

# Quantium Retail Analytics Final Report

## Customer Patterns and Trial Layout Performance

Nguyen Tran

2023-06-20

### Table of contents

<b>1</b>	<b>Abstract</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>2</b>
<b>3</b>	<b>Methods</b>	<b>3</b>
3.1	Task 1 . . . . .	3
3.2	Task 2 . . . . .	9
3.3	Task 3 . . . . .	18
<b>4</b>	<b>Conclusions</b>	<b>18</b>
4.1	Task 1 . . . . .	18
4.2	Task 2 . . . . .	19
4.3	Task 3 . . . . .	20
<b>5</b>	<b>References</b>	<b>20</b>
<b>6</b>	<b>Appendix A: Raw Data</b>	<b>20</b>
<b>7</b>	<b>Appendix B: Magnitude distance</b>	<b>20</b>

Note: This is a [virtual experience project](#) provided by theforage.com, done by Nguyen Tran.  
For the full analysis with code, visit [Task 1](#) and [Task 2](#).

### 1 Abstract

As part of Quantum’s retail analytics team, we have been approached by our client Julia, the Category Manager for Chips, who wants to better understand the types of customers who

purchase Chips and their purchasing behaviour within the region. She has also asked us to test the impact of the new trial layouts with a data driven recommendation to whether or not the trial layout should be rolled out to all their stores. We generated insights and provided commercial recommendations based on the provided data sets, identified benchmark stores based on two statistical measures to test the impact of the trial store layouts on customer sales, and presented a comprehensive client report in PowerPoint, introducing the actionable results and then getting into the why's and how's. In completion, we communicated statistical evidence to client, confirming the success of the trial layout, and recommended its deployment across all stores.

## 2 Introduction

Quantium has had a data partnership with a large supermarket brand for the last few years who provide transaction and customer data. We are responsible for delivering highly valued data analytics and insights to help the business make strategic decisions.

Supermarkets will regularly change their store layouts, product selections, prices and promotions. This is to satisfy their customer's changing needs and preferences, keep up with the increasing competition in the market or to capitalize on new opportunities. The Quantum analytics team is engaged in these processes to evaluate and analyse the performance of change and recommend whether it has been successful.

Our client Julia has provided two data sets: customer transaction data and purchase behavior, spanning from July 2018 to June 2019 (to see what the data sets look like, visit [Appendix A: Raw Data](#)). We need to present a strategic recommendation to Julia that is supported by data which she can then use for the upcoming category review however to do so we need to analyse the data to understand the current purchasing trends and behaviors. The client is particularly interested in customer segments and their chip purchasing behavior. She has also asked us to evaluate the performance of a store trial which was performed in stores 77, 86 and 88, and provide data-driven recommendations on whether they should be implemented across all their stores. The trial period is from February 2019 to the end of April 2019.

Our approach for this client problem is broken down into three tasks:

### 1. Data preparation and customer analytics

Conduct analysis on client's transaction data set and identify customer purchasing behaviors to generate insights and provide commercial recommendations.

### 2. Experimentation and uplift testing

Extend analysis from Task 1 to identify benchmark stores that allow us to test the impact of the trial store layouts on customer sales.

### 3. Analytics and commercial application

Use analytics and insights from Task 1 and 2 to prepare a report for the client.

## 3 Methods

In this section, we go over what we did for each task.

### 3.1 Task 1

#### 3.1.1 Data Cleaning

The first step in any analysis is to first understand the data. We looked at each of the data sets provided and found some data cleaning to do before exploring.

**Examining transaction data:** identified inconsistencies, no missing data, one outlier where 200 packets of chips are bought in one transaction (this customer has only had the two transactions over the year and is not an ordinary retail customer; the customer might be buying chips for commercial purposes instead; Thus, we removed these two transactions from further analysis), and items other than chips in product's name column such as salsa products.

**Examining customer data:** checked for similar issues in the customer data, no missing data, merged the transaction and customer data together so it's ready for the analysis.

Next, we dived into exploring the data.

#### 3.1.2 Analysis on customer segments

- Defined the metrics of interest to the client

- Who spends the most on chips (total sales), describing customers by lifestage and how premium their general purchasing behavior is?
- How many customers are in each segment?
- How many chips are bought per customer by segment?
- What's the average chip price by customer segment?

- Raised questions to our data team
  - The customer's total grocery spend over the period and total spend for each chips transaction to understand what proportion of their grocery spend is on chips.
  - Proportion of customers in each customer segment overall to compare against the mix of customers who purchase chips.
- Created charts to visualize our findings

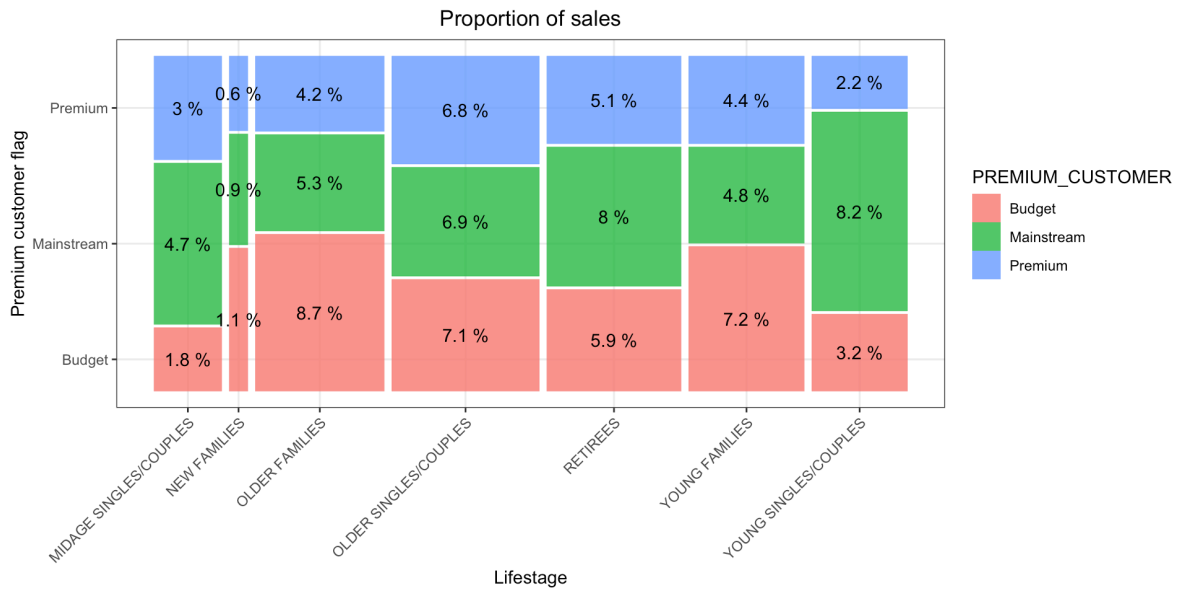
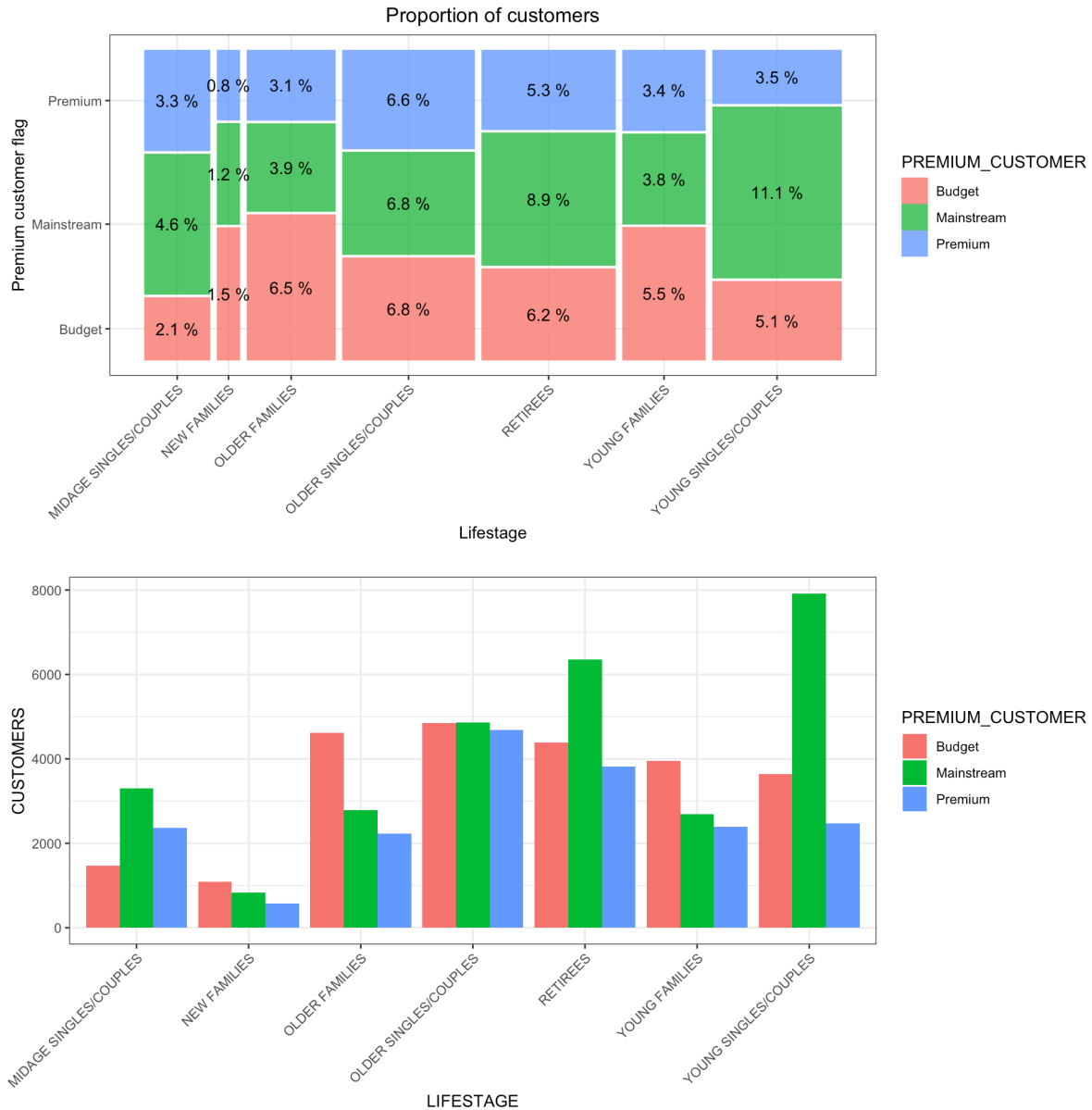


Figure 1: Note that the percentage of each tile is relative to the total sales.

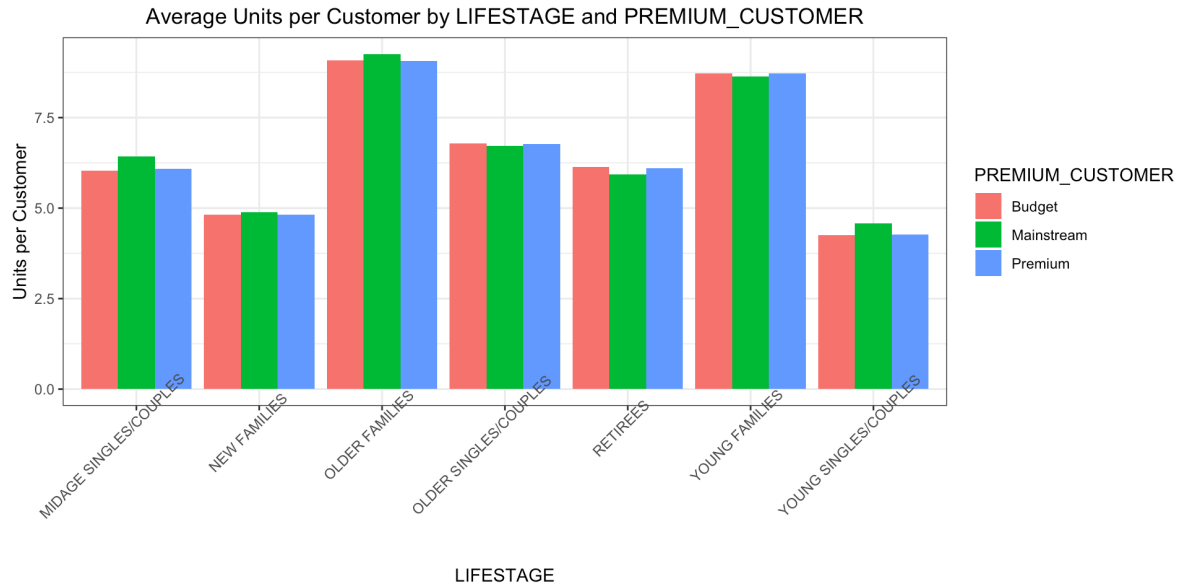
Sales are coming mainly from Budget - older families, Mainstream - young singles/couples, and Mainstream - retirees

Let's see if the higher sales are due to there being more customers who buy chips.



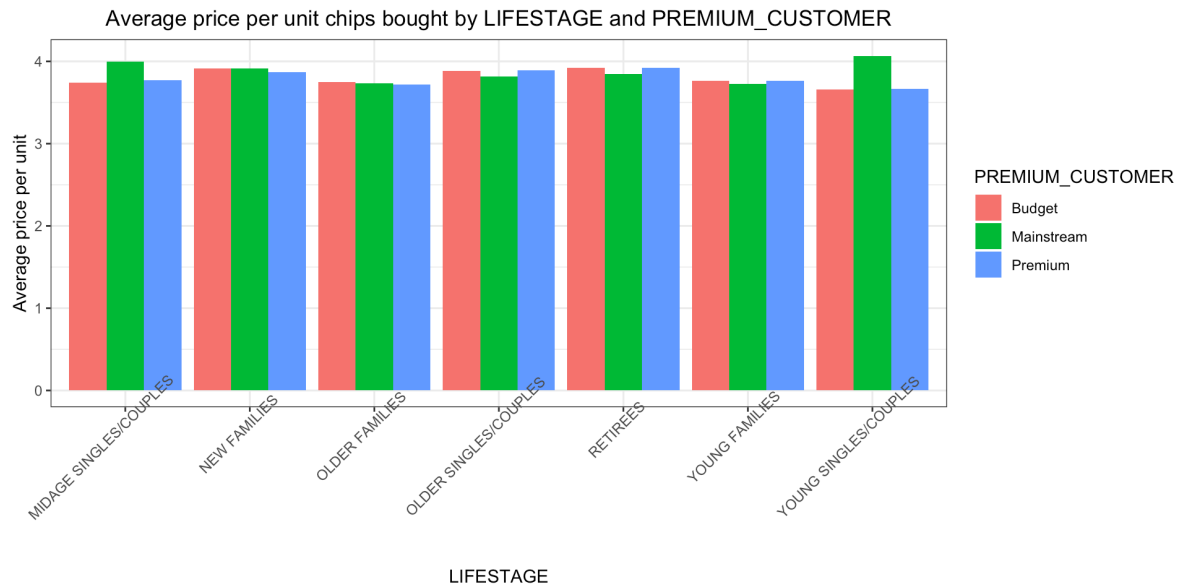
There are more Mainstream - young singles/couples and Mainstream - retirees who buy chips. This contributes to there being more sales to these customer segments but this is not a major driver for the Budget - Older families segment.

Higher sales may also be driven by more units of chips being bought per customer.



Older families and young families in general buy more chips per customer

What about the average price per unit chips bought for each customer segment as this is also a driver of total sales.



Mainstream midage and young singles and couples are more willing to pay more per packet of chips compared to their budget and premium counterparts. This may be due to premium shoppers being more likely to buy healthy snacks and when they buy chips, this is mainly for entertainment purposes rather than their own habit of consumption. This is also supported

by there being fewer premium midage and young singles and couples buying chips compared to their mainstream counterparts.

As the difference in average price per unit isn't large, we perform a t-test to see if this difference was statistically different.

Let  $\mu_m$  be the mean unit price of mainstream midage and young singles and couples Let  $\mu_{pb}$  be the mean unit price of premium and budget midage and young singles and couples

$H_0$  : true difference in  $\mu_m$  and  $\mu_{pb}$  is equal to 0

$H_a$  : true difference in  $\mu_m$  and  $\mu_{pb}$  is greater than 0

The t-test results in a p-value of 3.483677e-306, i.e. the unit price for mainstream, young and mid-age singles and couples ARE significantly higher than that of budget or premium, young and midage singles and couples.

### **3.1.3 Insights for specific customer segments**

We have found quite a few interesting insights that we can dive deeper into.

We might want to target customer segments that contribute the most to sales to retain them or further increase sales. Let's look at Mainstream - young singles/couples.

Do they tend to buy a particular brand of chips compared to other segments? We answered this with affinity analysis.

	BRAND	TargetSeg	otherSegs	AffinityToBrand
1:	Tyrrells	0.031552795	0.025692464	1.2280953
2:	Twisties	0.046183575	0.037876520	1.2193194
3:	Doritos	0.122760524	0.101074684	1.2145526
4:	Kettle	0.197984817	0.165553442	1.1958967
5:	Tostitos	0.045410628	0.037977861	1.1957131
6:	Infzns	0.014934438	0.012573300	1.1877898
7:	Pringles	0.119420290	0.100634769	1.1866703
8:	Grain	0.029123533	0.025121265	1.1593180
9:	Cobs	0.044637681	0.039048861	1.1431238
10:	Infuzions	0.049744651	0.044491379	1.1180739
11:	Thins	0.060372671	0.056986370	1.0594230
12:	Cheezels	0.017971014	0.018646902	0.9637534
13:	Smiths	0.096369910	0.124583692	0.7735355
14:	French	0.003947550	0.005758060	0.6855694
15:	Cheetos	0.008033126	0.012066591	0.6657329
16:	RRD	0.043809524	0.067493678	0.6490908
17:	Natural	0.015955832	0.024980768	0.6387246
18:	NCC	0.003643892	0.005873221	0.6204248
19:	CCs	0.011180124	0.018895650	0.5916771
20:	GrnWves	0.003588682	0.006066692	0.5915385
21:	Sunbites	0.006349206	0.012580210	0.5046980
22:	Woolworths	0.024099379	0.049427188	0.4875733
23:	Burger	0.002926156	0.006596434	0.4435967
	BRAND	TargetSeg	otherSegs	AffinityToBrand

Mainstream young singles/couples are 23% more likely to purchase Tyrrells chips compared to the rest of the population

Mainstream young singles/couples are 56% less likely to purchase Burger Rings compared to the rest of the population

Let's also find out if our target segment tends to buy larger packs of chips.



	PACK_SIZE	TargetSeg	otherSegs	AffinityToPack
1:	270	0.031828847	0.025095929	1.2682873
2:	380	0.032160110	0.025584213	1.2570295
3:	330	0.061283644	0.050161917	1.2217166
4:	134	0.119420290	0.100634769	1.1866703
5:	110	0.106280193	0.089791190	1.1836372
6:	210	0.029123533	0.025121265	1.1593180
7:	135	0.014768806	0.013075403	1.1295106
8:	250	0.014354727	0.012780590	1.1231662
9:	170	0.080772947	0.080985964	0.9973697
10:	150	0.157598344	0.163420656	0.9643722
11:	175	0.254989648	0.270006956	0.9443818
12:	165	0.055652174	0.062267662	0.8937572
13:	190	0.007481021	0.012442016	0.6012708
14:	180	0.003588682	0.006066692	0.5915385
15:	160	0.006404417	0.012372920	0.5176157
16:	90	0.006349206	0.012580210	0.5046980
17:	125	0.003008972	0.006036750	0.4984423
18:	200	0.008971705	0.018656115	0.4808989
19:	70	0.003036577	0.006322350	0.4802924
20:	220	0.002926156	0.006596434	0.4435967

Mainstream young singles/couples are 27% more likely to purchase pack size of 270g compared to the rest of the population.

Mainstream young singles/couples are 56% less likely to purchase pack size of 220g compared to the rest of the population.

Furthermore, Twisties is the only brand that offers 270g pack size. This brand also comes in 2nd in our Brand Affinity analysis. This contributes to why our target segment buys a lot of this pack size.

## 3.2 Task 2

For this part of the project we examined the performance in trial vs control stores to provide a recommendation for each location based on our insight.

### 3.2.1 Control Store Selections

The client has selected store numbers 77, 86 and 88 as trial stores and want control stores to be established stores that are operational for the entire observation period.

We would want to match trial stores to control stores that are similar to the trial store prior to the trial period of Feb 2019 in terms of :

- Monthly overall sales revenue
- Monthly number of customers
- Monthly number of transactions per customer

Thus, we created these metrics and filter to stores that are present throughout the pre-trial period.

Next, we needed to work out a way of ranking how similar each potential control store is to the trial store. The two metrics we decided to use are:

- Pearson correlations: how correlated the performance of each store is to the trial store.
- Magnitude distance: a standardized metric based on the absolute difference between the trial store's performance and each control store's performance. For more information on this metric, visit [Appendix B: Magnitude distance](#).

We selected control stores based on how similar monthly total sales in dollar amounts and monthly number of customers are to the trial stores. So as we applied the two metrics, we would get four scores, two for each of total sales and total customers.

We then combined all the scores calculated by our two metrics to create a composite score to rank on.

We combined by taking a simple average of the correlation and magnitude scores for each driver. Note that if we consider it more important for the trend of the drivers to be similar, we can increase the weight of the correlation score (a simple average gives a weight of 0.5 to each metric) or if we consider the absolute size of the drivers to be more important, we can lower the weight of the correlation score.

Now, we would have a score for each of total number of sales and number of customers. We combined the two via a simple average.

The store with the highest score is then selected as the control store since it is most similar to the trial store.

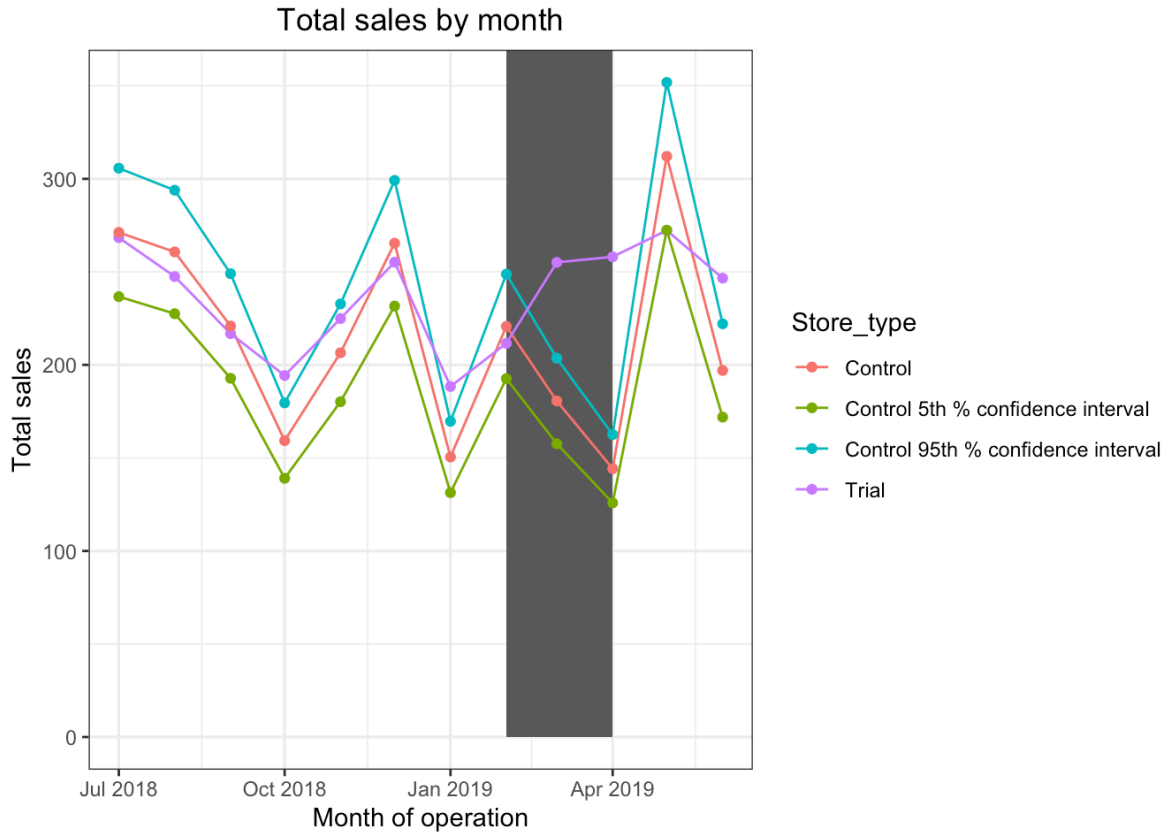
### 3.2.2 Store 77 Assessment

For trial store 77, the control store is 233. We visualized to check if this is reasonable.

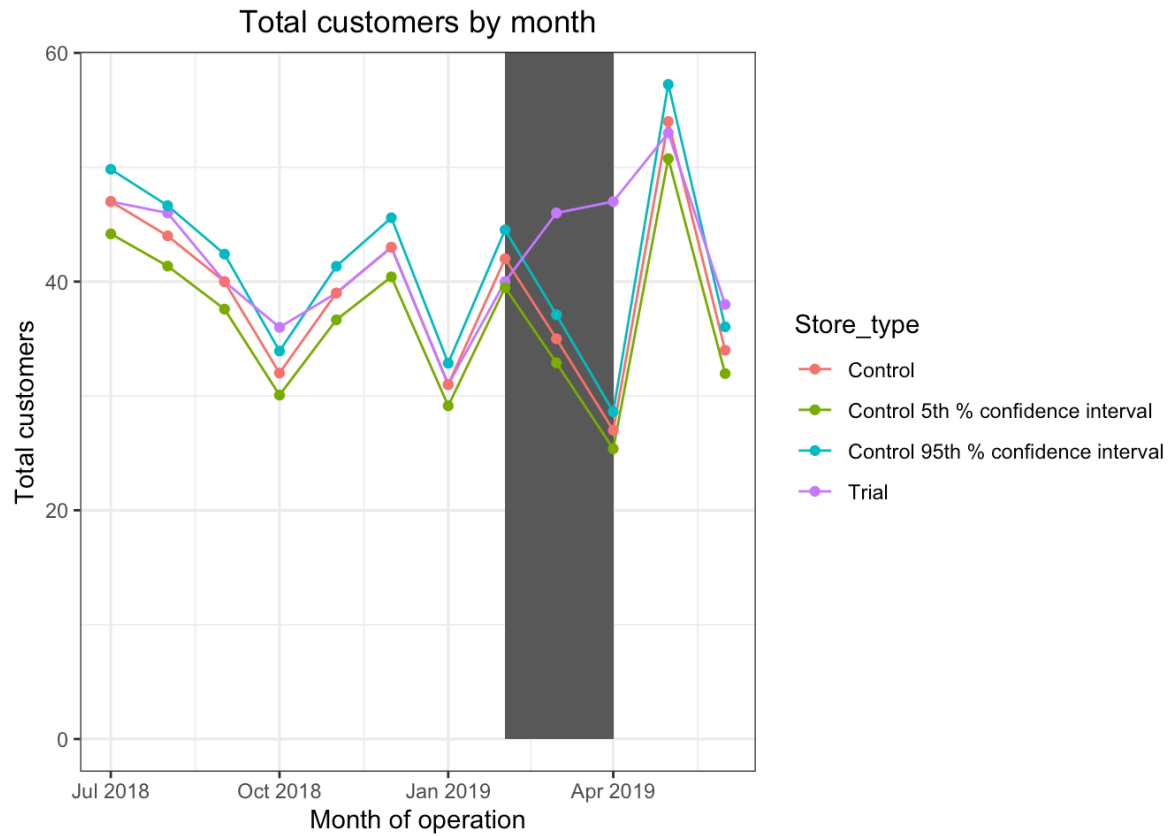




We plotted the sales of the control store, the sales of the trial stores and the 95th percentile value of sales of the control store.

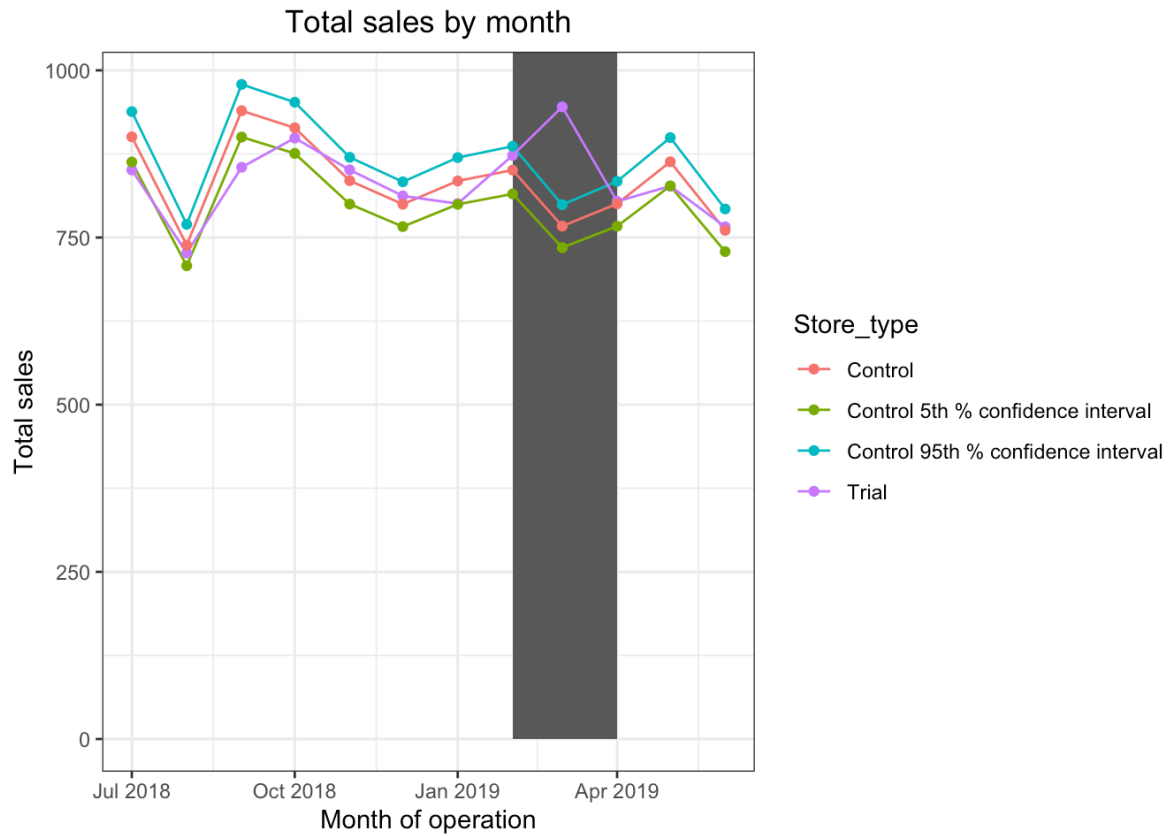


The results show that the trial in store 77 is significantly different to its control store in the trial period as the trial store performance lies outside the 5% to 95% confidence interval of the control store in two of the three trial months.



### 3.2.3 Store 86 Assessment

For trial store 86, the control store is 155.



The results show that the trial in store 86 is not significantly different to its control store in the trial period as the trial store performance lies inside the 5% to 95% confidence interval of the control store in the two of the three trial months.

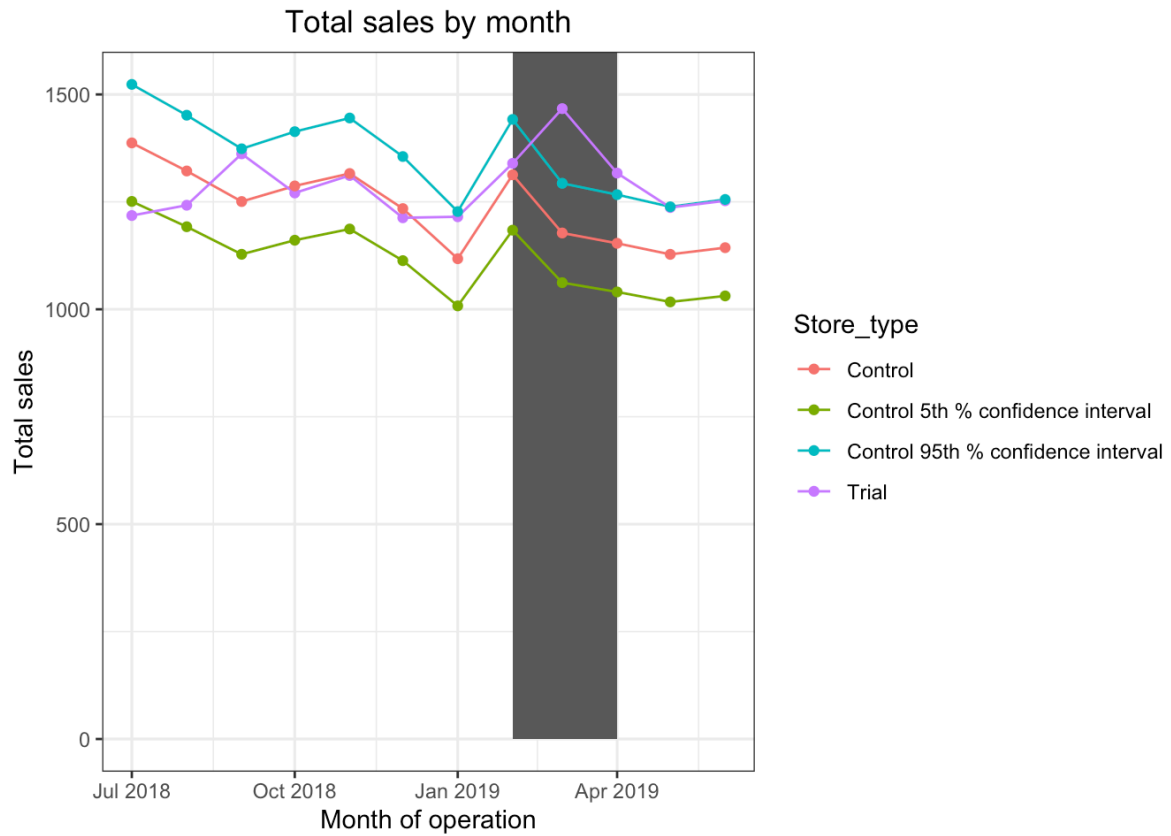


It looks like the number of customers is significantly higher in all of the three months. This seems to suggest that the trial had a significant impact on increasing the number of customers in trial store 86 but as we saw, sales were not significantly higher. We should check with the Category Manager if there were special deals in the trial store that were may have resulted in lower prices, impacting the results.

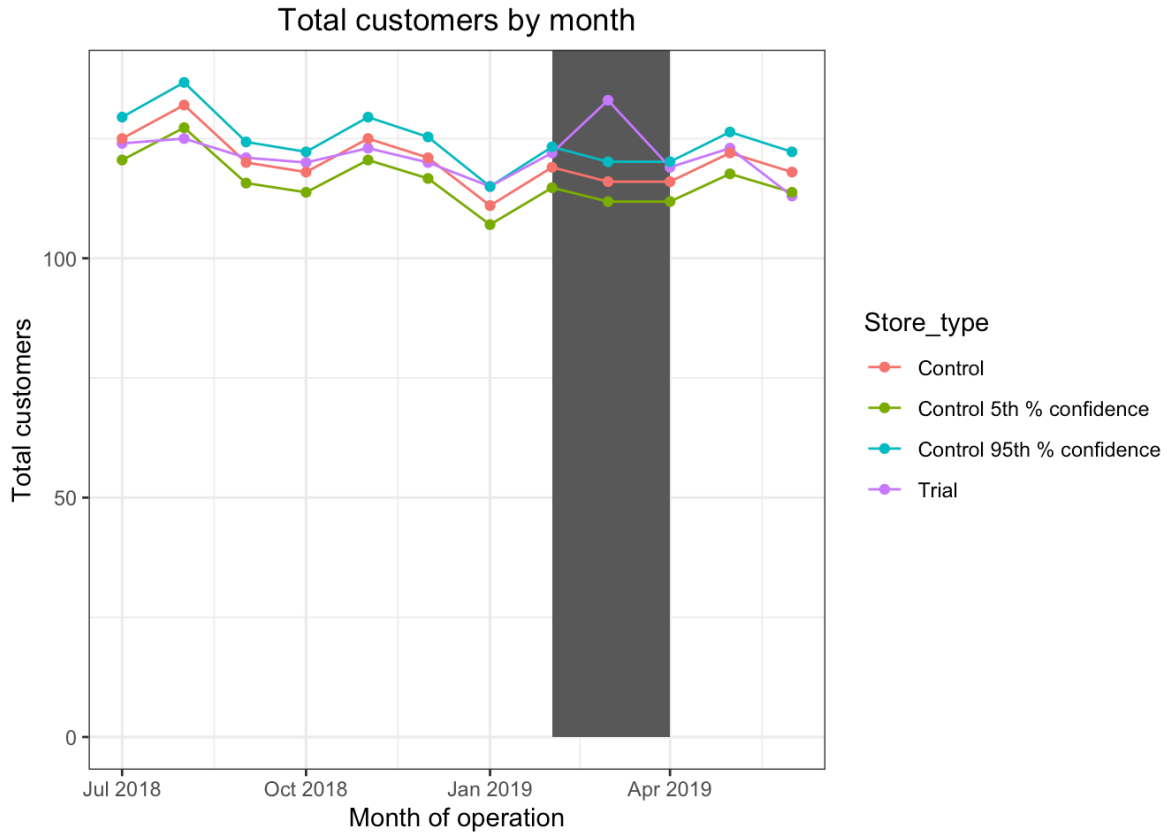
### 3.2.4 Store 88 Assessment

For trial store 88, the control store is 237.





The results show that the trial in store 88 is significantly different to its control store in the trial period as the trial store performance lies outside of the 5% to 95% confidence interval of the control store in two of the three trial months.



Total number of customers in the trial period for the trial store is not significantly higher than the control store for two out of three months, which indicates little to no effect by the trial.

### 3.3 Task 3

We created a [client report in PowerPoint](#) to provide commercial, actionable insights from our analysis and displayed it in a clear and concise way for your client, with minimal jargon. We followed the Pyramid Principles, getting at the actionable findings and what to implement first, then backing them up with explanations later on.

## 4 Conclusions

### 4.1 Task 1

Task 1 analysis serves as a precursor to more formal testing and analysis.

- Sales have mainly been due to Budget - older families, Mainstream - young singles/couples, and Mainstream - retirees shoppers.
- Although Budget - older families accounts for the largest proportion of sales (8.7%), the amount of customers in this segment (6.5%) is not as great as the Mainstream - young singles/couples (11.1%) and Mainstream - retirees shoppers (8.9%)
- This implies the sales can be further improved in the two Mainstream - young singles/couples and Mainstream - retirees shoppers segments
- We also found out the Mainstream - young singles/couples are among the segments that bought the least units of chips per customer. However, they are willing to spend the most on chips.
- Notable mention is Mainstream - young singles/couples, who are also willing to spend more on chips than others.
- This implies that the Mainstream - young singles/couples prefer higher quality chips and/or more popular brands. Their most preferred brands are Tyrrells, Twisties, Doritos, Kettle, and Tostitos
- Mainstream young singles/couples are 27% more likely to purchase pack size of 270g
- 270g seems to be the ideal pack size for Mainstream young singles/couples. However, the only brand that provides this pack size is Twisties in this store. Perhaps, including more brands with 270g pack size is a potential step towards more sales.
- Overall, we suggest an emphasis on the Mainstream - young singles/couples. One suggestion would be to have some Tyrrells chips near areas in the store that this target segment tends to visit the most.

## 4.2 Task 2

We've found control stores 233, 155, 237 for trial stores 77, 86 and 88 respectively.

The results for trial stores 77 during the trial period show a significant difference in at least two of the three trial months for both the total sales and number of customers metrics. Trial store 86 saw a statistically significant increase in customers but not total sales. Trial store 88 saw a statistically significant increase in sales but not customers.

Summary of Significant Difference in our Metrics		
Trial.Stores	Sales	Customers
77	Yes	Yes
86	No	Yes
88	Yes	No

We can check with the client if the implementation of the trial was different in trial stores 86 and 88 but overall, the trial shows a positive effect on our metrics; two out of three stores had higher sales and two out of three stores also had higher number of customers. Now that we have finished our analysis, we can prepare our presentation to the Category Manager.

### 4.3 Task 3

In completion, we communicated statistical evidence to client, confirming the success of the trial layout, and recommended its deployment across all stores.

## 5 References

- [Affinity Analysis](#)
- [Data table in R](#)
- [Measures of Distance in Data Mining](#)
- [The Pyramid Principles](#)
- [Two sample t-test](#)

## 6 Appendix A: Raw Data

## 7 Appendix B: Magnitude distance

i.e.  $1 - (\text{Observed distance} - \text{minimum distance}) / (\text{Maximum distance} - \text{minimum distance})$

- The fraction  $(\text{Observed distance} - \text{Minimum distance}) / (\text{Maximum distance} - \text{Minimum distance})$  represents the relative position of the observed distance within the entire range of possible distances.

	DATE	STORE_NBR	LYLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
1	43390	1	1000	1	5	Natural Chip Compny SeaSalt175g	2	6.00
2	43599	1	1307	348	66	CCs Nacho Cheese 175g	3	6.30
3	43605	1	1343	383	61	Smiths Crinkle Cut Chips Chicken 170g	2	2.90
4	43329	2	2373	974	69	Smiths Chip Thinly S/Cream&Onion 175g	5	15.00
5	43330	2	2426	1038	108	Kettle Tortilla ChpsHny&Jlpno Chili 150g	3	13.80
6	43604	4	4074	2982	57	Old El Paso Salsa Dip Tomato Mild 300g	1	5.10
7	43601	4	4149	3333	16	Smiths Crinkle Chips Salt & Vinegar 330g	1	5.70
8	43601	4	4196	3539	24	Grain Waves Sweet Chilli 210g	1	3.60
9	43332	5	5026	4525	42	Doritos Corn Chip Mexican Jalapeno 150g	1	3.90
10	43330	7	7150	6900	52	Grain Waves Sour Cream&Chives 210G	2	7.20
11	43602	7	7215	7176	16	Smiths Crinkle Chips Salt & Vinegar 330g	1	5.70
12	43332	8	8294	8221	114	Kettle Sensations Siracha Lime 150g	5	23.00
13	43603	9	9208	8634	15	Twisties Cheese 270g	2	9.20
14	43329	13	13213	12447	92	WW Crinkle Cut Chicken 175g	1	1.70
15	43600	19	19272	16686	44	Thins Chips Light& Tangy 175g	1	3.30
16	43604	20	20164	17136	54	CCs Original 175g	1	2.10
17	43330	20	20418	17413	94	Burger Rings 220g	4	9.20
18	43326	22	22411	18646	98	NCC Sour Cream & Garden Chives 175g	1	3.00
19	43329	22	22456	18696	93	Doritos Corn Chip Southern Chicken 150g	1	3.90
20	43601	23	23067	19162	56	Cheezels Cheese Box 125g	1	2.10
21	43604	25	25105	21815	7	Smiths Crinkle Original 330g	1	5.70
22	43328	33	33081	29949	98	NCC Sour Cream & Garden Chives 175g	1	3.00
23	43328	36	36012	32077	31	Infzns Crn Crnchers Tangy Gcamole 110g	1	3.80
24	43331	36	36302	33188	32	Kettle Sea Salt And Vinegar 175g	1	5.40
25	43327	38	38142	34181	108	Kettle Tortilla ChpsHny&Jlpno Chili 150g	2	9.20
26	43600	39	39144	35506	57	Old El Paso Salsa Dip Tomato Mild 300g	1	5.10
27	43331	39	39167	35638	111	Smiths Chip Thinly Cut Original 175g	2	6.00
28	43600	41	41423	38393	46	Kettle Original 175g	1	5.40

Figure 2: Transaction data

	LYLTY_CARD_NBR	LIFESTAGE	PREMIUM_CUSTOMER
1	1000	YOUNG SINGLES/COUPLES	Premium
2	1002	YOUNG SINGLES/COUPLES	Mainstream
3	1003	YOUNG FAMILIES	Budget
4	1004	OLDER SINGLES/COUPLES	Mainstream
5	1005	MIDAGE SINGLES/COUPLES	Mainstream
6	1007	YOUNG SINGLES/COUPLES	Budget
7	1009	NEW FAMILIES	Premium
8	1010	YOUNG SINGLES/COUPLES	Mainstream
9	1011	OLDER SINGLES/COUPLES	Mainstream
10	1012	OLDER FAMILIES	Mainstream
11	1013	RETIREEES	Budget
12	1016	OLDER FAMILIES	Mainstream
13	1018	YOUNG SINGLES/COUPLES	Mainstream
14	1019	OLDER SINGLES/COUPLES	Premium
15	1020	YOUNG SINGLES/COUPLES	Mainstream
16	1022	OLDER FAMILIES	Budget
17	1023	MIDAGE SINGLES/COUPLES	Premium
18	1024	YOUNG SINGLES/COUPLES	Premium
19	1025	YOUNG FAMILIES	Budget
20	1026	MIDAGE SINGLES/COUPLES	Premium
21	1027	OLDER FAMILIES	Premium
22	1028	YOUNG SINGLES/COUPLES	Budget
23	1030	RETIREEES	Mainstream
24	1034	RETIREEES	Premium
25	1038	OLDER FAMILIES	Mainstream
26	1039	YOUNG FAMILIES	Mainstream
27	1042	YOUNG SINGLES/COUPLES	Premium
28	1043	YOUNG FAMILIES	Budget

Figure 3: Customer data

- By subtracting the fraction  $(\text{Observed distance} - \text{Minimum distance}) / (\text{Maximum distance} - \text{Minimum distance})$  from 1, we obtain a value that reflects the similarity between the trial store and the control store.
  - With  $(\text{Observed distance} - \text{Minimum distance}) / (\text{Maximum distance} - \text{Minimum distance})$ , the lowest measure means highest similarity
- A value of 1 indicates maximum similarity (when the observed distance equals the minimum distance), while a value of 0 indicates minimum similarity (when the observed distance equals the maximum distance).