

```
# =====
# ROUND 2 – PREPROCESSING & VISUALIZATION
# =====
```

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.cluster import KMeans
from google.colab import files
```

```
plt.style.use("seaborn-v0_8")
```

```
# =====
# 1. LOAD RAW DATA
# =====
```

```
df = pd.read_csv("/content/snapdeal_raw_data.csv")
print("Raw dataset shape:", df.shape)
df.head()
```

```
Raw dataset shape: (100, 9)
```

```
{
  "summary": {
    "name": "df",
    "rows": 100,
    "fields": [
      {
        "column": "product_name",
        "properties": {
          "dtype": "string",
          "num_unique_values": 89,
          "samples": [
            "OLIVE OPS Sports Edition Smartwatch | Heart Rate | Steps | Wireless Charging Smartwatch (Black-Grey Strap, 1.73 Inch Big Sunlight Proof Display)",
            "Ramsons - MIDNIGHT Eau De Parfum Perfume For Men Long Lasting Premium Perfume 40ml (Pack of 1)",
            "Zebronics 10000 -mAh 22.5W Li-Polymer Power Bank"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "category",
        "properties": {
          "dtype": "category",
          "num_unique_values": 5,
          "samples": [
            "power_bank",
            "men_shoes",
            "smart_watch"
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "price",
        "properties": {
          "dtype": "number",
          "std": 291,
          "min": 129,
          "max": 1248,
          "num_unique_values": 77,
          "samples": [
            318,
            673,
            488
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "mrp",
        "properties": {
          "dtype": "number",
          "std": 1324,
          "min": 185,
          "max": 6999,
          "num_unique_values": 34,
          "samples": [
            349,
            799,
            2299
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "rating",
        "properties": {
          "dtype": "number",
          "std": 0.39630629146429436,
          "min": 3.5,
          "max":

```

```

5.0,\n          \"num_unique_values\": 16,\n          \"samples\": [\n
4.2,\n          3.9,\n          4.4\n          ],\n
\"semantic_type\": \"\", \n          \"description\": \"\"\n          }\n
n      },\n      {\n          \"column\": \"review_count\", \n
\"properties\": {\n          \"dtype\": \"number\", \n          \"std\":
1212.122512071145,\n          \"min\": 1.0,\n          \"max\": 7954.0,\n
\"num_unique_values\": 68,\n          \"samples\": [\n          272.0,\n
134.0,\n          7954.0\n          ],\n          \"semantic_type\":
\"\", \n          \"description\": \"\"\n          }\n      },\n      {\n
\"column\": \"product_url\", \n          \"properties\": {\n
\"dtype\": \"string\", \n          \"num_unique_values\": 100,\n
\"samples\": [\n          \"https://www.snapdeal.com/product/hotstyle-
gray-mens-sports-running/623981372905#breadcrumbSearch:men%20shoes\", \n
\"https://www.snapdeal.com/product/olive-ops-sports-edition-
smartwatch/655462266845#breadcrumbSearch:smart%20watch\", \n
\"https://www.snapdeal.com/product/st-john-one-men-army/659679519339#b
readcrumbSearch:perfume\"\n          ],\n          \"semantic_type\": \"\", \n
\"description\": \"\"\n          }\n      },\n      {\n          \"column\":
\"seller_name\", \n          \"properties\": {\n          \"dtype\":
\"category\", \n          \"num_unique_values\": 24,\n
\"samples\": [\n          \"P00JA ENTERPRISES\", \n          \"HELIOS
LIFESTYLE PVT LTD\", \n          \"Varni Enterprise\"\n          ],\n
\"semantic_type\": \"\", \n          \"description\": \"\"\n          }\n
n      },\n      {\n          \"column\": \"seller_rating\", \n
\"properties\": {\n          \"dtype\": \"number\", \n          \"std\":
0.3104070941563096,\n          \"min\": 3.5,\n          \"max\": 5.0,\n
\"num_unique_values\": 9,\n          \"samples\": [\n          4.3,\n
3.9,\n          3.5\n          ],\n          \"semantic_type\": \"\", \n
\"description\": \"\"\n          }\n      }\n      ]\n
n}], \"type\": \"dataframe\", \"variable_name\": \"df\"}

```

```

# Missing values

```

```

df.isnull().sum().sort_values(ascending=False)

```

```

plt.figure(figsize=(10,4))
sns.heatmap(df.isnull(), cbar=False)
plt.title("Missing Value Heatmap")
plt.show()

```

```

# Category distribution

```

```

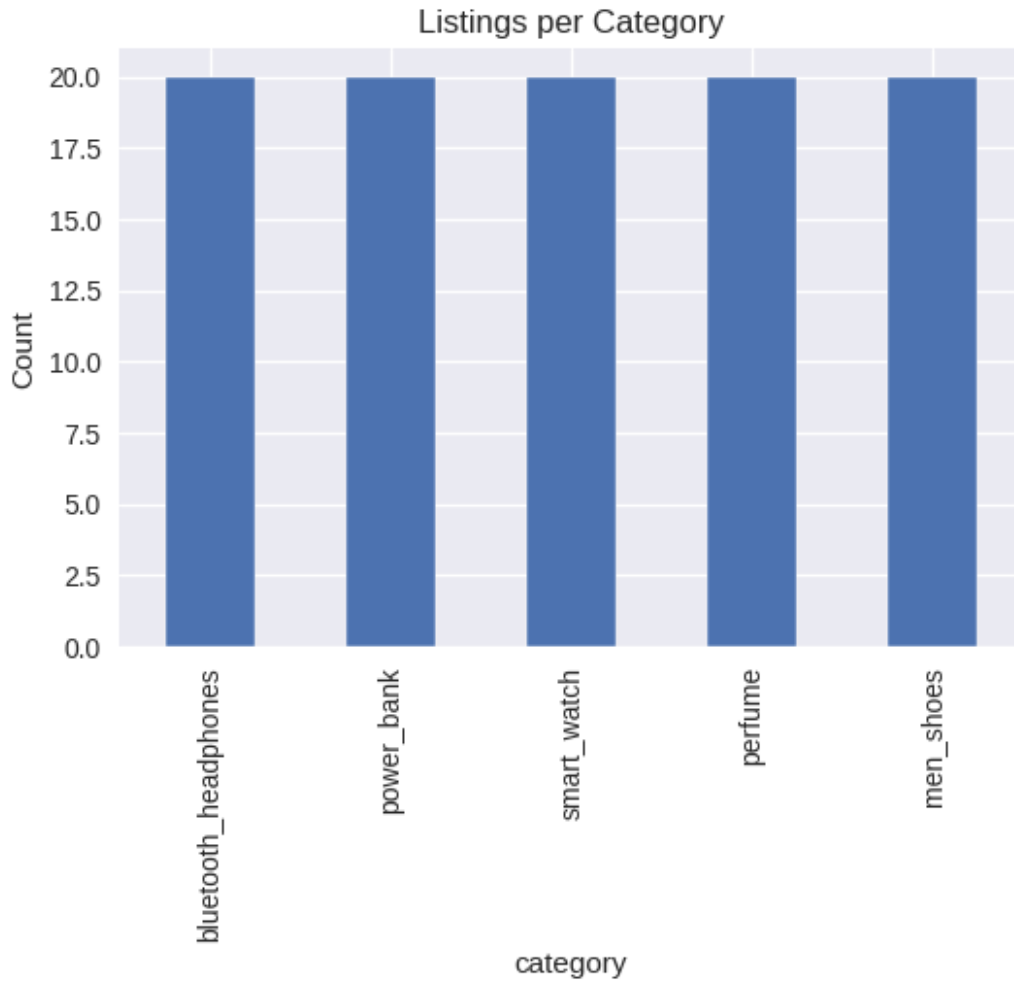
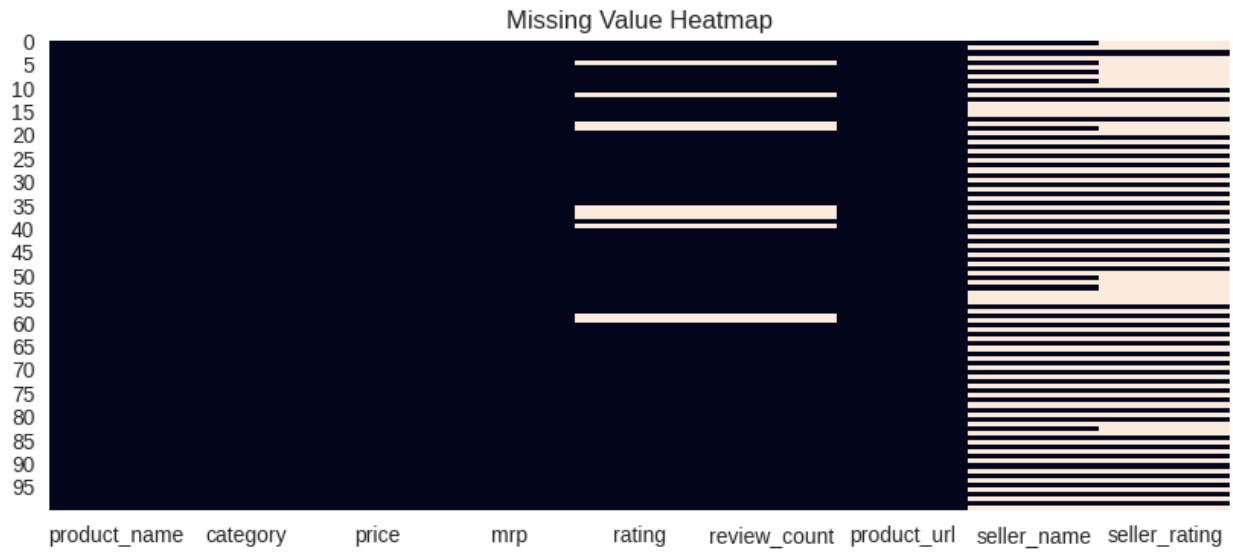
plt.figure(figsize=(6,4))
df[\"category\"].value_counts().plot(kind=\"bar\")
plt.title("Listings per Category")
plt.ylabel("Count")
plt.show()

```

```

df.dtypes

```



product_name	object
category	object

```
price          int64
mrp            int64
rating         float64
review_count   float64
product_url    object
seller_name    object
seller_rating  float64
dtype: object
```

```
before = df.shape[0]
```

```
df["price"] = pd.to_numeric(df["price"], errors="coerce")
df["mrp"] = pd.to_numeric(df["mrp"], errors="coerce")
df["rating"] = pd.to_numeric(df["rating"], errors="coerce")
df["review_count"] = pd.to_numeric(df["review_count"],
errors="coerce").fillna(0)
df["seller_rating"] = pd.to_numeric(df["seller_rating"],
errors="coerce")
```

```
df = df.dropna(subset=["price", "mrp", "rating"])
```

```
after = df.shape[0]
print("Rows before cleaning:", before)
print("Rows after cleaning:", after)
```

```
plt.figure(figsize=(6,4))
sns.boxplot(x=df["price"])
plt.title("Price Distribution After Cleaning")
plt.show()
```

```
Rows before cleaning: 100
```

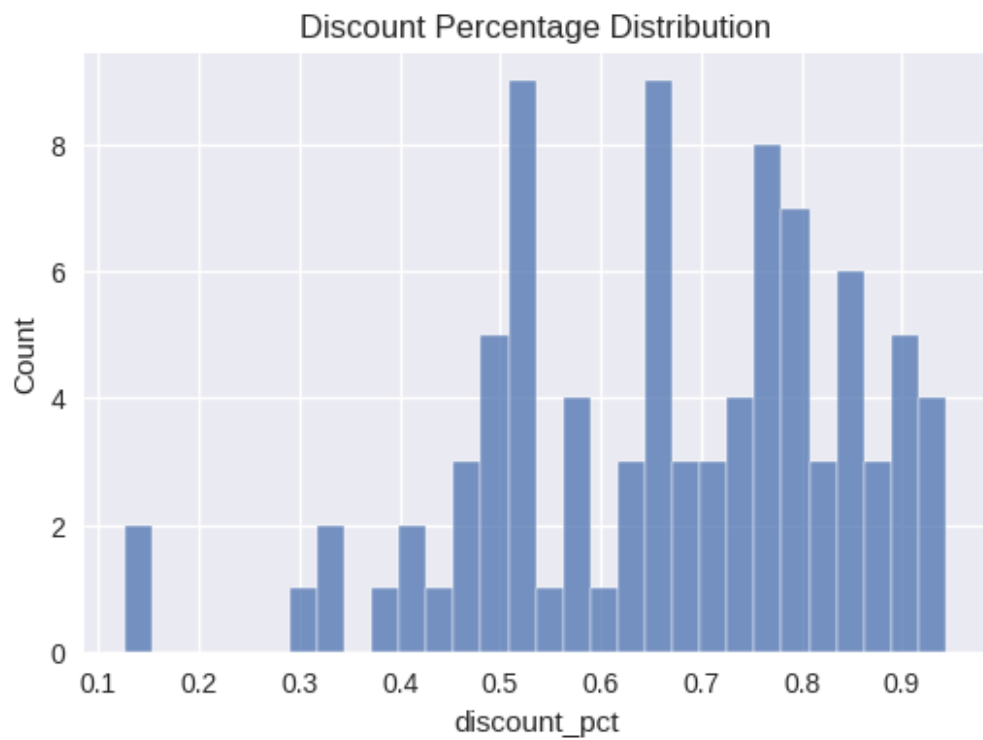
```
Rows after cleaning: 90
```

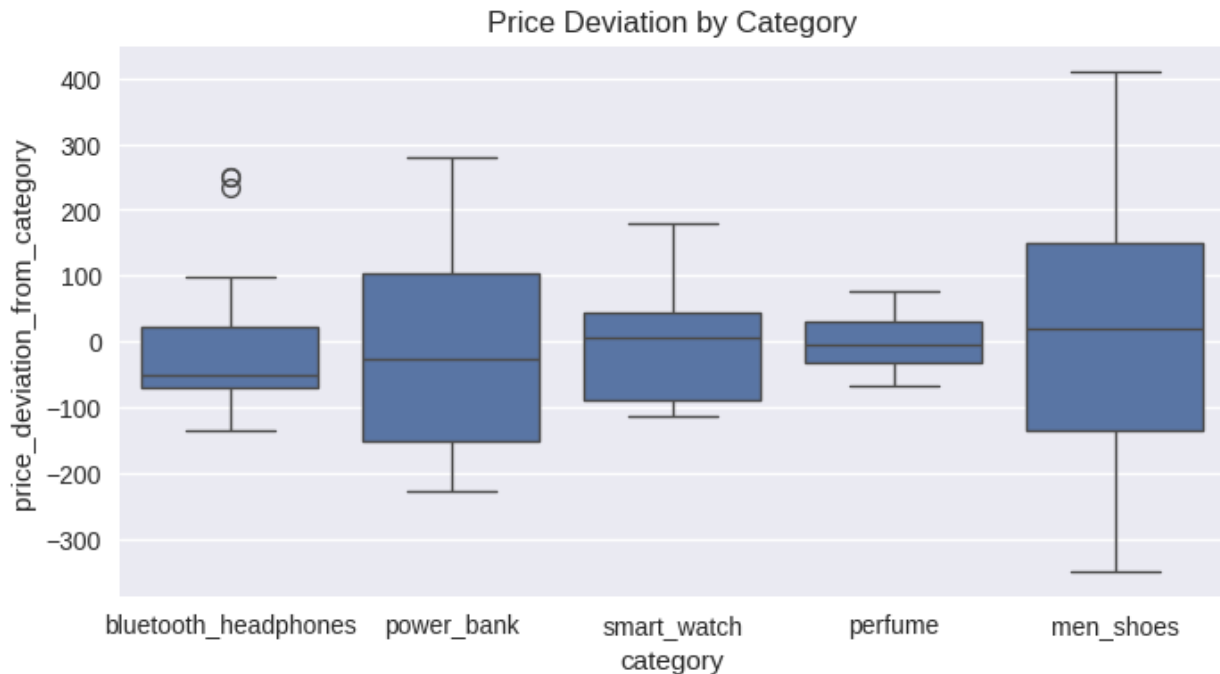


```
plt.figure(figsize=(6,4))
sns.histplot(df["discount_pct"], bins=30)
plt.title("Discount Percentage Distribution")
plt.show()

plt.figure(figsize=(6,4))
sns.scatterplot(data=df, x="price", y="rating", hue="category")
plt.title("Price vs Rating")
plt.show()

plt.figure(figsize=(8,4))
sns.boxplot(data=df, x="category", y="price_deviation_from_category")
plt.title("Price Deviation by Category")
plt.show()
```





```

df["extreme_discount_flag"] = (df["discount_pct"] > 0.6).astype(int)
df["suspicious_trust_flag"] = (df["rating_review_ratio"] <
1.2).astype(int)
df["aggressive_listing_flag"] = (
    (df["generic_word_count"] >= 3) | (df["uppercase_ratio"] > 0.25)
).astype(int)

df["risk_score"] = (
    df["extreme_discount_flag"] +
    df["suspicious_trust_flag"] +
    df["aggressive_listing_flag"]
)

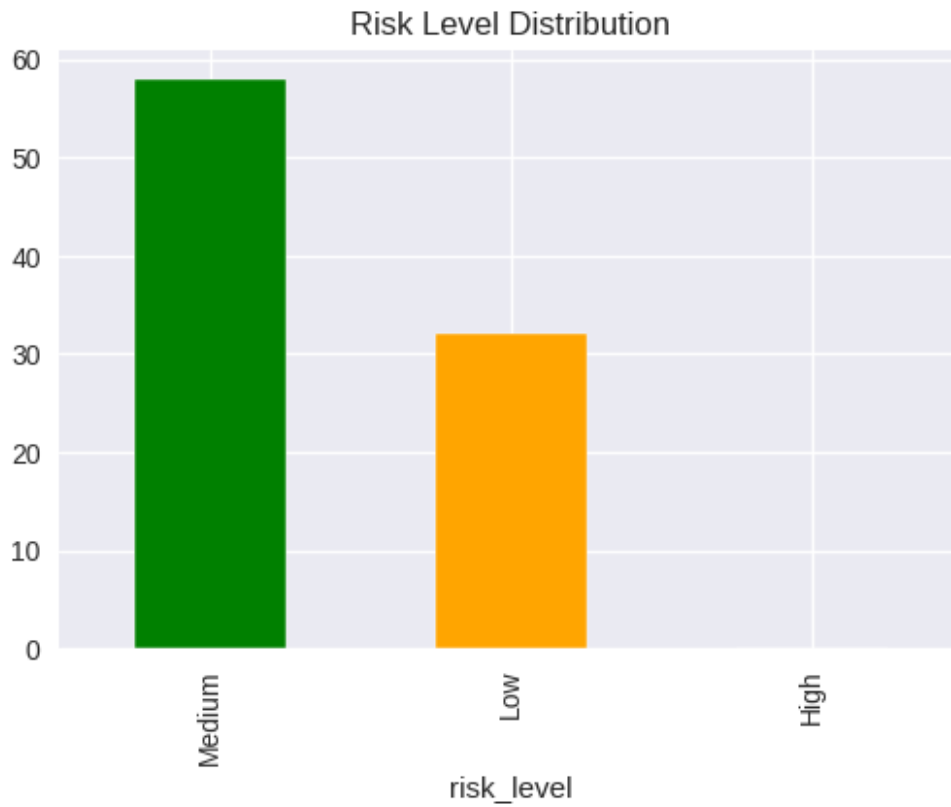
df["risk_level"] = pd.cut(
    df["risk_score"],
    bins=[-1, 1, 3, 10],
    labels=["Low", "Medium", "High"]
)

plt.figure(figsize=(6, 4))
df["risk_level"].value_counts().plot(kind="bar",
color=["green", "orange", "red"])
plt.title("Risk Level Distribution")
plt.show()

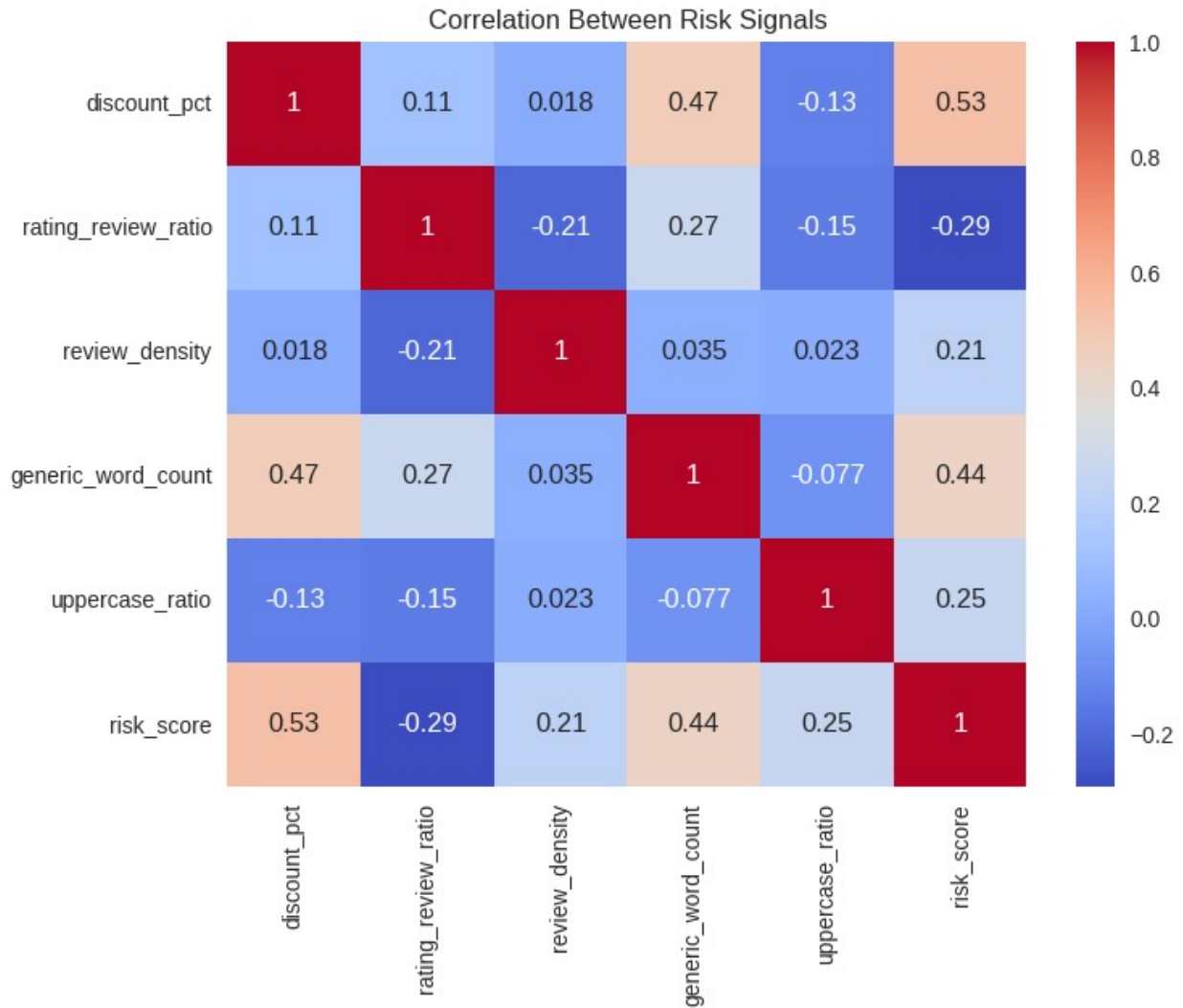
risk_features = [
    "discount_pct", "rating_review_ratio", "review_density",
    "generic_word_count", "uppercase_ratio", "risk_score"
]

```

```
plt.figure(figsize=(8,6))
sns.heatmap(df[risk_features].corr(), annot=True, cmap="coolwarm")
plt.title("Correlation Between Risk Signals")
plt.show()
```







```
cluster_features = df[
    ["price", "discount_pct", "rating_review_ratio",
     "generic_word_count", "uppercase_ratio"]
].fillna(0)

X = StandardScaler().fit_transform(cluster_features)
df["seller_cluster_id"] = KMeans(n_clusters=5,
random_state=42).fit_predict(X)

scaler = MinMaxScaler()
scale_cols = [
    "price", "discount_pct", "rating_review_ratio",
    "review_density", "price_deviation_from_category"
]

df[scale_cols] = scaler.fit_transform(df[scale_cols])
```

```
final_file = "round2_processed_dataset.csv"
df.to_csv(final_file, index=False)
files.download(final_file)
```

```
print("📊 FINAL DATASET SHAPE:", df.shape)
```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.Javascript object>
```

```
📊 FINAL DATASET SHAPE: (90, 23)
```