

✓ Video Game Sales Prediction - Machine Learning Lab

✓ Introduction and Setup

Welcome to this machine learning lab where we'll build a model to predict whether a video game will be a "hit" based on its characteristics and sales data. This notebook will guide you through the entire process, from data loading to model evaluation and optimization.

Learning objectives:

1. Learn to preprocess and explore a real-world dataset
2. Build and evaluate a decision tree classifier
3. Optimize a model through hyperparameter tuning
4. Interpret model results and feature importance

```
1 #install libraries if necessary
2
3 # Import the necessary libraries
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Machine learning libraries
10 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.preprocessing import LabelEncoder
13 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
14 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, RocCurveDisplay
15
```


```
1 # Set random seed for reproducibility
2 np.random.seed(42)
```

```
1 # Configure visualizations
2 plt.style.use('seaborn-v0_8-whitegrid')
3 sns.set_palette("colorblind")
```

```
1 # Display settings for better output
2 pd.set_option('display.max_columns', None)
3 pd.set_option('display.width', 1000)
```

✓ Download the Dataset

```
1 # You can run this cell to download the dataset directly, or upload it manually!
2 import requests
3
4 url = 'https://www.kaggle.com/datasets/gregorut/videogamesales/download'
5 response = requests.get(url)
6
7 with open('videogamesales.zip', 'wb') as f:
8     f.write(response.content)
9
10 print("Dataset downloaded successfully.")
11
```

 Dataset downloaded successfully.

✓ Load the Dataset

```
1 # Load the dataset
2 df = pd.read_csv('vgsales.csv')
```

```

3
4 # Let's take a look at the first few rows of the dataset
5 print("First 5 rows of the dataset:")
6 print(df.head())
7

```

```

First 5 rows of the dataset:

```

	Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
0	1	Wii Sports	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74
1	2	Super Mario Bros.	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24
2	3	Mario Kart	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	35.82
3	4	Wii Sports Resort	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	33.00
4	5	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37

Dataset Information

```

1 # Get basic information about the dataset
2 print("\nDataset basic information:")
3 print(df.info())
4
5 # Get descriptive statistics
6 print("\nDescriptive statistics:")
7 print(df.describe())

```

```

Dataset basic information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16598 entries, 0 to 16597
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Rank                 16598 non-null  int64
1   Name                 16598 non-null  object
2   Platform             16598 non-null  object
3   Year                 16327 non-null  float64
4   Genre                16598 non-null  object
5   Publisher            16540 non-null  object
6   NA_Sales              16598 non-null  float64
7   EU_Sales              16598 non-null  float64
8   JP_Sales              16598 non-null  float64
9   Other_Sales          16598 non-null  float64
10  Global_Sales          16598 non-null  float64
dtypes: float64(6), int64(1), object(4)
memory usage: 1.4+ MB
None

Descriptive statistics:

```

	Rank	Year	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
count	16598.000000	16327.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000
mean	8300.605254	2006.406443	0.264667	0.146652	0.077782	0.048063	0.537441
std	4791.853933	5.828981	0.816683	0.505351	0.309291	0.188588	1.555028
min	1.000000	1980.000000	0.000000	0.000000	0.000000	0.000000	0.010000
25%	4151.250000	2003.000000	0.000000	0.000000	0.000000	0.000000	0.060000
50%	8300.500000	2007.000000	0.080000	0.020000	0.000000	0.010000	0.170000
75%	12449.750000	2010.000000	0.240000	0.110000	0.040000	0.040000	0.470000
max	16600.000000	2020.000000	41.490000	29.020000	10.220000	10.570000	82.740000

```

1 # Cell 5: Check for missing values
2 print("\nMissing values per column:")
3 print(df.isnull().sum())

```

```

Missing values per column:
Rank      0
Name      0
Platform  0
Year     271
Genre     0
Publisher  58
NA_Sales  0
EU_Sales  0
JP_Sales  0
Other_Sales  0
Global_Sales  0
dtype: int64

```

```

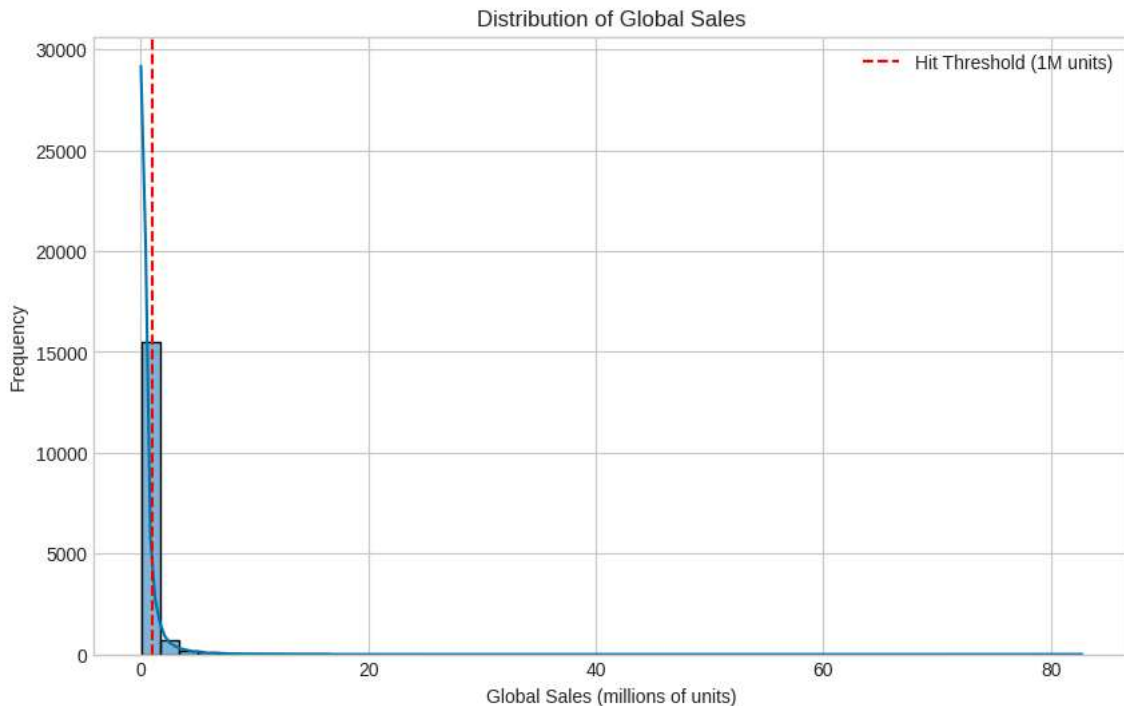
1 # Cell 6: Data Visualization - Global Sales Distribution
2 # =====
3 # Visualizing the distribution of global sales

```

```

3 # visualize the distribution of global sales
4 plt.figure(figsize=(10, 6))
5 sns.histplot(df['Global_Sales'], bins=50, kde=True)
6 plt.title('Distribution of Global Sales')
7 plt.xlabel('Global Sales (millions of units)')
8 plt.ylabel('Frequency')
9 plt.axvline(x=1, color='red', linestyle='--', label='Hit Threshold (1M units)')
10 plt.legend()
11 plt.show()

```



```

1 # Cell 7: Create Target Variable
2 # =====
3 # TASK: Create a binary target variable for "hit" games
4 # A game is considered a hit if it sold more than 1 million units (Global_Sales > 1)
5 # YOUR CODE HERE
6
7 df['Hit'] = (df['Global_Sales'] > 1).astype(int)
8
9 df[['Name', 'Global_Sales', 'Hit']].head()
10

```



	Name	Global_Sales	Hit
0	Wii Sports	82.74	1
1	Super Mario Bros.	40.24	1
2	Mario Kart Wii	35.82	1
3	Wii Sports Resort	33.00	1
4	Pokemon Red/Pokemon Blue	31.37	1

```

1 # Cell 8: Analyze Target Distribution
2 # =====
3 # Let's see the proportion of hits in our dataset
4 # YOUR CODE HERE
5
6 hit_counts = df['Hit'].value_counts()
7 hit_proportion = df['Hit'].value_counts(normalize=True)
8
9 print("Hit Counts:\n", hit_counts)
10 print("\nHit Proportion:\n", hit_proportion)
11
12 plt.figure(figsize=(6, 4))
13 sns.barplot(x=hit_counts.index, y=hit_counts.values)
14 plt.title('Distribution of Hit Games')

```

```

15 plt.xlabel('Hit (1 = Yes, 0 = No)')
16 plt.ylabel('Number of Games')
17 plt.xticks([0, 1], ['Not Hit', 'Hit'])
18 plt.show()
19

```

↗ Hit Counts:

```

Hit
0    14544
1     2054
Name: count, dtype: int64

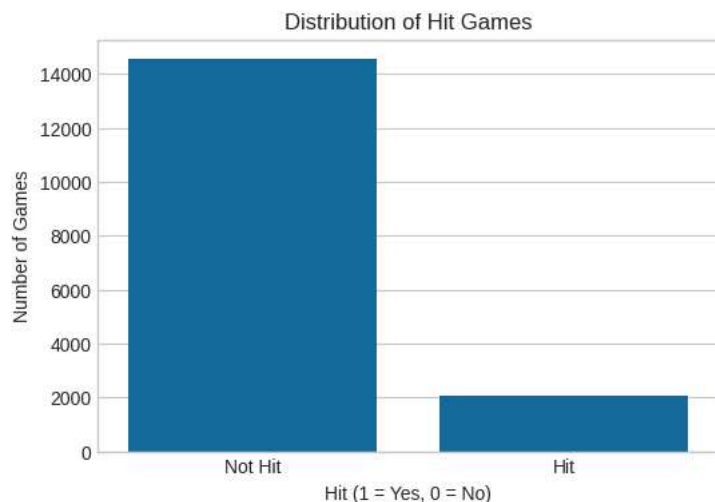
```

Hit Proportion:

```

Hit
0    0.87625
1    0.12375
Name: proportion, dtype: float64

```



```

1 # Cell 9: Drop Non-Informative Columns
2 # =====
3 # TASK: Drop non-informative columns
4 # Think about which columns won't help with prediction
5 # YOUR CODE HERE
6
7 # Drop columns that won't help with prediction
8 df.drop(columns=['Rank', 'Name'], inplace=True)
9
10 # Preview updated dataframe
11 df.head()
12

```

↗

	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	Hit
0	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74	1
1	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24	1
2	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	35.82	1
3	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	33.00	1
4	GB	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37	1

📊

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

```

1 # Cell 10: Missing Value Analysis
2 # =====
3 # Examine the 'Year' column which might have missing values
4
5 # Check for missing values in the 'Year' column
6 missing_year_count = df['Year'].isnull().sum()
7 missing_year_percentage = (missing_year_count / len(df)) * 100
8
9 print(f"Missing 'Year' values: {missing_year_count}")
10 print(f"Percentage missing: {missing_year_percentage:.2f}%")
11

```

Missing 'Year' values: 271
Percentage missing: 1.63%

```
1 # Cell 11: Handle Missing Values
2 # =====
3 # TASK: Handle missing values
4 # Option 1: Drop rows with missing values
5 # YOUR CODE HERE
6
7 # Option 2: Fill missing values with median or mean
8 # YOUR CODE HERE
9 df_clean['Year'] = df_clean['Year'].fillna(df_clean['Year'].median())
10
11 # Option 1: Drop rows with missing values
12 df_drop = df.dropna(subset=['Year'])
13
14 # Option 2: Fill missing values with the median year
15 df_fill = df.copy()
16 df_fill['Year'] = df_fill['Year'].fillna(df_fill['Year'].median())
17
18 df_clean = df_fill
19
20 df_clean[['Year']].isnull().sum()
21
```

0

Year
0

dtype: int64

```
1 # Cell 12: Categorical Variable Analysis
2 # =====
3 # Let's identify categorical columns
4 categorical_columns = df_clean.select_dtypes(include=['object']).columns.tolist()
5 print("\nCategorical columns:", categorical_columns)
```

Categorical columns: ['Platform', 'Genre', 'Publisher']

```
1 # Cell 13: Encode Categorical Variables
2 # =====
3 # TASK: Encode categorical variables using LabelEncoder
4 # Label Encoder transforms categorical variables into numerical ones
5 # YOUR CODE HERE
6
7 from sklearn.preprocessing import LabelEncoder
8
9 le = LabelEncoder()
10
11 for col in ['Platform', 'Genre', 'Publisher']:
12     df_clean[col] = le.fit_transform(df_clean[col].astype(str)) # Convert to string in case of NaNs
13
14 df_clean[['Platform', 'Genre', 'Publisher']].head()
15
```

Platform Genre Publisher

0	26	10	359
1	11	4	359
2	26	6	359
3	26	10	359
4	5	7	359

```
1 # Cell 14: Feature Engineering (Optional)
2 # =====
3 # BONUS TASK: Feature Engineering
4 # Creating new features might improve model performance
5 # Example: Total regional sales besides global
```

```

6 # YOUR CODE HERE
7
8 df_clean['Total_Regional_Sales'] = df_clean['NA_Sales'] + df_clean['EU_Sales'] + df_clean['JP_Sales'] + df_clean['Other_Sales']
9
10 df_clean[['Total_Regional_Sales', 'Global_Sales']].head()
11

```

	Total_Regional_Sales	Global_Sales
0	82.74	82.74
1	40.24	40.24
2	35.83	35.82
3	33.00	33.00
4	31.38	31.37

```

1 # Cell 15: Explore Processed Dataset
2 # =====
3 # Let's look at the processed dataset
4 # YOUR CODE HERE
5
6
7 print("Shape of dataset:", df_clean.shape)
8 print("\nData types:\n", df_clean.dtypes)
9 print("\nFirst 5 rows:\n", df_clean.head())
10 print("\nSummary statistics:\n", df_clean.describe())
11

```

Shape of dataset: (16598, 11)

Data types:

Platform	int64
Year	float64
Genre	int64
Publisher	int64
NA_Sales	float64
EU_Sales	float64
JP_Sales	float64
Other_Sales	float64
Global_Sales	float64
Hit	int64
Total_Regional_Sales	float64

dtype: object

First 5 rows:

	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	Hit	Total_Regional_Sales
0	26	2006.0	10	359	41.49	29.02	3.77	8.46	82.74	1	82.74
1	11	1985.0	4	359	29.08	3.58	6.81	0.77	40.24	1	40.24
2	26	2008.0	6	359	15.85	12.88	3.79	3.31	35.82	1	35.83
3	26	2009.0	10	359	15.75	11.01	3.28	2.96	33.00	1	33.00
4	5	1996.0	7	359	11.27	8.89	10.22	1.00	31.37	1	31.38

Summary statistics:

	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	Hit
count	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000	16598.000000
mean	15.797988	2006.416134	4.928124	294.599169	0.264667	0.146652	0.077782	0.048063	0.537441	0.537441
std	8.392298	5.781686	3.762015	178.082372	0.816683	0.505351	0.309291	0.188588	1.555028	1.555028
min	0.000000	1980.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.010000	0.010000
25%	7.000000	2003.000000	1.000000	137.000000	0.000000	0.000000	0.000000	0.000000	0.060000	0.060000
50%	16.000000	2007.000000	5.000000	323.000000	0.080000	0.020000	0.000000	0.010000	0.170000	0.170000
75%	21.000000	2010.000000	8.000000	461.000000	0.240000	0.110000	0.040000	0.040000	0.470000	0.470000
max	30.000000	2020.000000	11.000000	578.000000	41.490000	29.020000	10.220000	10.570000	82.740000	82.740000

```

1 # Cell 16: Split Features and Target
2 # =====
3 # TASK: Split the data into features (X) and target (y)
4 # YOUR CODE HERE
5
6 X = df_clean.drop(columns=['Hit']) # Features
7 y = df_clean['Hit']               # Target
8

```

```

1 # Cell 17: Train-Test Split
2 # =====
3 # TASK: Split the data into training and testing sets (80/20 split)

```

```

4
5 # X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
6 # Print the shapes to confirm the split
7
8 from sklearn.model_selection import train_test_split
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
11
12 print("Training features shape:", X_train.shape)
13 print("Testing features shape:", X_test.shape)
14 print("Training target shape:", y_train.shape)
15 print("Testing target shape:", y_test.shape)
16

```

```

↗ Training features shape: (13278, 10)
Testing features shape: (3320, 10)
Training target shape: (13278,)
Testing target shape: (3320,)

```

```

1 # Cell 18: Train Initial Model
2 # =====
3 # TASK: Train a Decision Tree classifier with default parameters
4 # YOUR CODE HERE
5
6 from sklearn.tree import DecisionTreeClassifier
7
8 clf = DecisionTreeClassifier(random_state=42)
9 clf.fit(X_train, y_train)
10

```

```

↗ DecisionTreeClassifier ⓘ ?
DecisionTreeClassifier(random_state=42)

```

```

1 # Cell 19: Make Predictions
2 # =====
3 # TASK: Make predictions on the test set
4 # YOUR CODE HERE
5
6 # Probability of being a hit
7
8 y_pred = clf.predict(X_test)
9
10 y_proba = clf.predict_proba(X_test)[: , 1] # Extract probability for class 1 (hit)

```

```

1 # Cell 20: Calculate Evaluation Metrics
2 # =====
3 # TASK: Calculate evaluation metrics
4 # YOUR CODE HERE
5
6
7 from sklearn.metrics import accuracy_score, precision_score, recall_score,
8 f1_score, roc_auc_score
9
10 accuracy = accuracy_score(y_test, y_pred)
11 precision = precision_score(y_test, y_pred, zero_division=0)
12 recall = recall_score(y_test, y_pred, zero_division=0)
13 f1 = f1_score(y_test, y_pred, zero_division=0)
14 roc_auc = roc_auc_score(y_test, y_proba)
15
16 print(f"Accuracy: {accuracy:.4f}")
17 print(f"Precision: {precision:.4f}")
18 print(f"Recall: {recall:.4f}")
19 print(f"F1 Score: {f1:.4f}")
20 print(f"ROC AUC Score: {roc_auc:.4f}")
21

```

```

↗ Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1 Score: 1.0000
ROC AUC Score: 1.0000

```

```

1 # Cell 21: Confusion Matrix Visualization

```

```

2 # =====
3 # TASK: Visualize the confusion matrix
4 # YOUR CODE HERE
5
6 from sklearn.metrics import ConfusionMatrixDisplay
7
8 ConfusionMatrixDisplay.from_estimator(clf, X_test, y_test, display_labels=['Not
Hit', 'Hit'], cmap='Blues')
9 plt.title('Confusion Matrix')
10 plt.show()
11

```

