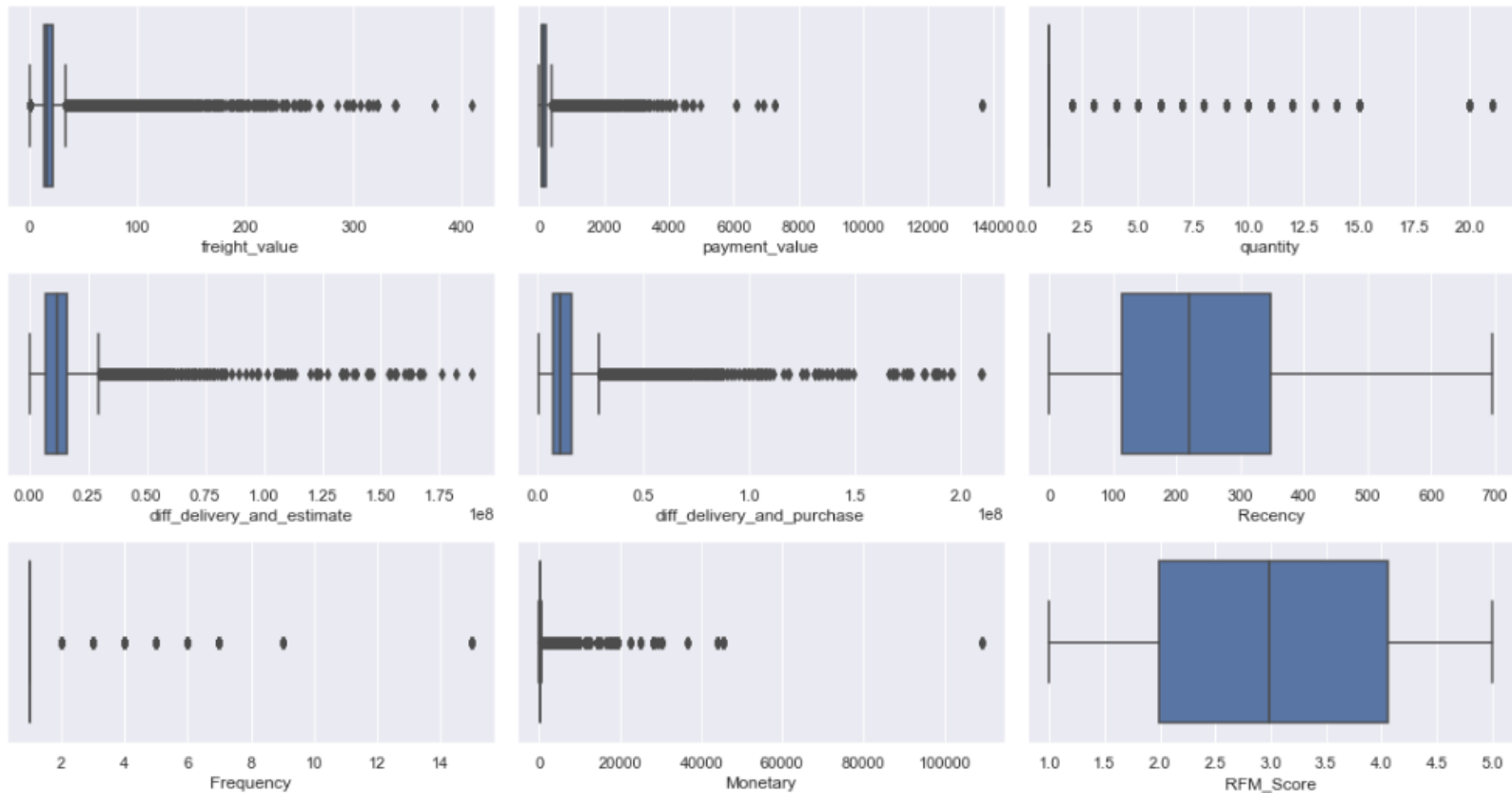
Antes da aplicação:

freight_value	payment_type	payment_installments	payment_value	review_score	quantity	diff_delivery_and_estimate	diff_delivery_and_purchase	customer_state	customer_region
13.29	credit_card	2	72.19	5	1	9.0	7.0	RJ	Sudeste
19.93	credit_card	3	259.83	4	1	3.0	16.0	SP	Sudeste
17.87	credit_card	5	216.87	5	1	14.0	8.0	MG	Sudeste
12.79	credit_card	2	25.78	4	1	6.0	6.0	SP	Sudeste
18.14	credit_card	3	218.04	5	1	16.0	25.0	SP	Sudeste

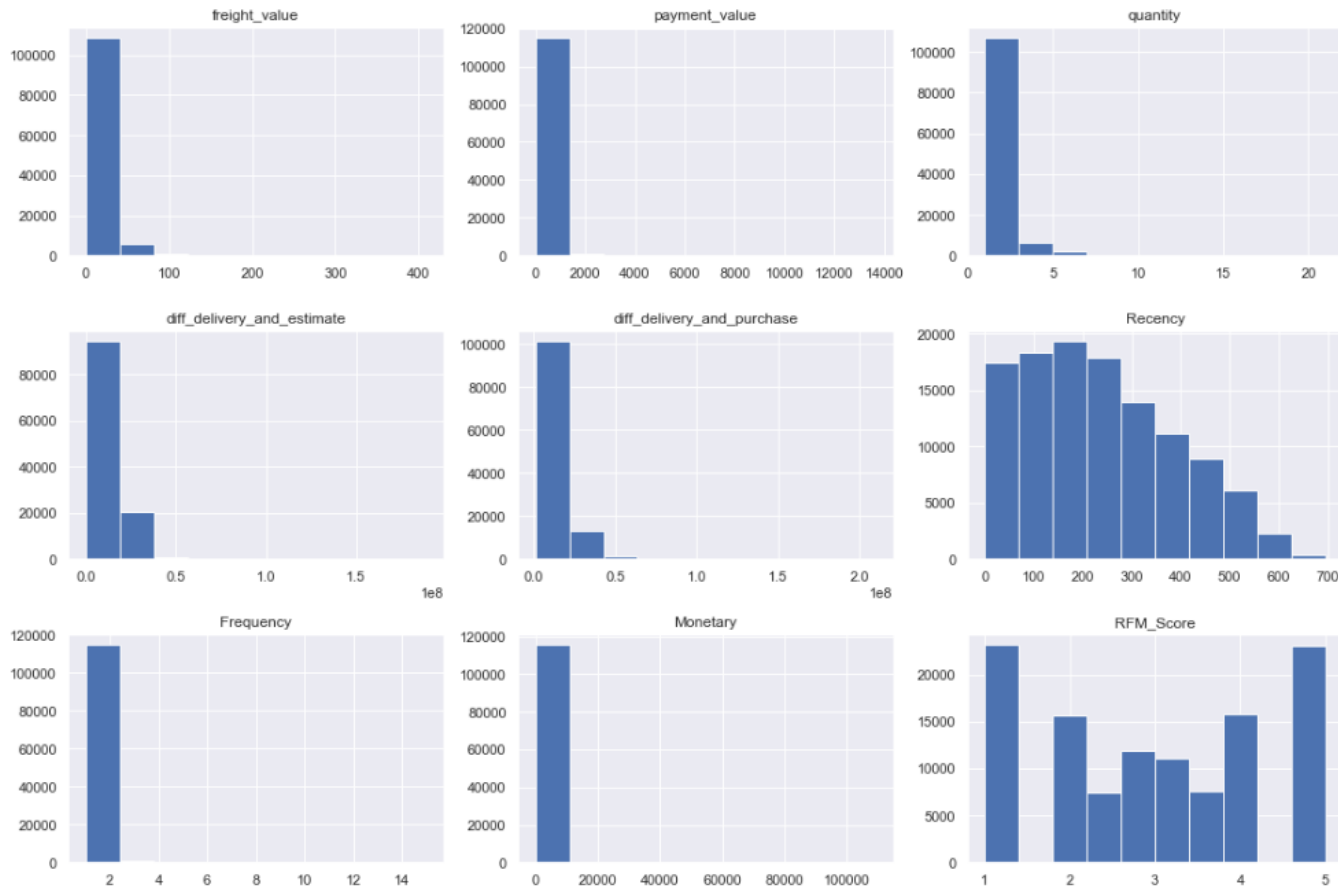
Depois da aplicação:

freight_value	payment_type	payment_installments	payment_value	review_score	quantity	diff_delivery_and_estimate	diff_delivery_and_purchase	customer_state	customer_region
13.29	1	2	72.19	5	1	9.0	7.0	18	3
19.93	1	3	259.83	4	1	3.0	16.0	25	3
17.87	1	5	216.87	5	1	14.0	8.0	10	3
12.79	1	2	25.78	4	1	6.0	6.0	25	3
18.14	1	3	218.04	5	1	16.0	25.0	25	3

# Boxplot dos atributos quantitativos:



# Histograma dos atributos quantitativos:



## Simetria dos atributos quantitativos:

Variável	Skewness	Kurtosis
freight_value	5.558939	58.300861
payment_value	14.448540	530.034479
quantity	6.331864	68.366469
diff_delivery_and_estimate	2.846043	32.528515
diff_delivery_and_purchase	3.912217	40.794495
Recency	0.447031	-0.656545
Frequency	11.008832	258.347088
Monetary	29.531411	1432.134347
RFM_Score	-0.003545	-1.271986

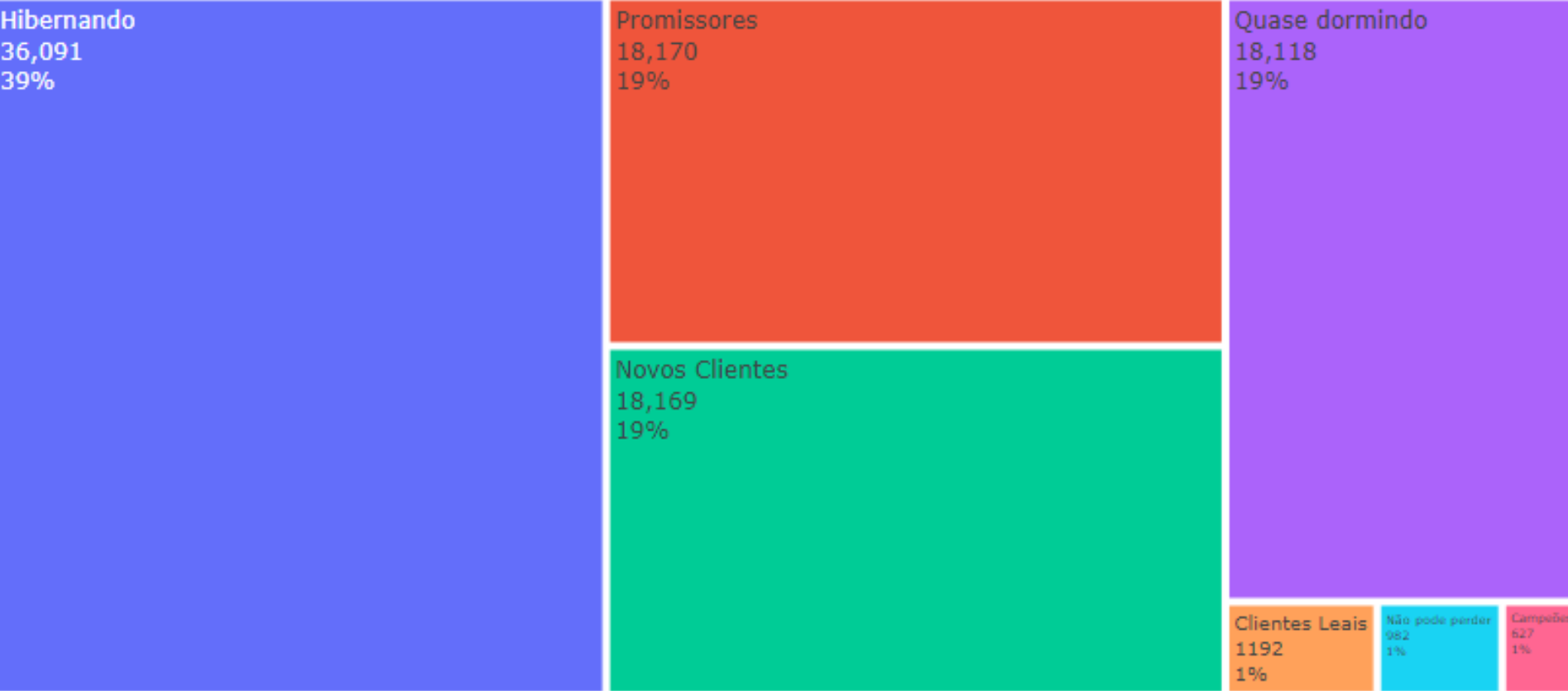
# Clusterização



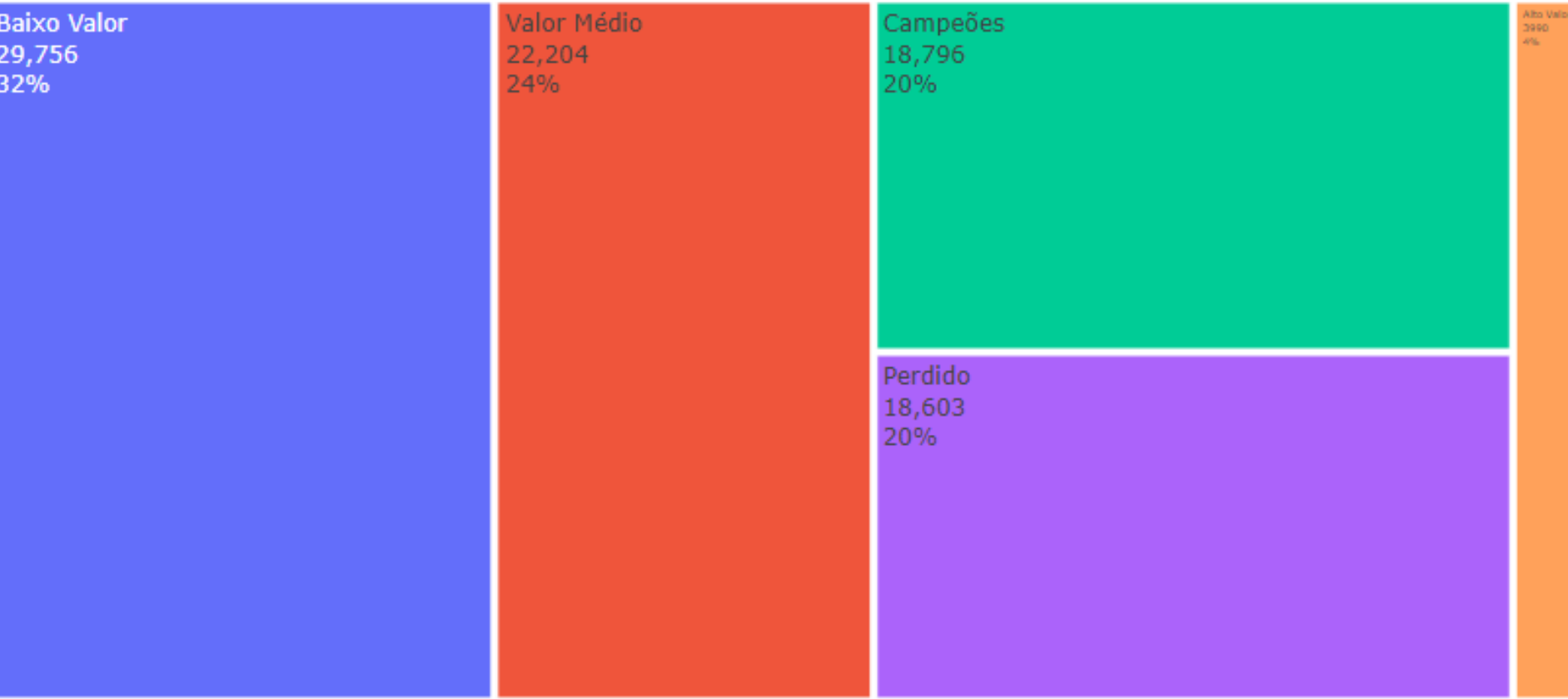
Categoria	Descrição
Recência	Quão recentemente o cliente fez uma transação. Por exemplo, o cliente realizou a sua última compra a 30 dias.
Frequência	Quão frequente o cliente realizar um pedido. Por exemplo, nos últimos 12 meses foram realizadas 4 compras.
Monetário	Quanto o cliente gastou ao total nos pedidos.

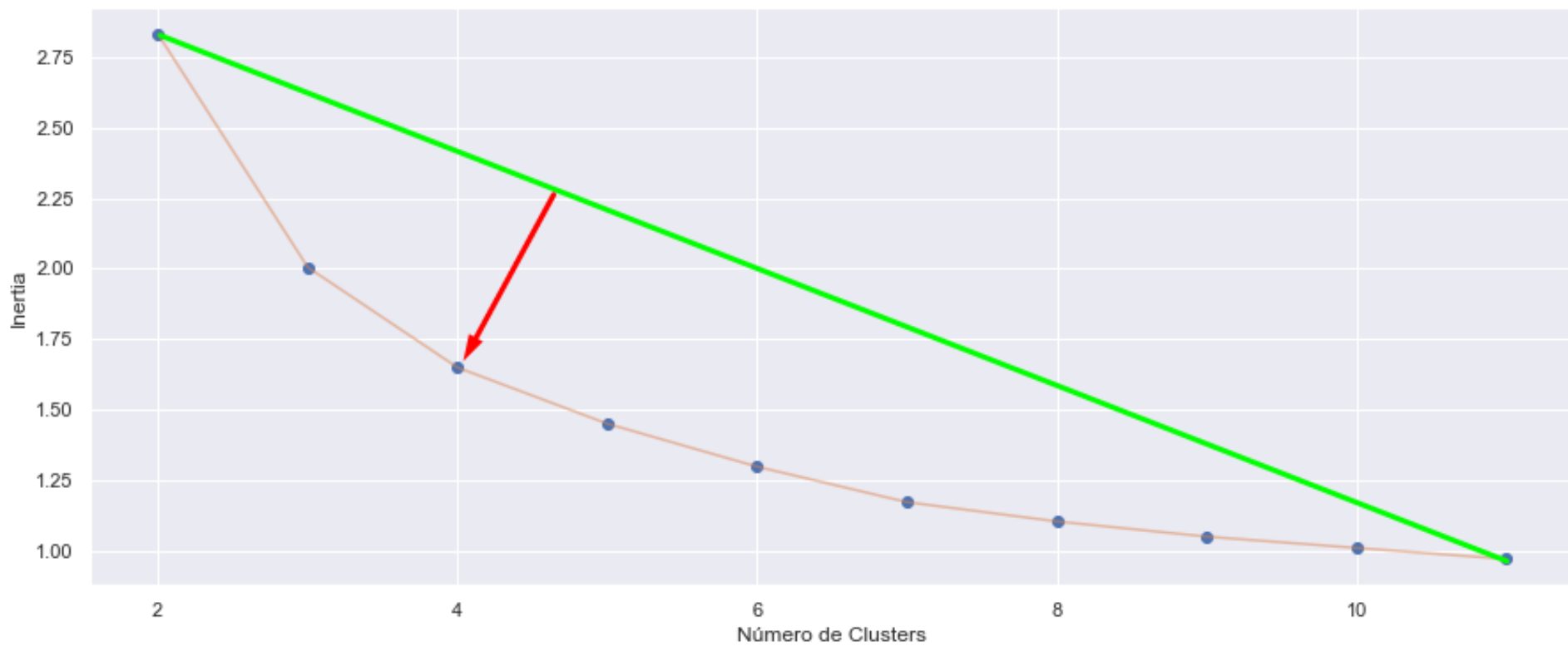


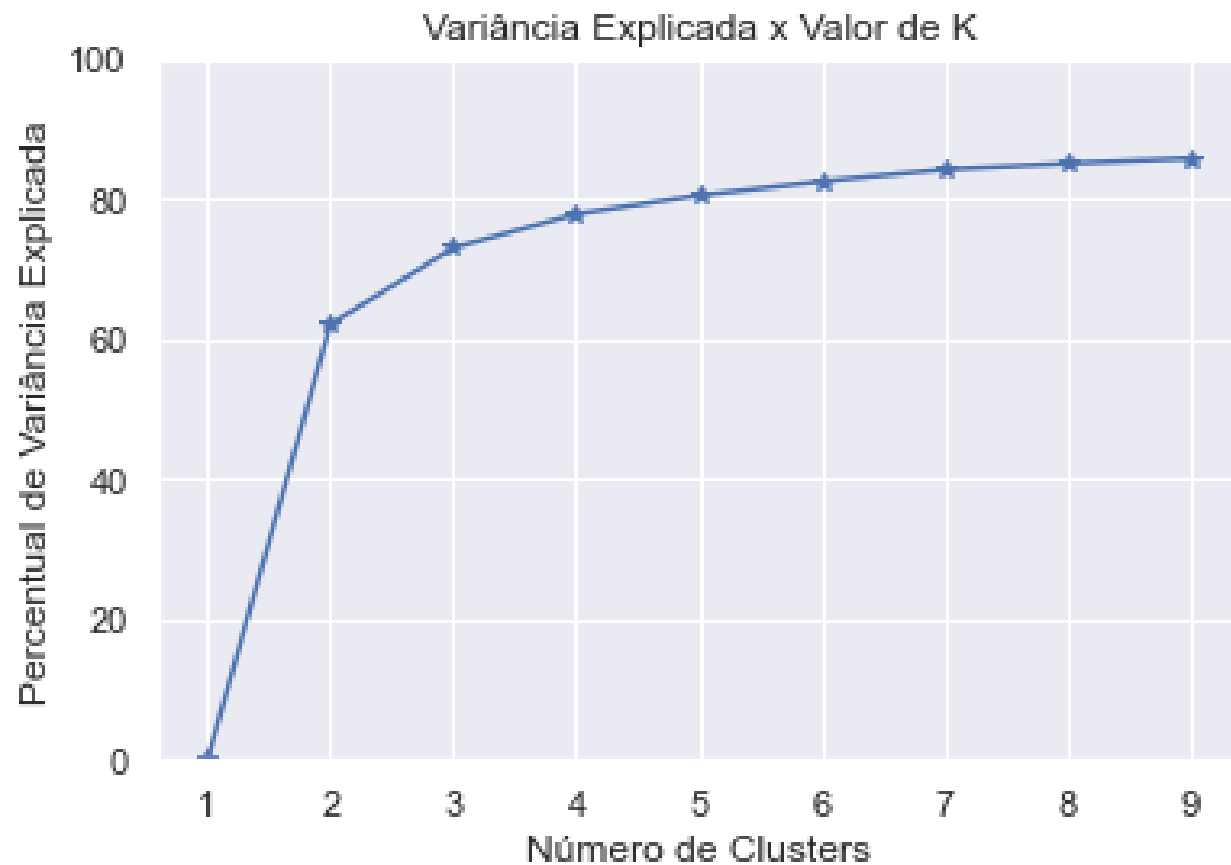
# Análise de RFM por Recência e Frequência



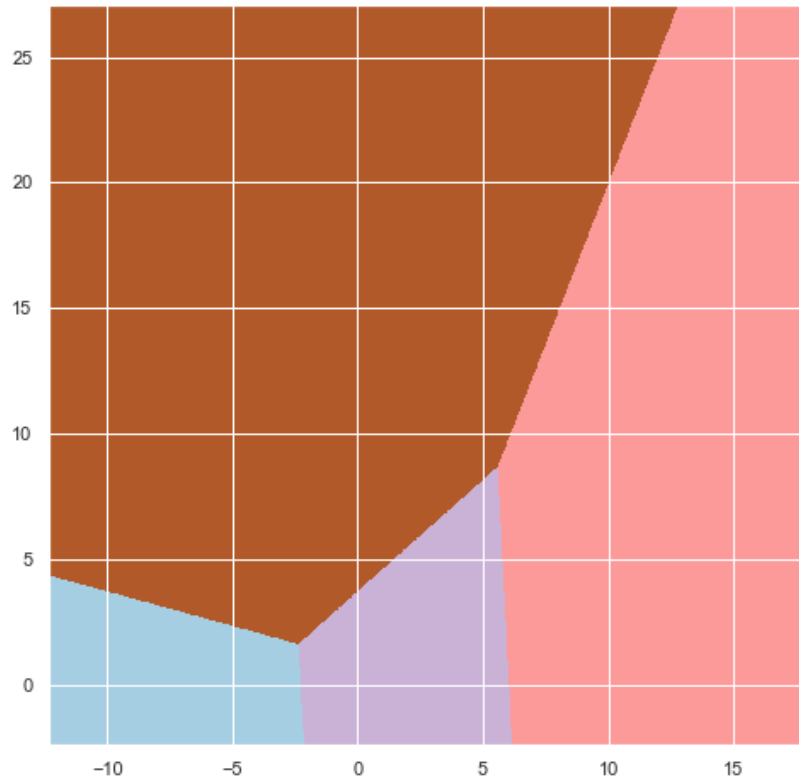
# Análise de RFM por Recência, Frequência e Valor Monetário



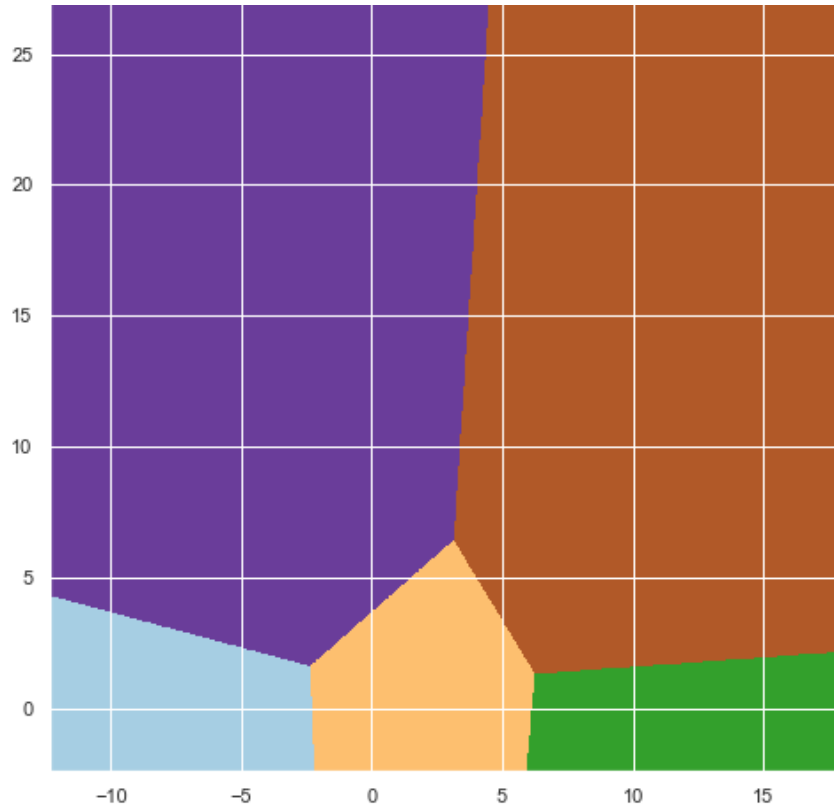




Meshgrid de 4 clusters com PCA



Meshgrid de 5 clusters com PCA

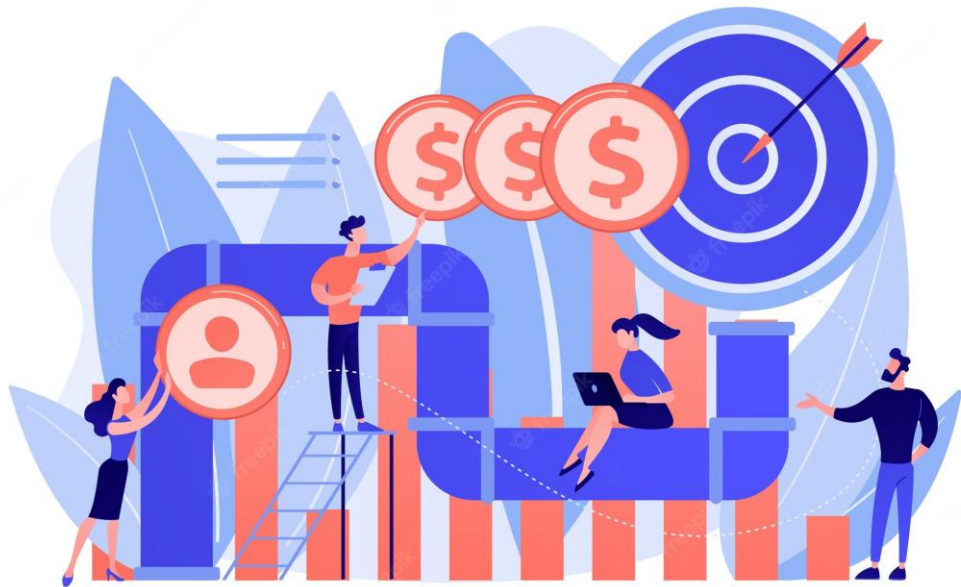


<b>K-Clusters</b>	<b>Silhouette Score com PCA</b>	<b>Silhouette Score sem PCA</b>
<b>4</b>	0.6094	0.4456
<b>5</b>	0.6025	0.4308
<b>6</b>	0.6114	0.4057
<b>8</b>	0.5799	0.3219

# Analise de RFM por Cluster



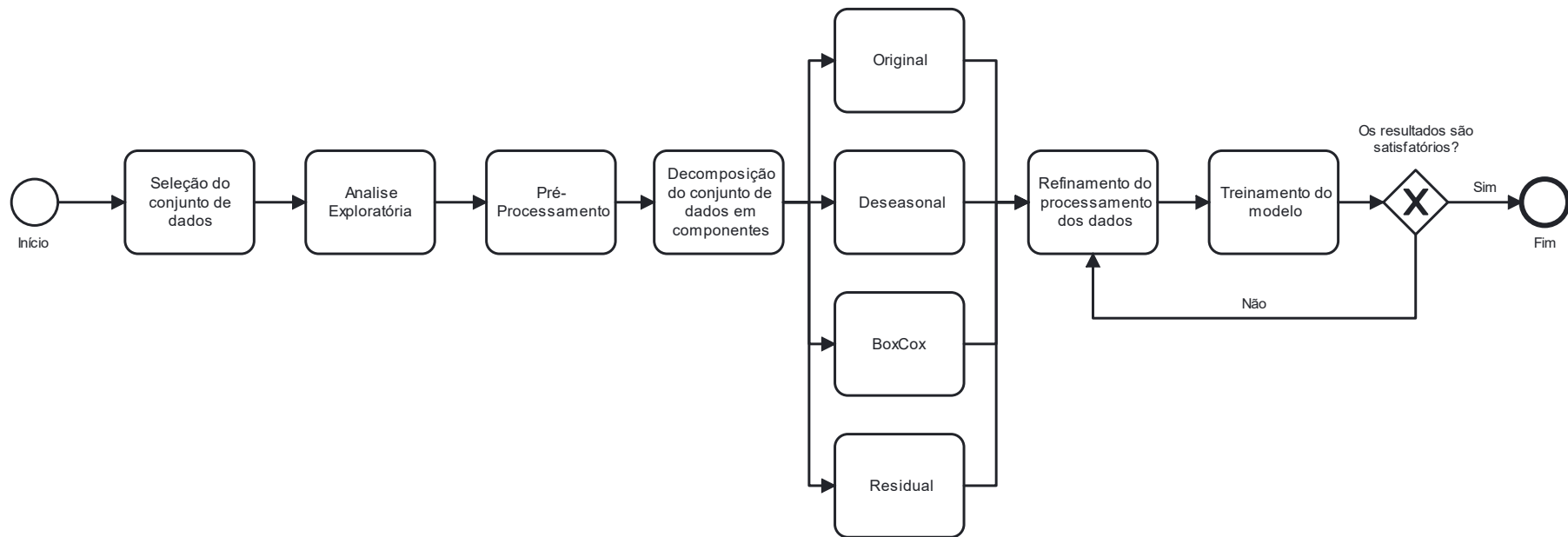
# Séries Temporais



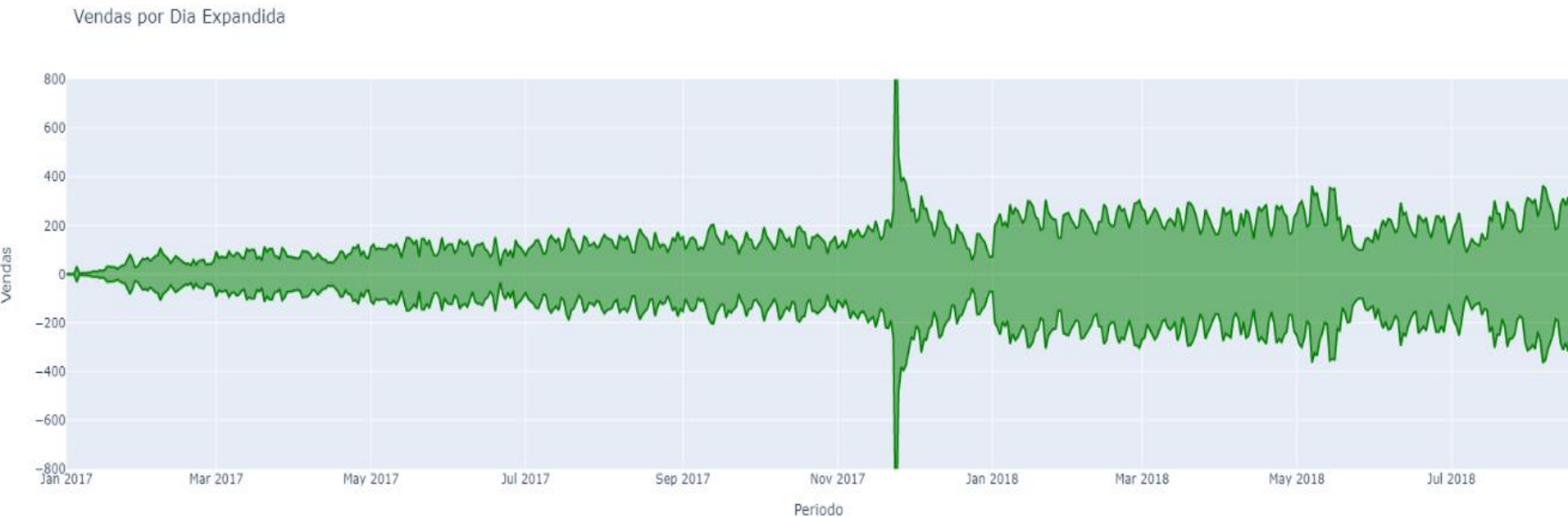


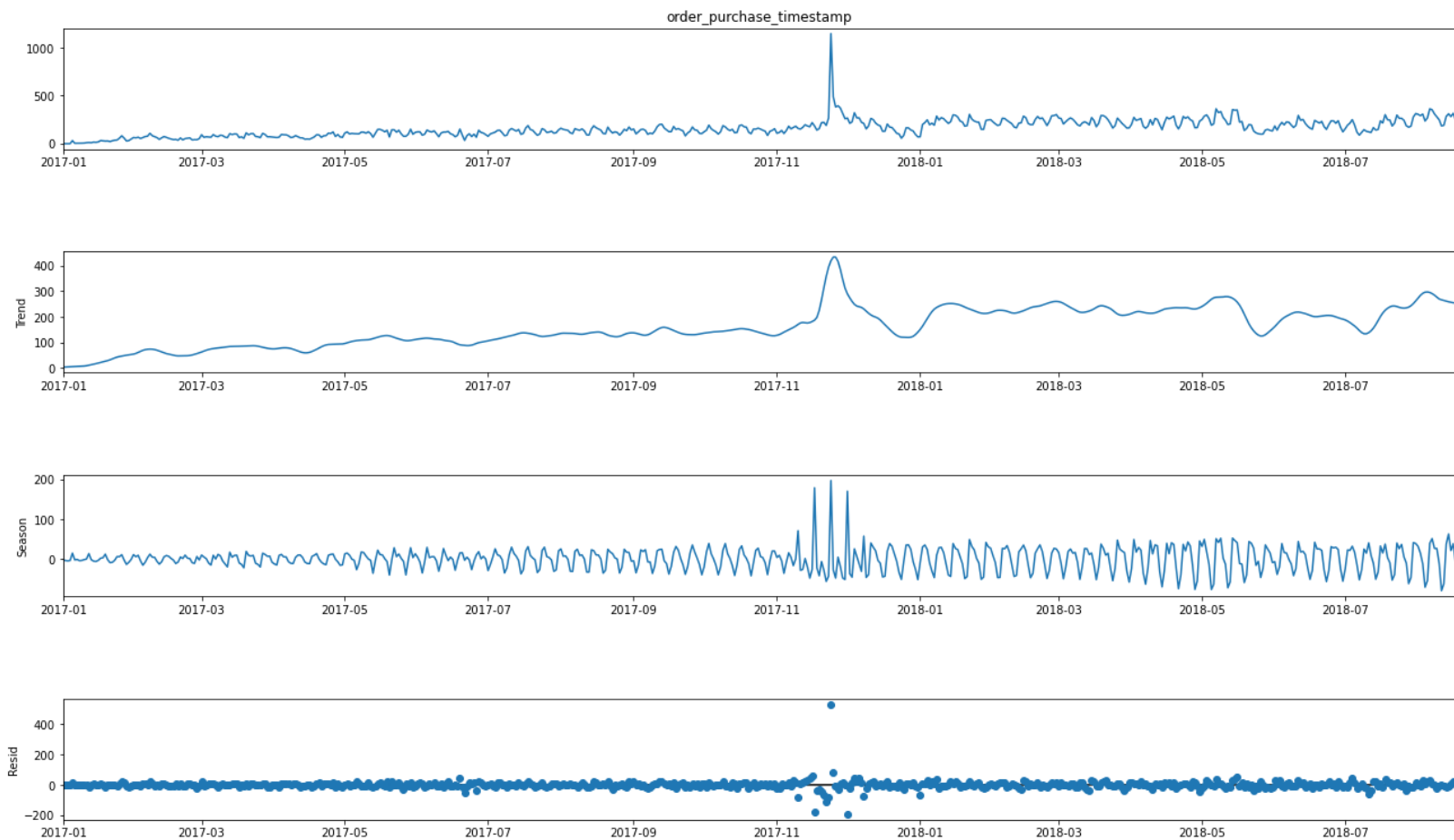


# 7. Séries Temporais



# Vendas por dia desconsiderando os outliers





Conjunto de Dado	Descrição
Original	Dados originais sem nenhuma transformação.
Deseasonal	Dados em que se aplica a remoção da sazonalidade.
BoxCox	Dados que aplicam a transformação de Box Cox, com objetivo de aproximar os dados de forma ótima de uma normal.
Residual	Variações aleatórias na série temporal.

Grupo	Nível de Confiança	Estacionário
Original	99%	Não
Deseasonal	99%	Não
BoxCox	99%	Não
Residual	99%	Sim



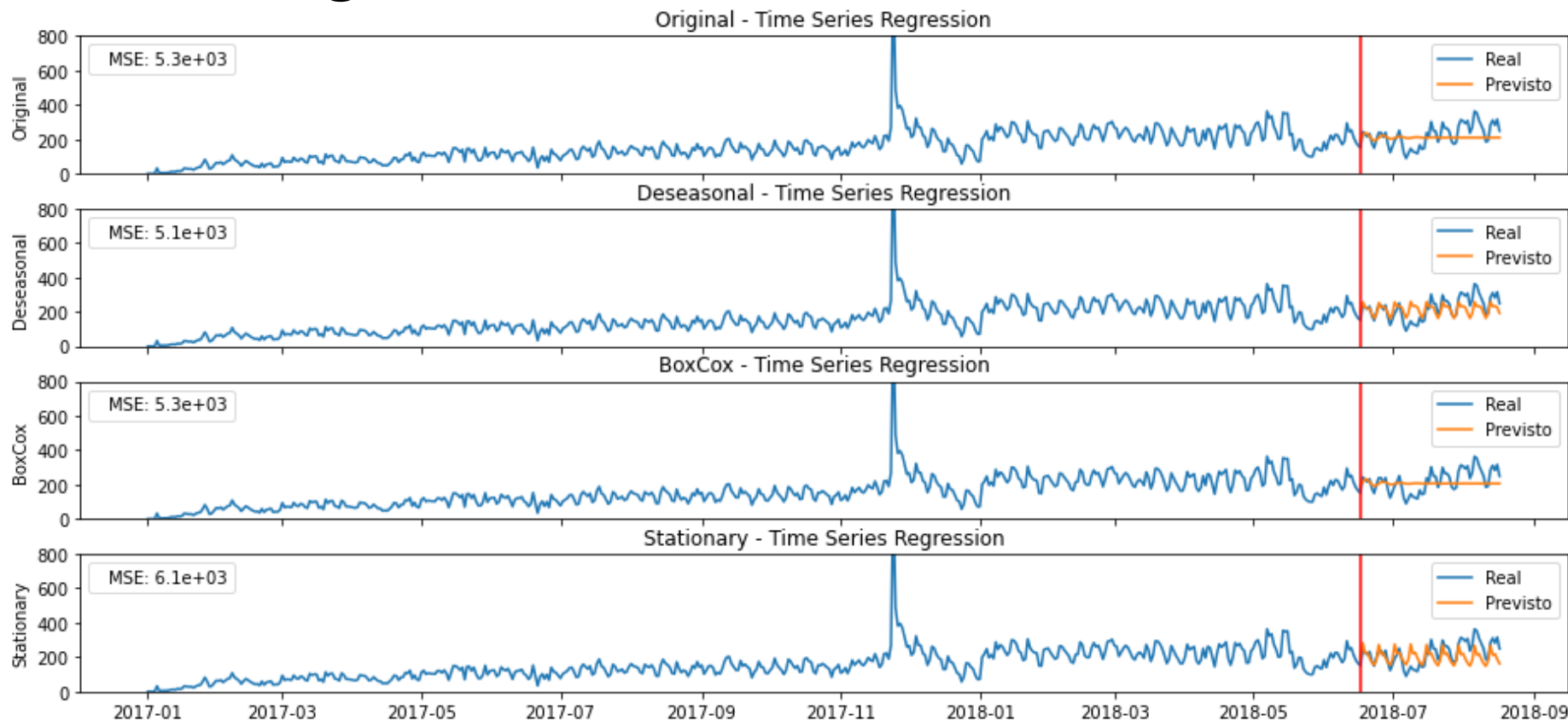
## 7.1 Algoritmos

- ⦿ Regressão Linear;
- ⦿ Triple Exponential Smoothing (TES);
- ⦿ Autoregressive integrated moving average (ARIMA);
- ⦿ Long short-term memory (LSTM);
- ⦿ XGBoost;



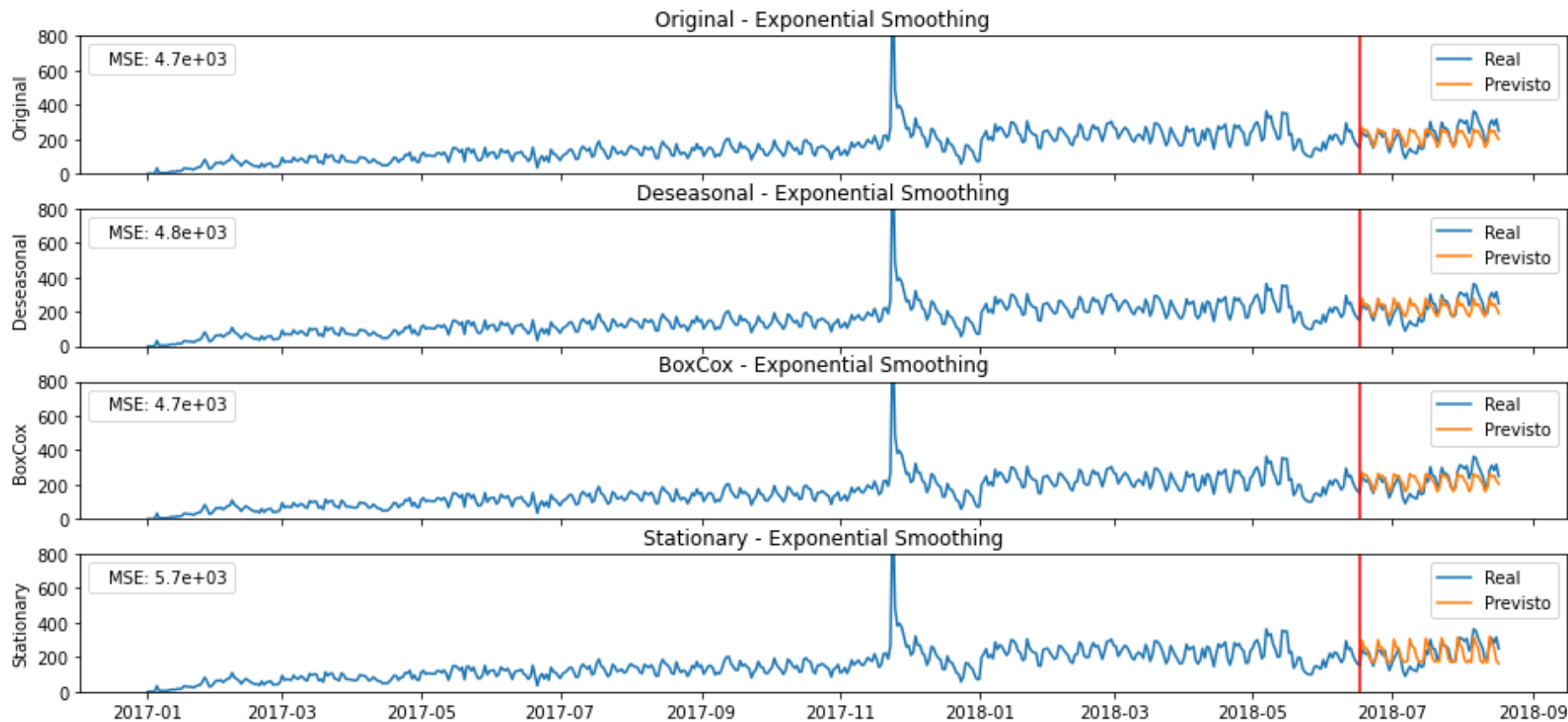


## 7.2 Regressão Linear





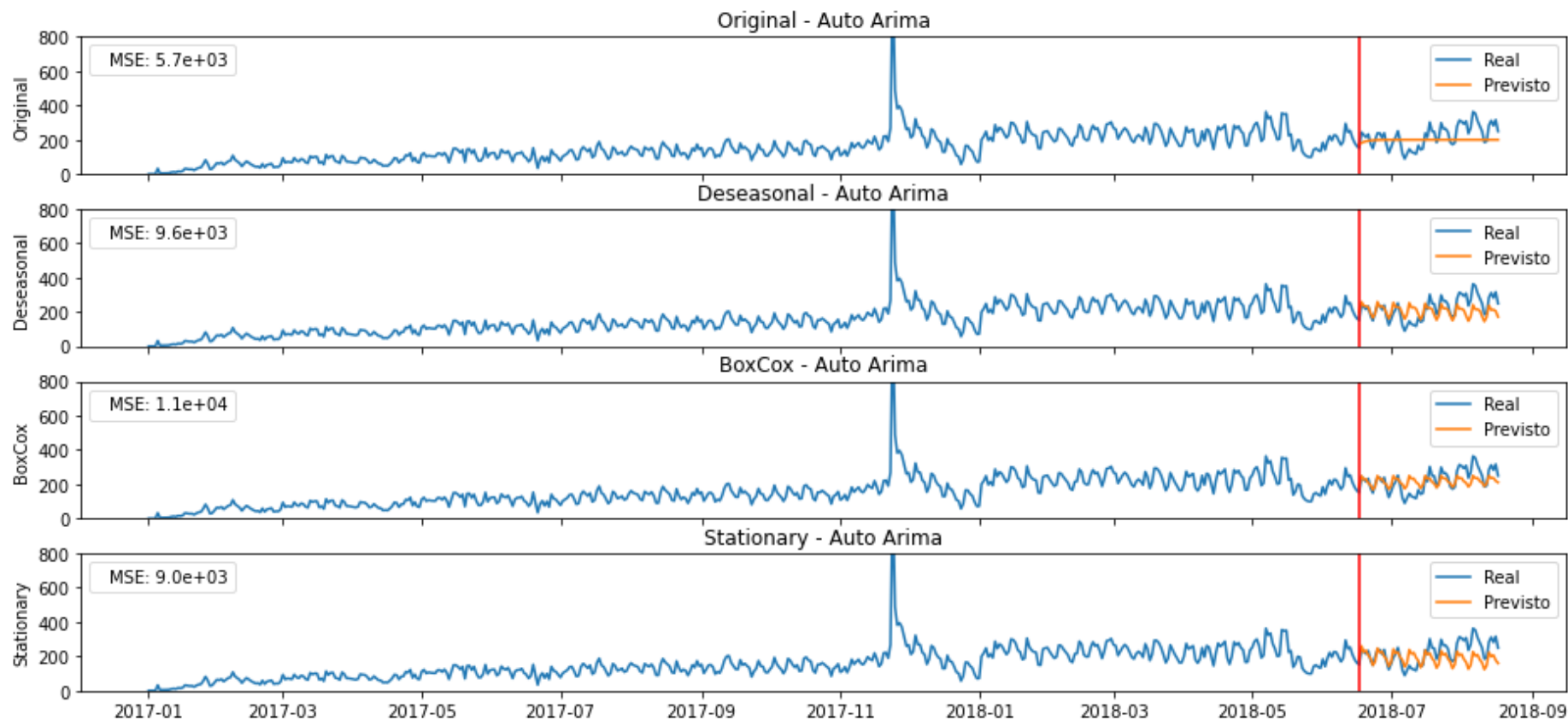
## 7.3 TES





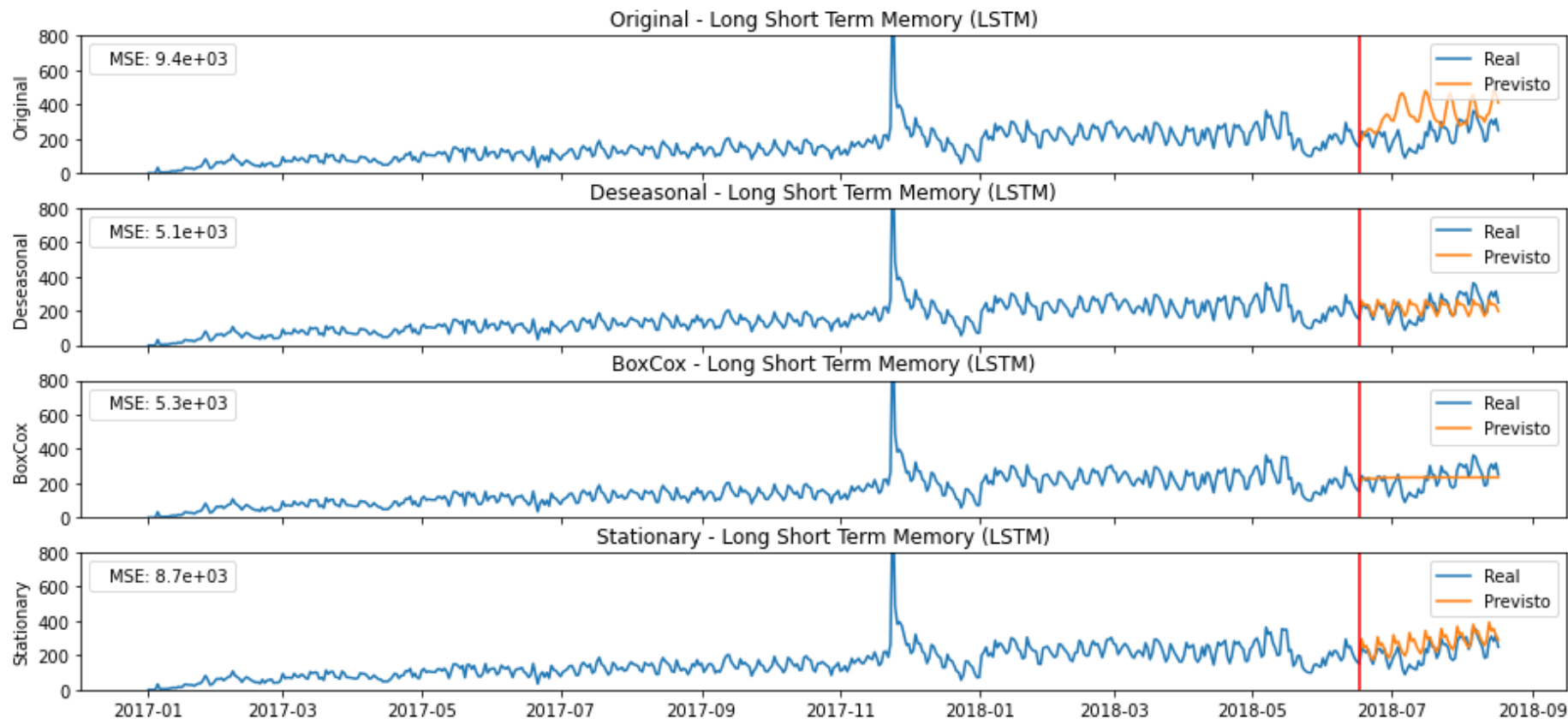


## 7.4 ARIMA





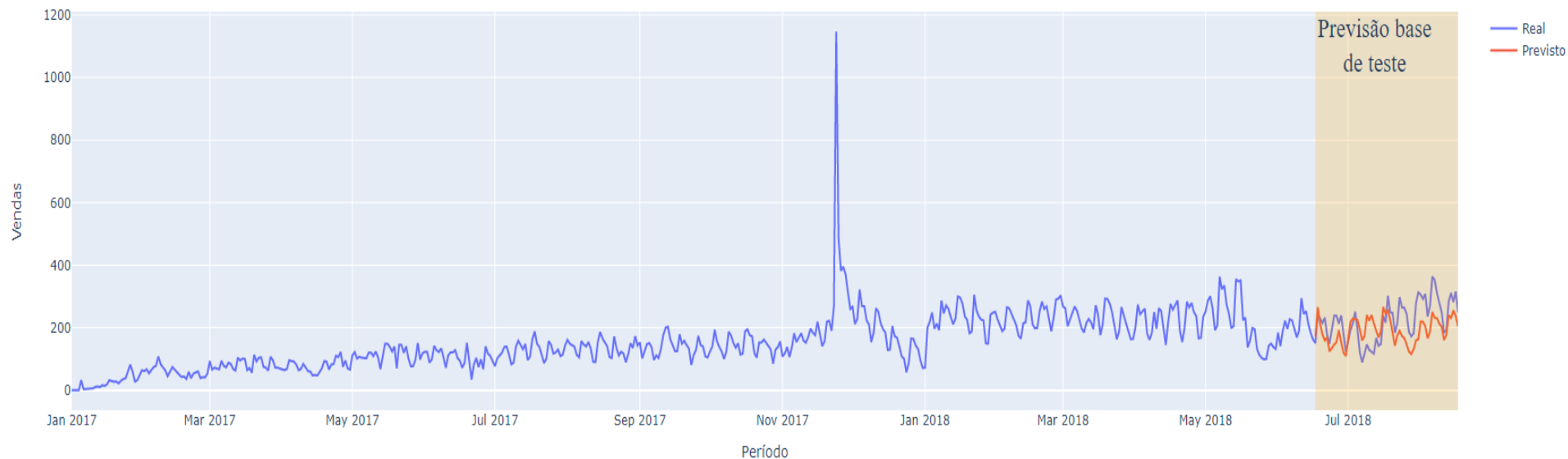
## 7.5 LSTM





## 7.6 XGBoost

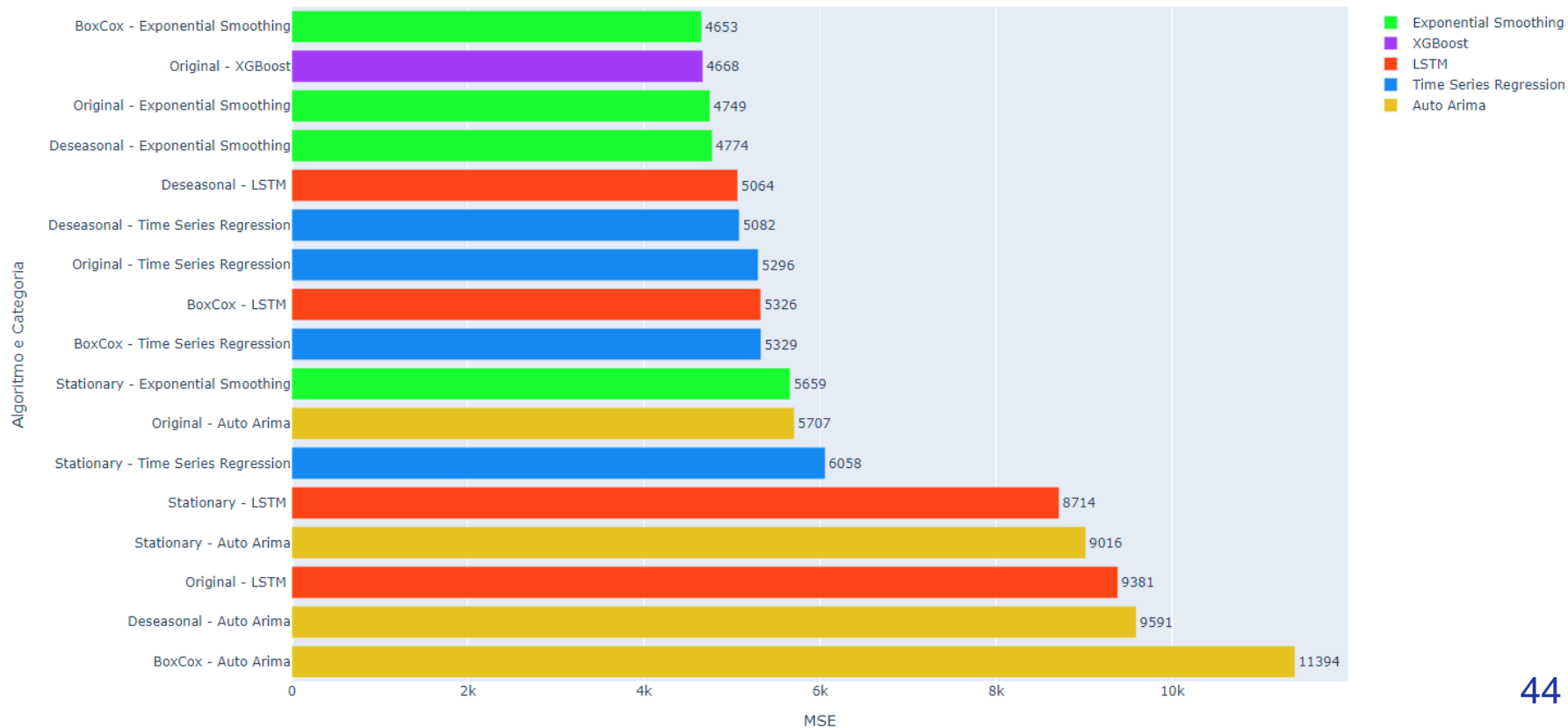
Original - XGBoost



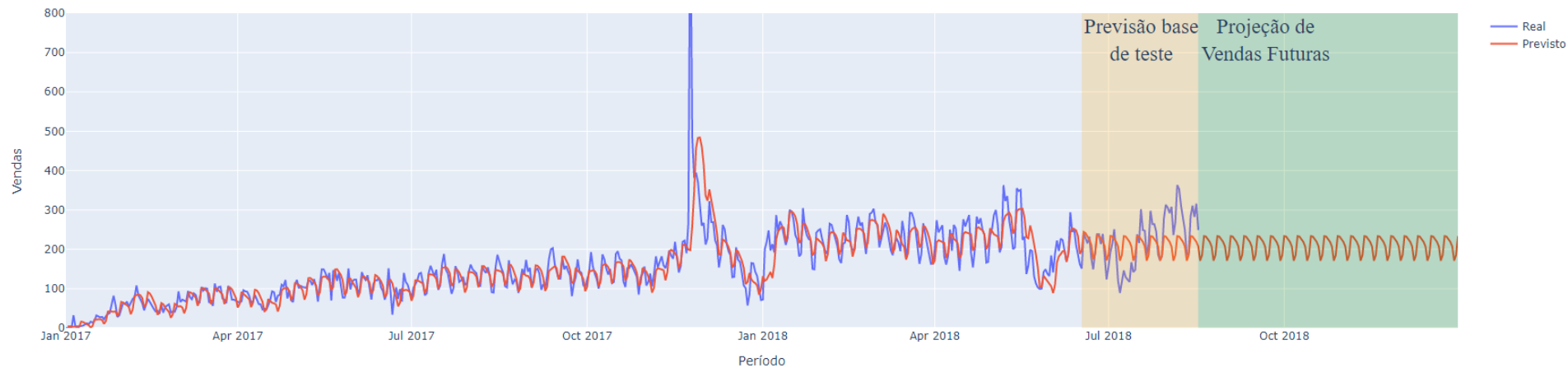


# 7.7 Compação Modelos

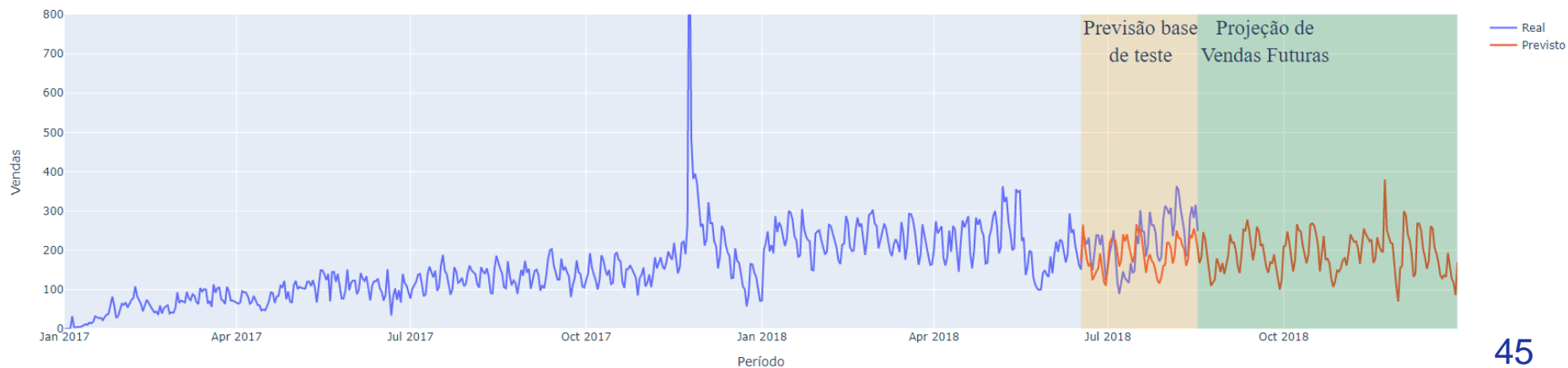
MSE por Algoritmo e Tipo de Dado



BoxCox - Exponential Smoothing



Original - XGBoost



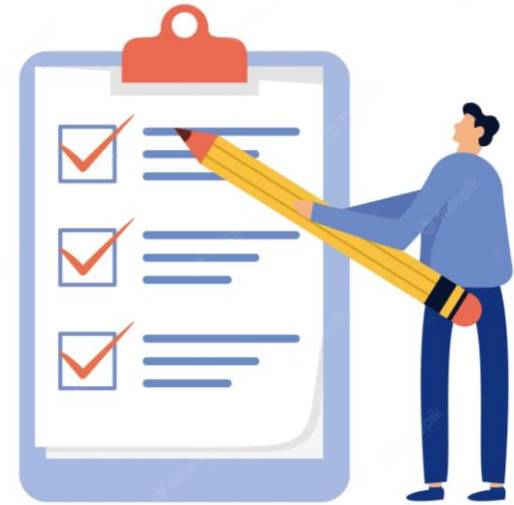
# Limites de confiança de 95% para os valores reais.

Original - XGBoost



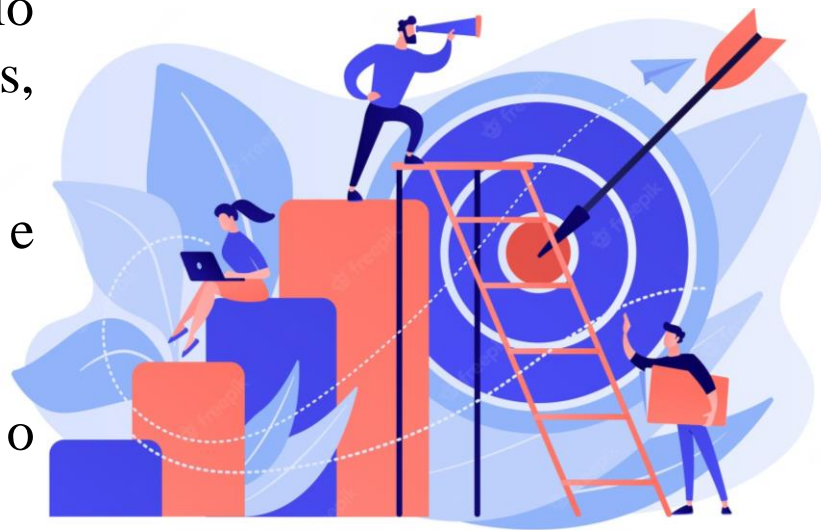
## 8. Conclusão

- ✓ Definição de atributos mais relevantes;
- ✓ Definição de algoritmos de Mineração de Dados;
- ✓ Aplicação de técnicas de Mineração de Dados;
- ✓ Avaliação de resultados obtidos.



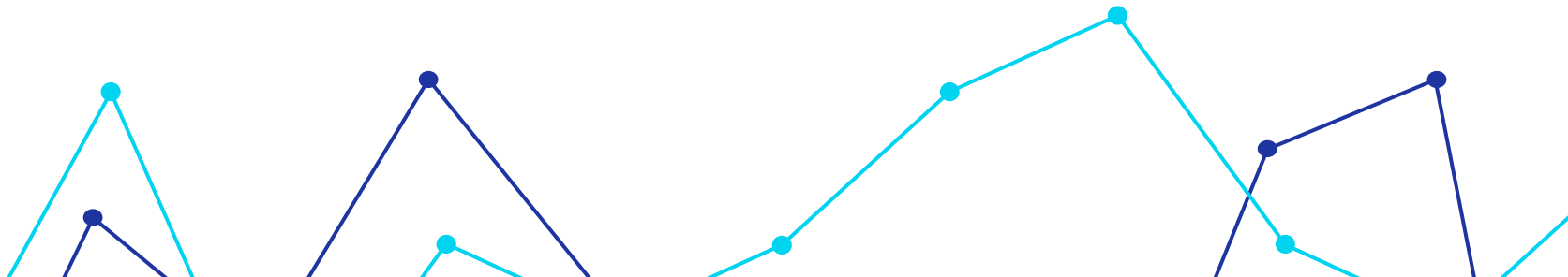
## 9. Trabalhos Futuros

- ✓ Continuação do estudo do comércio eletrônico, abordando outros algoritmos, fontes de dados e métricas;
- ✓ Ampliação do espaço de hipótese e parâmetros testados;
- ✓ Busca de soluções para manter o consumidor mais ativo no comércio.





# OBRIGADO





# Metodologia

---

O presente trabalho utiliza o método indutivo. Tratando-se de uma pesquisa aplicada, realizada com objetivo de obter novos conhecimentos para responder as perguntas de pesquisa do projeto.

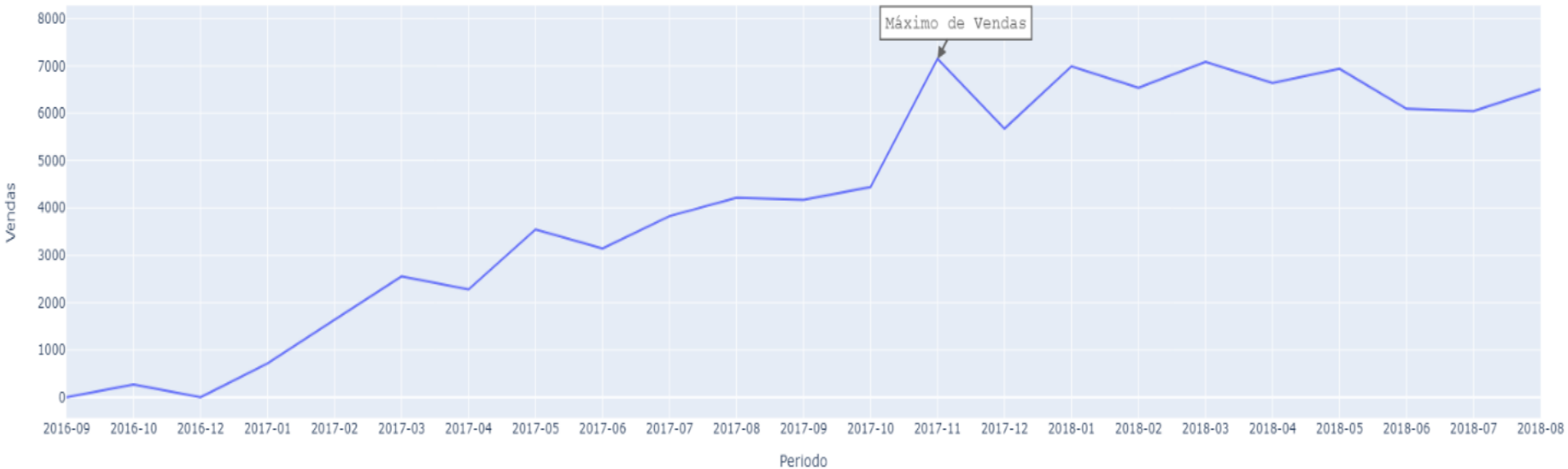


## 4. Trabalhos Relacionados

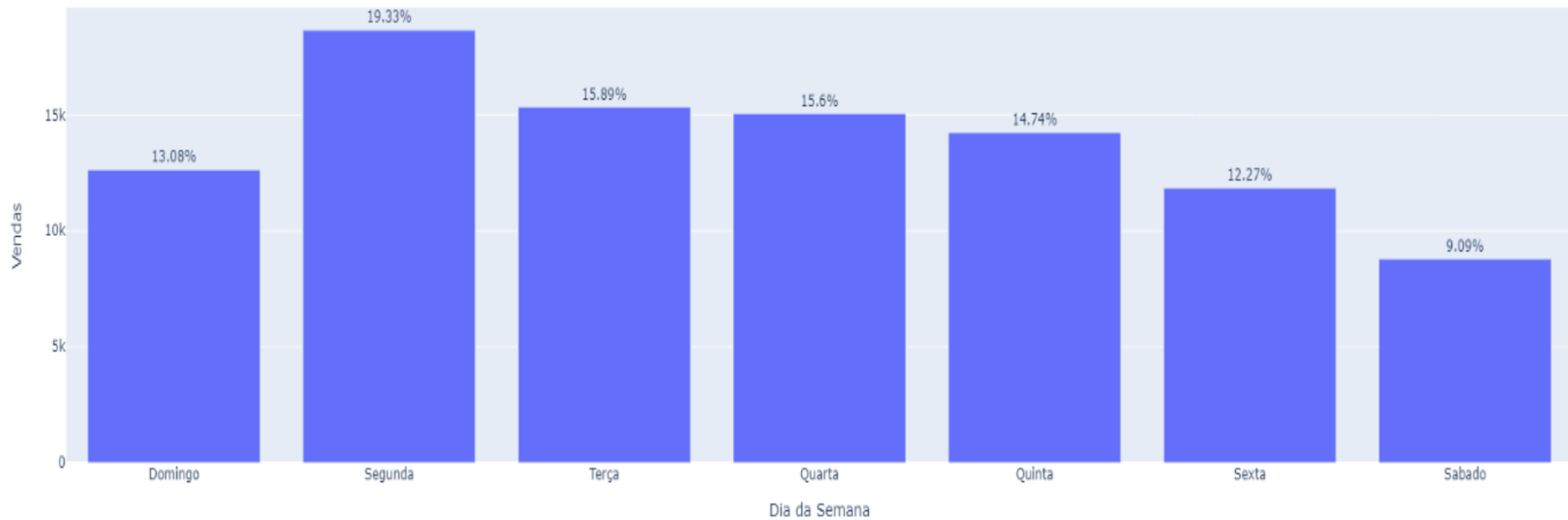
ID	Autores	Ano	Título	Fonte de dados
1	Seonghoon Moon, Suman Bae, Songkuk Kim	2014	Predicting the Near-Weekend Ticket Sales Using Web-Based External Factors and Box-Office Data	ACM Digital Libraty
2	Patcharin Ponyiam, Somjit Arch-int	2018	Customer Behavior Analysis Using Data Mining Techniques	IEEEExplore
3	Guixiang Zhua, Zhiang Wub, Youquan Wangb, Shanshan Caob, Jie Cao	2019	Online purchase decisions for tourism e-commerce	ScienceDirect

O período com maior quantidade de vendas foi Novembro de 2017.

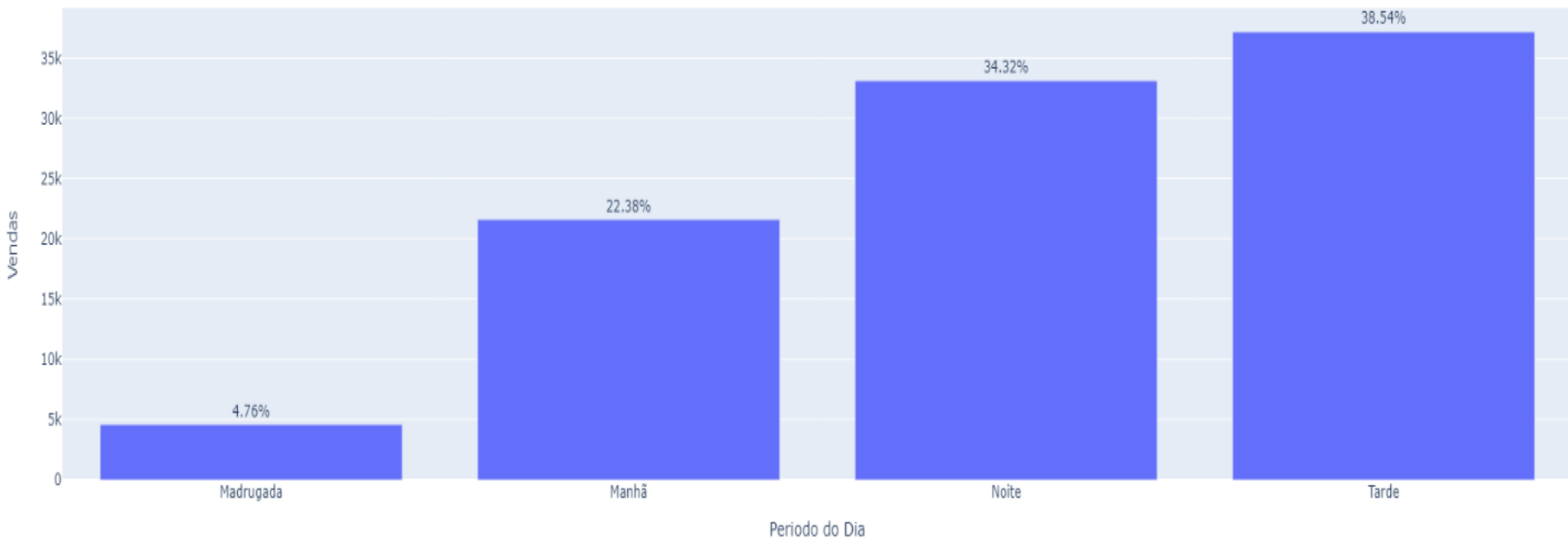
Evolução das Vendas Mensais



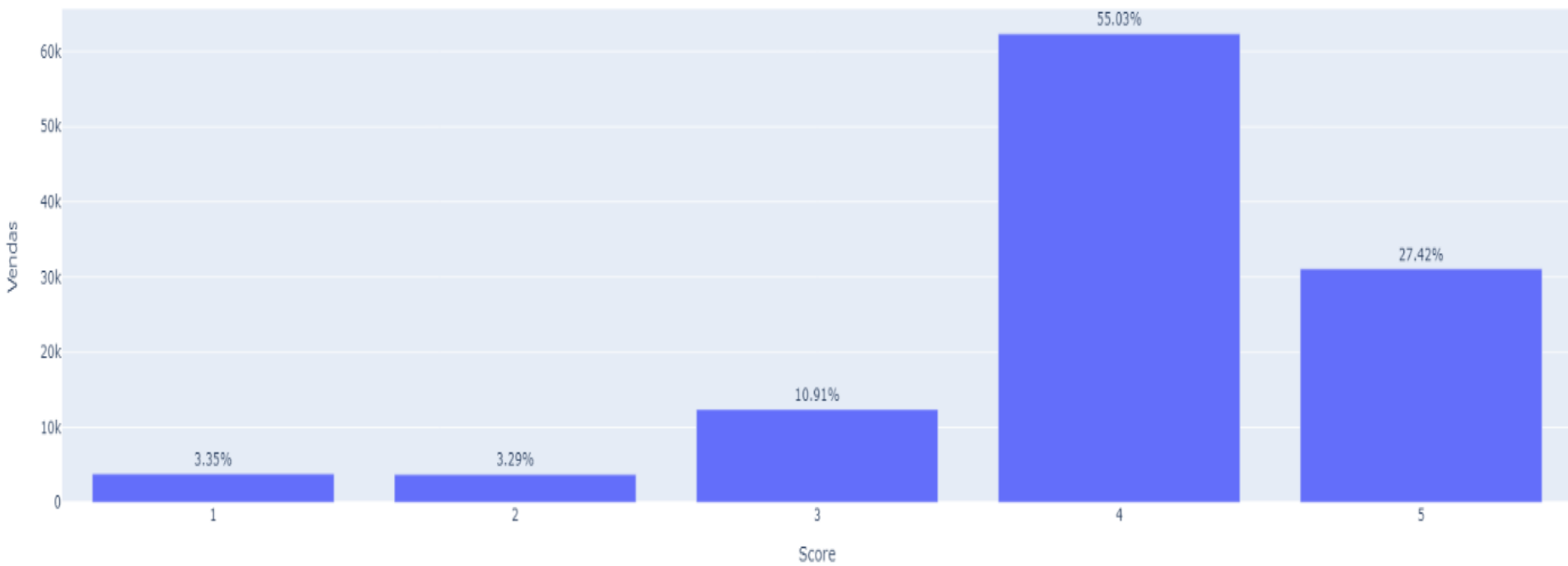
Total de Vendas por Dia da Semana



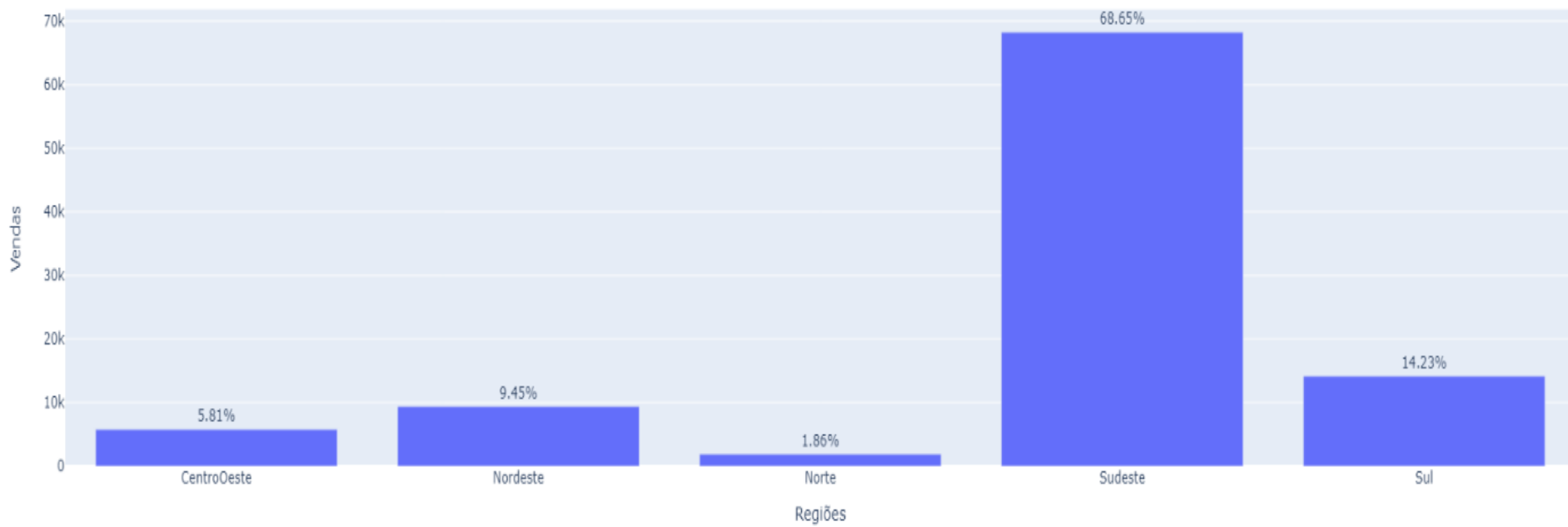
Total de Vendas por Período do Dia



Vendas por Score (Score Medio por Produto)

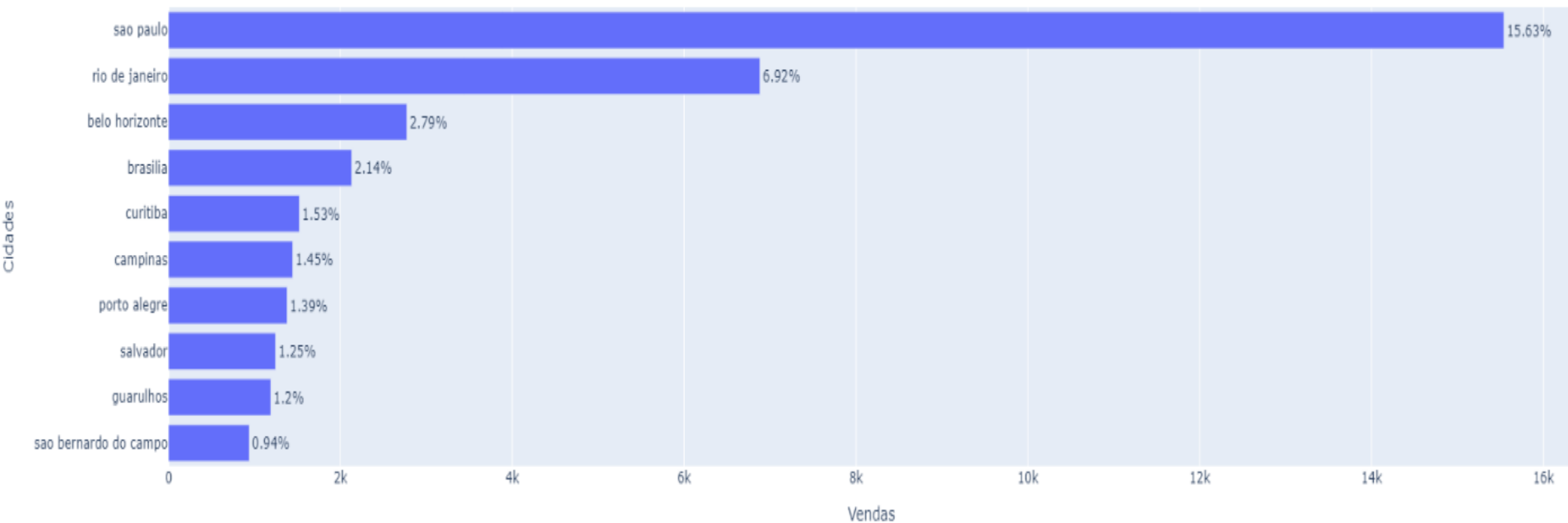


Total de Vendas por Região no Brasil

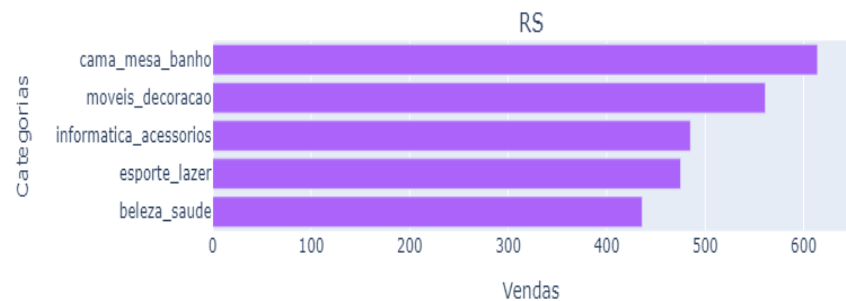
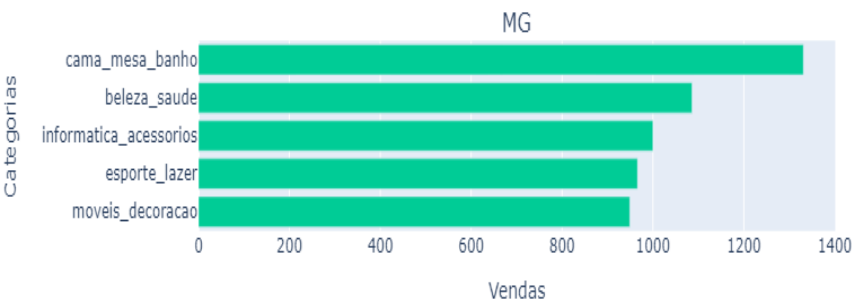
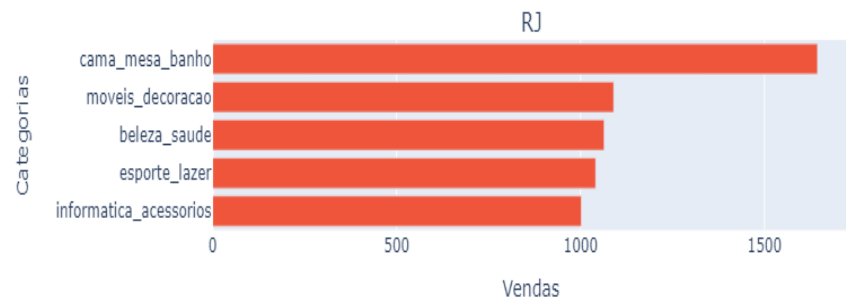
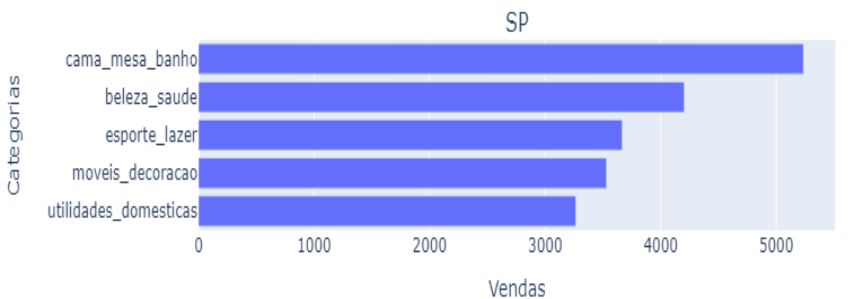




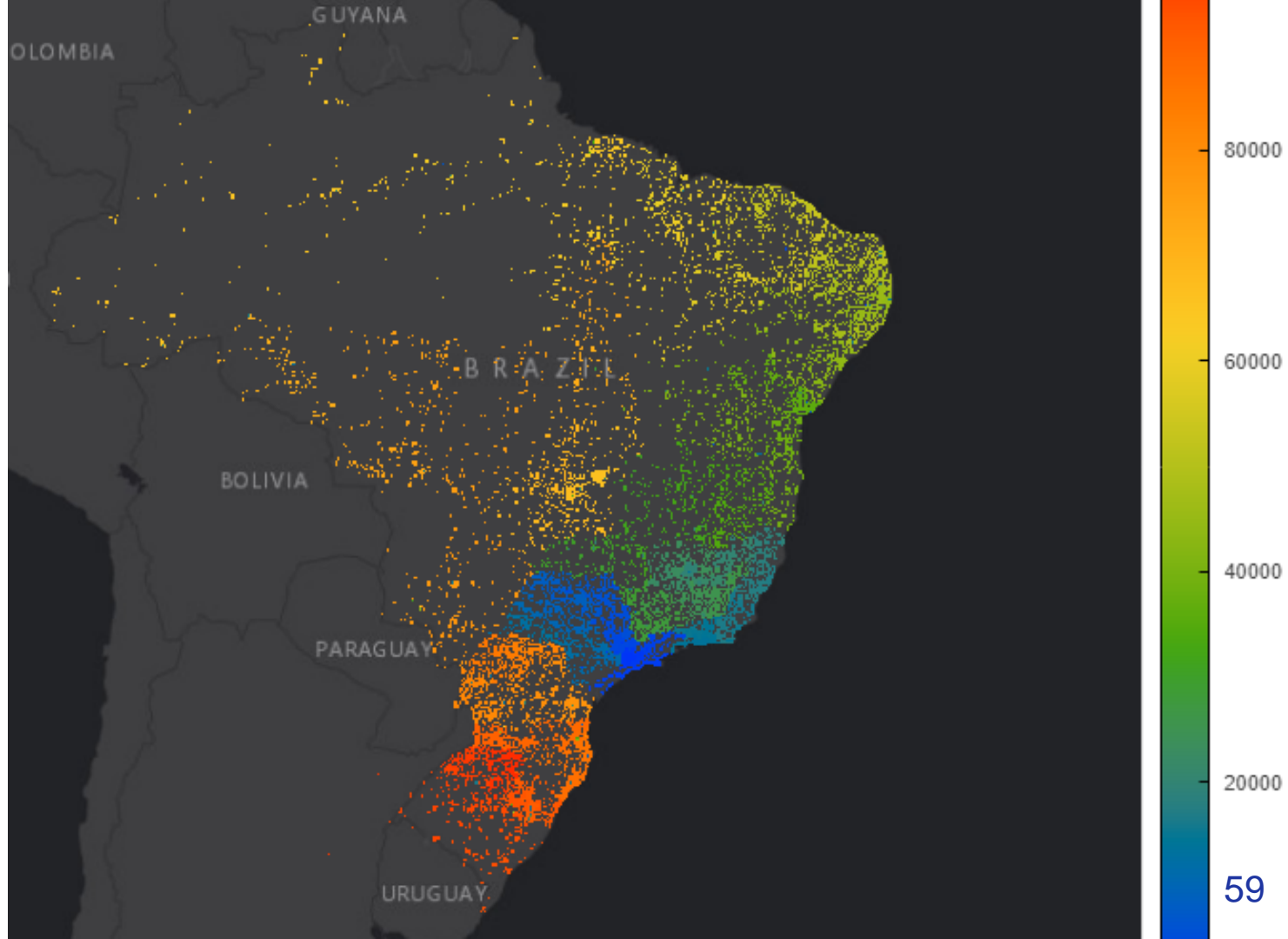
## Top 10 Cidades com Maior Numero de Vendas



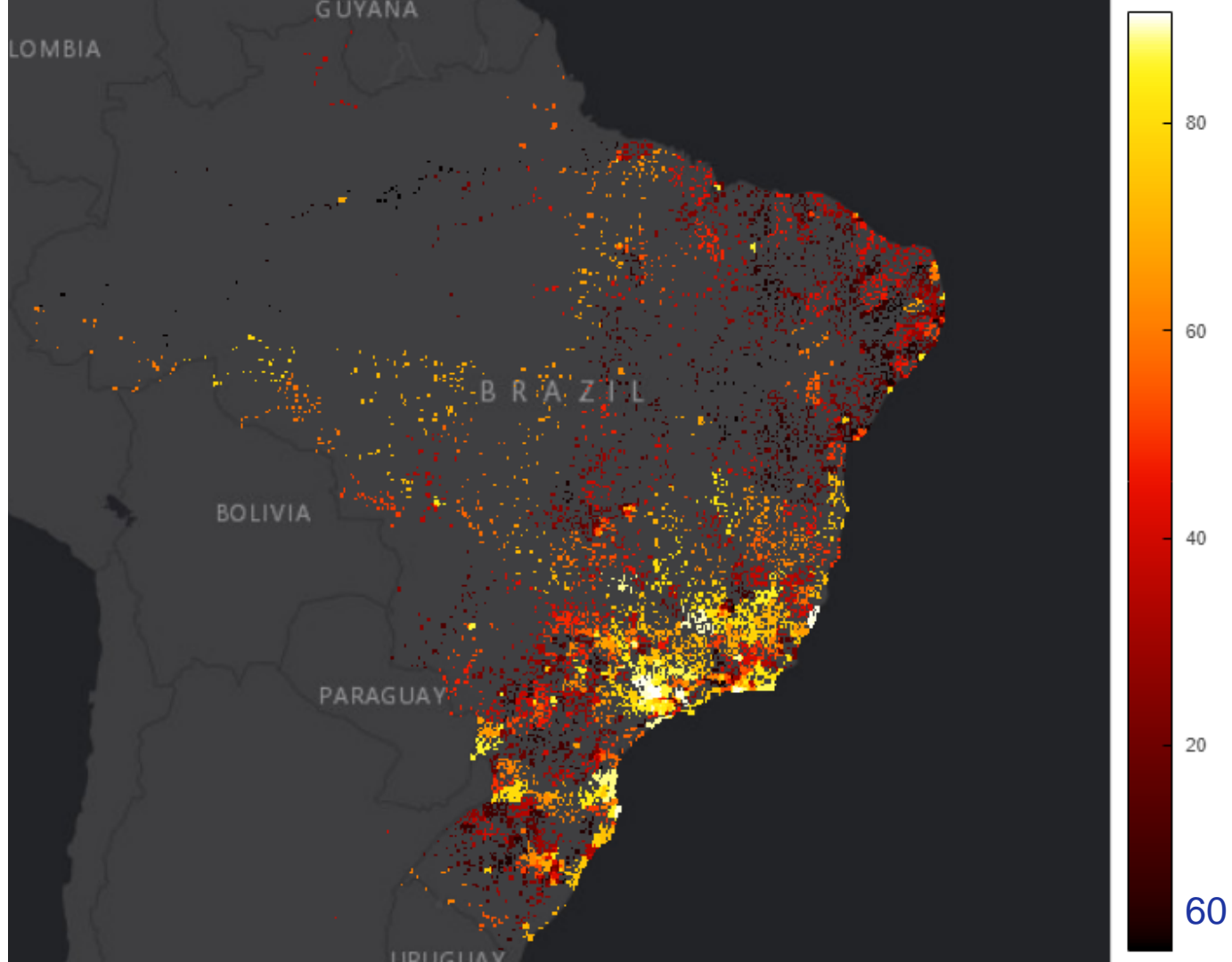
## Categorias Mais Vendidas por Estado



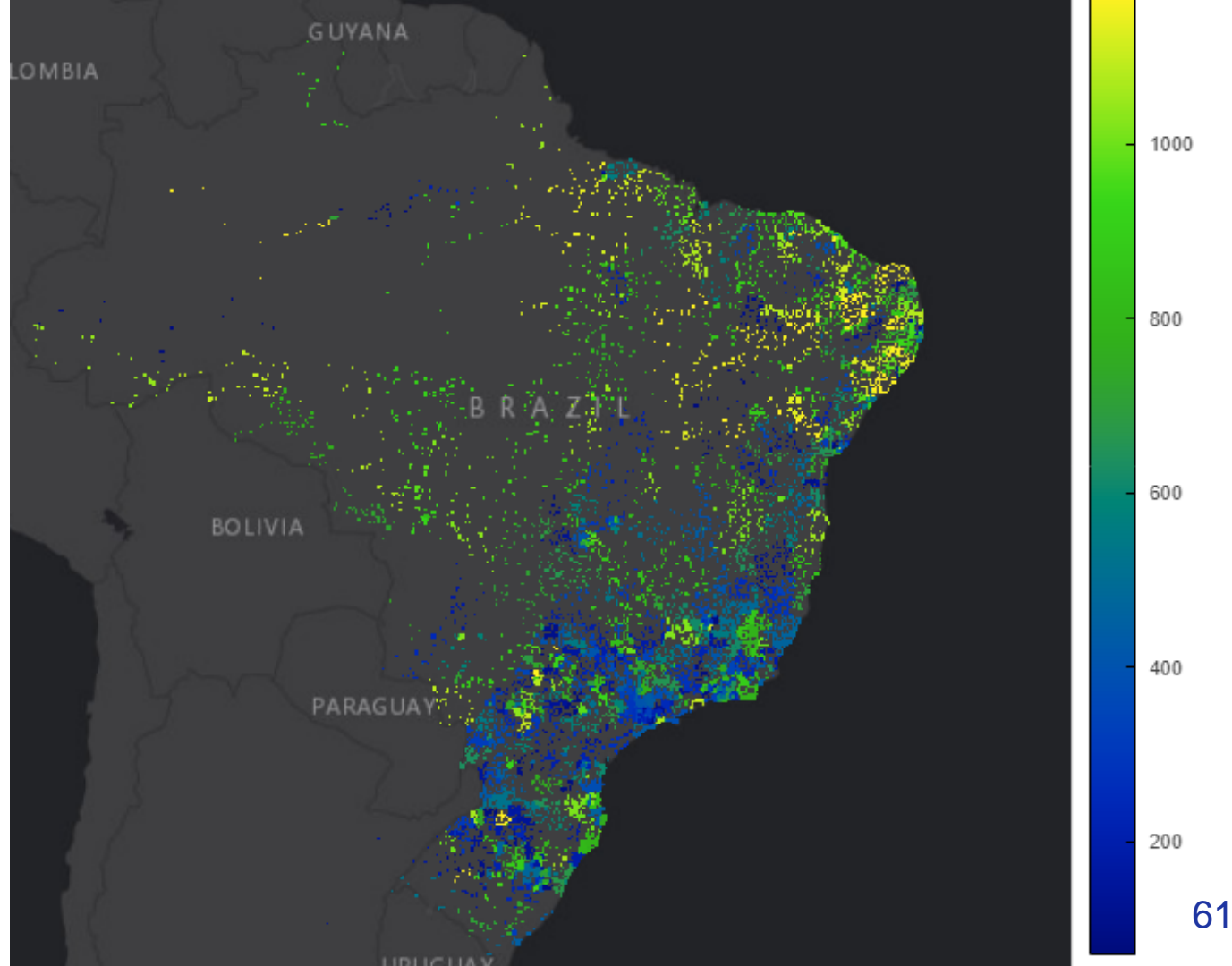
# Pedidos por CEP do Brasil



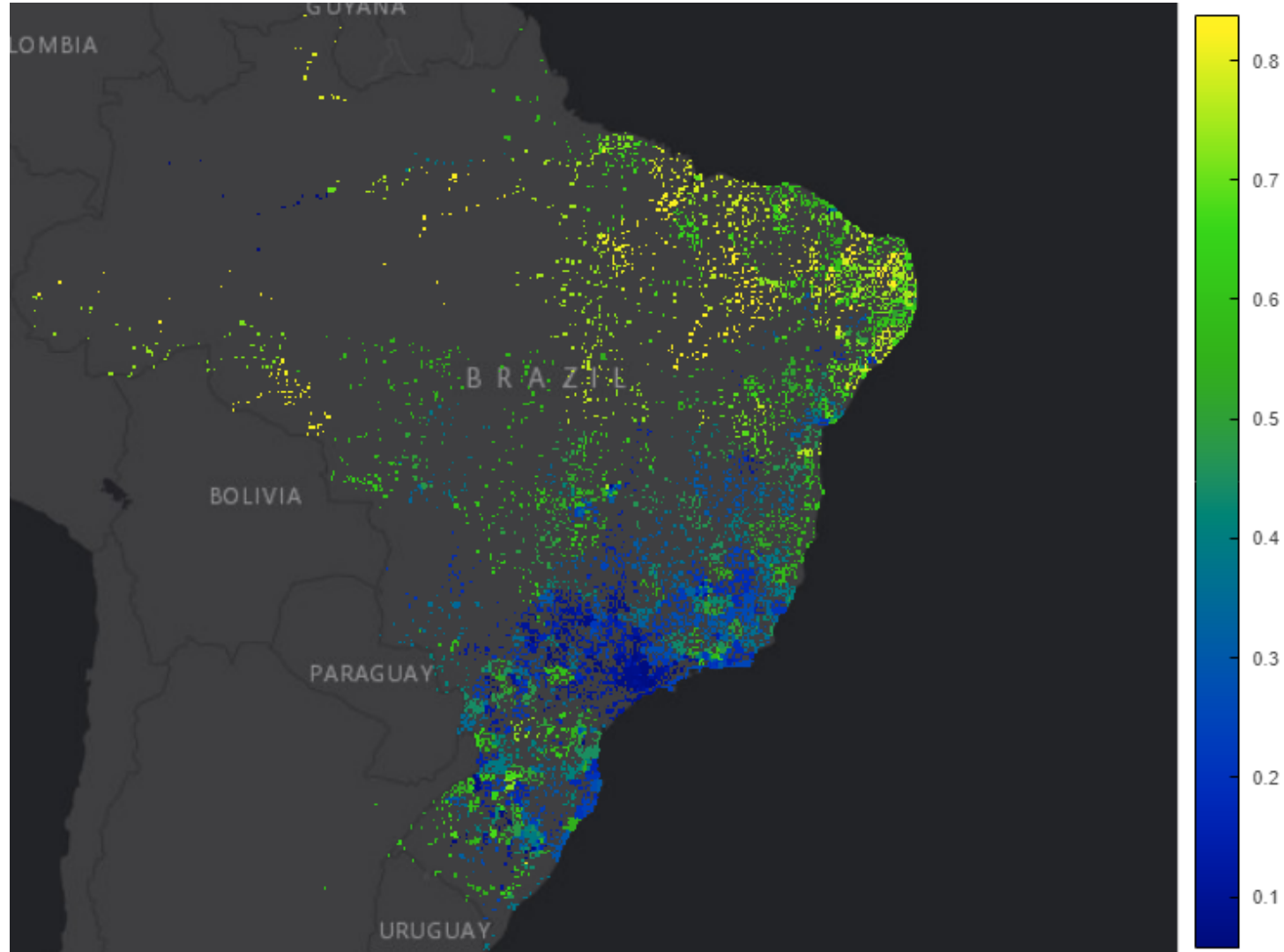
# Faturamento por região no Brasil



# Ticket médio por região no Brasil

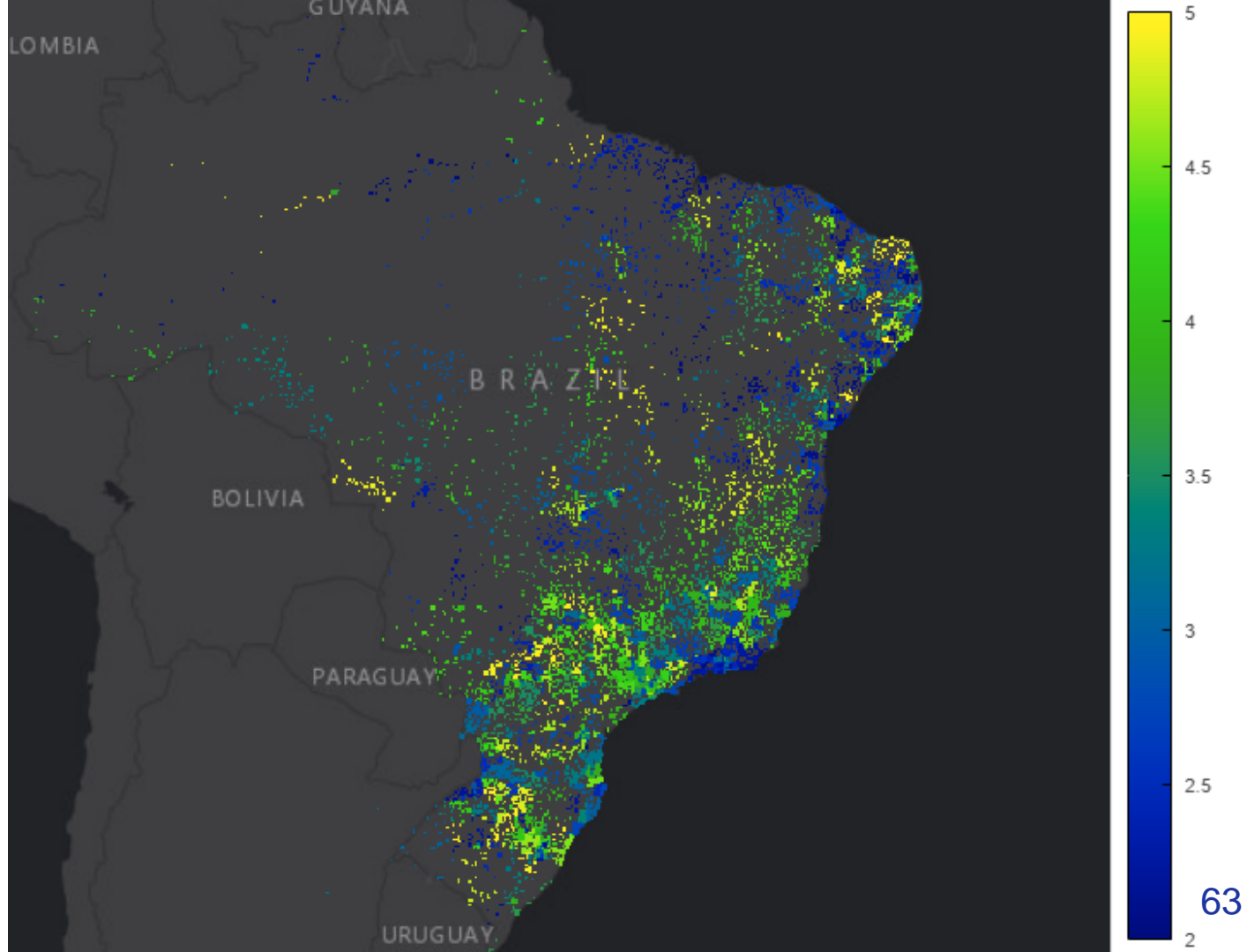


# Média de Frete por Região no Brasil

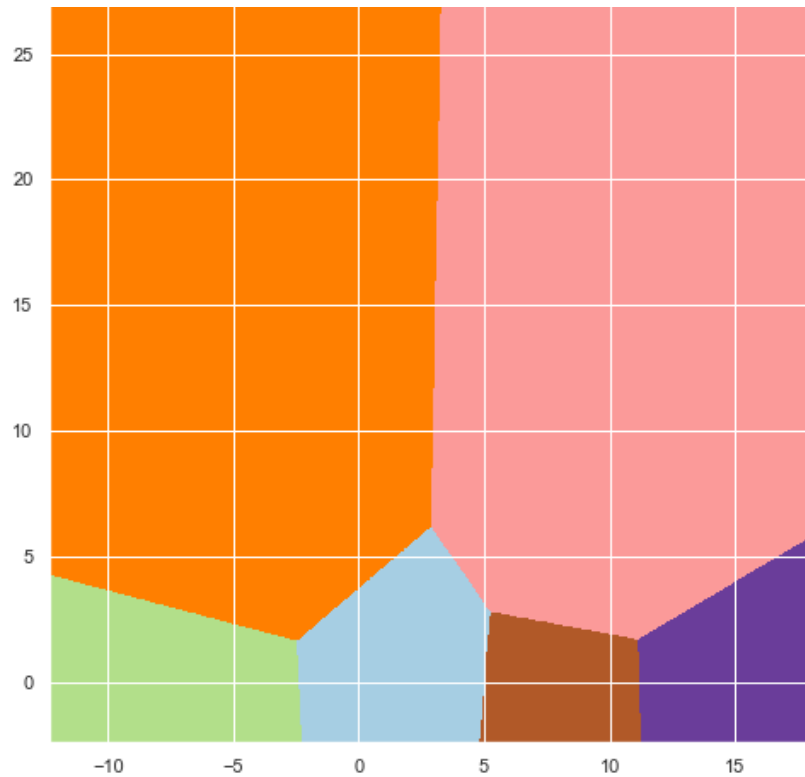


# Review

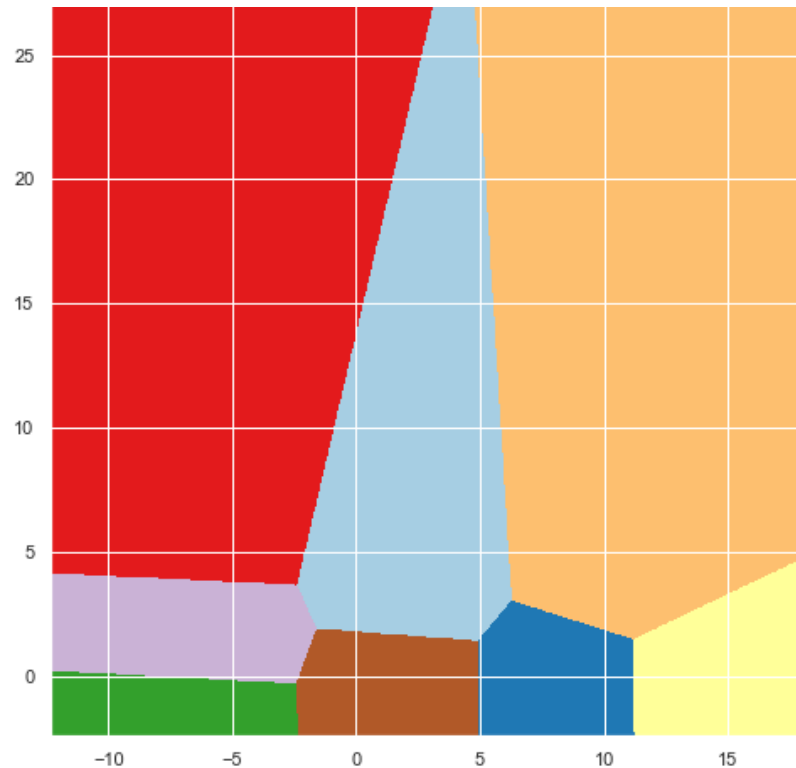
## Médio por Região no Brasil



Meshgrid de 6 clusters com PCA



Meshgrid de 8 clusters com PCA





# Clusterização Simples

Cluster	Preço	Frete	Parcelas	Score	Itens	RFM	Pagamento	Estado	Região
0	165.98	23.13	2.4	3.96	1.33	2.95	Crédito	RJ	Sudeste
1	130.2	15.62	1.77	4.16	1.37	3.05	Crédito	SP	Sudeste
2	187.74	23.64	3.16	4.01	1.38	2.97	Crédito	MG	Sudeste
3	320.17	23.27	8.41	3.98	1.67	3.03	Crédito	SP	Sudeste

# Clusterização por RFM

Cluster	Recency	Frequency	Monetary	RFM
0	395.67	1	296.31	1.64
1	133.88	1	296.32	3.9
2	221	2.23	682.74	3.27
3	221.85	1.49	26987.42	3.2



## 6.4 LSTM

Atributo	Valor / Parâmetro
Camada de entrada LSTM	Input_shape: (7, 1)
Camada Densa com 1 neurônio de entrada	Initializer: GlorotUniform
Camada de BatchNormalization	N/A
Camada Densa com 1 neurônio de entrada	Initializer: GlorotUniform
EarlyStopping	Monitor: loss Patience: 3 Mode: min
Otimizador Adam	Taxas de aprendizado testadas: 0.00001; 0.0001; 0.001; 0.01
Função de perda	MSE
Épocas	1000
StandarScaler	Pré-treinamento