



# Understanding Customer Reviews

## Team Cape Town

Anan



Teresa



Shuyue



Andy



Fannie



Cole



Akif





- **Overview**
- **Data Driven Insights**
- **Modelling Journey**
- **Model Evaluation**
- **LLM Integration**
- **Key Takeaways**



## Overview



### Challenge

Lots of reviews per product, making it difficult for customers to find useful, relevant information, to assist with their purchase decisions

### Current Situation



#### Poorly Written Reviews

Language & grammar errors make feedback difficult to decipher



#### Biased Reviews

Incentivized by discounts or brand affiliation



#### Unverified Reviews

No authentication for whether reviewer purchased the product



#### Duplicated Reviews

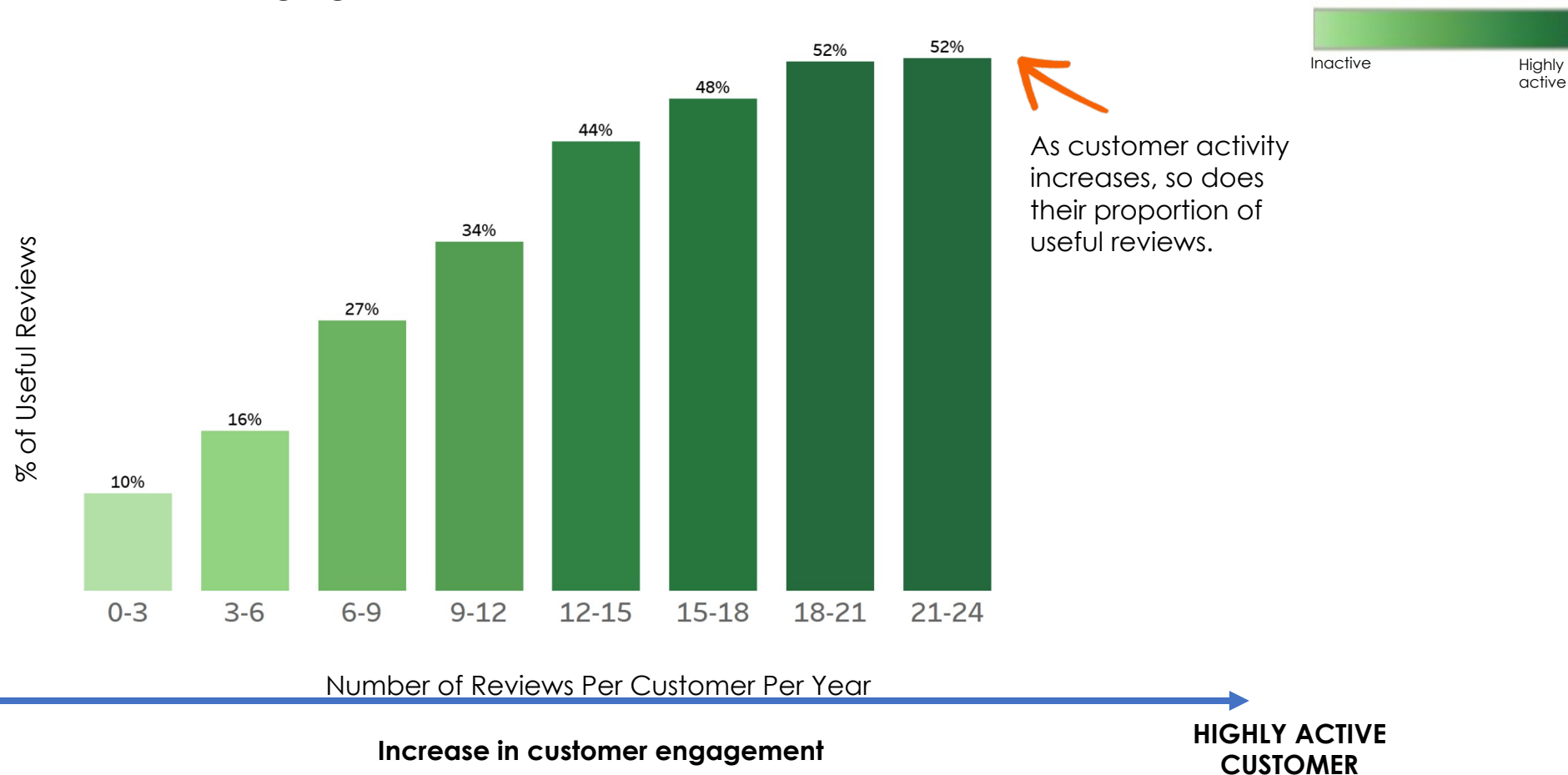
Reposting positive or negative comments

### Goal

Identify reviews which are helpful to **customers** in their purchase decisions and to **Amazon** to improve their product assortment

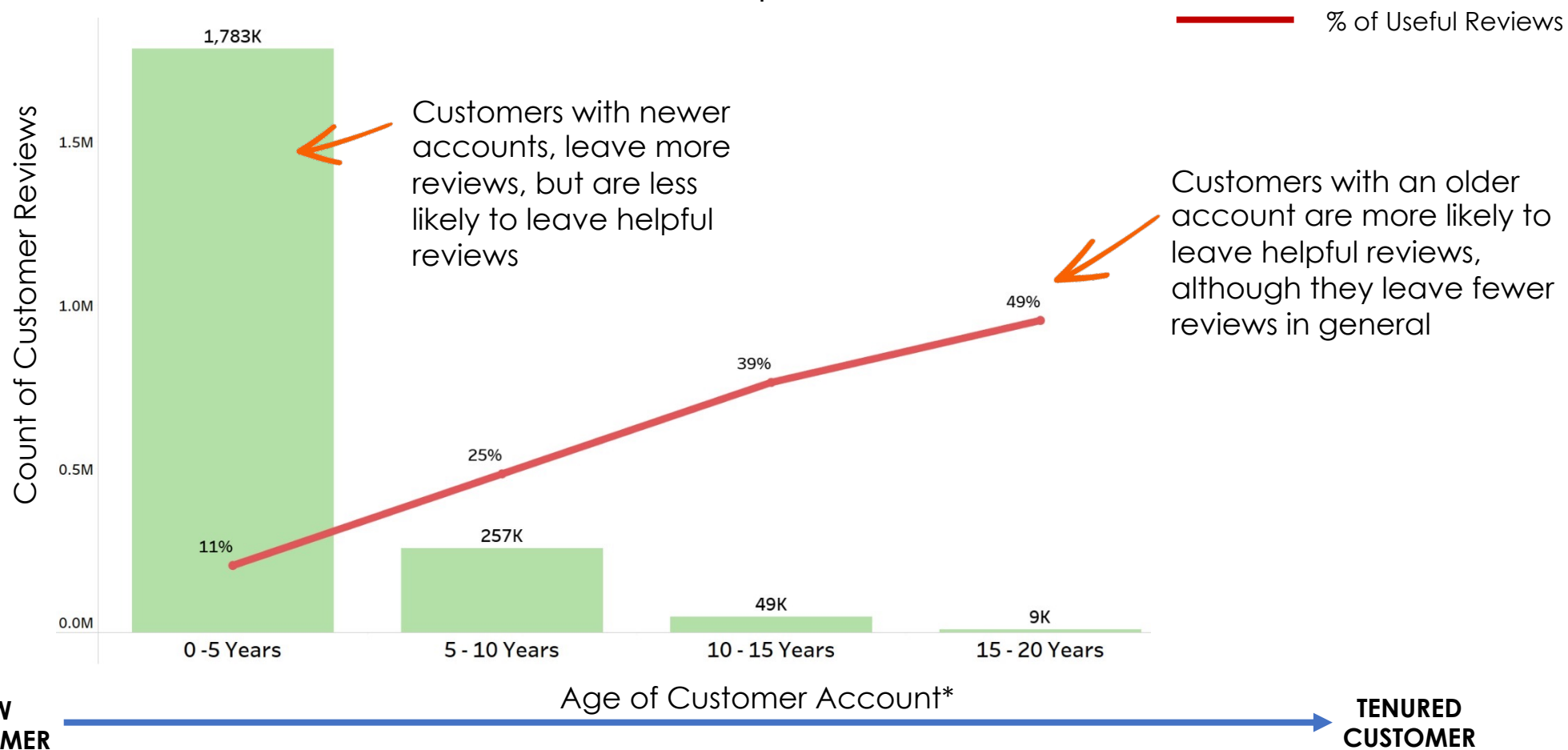


## Customer Engagement: More engaged customers leave useful reviews





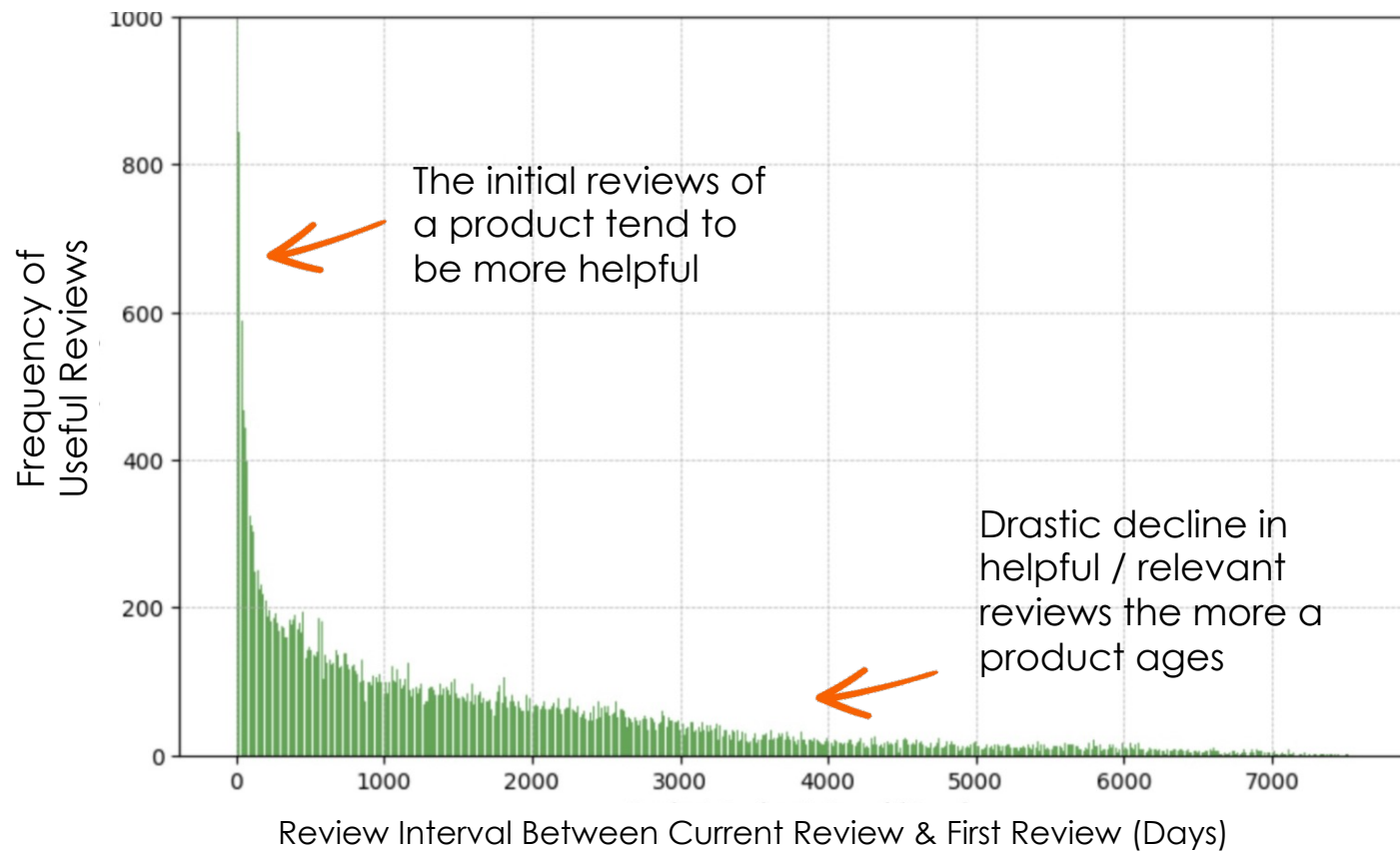
## Age of Customer Account\*: Older accounts leave helpful reviews



\* This is a proxy for the age of a customer account, calculated as the time since the first review of a customer.



**Reviews Per Product:** The initial reviews of a product are usually the helpful ones





### Customer Sentiment:





**Objective:** To flag customer reviews as helpful or not



**Approach:** Assess the recency, frequency & content of labelled, historically helpful reviews to build a predictive model, to spot similar patterns in recently generated reviews

## CRISP-DM Modelling Process



Data Cleaning  
& Pre-Processing



Pipeline



Helpful



Not Helpful



01

02

03

04

### Data Generation

Customers reviewing products on the Amazon website

### Feature Engineering

Features created at the customer, product & review levels

### Classification

The winning model for the Binary Classifier is **Light GBM**.

### Model Tuning & Evaluation

Grid Search to identify best hyper parameters, optimizing for AUC





## Feature Engineering

### CUSTOMER LEVEL

- Total reviews per customer
- Unique products reviewed
- Recency & Frequency of reviewing



### TEXT PREPROCESSING

- Tokenizing words
- Removing stop words
- Count Vector (Word Frequency)
- TF-IDF of frequently occurring words



### PRODUCT LEVEL

- Reviews generated per product
- General sentiment of product
- Recency & Frequency of reviews



### CONTENT OF REVIEW

- Word count of each review
- Sentiment of review
- Count of question marks, caps lock & exclamation points in review



### SEASONALITY

- Day of week / month of review creation
- Black Friday, Boxing Day & Prime Week



### DIVERGENCE METRICS

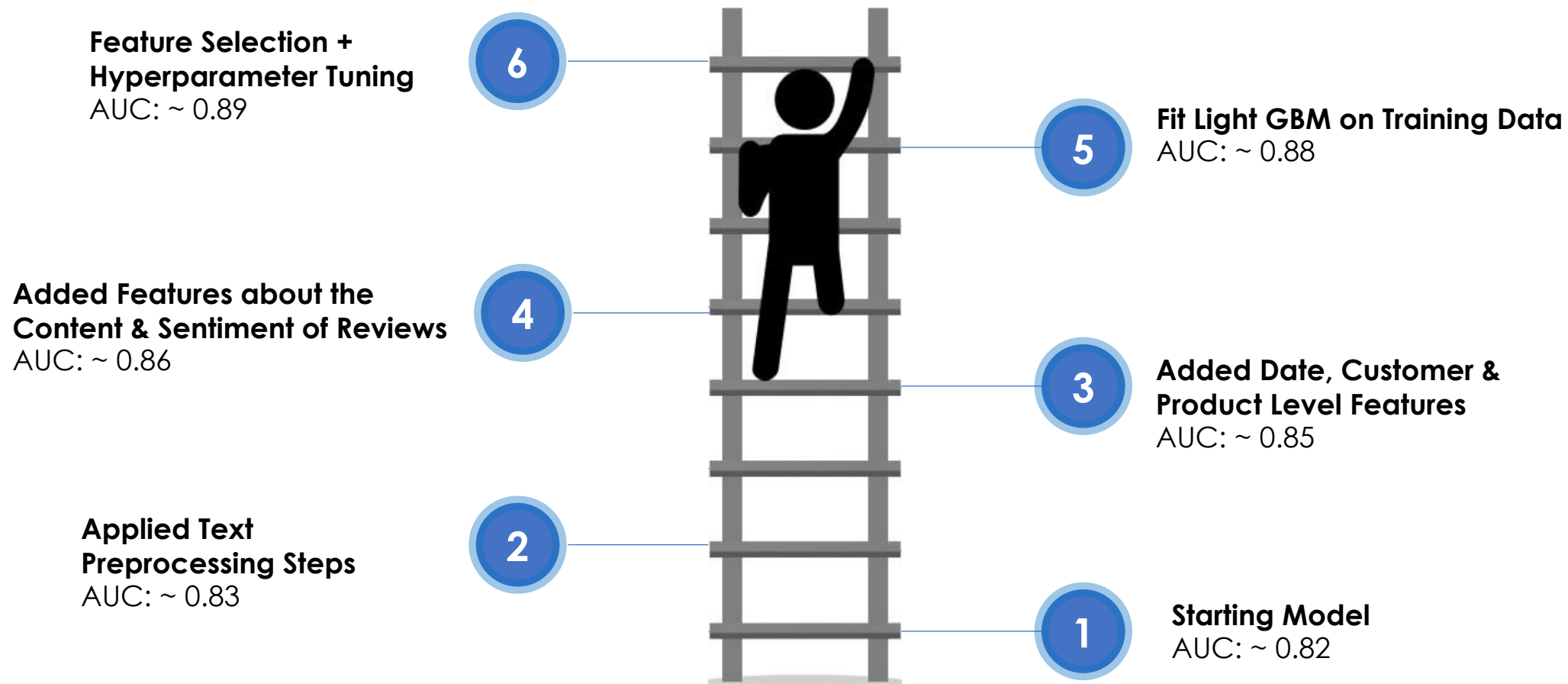
- Does the review go against majority opinion?
- Does the customer only review when they have a negative experience?



**34 Features that  
Help Identify  
What Makes a  
Review Useful**

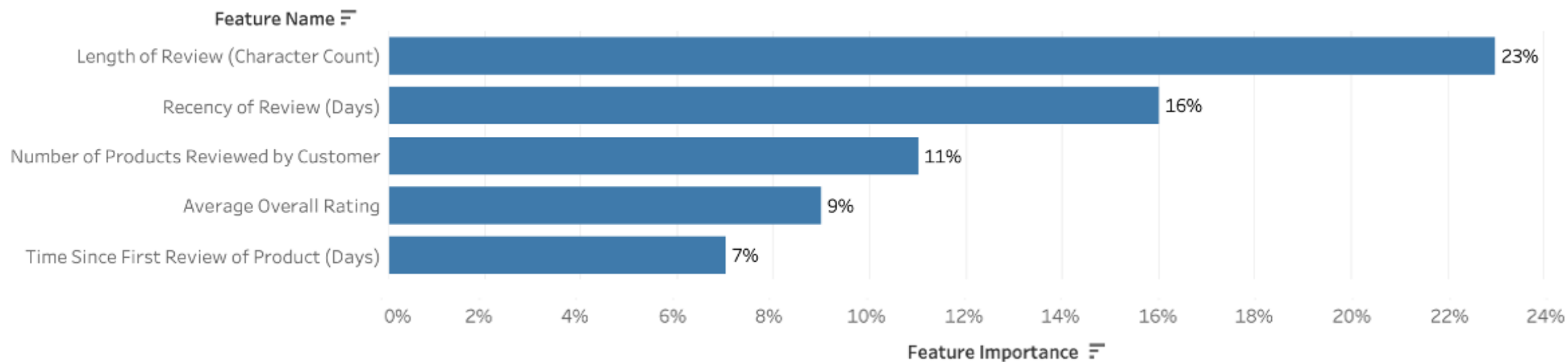


## Model Improvement





## Feature Importance Chart of Winning Model



Only top 5 features shown



## LLM Integration

- ✓ **Debugging Code:** Pasting the code and the error from the code provided a debugged version of the code, especially correcting for syntax issues.

### Error Message:

⊕ **AnalysisException:** [AMBIGUOUS\_REFERENCE] Reference `avg\_sentiment\_score\_per\_product` is ambiguous, could be: [`avg\_sentiment\_score\_per\_product`, `avg\_sentiment\_score\_per\_product`].

### LLM Solution:

M can you debug the following code for me?

[AMBIGUOUS\_REFERENCE] Reference `avg\_sentiment\_score\_per\_product` is ambiguous, could be: [`avg\_sentiment\_score\_per\_product`, `avg\_sentiment\_score\_per\_product`].

The error message "[AMBIGUOUS\_REFERENCE] Reference 'avg\_sentiment\_score\_per\_product' is ambiguous" typically indicates that there are two columns in the dataframe with the same name, which is causing confusion during operations.

This situation can arise in a couple of scenarios:

1. **Joining Two DataFrames:** If you join two DataFrames and both have a column with the same name, then the resulting DataFrame will have ambiguous references for that column.
2. **Adding a Column Twice:** You might have accidentally added a column with the same name to the DataFrame more than once.

To fix the error, you'll need to resolve the ambiguous reference. Here's a potential solution based on a general scenario:

Regenerate

Explanation  
of error  
message

#### Scenario: Joining Two DataFrames

Let's assume the ambiguity arose because of joining two DataFrames `df1` and `df2`, both of which have the column `avg\_sentiment\_score\_per\_product`.

When joining the DataFrames, you can rename the columns in one of the DataFrames to avoid ambiguity:

```
python Copy code

# Renaming the column in one DataFrame before joining
df2 = df2.withColumnRenamed("avg_sentiment_score_per_product", "avg_sentiment_score_per_product_df2")

# Now perform the join
joined_df = df1.join(df2, on="common_column", how="inner")
```

Now, the joined DataFrame `joined\_df` will have two distinct columns: `avg\_sentiment\_score\_per\_product` (from `df1`) and `avg\_sentiment\_score\_per\_product\_df2` (from `df2`).

Code  
sample

Now, the joined DataFrame `joined\_df` will have two distinct columns: `avg\_sentiment\_score\_per\_product` (from `df1`) and `avg\_sentiment\_score\_per\_product\_df2` (from `df2`).

Explanation  
of code  
snippet /  
solution

Remember to adjust the column operations accordingly after renaming. If you need to access the value from the original column in `df1`, use `avg\_sentiment\_score\_per\_product`. If you need to access the value from the renamed column in `df2`, use `avg\_sentiment\_score\_per\_product\_df2`.

In your given code, you might need to look at where the `df` dataframe is joined with `sentiment\_product` and `result` dataframes. Ensure that none of the joining DataFrames have the column `avg\_sentiment\_score\_per\_product` before joining or, if they do, resolve the ambiguity as shown above.



## LLM Integration

- ✓ **Code Refactoring:** ChatGPT made the code better formatted, legible and easier to edit, increasing efficiency, as the code becomes more “readable”.

### Pre LLM: 24 lines

```
1 def review_counts_per_product(data_input):
2     |
3     review_counts = data_input.groupBy("asin").agg(
4         countDistinct("reviewID").alias("reviews_per_product"),
5         countDistinct('reviewerID').alias("reviewers_per_product"),
6         min('reviewDate').alias('product_earliest_review'),
7         max('reviewDate').alias('product_latest_review'))
8
9     data_input = data_input.join(review_counts, on="asin", how="inner")
10
11    data_input = data_input.withColumn('product_review_interval',
12                                     datediff(data_input['product_latest_review'],
13                                               data_input['product_earliest_review']))
14
15    data_input = data_input.withColumn("product_review_interval",
16                                     when(col("product_latest_review") == col
17                                         ("product_earliest_review"), 0).otherwise(col
18                                         ("product_review_interval")))
19
20    data_input = data_input.withColumn('product_earliest_review',
21                                     datediff(current_date(), data_input
22                                               ['product_earliest_review']))
23
24    data_input = data_input.withColumn('product_latest_review',
25                                     datediff(current_date(), data_input
26                                               ['product_latest_review']))
27
28    return data_input
```

### Post LLM: 17 lines (~29% reduction)

```
1 def review_counts_per_product(data_input):
2
3     review_counts = (data_input.groupBy("asin")
4                     .agg(F.countDistinct("reviewID").alias("reviews_per_product"),
5                          F.countDistinct('reviewerID').alias("reviewers_per_product"),
6                          F.min('reviewDate').alias('product_earliest_review'),
7                          F.max('reviewDate').alias('product_latest_review')))
8
9     review_counts = (review_counts.withColumn('product_review_interval',
10                                              F.when(F.col("product_latest_review") == F.col
11                                                  ("product_earliest_review"), 0)
12                                                  .otherwise(F.datediff(F.col('product_latest_review'),
13                                                                    F.col('product_earliest_review'))))
14                    .withColumn('product_earliest_review_age', F.datediff(F
15                                current_date(), F.col('product_earliest_review')))
16                    .withColumn('product_latest_review_age', F.datediff(F.current_date
17                                (), F.col('product_latest_review'))))
18
19    data_input = data_input.join(review_counts, on="asin", how="inner")
20
21    return data_input
```

Compact code

Multiple new columns created all at once



## Learnings

- ✓ **Feature Engineering:** The biggest improvement in model performance was brought about by adding features on how reviews are perceived by customers.
- ✓ **Feature Selection:** Relevant features need to be included in the model. Too many / too little features can erode performance.
- ✓ **Fitting Different Models:** Switching from Logistic Regression to Boosting models (Light GBM) improved model performance.
- ✓ **Cross Validation:** Cross-validation to optimize for hyper parameter values (using Optuna, Grid Search) improved model performance.
- ✓ **Sampling Data:** Data must be sampled to complete model training within a reasonable amount of time.



**Thanks!**





 Our team: **Team Cape Town**

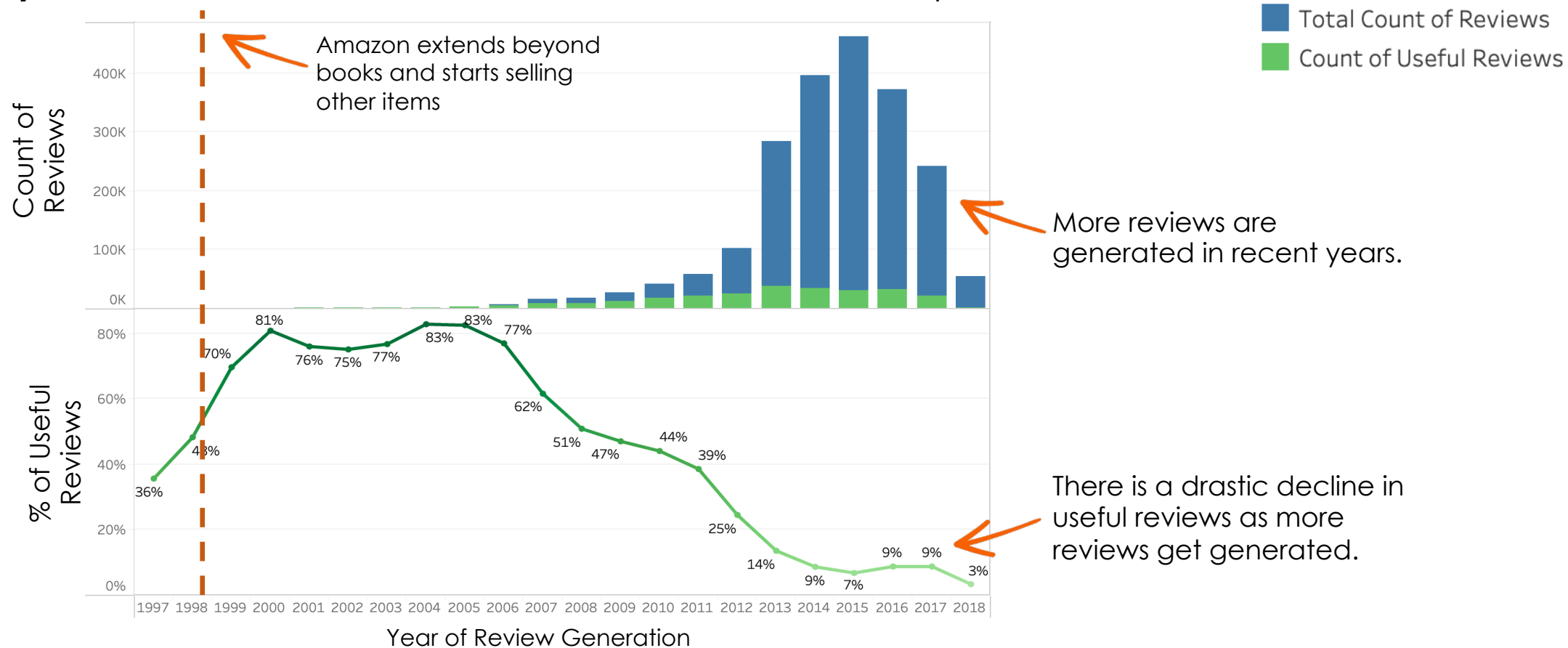


# APPENDIX



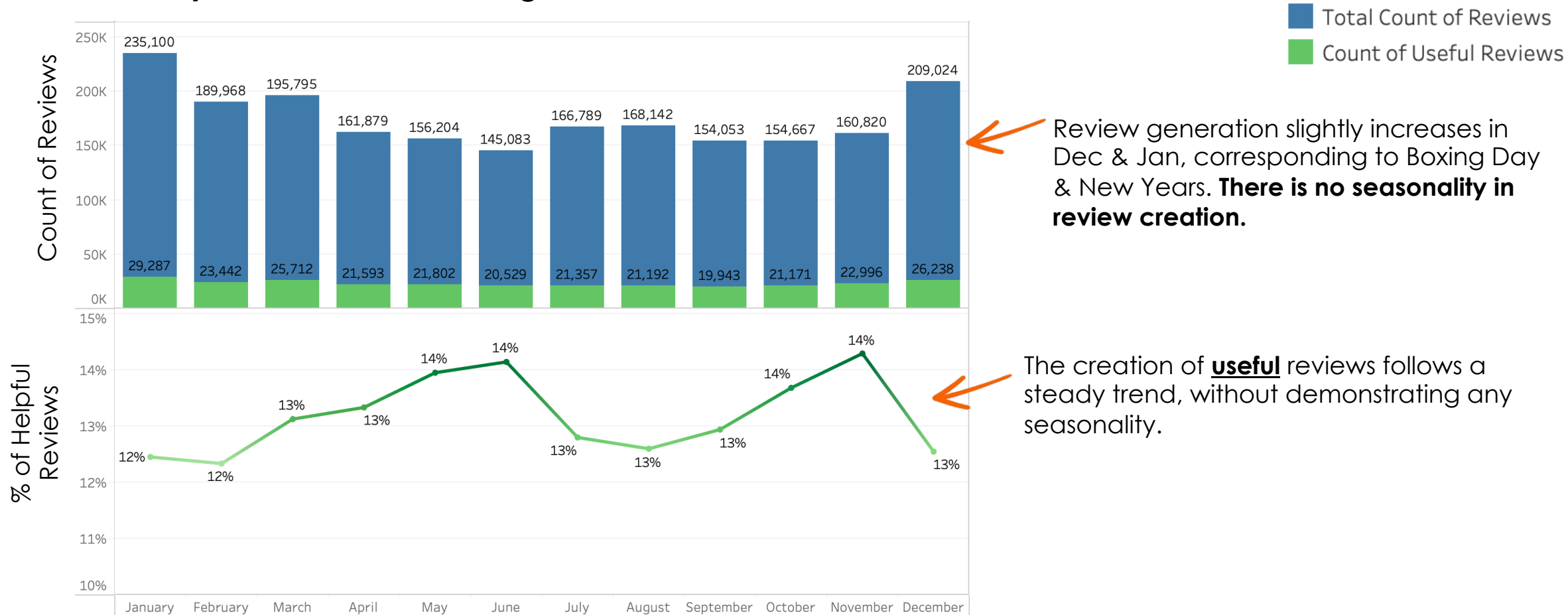


## Recency of Reviews: A recent review does not mean it is a helpful one!



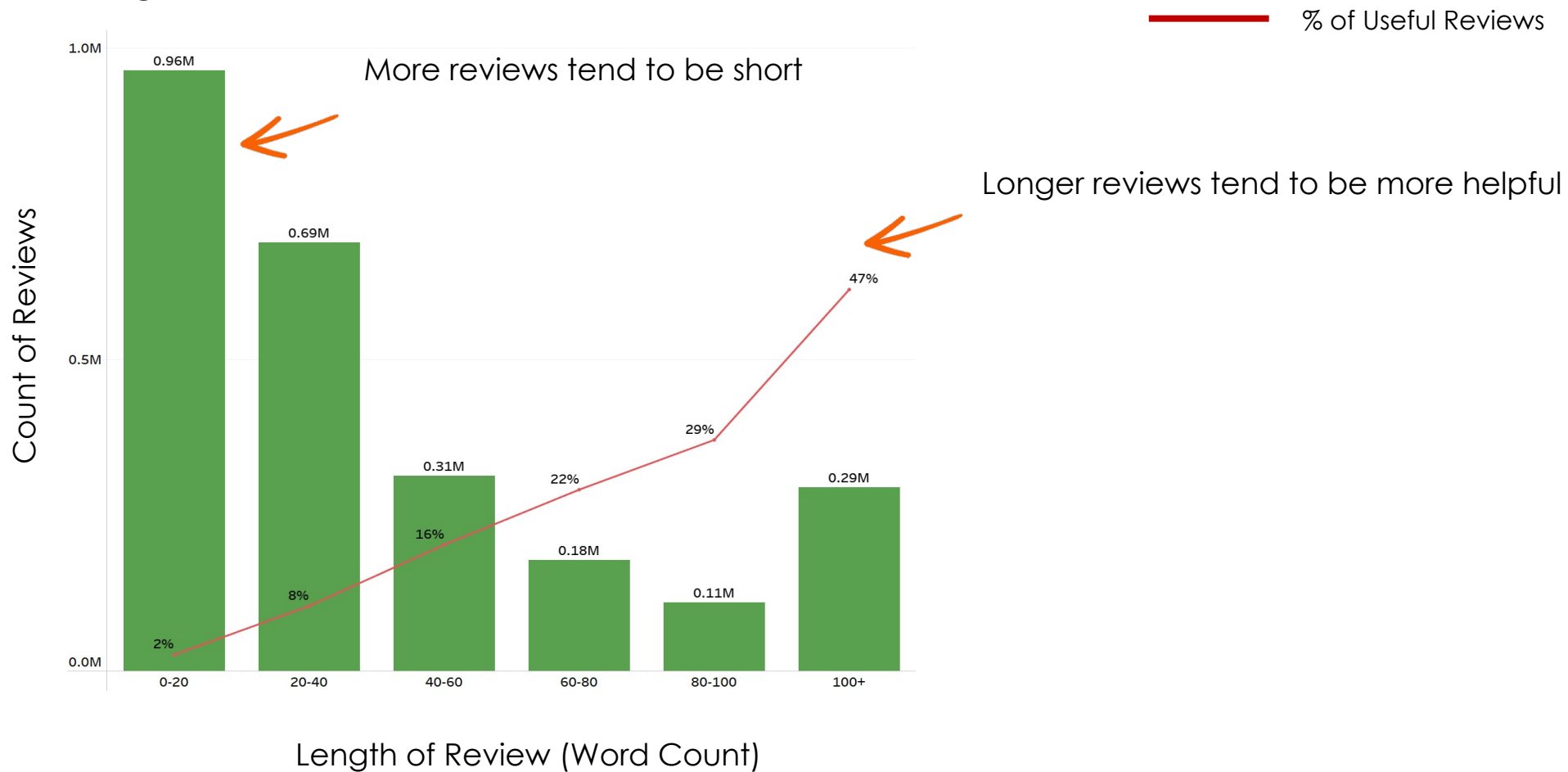


## Seasonality of Reviews: Review generation does not follow a seasonal trend





## Length of Reviews: Longer reviews tend to be more helpful





## Feature Importance Chart of Winning Model

