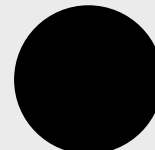
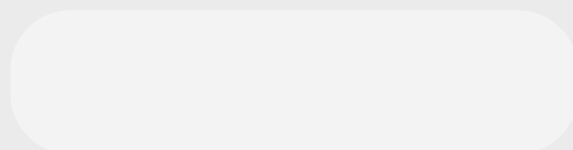
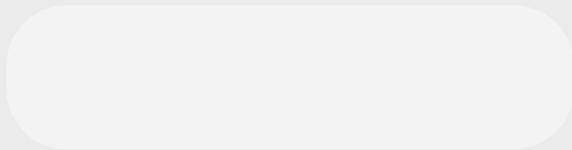
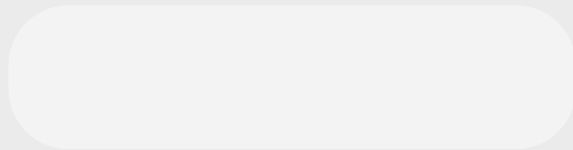
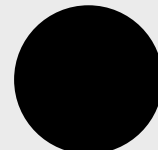
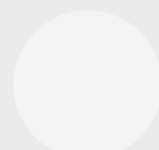
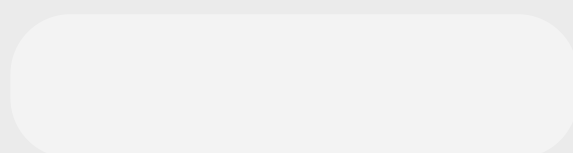
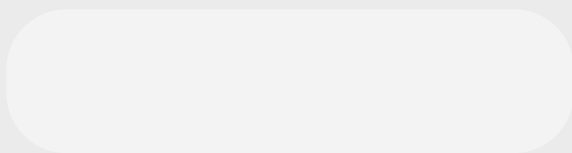
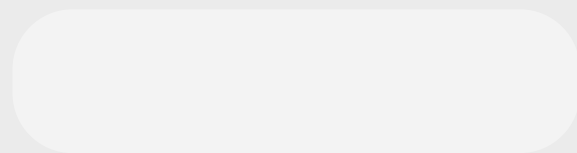


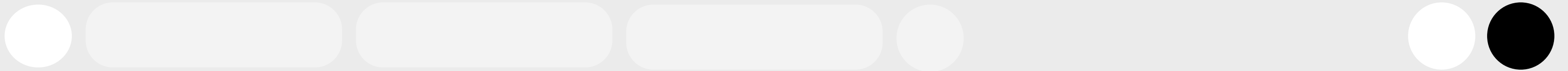
Ask anything



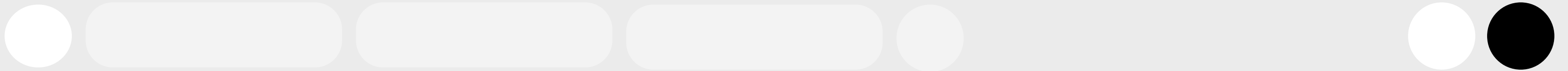
yo chat give me a list of laptops with these specifications



bro no something within a reasonable price range



ok but which ones better in terms of what i want



we've got  
a *better*  
solution  
than  
chatgpt

A black silhouette of a person sitting at a desk with a laptop, positioned on the left side of the image. The person is facing right, and the laptop is open in front of them. The background is split vertically: white on the left and dark gray on the right.

TECH GENIE  
AT YOUR SERVICE

# PC PART PICKER 30000

# AN OVERVIEW

Ever wish you had a  
genie who could  
instantly tell you the  
price of your **dream**  
**laptop** and **show you the**  
**best match** in the  
market?

We deliver three core ML functions:

- **Descriptive:** K-Means Clustering to segment computers based on price, RAM, and other specs. PCA is a way of showing the clusters of the K-Means
- **Predictive:** LightGBM Regression to estimate **prices** from user inputs
- **Prescriptive:** K-Nearest Neighbors (KNN) to recommend similar listings with ranked similarity

# DATA COLLECTION & PREP





- **CSV file:** 8,064 marketplace listings (rows) x **135 raw Spanish-language columns**
- Encoded in **UTF-8-SIG** with mixed metrics, units, and labels; **.CSV file with 135 columns**.
- **No scraping or APIs; data ingested directly via `pandas.read_csv`.**

# RAW DATA

# ● CLEANED DATA

**1. Dropped Duplicates →**  
`df.duplicated().sum()`

**2. Standardized Column Names with custom slugify function**

- → removed accents, lowercase, dropped stopwords (e.g., **Pantalla\_Tamaño** → **pantalla\_tamano**)

**3. Dropped Unnamed Columns:**

- `df.drop(columns=['unnamed_0'])`
- Full null or >70% null columns

#### 4. Price Normalization

- Parsed "Precio\_Rango" (e.g., "1.026,53 € – 2.287,17 €") into:
- `precio_min`, `precio_max`, and `precio_mean`
- Dropped original string after parsing

#### 5. Numerical Extraction

- Created functions to **extract float from strings (e.g., RAM, CPU speed)**
- Remove thousands separators.
- Apply `apply_cleaning_to_column()` across many dirty fields

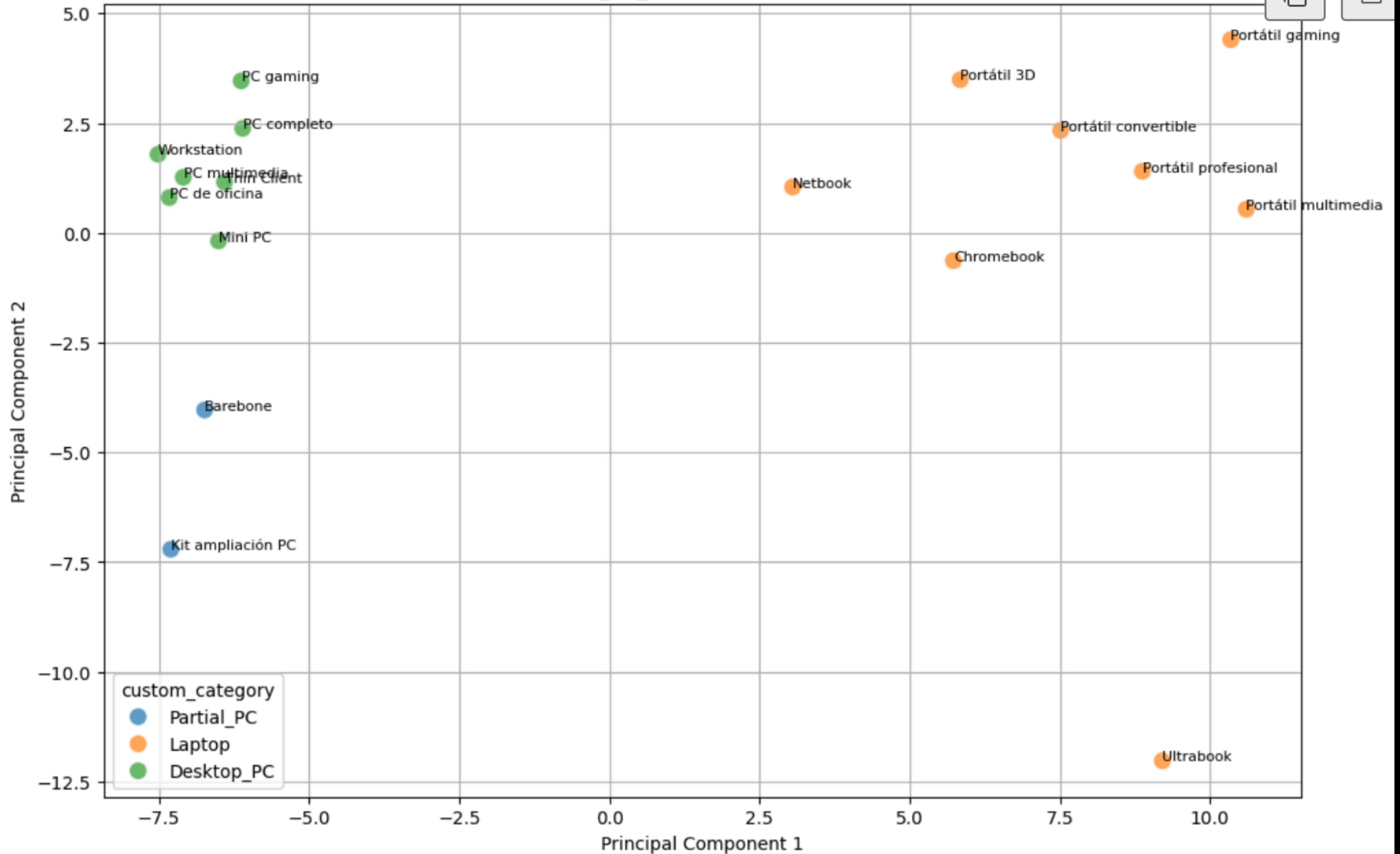
#### 6. Standardized Screen Resolution

- Used regex to convert inconsistent resolution strings to "WIDTHXHEIGHT"
- E.g., "4K (3.840 x 2.160)" → "3840x2160"

#### 7. Offers Cleaning

- Convert strings like "200 ofertas" to 200.0 (float) for numeric ops.

PCA of tipo\_de\_producto (2 Components)



# HANDLING MISSING DATA

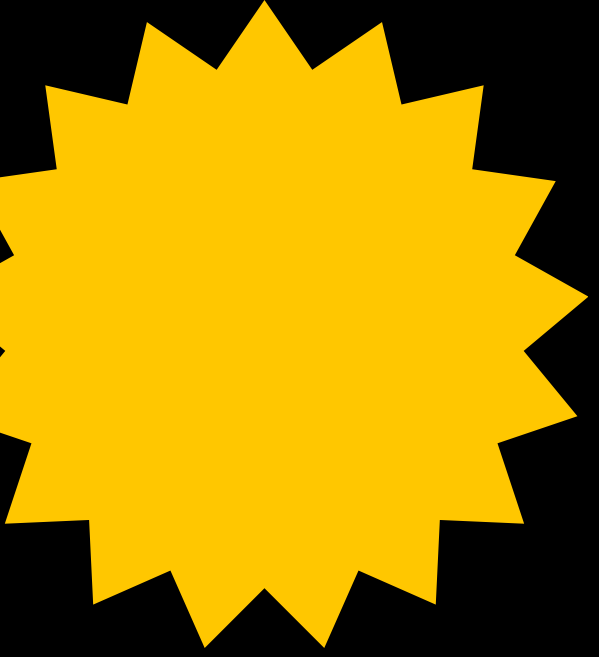
Used `df.isnull().sum()` and `missingno`  
heatmaps

• **Aware of Missing-Not-At-Random (MNAR)** issues (e.g., screens missing in desktops). To solve, we handled it by isolating category-specific structures and then:

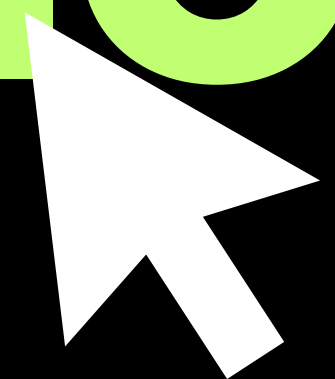



## STRATEGY?

- | 70% missing: dropped
- **30–70%:** conditional imputation or dropped
- **<30%:** imputed by product category using mean/mode



# FEATURE ENGINEERING & SELECTION





# HALL OF FAME

## FEATURE ENGINEERING & SELECTION

1

### ONE-HOT ENCODING

for low-cardinality  
categorical fields.

2

### ORDINAL ENCODING

for ordered  
features *like*  
processor  
generation

3


### PCA & CORRELATION ANALYSIS

**PCA** to retain  
features explaining  
90%+ variance  
+ Removed highly  
correlated  
variables  
(Pearsons).

4

### FINAL MATRIX

Final feature  
matrix optimized  
for model  
performance &  
interpretability.



**CATEGORICAL  
HANDLING**

**MEAN PRICE:** Extracted from raw price range string

```
def process_price_range(price_str)
```

**Volume (cm<sup>3</sup>)**= height x width x depth

**Category Mapping:** Mapped devices to **English Classes (Ultrabook, Tower, All-in-One)**



**FEATURE ENGINEERING**

# MODEL TRAINING & VALIDATION

TECHNICAL APPROACH FOR SOLVING  
FUNCTIONALITIES



# DESCRIPTIVE

## K-MEANS CLUSTERING

- **Business value:** clustering different product segments.
- **Inputs:** features in the `df_engineered.csv` dataset.
- **Output:** k=2 clusters.
- **Evaluation:** PCA dimensionality reduction for visual inspection of clusters, plotted clusters and `tipo_de_producto` feature against same principal component axes.

# PREDICTIVE

## LIGHTGBM REGRESSION

- **Business value:**
  - Price prediction of computer given user selected specifications.
  - Feature importance on price of computer.
- **Inputs:** features in `df_engineered_desktop_pc.csv` or `df_engineered_laptop.csv`.
- **Output:** predicted price and feature importance.
- **Target:** `precio_mean` feature.
- **Validation:** cross-validation hyperparameter tuning.
- **Evaluation:** RMSE  $\approx$  520 EUR,  $R^2 \approx$  0.81.

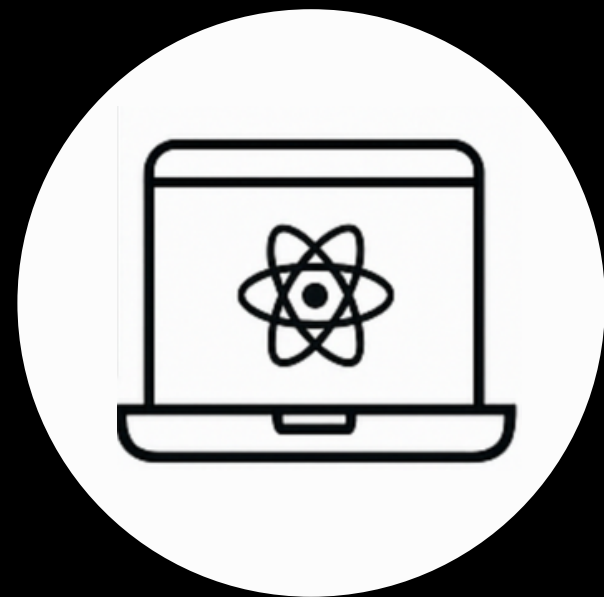
# PRESCRIPTIVE

## K-NEAREST NEIGHBORS

- **Business value:**
  - K-similar product offers to user selected specifications.
- **Inputs:** features in `df_engineered_desktop_pc.csv` or `df_engineered_laptop.csv`.
- **Output:** k=5 neighbours.
- **Evaluation:** KNN regressor predicted versus target price scatter plot and residual plot as baseline performance metrics.

# APP ARCHITECTURE & DEPLOYMENT

# FRONTEND STACK



REACT-BASED UI

# BACKEND APIS

*PYTHON:  
SCIKIT-LEARN  
PANDAS,  
MATPLOTLIB*

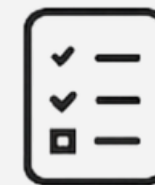
FOR THE EDA &  
TRAINING



Google Cloud



API HOSTING: DEPLOYED VIA GIT  
HUB → **GOOGLE CLOUD RUN  
FUNCTIONS**



ML MODELS: LIGHTGBM, KMEANS,  
KNN IN PYTHON (**JOBLIB  
SERIALIZED**)



MODEL STORAGE: **GOOGLE CLOUD  
STORAGE**



**CI/CD AUTOMATION:** GITHUB  
ACTIONS - TRIGGERED ON PUSH TO  
MAIN FOR THE MODELS



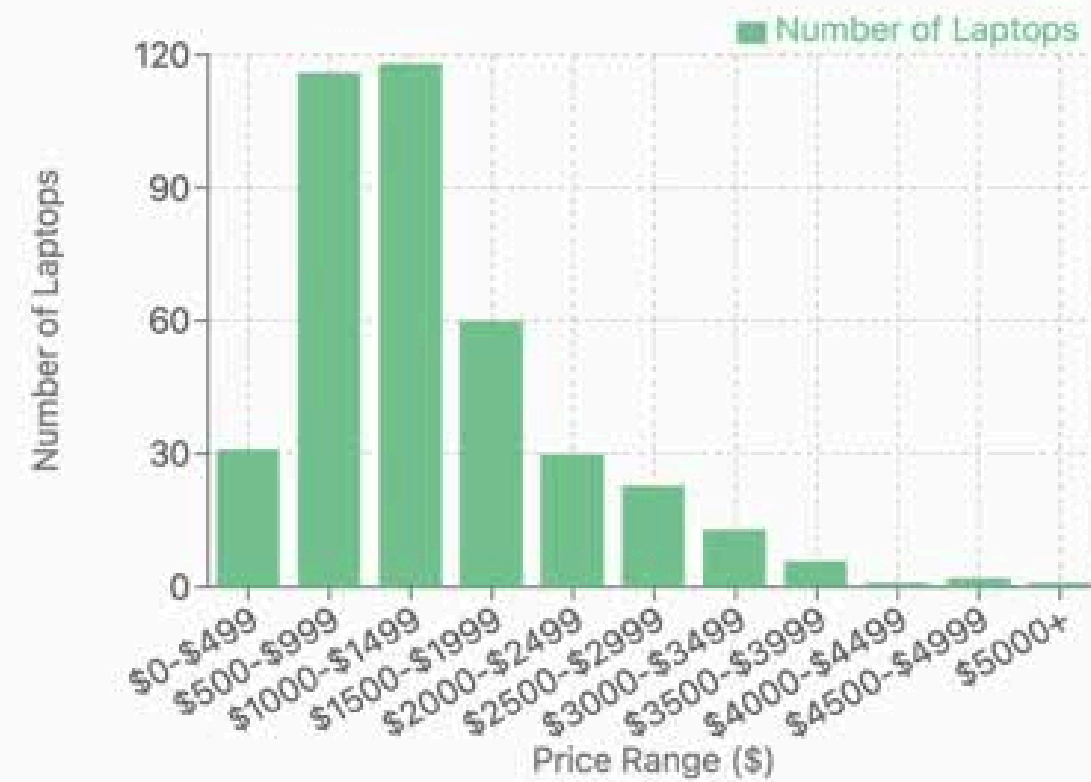
# **LIVE DEMONSTRATION**



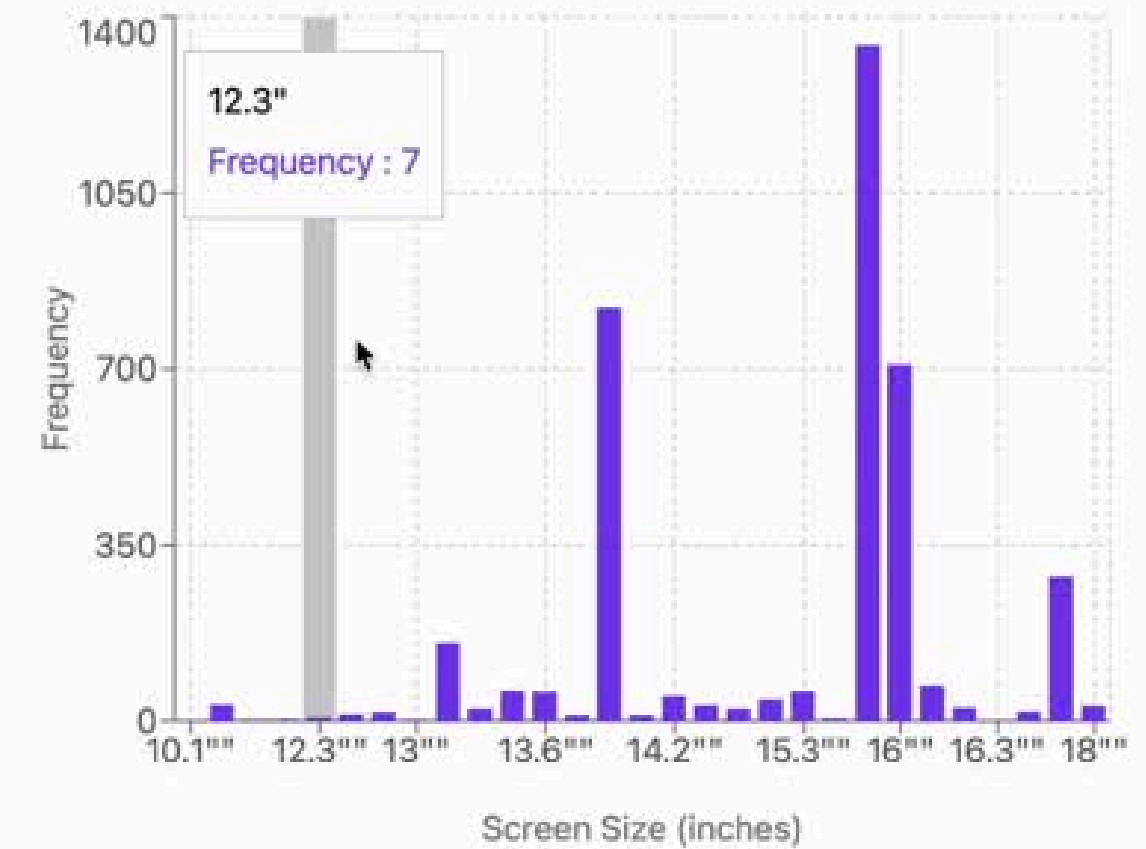
Segmentation

Prediction

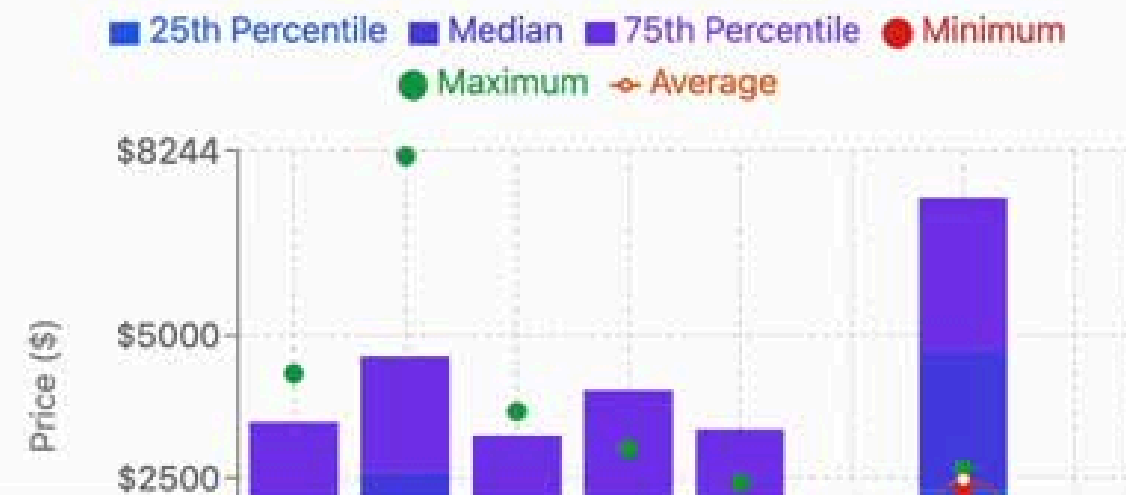
## Overall Price Distribution



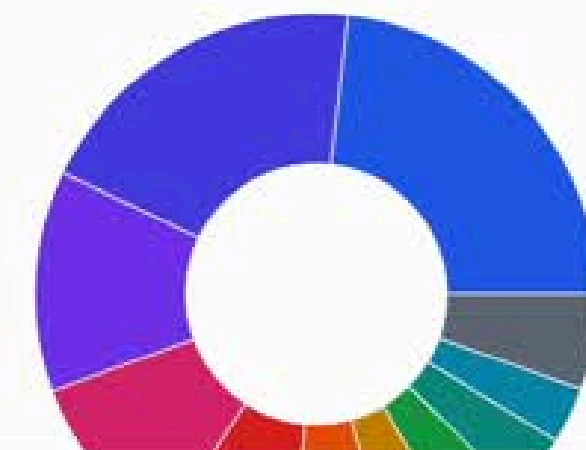
## Screen Size Distribution of Laptops



## Price Distribution by Top 7 Product Types



## Top 10 Popular Brands



**IMPROVEMENTS  
& NEXT STEPS**

# GENIE'S NEXT EVOLUTION

- **Live data integration** via APIs to keep listings up to date
- **Prediction confidence intervals** to show uncertainty
- **User-based personalization** using historical preferences
- **Model retraining** via feedback log ingestion
- **Multilingual toggle** to support Spanish/English UIs
- **Domain expansion** to peripherals, monitors, GPUs
- **Feature Feedback** to allow for constant improvements of model & the display of processed data.

The background features a white field with several dark grey, organic, and somewhat jagged shapes. These shapes are positioned around the central text, with some extending towards the edges of the frame. The overall effect is a modern, minimalist design.

**THANK YOU**