

****Proyecto 2 - Entrega 6**

Integrantes:

- **Diederich Solis** (22952)
 - **Gabriel Paz** (221087)
-

Uso del conjunto de datos de entrenamiento y prueba

En esta sección se utiliza el mismo conjunto de datos `train.csv` empleado en entregas anteriores. La separación de datos en entrenamiento y prueba se mantiene constante para garantizar la validez de las comparaciones entre modelos.

Se cargan los datos utilizando `pandas` y se verifica su correcta estructura para preparar el preprocesamiento necesario para los modelos SVM.

```
# Librerías necesarias
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler

# Cargar el dataset
data = pd.read_csv('train.csv')

# Mostrar primeras filas
data.head()

{"type": "dataframe", "variable_name": "data"}
```

Exploración y Transformación de Datos

Se realiza un análisis exploratorio inicial para comprender la estructura del dataset, identificar valores faltantes y analizar las variables relevantes.

Posteriormente, se aplican transformaciones necesarias como imputación de datos, codificación de variables categóricas y escalado de variables numéricas para preparar el dataset para su uso en modelos de Máquinas de Vectores de Soporte (SVM).

```
# Información general del dataset
data.info()

# Verificar valores nulos
missing_values = data.isnull().sum()
missing_values[missing_values > 0]
```

```
# Rellenar valores nulos para simplicidad (puedes cambiar si quieres ser más sofisticado)
```

```
data = data.fillna(data.median(numeric_only=True))
```

```
# Eliminar columnas no numéricas o altamente categóricas para este experimento
```

```
data = data.select_dtypes(include=[np.number])
```

```
# Confirmar limpieza
```

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1460 entries, 0 to 1459
```

```
Data columns (total 81 columns):
```

#	Column	Non-Null Count	Dtype
0	Id	1460 non-null	int64
1	MSSubClass	1460 non-null	int64
2	MSZoning	1460 non-null	object
3	LotFrontage	1201 non-null	float64
4	LotArea	1460 non-null	int64
5	Street	1460 non-null	object
6	Alley	91 non-null	object
7	LotShape	1460 non-null	object
8	LandContour	1460 non-null	object
9	Utilities	1460 non-null	object
10	LotConfig	1460 non-null	object
11	LandSlope	1460 non-null	object
12	Neighborhood	1460 non-null	object
13	Condition1	1460 non-null	object
14	Condition2	1460 non-null	object
15	BldgType	1460 non-null	object
16	HouseStyle	1460 non-null	object
17	OverallQual	1460 non-null	int64
18	OverallCond	1460 non-null	int64
19	YearBuilt	1460 non-null	int64
20	YearRemodAdd	1460 non-null	int64
21	RoofStyle	1460 non-null	object
22	RoofMatl	1460 non-null	object
23	Exterior1st	1460 non-null	object
24	Exterior2nd	1460 non-null	object
25	MasVnrType	588 non-null	object
26	MasVnrArea	1452 non-null	float64
27	ExterQual	1460 non-null	object
28	ExterCond	1460 non-null	object
29	Foundation	1460 non-null	object
30	BsmtQual	1423 non-null	object
31	BsmtCond	1423 non-null	object

32	BsmtExposure	1422	non-null	object
33	BsmtFinType1	1423	non-null	object
34	BsmtFinSF1	1460	non-null	int64
35	BsmtFinType2	1422	non-null	object
36	BsmtFinSF2	1460	non-null	int64
37	BsmtUnfSF	1460	non-null	int64
38	TotalBsmtSF	1460	non-null	int64
39	Heating	1460	non-null	object
40	HeatingQC	1460	non-null	object
41	CentralAir	1460	non-null	object
42	Electrical	1459	non-null	object
43	1stFlrSF	1460	non-null	int64
44	2ndFlrSF	1460	non-null	int64
45	LowQualFinSF	1460	non-null	int64
46	GrLivArea	1460	non-null	int64
47	BsmtFullBath	1460	non-null	int64
48	BsmtHalfBath	1460	non-null	int64
49	FullBath	1460	non-null	int64
50	HalfBath	1460	non-null	int64
51	BedroomAbvGr	1460	non-null	int64
52	KitchenAbvGr	1460	non-null	int64
53	KitchenQual	1460	non-null	object
54	TotRmsAbvGrd	1460	non-null	int64
55	Functional	1460	non-null	object
56	Fireplaces	1460	non-null	int64
57	FireplaceQu	770	non-null	object
58	GarageType	1379	non-null	object
59	GarageYrBlt	1379	non-null	float64
60	GarageFinish	1379	non-null	object
61	GarageCars	1460	non-null	int64
62	GarageArea	1460	non-null	int64
63	GarageQual	1379	non-null	object
64	GarageCond	1379	non-null	object
65	PavedDrive	1460	non-null	object
66	WoodDeckSF	1460	non-null	int64
67	OpenPorchSF	1460	non-null	int64
68	EnclosedPorch	1460	non-null	int64
69	3SsnPorch	1460	non-null	int64
70	ScreenPorch	1460	non-null	int64
71	PoolArea	1460	non-null	int64
72	PoolQC	7	non-null	object
73	Fence	281	non-null	object
74	MiscFeature	54	non-null	object
75	MiscVal	1460	non-null	int64
76	MoSold	1460	non-null	int64
77	YrSold	1460	non-null	int64
78	SaleType	1460	non-null	object
79	SaleCondition	1460	non-null	object
80	SalePrice	1460	non-null	int64

```

dtypes: float64(3), int64(35), object(43)
memory usage: 924.0+ KB
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 38 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Id                    1460 non-null   int64
1   MSSubClass            1460 non-null   int64
2   LotFrontage           1460 non-null   float64
3   LotArea               1460 non-null   int64
4   OverallQual           1460 non-null   int64
5   OverallCond           1460 non-null   int64
6   YearBuilt             1460 non-null   int64
7   YearRemodAdd          1460 non-null   int64
8   MasVnrArea           1460 non-null   float64
9   BsmtFinSF1            1460 non-null   int64
10  BsmtFinSF2            1460 non-null   int64
11  BsmtUnfSF             1460 non-null   int64
12  TotalBsmtSF           1460 non-null   int64
13  1stFlrSF              1460 non-null   int64
14  2ndFlrSF              1460 non-null   int64
15  LowQualFinSF          1460 non-null   int64
16  GrLivArea             1460 non-null   int64
17  BsmtFullBath          1460 non-null   int64
18  BsmtHalfBath          1460 non-null   int64
19  FullBath              1460 non-null   int64
20  HalfBath              1460 non-null   int64
21  BedroomAbvGr          1460 non-null   int64
22  KitchenAbvGr          1460 non-null   int64
23  TotRmsAbvGrd          1460 non-null   int64
24  Fireplaces            1460 non-null   int64
25  GarageYrBlt           1460 non-null   float64
26  GarageCars            1460 non-null   int64
27  GarageArea            1460 non-null   int64
28  WoodDeckSF            1460 non-null   int64
29  OpenPorchSF           1460 non-null   int64
30  EnclosedPorch         1460 non-null   int64
31  3SsnPorch             1460 non-null   int64
32  ScreenPorch           1460 non-null   int64
33  PoolArea              1460 non-null   int64
34  MiscVal               1460 non-null   int64
35  MoSold                1460 non-null   int64
36  YrSold                1460 non-null   int64
37  SalePrice              1460 non-null   int64
dtypes: float64(3), int64(35)
memory usage: 433.6 KB

```

Creación de la variable categórica de precios (Barata, Media, Cara)

Se genera una variable categórica basada en el valor de `SalePrice` para clasificar las propiedades en "baratas", "medias" y "caras" usando los percentiles 33% y 66% como umbrales de segmentación.

```
# Crear variable categórica
percentiles = np.percentile(data['SalePrice'], [33, 66])

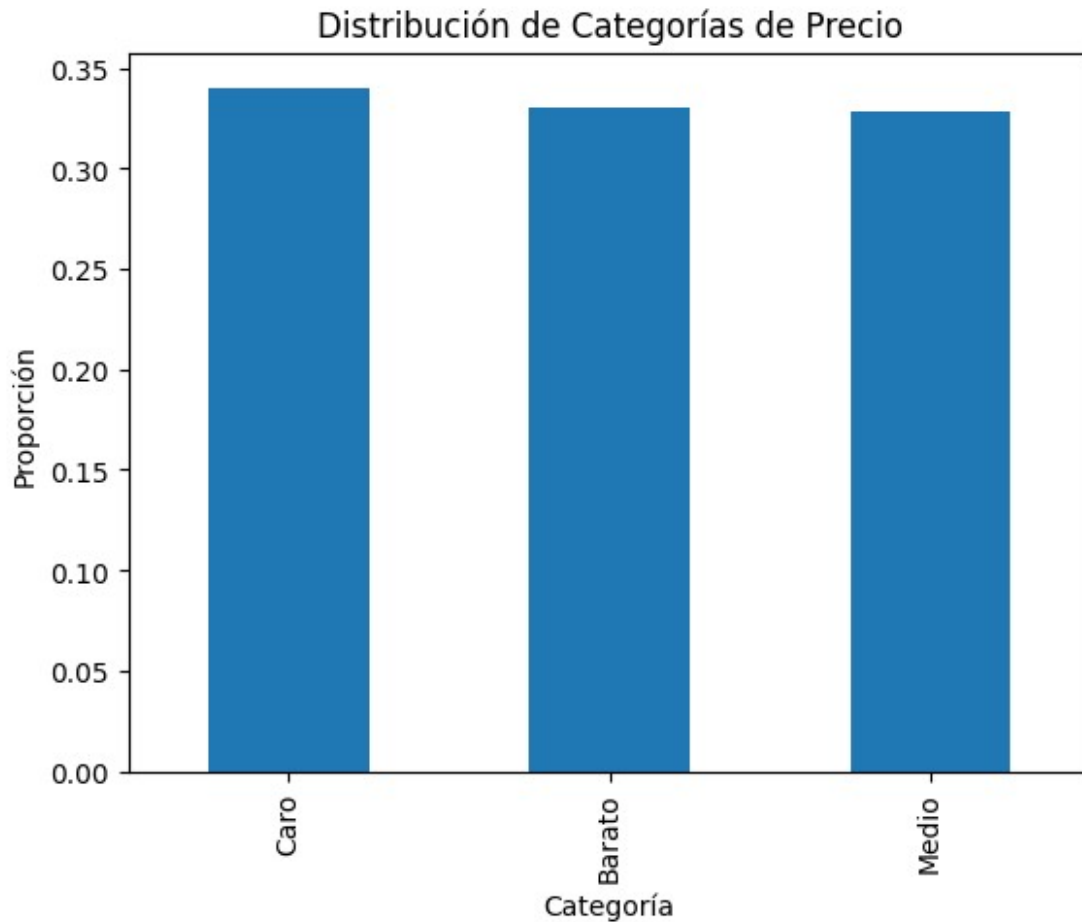
def categorizar_precio(precio):
    if precio <= percentiles[0]:
        return 'Barato'
    elif precio <= percentiles[1]:
        return 'Medio'
    else:
        return 'Caro'

data['PrecioCategoria'] = data['SalePrice'].apply(categorizar_precio)

# Visualización del balance de clases
balance = data['PrecioCategoria'].value_counts(normalize=True)
print(balance)

balance.plot(kind='bar', title='Distribución de Categorías de Precio')
plt.xlabel('Categoría')
plt.ylabel('Proporción')
plt.show()

PrecioCategoria
Caro      0.340411
Barato    0.330822
Medio     0.328767
Name: proportion, dtype: float64
```



Creación de Modelos SVM con diferentes kernels y parámetros

Se crean múltiples modelos SVM utilizando diferentes configuraciones de kernels: `lineal`, `rbf` (gaussiano) y `polinomial`.

Se ajustan también hiperparámetros como `C`, `gamma` y `degree` para explorar su impacto en el desempeño del modelo.

```
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV

# Variables predictoras y respuesta
X = data.drop(columns=['SalePrice', 'PrecioCategoria'])
y = data['PrecioCategoria']

# Escalado
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Partición
```

```

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y,
test_size=0.3, random_state=42)

# Definimos modelos básicos
models = {
    "SVM Lineal": SVC(kernel='linear'),
    "SVM RBF": SVC(kernel='rbf'),
    "SVM Polinomial": SVC(kernel='poly')
}

# Hiperparámetros a buscar
param_grid = {
    'linear': {'C': [0.1, 1, 10]},
    'rbf': {'C': [0.1, 1, 10], 'gamma': [0.01, 0.1, 1]},
    'poly': {'C': [0.1, 1, 10], 'degree': [2, 3, 4]}
}

# GridSearchCV
best_models = {}
for name, model in models.items():
    if model.kernel == 'linear':
        grid = GridSearchCV(model, param_grid['linear'], cv=5,
n_jobs=-1)
    elif model.kernel == 'rbf':
        grid = GridSearchCV(model, param_grid['rbf'], cv=5, n_jobs=-1)
    else: # polinomial
        grid = GridSearchCV(model, param_grid['poly'], cv=5, n_jobs=-
1)

    grid.fit(X_train, y_train)
    best_models[name] = grid.best_estimator_

# Mostrar mejores hiperparámetros
for name, model in best_models.items():
    print(f"Mejor modelo {name}: {model}")

Mejor modelo SVM Lineal: SVC(C=0.1, kernel='linear')
Mejor modelo SVM RBF: SVC(C=1, gamma=0.01)
Mejor modelo SVM Polinomial: SVC(C=10, kernel='poly')

```

Predicción de la Variable Respuesta con los Modelos SVM

Se realizan las predicciones en el conjunto de prueba para evaluar el desempeño de cada modelo SVM utilizando diferentes configuraciones de kernel.

```

# Predicciones
y_pred_linear = svm_linear.predict(X_test)

```

```
y_pred_rbf = svm_rbf.predict(X_test)
y_pred_poly = svm_poly.predict(X_test)
```

Evaluación: Matrices de Confusión

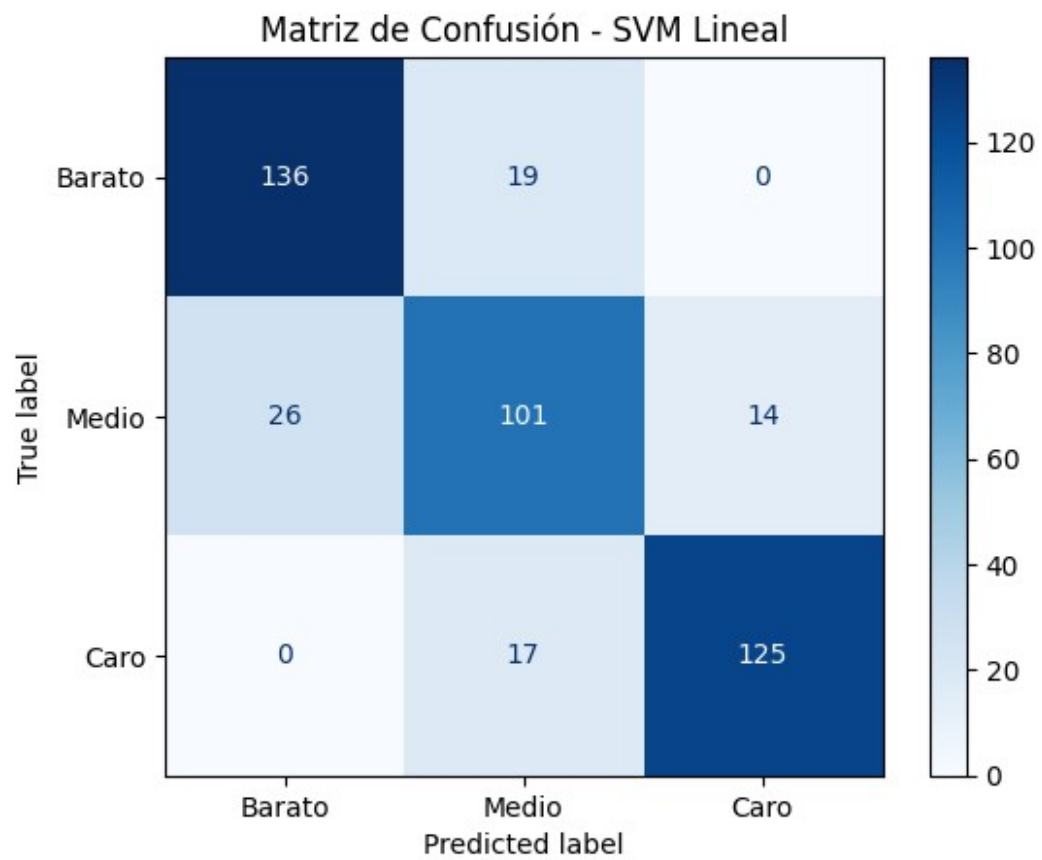
Se presentan las matrices de confusión para los diferentes modelos SVM creados. Estas matrices permiten visualizar el desempeño del modelo en términos de predicciones correctas e incorrectas en cada clase.

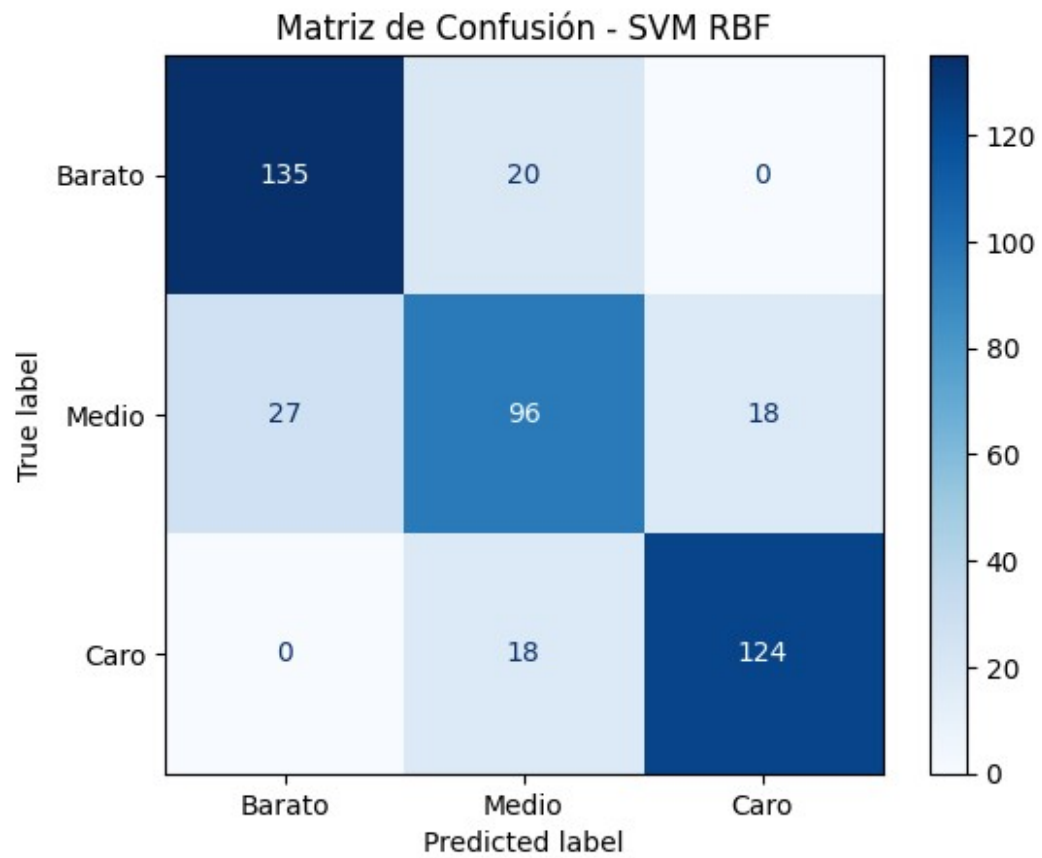
```
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

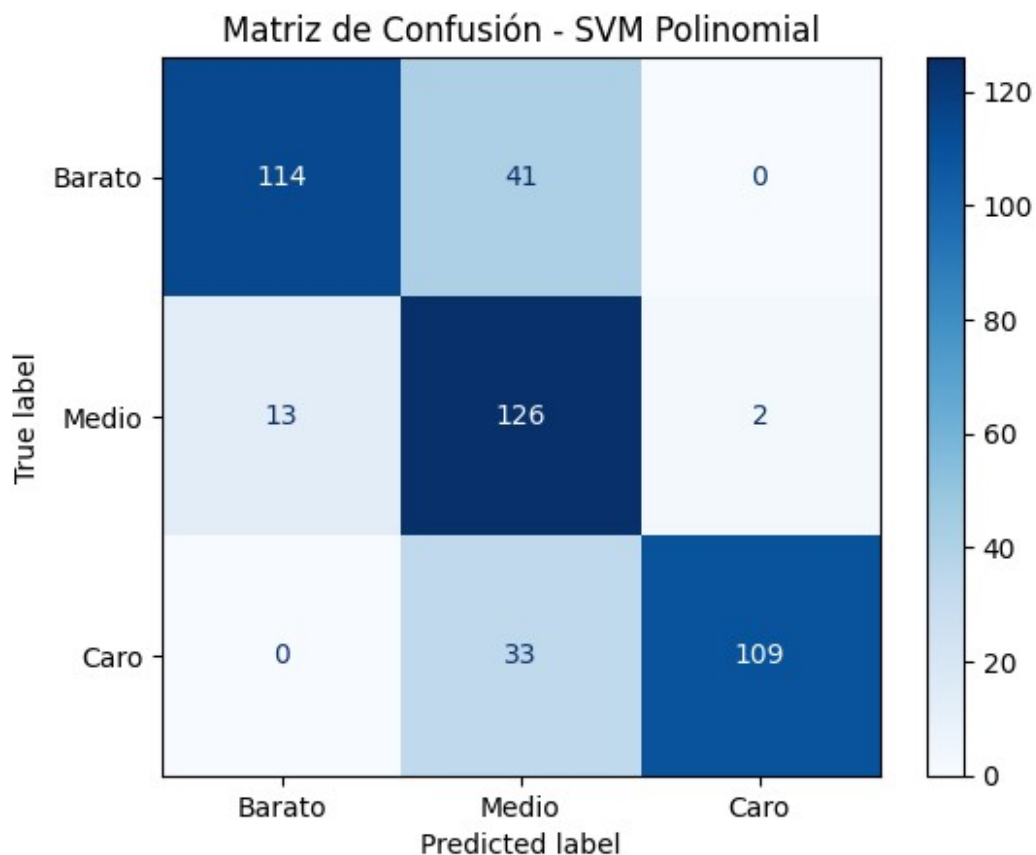
# Función para graficar matrices
def plot_confusion(y_true, y_pred, title):
    cm = confusion_matrix(y_true, y_pred, labels=['Barato', 'Medio',
'Caro'])
    disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=['Barato', 'Medio', 'Caro'])
    disp.plot(cmap='Blues')
    plt.title(title)
    plt.show()

import matplotlib.pyplot as plt

# Mostrar matrices
plot_confusion(y_test, y_pred_linear, 'Matriz de Confusión - SVM
Lineal')
plot_confusion(y_test, y_pred_rbf, 'Matriz de Confusión - SVM RBF')
plot_confusion(y_test, y_pred_poly, 'Matriz de Confusión - SVM
Polinomial')
```





Análisis de Sobreajuste o Subajuste

Se analiza el posible sobreajuste (overfitting) o subajuste (underfitting) de los modelos generados observando su desempeño en entrenamiento y prueba.

Además, se discuten estrategias de ajuste de hiperparámetros para mejorar el balance entre sesgo y varianza de los modelos.

```
# Scores
for name, model in best_models.items():
    print(f"{name} - Train Score: {model.score(X_train, y_train):.4f},
    Test Score: {model.score(X_test, y_test):.4f}")
```

SVM Lineal - Train Score: 0.8630, Test Score: 0.8151

SVM RBF - Train Score: 0.8796, Test Score: 0.8105

SVM Polinomial - Train Score: 0.9746, Test Score: 0.8059