# COUNTERFEIT DETECTION ON E-COMMERCE PLATFORMS
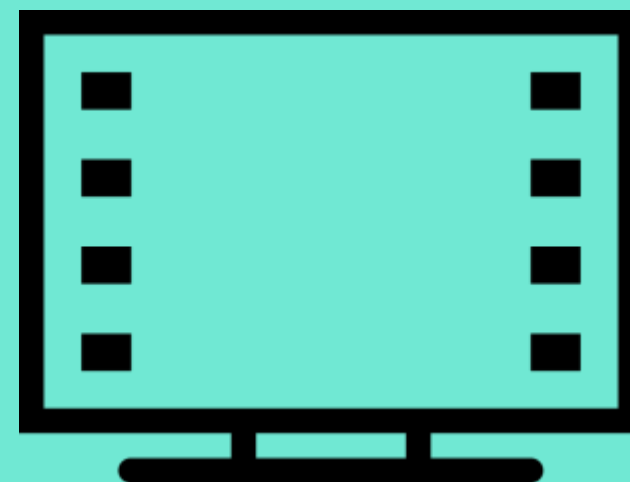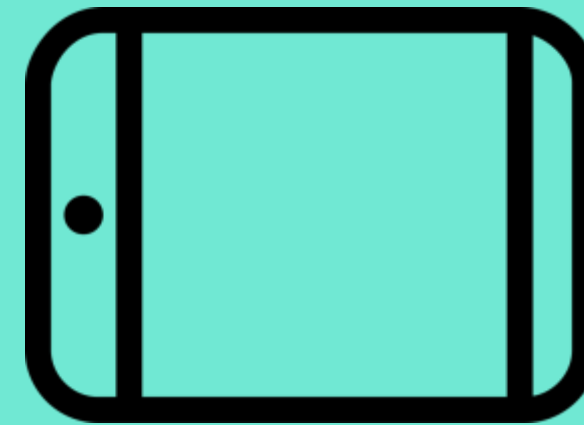
Data Science Bootcamp
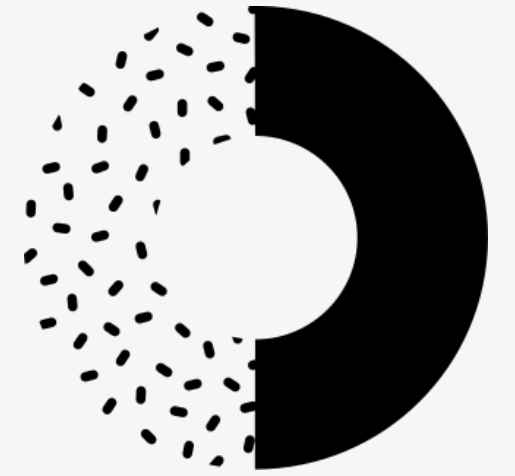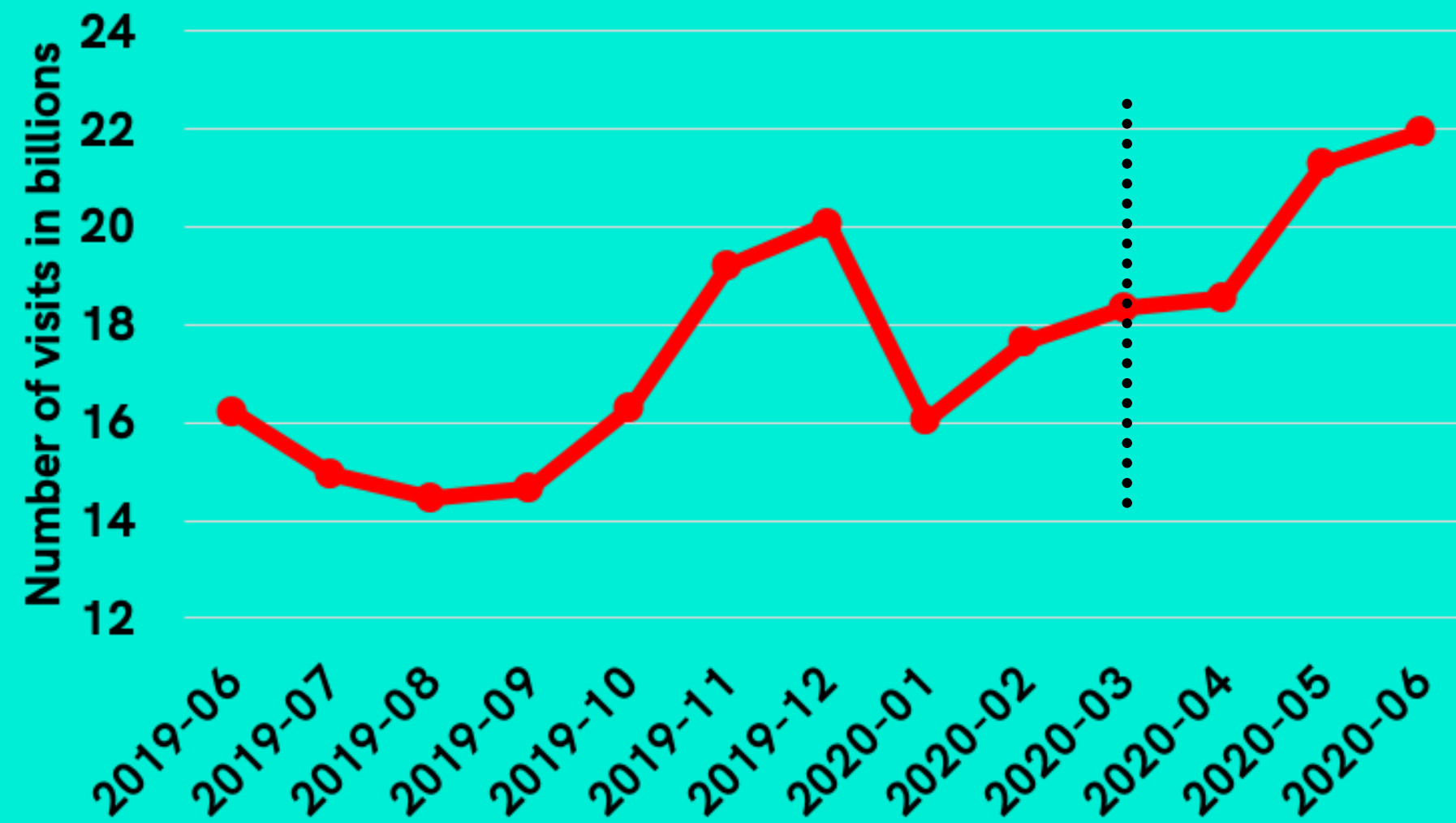
Analytiks Inc.

October 3, 2020

ONLINE
SHOPPING

# E-commerce

*n.* buying and selling of goods or services using the internet

## Coronavirus impact on retail e-commerce website traffic worldwide

Number of visits in billions

## New Normal

By 2040, an estimated 95% of purchases will be made online.

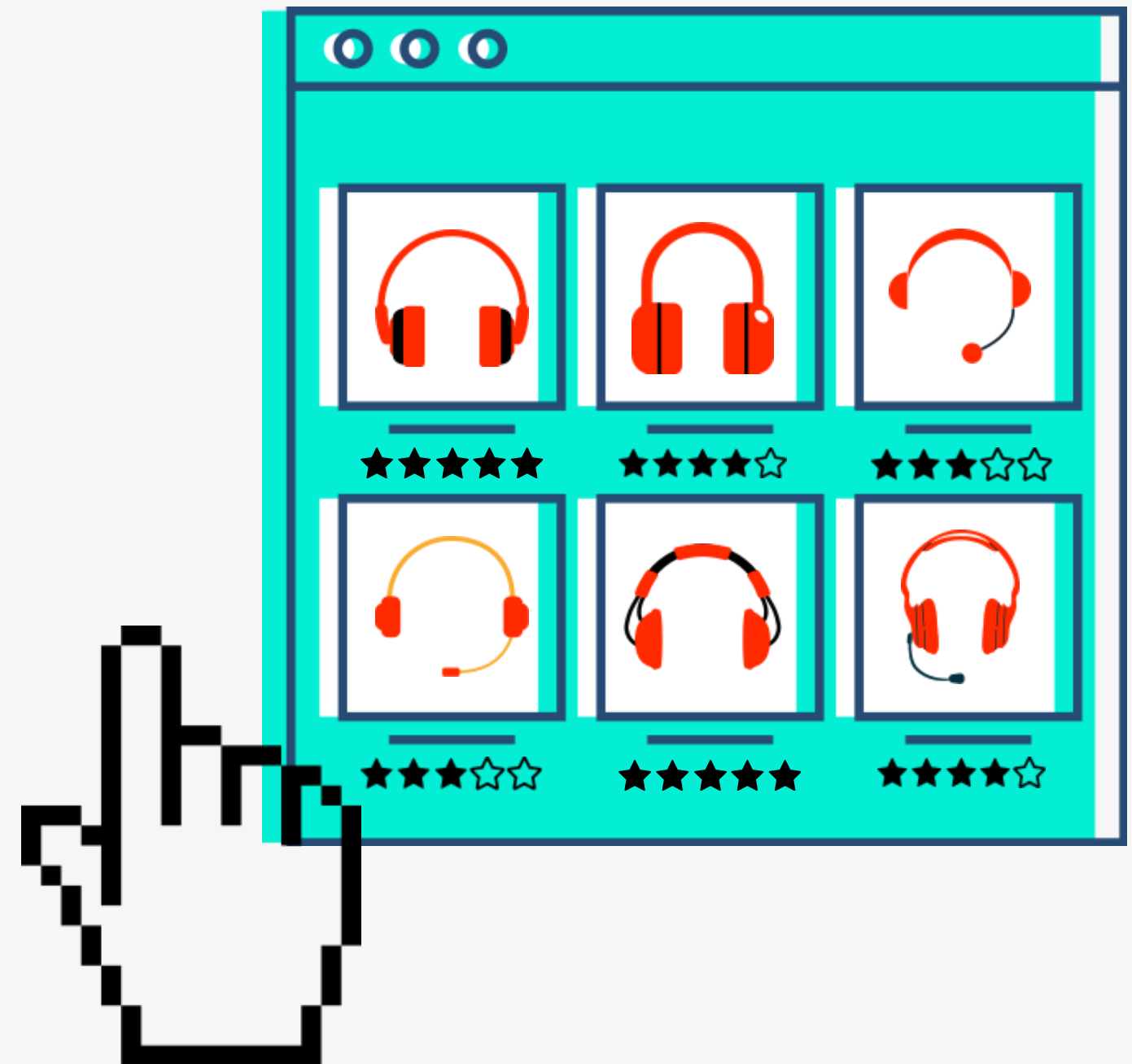# Advantages of Online Shopping

## Convenience
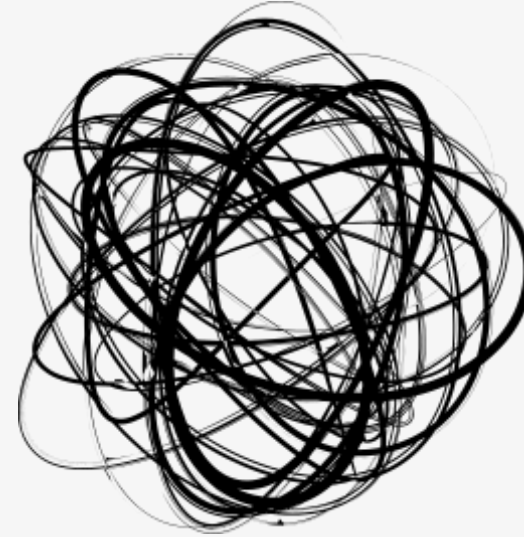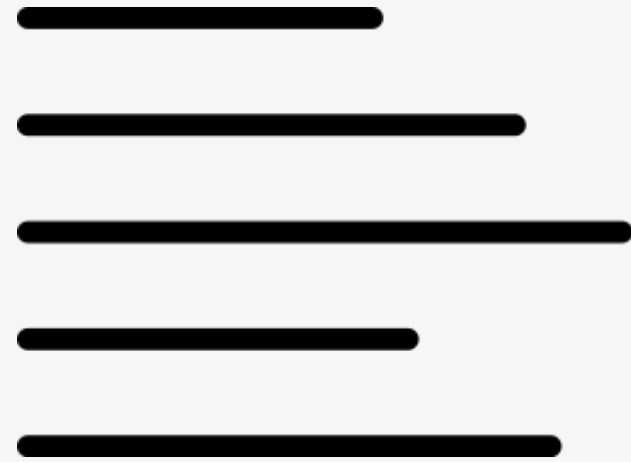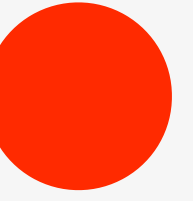No need to go to stores

## Reviews
Testimonials help in decision-making

## Diversity
Providers from all over the world

# COUNTERFEITING

- fraudulent imitation of a product

Amount of total counterfeiting globally is bound to reach

$ 1.82 Trillion

by the year 2020.

## Key Concerns of Counterfeiting:

- Introduce dangerous products into the market
- Weaken environmental, health, and safety regulations
- Diminish tax revenues
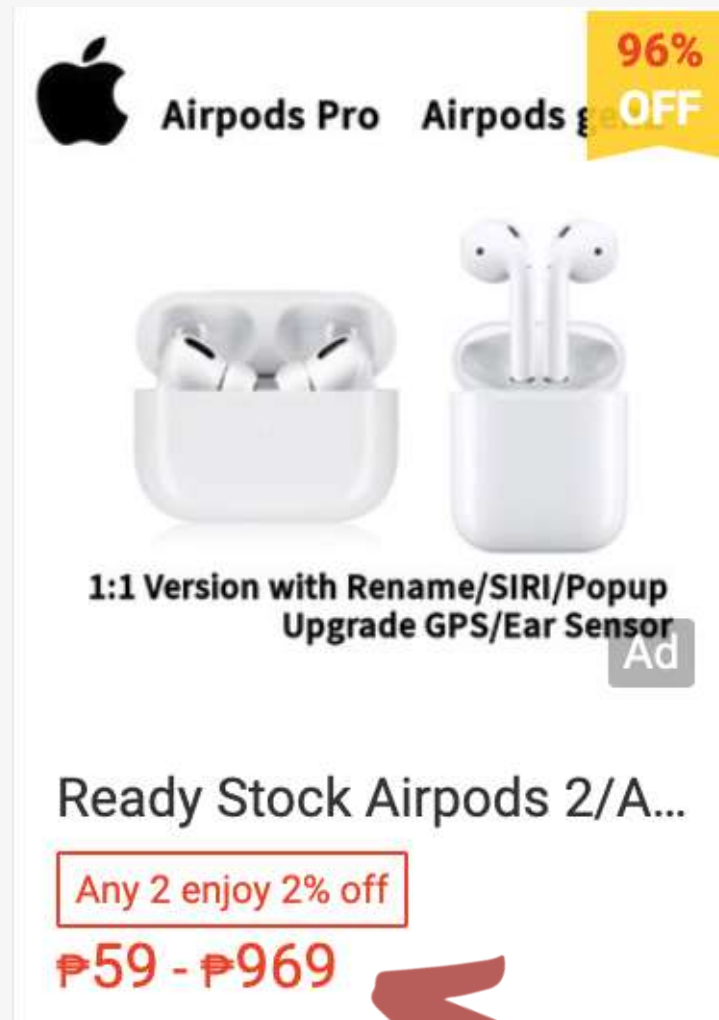- Support organized crime
- Promote child labor

# PROBLEM

How can we detect counterfeit products on e-commerce platforms?
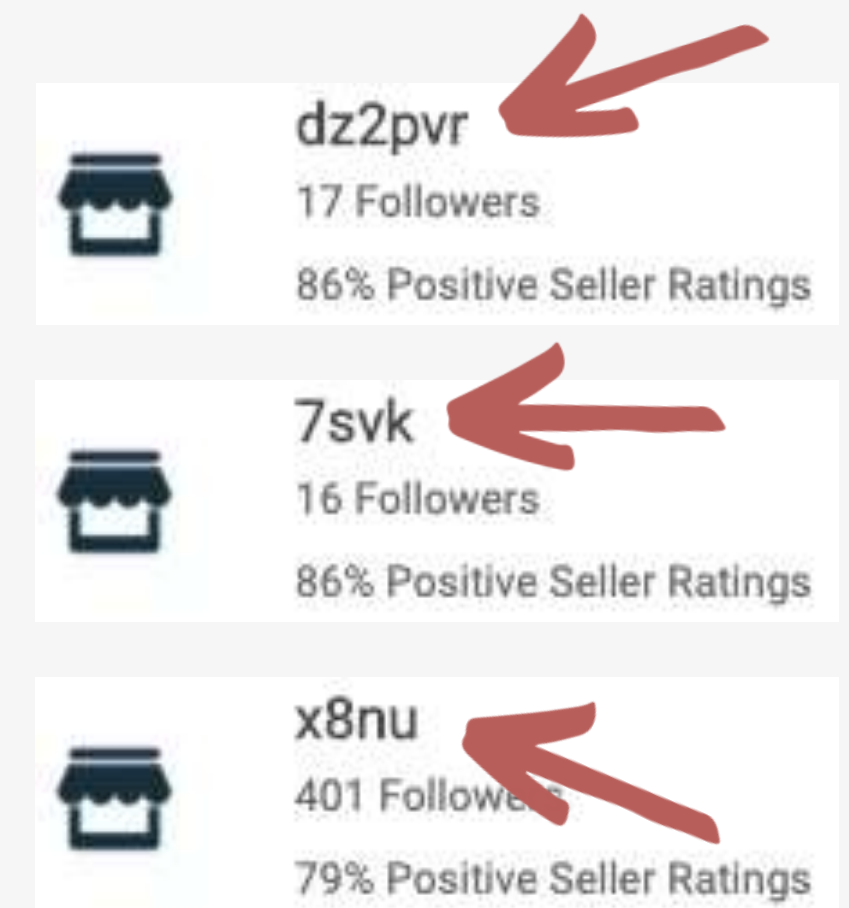
# Traditional ways to spot counterfeits
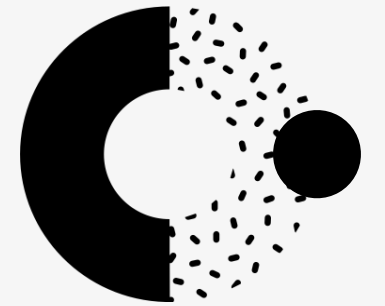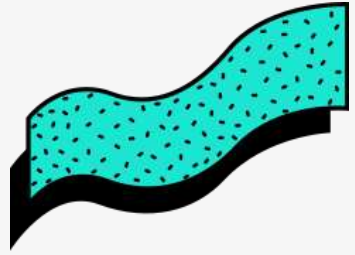
## Price
### Way below SRP



## Packaging
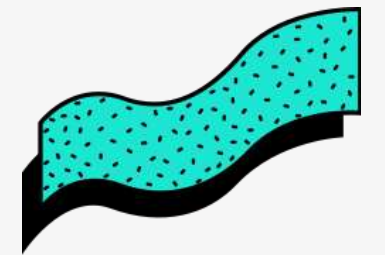### Suspicious branding



## Place
### Sketchy retailers

# Use machine learning!

Price, Packaging, and Place are <u>features</u> that can be used to classify products using a predictive model.

# OBJECTIVES

1. Use machine learning to **predict** whether or not an e-commerce listing is authentic.
2. **Determine** how the features affect the decision-making of the model.

# PROJECT FLOW

Step 1
Assessment

Step 2
Collect data

Step 3
Prepare data

Step 4
Build model

Step 5
Evaluate model

# Evaluate options

## 1. Platforms

Lazada

Shopee

- Top 2 e-commerce sites based on traffic
- Diverse selection of products
- Massive user base

# Evaluate options

## 2. Products

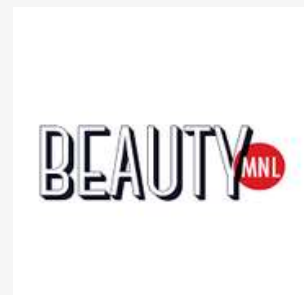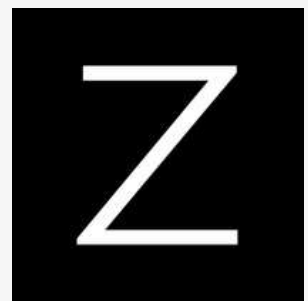- In the Philippines, the e-market's largest segment is Electronics & Media.
- In 2016, electronics accounted for 35% in the global trade of fake goods.

airpods pro

apple watch

ps4 controller

# How to collect?

## Dynamic web scraping

Use Selenium

- Automated
- Simulate "human-like" behavior

Page-by-page collection

- Lazada: 40 listings/page
- Shopee: 50 listings/page

search results     list of product URLs     product page     database

**get individual links**     **open in browser**     **scrape the data**

# How much data?

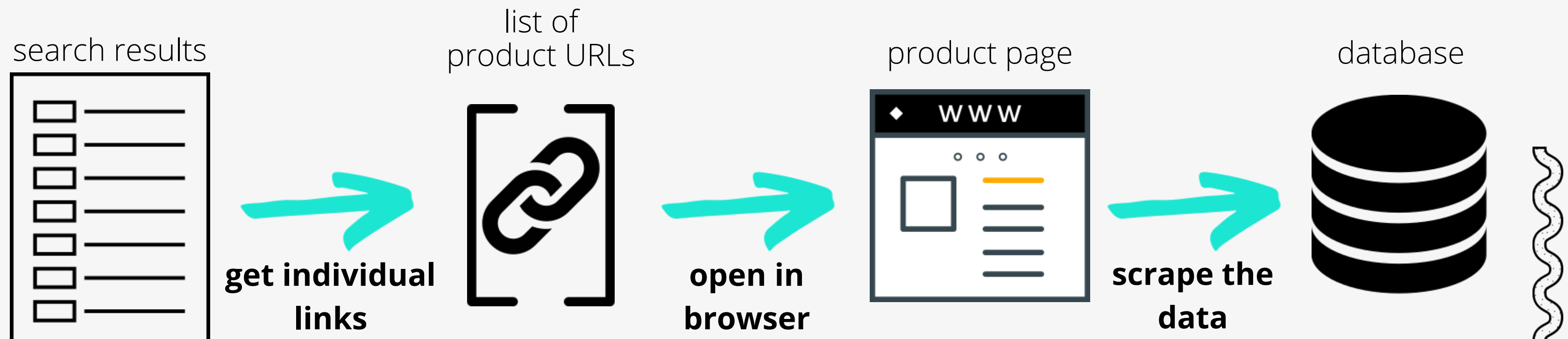## All of it, of course!

### Challenge #1

- Unrealistic! Customer cannot possibly view all the listings.

104243 items found for "gopro"

1/100   <   >

### Challenge #2

- Need to limit request rate.

Sorry, we have detected unusual traffic from your network.
Please slide to continue.

>>    Please slide to verify

6e60edfe5d95276a6fdb8a2ee8a541de
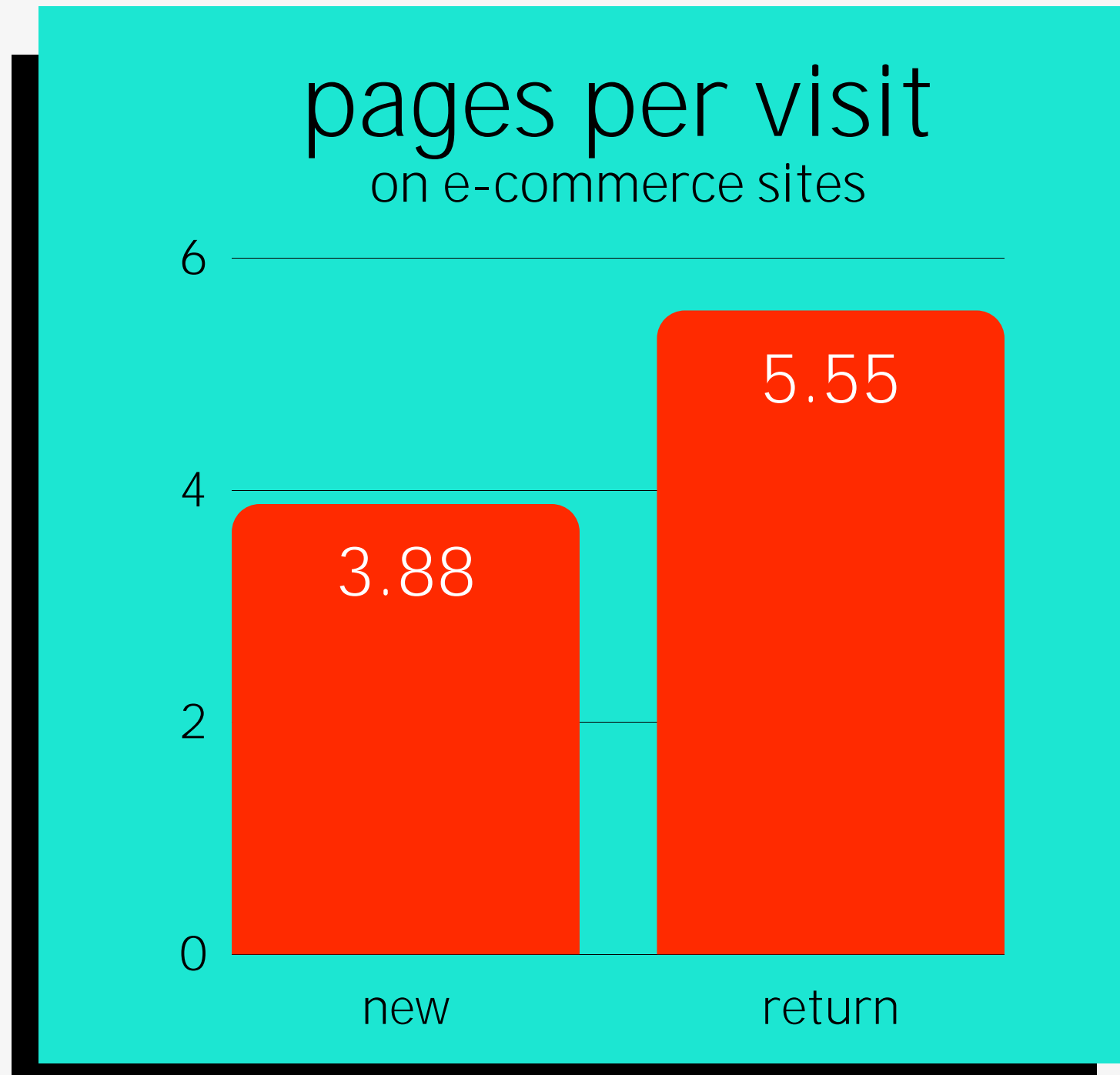
# How much data?

**Simulate real-life scenario**

Page-by-page collection

- Lazada: 40 listings/page
- Shopee: 50 listings/page

Target:

- Lazada: 4 pages = 160 listings
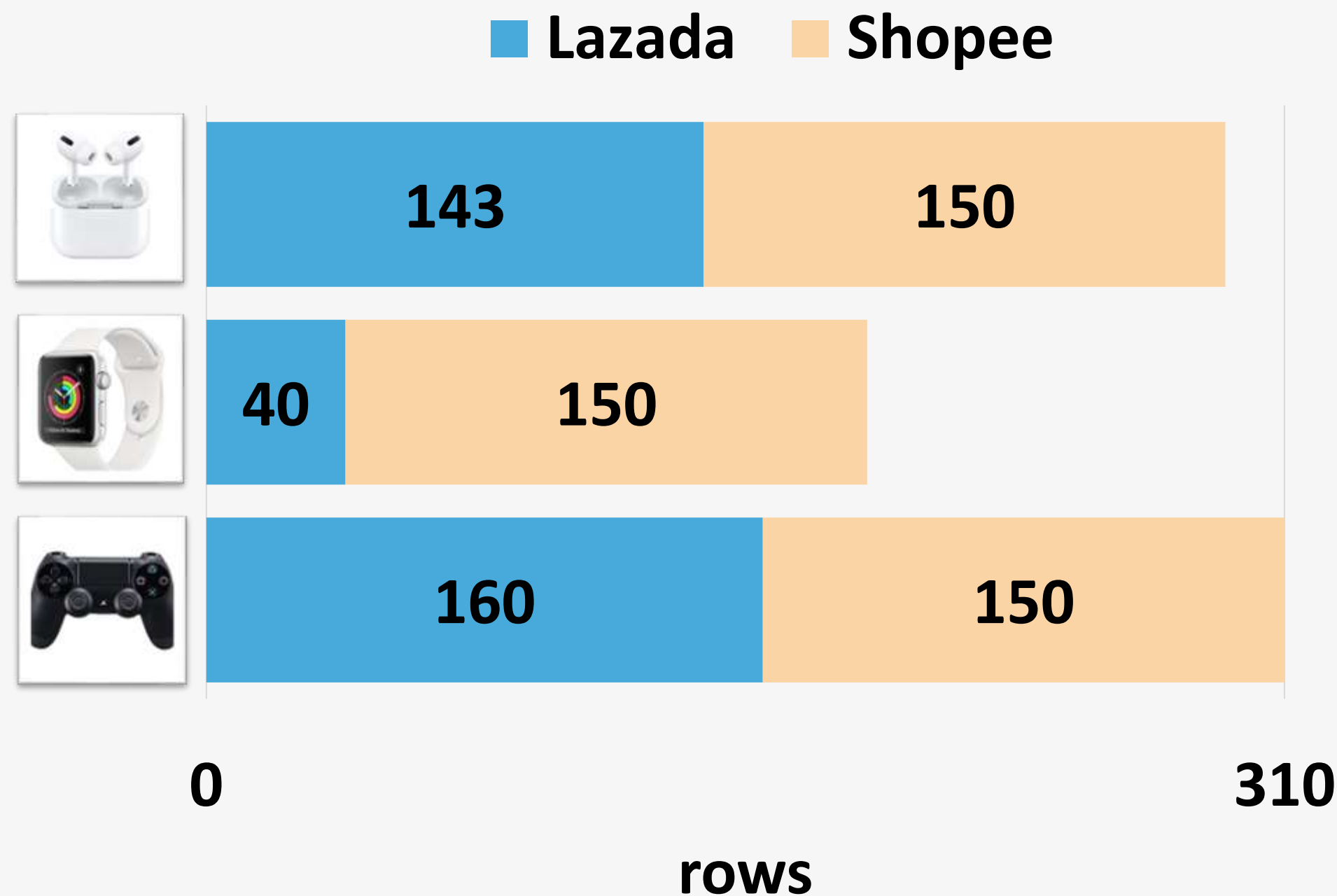- Shopee: 3 pages = 150 listings
- **Total: 310 per product**

## pages per visit
on e-commerce sites

| | new | return |
|---|---|---|
| value | 3.88 | 5.55 |

# We now have data!

## We collected six datasets.

**Lazada**  **Shopee**

| | Lazada | Shopee |
|---|---|---|
| AirPods | 143 | 150 |
| Apple Watch | 40 | 150 |
| Controller | 160 | 150 |

0 — 310

**rows**

# Why are there missing data?

- Listing is displayed but product page is unavailable.
- Not enough listings to display.

# We now have data!

**We collected twelve features.**

**Categorical (4)**
'ProductURL'
'ProductTitle'
'SellerName'
'ChatResponseTime'

**Numerical (7)**
'IsMall'                                      'SellerRating'
'Price'                                       'ShipOnTime'
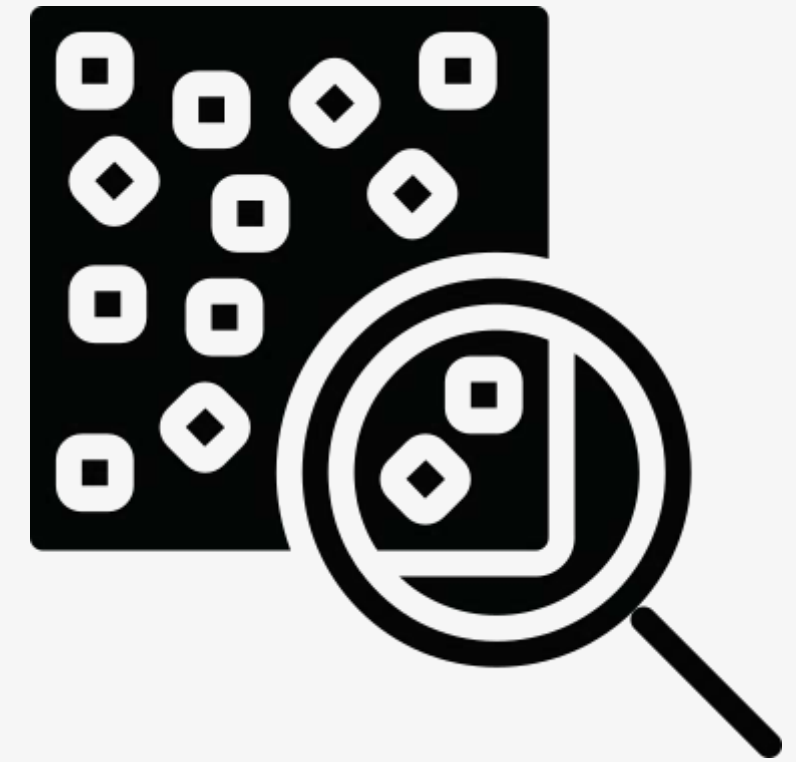'ProductRatingAverage'        'ChatResponseRate'
'ProductRatingCount'

**Target**
'Class'

# Target variable

**'Class', dtype: int**

We manually labeled each row of data, as follows:

- 0 - Counterfeit
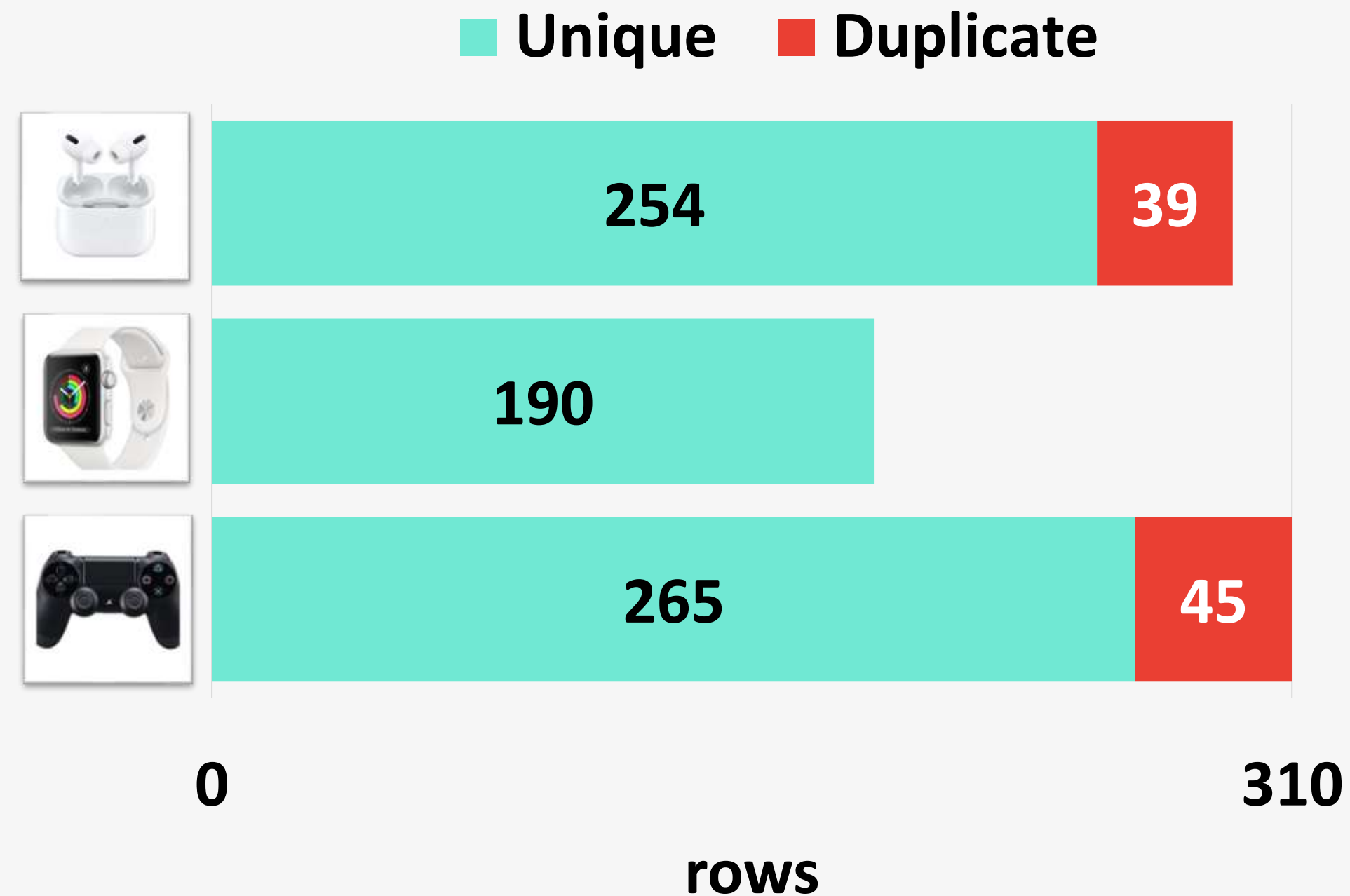- 1 - Authentic
- 2 - Unsure
- 3 - Different Product

Each entry was reviewed by another member
at least once.

# Data Preparation

**Drop duplicates**

# Data Preparation

## Handle missing values

**ChatResponseTime** : unavailable on Lazada product pages; difficult to scale

**ShipOnTime** : not provided by Shopee

**Resolution** : drop both features

    (1) high missing value count

    (2) arguably irrelevant

within minutes
within hours
Active in: 1 hours
Active in: 10 mins
Active in: hours

# Data Preparation

## Handle incorrect column dtype

These features should be float-type, but contain object-type data.

ChatResponseRate : 'not enough data'

ProductRatingAverage : 'No ratings yet'

SellerRating : 'New Seller', 'No ratings yet', '0 R'

# Data Preparation

## Handle incorrect column dtype

**Resolution** : replace with 0

| ProductRatingAverage | SellerRating | ChatResponseRate |
|---|---|---|
| 4.6 | New Seller | 1 |
| 5 | No ratings yet | 0.97 |
| No ratings yet | 0 R | 0.57 |
| 4.9 | 0.94 | not enough data |
| No ratings yet | 4.6 | 0.97 |

→

| ProductRatingAverage | SellerRating | ChatResponseRate |
|---|---|---|
| 4.6 | 0 | 1 |
| 5 | 0 | 0.97 |
| 0 | 0 | 0.57 |
| 4.9 | 0.94 | 0 |
| 0 | 4.6 | 0.97 |

# EDA

## 'Class' distribution of each product



0 – Counterfeit | 1 – Authentic | 2 – Unsure | 3 - Different Product

# EDA

## Word cloud of common terms in counterfeit product labels





LazMall
Apple AirPods Pro

# EDA

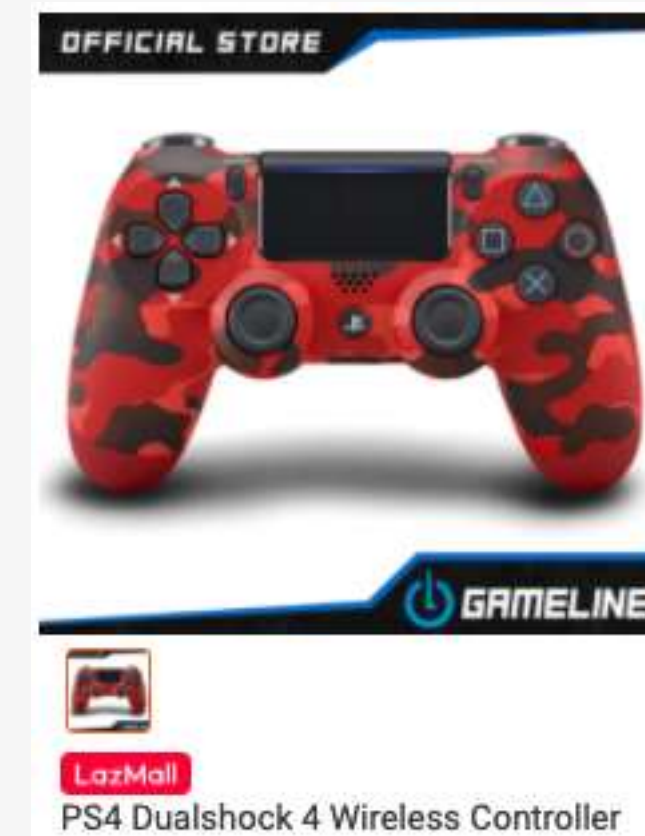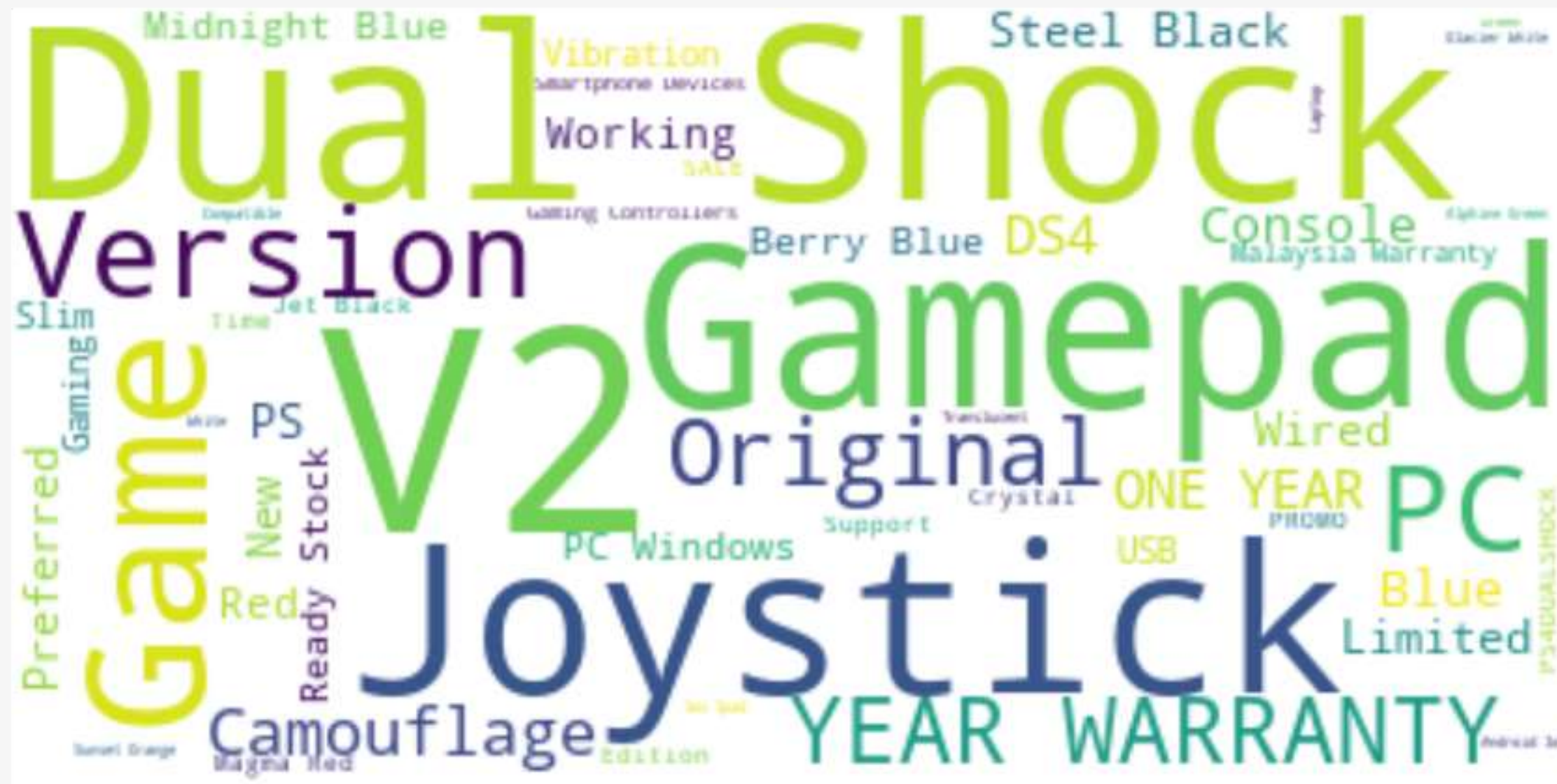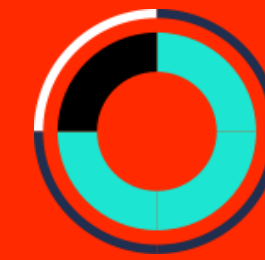**Word cloud of common terms in counterfeit product labels**

# EDA

**Word cloud of common terms in counterfeit product labels**

# Modelling

## Getting Started

**Data:**

- 6 datasets, 6 (numerical) features, 1 target

**Model:**

- Supervised
- Classification
  - ➤ Voting: Hard, Soft (Uniform), Soft with GridSearchCV

➤ Logistic Regression (scaled)

➤ K-Nearest Neighbors

➤ Gaussian Naïve-Bayes (scaled)

➤ Decision Tree

➤ Random Forest

- all default except for max_depth = 10 for DT
- 5-fold cross validation

# Metrics

**Lazada 'airpods pro': best parameters**

```
param_grid = {
    'knn__n_neighbors' : [1,3],
    'rf__max_depth': [1, 4, 7],
    'rf__criterion' : ['gini', 'entropy']
    }
```

```
Best parameters = {
    'knn__n_neighbors': 1,
    'rf__max_depth': 4
    'rf__criterion': 'gini',
    }
```
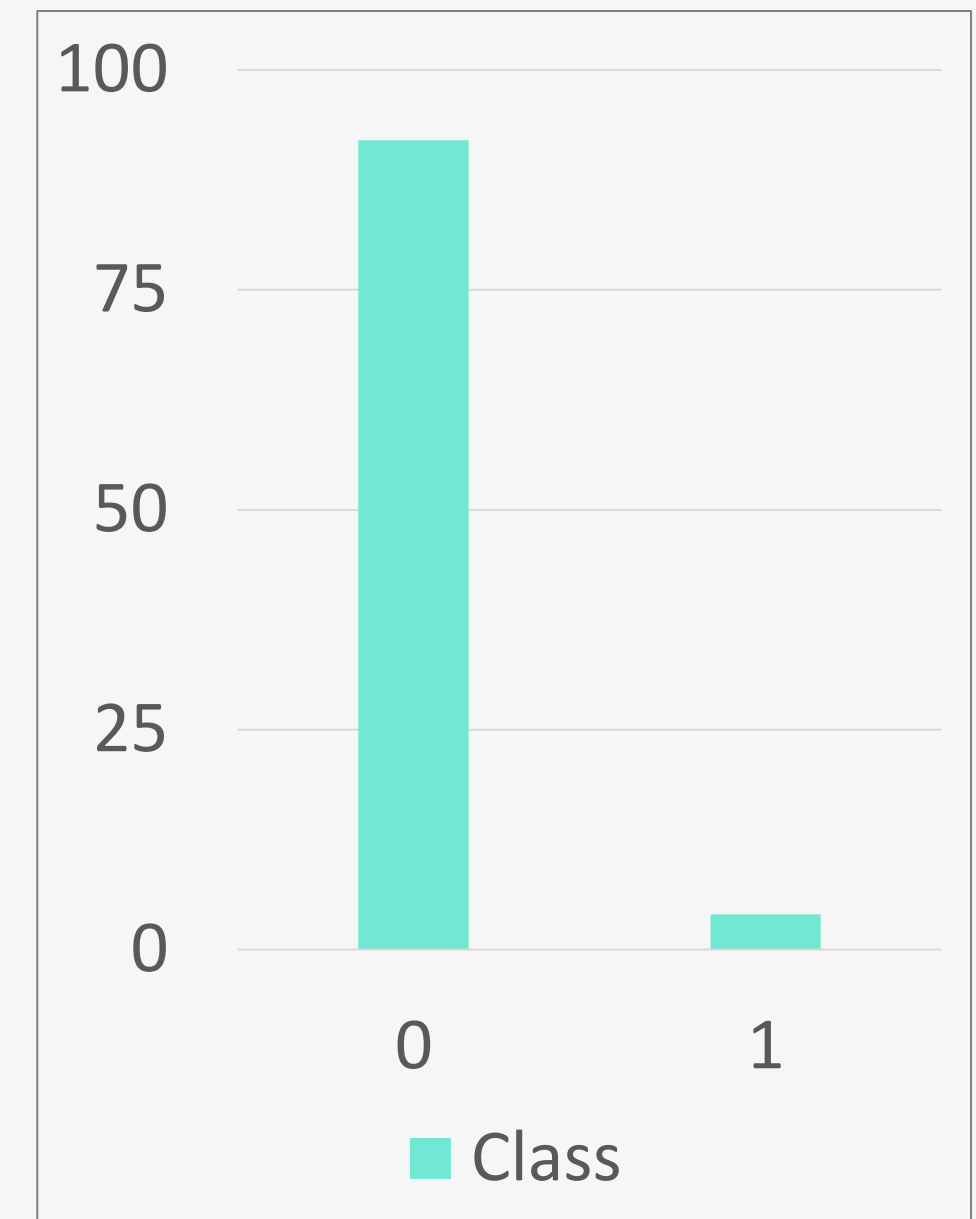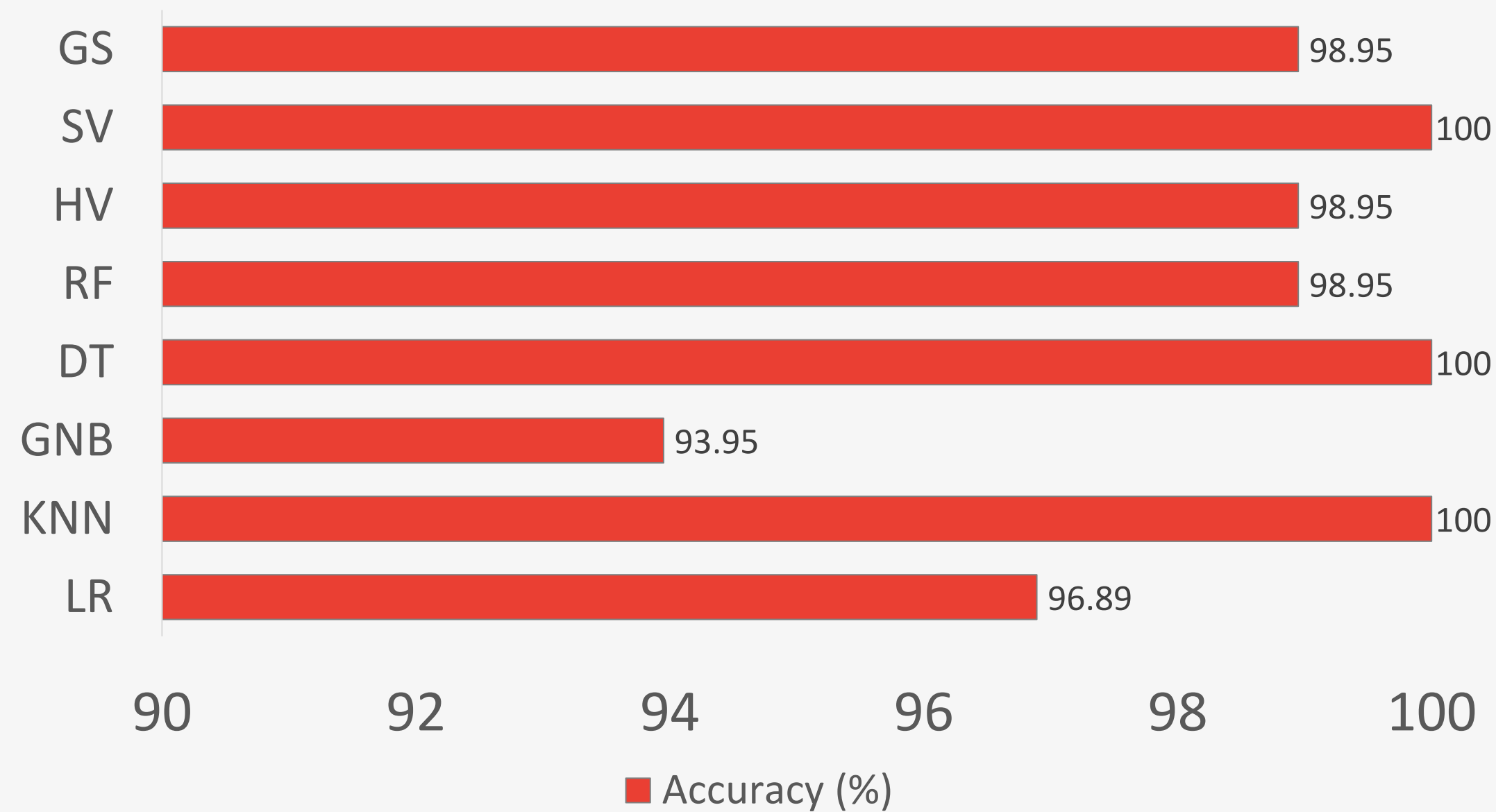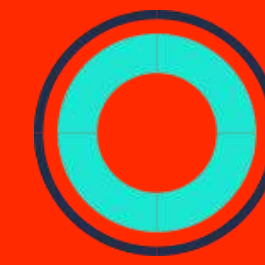
# Metrics

## Lazada 'airpods pro': accuracy

# Metrics

## Lazada 'airpods pro': other scores

'Class' = 0

'Class' = 1

F1 Score ■ Recall ■ Precision

# Metrics

**Lazada 'airpods pro': best estimator**



'Class' = 1

# Metrics

## All products: best estimator

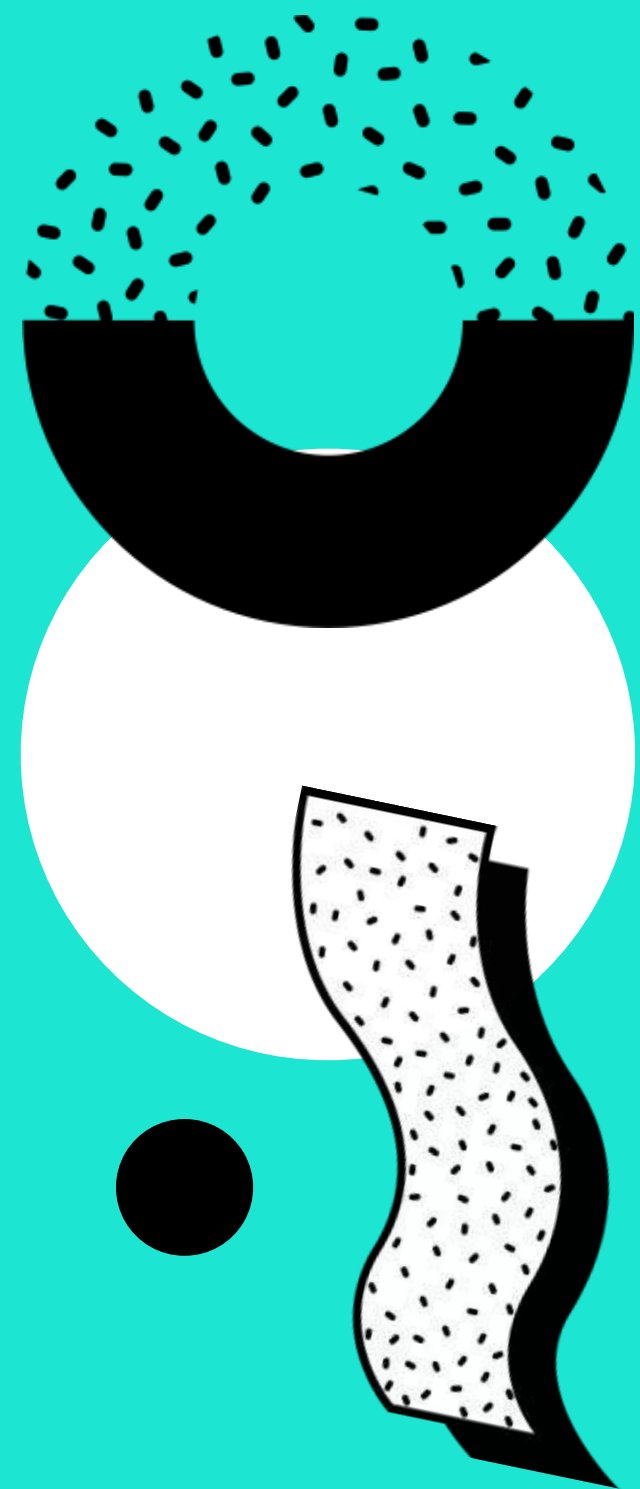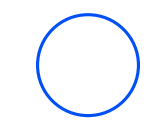| Dataset | Best estimator | F1 score |
| --- | --- | --- |
| Lazada 'airpods pro' | Soft Voting with GridSearchCV | 0.81 |
| Lazada 'apple watch' | LR | 1 |
| Lazada 'ps4 controller' | Soft Voting with GridSearchCV | 0.77 |
| Shopee 'airpods pro' | Soft Voting with GridSearchCV | 0.78 |
| Shopee 'apple watch' | Soft Voting with GridSearchCV | 0.85 |
| Shopee 'ps4 controller' | Soft Voting with GridSearchCV | 0.86 |

# Metrics
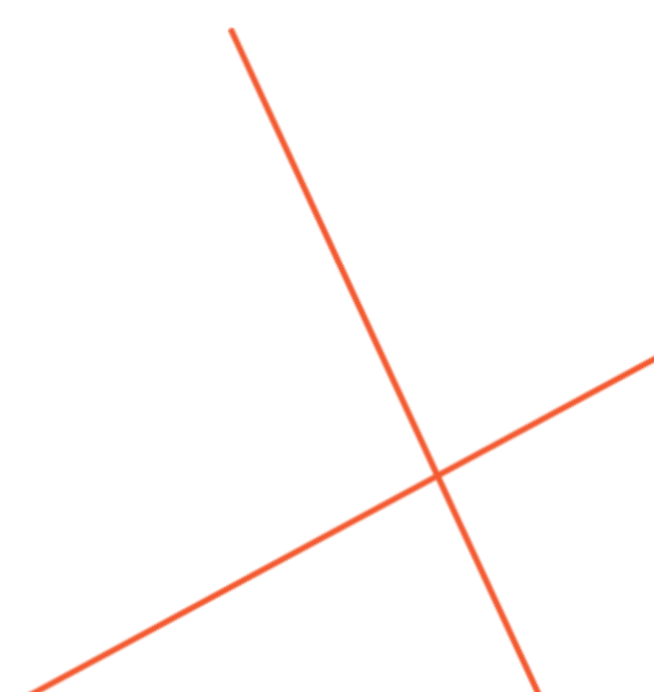
**All products: Feature Importance**

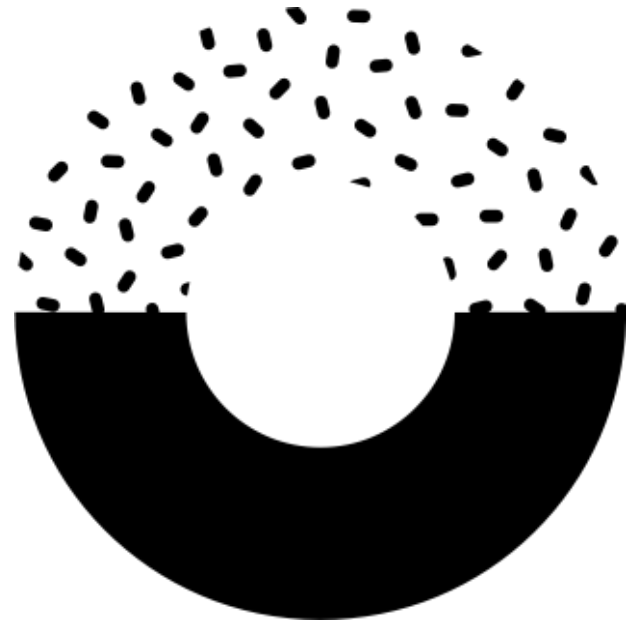# Summary

and Conclusions

# In summary,

- **Detection of counterfeit items** on e-commerce platforms can be achieved through machine learning

- **Price, Seller Rating, and authentic certification** (e.g., LazMall and ShopeeMall) are key indicators of authenticity
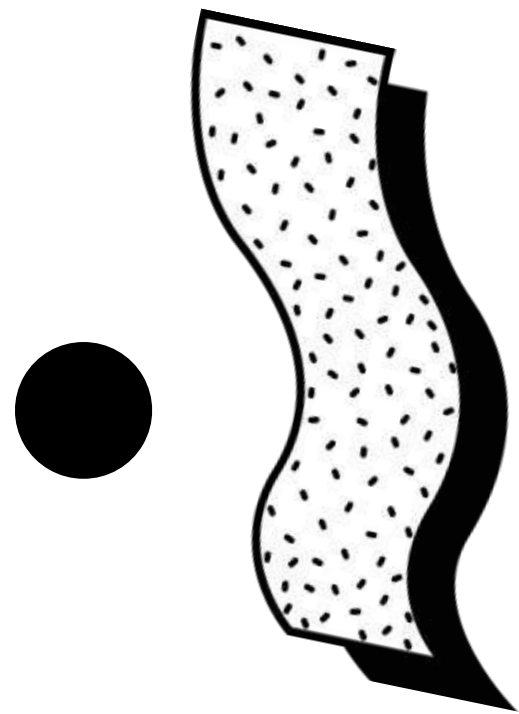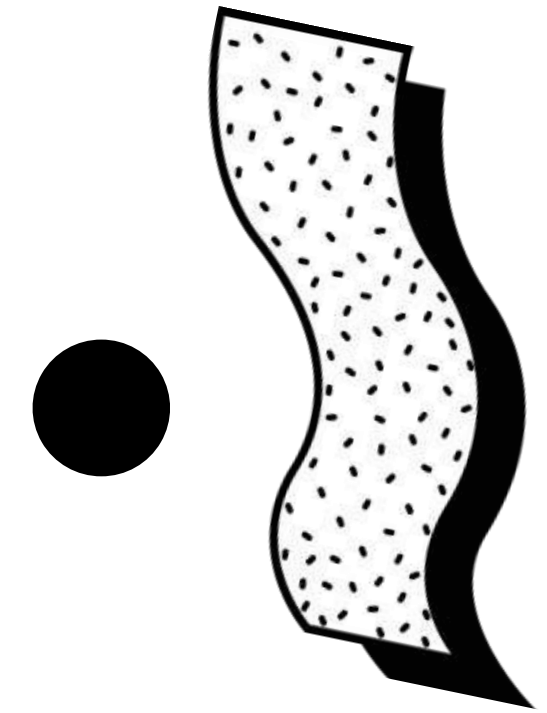
1. How will we scale this in the real world?

2. Points of improvement

# Thank you for listening!

Presented by:

Romar Acejo

Chris Dion Bautista

Ladie Guevarra