

Aprendizado de Máquina

...

Alexandre Furtado Violante, RA: 632880

House Prices - Advanced Regression Techniques

- **Objetivo:** estimar o preço de casas (SalePrice) em Ames, Iowa, usando dados históricos de vendas e atributos detalhados das casas.
- **Kaggle fornece 2 planilhas:**
 - **train.csv**, com 1460 casas e 79 colunas de características e uma SalePrice.
 - **test.csv**, com 1459 casas e as mesmas 79 colunas, sem a coluna SalePrice

Desafios

- O conjunto de dados possui variáveis numéricas e categóricas, muitas com valores ausentes.
- O modelo deve lidar com dados heterogêneos, diferentes escalas, outliers e relações complexas entre atributos.

Modelo

- **RandomForest Regressor:** método robusto para regressão em dados tabulares, lidando bem com muitas variáveis e relações não lineares.
- **Funcionamento:**
 - Algoritmo de ensemble que combina centenas de árvores de decisão.
 - Cada árvore recebe uma amostra aleatória dos dados e um subconjunto aleatório de features.
 - A predição final é a média das previsões de todas as árvores.

Bibliotecas

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import RandomForestRegressor
```

- *Pandas*: manipulação de dados
- *Numpy*: operações numéricas
- *Matplotlib e Seaborn*: visualização gráfica
- *RandomForestRegressor*: modelo utilizado

Carregamento dos dados

```
train = pd.read_csv("train.csv")  
test  = pd.read_csv("test.csv")  
print(f"Train: {train.shape} | Test: {test.shape}")
```

```
Train: (1460, 81) | Test: (1459, 80)
```

- Os arquivos train.csv e test.csv são carregados.
- A inspeção com shape e info() mostra quantidade de dados, tipos de variáveis e presença de valores ausentes.



```
train.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

Valores ausentes

```
for df in [train, test]:  
    df.fillna(df.select_dtypes(include=np.number).median(), inplace=True)  
    df.fillna('Missing', inplace=True)  
  
print("Missing values tratados → nulos:", train.isnull().sum().sum())  
  
Missing values tratados → nulos: 0
```

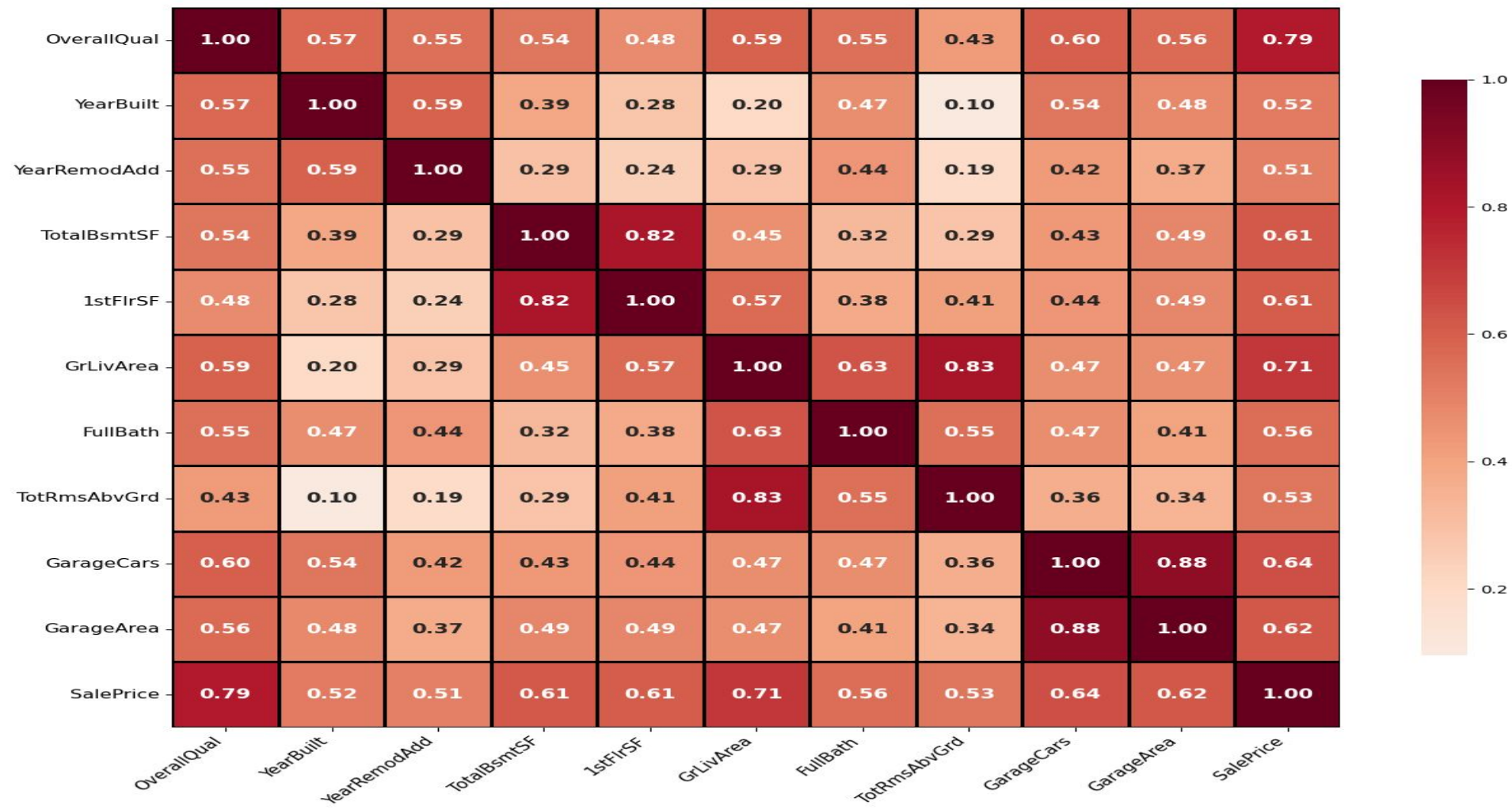
- Numéricos → preenchidos com a mediana (menos sensível a outliers)
- Categóricos → preenchidos com “Missing”
- Essa etapa evita falhas no modelo e permite que ele interprete ausência como informação.

Análise de Correlação

```
numeric_cols = train.select_dtypes(include='number')
corr = numeric_cols.corr()
top_corr = corr.index[abs(corr["SalePrice"]) > 0.5].tolist()

plt.figure(figsize=(14, 11))
sns.heatmap(
    train[top_corr].corr(),
    annot=True,
    fmt=".2f",
    cmap="RdBu_r",
    center=0,
    linewidths=1,
    linecolor='black',
    cbar_kws={"shrink": .8},
    square=True,
    annot_kws={"size": 12, "weight": "bold"}
)
plt.title('Relacao de Features com SalePrice (>0.5)', fontsize=20, fontweight='bold', pad=20)
plt.xticks(rotation=45, ha='right', fontsize=11)
plt.yticks(fontsize=11)
plt.tight_layout()
plt.show()
```

Relacao de Features com SalePrice (>0.5)

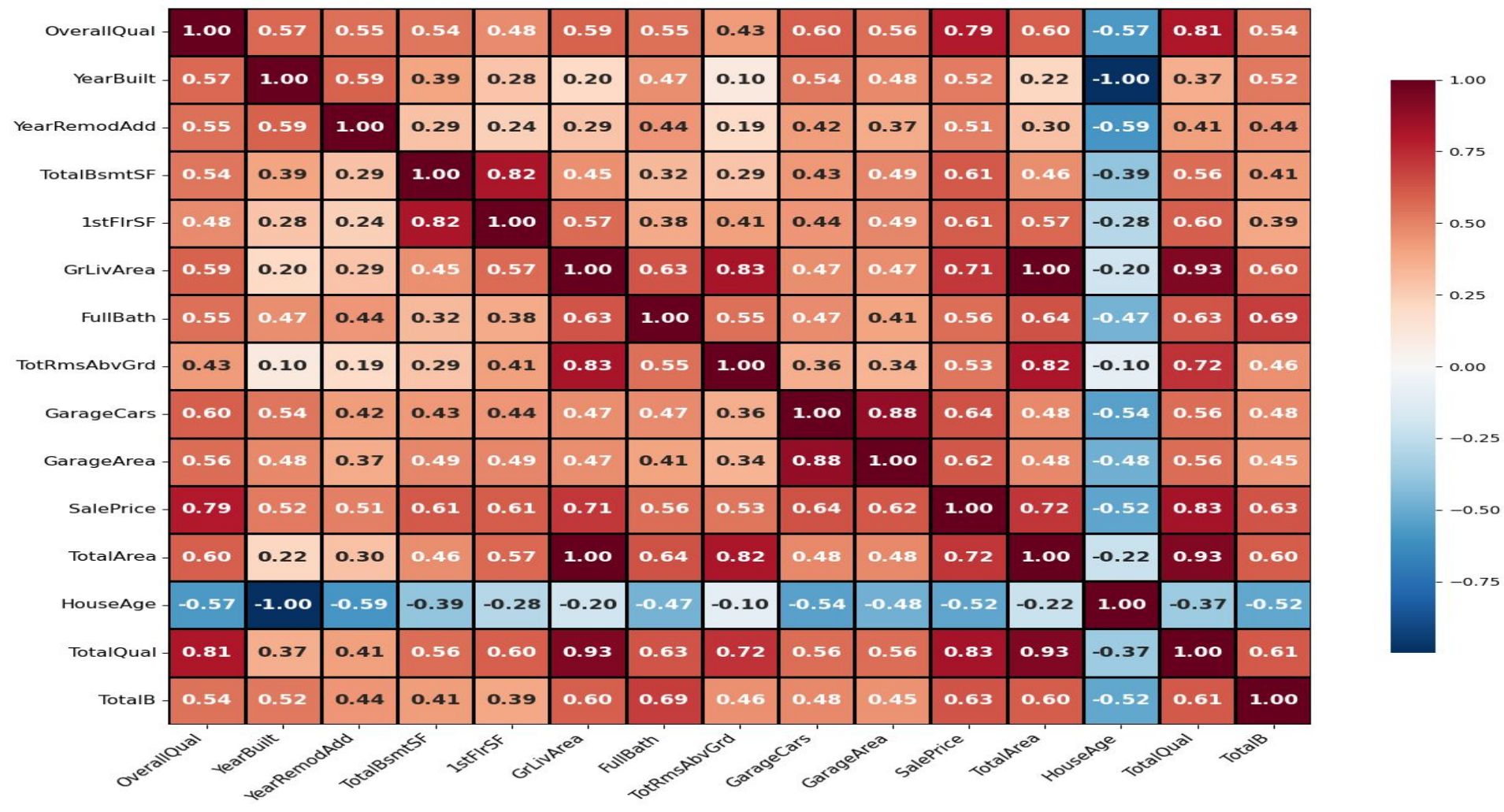


Criação das Features

```
numeric_cols = train.select_dtypes(include='number')
for df in [train, test]:
    df['TotalArea'] = df['1stFlrSF'] + df['2ndFlrSF']
    df['HouseAge'] = df['YrSold'] - df['YearBuilt']
    df['TotalQual'] = df['OverallQual'] * df['GrLivArea']
    df['HasGarage'] = (df['GarageArea'] > 0).astype(int)
    df['TotalB'] = df['FullBath'] + 0.5*df['HalfBath'] + df['BsmtFullBath'] + + 0.5*df['BsmtHalfBath']
```

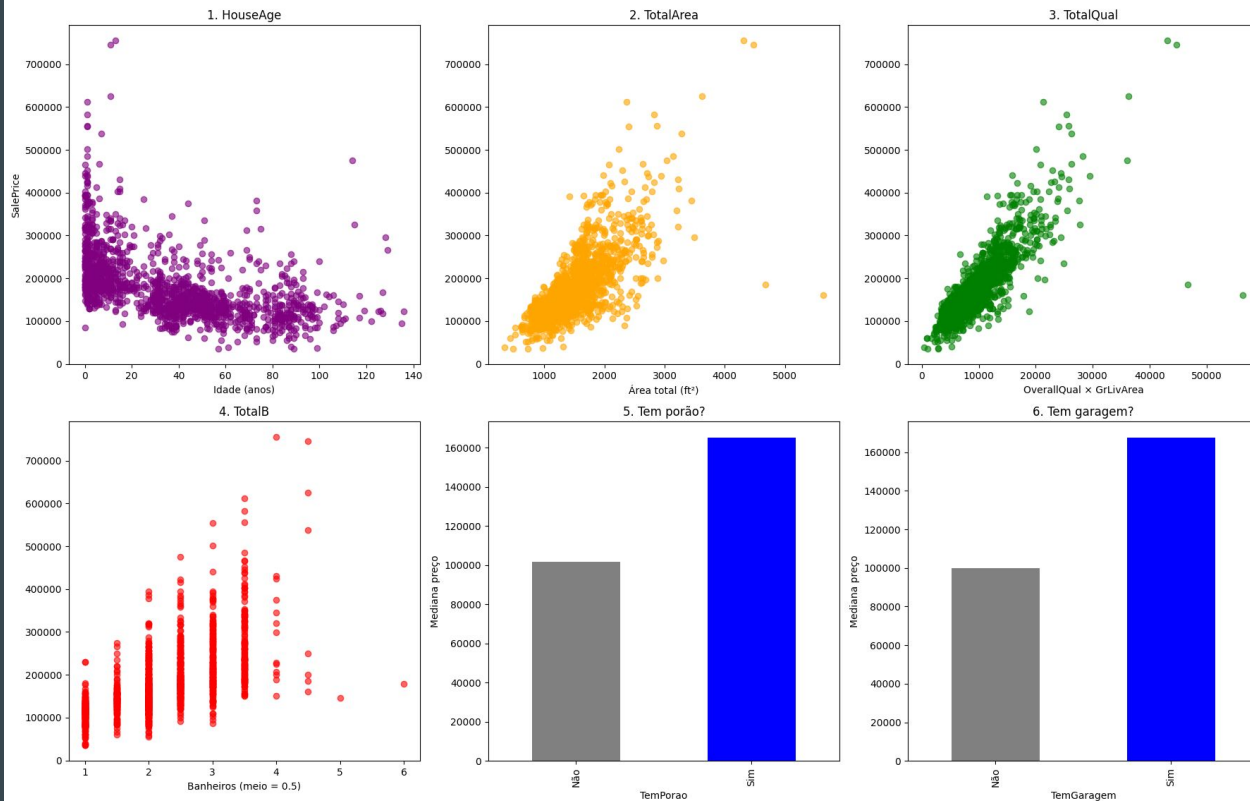
- Área total = 1º andar + 2º andar
- Idade da casa = Ano da venda - Ano da construção
- Qualidade total = Qualidade * Área
- Tem garagem?
- Quantidade de banheiros

Relacao de Features com SalePrice (>0.5)



Visualização

Features Criadas



Remoção das colunas

```
cols_to_drop = [  
    '1stFlrSF',  
    '2ndFlrSF',  
    'YearBuilt',  
    'OverallQual',  
    'GrLivArea',  
    'GarageArea',  
    'FullBath',  
    'HalfBath',  
    'BsmtFullBath',  
    'BsmtHalfBath'  
]  
  
train = train.drop(columns=cols_to_drop, errors='ignore')  
test  = test.drop(columns=cols_to_drop, errors='ignore')
```

- Remoção das colunas utilizadas na criação das novas features.

```
y_train = train['SalePrice']  
  
X_train = train.drop(['SalePrice'], axis=1)  
X_test = test.copy()
```

- Extrai a coluna alvo SalePrice do DataFrame train e armazena em y_train
- Cria X_train removendo a coluna alvo de train
- Faz uma cópia do DataFrame test para usar como X_test

One-hot encoding

```
X_train = pd.get_dummies(X_train, drop_first=True)  
X_test  = pd.get_dummies(X_test, drop_first=True)
```

- O modelo não entende texto, só números.
- Todas as colunas categóricas viram colunas binárias (dummies).
- `drop_first=True` evita multicolinearidade (remove uma das categorias).

Alinhamento

```
X_train, X_test = X_train.align(X_test, join='left', axis=1, fill_value=0)  
  
print(f"Shapes finais → Train: {X_train.shape} | Test: {X_test.shape}")
```

- align() faz os dois datasets terem exatamente as mesmas colunas.

```
Shapes finais → Train: (1460, 258) | Test: (1459, 258)
```

Treino do Modelo Random Forest

```
rf = RandomForestRegressor(n_estimators=1000, random_state=50, n_jobs=-1)
rf.fit(X_train, y_train)
print("Modelo treinado!")
```

- Criação de um Random Forest Regressor:
 - `n_estimators=1000` → 1000 árvores (mais estabilidade).
 - `random_state=50` → reprodução consistente dos resultados.
 - `n_jobs=-1` → usa todos os núcleos da máquina (mais rápido).
- `fit()` treina o modelo com as features X e os preços y.

Previsão e criação do submission.csv

```
pred = rf.predict(X_test)

submission = pd.DataFrame({'Id': test['Id'], 'SalePrice': pred})
submission.to_csv('submission.csv', index=False)
```

- predict() gera o preço previsto para cada casa do dataset de teste.
- O arquivo final é salvo como submission.csv.

Resultado

```
print(submission.head(15))
```

	Id	SalePrice
0	1461	122897.004
1	1462	165727.655
2	1463	178779.357
3	1464	183641.636
4	1465	190005.514
5	1466	179902.925
6	1467	166941.752
7	1468	178063.812
8	1469	184326.727
9	1470	124258.614
10	1471	192269.158
11	1472	92721.317
12	1473	93871.900
13	1474	153041.190
14	1475	110081.410

Obrigado!