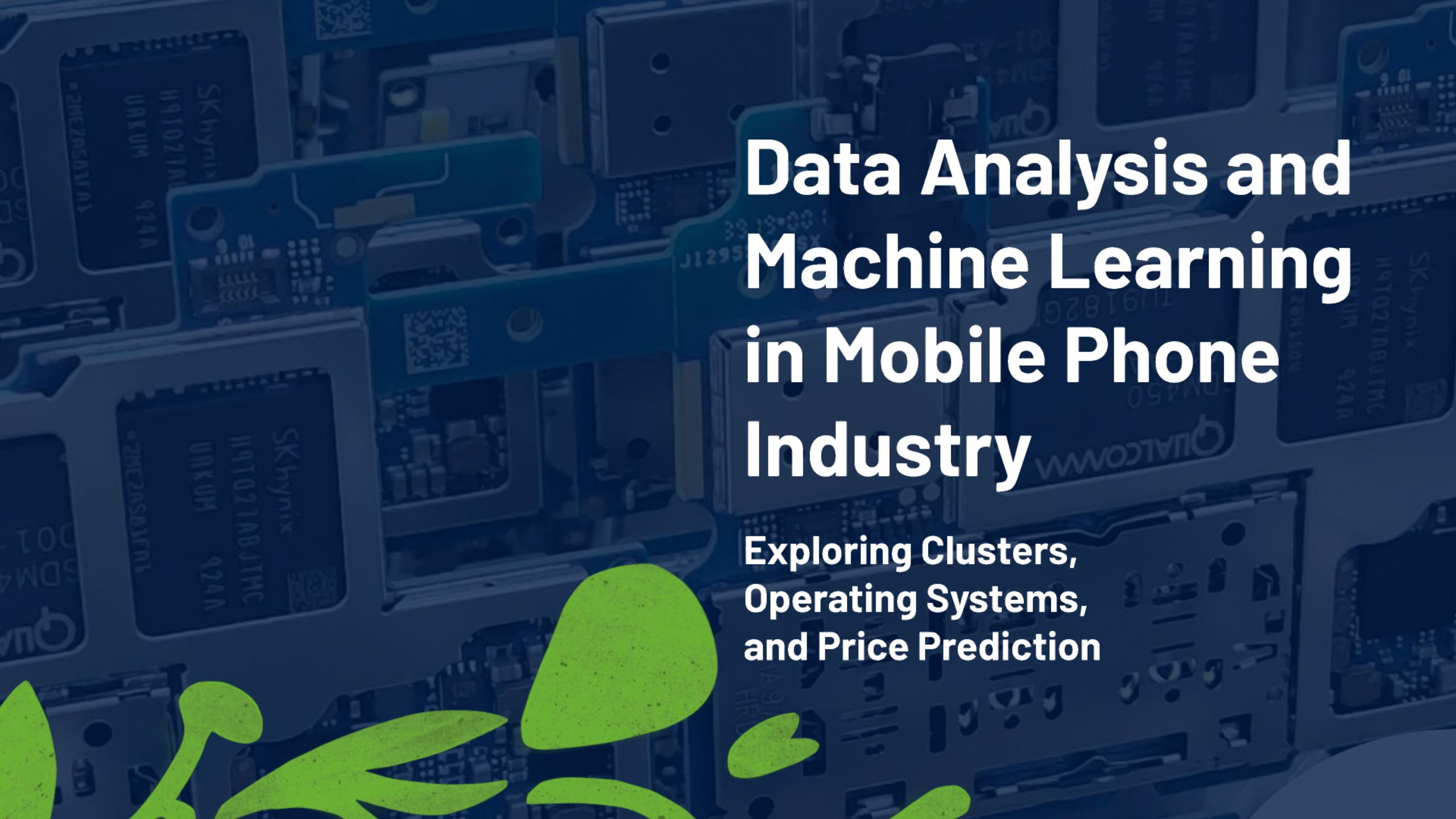


# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction





# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction

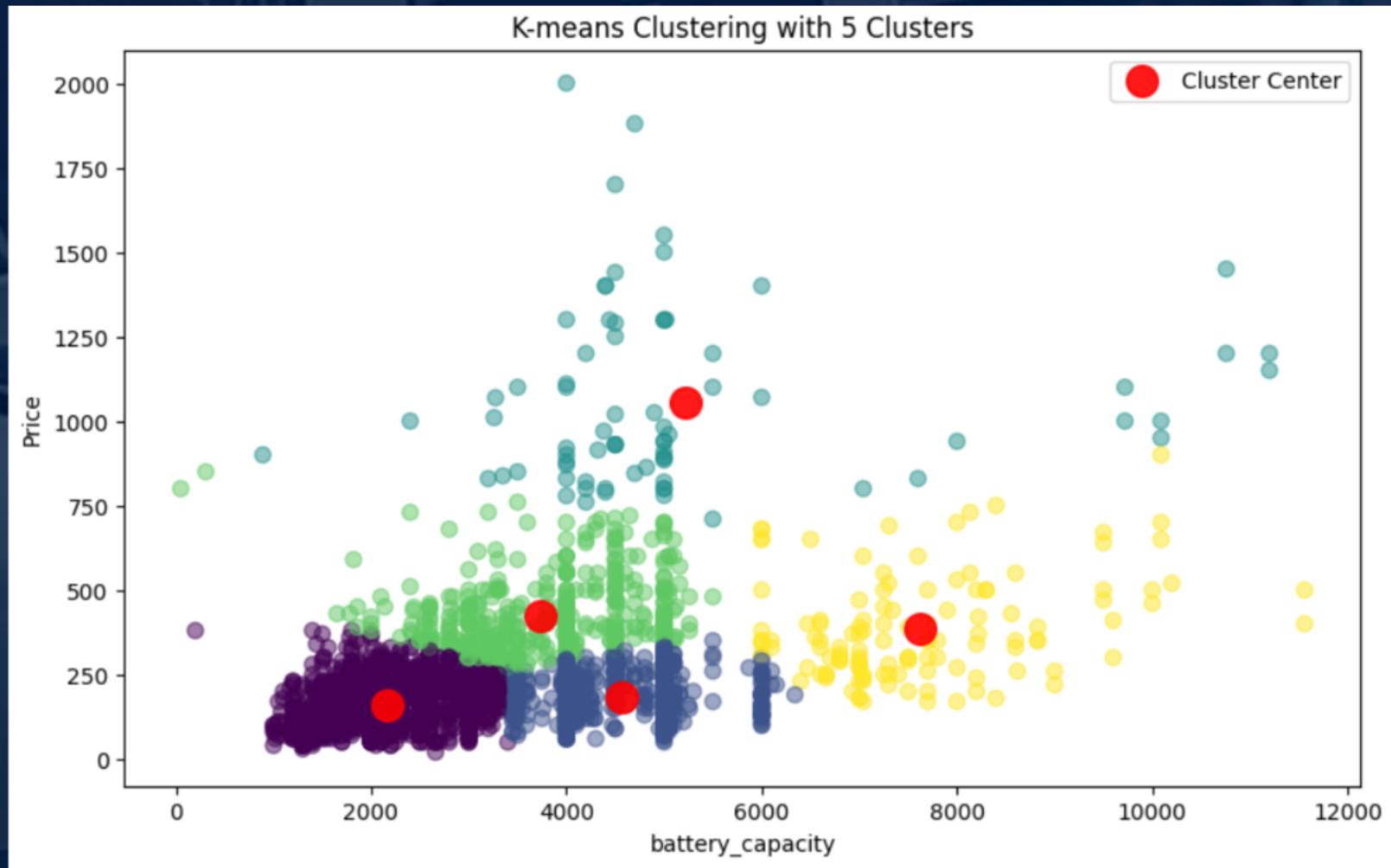
# Clustering based on Battery Capacity and Price



Explore K-means and DBScan clustering for mobile phone data based on battery capacity and price.

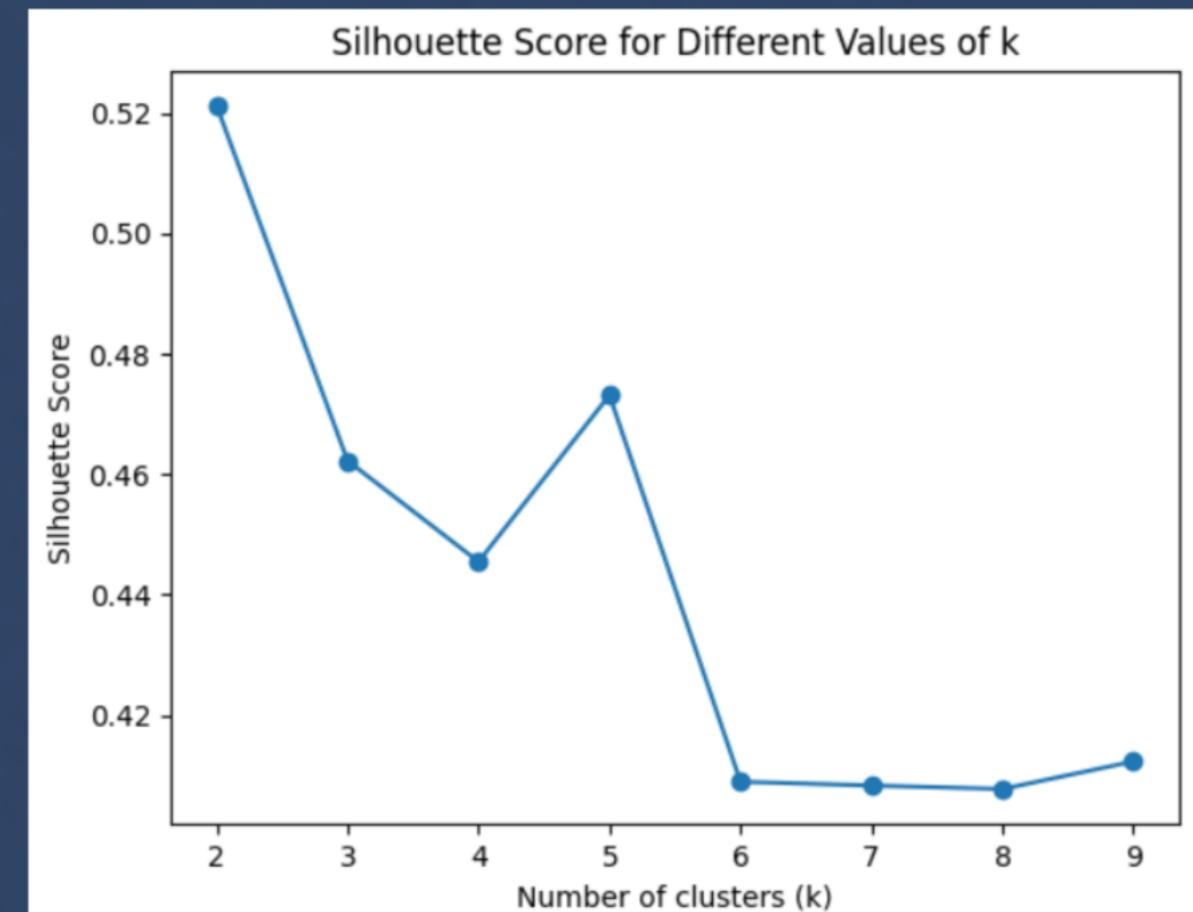
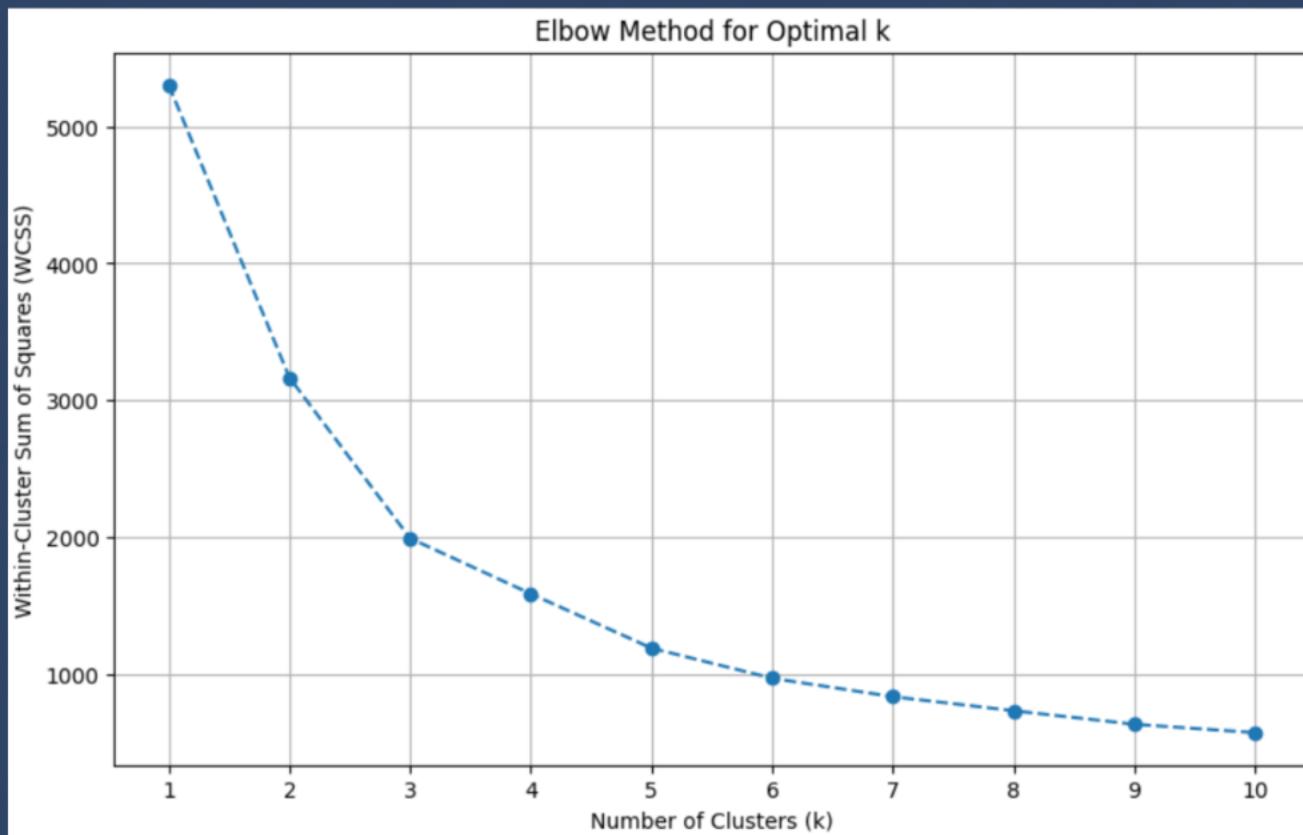
# K-means Clustering

Executing K-means algorithm with 5 clusters  
scatter plot and cluster centers.



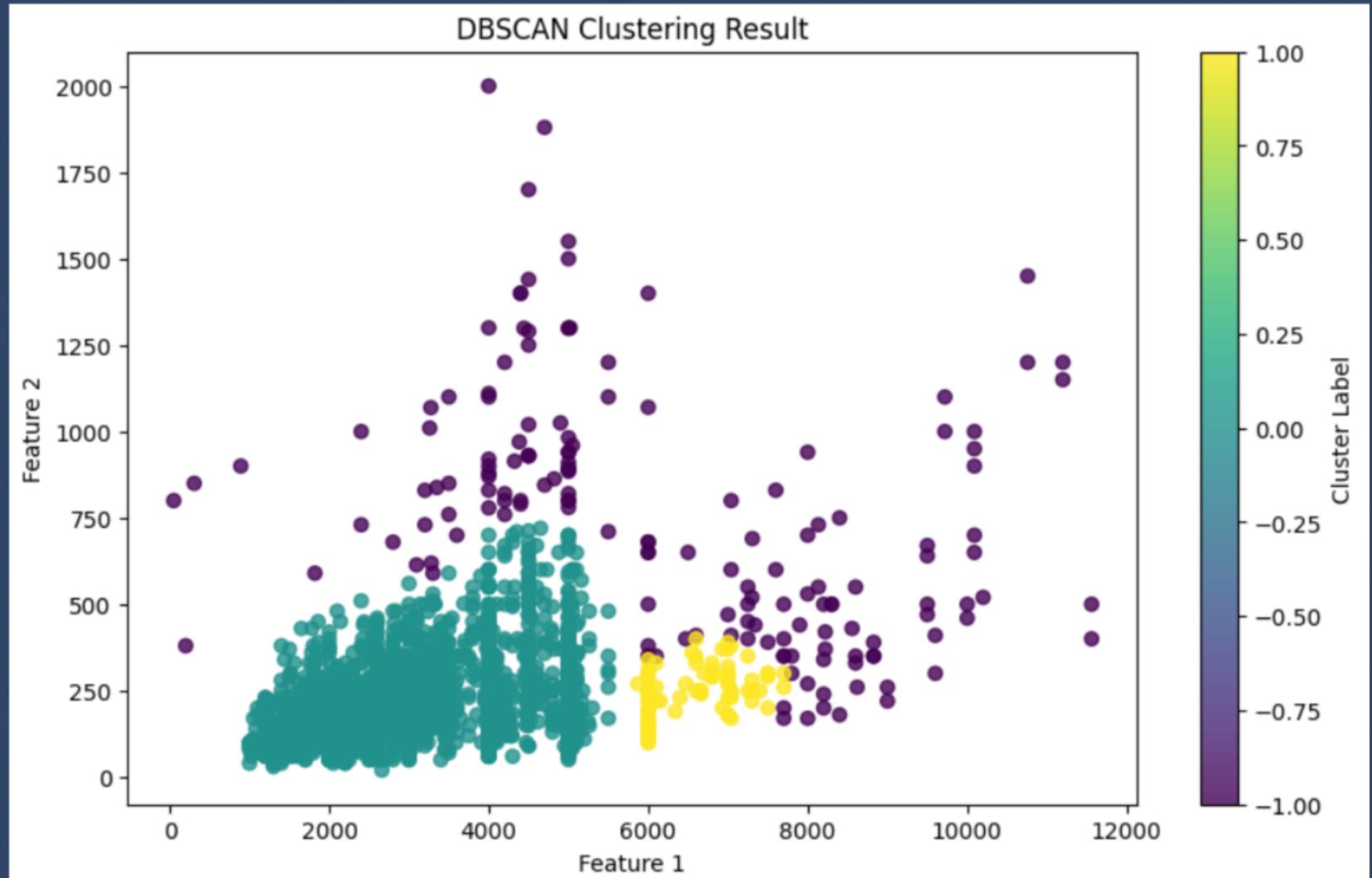
# Optimal K-value Selection

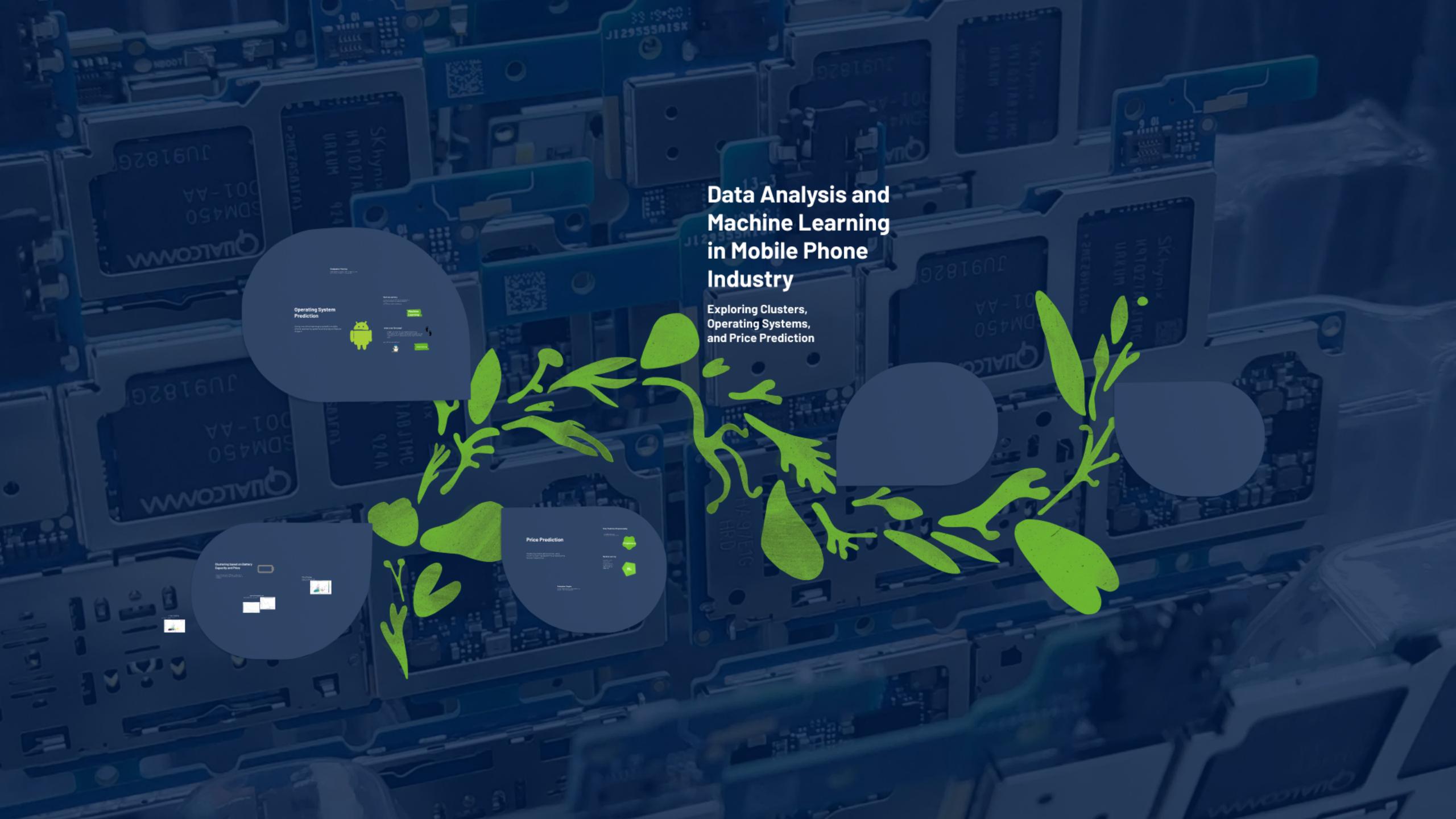
Running K-means for k=1 to k=10 to choose optimal k based on Elbow Method and Silhouette Score.



# DBScan Clustering

Clustering data using DBScan with modified hyperparameters to produce 3 meaningful clusters(including noise).





# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction



Price Prediction



The Predicted Prices

# Operating System Prediction

Using machine learning to predict mobile phone operating systems and analyze feature impact.



# what is our first step?

- Alright, in order to use machine learning techniques for predicting operating systems, we need to perform data preprocessing in the first step.



SO LET'S DIVE INTO IT



preprocessing

Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone



We want to show you what we  
have done with our data and  
what insights we have gained.

I'm trusting you

Let's begin!

Columns

2G 3G 4G 5G

Tools

Weight

Length, Width,  
Diameter

SIM

Display

ppi

Body Ratio

And more

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

# We have bunch of columns

```
data.columns
```

✓ 0.0s

```
Index(['name', 'brand', '2G', '3G', '4G', '5G', 'Announced', 'Status',
       'Weight', 'Length', 'Width', 'Diameter', 'SIM', 'Display Type',
       'Display Size', 'ppi', 'body ratio', 'os', 'battery_capacity', 'Price',
       'CPU', 'ratio', 'pixel', 'WLAN', 'Colors', 'Sensors', 'Bluetooth',
       'GPU', 'Loudspeaker', '3.5mm jack', 'Chipset', 'Network', 'Internal',
       'Card slot', 'RAM', 'Storage'],
      dtype='object')
```

Our expert  
team has  
decided to  
move this  
comment in  
more detail.

so we have bunch of OS

```
OS
```

```
data['OS'] = data['os'].str.extract(r'^(.+)\s')
data['OS'].unique()
✓ 0.0s
array(['Android', 'iOS', 'Windows', 'Firefox', 'OS', 'iOS', 'Windows',
       'Phone', 'Mobile', 'Macintosh', 'HTML', 'platform', 'Hello',
       'None', 'System', 'OS', 'Norton', 'Linux', 'iTunes'], dtype=object)
```

so we have bunch of OS

## OS

```
data['os'] = data['os'].str.extract(r'^(\w+)')
```

```
data['os'].unique()
```

✓ 0.1s

```
array(['Android', 'KaiOS', 'Windows', 'Firefox', 'os', 'ios', 'iPadOS',
       'Phone', 'Mobile', 'HarmonyOS', 'EMUI', 'platform', 'Belle',
       'Anna', 'Symbian', 'FP1', 'MeeGo', 'Linux', 'Tizen'], dtype=object)
```

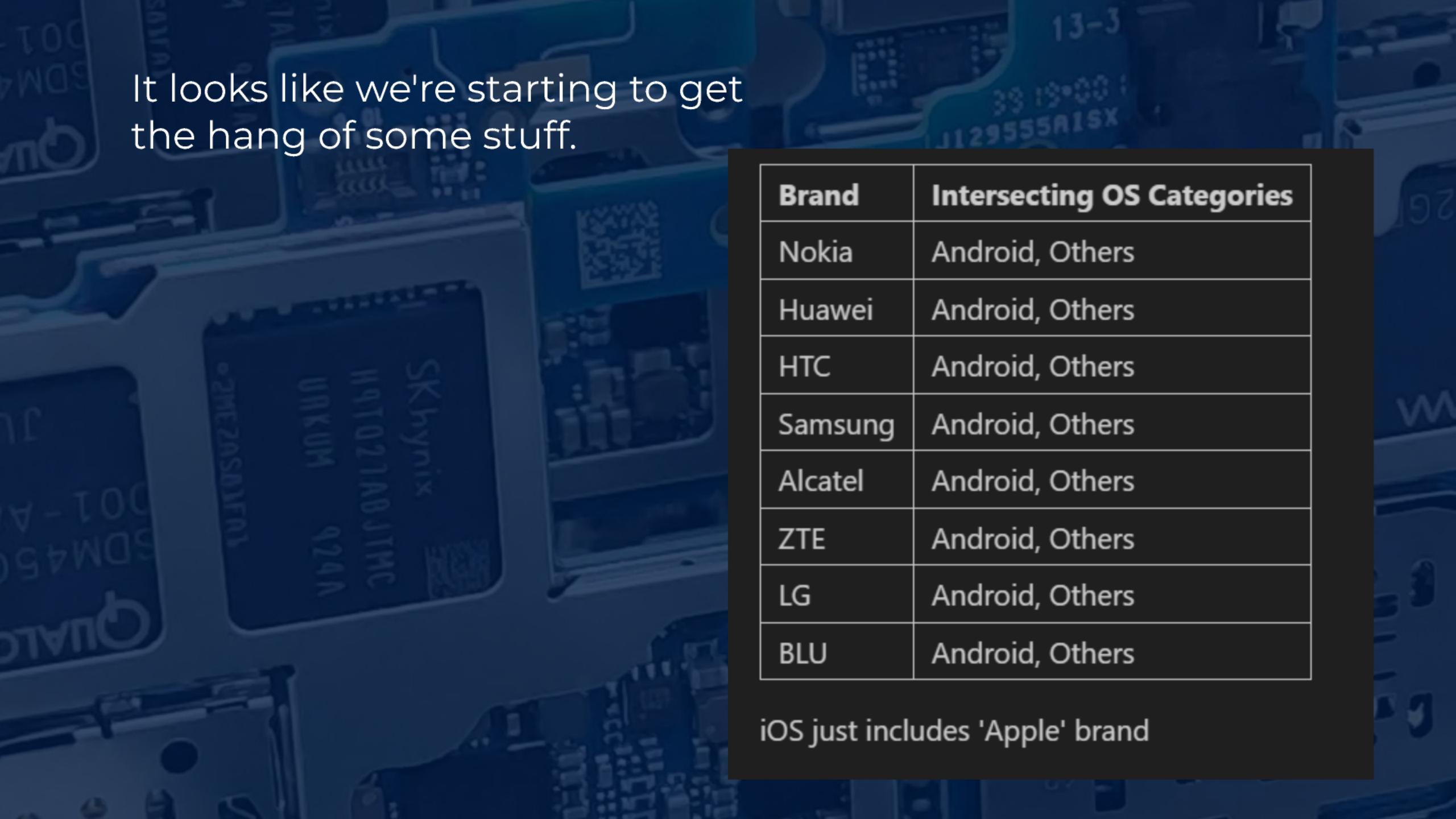
and bunch of brands too

brand

```
data.groupby('os')['brand'].value_counts()
```

| os      | brand   | count |
|---------|---------|-------|
| Android | samsung | 521   |
|         | xiaomi  | 285   |
|         | huawei  | 250   |
|         | lg      | 226   |
|         | zte     | 202   |
|         | htc     | 187   |
|         | lenovo  | 174   |
|         | blu     | 146   |
|         | sony    | 122   |
|         | asus    | 113   |
| Others  | alcatel | 103   |
|         | nokia   | 82    |
|         | infinix | 38    |
|         | nokia   | 67    |
|         | huawei  | 27    |
|         | htc     | 17    |
|         | samsung | 11    |
| ios     | alcatel | 8     |
|         | zte     | 4     |
|         | lg      | 2     |
|         | blu     | 1     |
| ios     | apple   | 65    |

Name: count, dtype: int64



It looks like we're starting to get  
the hang of some stuff.

| Brand   | Intersecting OS Categories |
|---------|----------------------------|
| Nokia   | Android, Others            |
| Huawei  | Android, Others            |
| HTC     | Android, Others            |
| Samsung | Android, Others            |
| Alcatel | Android, Others            |
| ZTE     | Android, Others            |
| LG      | Android, Others            |
| BLU     | Android, Others            |

iOS just includes 'Apple' brand



Our expert team has identified the need to delve into these columns in more detail.

Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

Columns

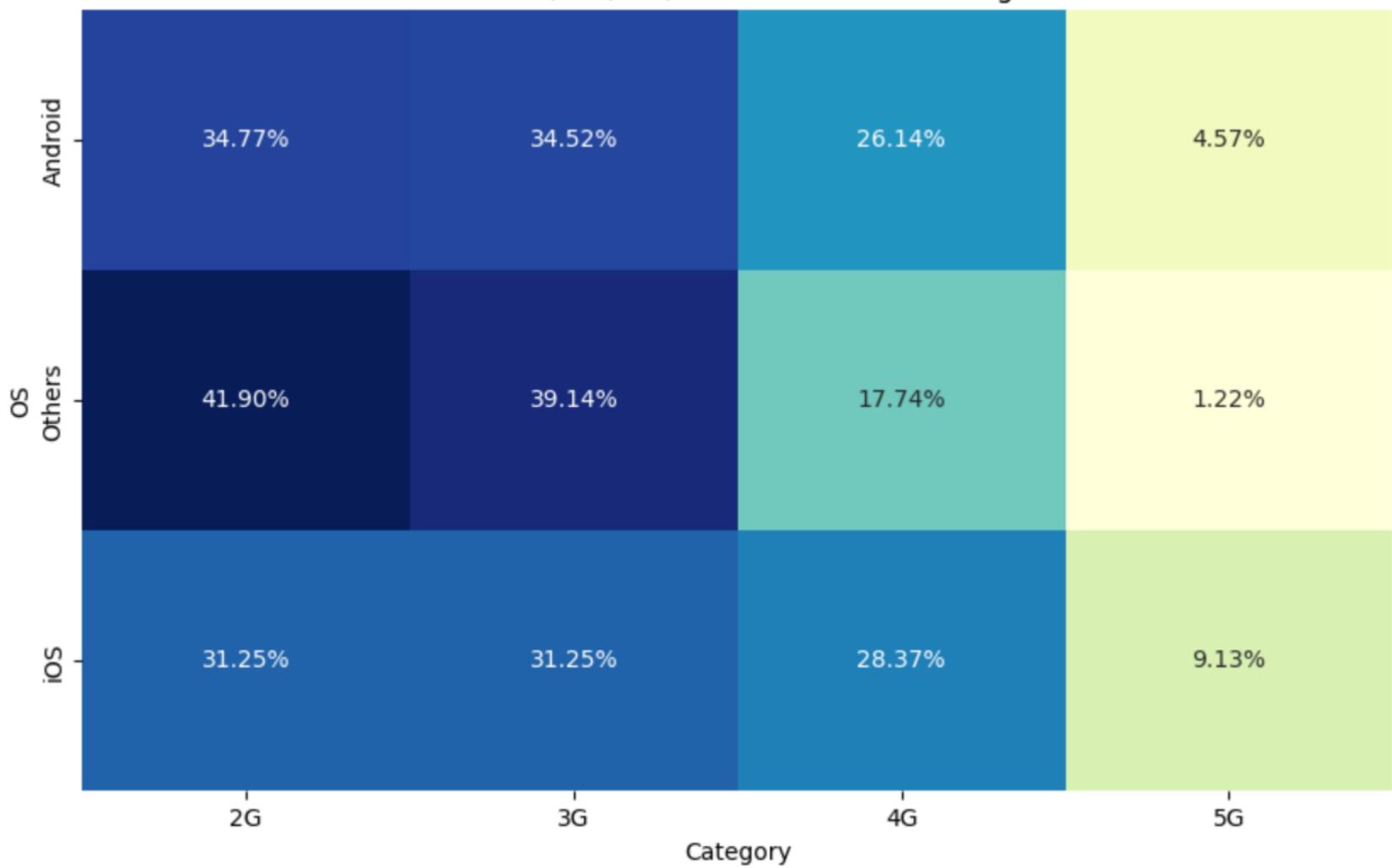
2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

Ratio of 2G, 3G, 4G, and 5G within OS categories



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

Along the way, we need to develop some tools for ourselves.

```
def calculate_and_plot_ratio(data, group_columns, value_column):  
    # calculate the total count  
  
    total_counts = data.groupby(group_columns[0])[value_column].transform('count')  
  
    # calculate the ratio  
  
    data['Ratio'] = data.groupby(group_columns)[value_column].transform('count') / total_counts  
  
    # Plotting  
  
    plt.figure(figsize=(12, 8))  
    sns.barplot(data=data, x=value_column, y='Ratio', hue=group_columns[0])  
    plt.title(f'Ratio of {value_column} by {group_columns[0]}')  
    plt.xlabel(value_column)  
    plt.ylabel('Ratio')  
    plt.legend(title=group_columns[0])  
    plt.show()
```

```
def plot_feature_distribution(data, feature):  
    data[feature] = pd.to_numeric(data[feature], errors='coerce')  
  
    plt.figure(figsize=(10, 6))  
    sns.boxplot(data=data, x='OS', y=feature)  
    plt.title(f'distribution of {feature} by OS')  
    plt.xlabel('OS')  
    plt.ylabel(feature)  
    plt.show()
```

```
def plot_feature_distribution(data, feature):
    data[feature] = pd.to_numeric(data[feature], errors='coerce')

    plt.figure(figsize=(10, 6))
    sns.boxplot(data=data, x='os', y=feature)
    plt.title(f'Distribution of {feature} by os')
    plt.xlabel('os')
    plt.ylabel(feature)
    plt.show()
```

Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

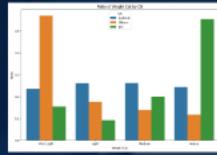
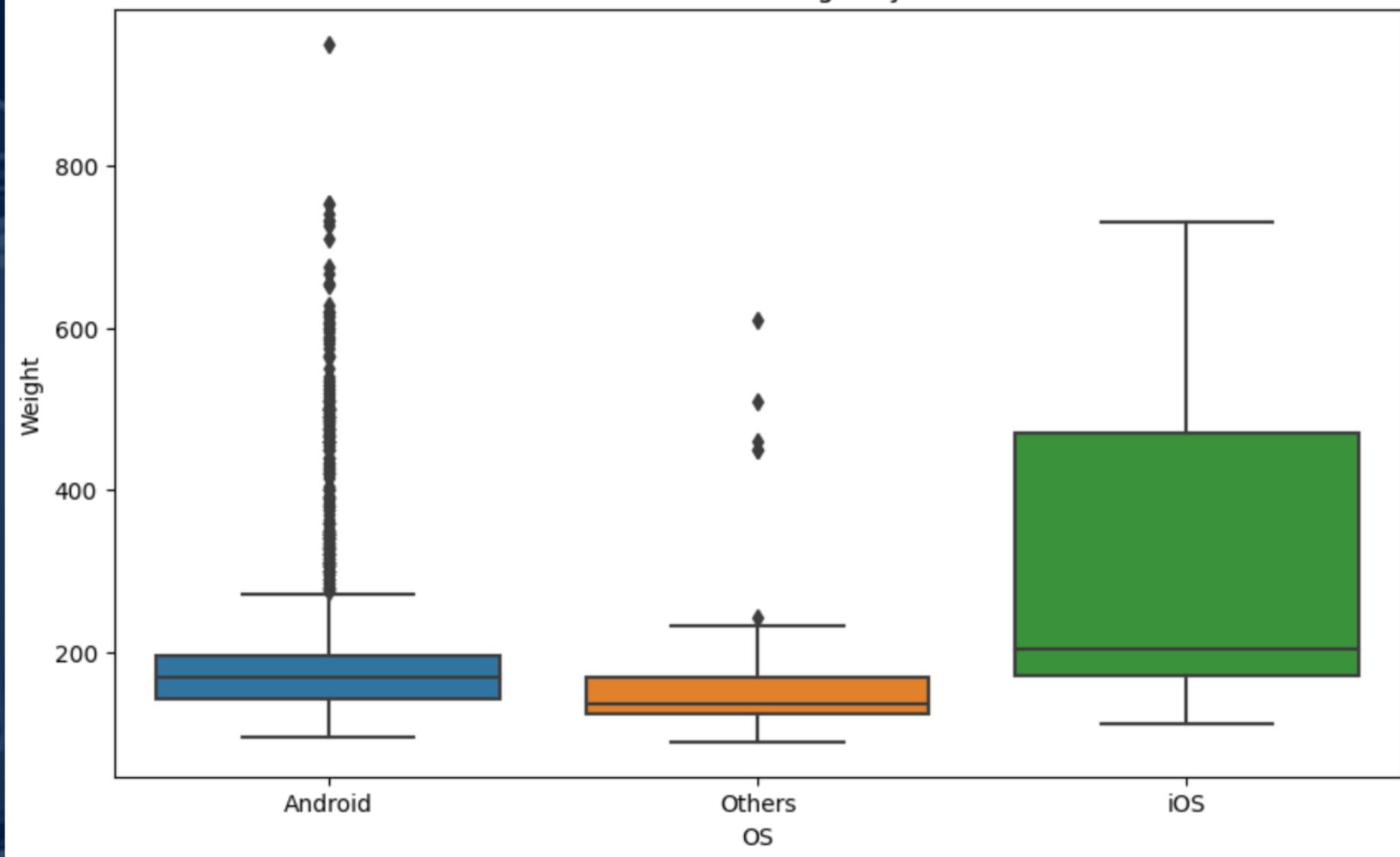
2G 3G 4G 5G

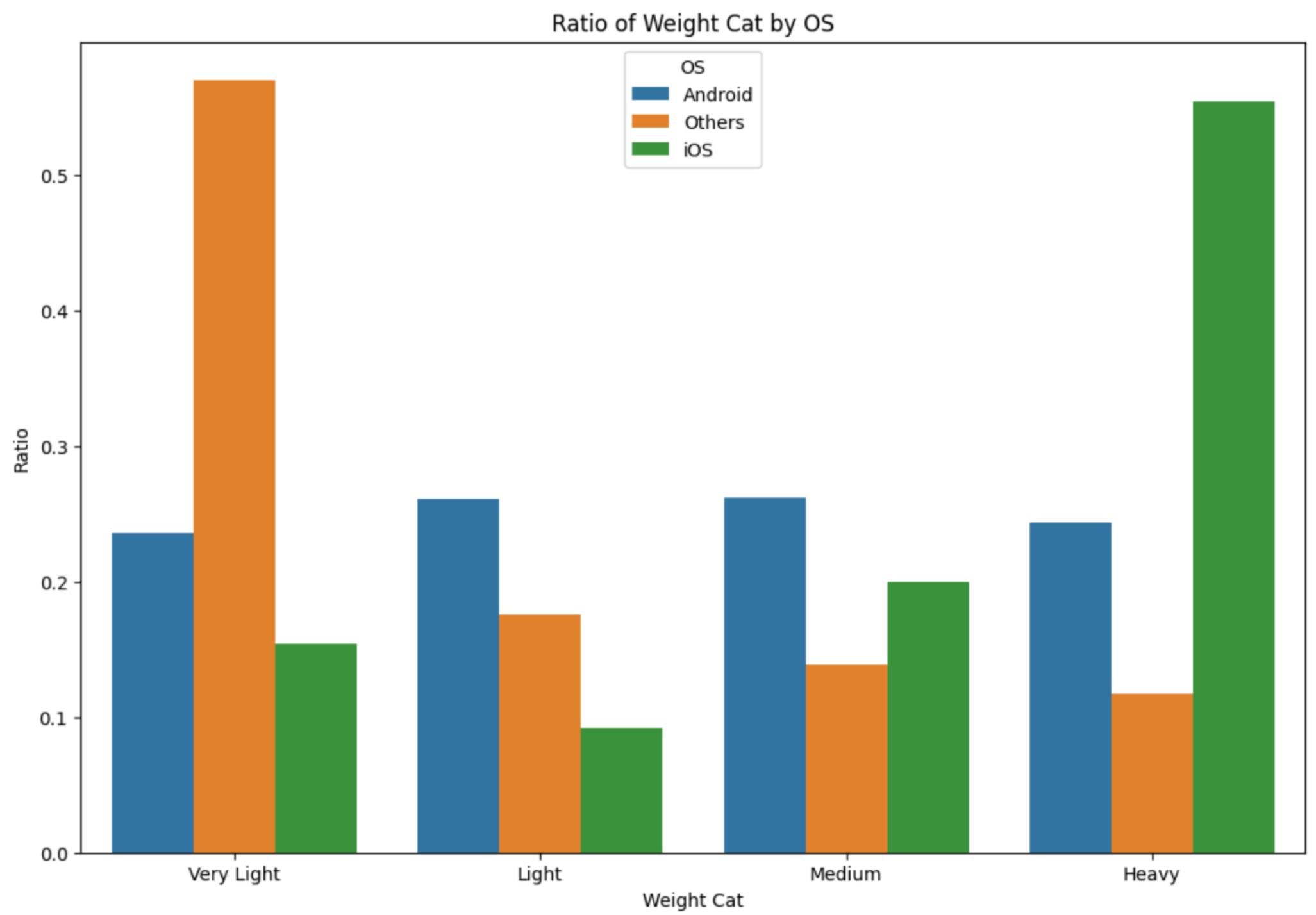
Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

## Distribution of Weight by OS





Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

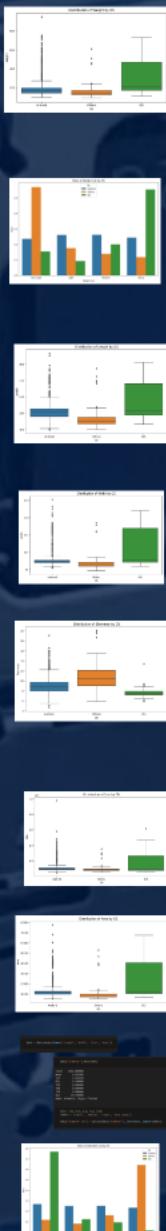
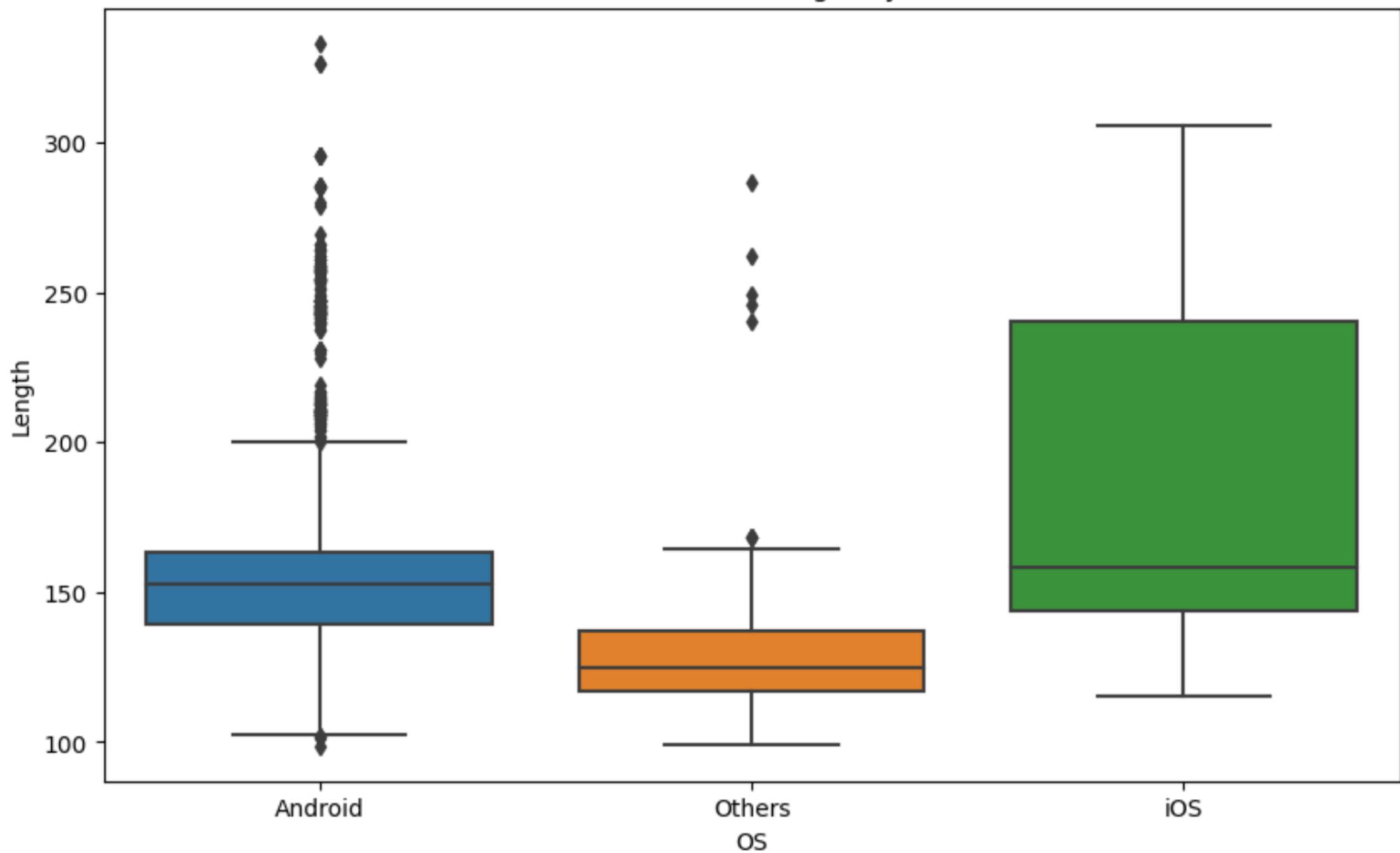
2G 3G 4G 5G

Tools

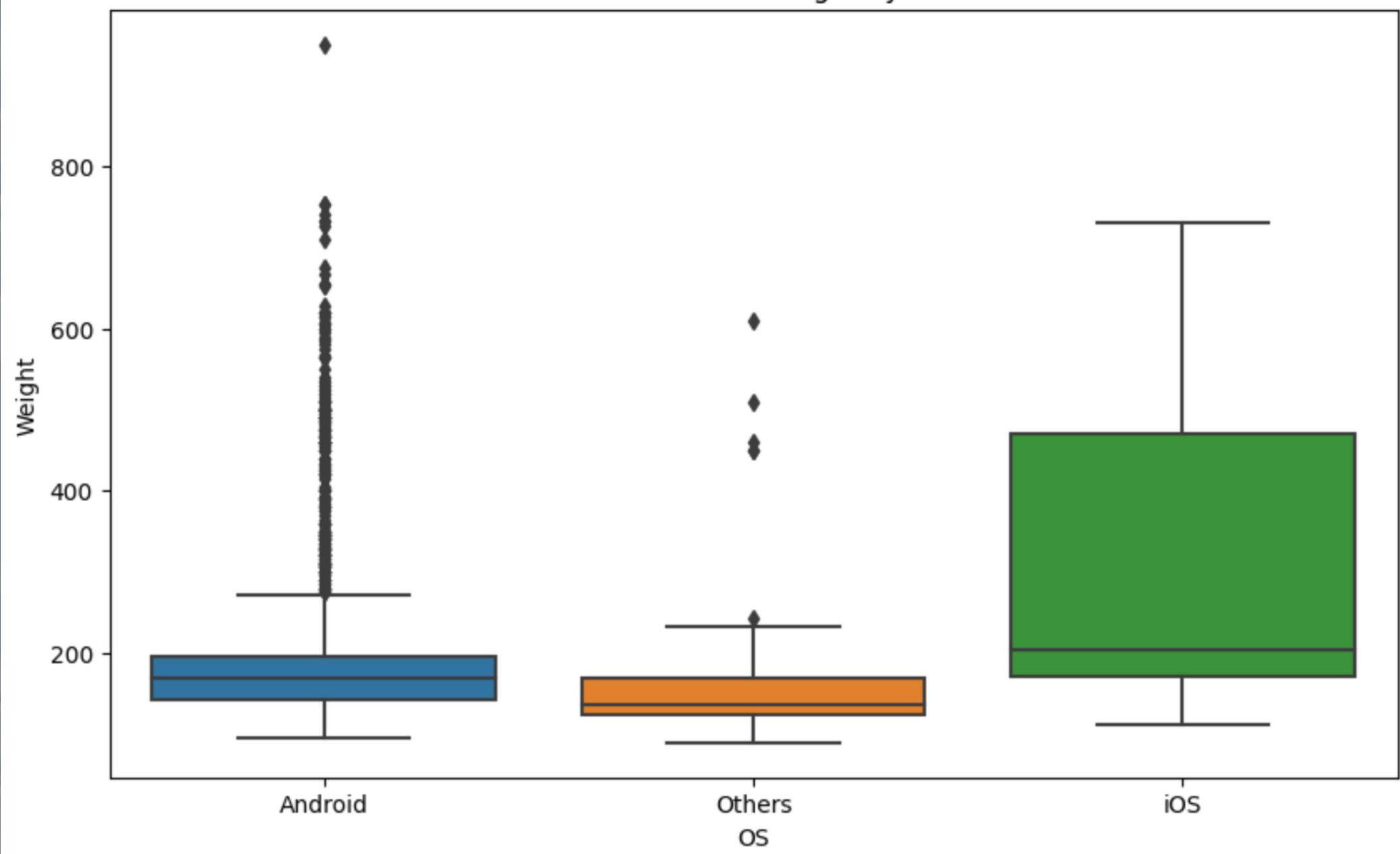
This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

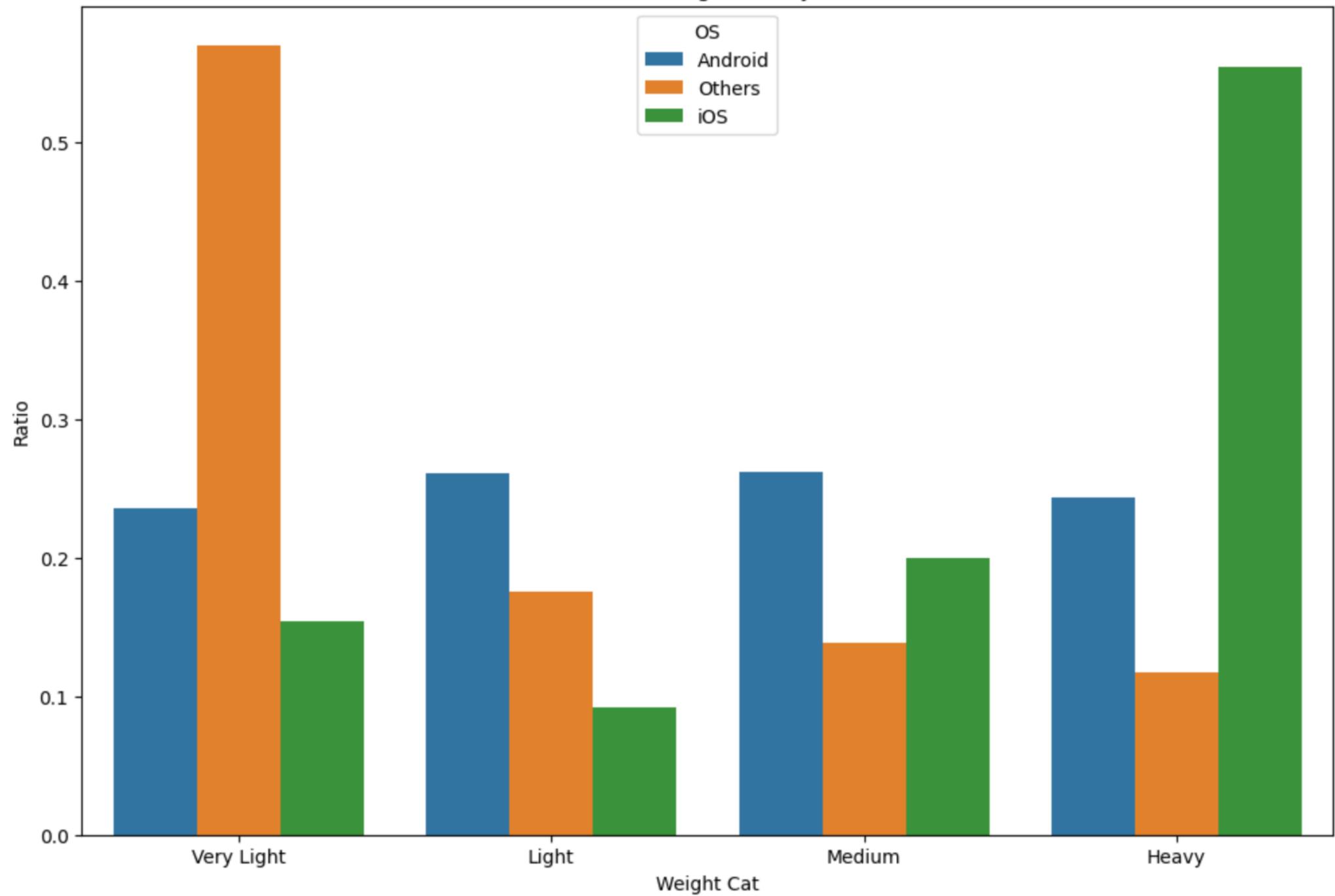
# Distribution of Length by OS



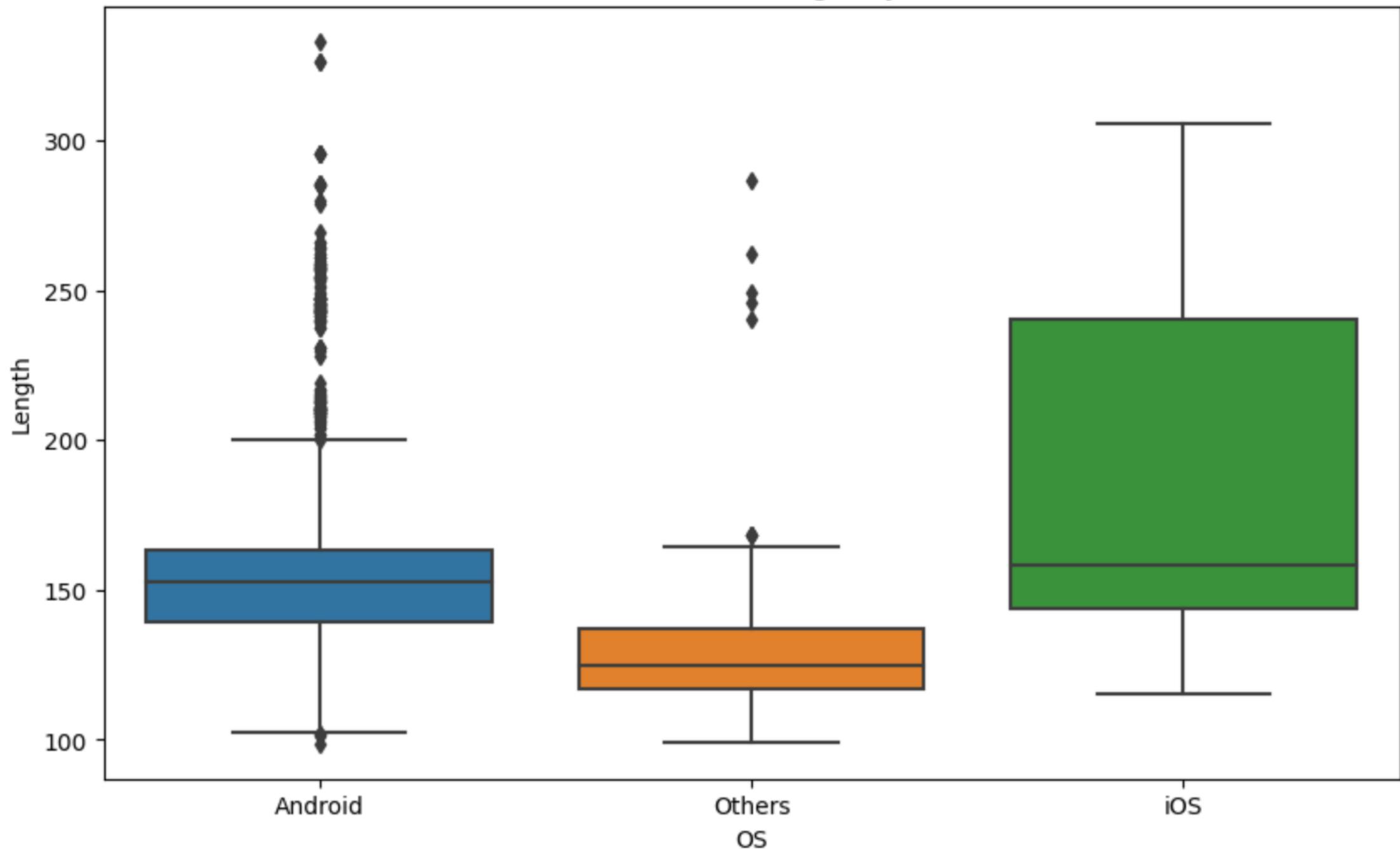
### Distribution of Weight by OS



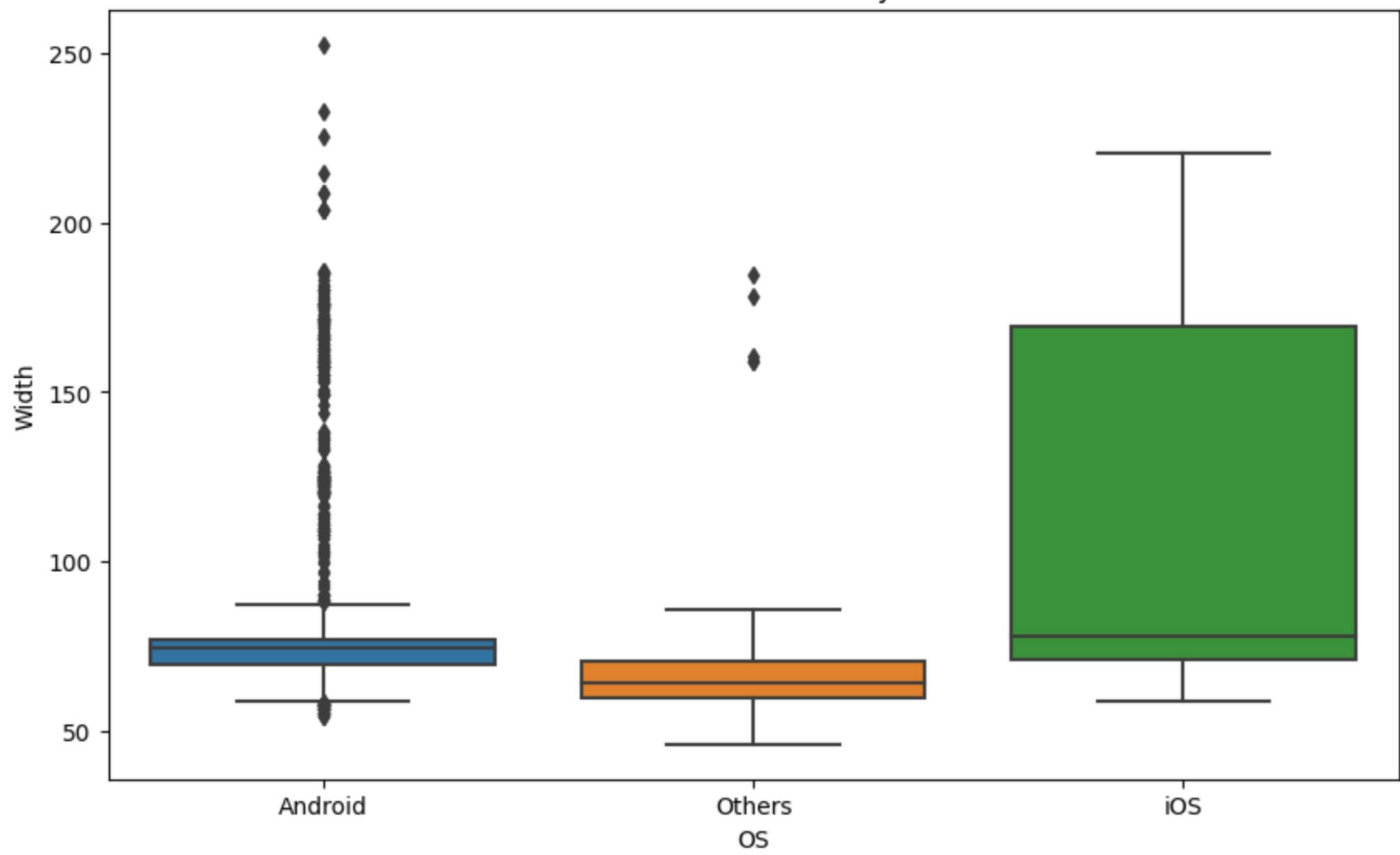
Ratio of Weight Cat by OS



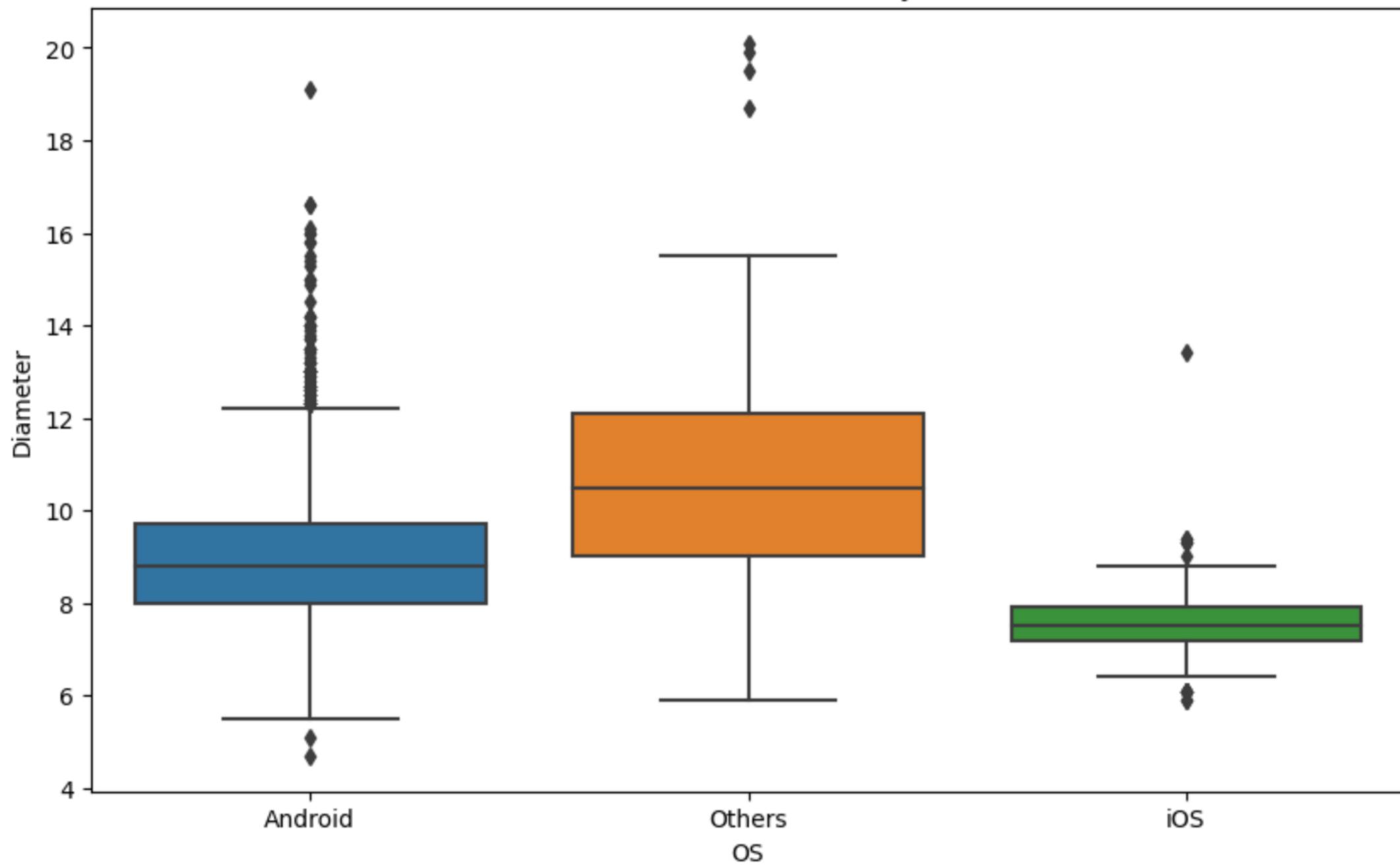
# Distribution of Length by OS



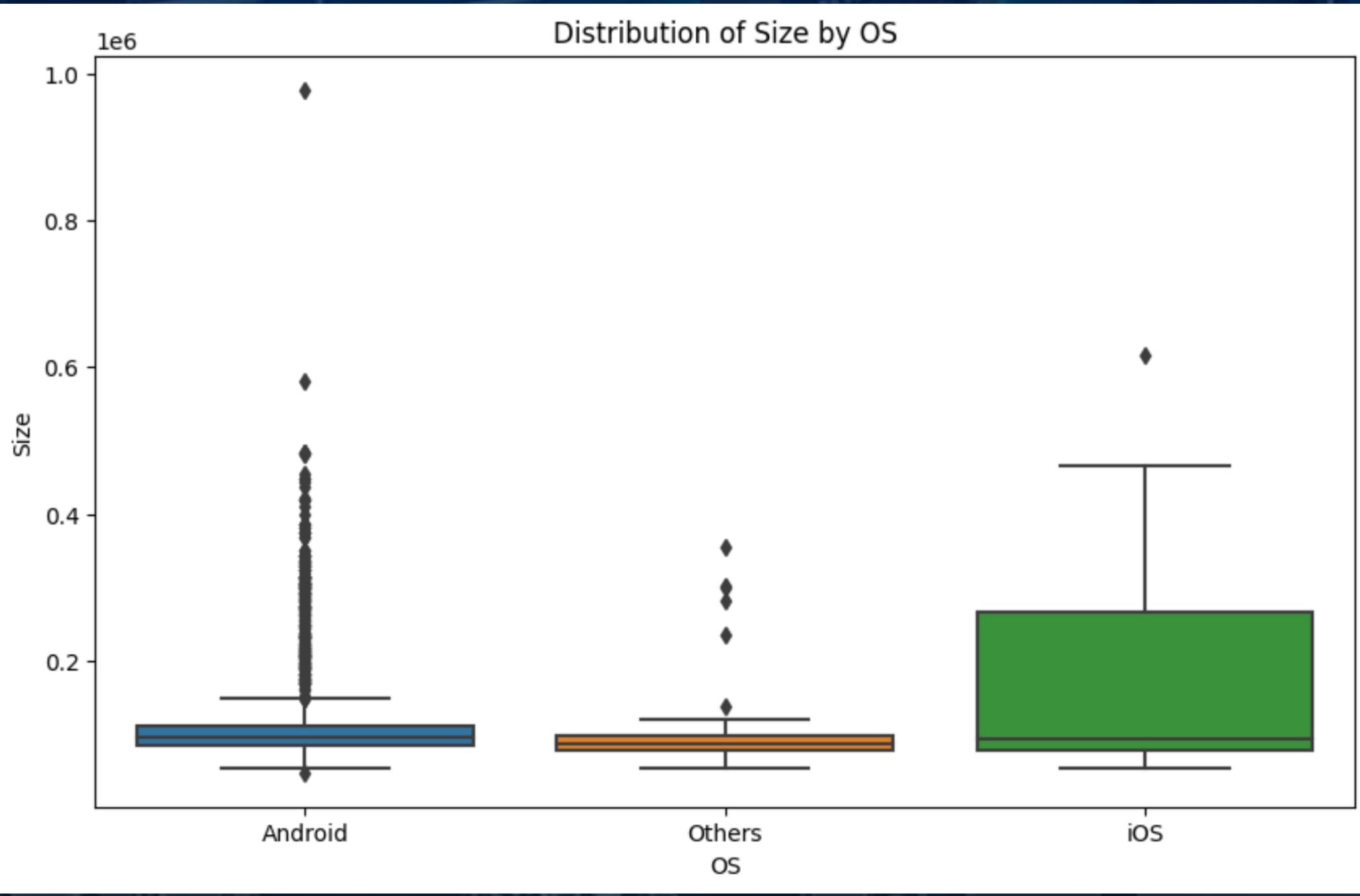
Distribution of Width by OS



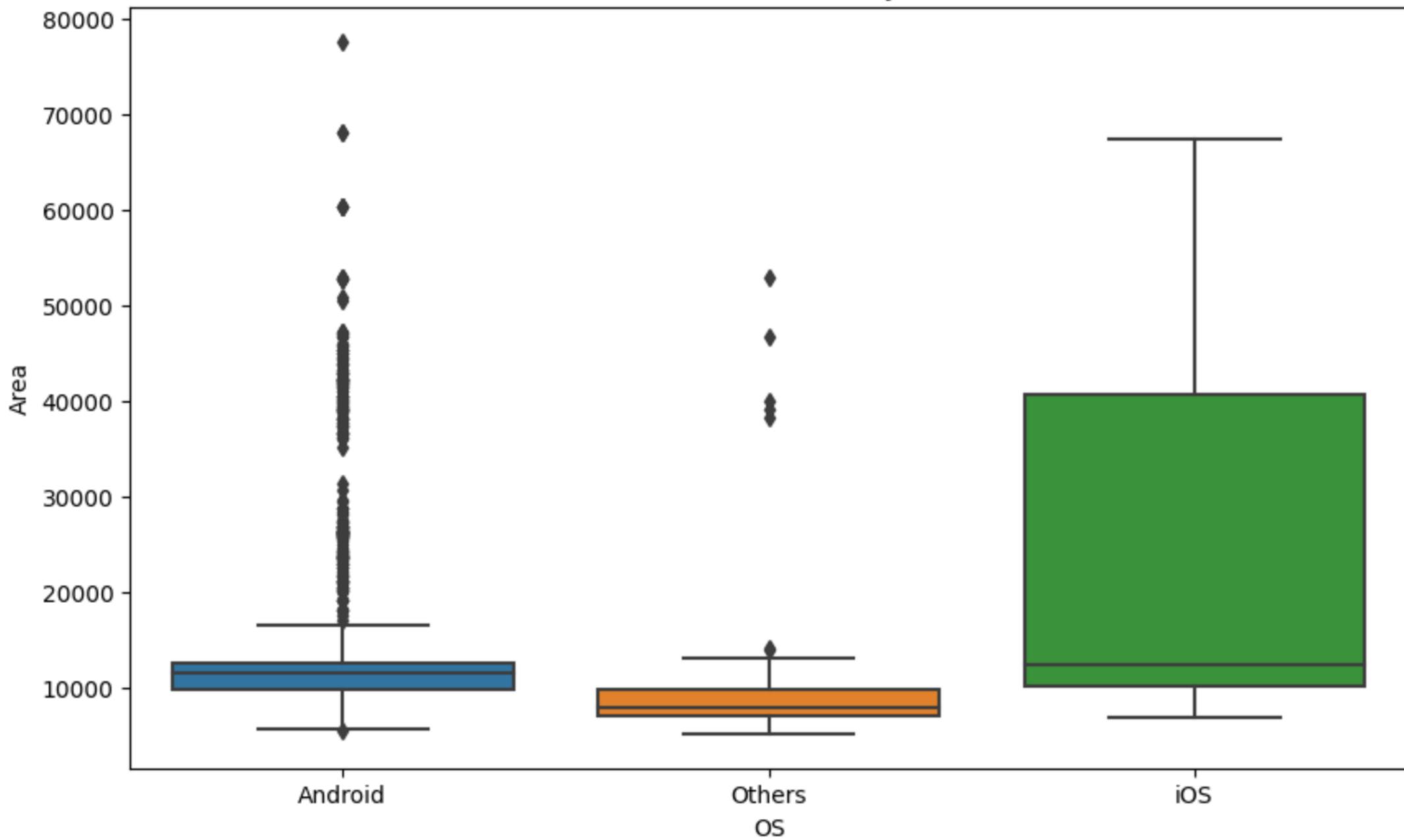
### Distribution of Diameter by OS



## Distribution of Size by OS



### Distribution of Area by OS



```
data = data.drop(columns=['Length', 'Width', 'Size', 'Area'])
```

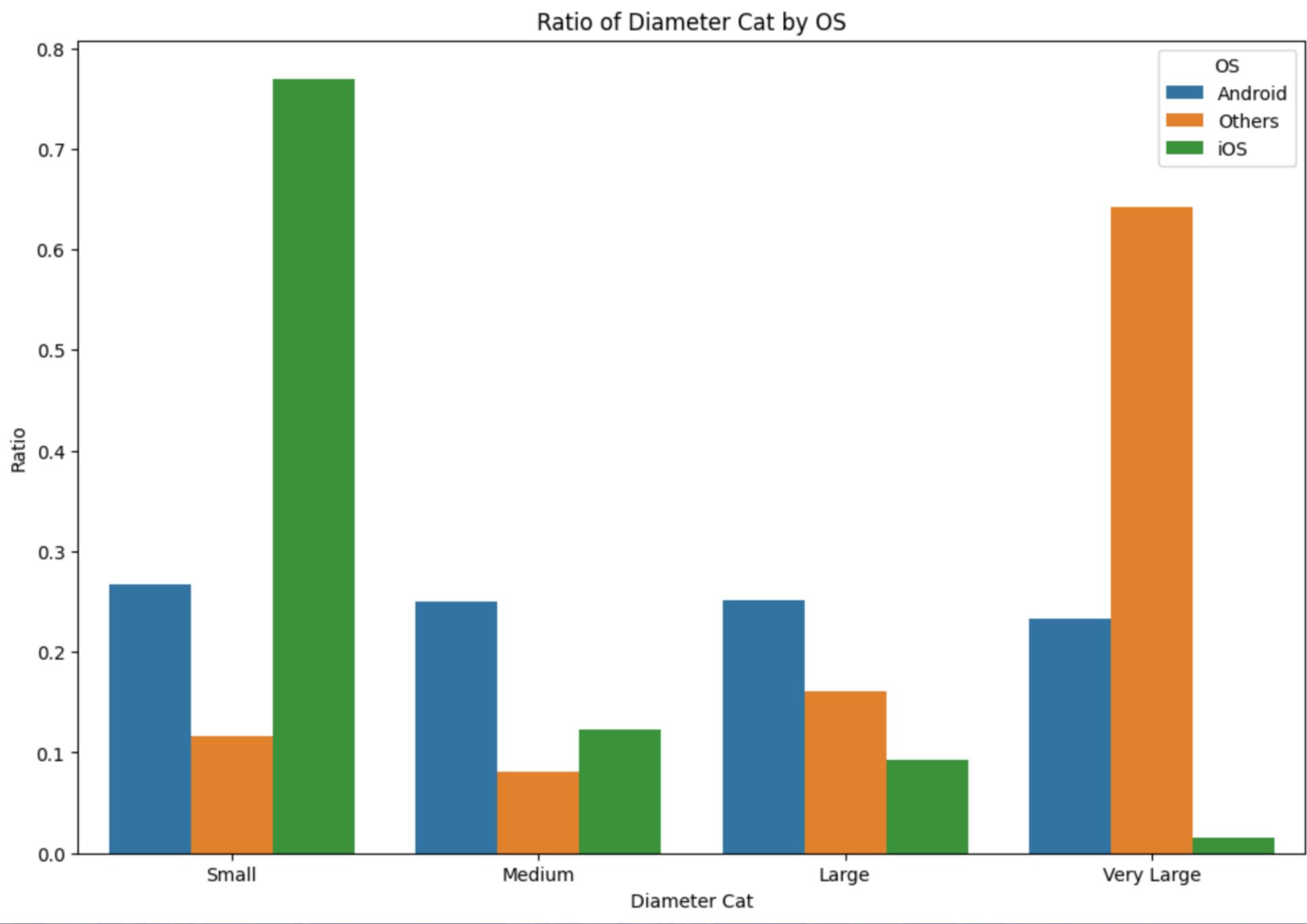
```
data['Diameter'].describe()
```

```
count    2646.000000
mean      9.101209
std       1.641620
min       4.700000
25%       8.000000
50%       8.800000
75%       9.800000
max      20.100000
Name: Diameter, dtype: float64
```

```
bins = [0, 8.0, 8.8, 9.8, 170]
```

```
labels = ['Small', 'Medium', 'Large', 'Very Large']
```

```
data['Diameter Cat'] = pd.cut(data['Diameter'], bins=bins, labels=labels)
```



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

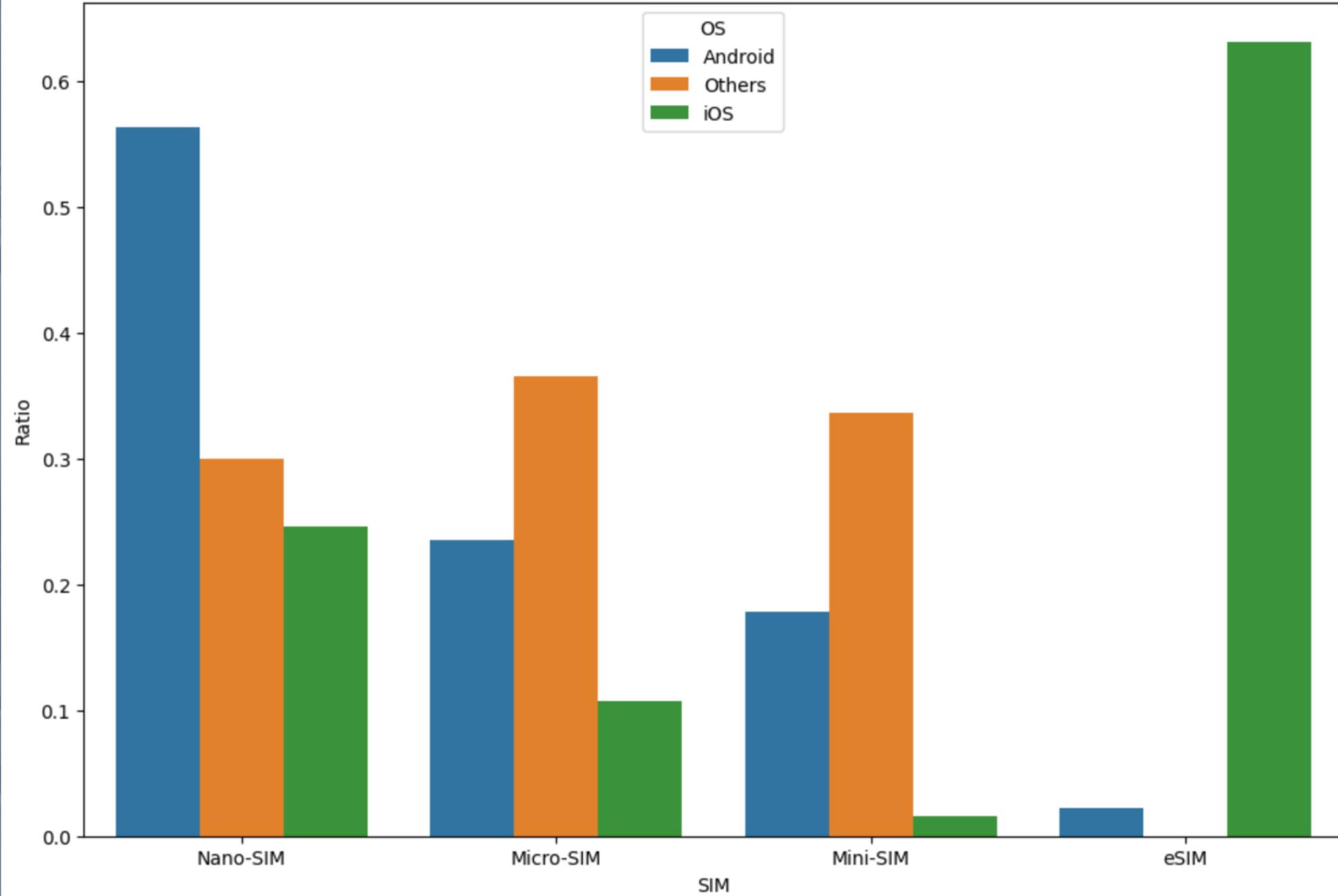
2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

Ratio of SIM by OS



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

Columns

2G 3G 4G 5G

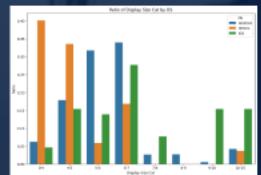
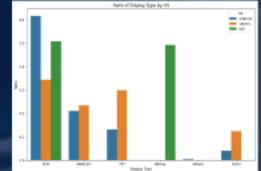
Tools

This is a magical place  
where we prepare our  
data for the next stage.

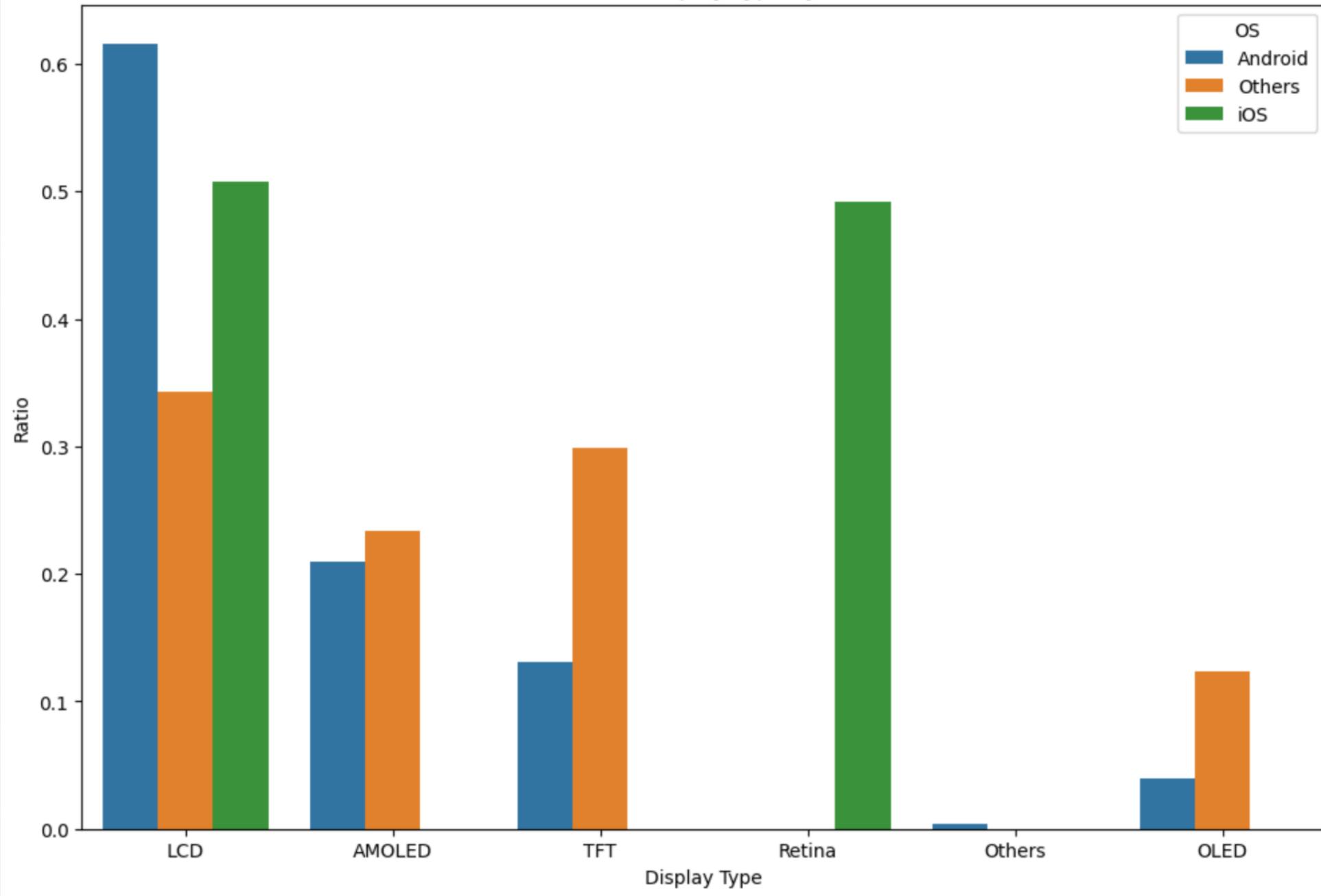
p.s: don't tell anyone

```
data['Display Type'].unique()
```

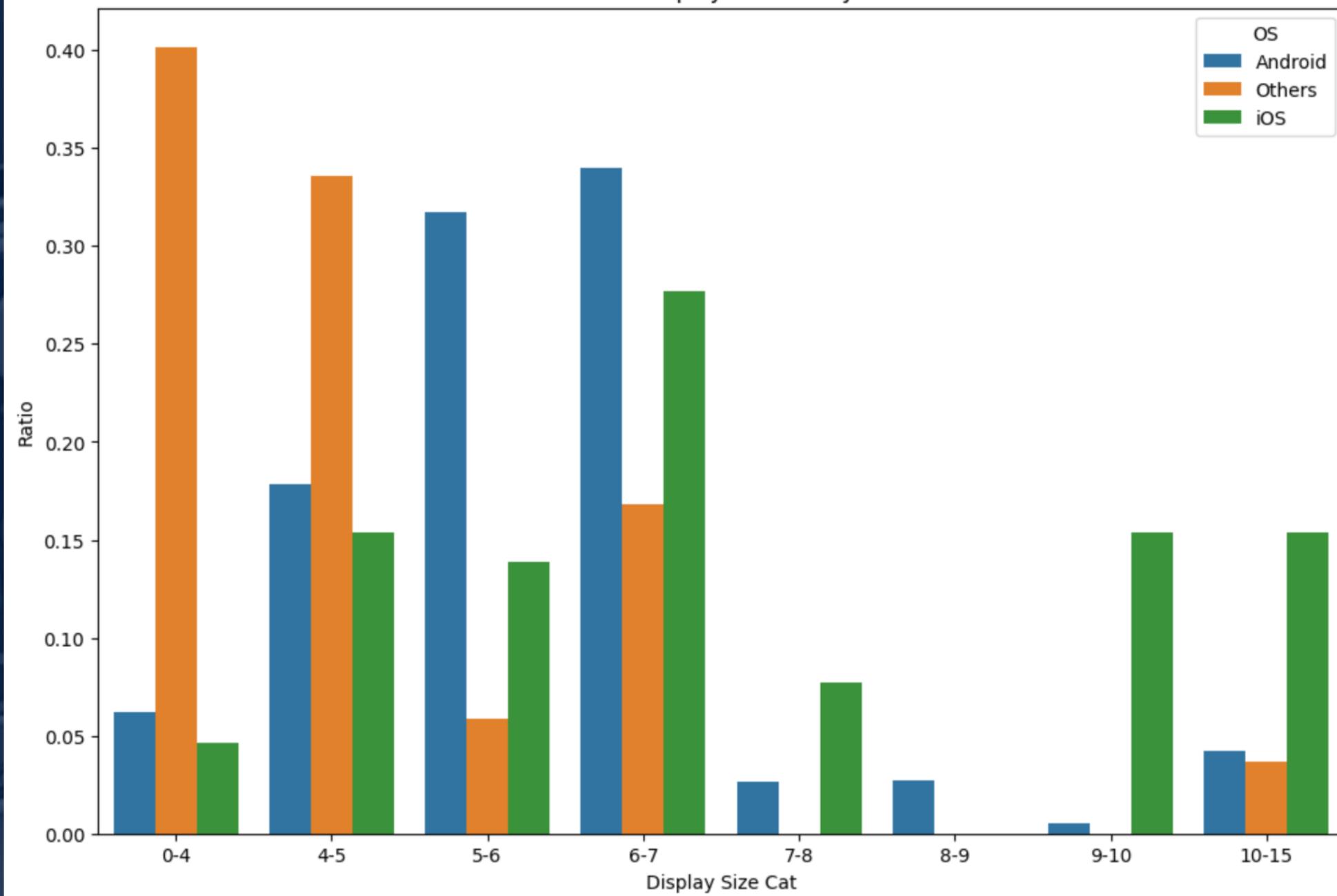
```
array(['TFT LCD', 'IPS LCD', 'TN TFT LCD', 'Super AMOLED', 'TFT',
       'AMOLED', 'LTPO Super Retina XDR OLED', 'Super Retina XDR OLED',
       'Liquid Retina XDR mini-LED LCD', 'Retina IPS LCD',
       'Liquid Retina IPS LCD', 'Super Retina OLED', 'LTPO AMOLED',
       'Super AMOLED ', 'S-IPS LCD', 'Super IPS+', 'IPS+ LCD',
       'Super IPS+ LCD', 'Super IPS LCD', 'LED-backlit IPS LCD', 'LCD',
       'Super IPS LCD2', 'TFT resistive touchscreen', 'OLED',
       'Super LCD6', 'Super LCD3', 'Super LCD', 'Super LCD5',
       'Super LCD2', 'S-LCD2', 'S-LCD', '3D LCD', '3D S-LCD',
       'AMOLED or SLCD', 'LTPO OLED', 'Foldable OLED', 'IPS-NEO LCD',
       'TDDI IPS LCD', 'TDDI', 'P-OLED', 'Curved P-OLED',
       'True HD-IPS + LCD', 'True IPS+ LCD', 'True Full HD IPS Plus LCD',
       'True HD IPS+', 'True HD-IPS LCD', 'AH IPS LCD',
       'True HD IPS Plus', 'HD-IPS LCD', '3D IPS LCD', 'AH-IPS LCD',
       'TFT LCD', 'LED-backlit LCD', 'Dual IPS LCD', 'TN',
       'IGZO IPS LCD', 'AMOLED or LCD', 'Samsung AMOLED', 'LTPO2 AMOLED',
       'Foldable AMOLED', 'Dynamic LTPO AMOLED 2X', 'Dynamic AMOLED 2X',
       'PLS LCD', 'Super AMOLED Plus', 'Dynamic AMOLED', 'PLS',
       'Super clear LCD', 'Super Flexible AMOLED', 'SC-LCD',
       'Super Clear LCD'], dtype=object)
```



### Ratio of Display Type by OS



Ratio of Display Size Cat by OS



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

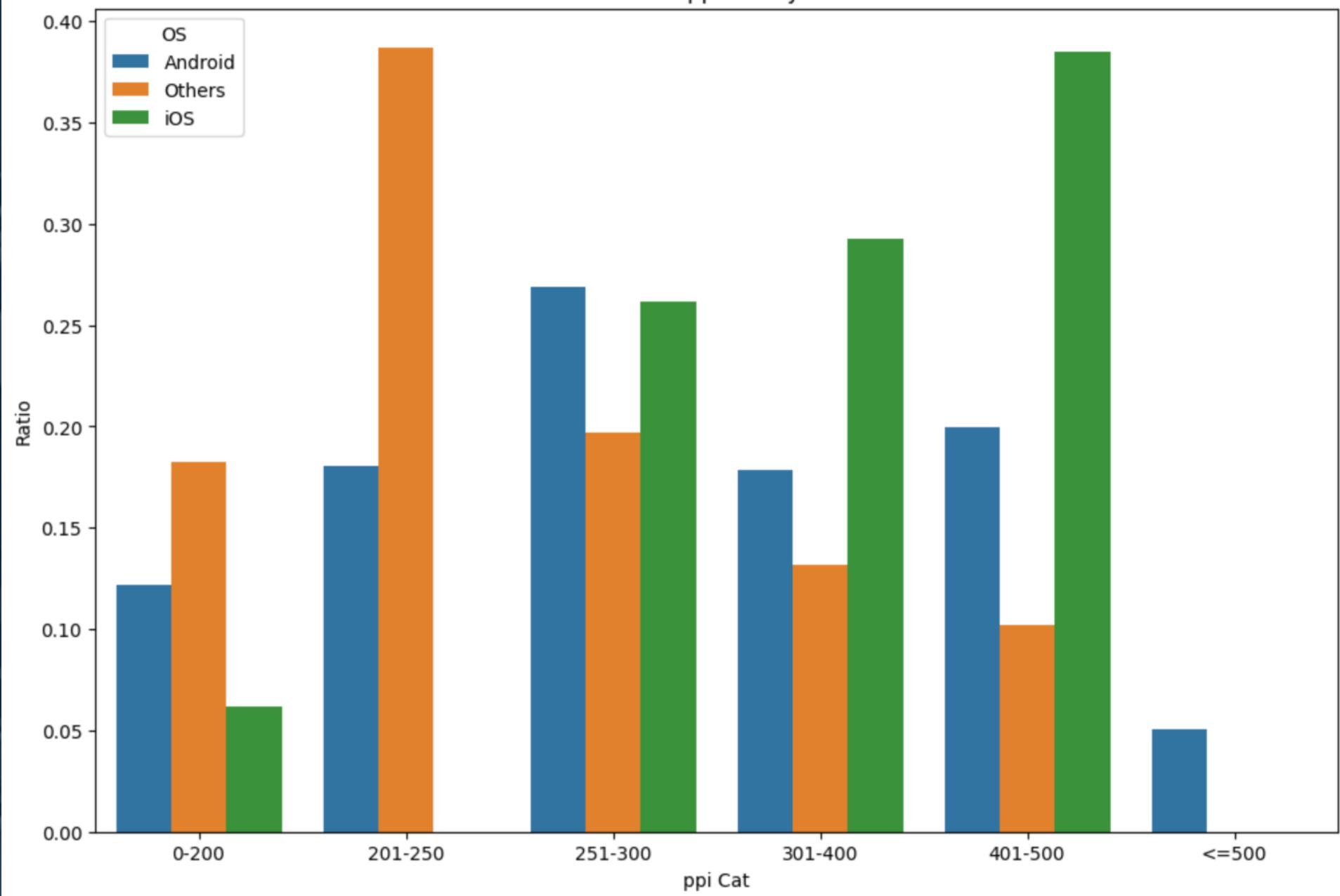
2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

Ratio of ppi Cat by OS



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

Columns

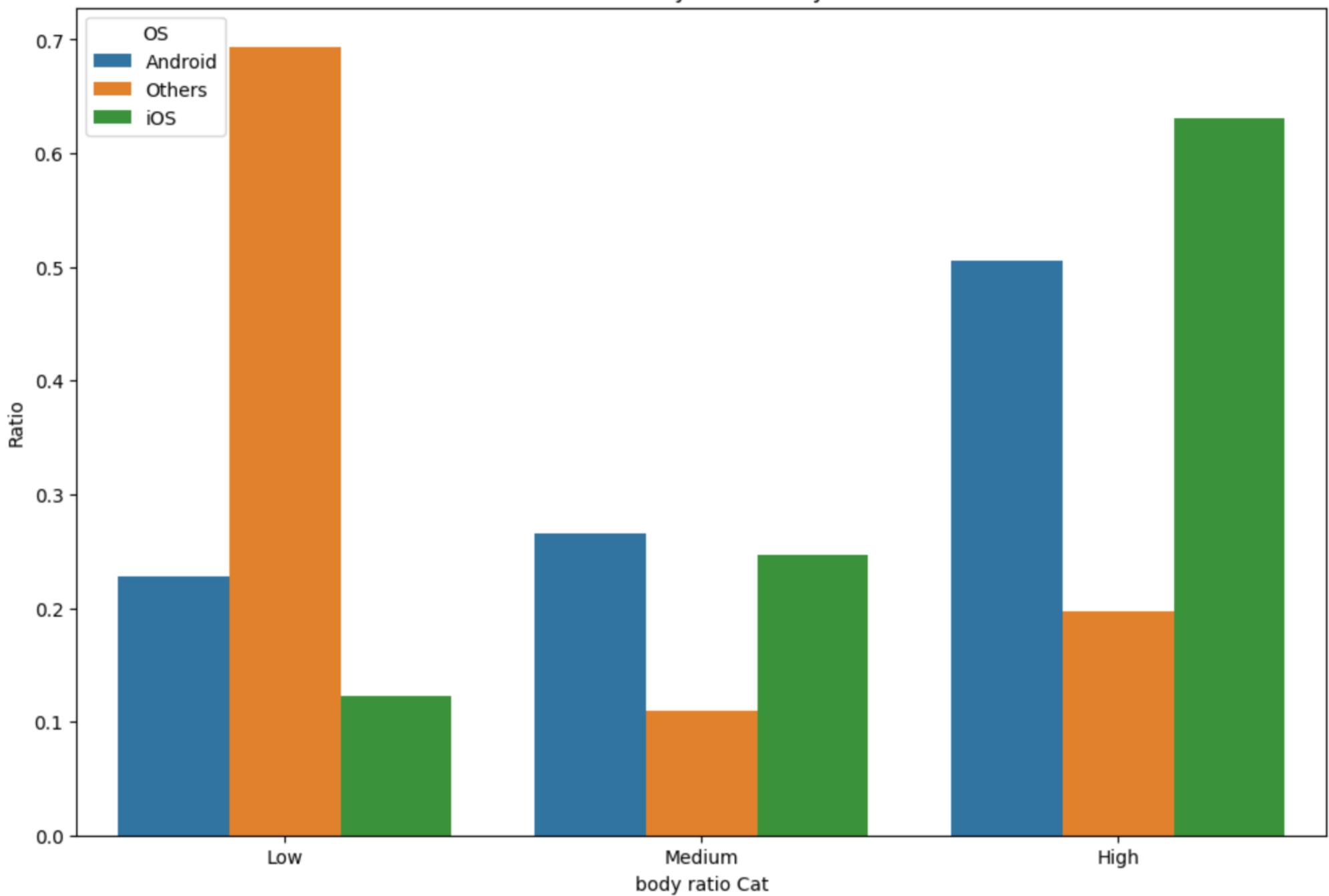
2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

Ratio of body ratio Cat by OS



Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

Weight

Columns

2G 3G 4G 5G

Tools

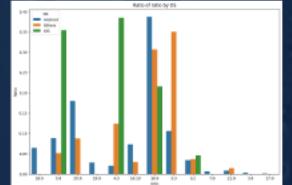
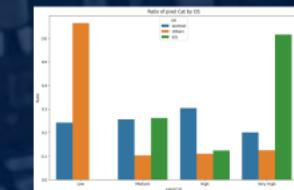
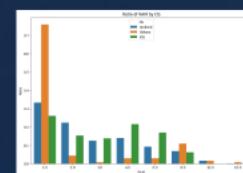
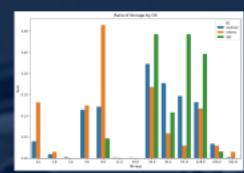
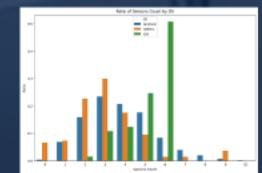
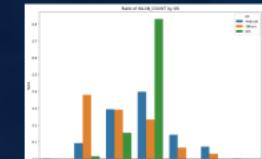
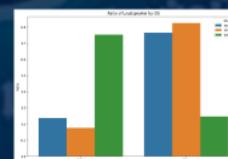
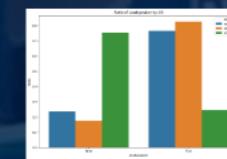
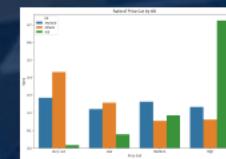
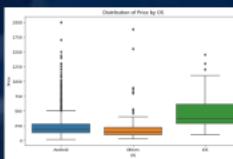
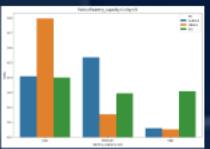
This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone

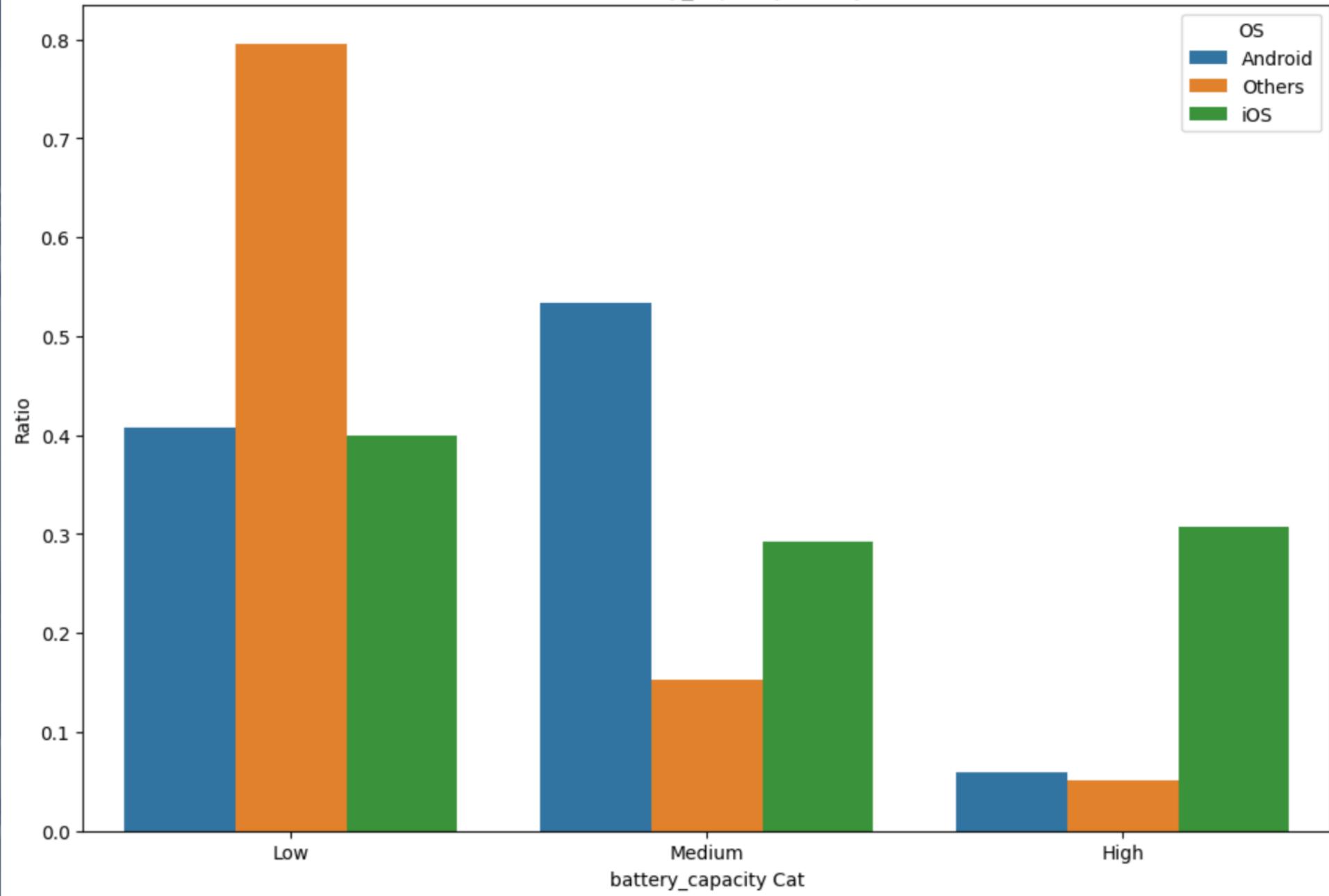
we also need to handle some nulls

| DELETED COLUMNS | MODIFIED COLUMNS | COLUMNS TO KEEP FOR 1M+ |
|-----------------|------------------|-------------------------|
| NAME            | NAME             | ANNOUNCED               |
| 2G              | 2G               | DISPLAY CAT             |
| 3G              |                  | DISPLAY SIZE CAT        |
| LENGTH          | DISPLAY CAT      | DISPLAY TYPE            |
| WIDTH           | DISPLAY SIZE CAT | DISPLAY TYPE CAT        |
| WLAN            | PPR CAT          | GPU CAT                 |
| COLORS          | BODY RATIO CAT   | BODY RATIO CAT          |
| SENSORS         | PRICE CAT        | PRICE CAT               |
| BLUETOOTH       | PPR CAT          | PPR CAT                 |
| LOUDPRAKER      | PPR CAT          | PPR CAT                 |
| 3.5MM JACK      | PPR CAT          | PPR CAT                 |
| INTERNAL        | PPR CAT          | PPR CAT                 |
| CARD SLOT       | PPR CAT          | PPR CAT                 |
| NETWORK         | PPR CAT          | PPR CAT                 |
|                 | PPR CAT          | PPR CAT                 |
| TOTAL COUNT: 14 | TOTAL COUNT: 6   | TOTAL COUNT: 6          |
|                 | TOTAL COUNT: 1   | TOTAL COUNT: 1          |

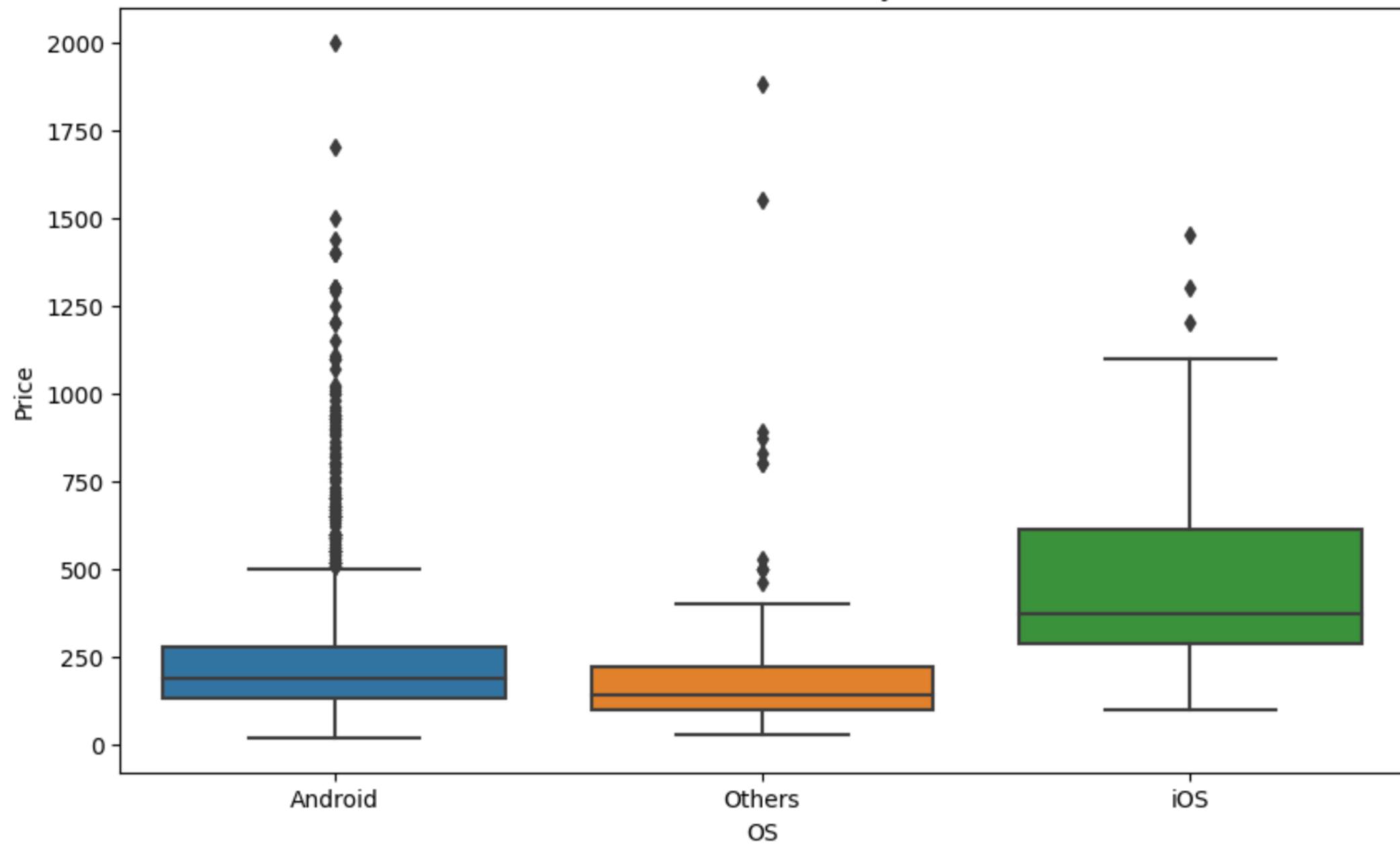
# Let's explore the rest quickly



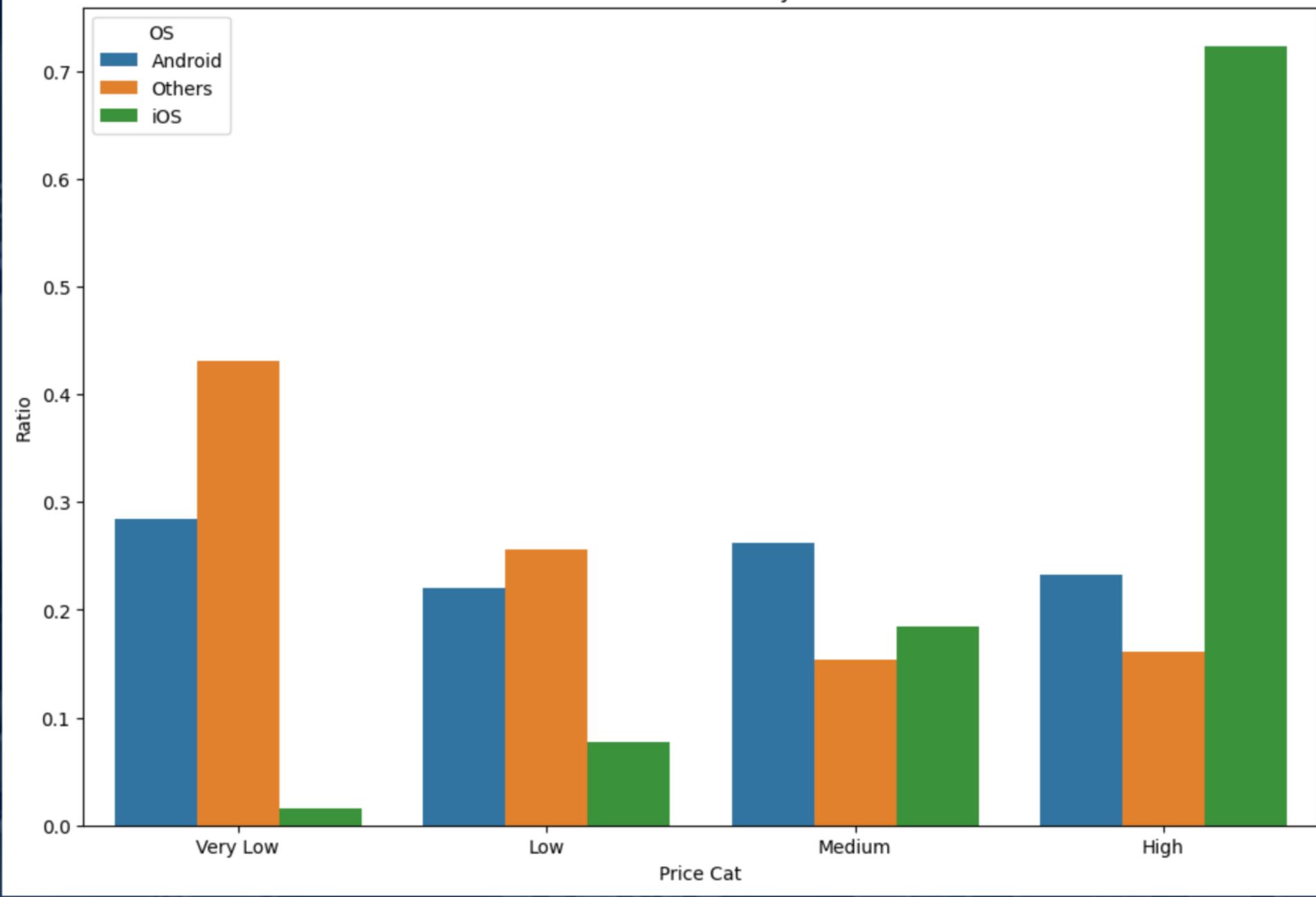
Ratio of battery\_capacity Cat by OS



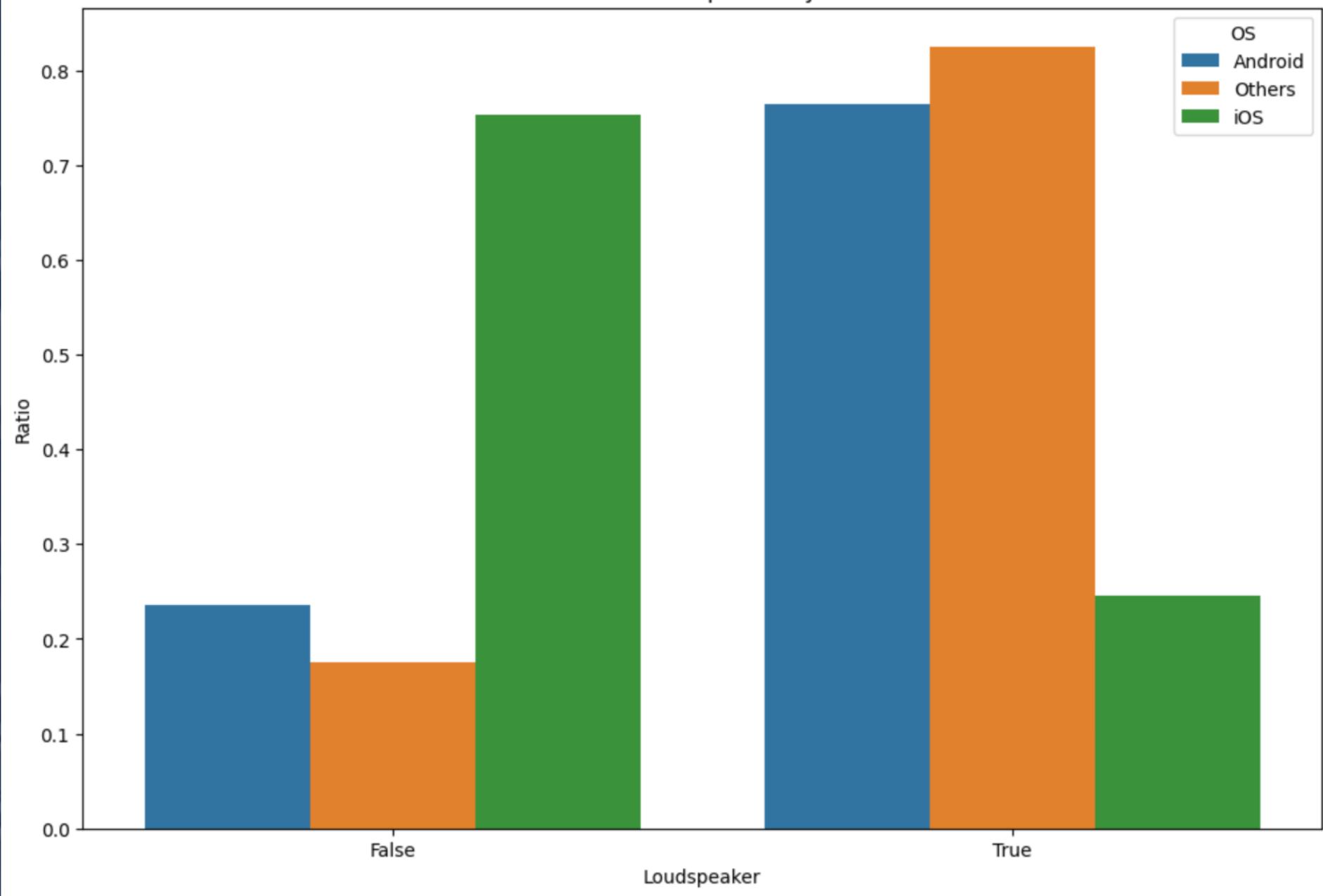
### Distribution of Price by OS



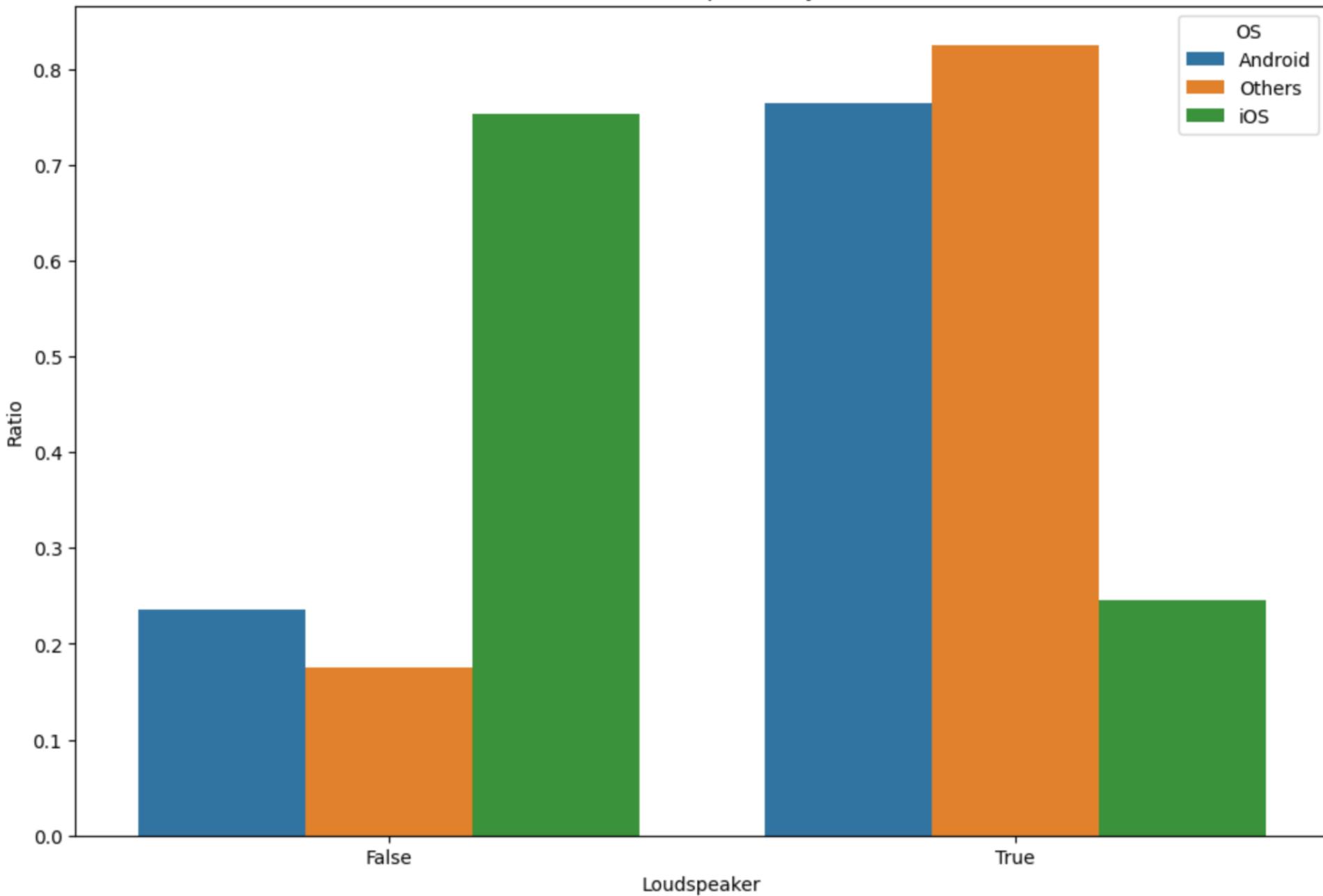
Ratio of Price Cat by OS



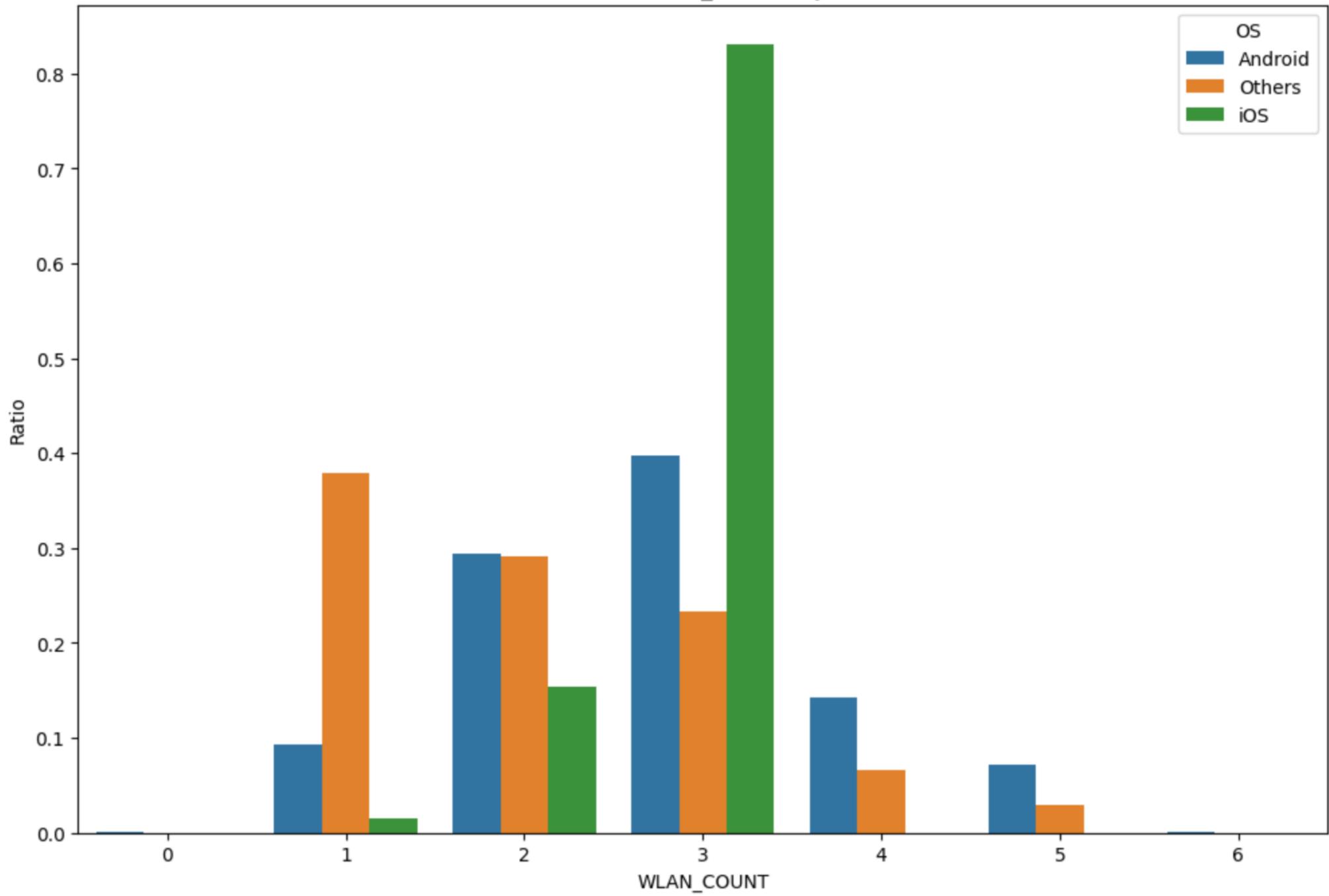
Ratio of Loudspeaker by OS



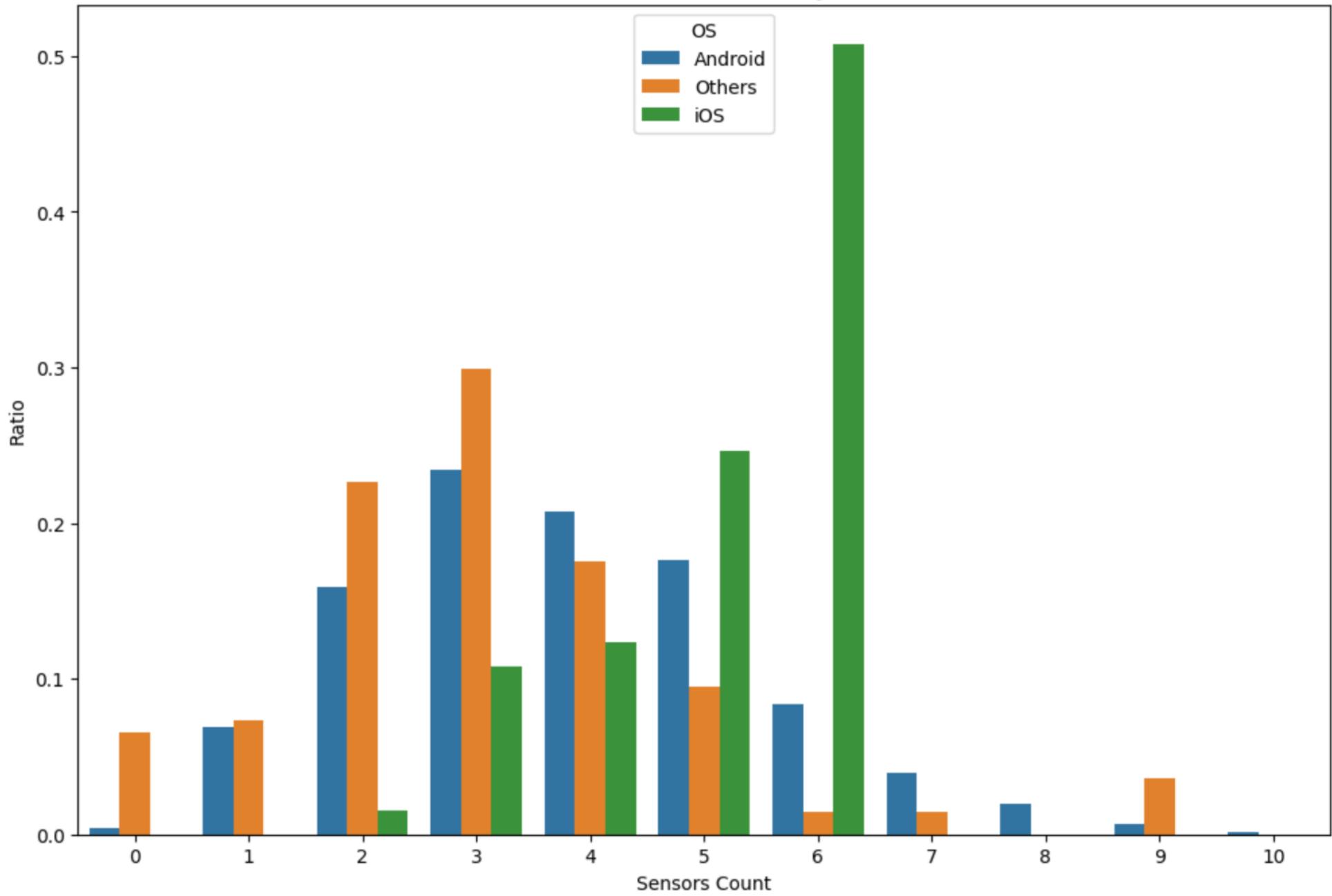
Ratio of Loudspeaker by OS



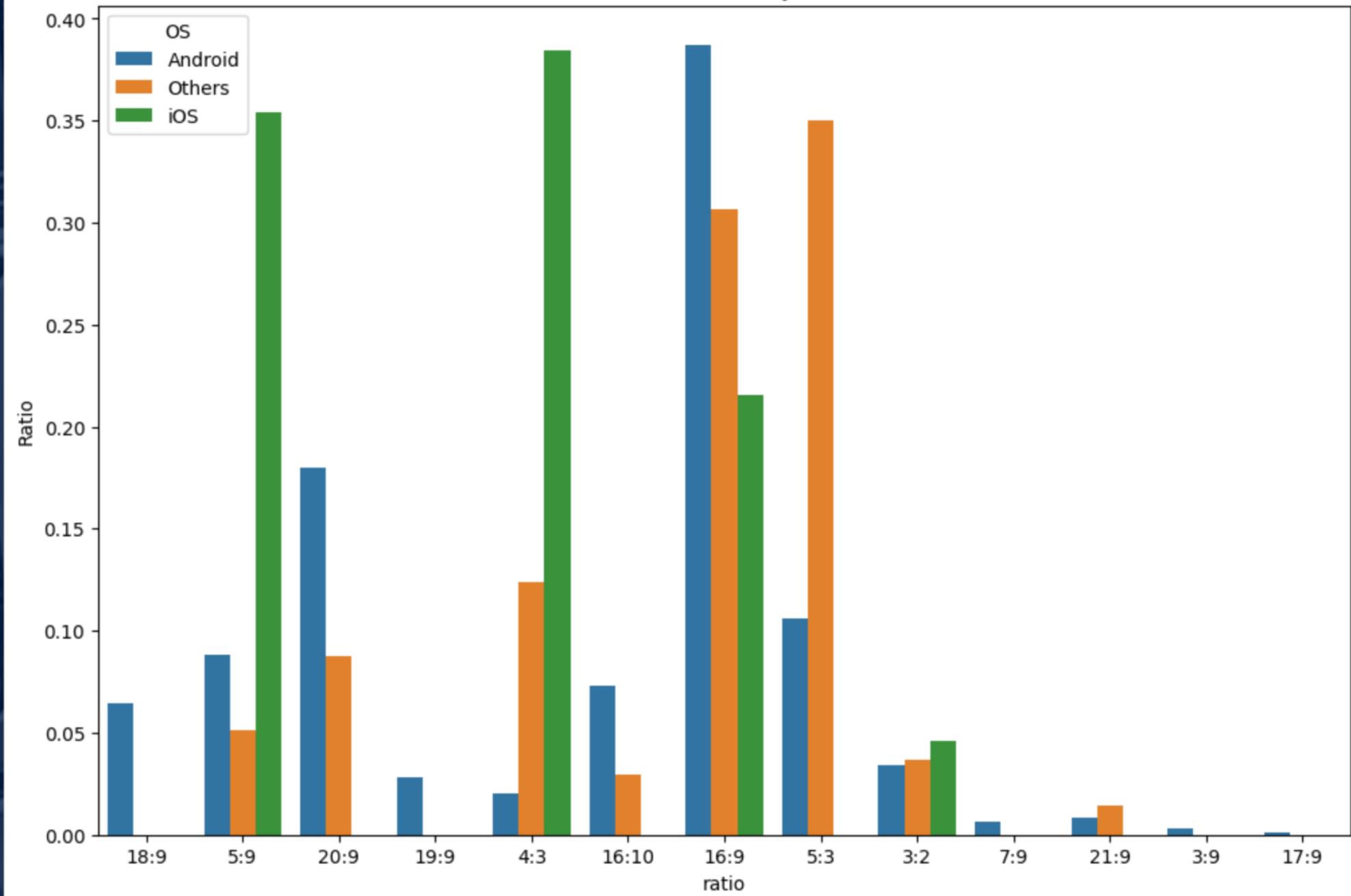
Ratio of WLAN\_COUNT by OS



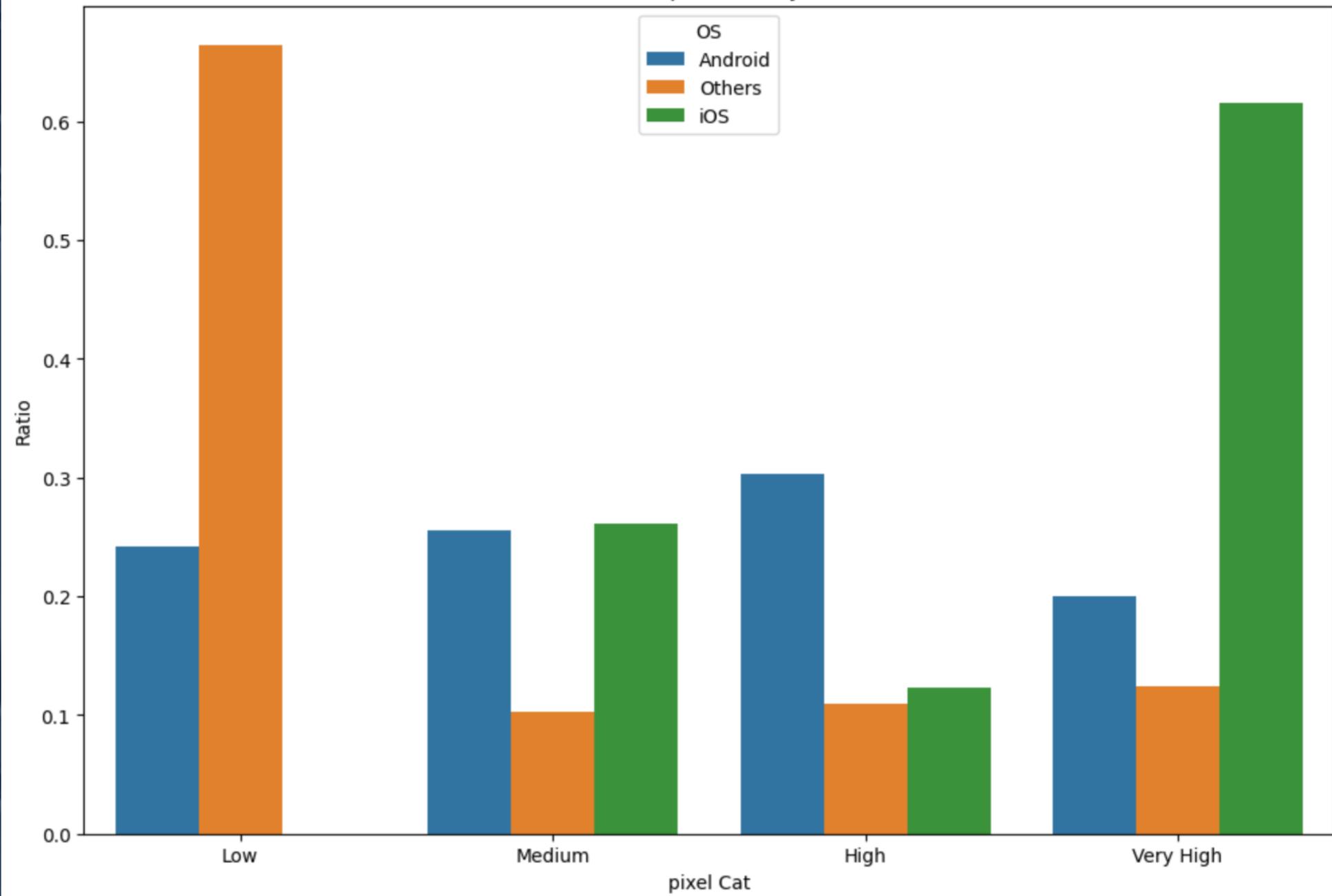
### Ratio of Sensors Count by OS



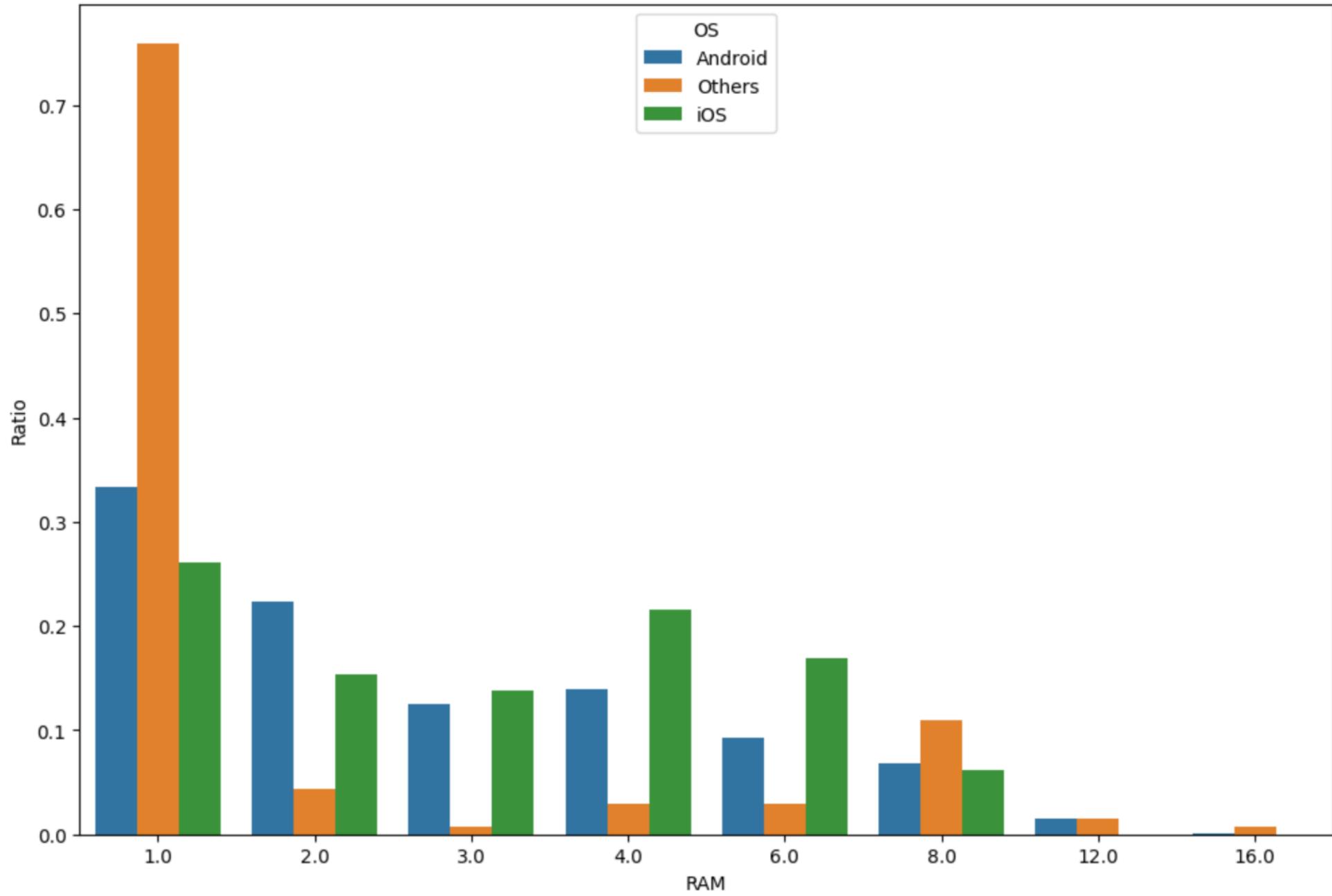
Ratio of ratio by OS



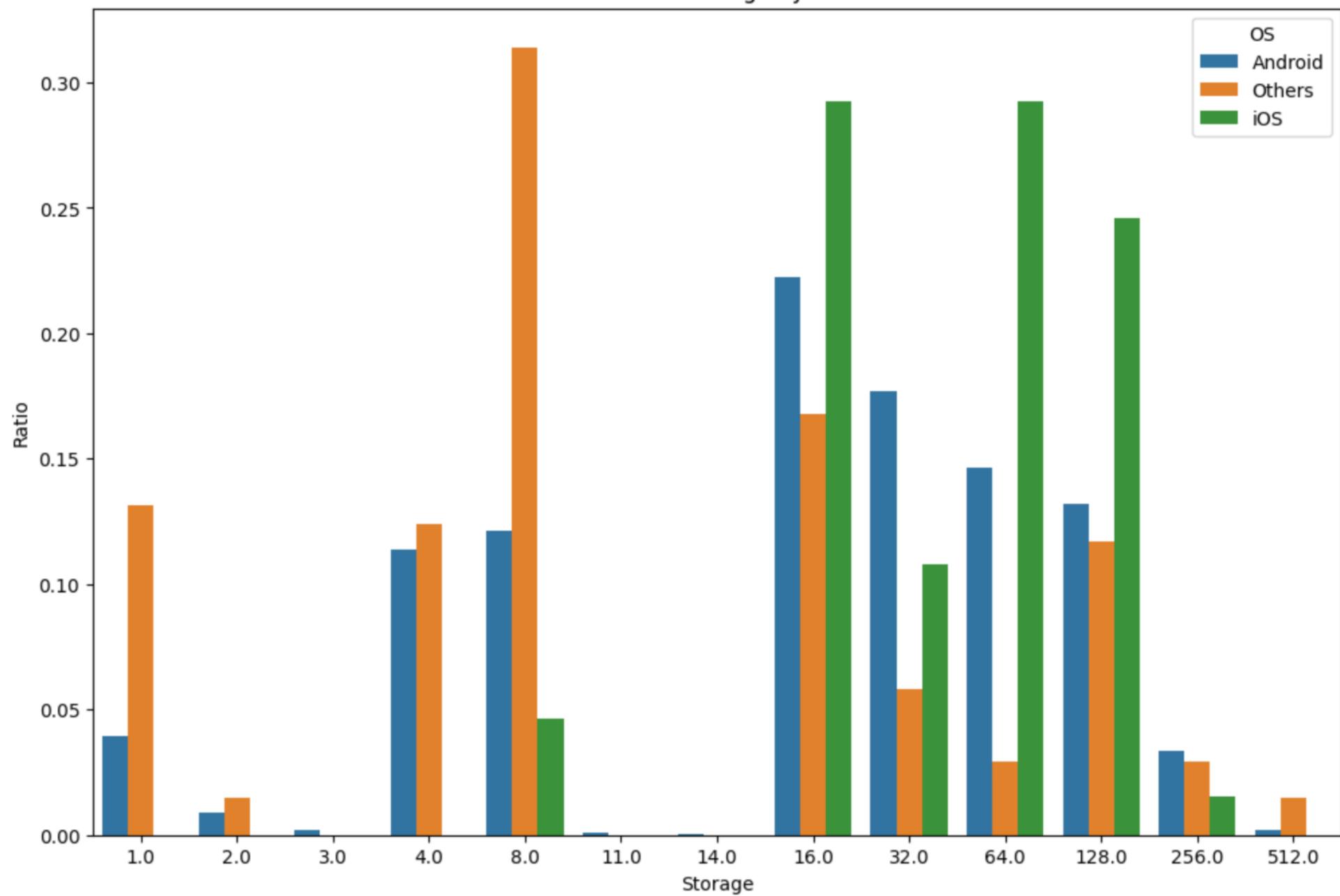
Ratio of pixel Cat by OS



Ratio of RAM by OS



### Ratio of Storage by OS



| DELETED COLUMNS |
|-----------------|
| NAME            |
| 2G              |
| 3G              |
| LENGTH          |
| WIDTH           |
| WLAN            |
| COLORS          |
| SENSORS         |
| BLUETOOTH       |
| LOUDSPEAKER     |
| 3.5MM JACK      |
| INTERNAL        |
| CARDSLOT        |
| NETWORK         |

TOTAL COUNT: 14

| MODIFIED COLUMNS     |
|----------------------|
| WEIGHT CAT           |
| DIAMETER CAT         |
| DISPLAY TYPE         |
| DISPLAY SIZE CAT     |
| PPI CAT              |
| BODY RATIO CAT       |
| BATTERY_CAPACITY CAT |
| PRICE CAT            |
| WLAN_COUNT           |
| PIXEL CAT            |

TOTAL COUNT: 10

| COLUMNS TO KEEP FOR 100% |
|--------------------------|
| BRAND                    |
| 5G                       |
| SIM                      |
| DISPLAY TYPE             |
| BATTERY CAPACITY         |
| CPU                      |

TOTAL COUNT: 6

| COLUMNS TO KEEP FOR LESS THAN 50% |
|-----------------------------------|
| STATUS                            |
| GPU ?                             |
| CHIPSET ?                         |

TOTAL COUNT: 3

| COLUMNS TO KEEP FOR 50% |
|-------------------------|
| ANNOUNCED               |
| WEIGHT CAT              |
| DIAMETER CAT            |
| DISPLAY SIZE CAT        |
| PPI CAT                 |
| BODY RATIO CAT          |
| PRICE CAT               |
| WLAN_COUNT              |
| RATIO                   |
| PIXEL CAT               |
| STORAGE                 |
| RAM                     |

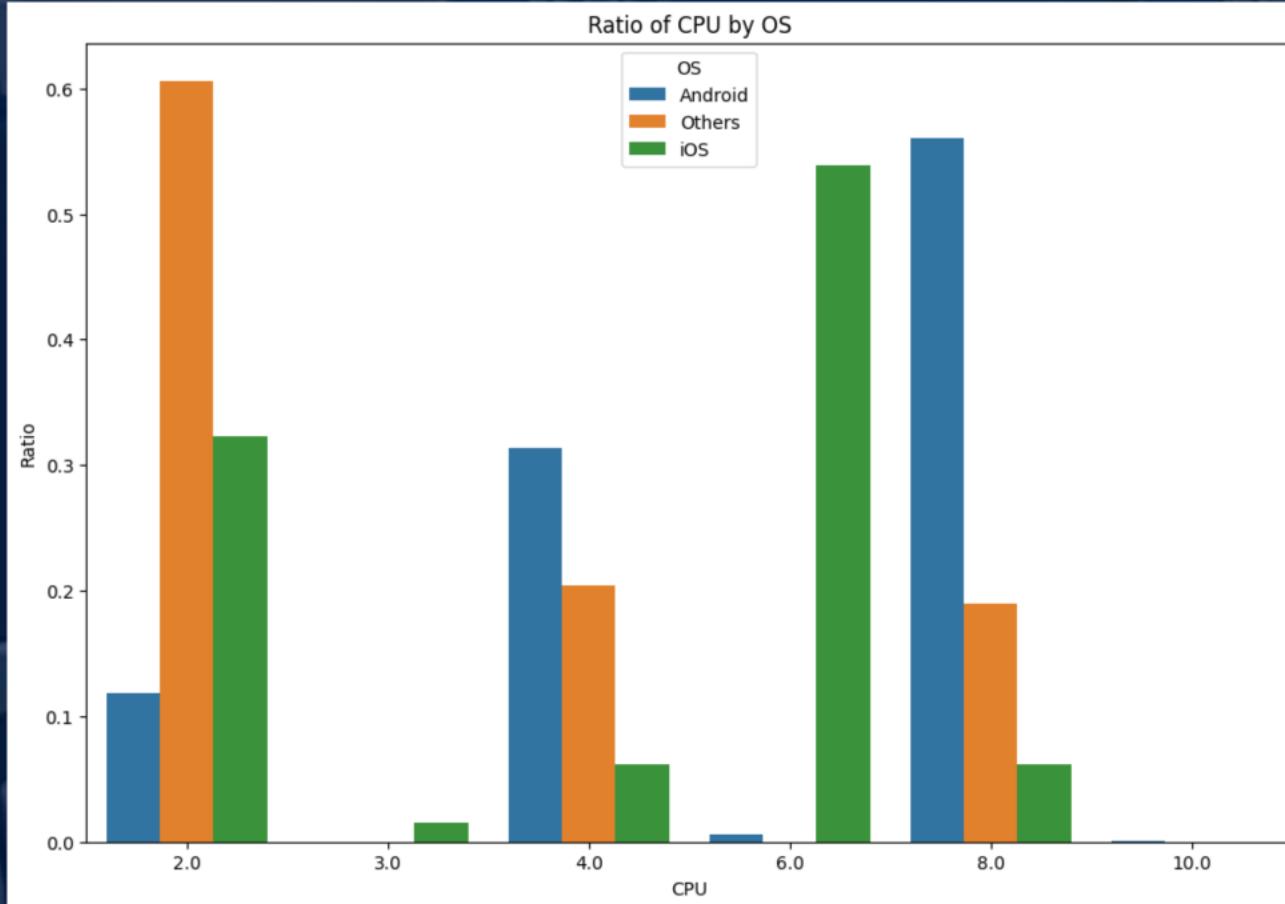
TOTAL COUNT: 12

**we also need  
to handle  
some nulls**

```
data.isnull().sum()
```

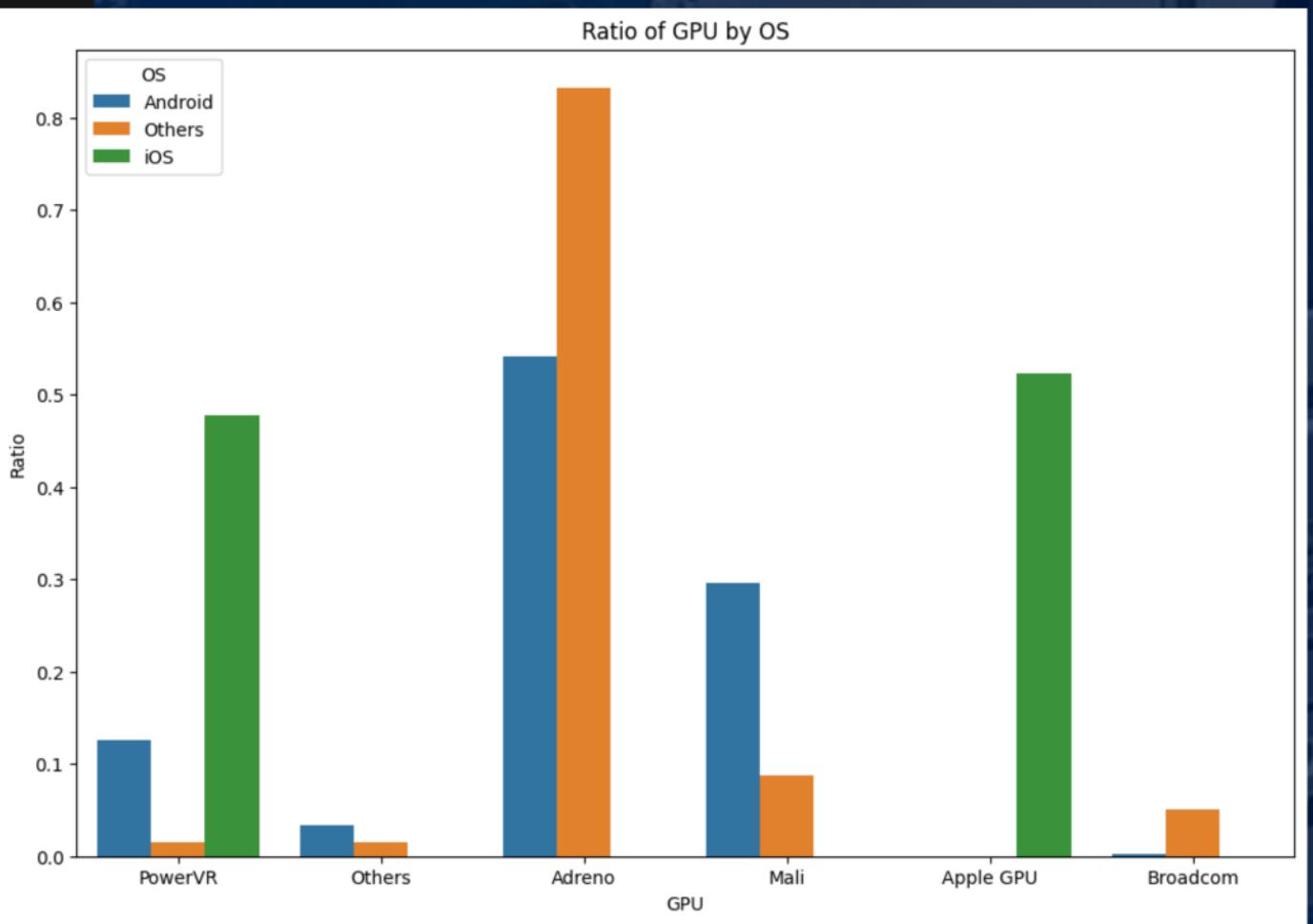
```
Diameter          0
SIM              0
Display Type     0
Display Size     0
ppi               0
body ratio        0
OS                0
battery_capacity 0
Price             0
CPU              279
ratio             0
pixel             0
Bluetooth         0
GPU               224
Chipset           257
Network           0
Internal          0
Card slot         0
RAM               0
Storage           0
Ratio              0
Weight Cat        0
Diameter Cat      0
Display Size cat  0
ppi Cat            0
body ratio Cat    0
battery_capacity Cat 0
Price Cat          0
WLAN_COUNT        0
dtype: int64
```

```
cpu_counts = data.groupby(['os', 'brand', 'CPU']).size().reset_index(name='COUNT')
max_cpu = cpu_counts.loc[cpu_counts.groupby(['os', 'brand'])['COUNT'].idxmax()]
max_cpu.drop(columns='COUNT', inplace=True)
max_cpu
```

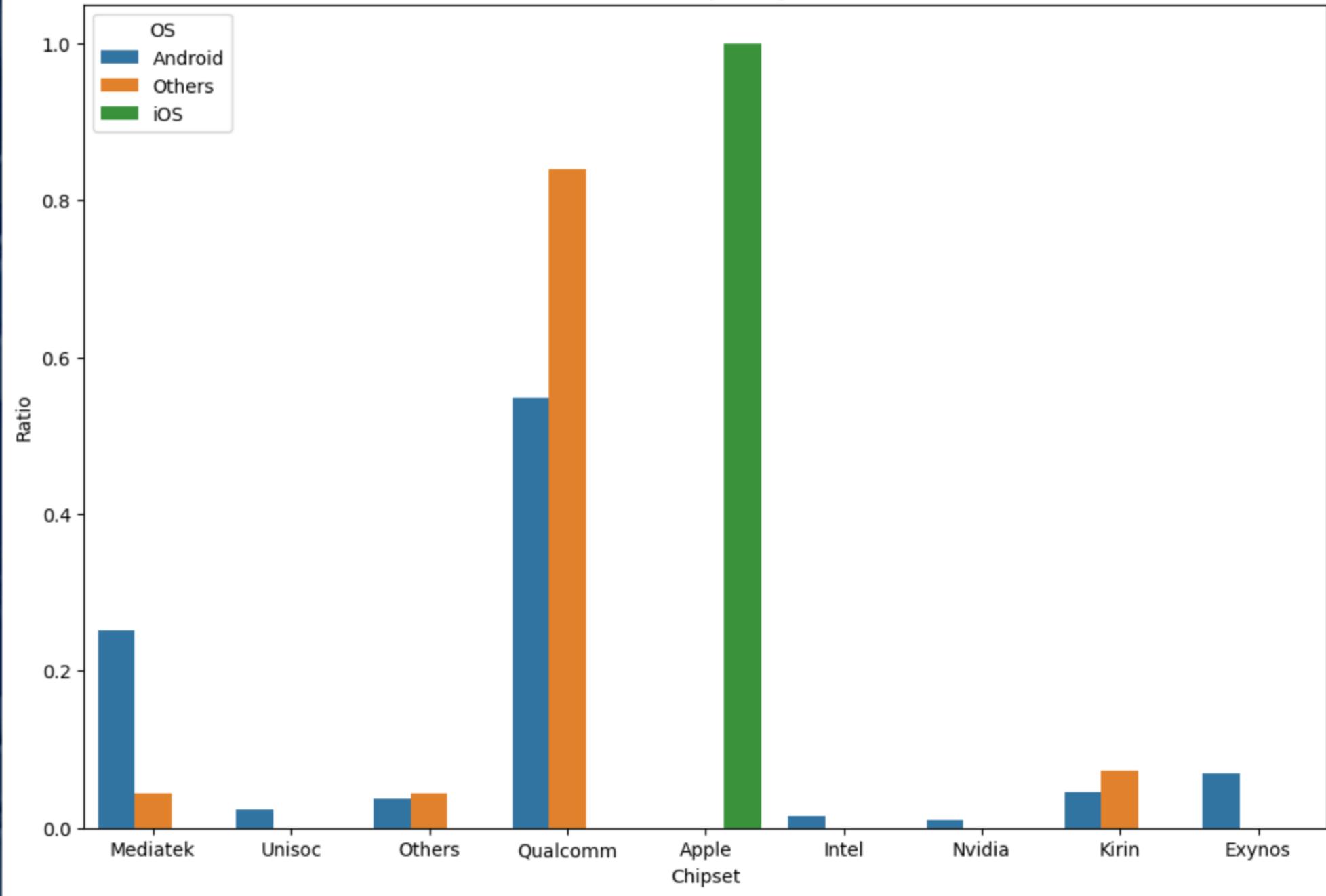


```
def preprocess_gpu(gpu):
    if '-' in gpu:
        return gpu.split('-')[0].strip()
    else:
        return gpu.strip()

data['GPU'] = data['GPU'].apply(preprocess_gpu)
data['GPU'].nunique()
```



### Ratio of Chipset by OS



## Handling Categorical Variables:

```
label_encoder = LabelEncoder()

columns_to_encode = ['brand', 'ratio', 'os', 'GPU', 'Chipset', 'Weight Cat', 'Diameter Cat', 'Display Size Cat',
| | | | | | | 'ppi Cat', 'body ratio Cat', 'battery_capacity Cat', 'Price Cat', 'pixel Cat', 'Status']

for column in columns_to_encode:
    data[column] = label_encoder.fit_transform(data[column])

columns_to_dummies = ['SIM', 'Display Type']
data = pd.get_dummies(data, columns=columns_to_dummies)
```

Let's begin!

And more

Body Ratio

ppi

Display

SIM

Length, Width,  
Diameter

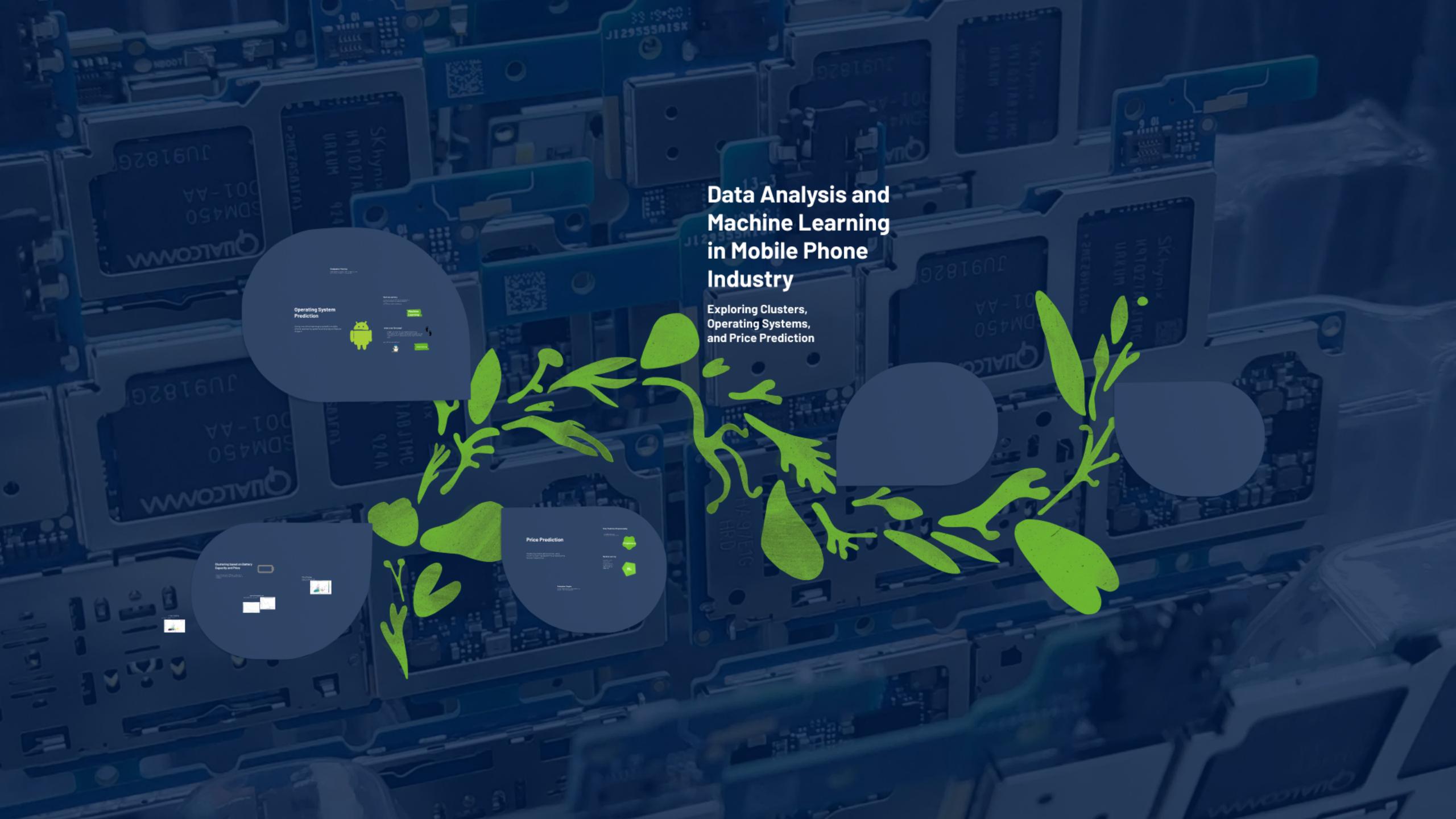
Weight

2G 3G 4G 5G

Tools

This is a magical place  
where we prepare our  
data for the next stage.

p.s: don't tell anyone



# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction



# Machine Learning

OK, now that our work with preprocessing is done,  
we can try to predict the operating system of  
phones  
using machine learning techniques.



Machine  
Learning

first, we needed to split our data

```
X = df.drop(['OS', 'size'], axis=1)
Y = df['OS']

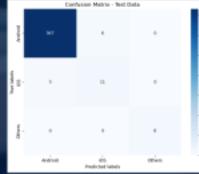
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.15, random_state=42)

X_tr, X_val, y_tr, y_val = train_test_split(X_train, Y_train, test_size=0.15, random_state=42)
```

```
Metrics for Class 0:
Precision: 0.96891091762805
Recall: 0.9691800109350899
F1 Score: 0.9691800109350899

Metrics for Class 1:
Precision: 0.6666666666666667
Recall: 0.4875
F1 Score: 0.6666666666666667

Metrics for Class 2:
Precision: 1.0
Recall: 1.0
F1 Score: 1.0
```



our model and scale didn't change

here's the results:

```
Validation Metrics:
Accuracy: 0.9615384615384616
Precision: 0.96891091762805
Recall: 0.9615384615384616
F1 Score: 0.9615384615384616
AUC Score: 0.9685459684581547
```

WOW THERE IS A PART 2, JUST WHAT WE NEEDED, MORE FUN!

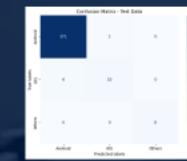
Ahright, so We just needed to change the 'battery\_capacity' in Selected\_features to 'battery\_capacity Cat'

we heard you asked for each class metrics, so here we go:

```
class 0: android
class 1: ios
class 2: others
```

```
Metrics for Class 0:
Precision: 0.968908908906948
Recall: 0.94638608904938
F1 Score: 0.9899333333333334

Metrics for Class 1:
Precision: 0.8333333333333334
Recall: 0.625
F1 Score: 0.714857142857143
```



then we balanced our data

#### Handling Imbalanced Data

```
# we can use imbalanced-learn library to handle this
from imblearn.over_sampling import SMOTE
smote = SMOTE()
X_tr_balanced, y_tr_balanced = smote.fit_resample(X_tr, y_tr)

print(X_tr_balanced.shape, y_tr_balanced.shape, np.unique(y_tr_balanced, return_counts=True))
print(y_tr_balanced)

# print(X_tr_balanced, y_tr_balanced, np.unique(y_tr_balanced, return_counts=True))

(1900, 30) (1900,) array([0, 1, 2], dtype=int64, order='fortran')
array([1000, 900, 1000], dtype=int64, order='fortran', max(1000, 900, 1000), shape=(3,))

Balanced
```

and then we tried to find the best combination of model and scaling using Gridsearch

```
classifiers = [
    ('Logistic Regression', LogisticRegression()),
    ('Random Forest', RandomForestClassifier()),
    ('Adaboost', AdaBoostClassifier()),
    ('SVM', SVC())
]

param_grids = {
    'Logistic Regression': {
        'C': [1, 10, 100, 1000]
    },
    'Random Forest': {
        'n_estimators': [50, 100, 200],
        'max_depth': [10, 20, 30]
    },
    'Adaboost': {
        'n_estimators': [50, 100, 200],
        'learning_rate': [0.05, 0.1, 0.2]
    },
    'SVM': {
        'C': [1, 10, 100, 1000],
        'gamma': [0.001, 0.0001]
    }
}

grid_search = GridSearchCV(classifiers, param_grids, cv=5, scoring='f1_weighted')

grid_search.fit(X_val_scaled, y_val)

grid_search.best_params_
grid_search.score(X_val_scaled, y_val)

grid_search.best_score_
grid_search.best_params_
print()
```

And we have found the best combination of model and scale that we can have as follows:

```
best_model_key = max(best_scores, key=best_scores.get)
best_model_key

('StandardScaler', 'XGBoost')
```

# This is where the real magic happens.

So we moved forward and built our model with the things we had understood and we reached these results.

Validation Metrics:  
Accuracy: 0.9615384615384616  
Precision: 0.9589180050718513  
Recall: 0.9615384615384616  
F1 Score: 0.9573792985164221  
AUC Score: 0.9685459684581547

Test Metrics:  
Accuracy: 0.9798488664987406  
Precision: 0.9783299592211258  
Recall: 0.9798488664987406  
F1 Score: 0.9784632361760827  
AUC Score: 0.9840589003117238

# first, we needed to split our data

```
X = df.drop('os' , axis=1)
Y = df['os']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.15, random_state=42)

X_tr_tr, X_val_tr, y_tr_tr, y_val_tr = train_test_split(X_train, Y_train, test_size=0.15, random_state=42)
```

then we balanced our data

## Handling Imbalanced Data

SMOTE

```
smt=SMOTE(random_state=42)

X_tr_tr_balanced_smote, y_tr_tr_balanced_smote = smt.fit_resample(X_tr_tr, y_tr_tr)

print(X_tr_tr.shape, y_tr_tr.shape, np.unique(y_tr_tr, return_counts=True))

print('-----')

print('Balanced :')
print(X_tr_tr_balanced_smote.shape, y_tr_tr_balanced_smote.shape, np.unique(y_tr_tr_balanced_smote, return_counts=True))

(1911, 38) (1911,) (array([0, 1, 2], dtype=int64), array([1765,    97,    49], dtype=int64))
-----
Balanced :
(5295, 38) (5295,) (array([0, 1, 2], dtype=int64), array([1765, 1765, 1765], dtype=int64))
```

and then we tried to find the best combination of model and scaling using Gridsearch

## Model Selection and Scaling using GridSearch

```
classifiers = {
    'Decision Tree': DecisionTreeClassifier(),
    'Random Forest': RandomForestClassifier(),
    'AdaBoost': AdaBoostClassifier(),
    'XGBoost': XGBClassifier()
}

param_grids = {
    'Decision Tree': {
        'max_depth': [3, 5, 7]
    },
    'Random Forest': {
        'n_estimators': [50, 100, 200]
    },
    'AdaBoost': {
        'n_estimators': [50, 100, 200],
        'learning_rate': [0.01, 0.1, 1]
    }
}

for scaler_name, scaler in [('StandardScaler', scaler_standard), ('MinMaxScaler', scaler_minmax)]:
    X_val_scaled = scaler.fit_transform(X_val_selected)
    for clf_name, clf in classifiers.items():
        grid_search = GridSearchCV(clf, param_grids[clf_name], cv=5, scoring='f1_weighted')
        grid_search.fit(X_val_scaled, y_val_tr)
        results[(scaler_name, clf_name)] = grid_search

        print(f"Scaler: {scaler_name}, Classifier: {clf_name}")
        print("Best Score:", grid_search.best_score_)
        print("Best Parameters:", grid_search.best_params_)
        print()
```

And we have found  
the best  
combination of  
model and scale that  
we can have as  
follows:

```
best_model_key = max(best_scores, key=best_scores.get)
best_model_key
('StandardScaler', 'XGBoost')
```

So we moved forward and built our model with the things we had understood and we reached these results.

#### Validation Metrics:

Accuracy: 0.9615384615384616

Precision: 0.9589180050718513

Recall: 0.9615384615384616

F1 Score: 0.9573792985164221

AUC Score: 0.9685459684581547

#### Test Metrics:

Accuracy: 0.9798488664987406

Precision: 0.9783299592211258

Recall: 0.9798488664987406

F1 Score: 0.9784632361760827

AUC Score: 0.9840589003117238

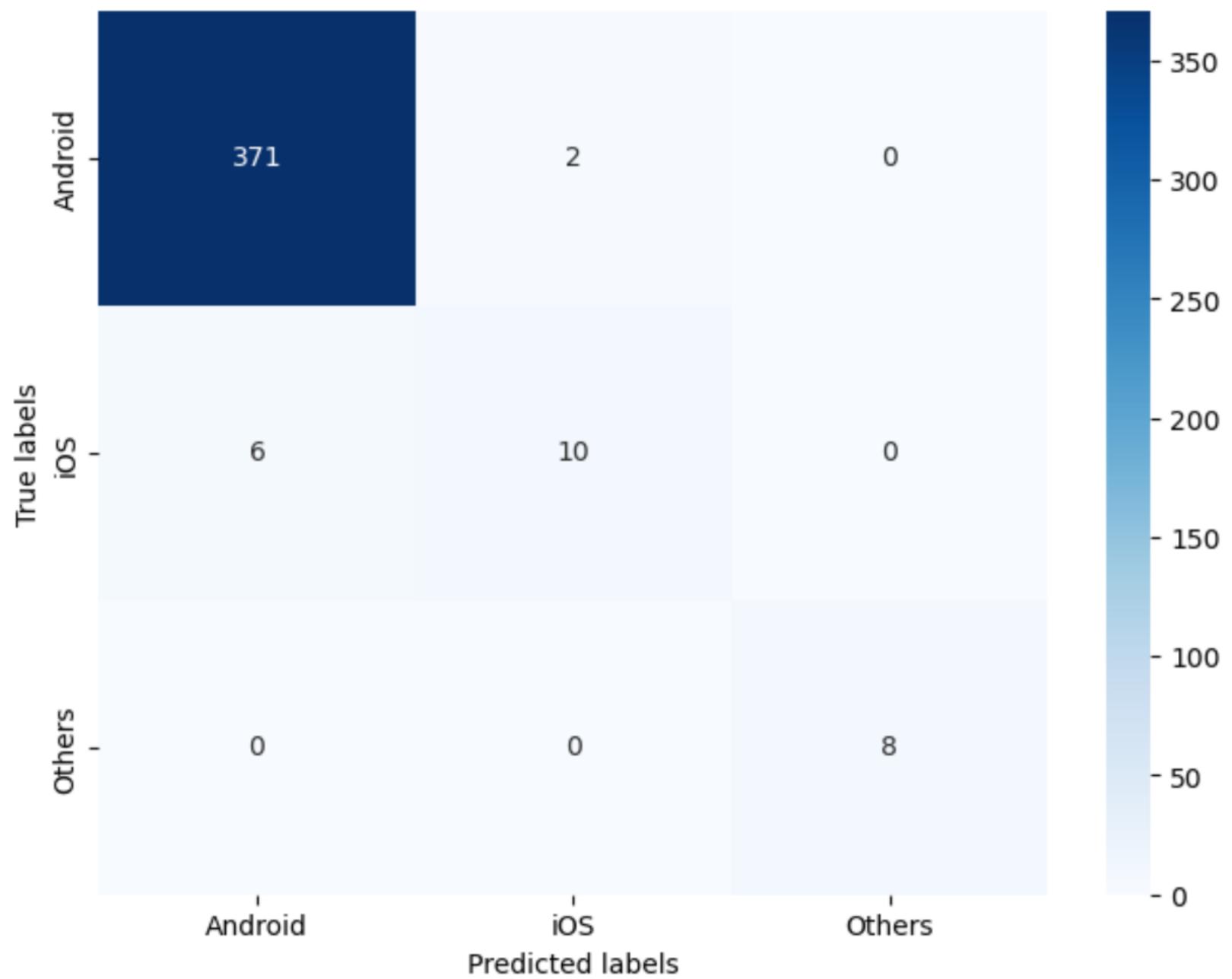
## Validation Metrics:

- **Accuracy:** The model correctly predicts about 96.15% of the validation data samples.
- **Precision:** When the model predicts a positive class label, it is correct about 95.89% of the time on average.
- **Recall:** The model identifies around 96.15% of the actual positive class instances in the validation dataset.
- **F1 Score:** The F1 score, which combines precision and recall, is approximately 95.74%.
- **AUC Score:** The high AUC score of 96.85% indicates that the model has strong discriminatory power between the positive and negative classes.

## Test Metrics:

- **Accuracy:** On the test dataset, the model achieves an accuracy of 97.98%.
- **Precision:** With a precision score of 97.83%, the model maintains a high level of accuracy when predicting positive class instances on the test data.
- **Recall:** The model captures about 97.98% of the actual positive class instances in the test dataset.
- **F1 Score:** The F1 score of 97.85% on the test data suggests a good balance between precision and recall.
- **AUC Score:** Despite being evaluated on unseen data, the AUC score of 98.41% on the test dataset indicates that the model's discriminatory power remains strong.

### Confusion Matrix - Test Data



we heard you  
asked for each  
class metrics, so  
here we go:

class 0: android  
class 1: ios  
class 2 :others

Metrics for Class 0:

Precision: 0.9840848806366048

Recall: 0.9946380697050938

F1 Score: 0.9893333333333334

Metrics for Class 1:

Precision: 0.8333333333333334

Recall: 0.625

F1 Score: 0.7142857142857143

Metrics for Class 2:

Precision: 1.0

Recall: 1.0

F1 Score: 1.0

**WOW THERE  
IS A PART 2,  
JUST WHAT  
WE NEEDED,  
MORE FUN!**

Alright, so We just needed to change the 'battery\_capacity' in Selected\_features to 'battery\_capacity Cat'

our model and  
scale didn't  
change!

here's the  
results:

#### Validation Metrics:

Accuracy: 0.9615384615384616

Precision: 0.9608341175672018

Recall: 0.9615384615384616

F1 Score: 0.9611606678004927

AUC Score: 0.9686820256234129

#### Test Metrics:

Accuracy: 0.9722921914357683

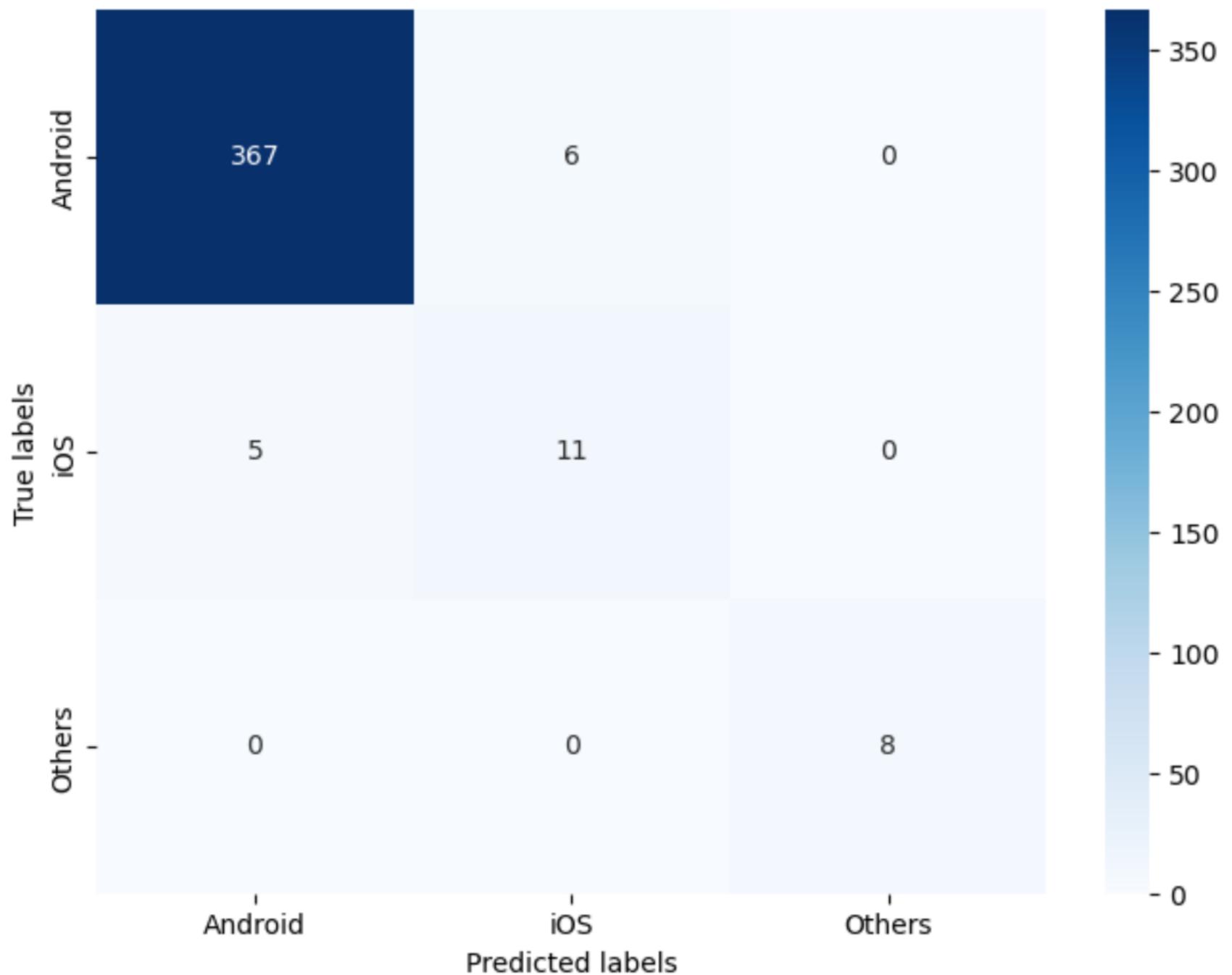
Precision: 0.9731473559603415

Recall: 0.9722921914357683

F1 Score: 0.9726934108723705

AUC Score: 0.9820786686181654

### Confusion Matrix - Test Data



Metrics for Class 0:

Precision: 0.9865591397849462

Recall: 0.9839142091152815

F1 Score: 0.985234899328859

Metrics for Class 1:

Precision: 0.6470588235294118

Recall: 0.6875

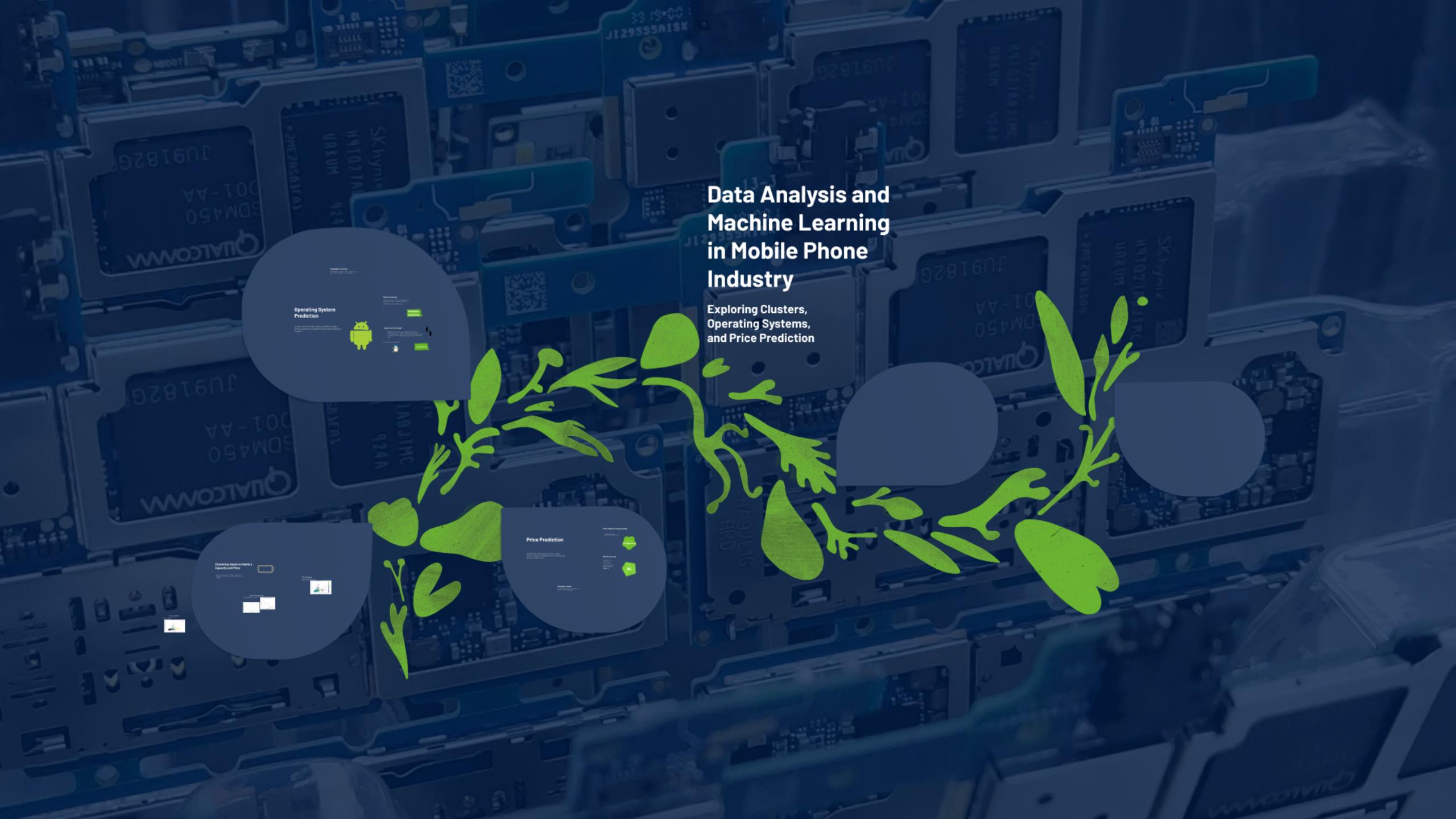
F1 Score: 0.6666666666666667

Metrics for Class 2:

Precision: 1.0

Recall: 1.0

F1 Score: 1.0



# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction



# Price Prediction

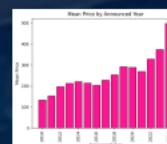
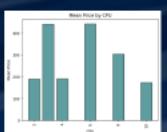
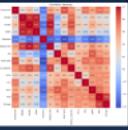
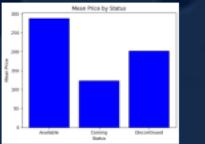
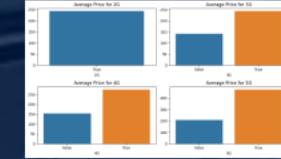
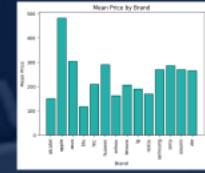
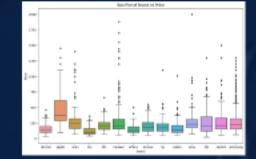
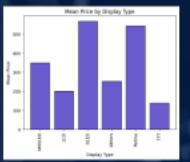
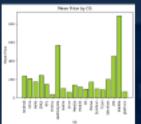
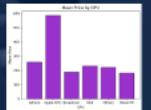
Predicting mobile phone prices using  
machine learning algorithms and analyzing  
feature importance.

# Price Prediction Preprocessing

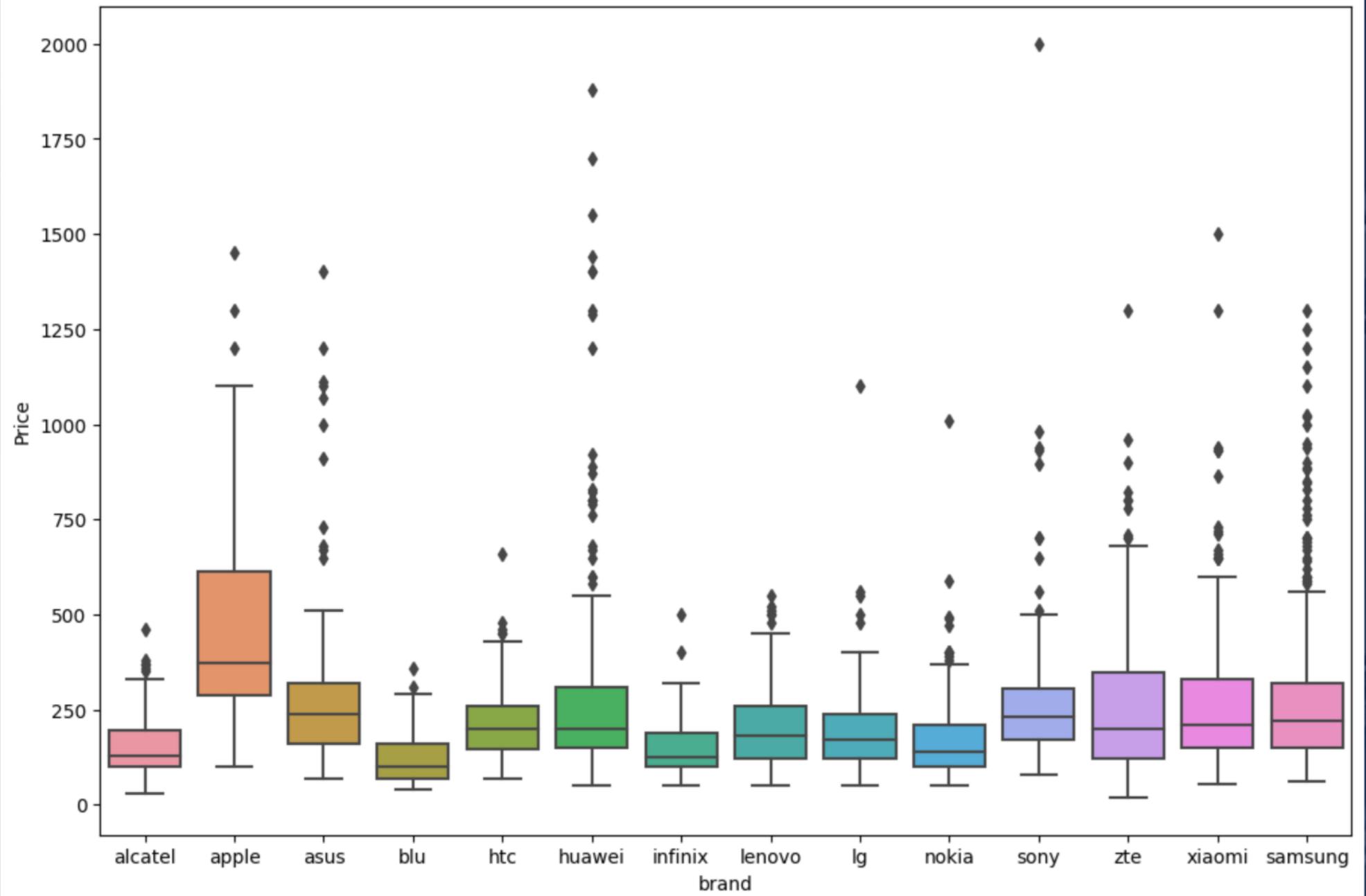
Once again, here we are  
just let's go for it without talking.



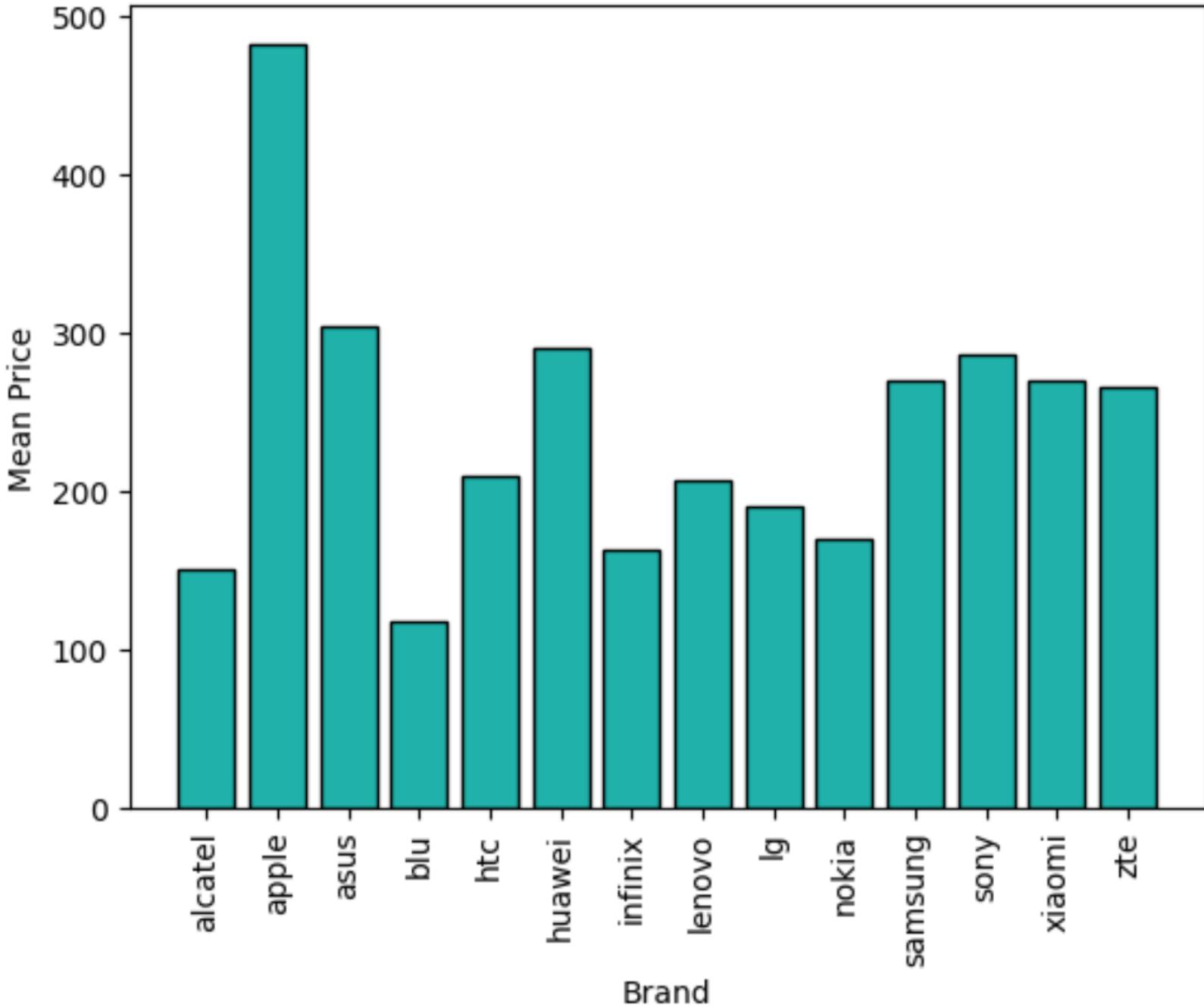
# Data Discovery Area



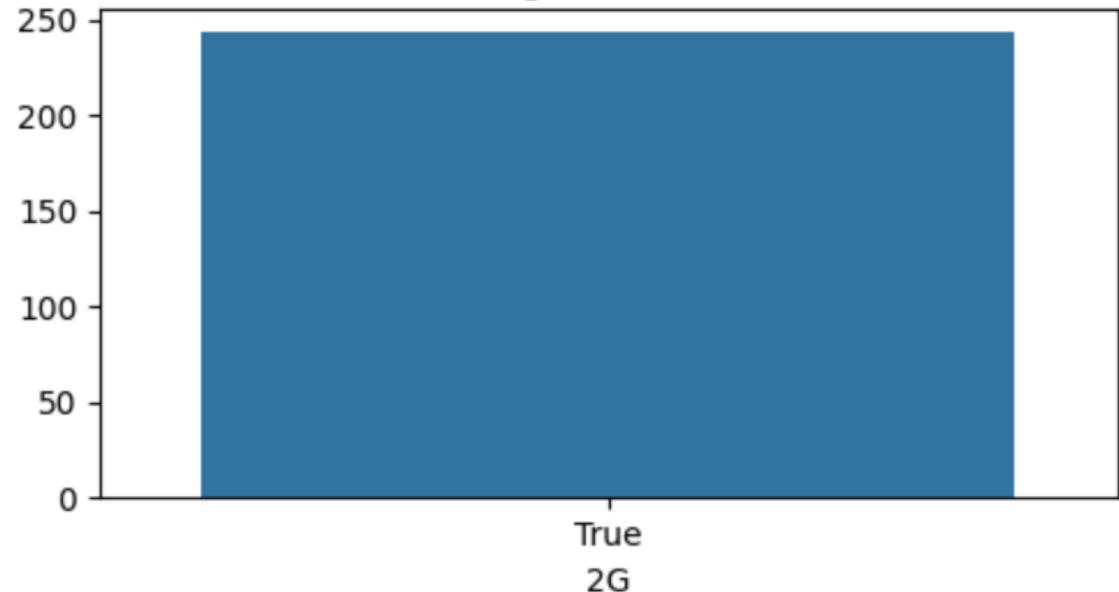
Box Plot of Brand vs Price



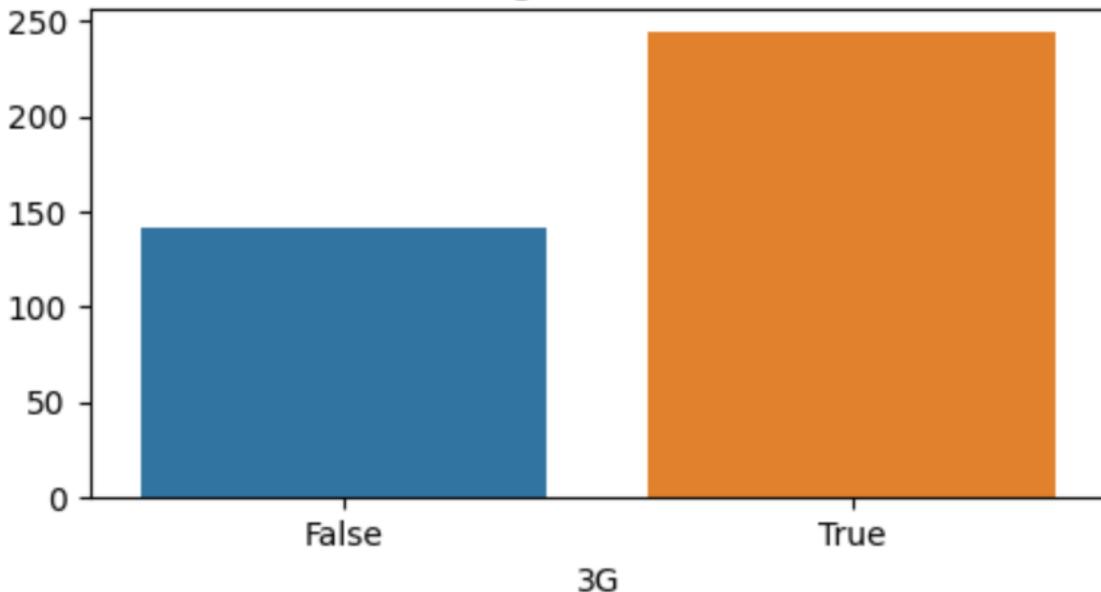
### Mean Price by Brand



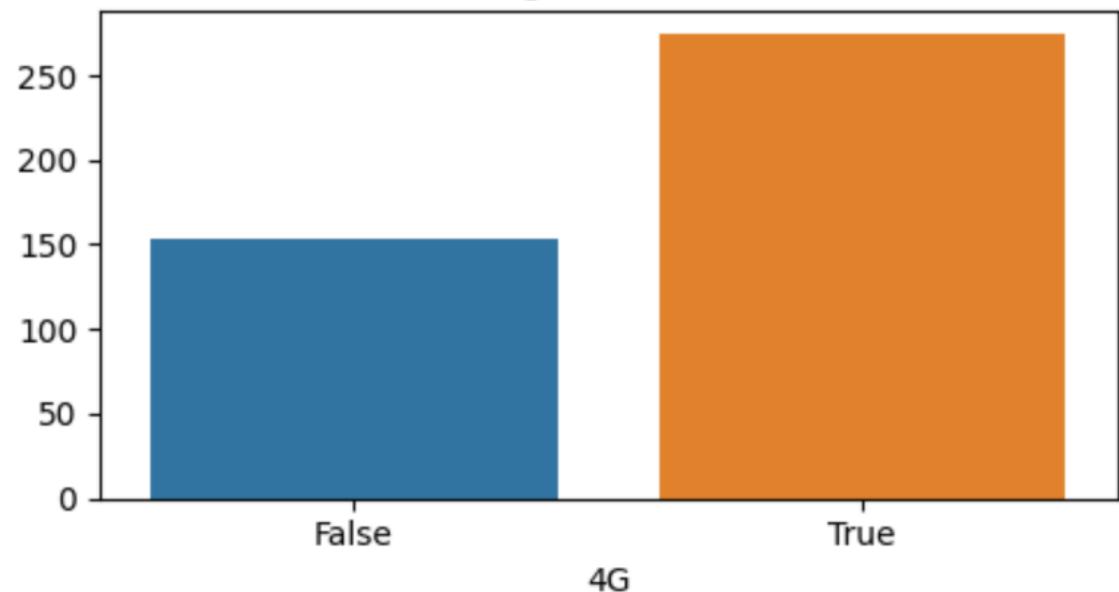
Average Price for 2G



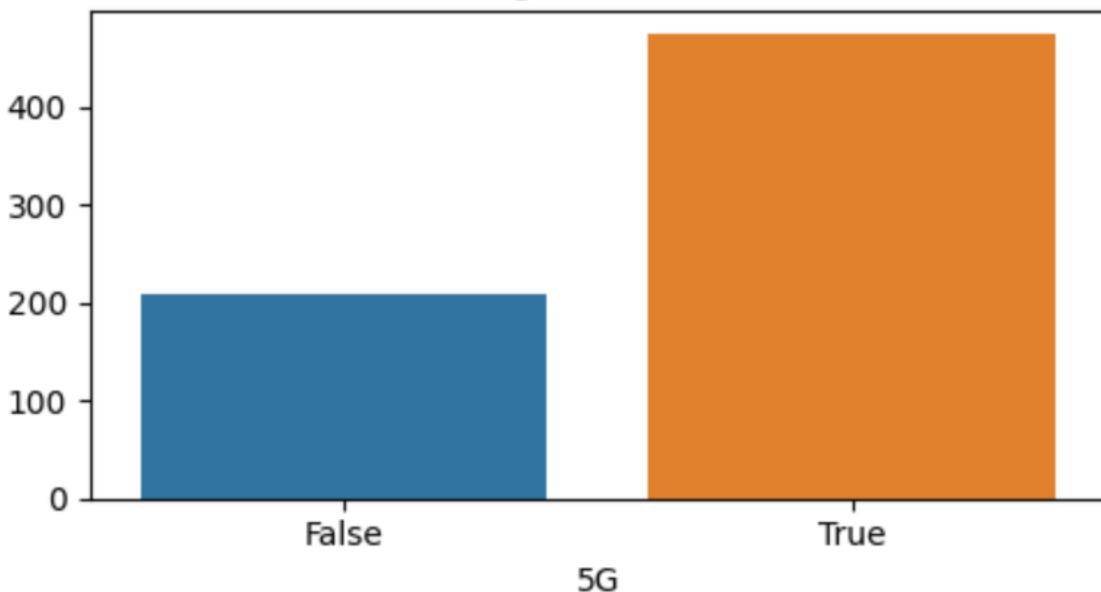
Average Price for 3G



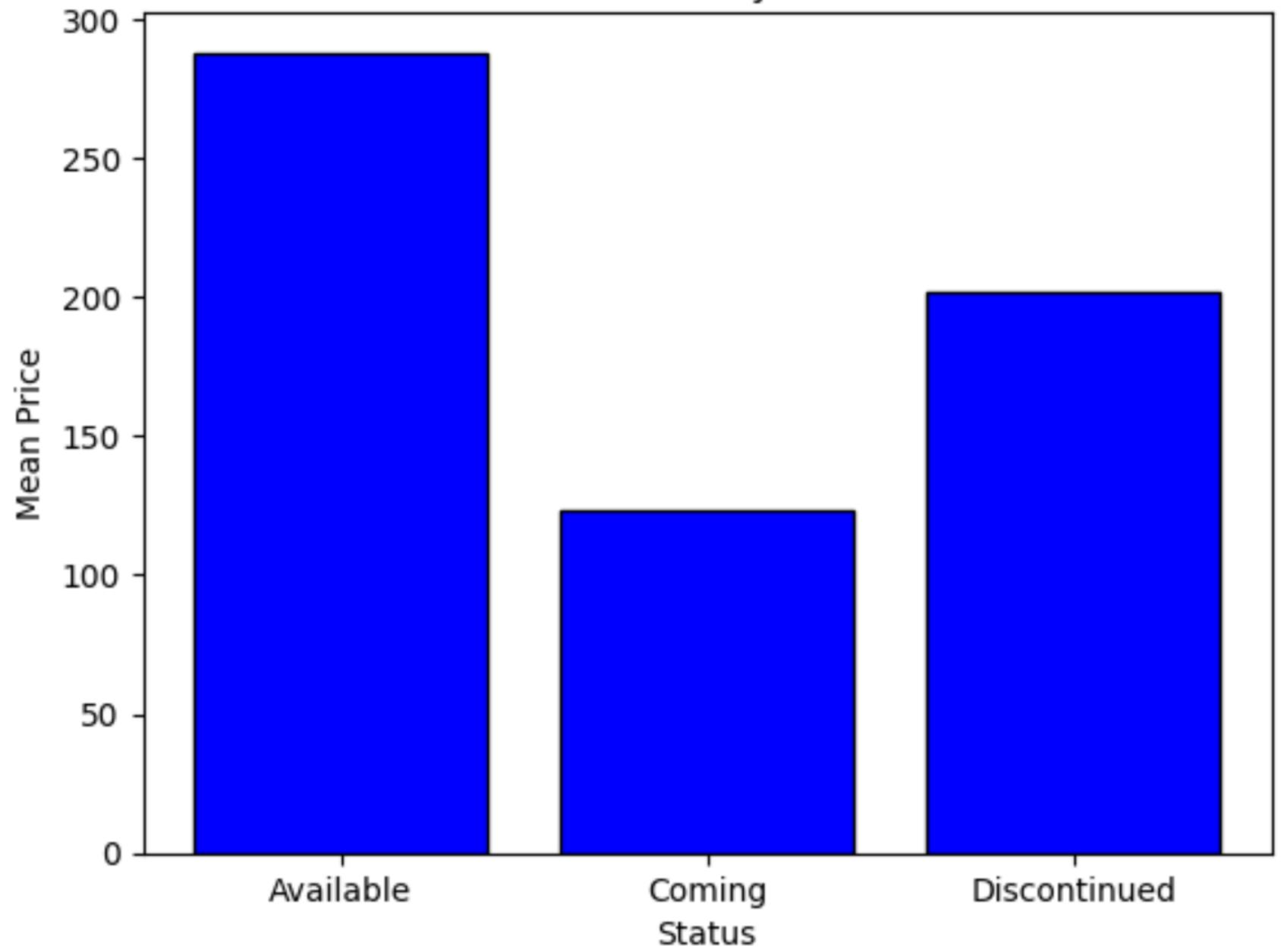
Average Price for 4G



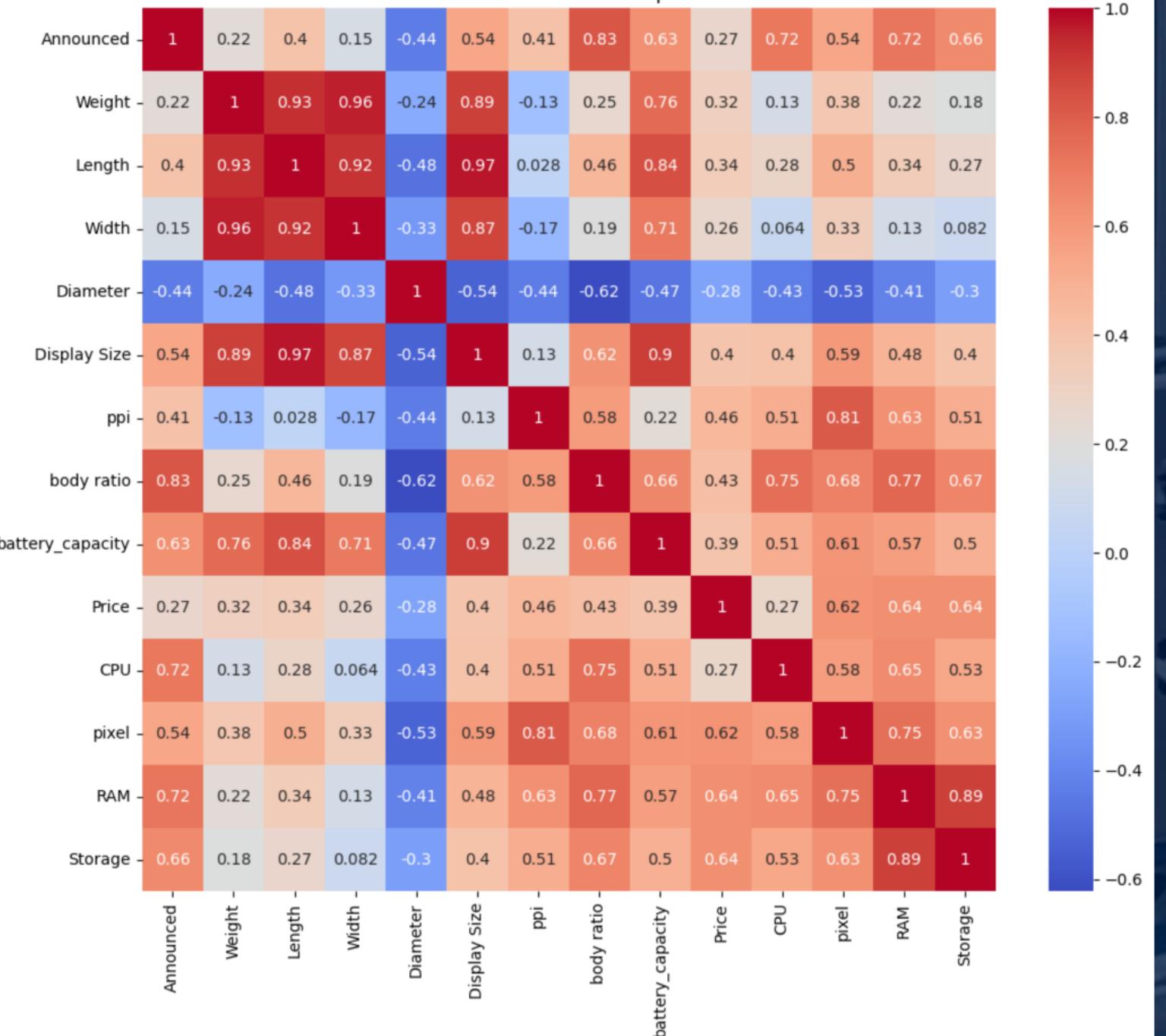
Average Price for 5G

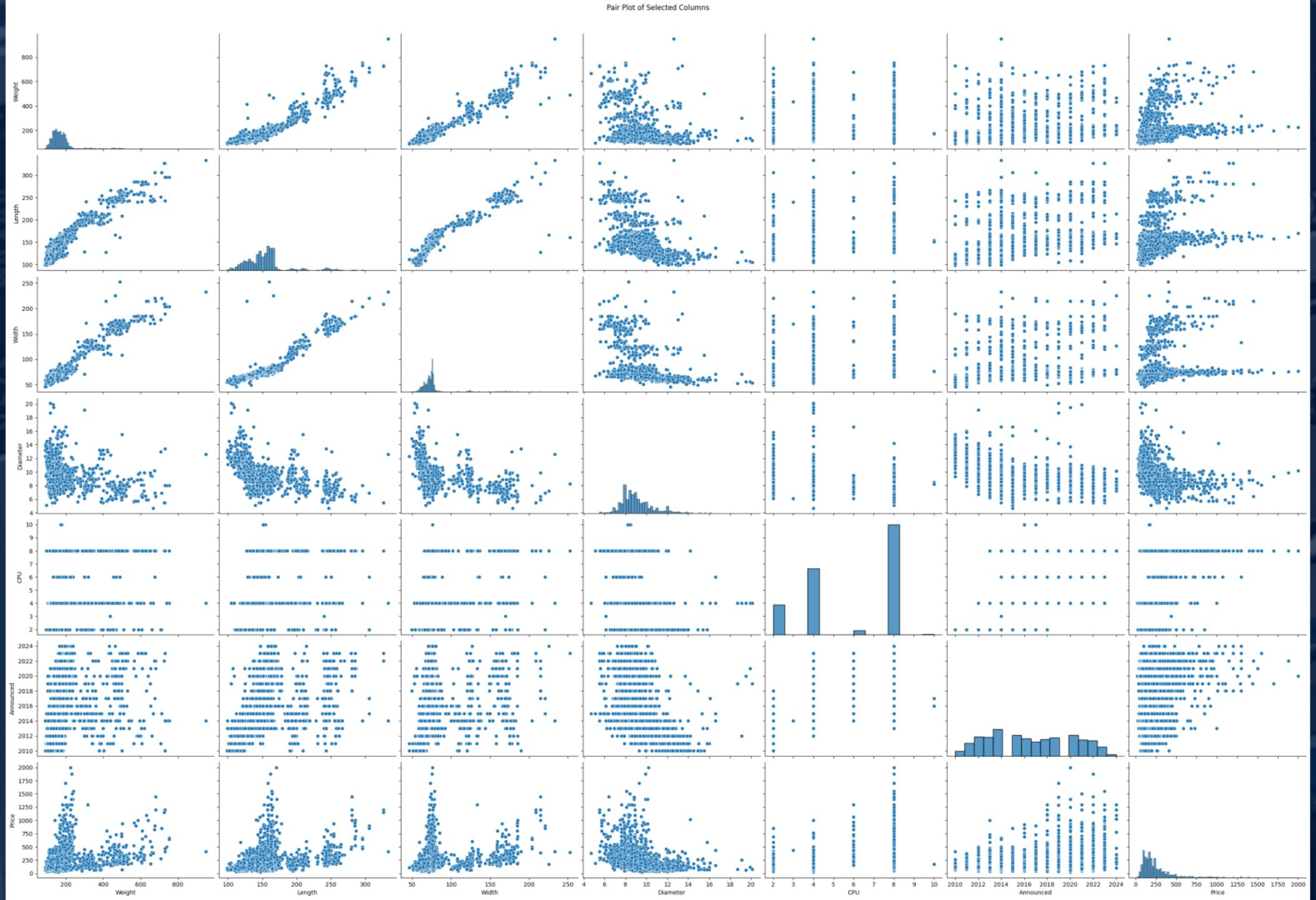


### Mean Price by Status

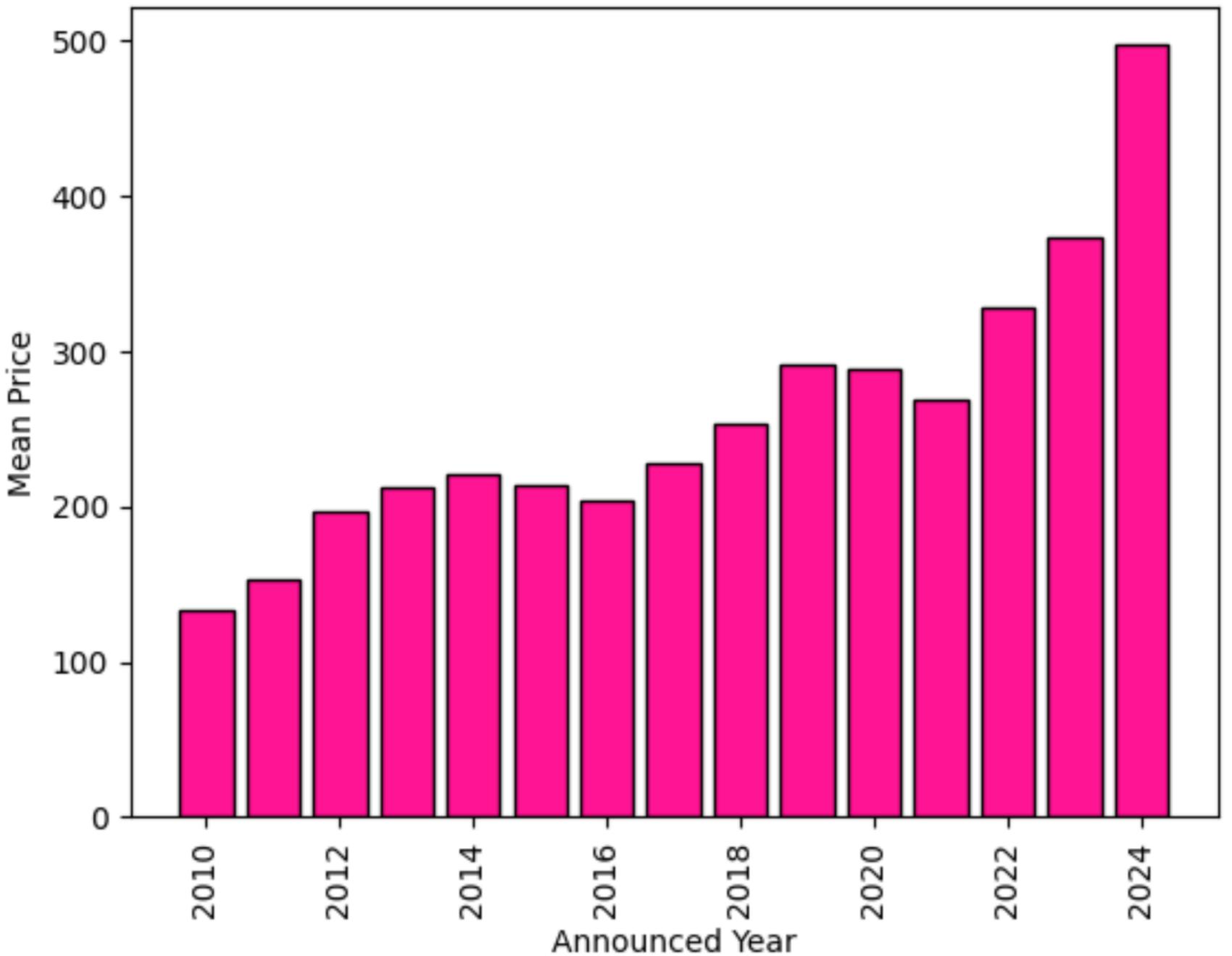


Correlation Heatmap

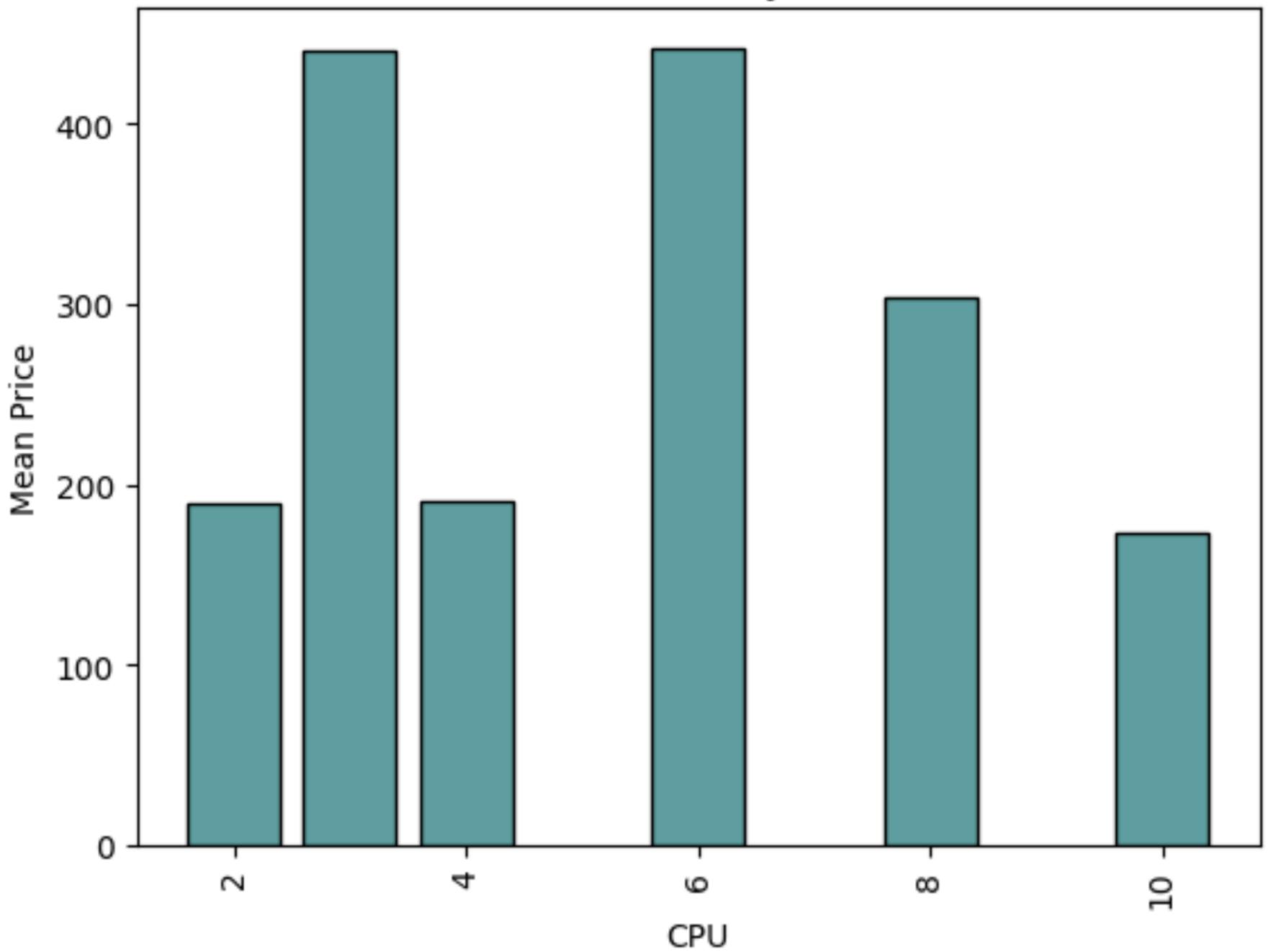




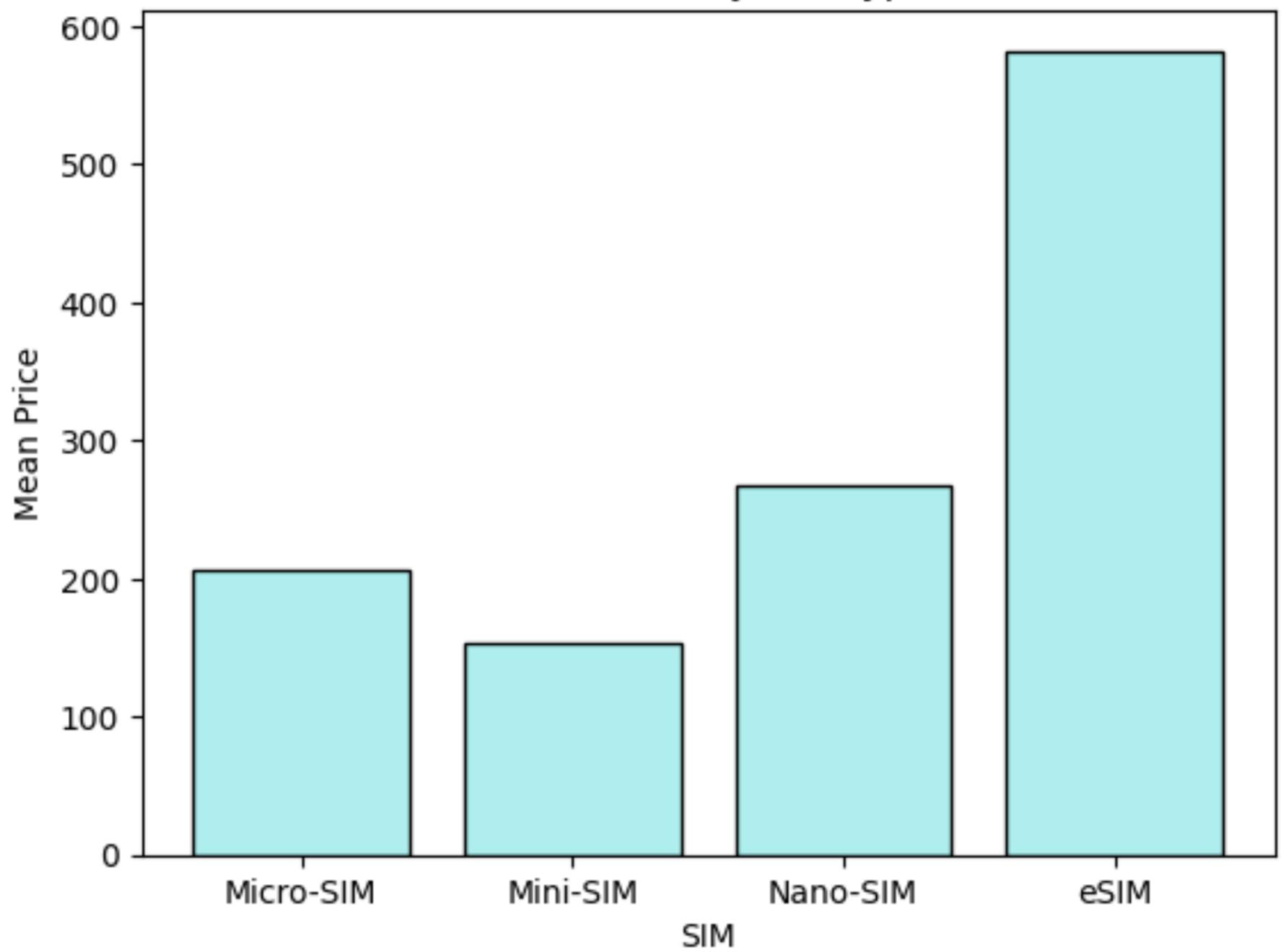
### Mean Price by Announced Year



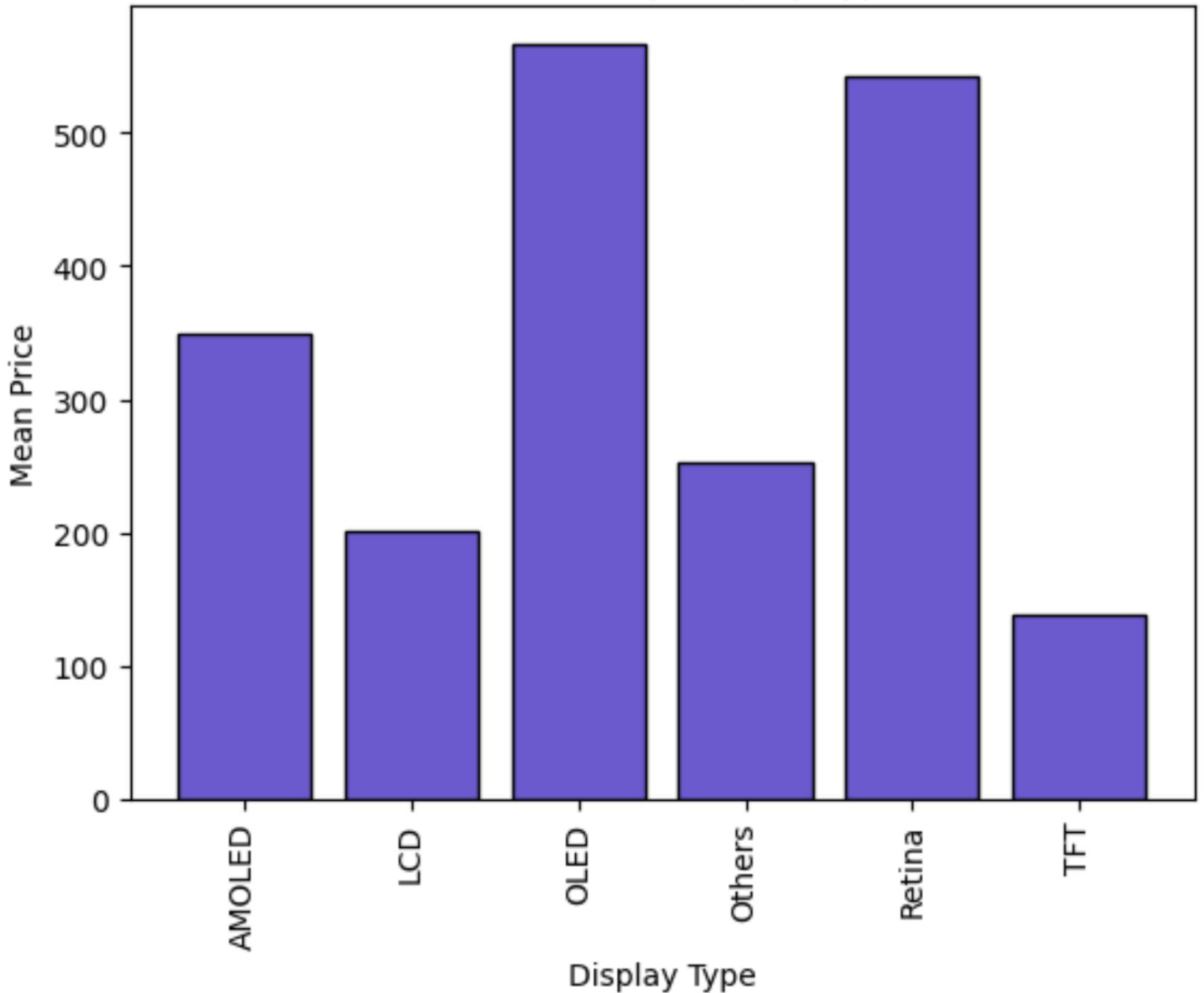
## Mean Price by CPU



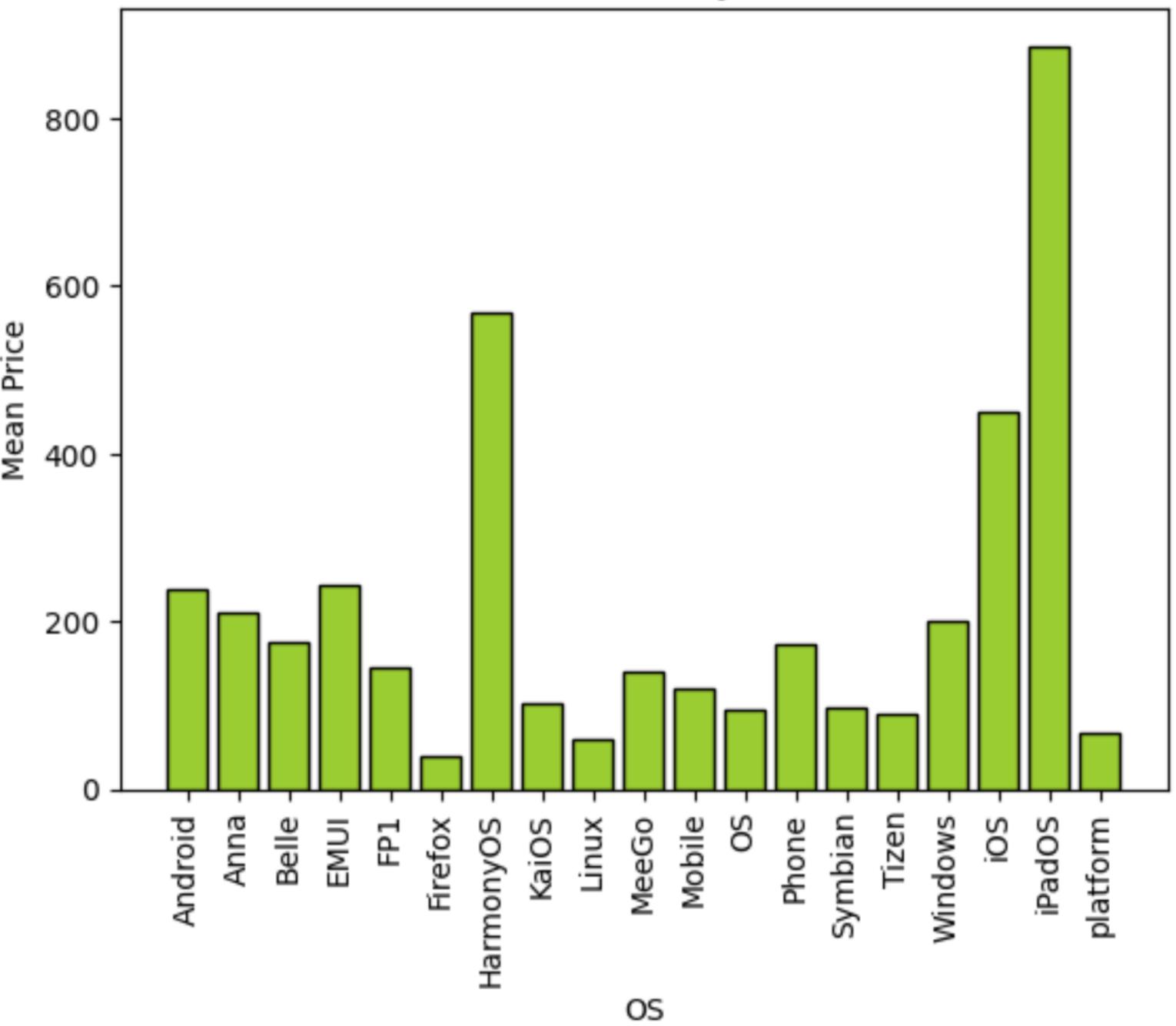
## Mean Price by SIM Type



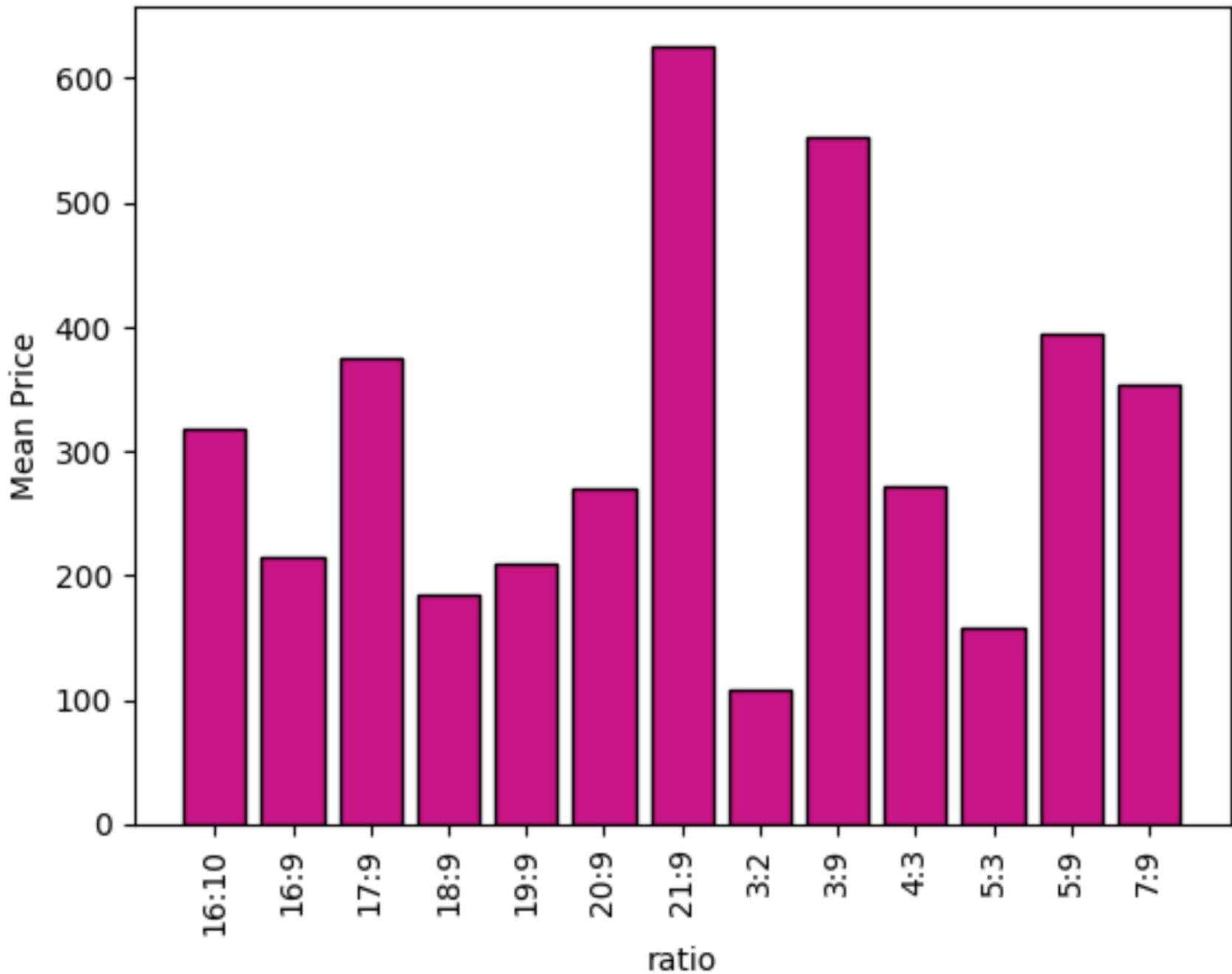
## Mean Price by Display Type



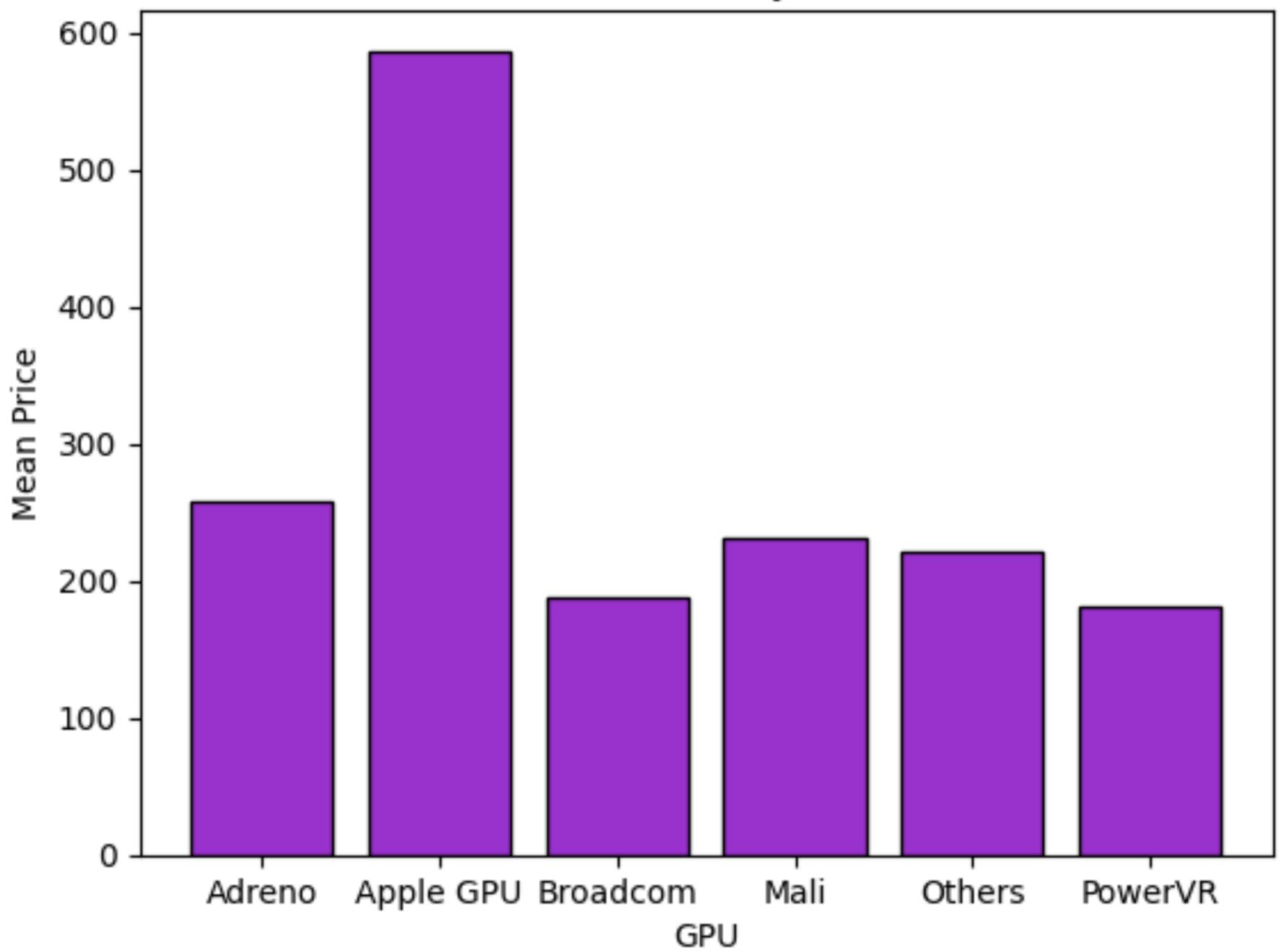
## Mean Price by OS



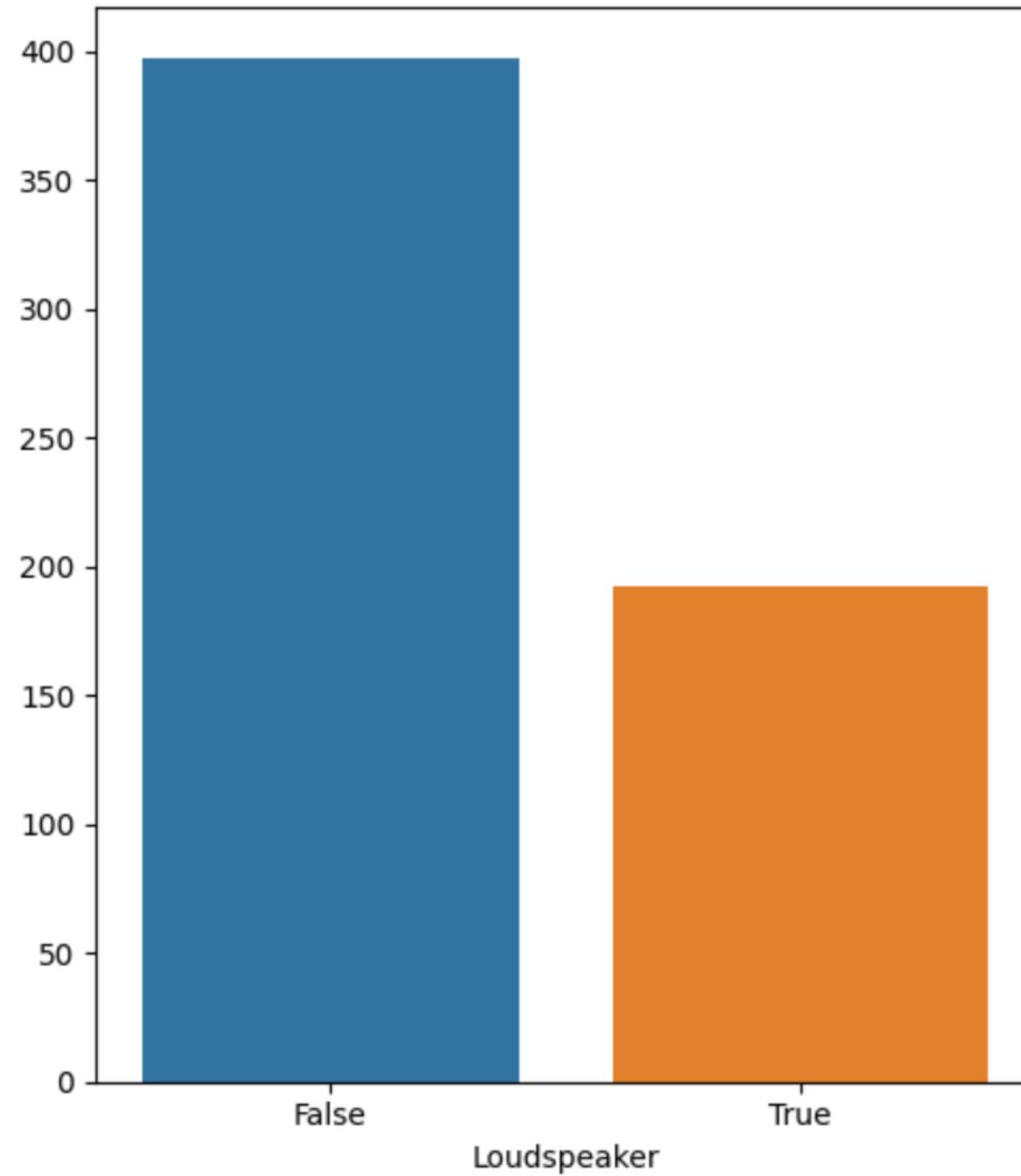
### Mean Price by ratio



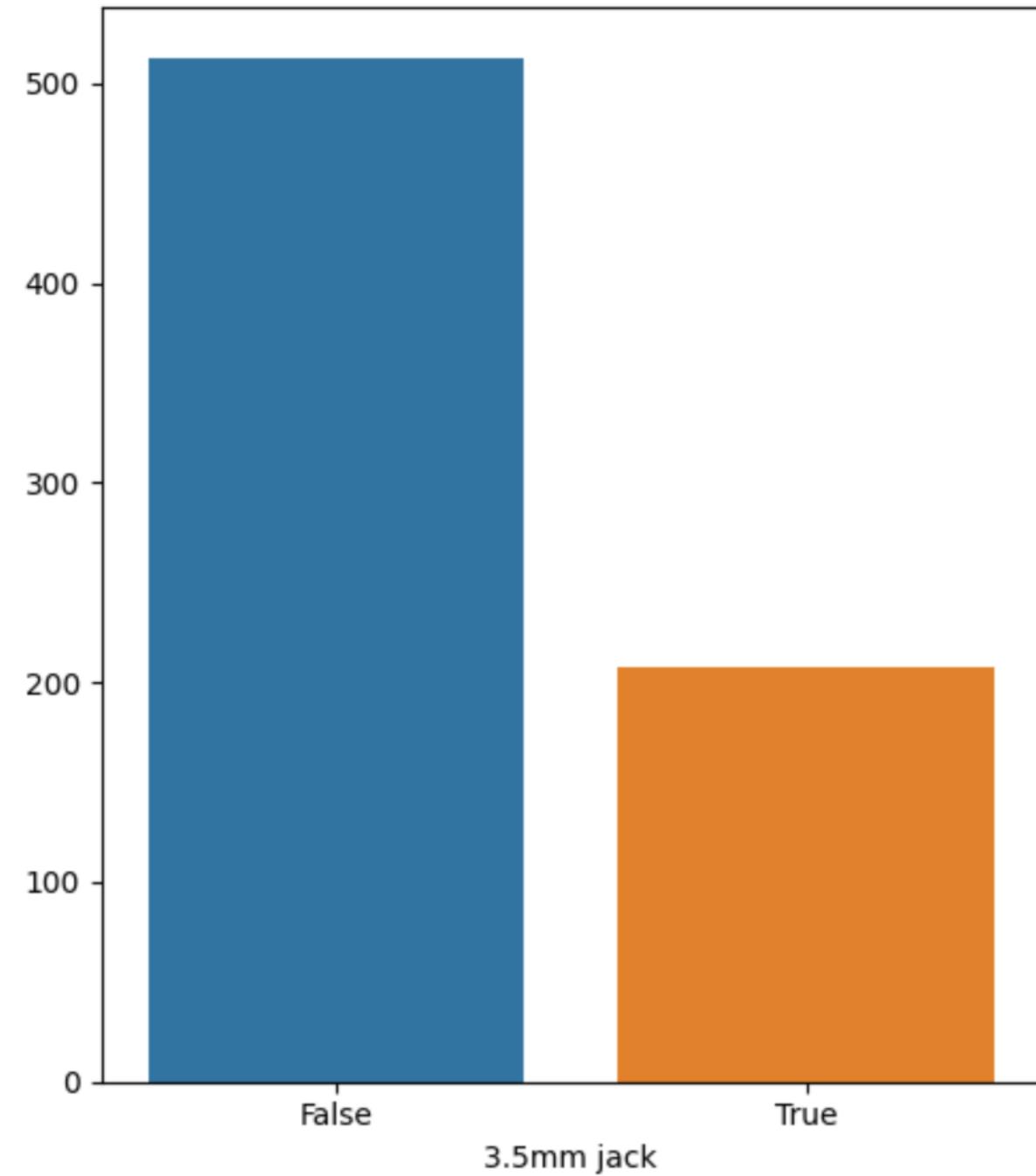
### Mean Price by GPU



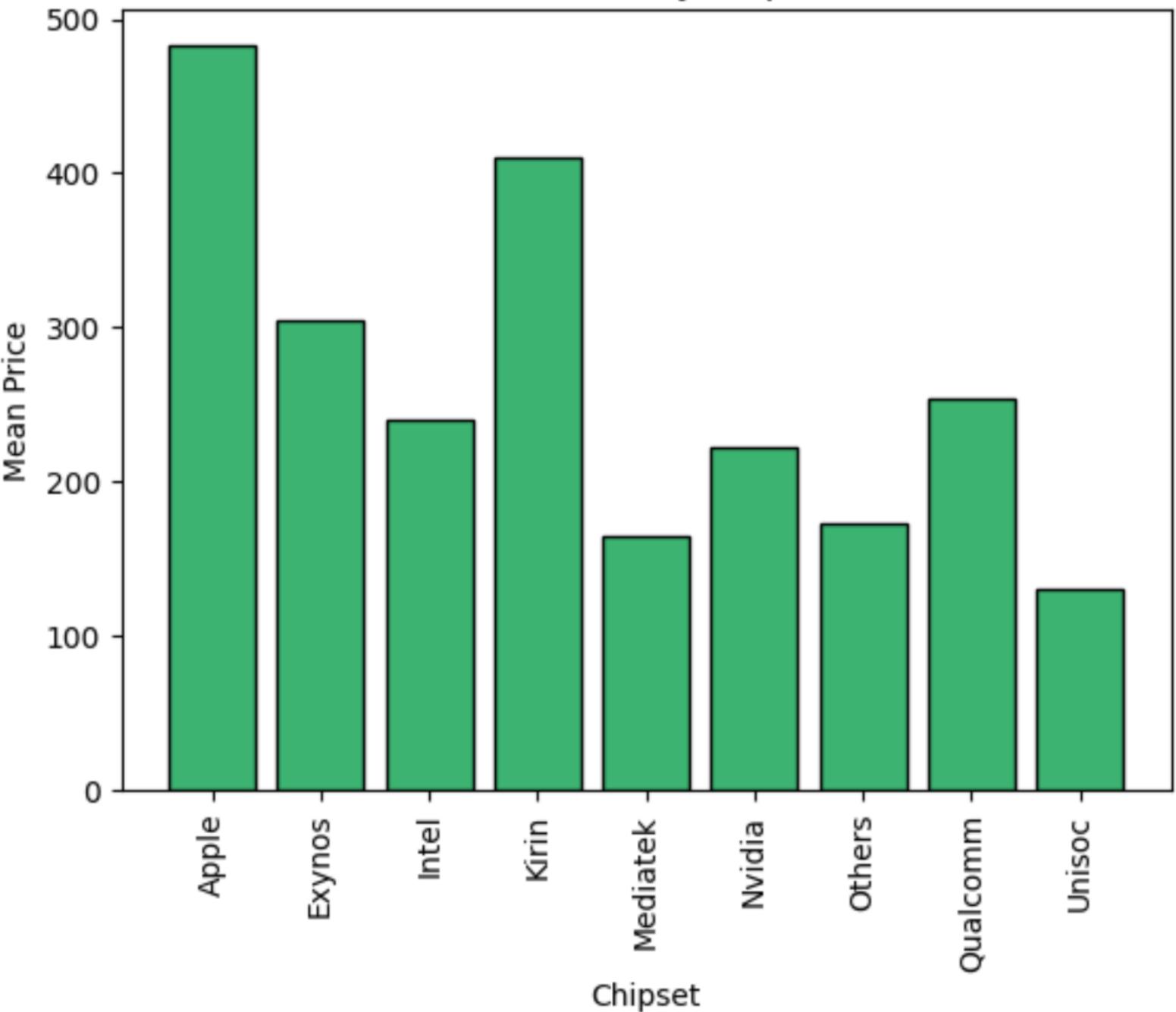
Average Price for Loudspeaker



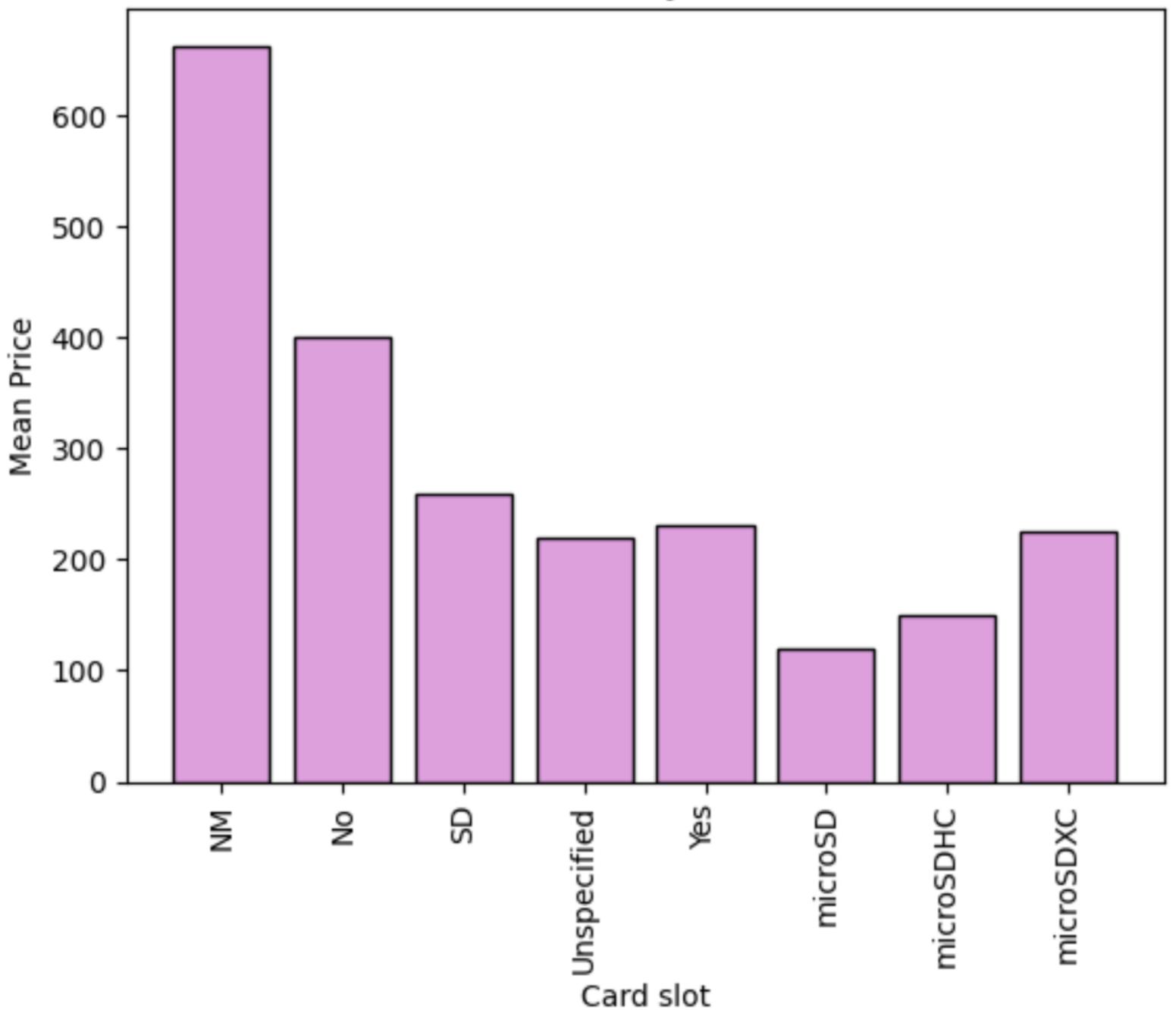
Average Price for 3.5mm jack



## Mean Price by Chipset



### Mean Price by Card slot



# drop wlan/color/network

```
df.drop(columns=['WLAN'], inplace=True)
df.drop(columns=['Colors'], inplace=True)
df.drop(columns=['Network'], inplace=True)
df.drop(columns=['Sensors'], inplace=True)
```

let's handle  
those missing  
values

## Length Width Diameter

```
mean_length = df['Length'].mean()  
mean_width = df['Width'].mean()  
mean_diameter = df['Diameter'].mean()  
  
df['Length'].fillna(mean_length, inplace=True)  
df['Width'].fillna(mean_width, inplace=True)  
df['Diameter'].fillna(mean_diameter, inplace=True)
```

## CPU

```
cpu_counts = df.groupby(['OS', 'brand', 'CPU']).size().reset_index(name='COUNT')
max_cpu = cpu_counts.loc[cpu_counts.groupby(['OS', 'brand'])['COUNT'].idxmax()]
max_cpu.drop(columns='COUNT', inplace=True)
max_cpu.head()
```

|    | OS      | brand   | CPU |
|----|---------|---------|-----|
| 1  | Android | alcatel | 4.0 |
| 4  | Android | asus    | 4.0 |
| 8  | Android | blu     | 4.0 |
| 12 | Android | htc     | 8.0 |
| 16 | Android | huawei  | 8.0 |

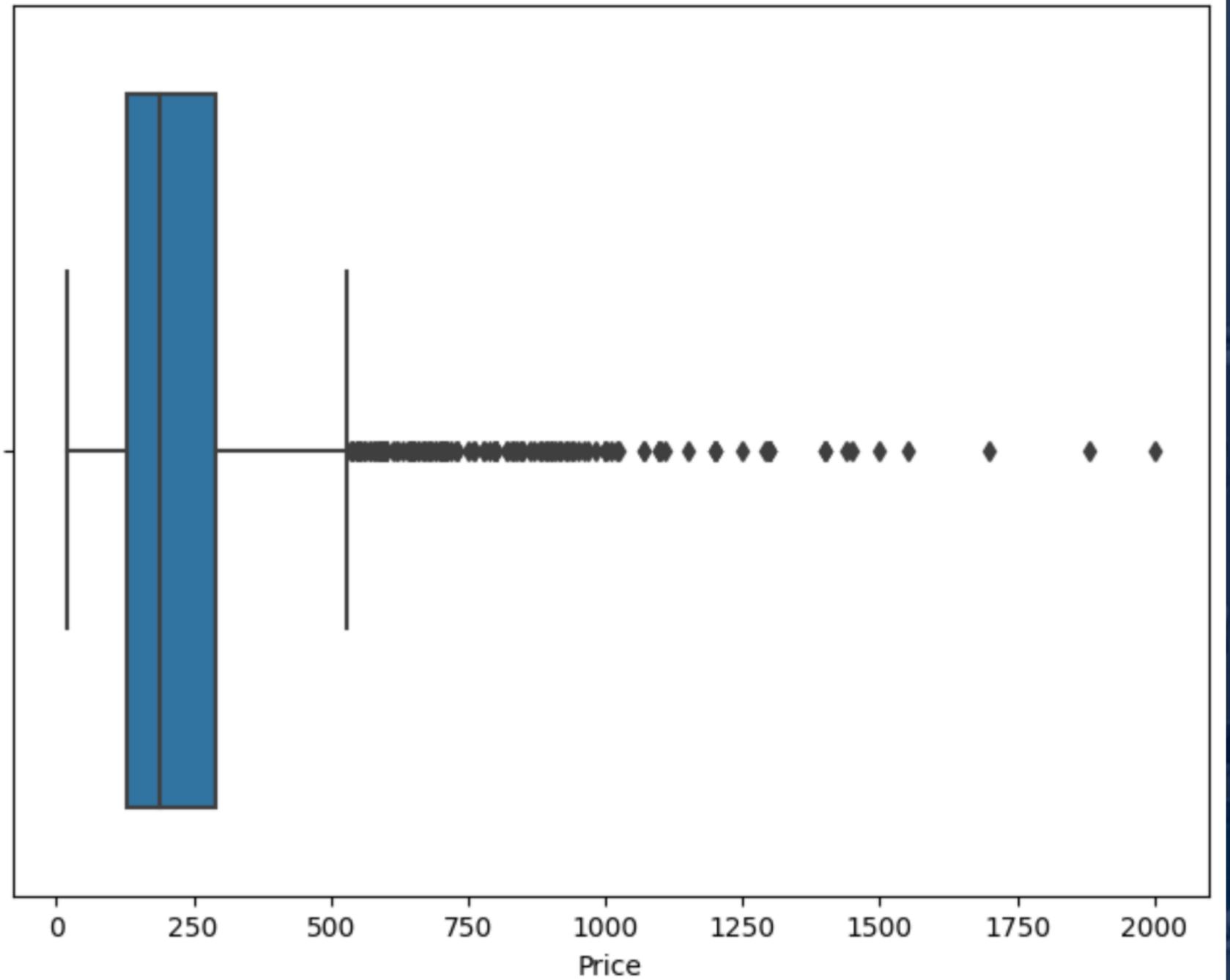
```
for index, row in df.iterrows():
    if pd.isnull(row['CPU']):
        os_value = row['OS']
        brand_value = row['brand']

        matching_rows = max_cpu[(max_cpu['OS'] == os_value) & (max_cpu['brand'] == brand_value)]

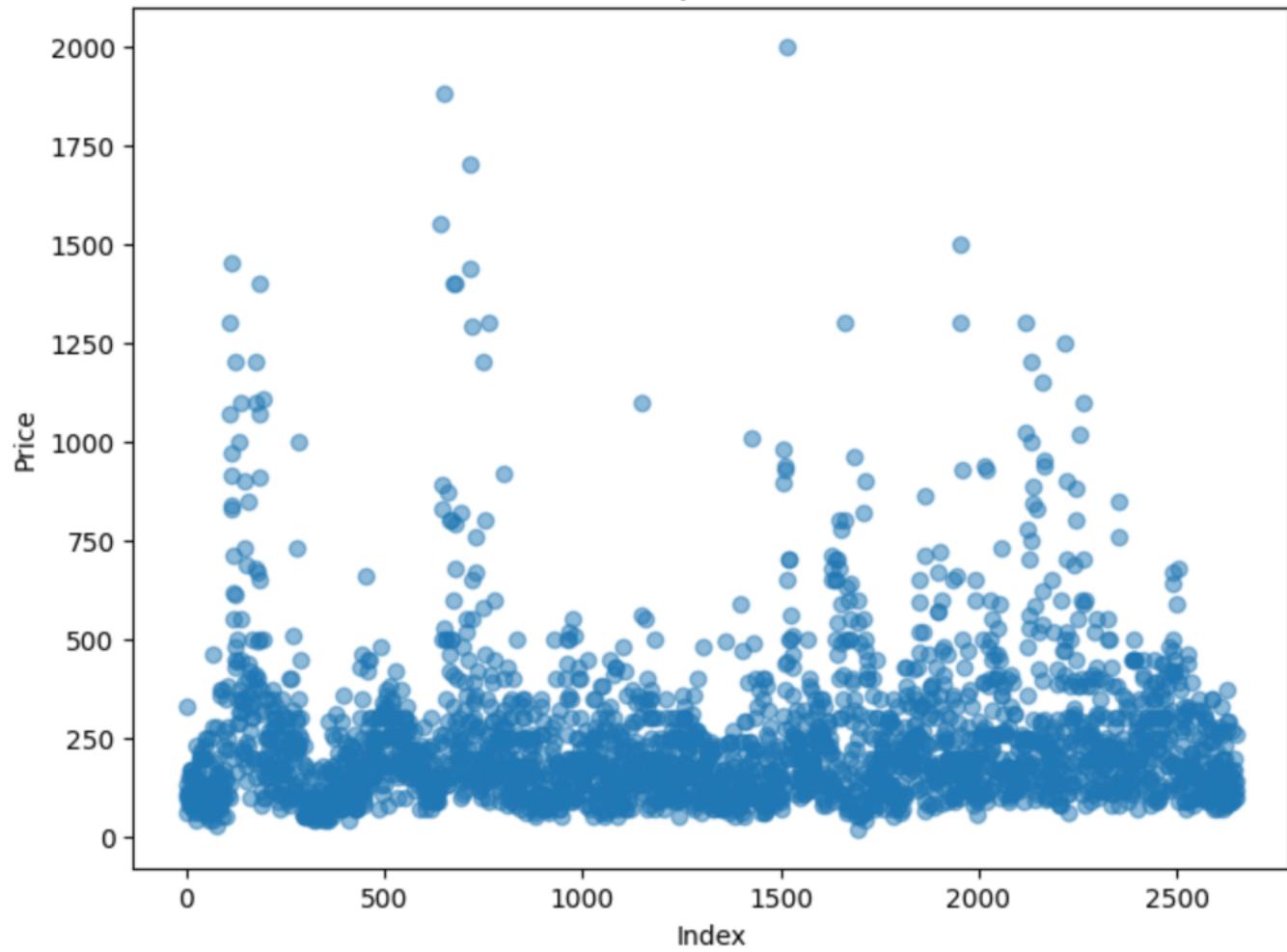
        if not matching_rows.empty:
            corresponding_cpu = matching_rows.iloc[0]['CPU']
            df.at[index, 'CPU'] = corresponding_cpu
        else:
            print(os_value, brand_value)
            df.at[index, 'CPU'] = 4

df['CPU'].isnull().sum()
```

### Boxplot of Price

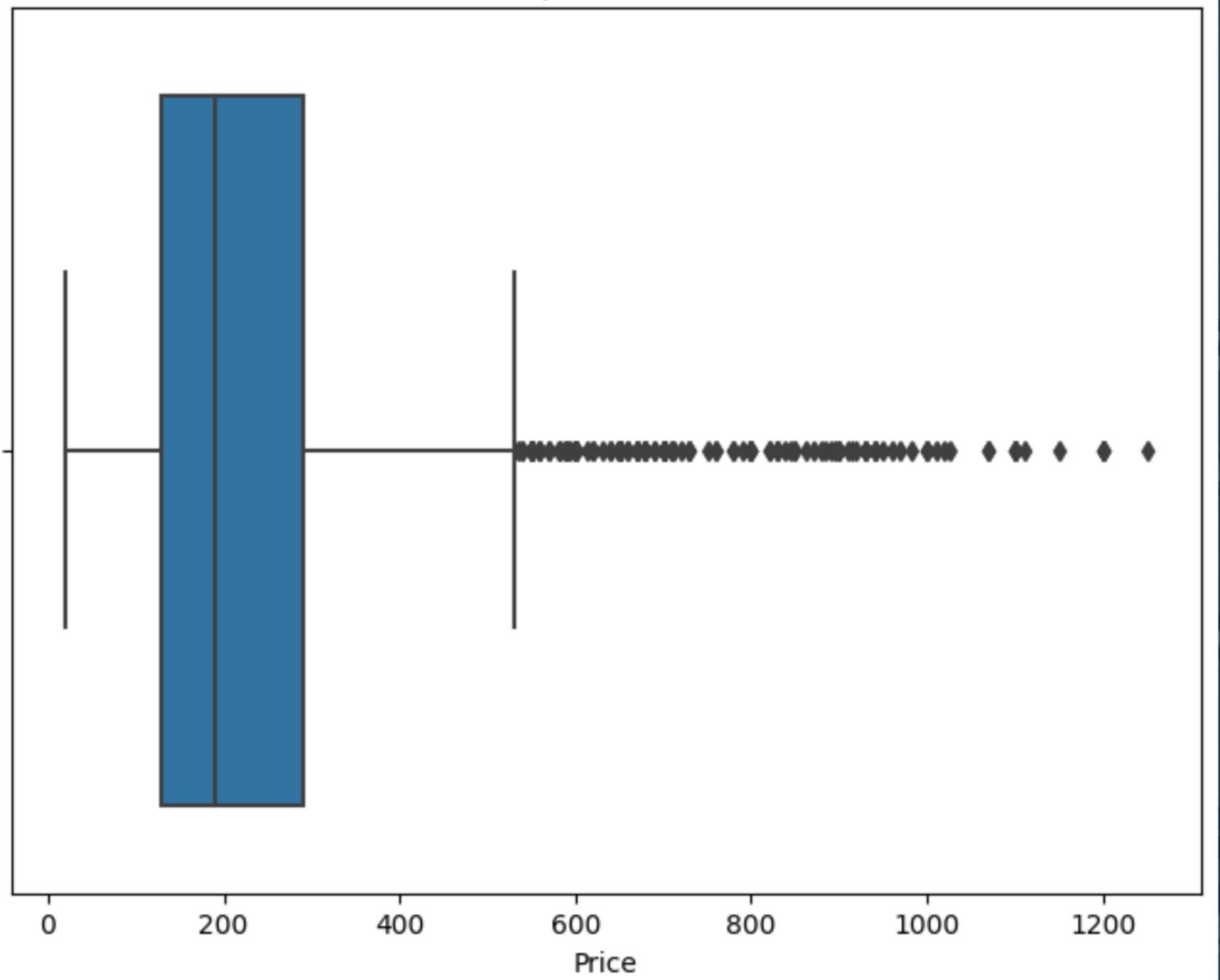


Scatter plot of Price



```
df['Price'] = df[df['Price'] <= 1250].Price
```

### Boxplot of Price



## Handling Categorical Variables:

```
label_encoder = LabelEncoder()  
  
columns_to_encode = ['brand', 'ratio', 'os', 'GPU', 'chipset', 'status', 'card slot', 'SIM', 'Display Type']  
  
for column in columns_to_encode:  
    df[column] = label_encoder.fit_transform(df[column])
```

# droped columns

| DELETED COLUMNS |
|-----------------|
| NAME            |
| Internal        |
| 2G              |
| WLAN            |
| COLORS          |
| SENSORS         |
| BLUETOOTH       |
| NETWORK         |

TOTAL COUNT: 8

# columns to keep

| COLUMNS TO KEEP FOR 100% |
|--------------------------|
| BRAND                    |
| DISPLAY TYPE             |
| OS                       |
| CPU                      |
| GPU                      |
| Chipset                  |
| RAM                      |
| Storage                  |

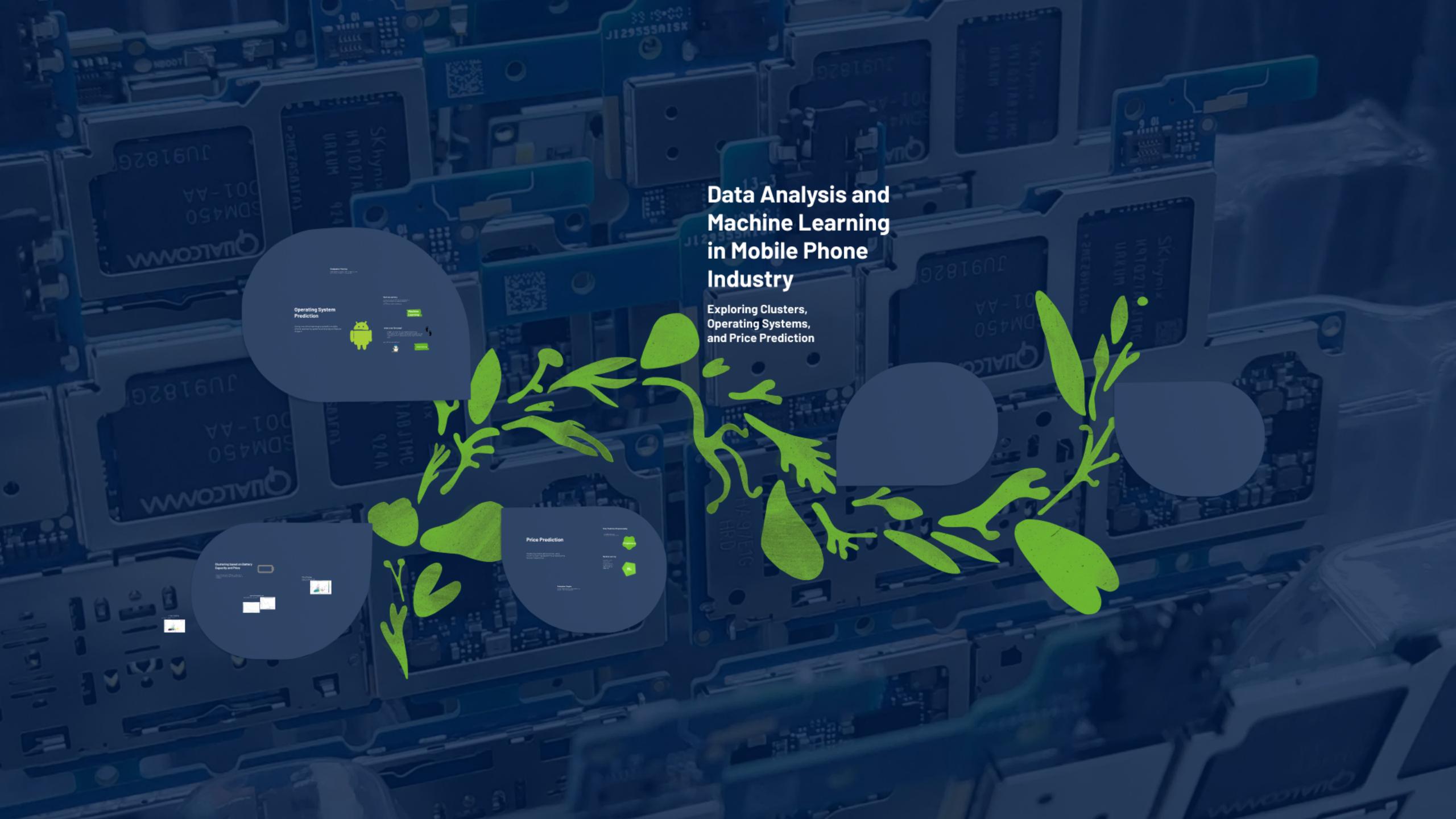
TOTAL COUNT: 8

| COLUMNS TO KEEP FOR 50% |
|-------------------------|
| 3G                      |
| 4G                      |
| 5G                      |
| SIM                     |
| Display Size            |
| ppi                     |
| body ratio              |
| battery_capacity        |
| pixel                   |

TOTAL COUNT: 6

| COLUMNS TO KEEP FOR LESS THAN 50% |
|-----------------------------------|
| Announced                         |
| Status                            |
| Weight                            |
| Length                            |
| Width                             |
| Loudspeaker                       |
| 3.5mm jack                        |
| Chipset                           |
| Card slot                         |
| ratio                             |

TOTAL COUNT: 10



# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction



Mobile Phone



Mobile Phone

The mobile phone  
industry is highly  
competitive, with  
numerous players  
offering a wide  
range of products.  
In this analysis,  
we will explore  
the use of data  
analysis and  
machine learning  
to predict price  
clusters and  
operating systems  
based on shared  
characteristics.

We will also  
examine the  
relationship  
between price  
and various  
factors such as  
processor type,  
RAM capacity,  
storage size,  
display size,  
camera quality,  
battery life,  
and software  
version.

By understanding  
these patterns,  
manufacturers  
can better  
position their  
products and  
target specific  
markets.

Overall, this  
analysis provides  
valuable insights  
into the mobile  
phone industry  
and highlights  
the potential  
of data analysis  
and machine  
learning to  
drive innovation  
and growth.

Stay tuned for  
more updates  
and analysis  
on the mobile  
phone industry.

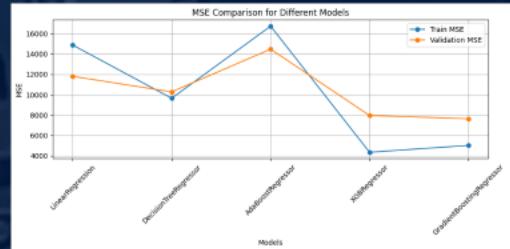
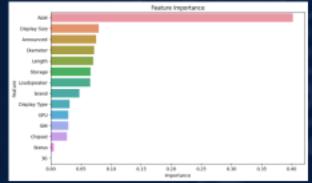
If you have any  
questions or  
comments,  
please feel  
free to  
leave a  
message  
below.

# Machine Learning

Now that our work with preprocessing is done, we can start machine learning



**we are here  
training the  
machines.**



```

once again, we
find the best
model using
Gridsearch

```

```
Best model: GradientBoostingRegressor with StandardScaler  
Best train R2 score: 0.809 / 0.809 (0.809)  
Best Validation R2 score: 0.809 / 0.804 (0.804)  
Best params: {'learning_rate': 0.05, 'max_depth': 3, 'n_estimators': 200}
```



```
x = df.drop('Price', axis=1)
Y=df['Price']

x_train, x_test, Y_train, Y_test = train_test_split(x, Y, test_size=0.15, random_state=42)

x_tr, x_val, y_tr, y_val= train_test_split(x_train, Y_train, test_size=0.15, random_state=42)
```

```
selected_features = ['brand', '3G','Announced','Length', 'Diameter', 'Loudspeaker','Chipset','GPU', 'RAM','Display Size', 'Storage','SIM' , 'Display Type' , 'Status' ]  
  
X_tr_selected = X_tr[selected_features]  
  
X_val_selected = X_val[selected_features]  
  
X_test_selected = X_test[selected_features]
```

Py



once again, we  
find the best  
model using  
Gridsearch

```
scaler = {  
    "StandardScaler": StandardScaler()  
}  
  
models = {  
    "LinearRegression": LinearRegression(),  
    "DecisionTreeRegressor": DecisionTreeRegressor(),  
    "AdaBoostRegressor": AdaBoostRegressor() ,  
    "XGBRegressor": xgb.XGBRegressor(),  
    "GradientBoostingRegressor": GradientBoostingRegressor()  
}  
  
params = {  
    "DecisionTreeRegressor": {  
        "max_depth": [3, 5, 7],  
        "min_samples_split": [2, 5, 10]  
    },  
    "XGBRegressor": {  
        "learning_rate": [0.05, 0.1, 0.15],  
        "max_depth": [3, 5, 7],  
        "min_child_weight": [1, 3, 5],  
        "subsample": [0.7, 0.8, 0.9],  
        "colsample_bytree": [0.7, 0.8, 0.9],  
        "n_estimators": [100, 200, 300]  
    },  
    "AdaBoostRegressor": {  
        "n_estimators": [50, 100, 200 , 400],  
        "learning_rate": [0.01, 0.05, 0.1, 0.5]  
    },  
    "LinearRegression": {},  
    "GradientBoostingRegressor": {
```

Best model: GradientBoostingRegressor with StandardScaler  
Best Train R2 score: 0.8295773398116387  
Best Validation R2 score: 0.8096783047739387  
Best params: {'learning\_rate': 0.05, 'max\_depth': 3, 'n\_estimators': 300}

# model validation with test data

```
best_model = GradientBoostingRegressor (learning_rate =0.05, max_depth =3, n_estimators =300)
best_scaler = StandardScaler()

X_train_scaled = best_scaler.fit_transform(X_tr_selected)
X_test_scaled = best_scaler.transform(X_test_selected)

best_model.fit(X_train_scaled, y_tr)

y_test_pred = best_model.predict(X_test_scaled)

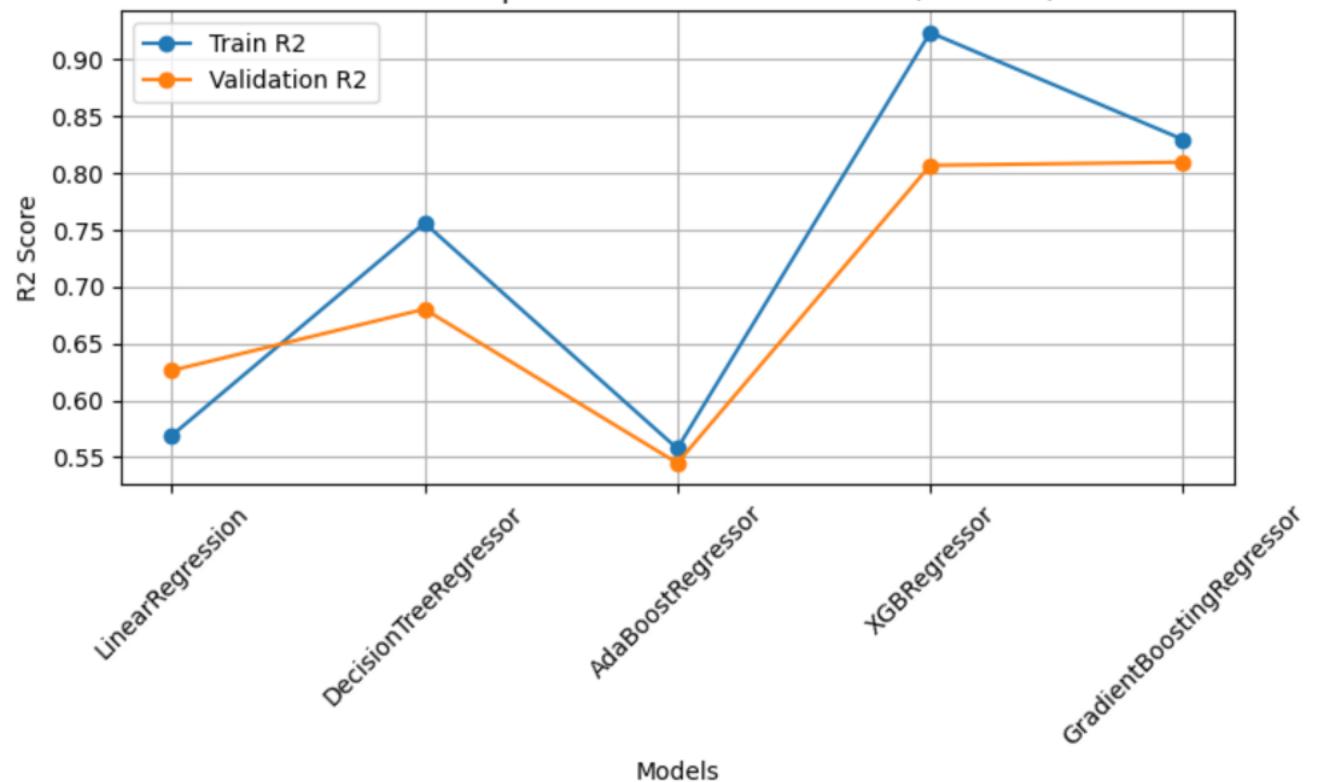
test_r2 = r2_score(Y_test, y_test_pred)
test_mse = mean_squared_error(Y_test, y_test_pred)

print(f"R2 Score on test set: {test_r2}")
print(f"MSE on test set: {test_mse}")
```

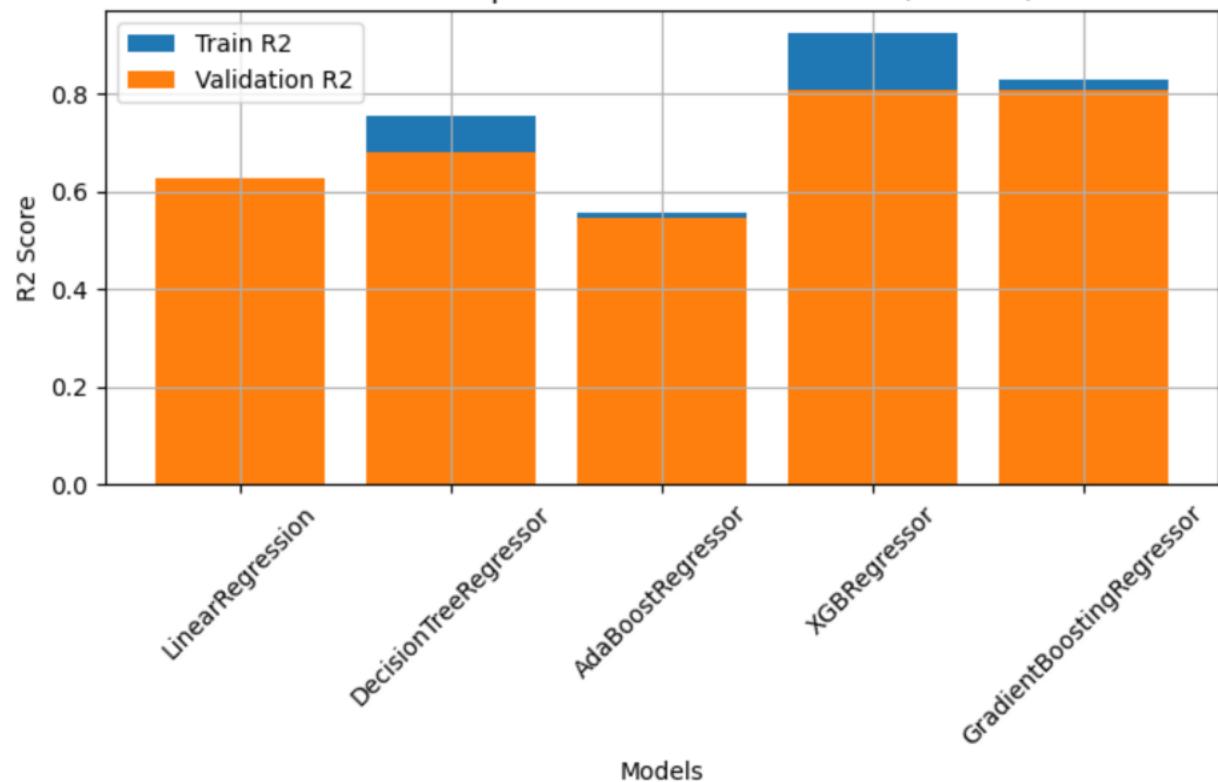
R2 Score on test set: 0.7914219500715205

MSE on test set: 6307.529155506736

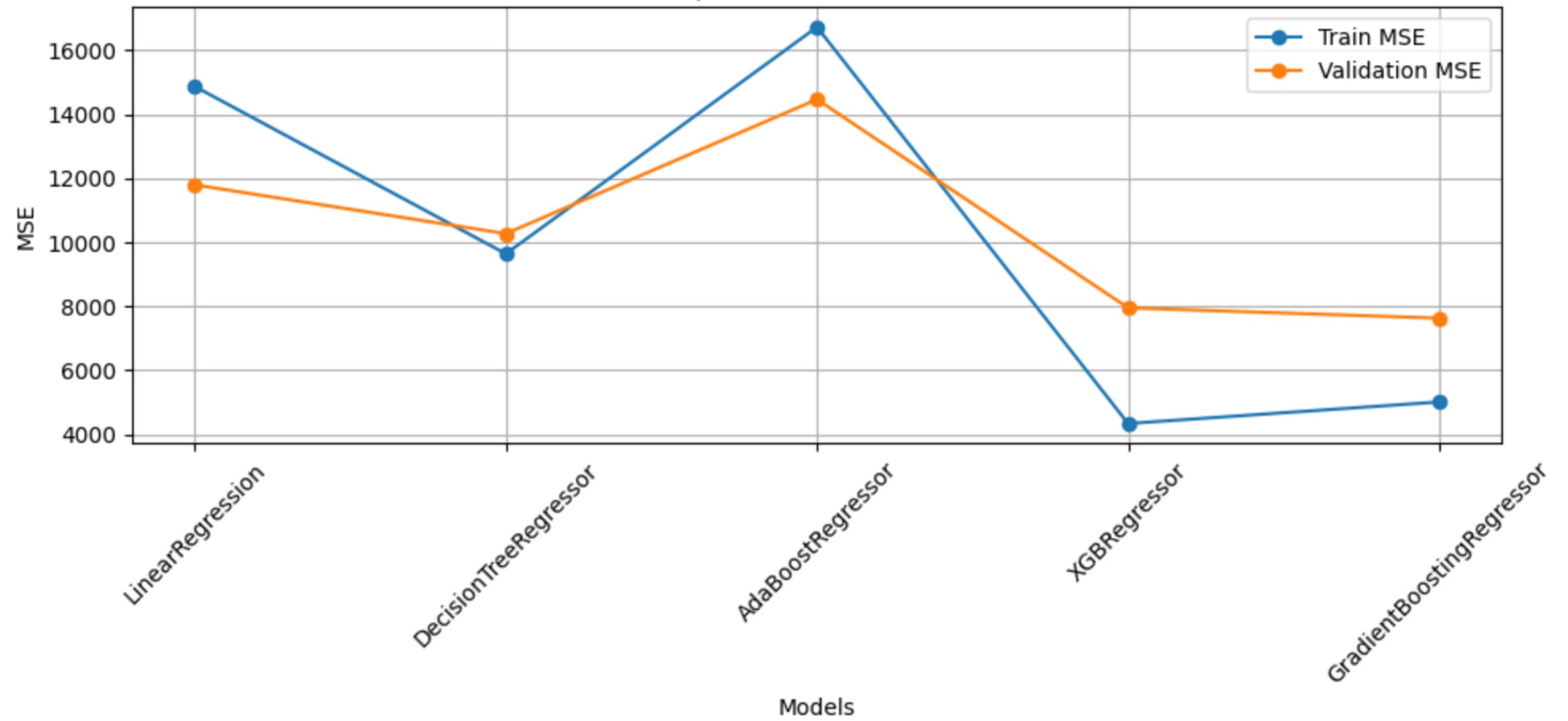
R2 Score Comparison for Different Models (Line Plot)



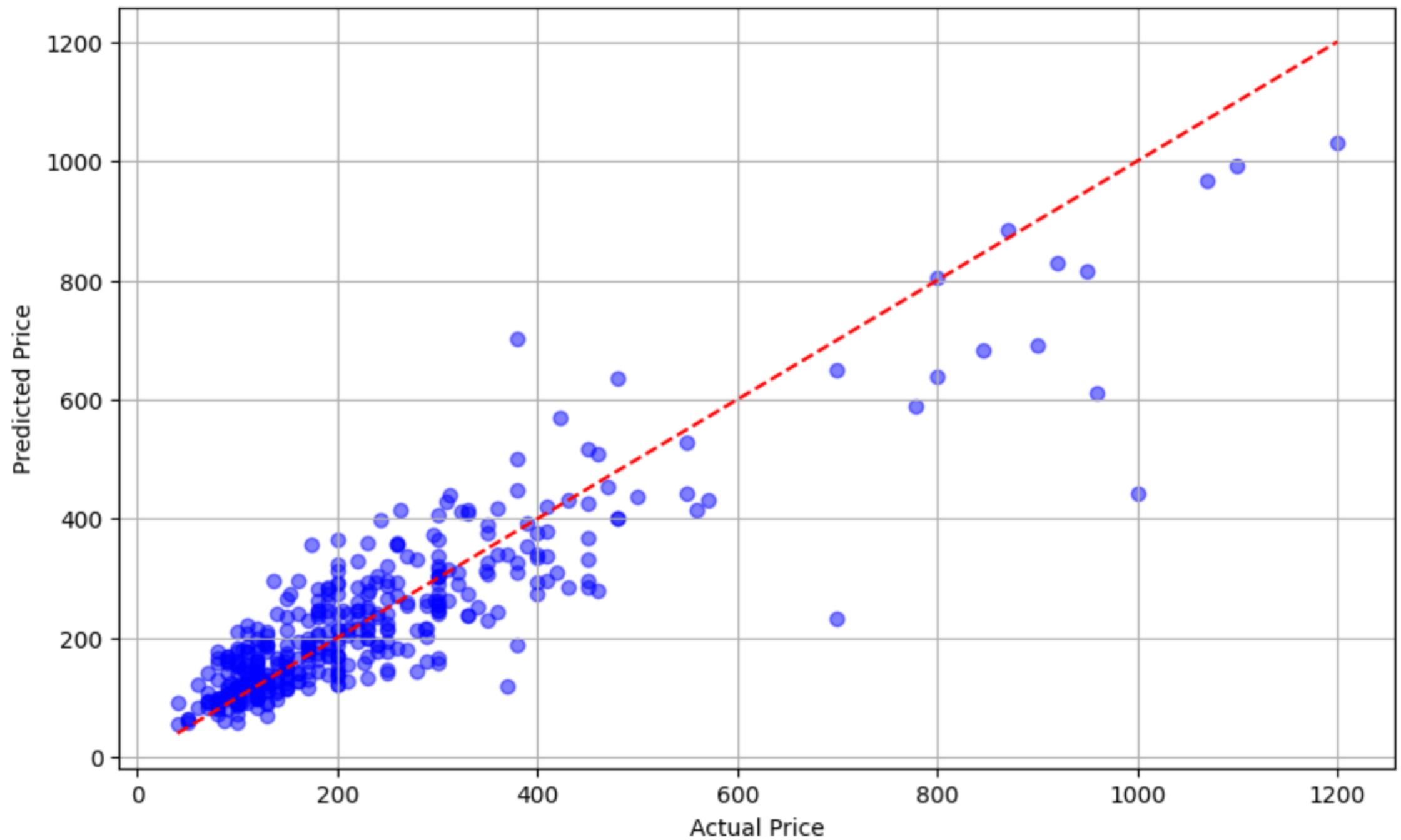
R2 Score Comparison for Different Models (Bar Plot)



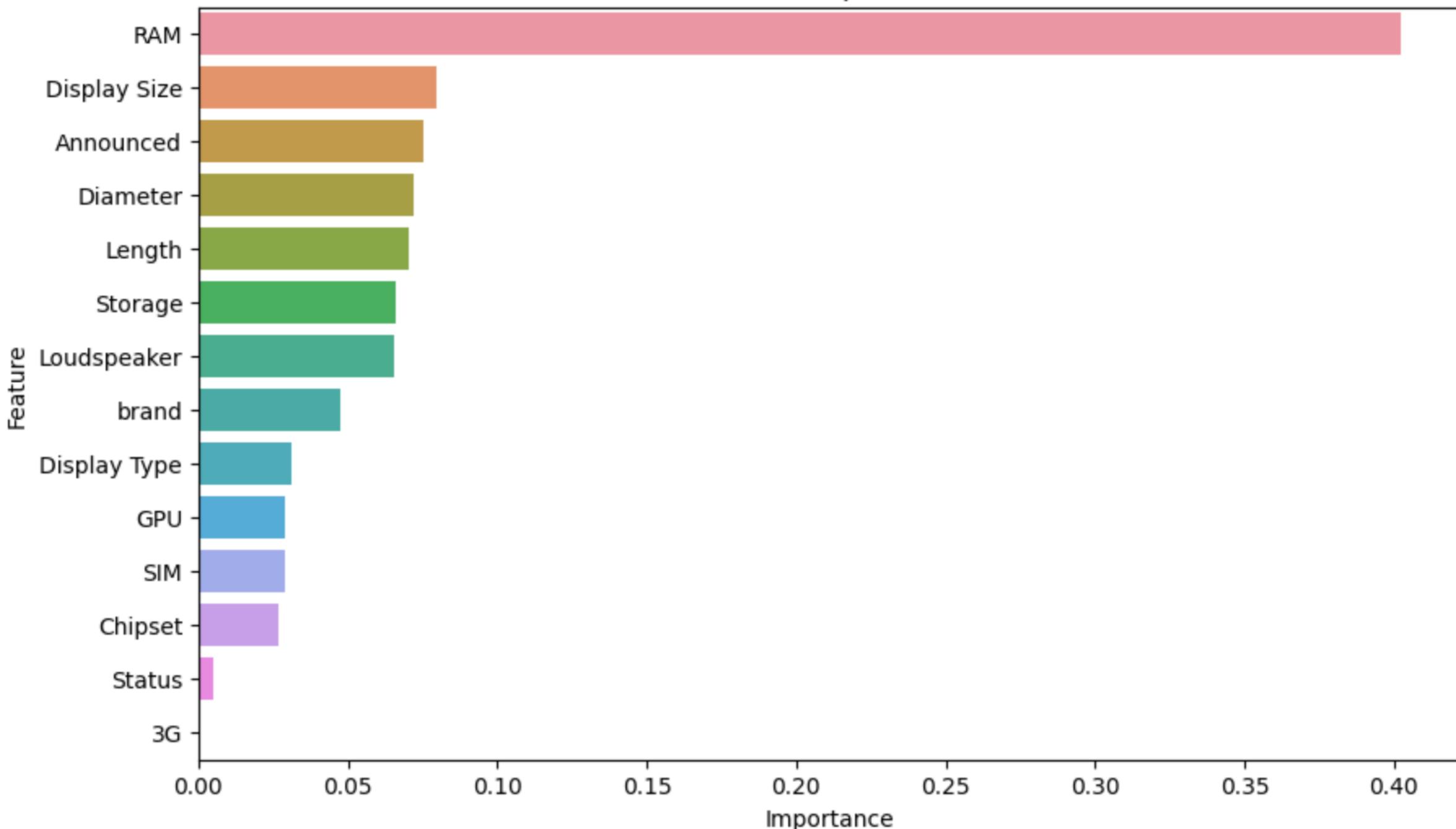
## MSE Comparison for Different Models

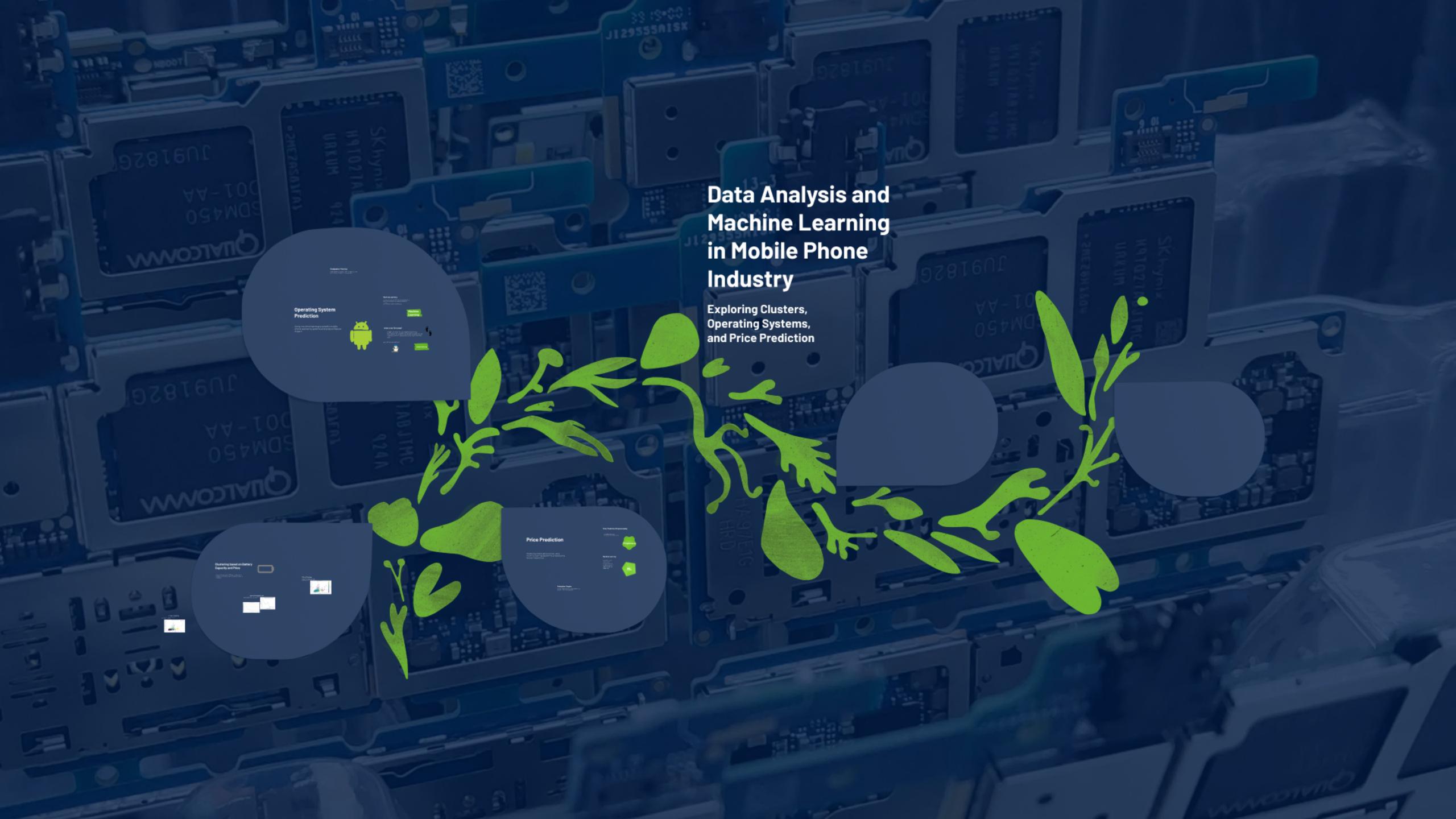


### Actual vs Predicted Prices for Test Data



### Feature Importance





# Data Analysis and Machine Learning in Mobile Phone Industry

Exploring Clusters,  
Operating Systems,  
and Price Prediction



Price Prediction

