

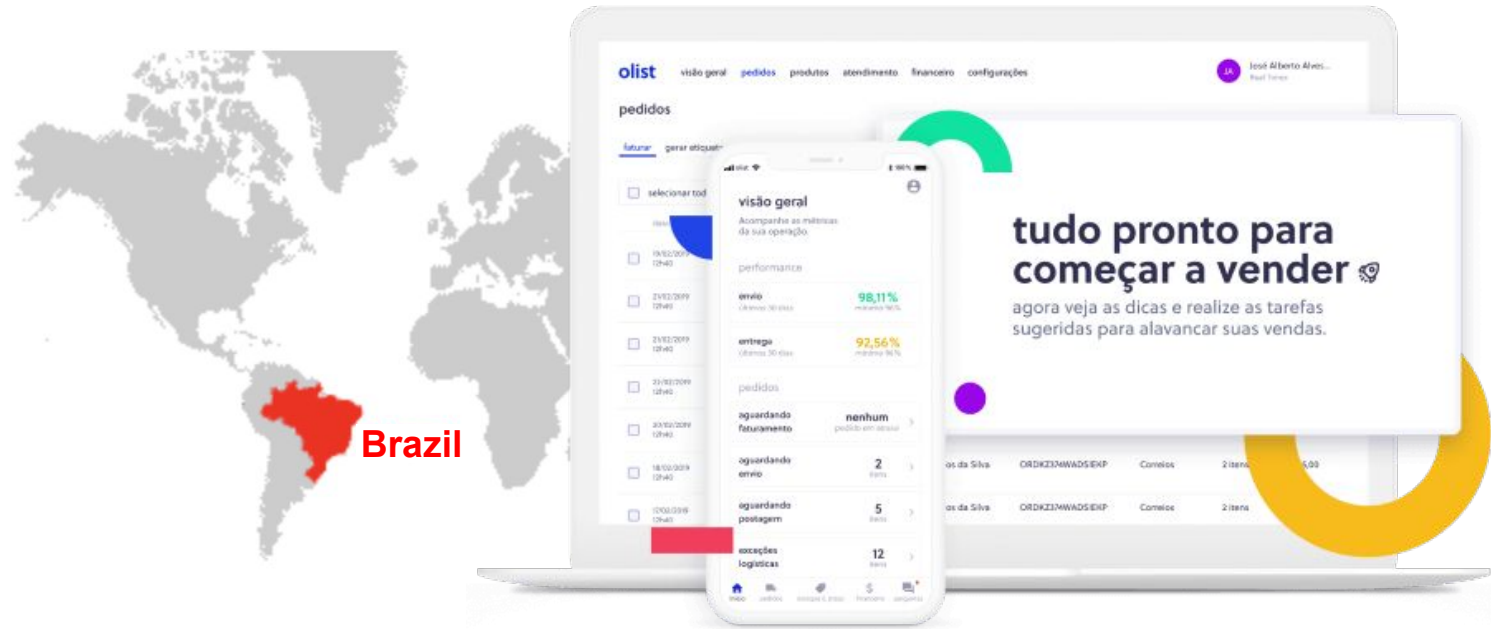
The background of the slide is a night-time photograph of a city skyline, likely New York City, with numerous skyscrapers illuminated. Overlaid on this image is a complex digital graphic consisting of many thin, vertical blue lines of varying heights. These lines are connected by a series of flowing, translucent blue waveforms that create a sense of movement and data flow. Small, bright blue dots are scattered throughout the scene, particularly along the vertical lines and within the waveforms, resembling data points or stars. The overall color palette is dominated by deep blues and teals, with the warm lights of the city providing a contrasting background.

Olist E-Commerce

Final Visualisation Project Presentation

Group 16

Recap: Dataset from Olist e-commerce in Brazil, to provide insights and recommendations to sellers on e-commerce



Selected dataset is on Olist orders, **an e-commerce platform in Brazil on orders between Oct 2016 to Oct 2018**. Our objective is to provide insights to improve e-commerce performance of our target audiences, i.e. **the potential and existing sellers** based on three proposed key hypotheses.

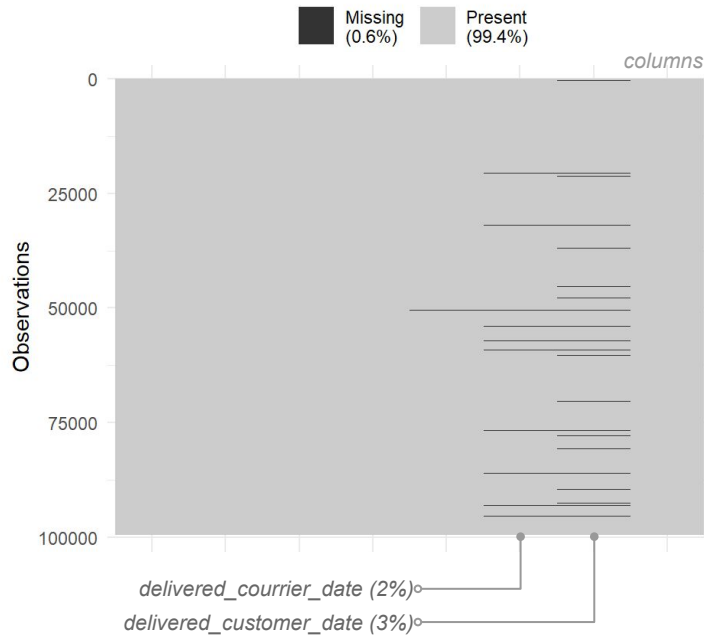
Data Wrangling & Pre-processing

- Missing values
- Orders data across time period
- Geolocation data
- Product category data

Missing Values: Found in Orders and Reviews data at ~0.6% and ~20.9% respectively; to remove missing Orders data only

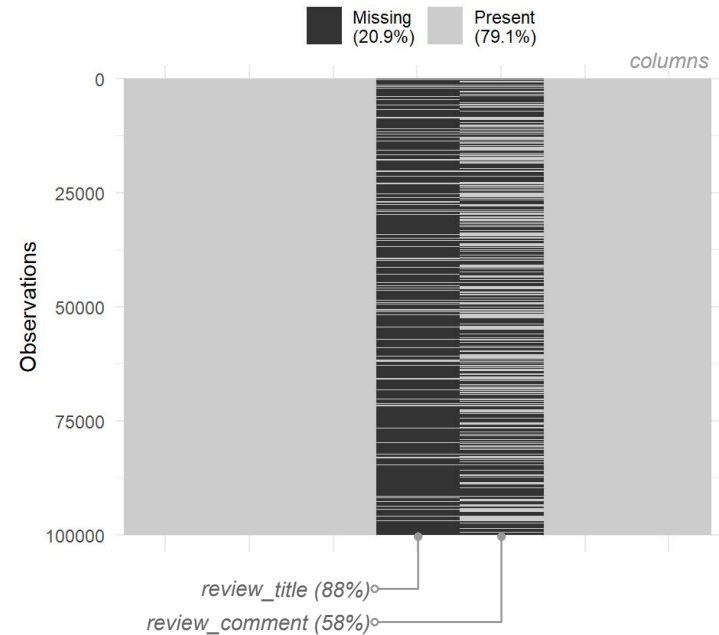
Missing values in Orders data found ~0.6%

Rows with missing values to be removed for data cleanliness

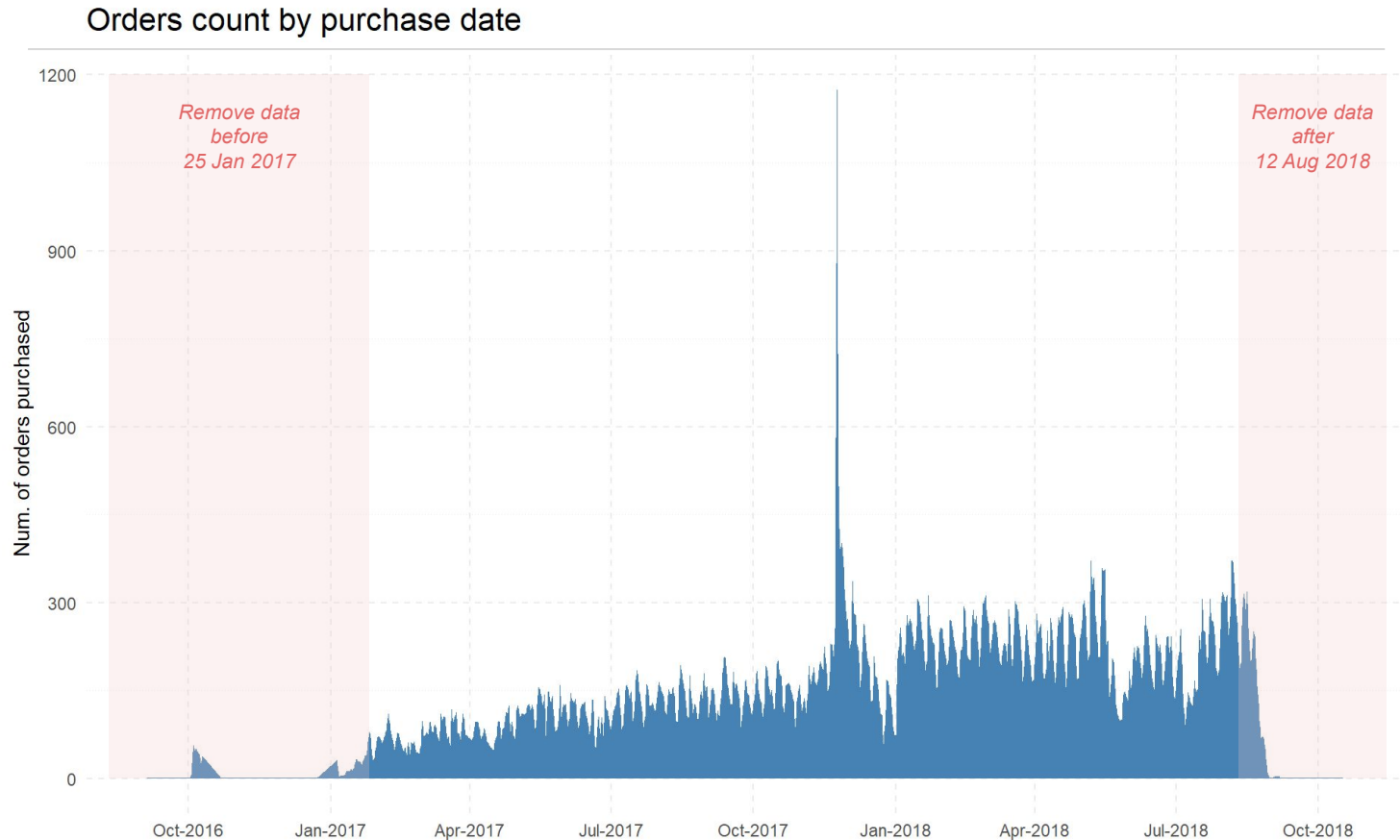


Missing values in Reviews data found ~20.9%

Missing values are expected and require no data cleaning



Orders: Suspicious order counts before 25 Jan 2017 and after 12 Aug 2018; to remove data points outside the time period (~4.0%)



Geolocation: Duplicated geolocation zip codes around ~98.1% and inconsistent spelling variations on city names

Zip Code Prefix	City	State	Latitude	Longitude
01001	são paulo	SP	-23.55064	-46.63441
01001	sao paulo	SP	-23.55050	-46.63434
01001	sao paulo	SP	-23.54978	-46.63396
01001	saopaulo	SP	-23.55050	-46.63434
01001	sao£ paulo	SP	-23.55143	-46.63407
01001	são paulo	SP	-23.55143	-46.63407

Geolocation data issues:

- Duplicated zip code (key) around ~98.1%
- Inconsistent city names such as **sao£ paulo** (random characters), **saopaulo** (no spacing) and **são paulo** (accent)

01001	saopaulo	SP	-23.55064	-46.63441
01001	saopaulo	SP	-23.55050	-46.63434
01001	saopaulo	SP	-23.54978	-46.63396
01001	saopaulo	SP	-23.55143	-46.63407

Clean geolocation data:













- Remove accents, special characters and strip the spaces
- Remove duplicated rows to prevent skew in centroids

01001	saopaulo	SP	-23.55059	-46.63420
-------	----------	----	-----------	-----------

Aggregate into centroid

Category: Product category data are re-categorised by referencing product categories from main competitors

Competitor 1: Mercado Libre				
				
Autos, Motos y Otros	Computación	Electrodomésticos y Aires Ac.	Deportes y Fitness	Inmuebles
				
Celulares y Teléfonos	Electrónica, Audio y Video	Ropa y Accesorios	Hogar, Muebles y Jardín	Accesorios para Vehículos

Competitor 2: Shopee					
					
Casa e Decoração	Celulares e Dispositivos	Brinquedos e Hobbies	Roupas Femininas	Beleza	Sapatos Femininos
					
Animais Domésticos	Câmeras e Drones	Acessórios de Moda	Roupas Masculinas	Papelaria	Sapatos Masculinos

Re-categorisation of product category

Drilled-down into 20 new category groups from 71

Original product category	Recategorised
home_appliances	Home Appliances
home_appliances_2	
small_appliances	
small_appliances_oven_and_coffee	
fashion_sport	Fashion
fashion_male_clothing	
fashion_female_clothing	

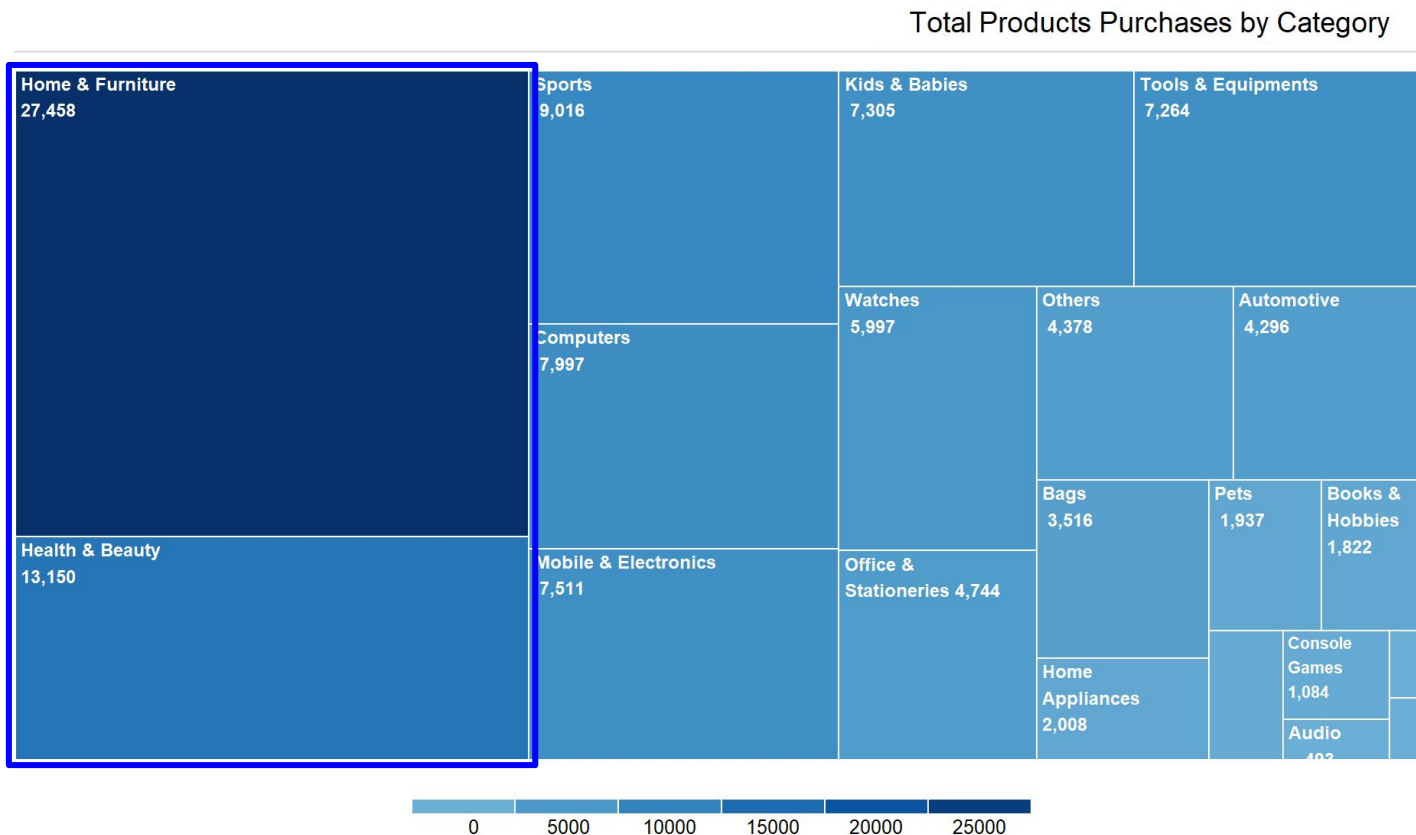
- New category groupings are referenced from **Olist's e-commerce competitors** in Brazil, i.e. Mercado Libre and Shopee
- Duplicated categories are grouped and recategorised manually from **71** into **20 new high-level categories**

Hypotheses & Insights

- 1 **Customers:** Customer's product preferences are different across states
- 2 **Orders:** Higher orders generated on products with greater product details
- 3 **Reviews:** Review ratings are significantly affected by delivery efficiency

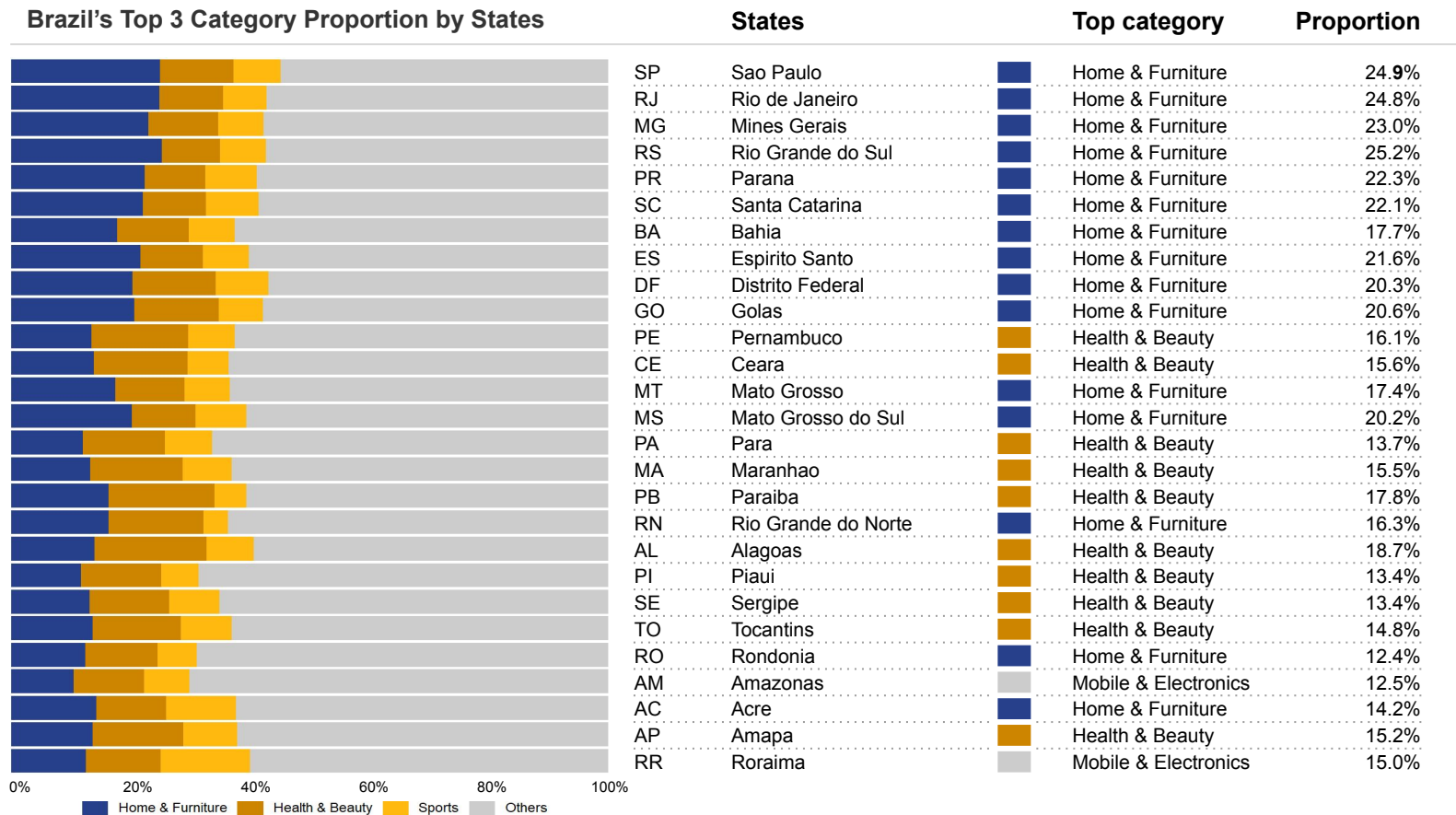
1

Customers: Home & Furniture is the most preferred, generating ~27.5k orders, followed by Health & Beauty, at ~13.1k orders



Customers: Most popular products are Home & Furniture in 15 states and Health & Beauty in 10 states, takes up to 25% of orders

Brazil's Top 3 Category Proportion by States

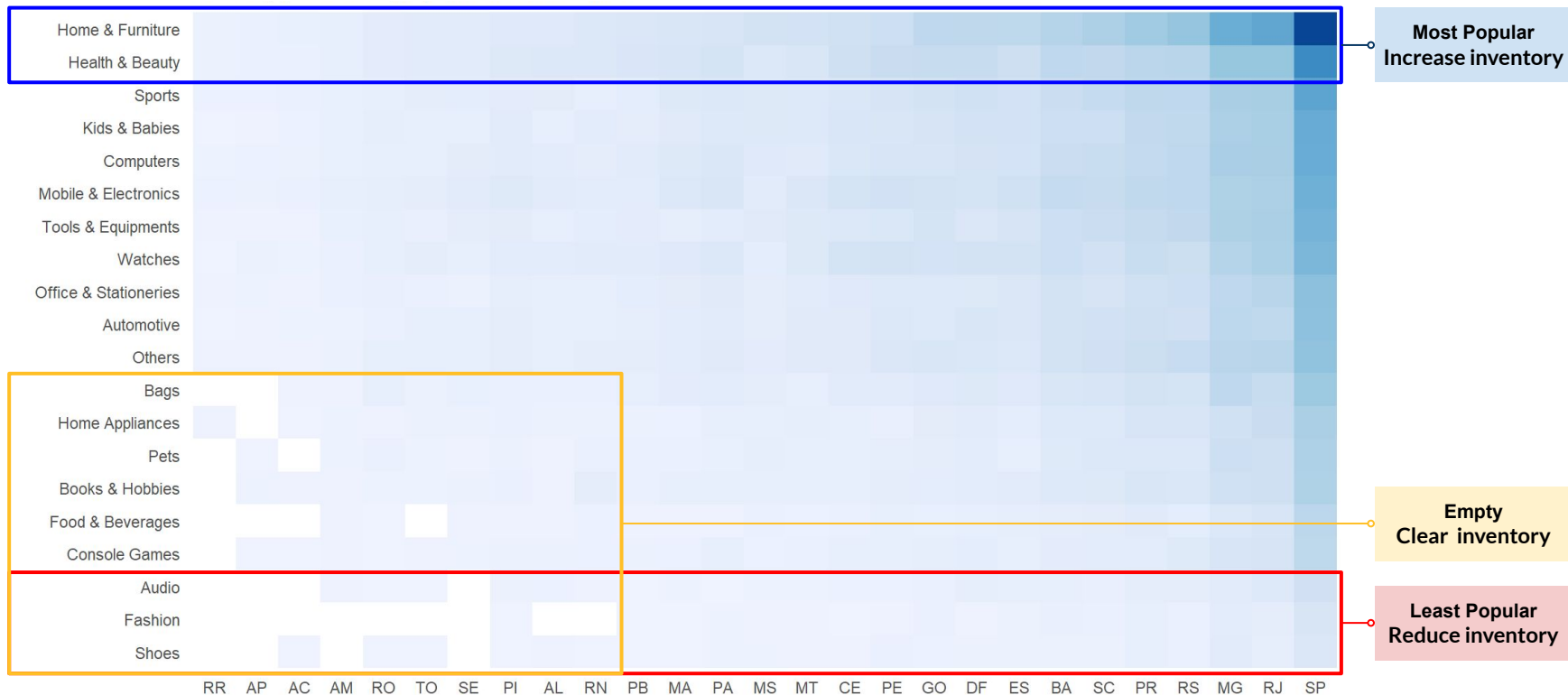


1

Customers: Similar preferences across states, most popular being Home & Furniture, and the least popular Audio, Fashion and Shoes

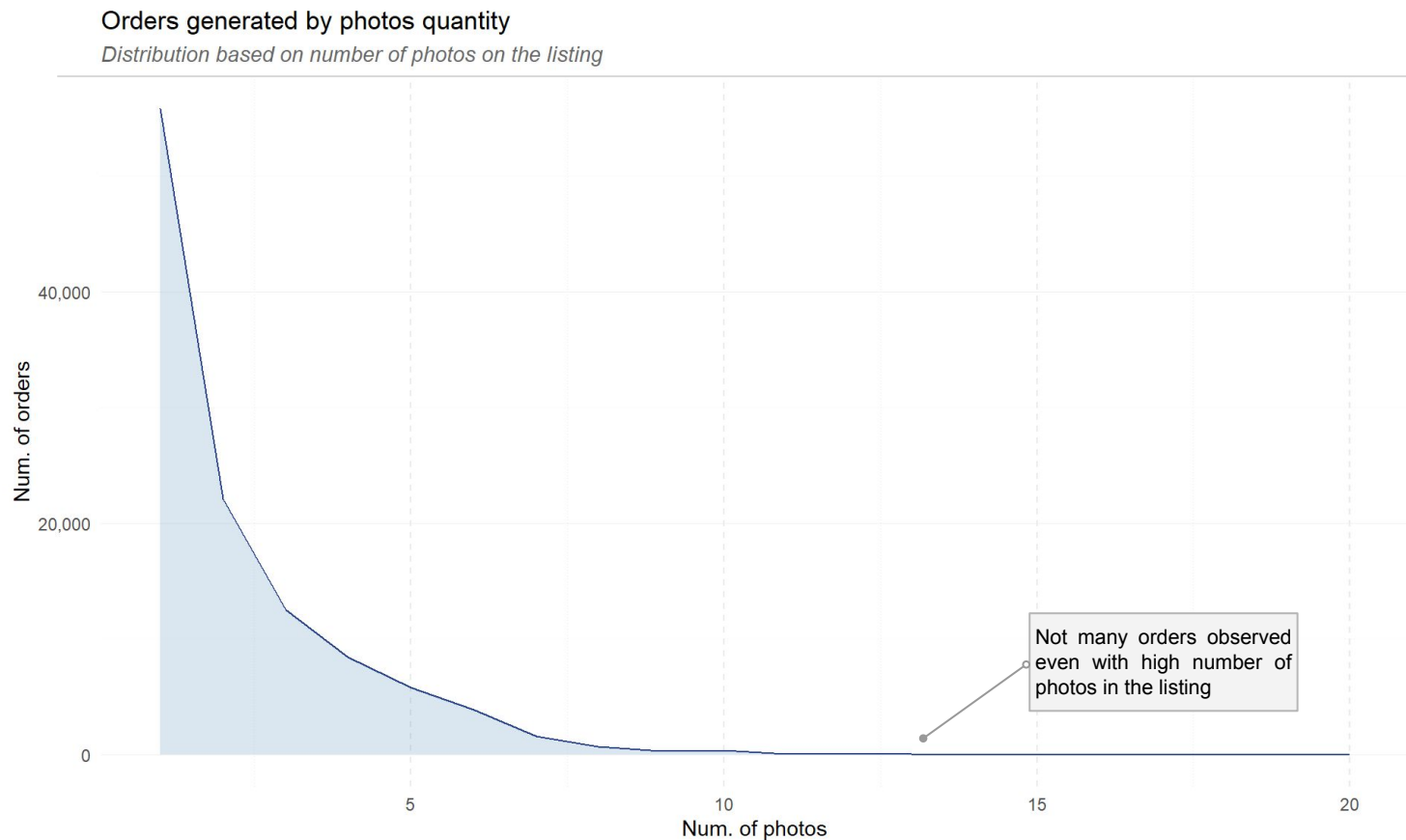
Total Orders Generated

Orders by product category and customer's state

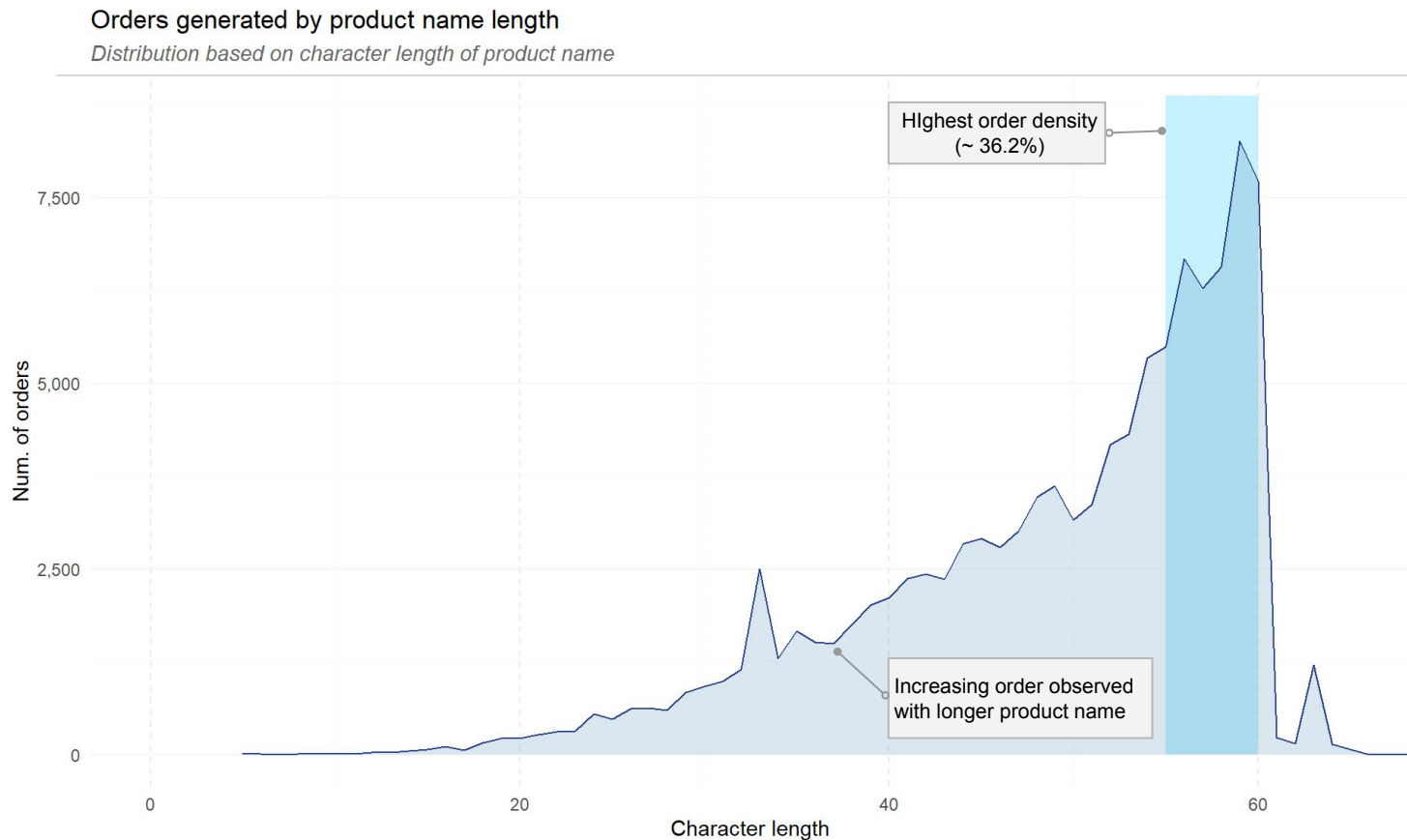


Data: Olist Public Dataset Feb 2017 to Aug 2018

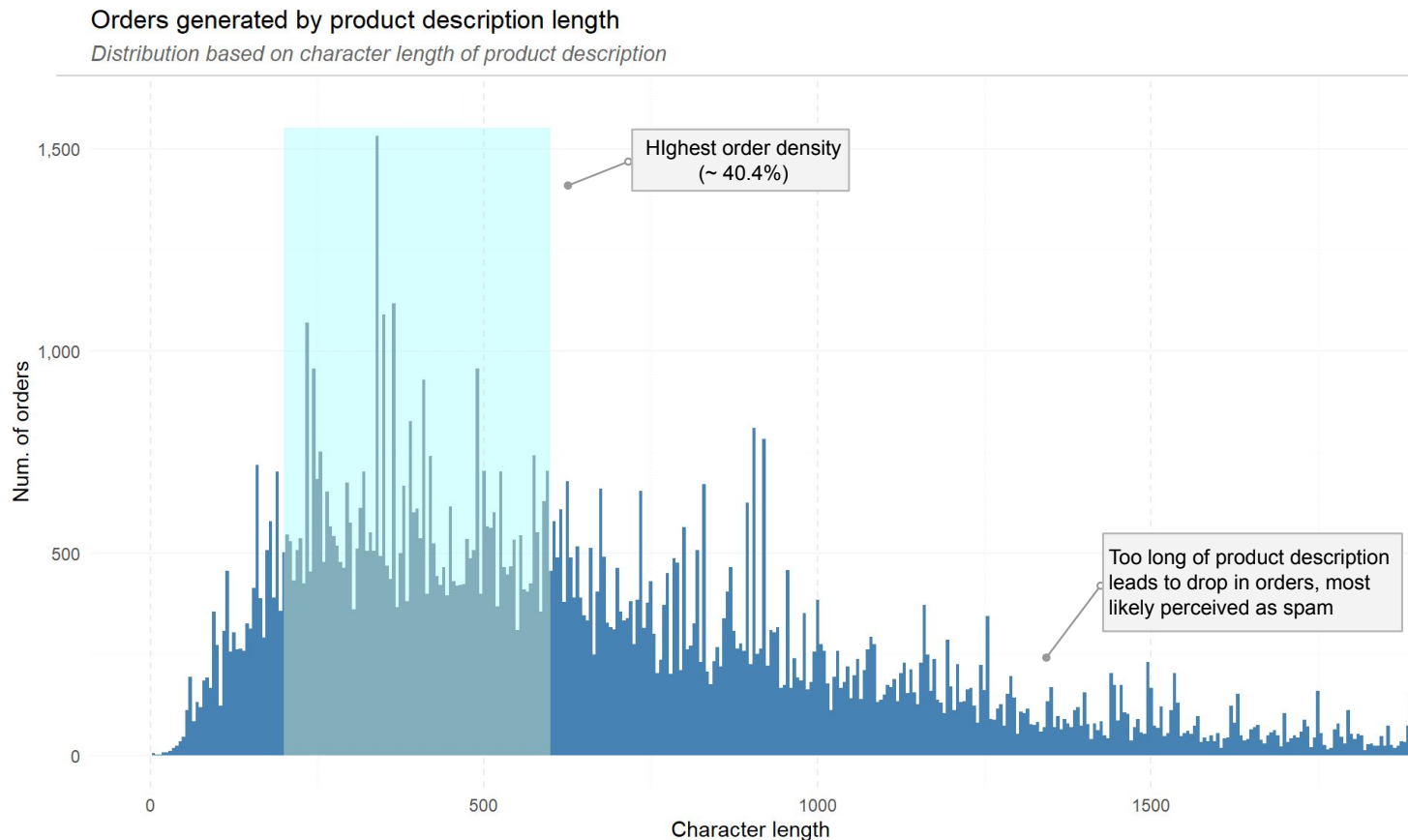
Orders: Increasing photos quantity does not affect the number of orders generated



Orders: Increase in orders with longer product name, highest orders observed at approximately 55 to 60 product name length



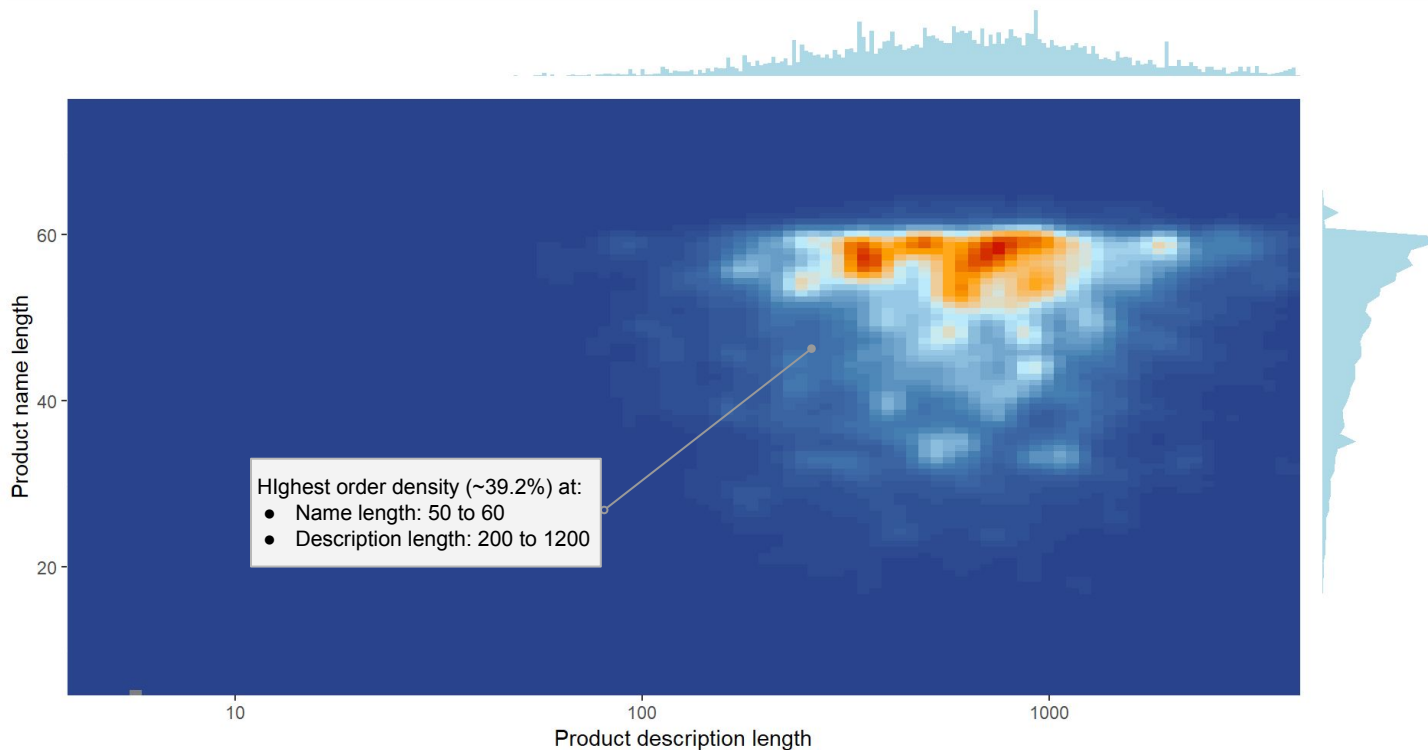
Orders: Highest orders observed at around 200 to 600 product description length, and drops when description gets too long



Orders: Highest order density observed when length of product name between 50-60 and description between 200-1200

Orders density by product attributes

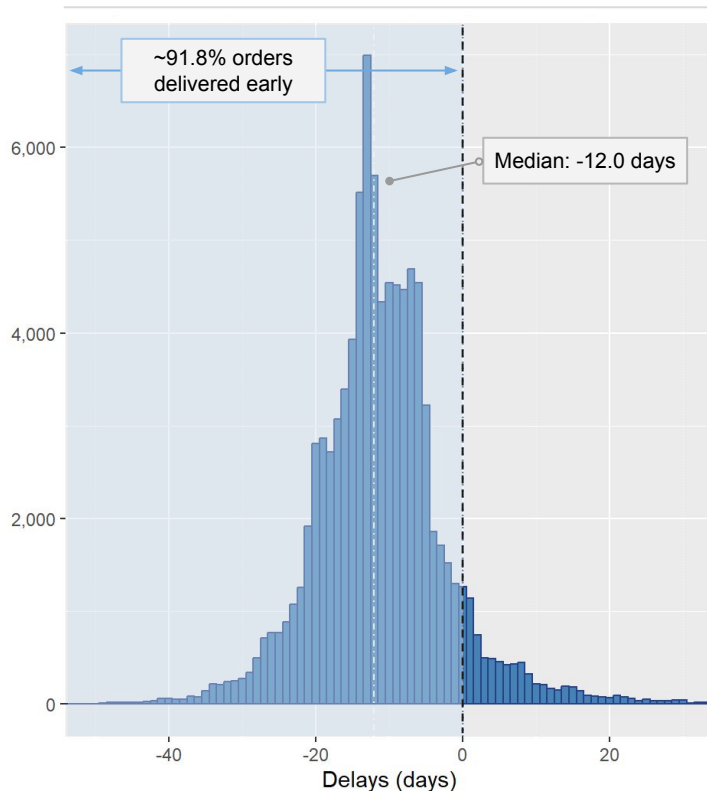
Based on product name and description length



Reviews: About ~92% orders arrived earlier than estimated, with median at ~12 days early; Median delivery lead time at ~10 days

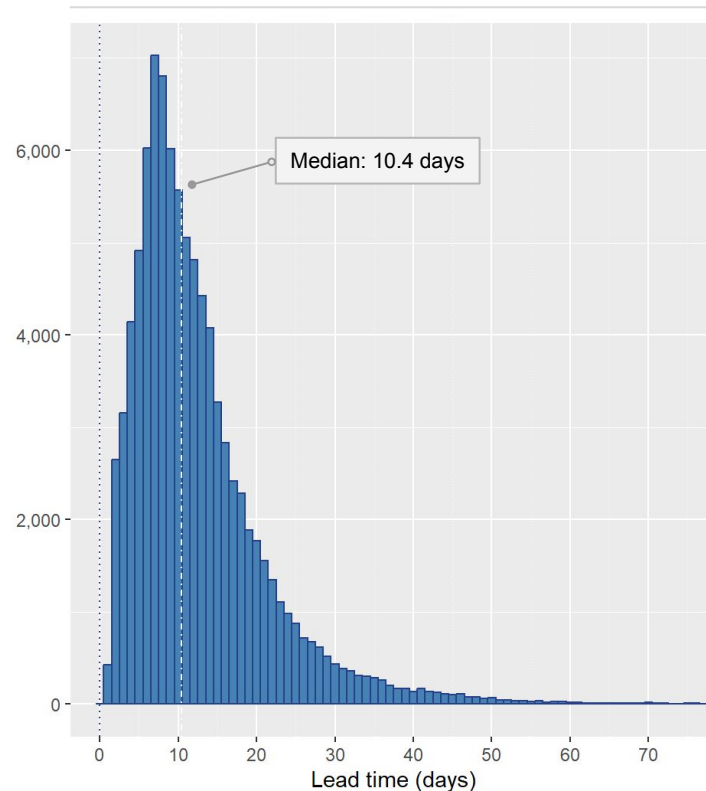
Distribution of delivery delays (in days)

Delays calculated from estimated delivery date, in days

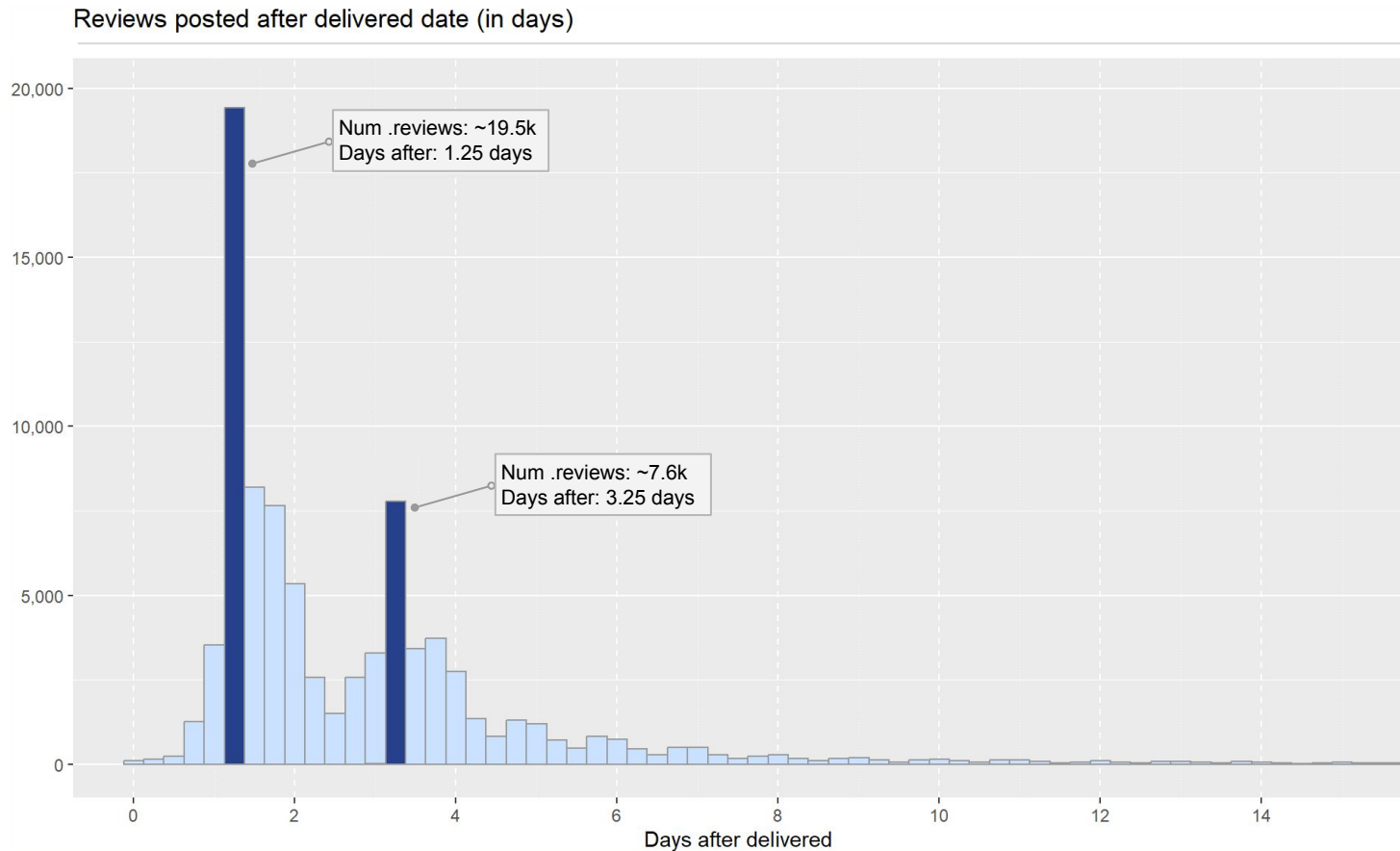


Distribution of delivery lead time (in days)

Lead time from date of purchase, in days



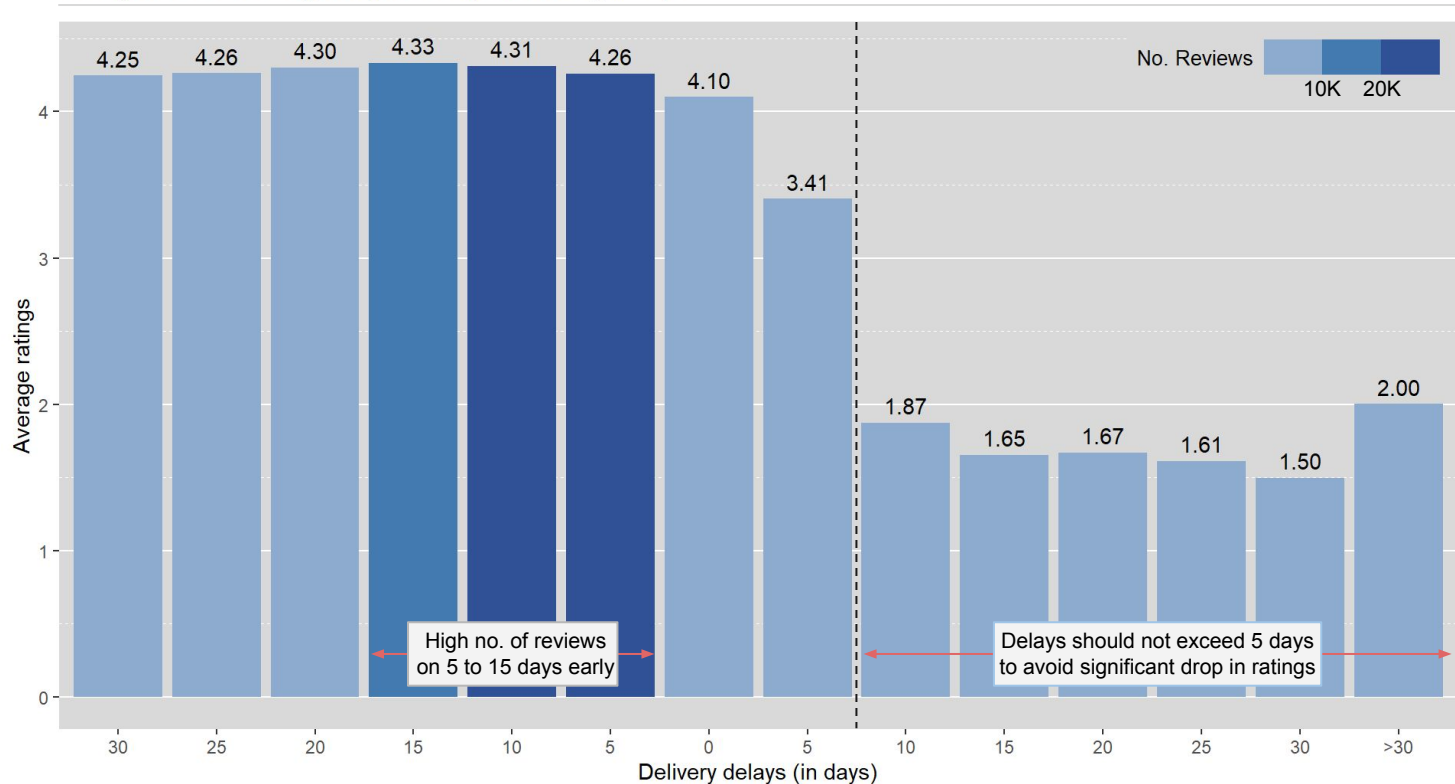
Reviews: Around 68% customers leave a review by day +3.25 after receiving the orders, peaking at +1.25 and +3.25 days post-delivery






Reviews: Review scores ranges between 4.1 - 4.3 when delivery is early, but drops to 1.5 - 2 when delayed longer than 5 days

Review score against delays (in days)

Average review scores against given delays in delivery, in days






Conclusion: Sellers and investors to consider recommendations derived from the key insights over customer, review and orders

Analysis Scope	Key Insights Gained	Recommendations
1  Customer	<ul style="list-style-type: none">• No big difference in preferences across states• Highest demand on Home & Furniture (H&F) and Health & Beauty (H&B) products• Lowest demand on Audio, Fashion and Shoes products	<ul style="list-style-type: none">• Sellers should expand or increase focus on selling H&F and H&B products• Minimise focus on selling Audio, Fashion and Shoes products
2  Order	<ul style="list-style-type: none">• Quantity of photos do not affect customer's purchasing decision• Highest order density observed on long product name (50 to 60) and medium description length (200 to 1200)	<ul style="list-style-type: none">• Sellers set up guidelines and best practices on product name and description for listing creation
3  Review	<ul style="list-style-type: none">• Majority of products ~92% arrived earlier than given estimated delivery time• Majority of customers left a review within ~3.25 days• Review ratings drops significantly when there are delays in delivery	<ul style="list-style-type: none">• Sellers should avoid delivery delays by partnering with reputable 3PLs and improve delivery estimation accuracy• Marketing campaigns to promotes timely feedback

Archived Slides I (from current slides)

Hypotheses: Aims to value-add our investors and sellers through key explorations on customer, review and orders data

Analysis Scope	Key Hypotheses / Questions	Actionable Insights
<div>1</div> <div> Customer</div>	<p>Customers have different product type preference in different states</p> <ul style="list-style-type: none">• States with high percentage of younger population have higher demand on electronics,• States with high percentage of older population have higher demand on health products	<p>Insight used to support sellers in finding the ideal ecommerce product mix for different target states</p>
<div>2</div> <div> Review</div>	<p>Fast-delivery products lead to better ratings and review scores</p> <ul style="list-style-type: none">• Many factors affect reviews, among which, low priced and faster delivery time may be the most important	<p>Advocate more reliable 3PL outsourcing to maintain their reputation and quality of commerce</p>
<div>3</div> <div> Order</div>	<p>More orders generated when products have complete attributes, i.e. longer titles, detailed descriptions and more photos quantity</p> <ul style="list-style-type: none">• More information on the products allow customers to make more confident purchases	<p>Provide guidelines for sellers to potentially generate more orders</p>

① Hypothesis: Product preferences are different across states

- Home & Furniture is the most popular category in Brazil, generating ~27.5k orders over 1.5-year period
- Home & Furniture is the most popular product category in 15 states, while Health & Beauty in 10 states
- Preferences are similar across states, with the top being Home & Furniture and Health & Beauty, while the least popular are Audio, Fashion and Shoes

② Hypothesis: Reviews are significantly affected by delivery efficiency

- 91.8% of orders arrived earlier than the estimated delivery time, with an average early arrival time of 13.4 days and median of 12.0 days
- 50% of shipments were received by customers within 10.4 days from date of purchase
- 68 % of customers tend to post reviews in 1.25 or 3.25 days after receiving their goods
- Review ratings average ~4.2 when delivered early, but drops to 3.4 when up to 5 days late and drops to approximately 1.7 when more than 5 days

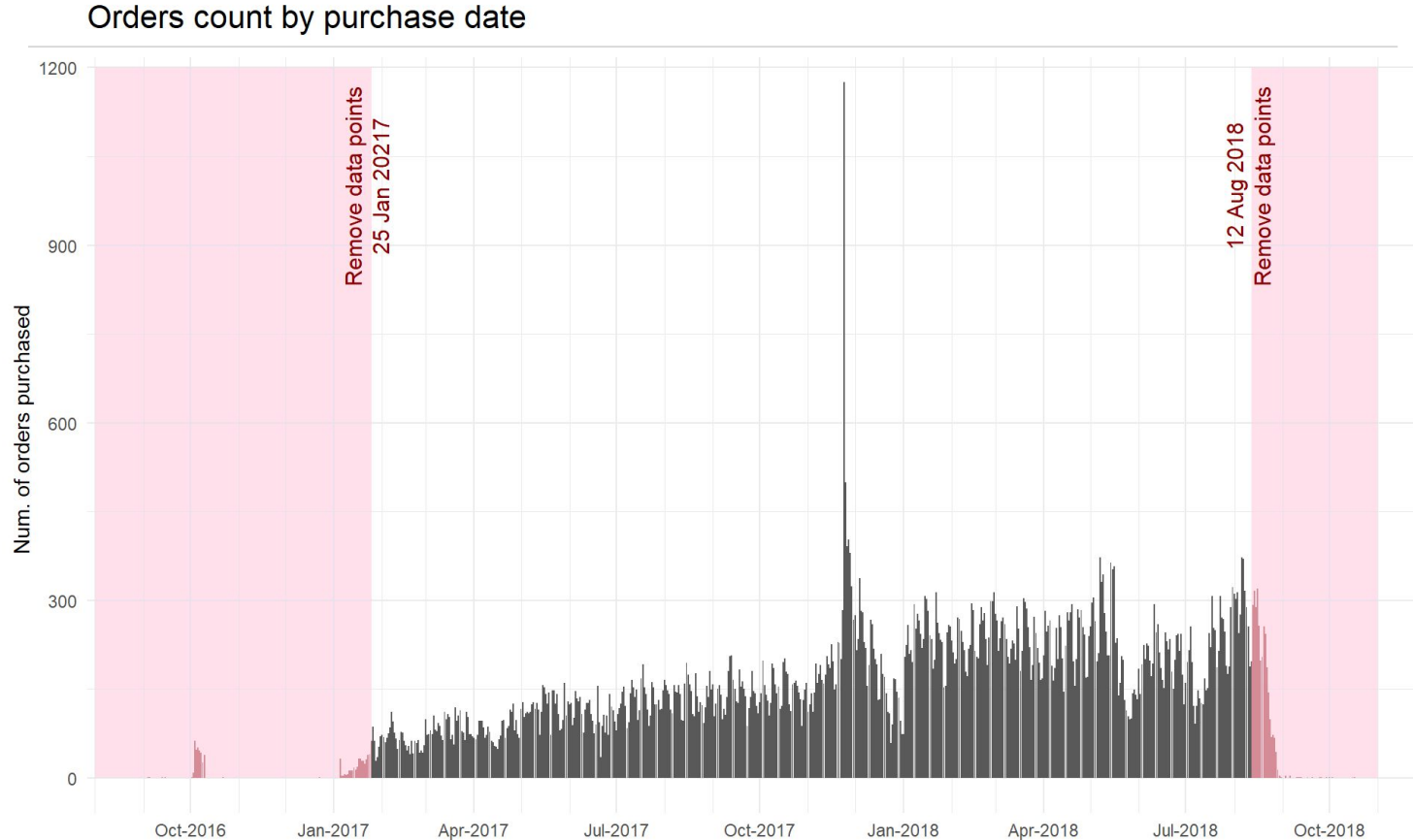
③ Hypothesis: Higher orders on products with greater details

- The number of photos included in the listing does not necessarily affect the orders generated, as we observe high orders with single-photo products
- Orders are generally higher with long product names around 50 to 60 characters and product descriptions at medium length at around 200 to 1200
- Product descriptions that are too long lead to lower orders

Conclusion & Recommendations

- **Customer:** Expand on selling **H&F and H&B products** in states like SP, RJ and MG, and avoid **Audio, Fashion and Shoes products** in states like RR, AP and AC to gain customer traction
- **Reviews:** Partner with reputable 3PLs and ensure proper pick-pack-and-ship SOP to avoid delivery delays
- **Orders:** Guidelines to ensure product name and description are within the desirable length

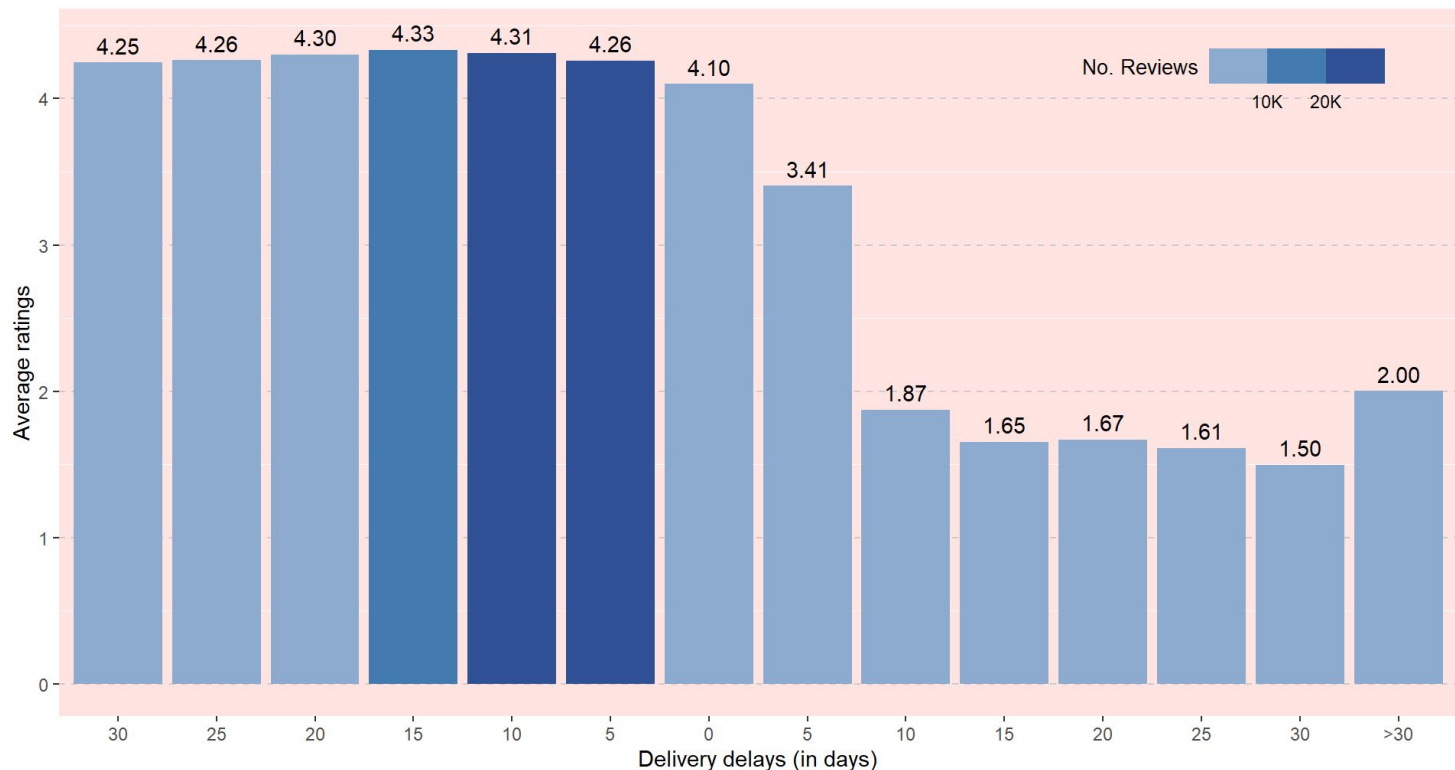
Orders Data: Suspicious order counts before 25 Jan 2017 and after 12 Aug 2018; to remove data points outside the time period (~4.0%)



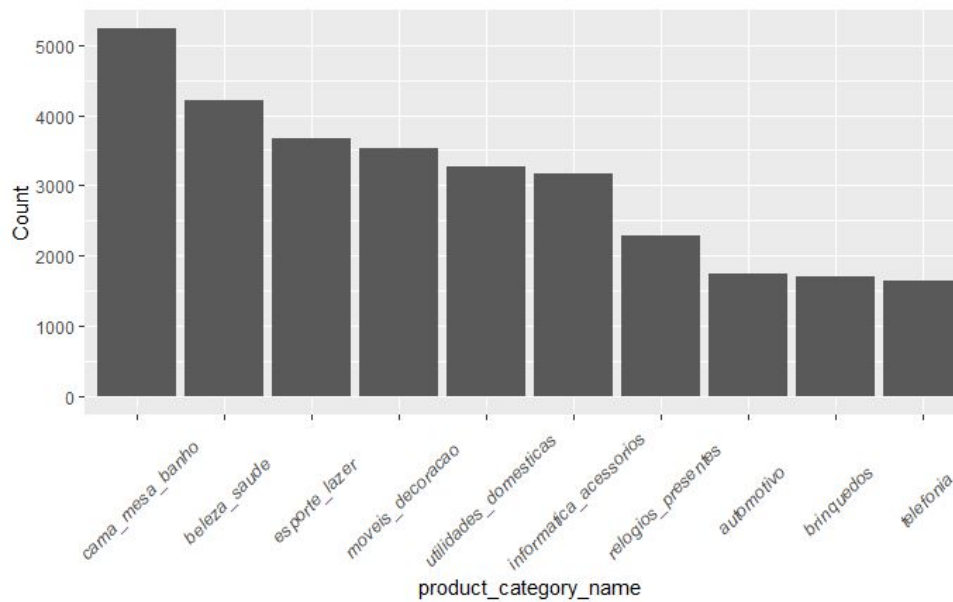
Insights: Review scores ranges between 4.1 - 4.3 when delivery is early, but drops to 1.5 - 2 when delayed longer than 5 days

Review score against delays (in days)

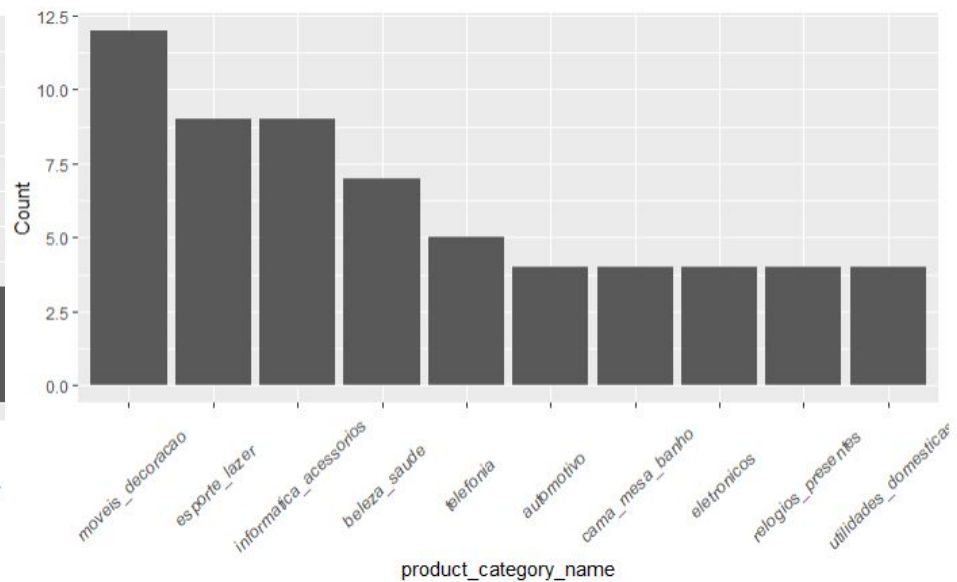
Average review scores against given delays in delivery, in days



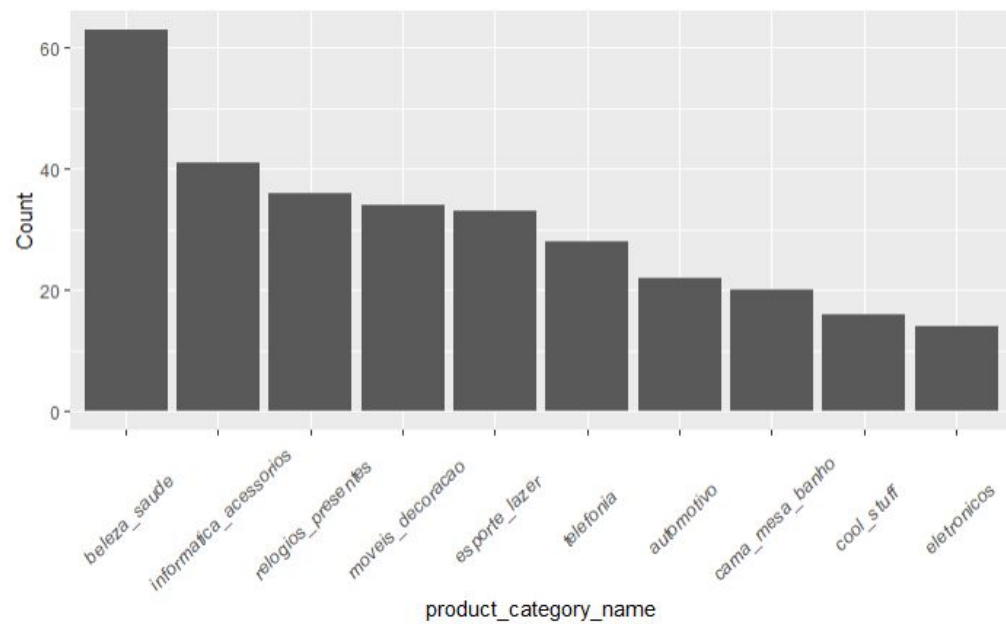
SP



AC



AL



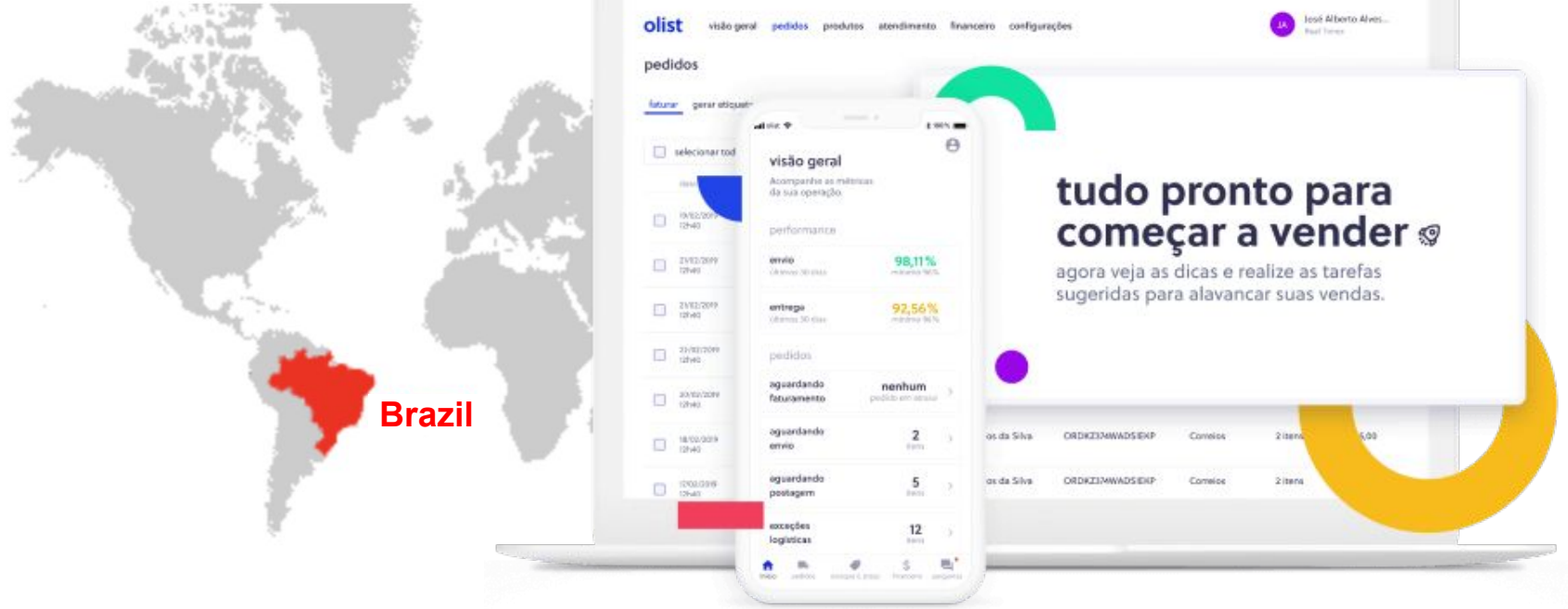
1. Overall most popular product

AC	AL	AM	AP	BA	CE	DF	ES	GO
furniture_de cor	beleza_sau de	beleza_sau de	beleza_sau de	beleza_sau de	beleza_sau de	beleza_sau de	cama_mess _banha	cama_mess _banha
MA	MG	MS	MT	PA	PB	PE	PI	PR
beleza_sau de	cama_mess _banha	esporte_laz er	beleza_sau de	beleza_sau de	beleza_sau de	beleza_sau de	beleza_sau de	moveis_dec oracao
RJ	RN	RO	RR	RS	SC	SE	SP	TO
cama_mess _banha	beleza_sau de	beleza_sau de	esporte_laz er	cama_mess _banha	esporte_laz er	beleza_sau de	cama_mess _banha	beleza_sau de

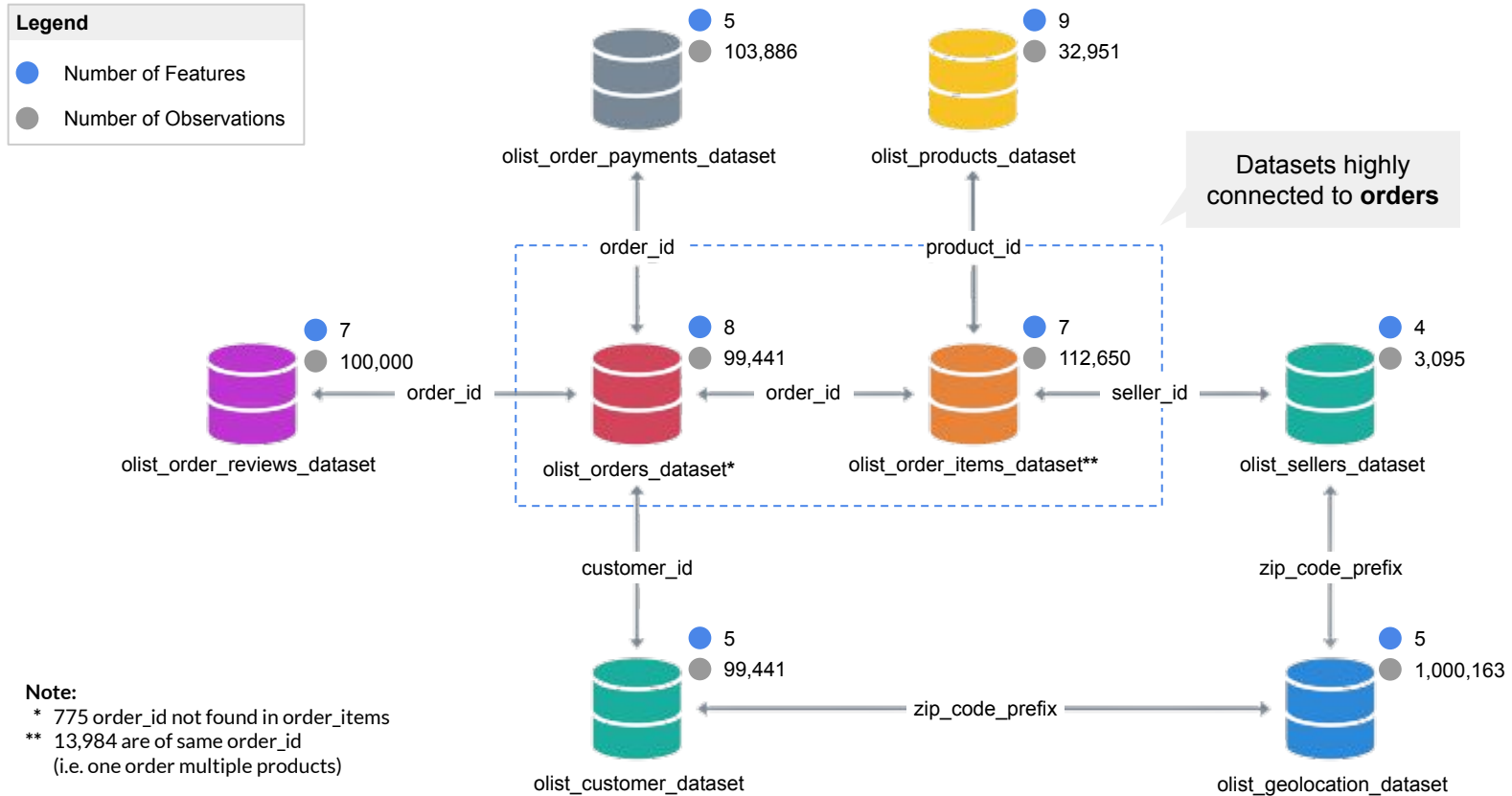
Archived Slides II

(from hypotheses proposal)

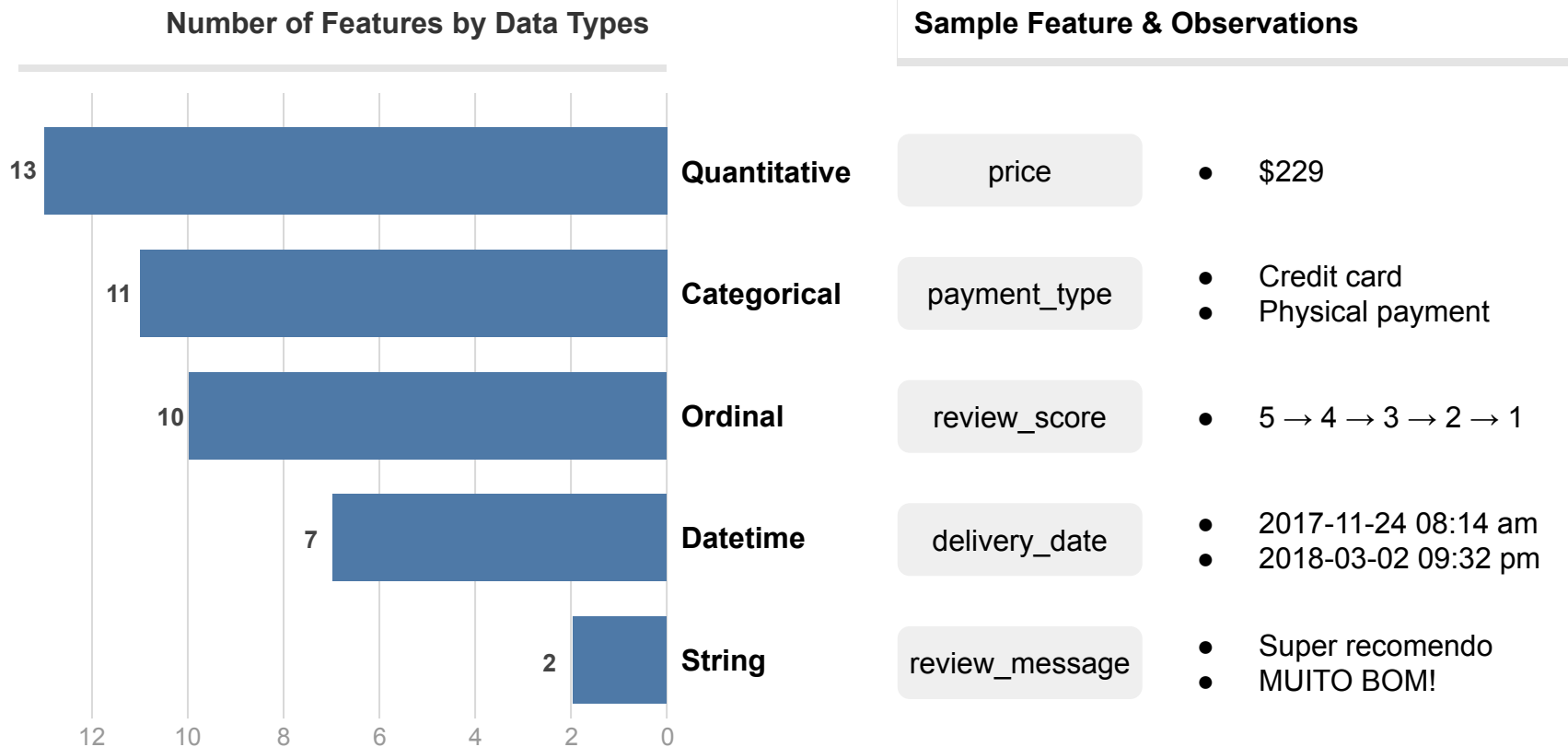
Context & Target Audience: Olist, a Brazilian e-commerce platform; target investors /sellers to improve performance through data



Data structure: Total 8 data sets focusing on customers, orders and products, but structure mainly centralised on orders



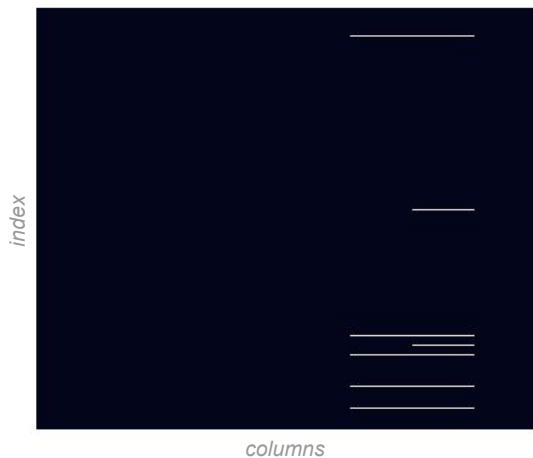
Data types: Total 5 types, with majority being quantitative data



Missing values: Found in 3 datasets through heatmaps; around 2-3% on orders and products, and around 89% on reviews

Orders

Around 3% missing data on *delivered dates*

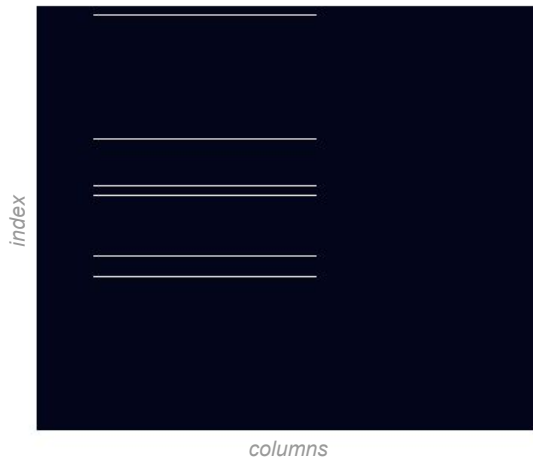


Potentially due to multiple reasons

1. **Canceled orders**
2. Incomplete data slicing
3. Software latency / API issue

Products

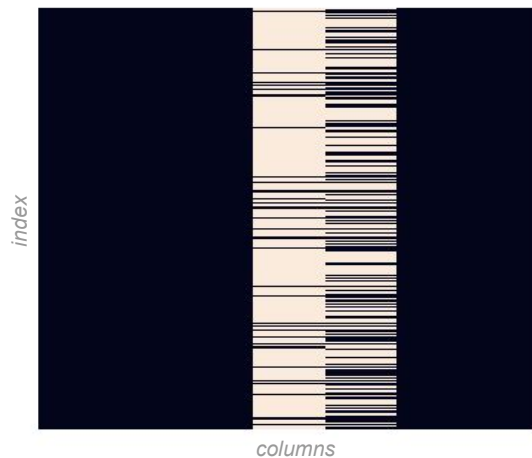
Around 2% missing data on *product category, title, description, image*



Potentially due to **incomplete/unpublished listings** commonly listed for testing purposes

Reviews

Around 89% missing data on *review titles and review messages*



Missing data expected as customers often leave ratings but do not write any review

Legend:

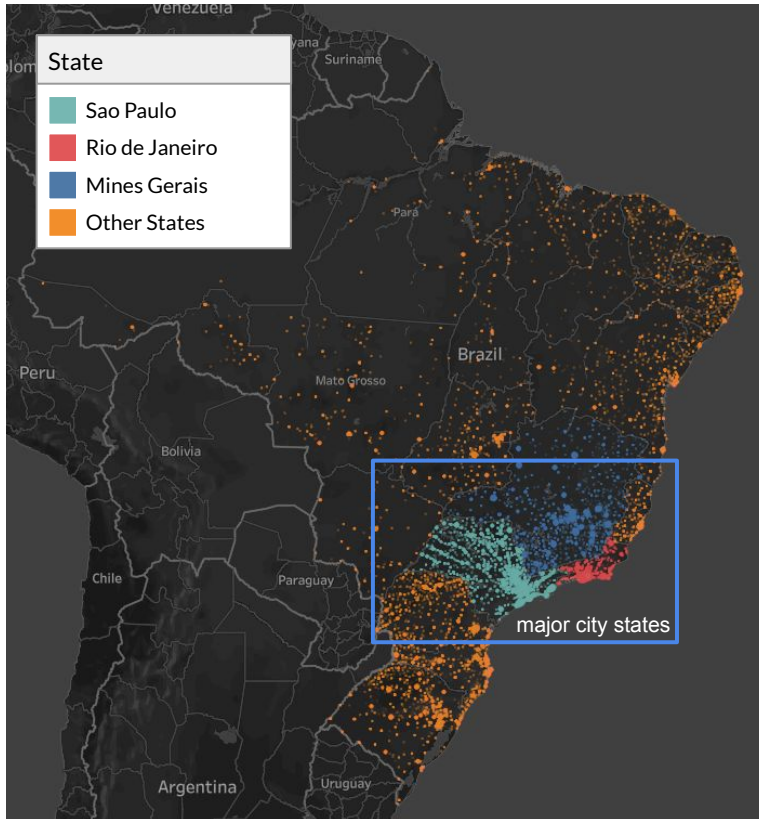


Null values

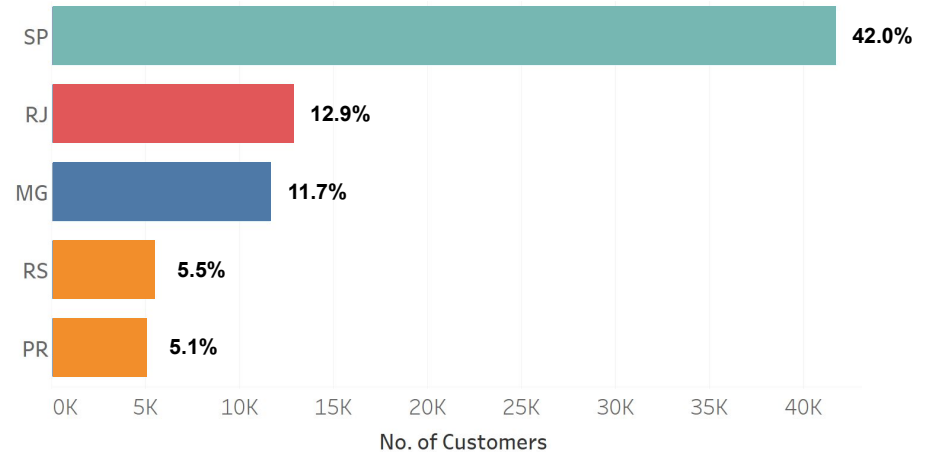


Non-null values

Customer: Buyer profiles heavily concentrated in major city states such as Sao Paulo, Rio de Janeiro and Minas Gerais



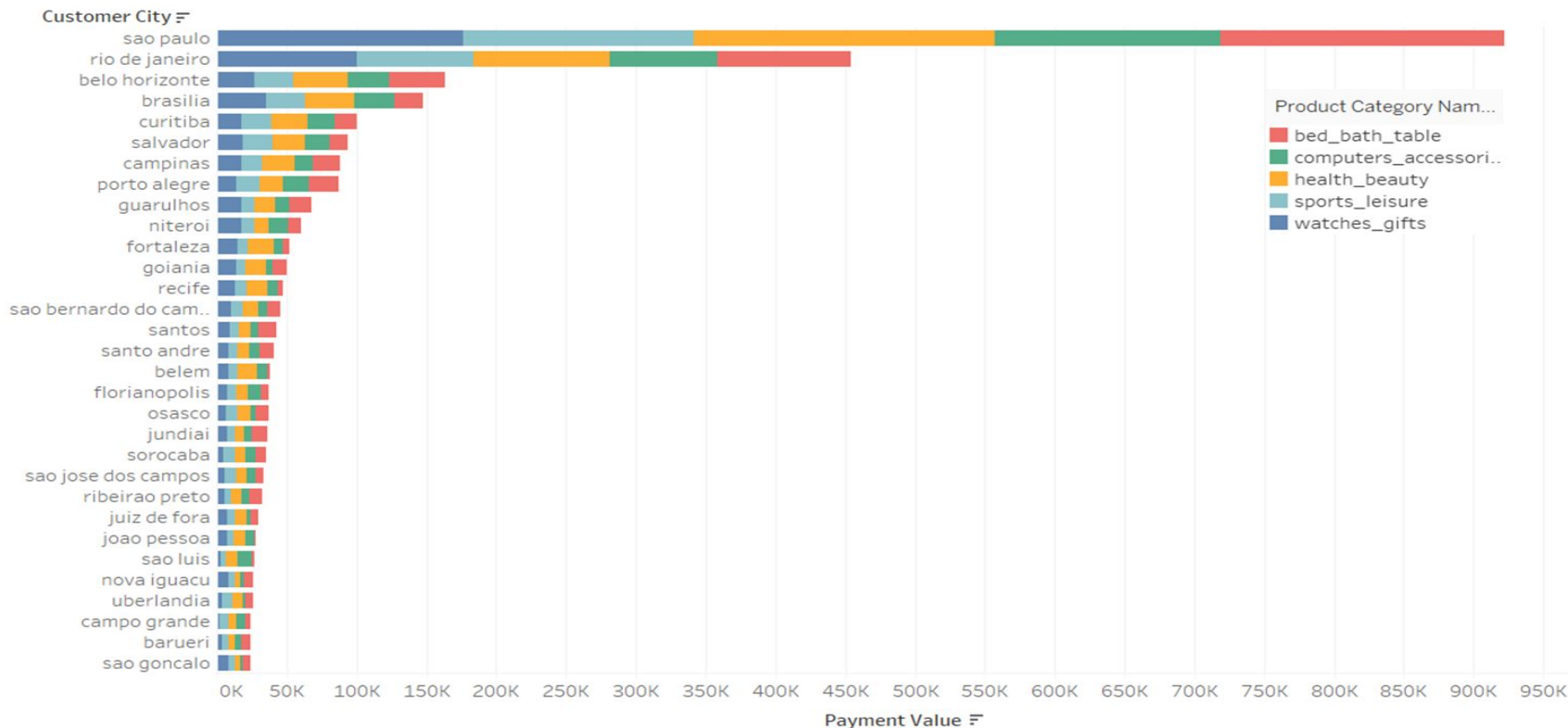
Customer distribution across the top 5 states in Brazil
(no. of customers in thousands)



Top 5 states out of 27 in Brazil make up to **~77.2%** of the total customers, with the highest density in **Sao Paulo at ~42.0%**

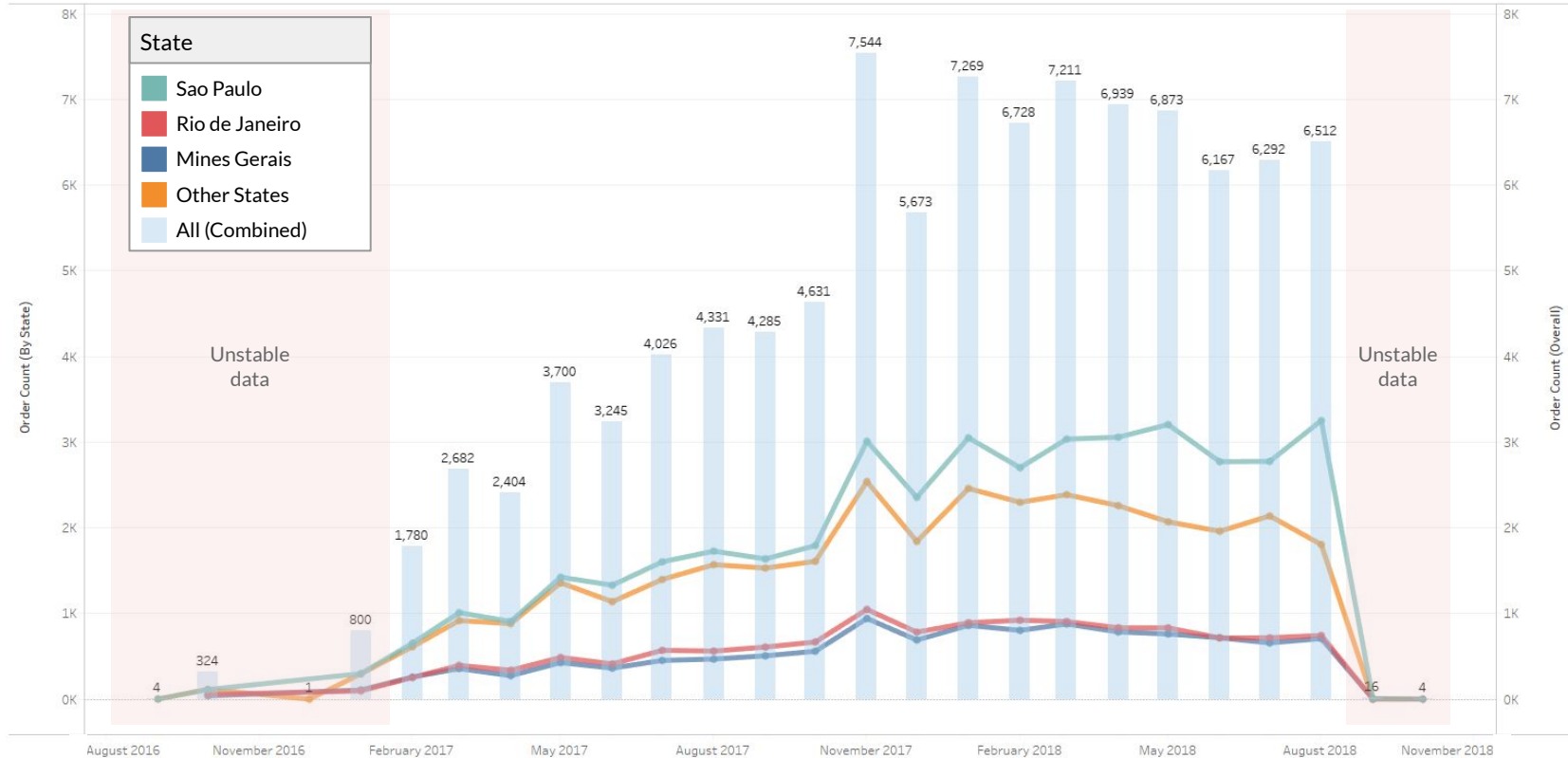
Customer: Brazilian consumers in different states have highly similar preferences for types of goods, but rank them differently

Customer Purchasing Preferences on Top 5 Categories (by States)



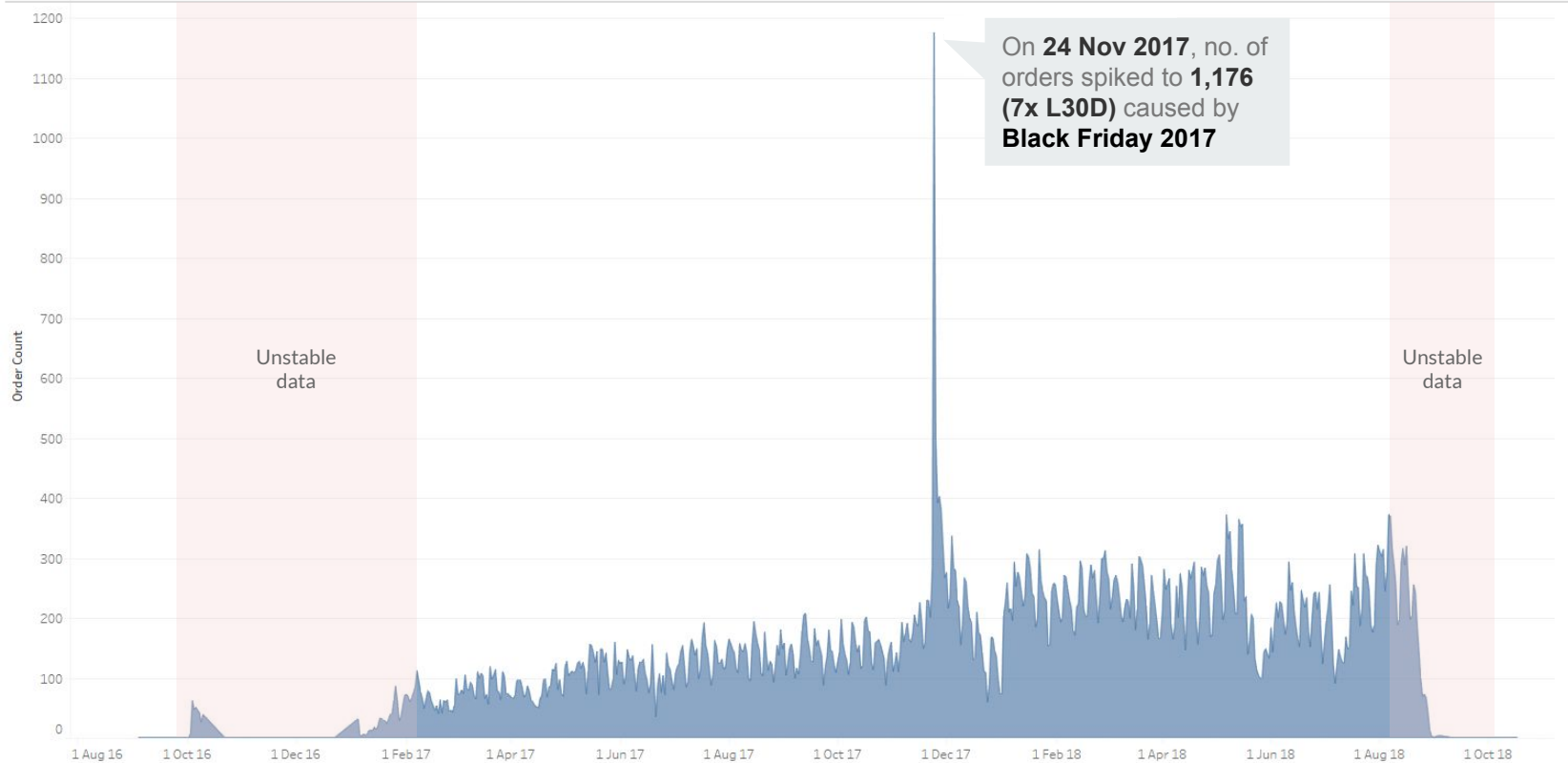
Orders: Upward trends observed on monthly orders; unreliable order counts before Feb 2017 and after Aug 2018 to be removed

Monthly Order Distribution (by states and combined)



Orders: Upon further investigation, order spike uncovered on the 24 Nov 2017 due to 2017 Black Friday

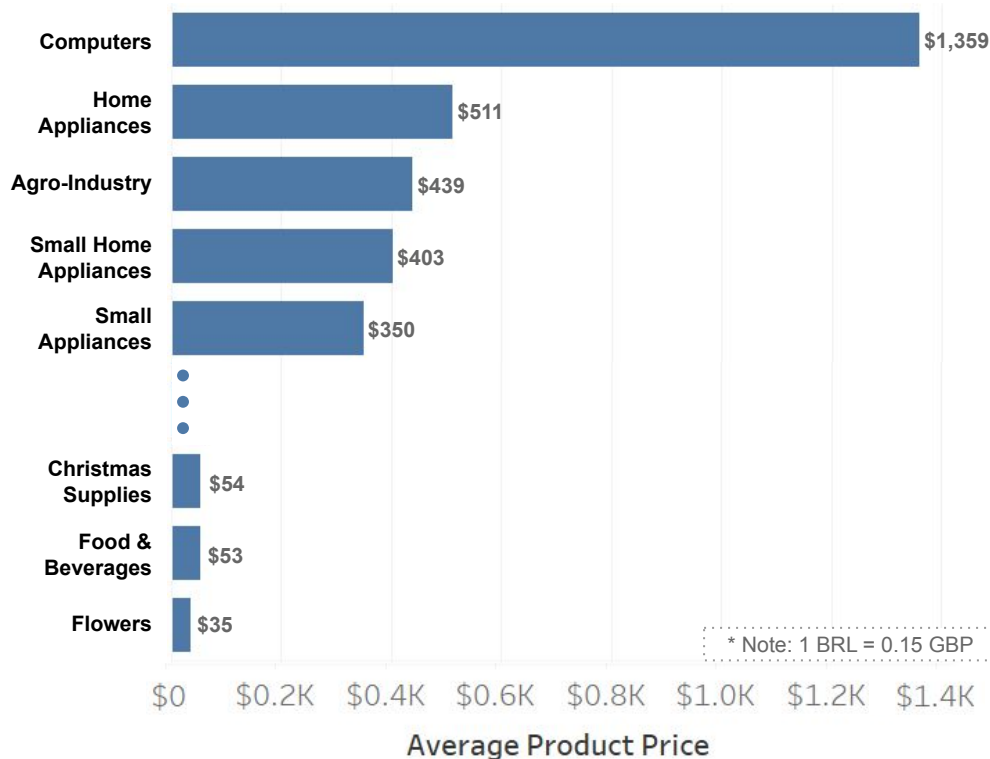
Daily Order Distribution (all states)



Products: High priced products dominated by electronic-related categories, further data cleaning required to properly categorise

Product price distribution across different categories

(Average price in thousands of Brazilian Reals)

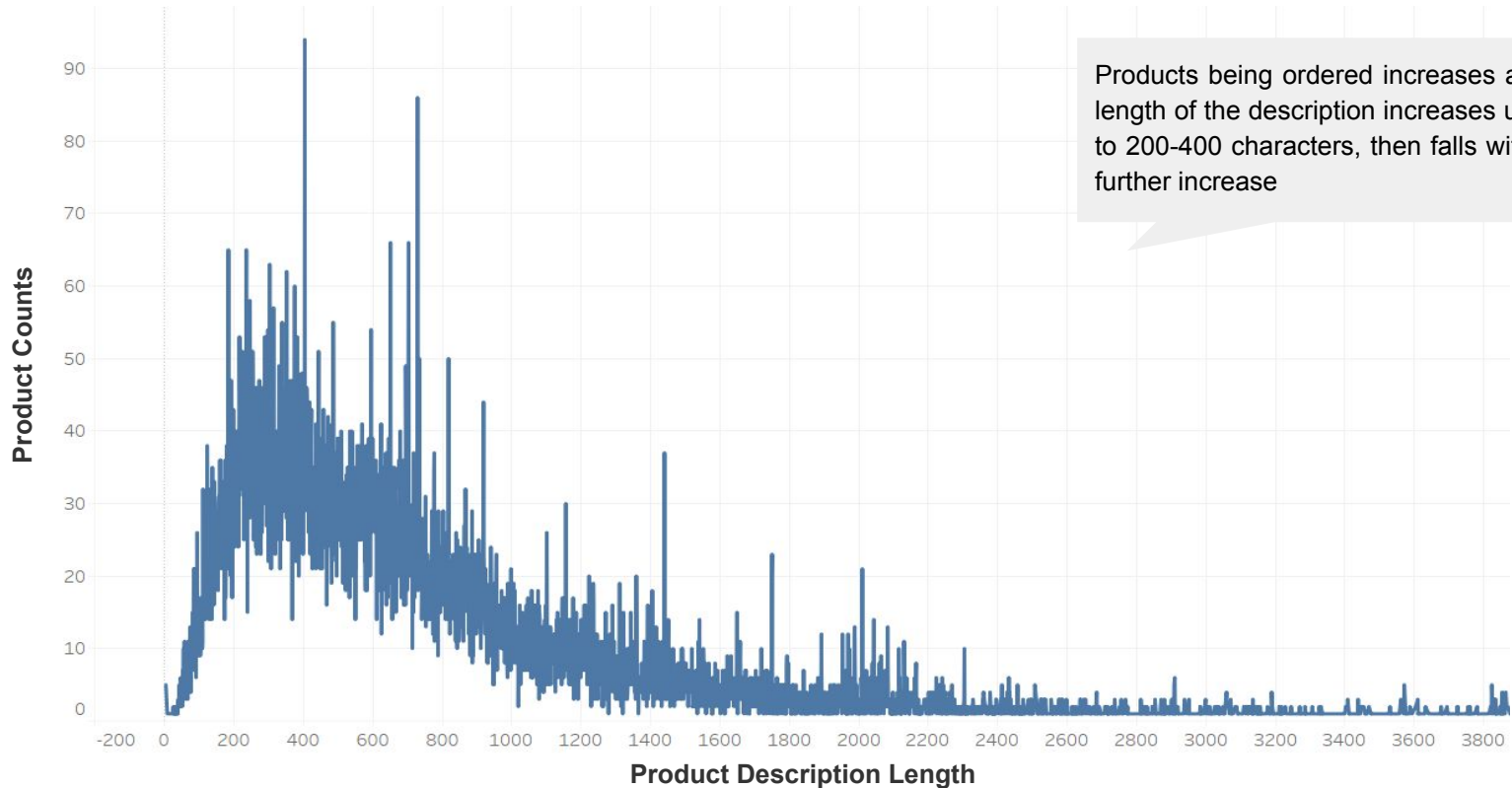


- Out of 77 available product categories, the highest priced are mostly **electronics**, with **computers** averaging at **~R\$1,359**
- Duplicated or similar categories observed, **data cleaning and re-categorising** will be required

Original product category	Cleaned
home_appliances	Home Appliances
home_appliances_2	
small_appliances	
small_appliances_oven_and_coffee	
fashion_sport	Fashion
fashion_male_clothing	
fashion_female_clothing	

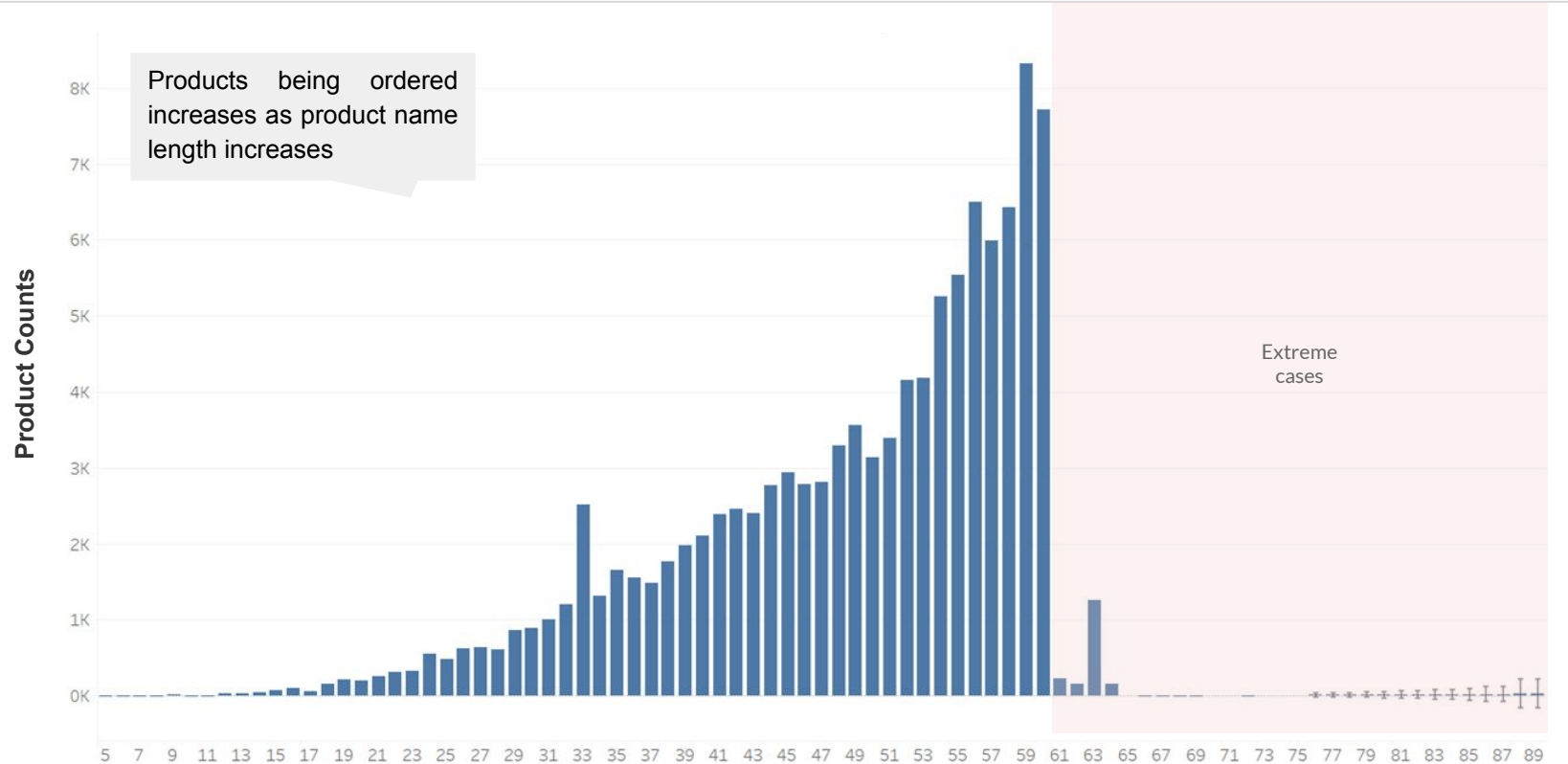
Products: Descriptions within 200-400 length range tend to generate more orders, but drops when too long

Distribution of Product Description Length



Products: Longer product names facilitate sales, further data cleaning required to eliminate effects of extreme cases

Distribution of Product Name Length



Derived features: Delivery lead time and freight-volumetric weight ratio can also be derived to better understand the order profiles

Delivery Lead Time

Delivered date - Purchased date



Freight-Volumetric Weight Ratio

Freight value ÷ (Volume × Weight)



Hypotheses: Aims to value-add our investors and sellers through key explorations on customer, review and orders data

Analysis Scope	Key Hypotheses / Questions	Actionable Insights
 Customer	<p>Customers prefer certain product types more in different city/ states</p> <ul style="list-style-type: none">• States with high percentage of younger population have higher demand on electronics,• States with high percentage of older population have higher demand on health products	Insight used to support sellers in finding the ideal product mix for different target states
 Review	<p>Customers tend to give better ratings and reviews on low-priced, fast-delivery products</p> <ul style="list-style-type: none">• Many factors affect reviews, among which, low priced and faster delivery time may be the most important	Advocate competitive pricing to sellers and reliable 3PL outsourcing to maintain their reputation and quality of commerce
 Order	<p>Products with longer titles, descriptions and more image quantity tend to generate more orders</p> <ul style="list-style-type: none">• More information on the products allow customers to make more confident purchases	Provide guidelines for sellers to potentially generate more orders