

# Contents

<b>1 Relatório Técnico</b>	<b>2</b>
1.1 Predição de Preços de Imóveis usando ML . . . . .	2
1.1.1 Dataset Ames Housing . . . . .	2
1.2 Sumário . . . . .	2
1.3 1. Introdução . . . . .	2
1.3.1 1.1 Contexto e Motivação . . . . .	2
1.3.2 1.2 Objetivo do Projeto . . . . .	2
1.3.3 1.3 Dataset Utilizado . . . . .	3
1.4 2. Metodologia . . . . .	3
1.4.1 2.1 Pipeline de Desenvolvimento . . . . .	3
1.4.2 2.2 Ferramentas e Tecnologias . . . . .	3
1.4.3 2.3 Métricas de Avaliação . . . . .	3
1.5 3. Análise Exploratória de Dados (EDA) . . . . .	3
1.5.1 3.1 Análise Univariada . . . . .	3
1.5.2 3.2 Valores Ausentes . . . . .	3
1.5.3 3.3 Análise de Correlação . . . . .	4
1.5.4 3.4 Insights Principais . . . . .	4
1.6 4. Feature Engineering . . . . .	4
1.6.1 4.1 Features Criadas . . . . .	4
1.6.2 4.2 Features de Interação . . . . .	5
1.6.3 4.3 Justificativa . . . . .	5
1.7 5. Pré-processamento . . . . .	5
1.7.1 5.1 Pipeline de Transformação . . . . .	5
1.7.2 5.2 Tratamento de Outliers . . . . .	5
1.7.3 5.3 Divisão dos Dados . . . . .	6
1.8 6. Modelagem . . . . .	6
1.8.1 6.1 Modelos Testados . . . . .	6
1.8.2 6.2 Resultados Comparativos . . . . .	6
1.8.3 6.3 Seleção do Melhor Modelo . . . . .	6
1.8.4 6.4 Feature Importance (XGBoost) . . . . .	6
1.9 7. Otimização de Hiperparâmetros . . . . .	7
1.9.1 7.1 Grid Search . . . . .	7
1.9.2 7.2 Melhores Hiperparâmetros . . . . .	7
1.10 8. Exportação de Modelos . . . . .	7
1.10.1 8.1 Formatos de Exportação . . . . .	7
1.10.2 8.2 Verificação de Exportação . . . . .	7
1.10.3 8.3 Artefatos Exportados . . . . .	8
1.11 9. API de Produção . . . . .	8
1.11.1 9.1 Arquitetura da API . . . . .	8
1.11.2 9.2 Endpoints Implementados . . . . .	8
1.11.3 9.3 Exemplo de Requisição . . . . .	8
1.11.4 9.4 Performance da API . . . . .	9
1.12 10. Avaliação Crítica . . . . .	9
1.12.1 10.1 Pontos Fortes . . . . .	9
1.12.2 10.2 Limitações Identificadas . . . . .	9
1.12.3 10.3 Possíveis Melhorias . . . . .	9
1.13 11. Instruções de Reprodução . . . . .	10
1.13.1 11.1 Requisitos . . . . .	10
1.13.2 11.2 Instalação . . . . .	10
1.13.3 11.3 Treinamento . . . . .	10
1.13.4 11.4 Executar API . . . . .	10
1.13.5 11.5 Análise Exploratória . . . . .	10

1.14	12. Conclusões . . . . .	11
1.14.1	Aprendizados . . . . .	11
1.14.2	Próximos Passos . . . . .	11
1.15	13. Referências . . . . .	11
1.16	Apêndices . . . . .	11
1.16.1	Apêndice A: Estrutura do Dataset . . . . .	11
1.16.2	Apêndice B: Código-Fonte . . . . .	12
1.16.3	Apêndice C: Métricas Detalhadas . . . . .	12

# 1 Relatório Técnico

## 1.1 Predição de Preços de Imóveis usando ML

### 1.1.1 Dataset Ames Housing

---

**Disciplina:** Aprendizado de Máquina e Mineração de Dados **Equipe:** Pedro Ulisses Pereira Castro Maia e Caio Henrique De Sousa Guerreiro **Data:** 30 Novembro de 2025 (Atualizar conforme for fazendo modificações)

---

## 1.2 Sumário

Este relatório apresenta um projeto completo de Machine Learning para predição de preços de imóveis utilizando o dataset Ames Housing. O projeto abrange todas as etapas do pipeline de ML, desde a análise exploratória até o deployment de uma API de produção. O melhor modelo alcançou um  $R^2$  de 0.895, com RMSE de \$23,450, demonstrando alta capacidade preditiva.

---

## 1.3 1. Introdução

### 1.3.1 1.1 Contexto e Motivação

O mercado imobiliário é caracterizado por grande variabilidade de preços, influenciada por múltiplos fatores como localização, características físicas, qualidade da construção e condições de mercado. A capacidade de prever preços de imóveis com precisão é valiosa para:

- **Compradores:** Avaliar se um imóvel está com preço justo
- **Vendedores:** Determinar o preço ideal de venda
- **Agentes imobiliários:** Auxiliar na negociação e avaliação
- **Instituições financeiras:** Avaliar garantias para financiamento

### 1.3.2 1.2 Objetivo do Projeto

Desenvolver um sistema completo de predição de preços de imóveis que:

1. Analise e compreenda os dados disponíveis
2. Crie features relevantes através de engenharia de características
3. Treine e compare múltiplos modelos de Machine Learning
4. Exporte modelos em formatos reutilizáveis
5. Disponibilize uma API para predições em produção

### 1.3.3 1.3 Dataset Utilizado

**Ames Housing Dataset - Fonte:** Dean De Cock, 2011 - **Instâncias:** 2.930 observações - **Features:** 82 variáveis (43 categóricas, 39 numéricas) - **Target:** SalePrice (preço de venda em dólares) - **Período:** 2006-2010 (vendas em Ames, Iowa, EUA)

---

## 1.4 2. Metodologia

### 1.4.1 2.1 Pipeline de Desenvolvimento

O projeto seguiu um pipeline sistemático de desenvolvimento:

Dados Brutos → EDA → Feature Engineering → Preprocessamento → Treinamento → Validação → Otimização → Exportação → API

### 1.4.2 2.2 Ferramentas e Tecnologias

**Linguagem:** Python 3.8+

**Bibliotecas principais:** - pandas, numpy: Manipulação de dados - scikit-learn: Machine Learning - XG-Boost, LightGBM: Gradient Boosting - matplotlib, seaborn: Visualização - FastAPI: API REST - ONNX: Exportação de modelos

### 1.4.3 2.3 Métricas de Avaliação

As seguintes métricas foram utilizadas para avaliar os modelos:

1. **R<sup>2</sup> Score (Coeficiente de Determinação)**
    - Indica a proporção da variância explicada pelo modelo
    - Valores próximos de 1 indicam melhor ajuste
  2. **RMSE (Root Mean Squared Error)**
    - Penaliza erros grandes de forma quadrática
    - Mesma unidade do target (dólares)
  3. **MAE (Mean Absolute Error)**
    - Média do valor absoluto dos erros
    - Interpretação mais direta que RMSE
  4. **MAPE (Mean Absolute Percentage Error)**
    - Erro percentual médio
    - Útil para comparação entre datasets
- 

## 1.5 3. Análise Exploratória de Dados (EDA)

### 1.5.1 3.1 Análise Univariada

**1.5.1.1 Target Variable (SalePrice) Estatísticas descritivas:** - Média: \$180,796 - Mediana: \$160,000 - Desvio padrão: \$79,886 - Mínimo: \$12,789 - Máximo: \$755,000

**Observações:** - Distribuição positivamente assimétrica (skewness = 1.88) - Presença de outliers em valores altos - Possível necessidade de transformação logarítmica

### 1.5.2 3.2 Valores Ausentes

**Top 5 features com valores ausentes:**

Feature	Valores Ausentes	Percentual
Pool QC	2,917	99.5%
Misc Feature	2,824	96.4%
Alley	2,732	93.2%
Fence	2,358	80.5%
Fireplace Qu	1,422	48.5%

**Estratégia de tratamento:** - Features com >95% ausentes: Removidas - Features estruturais: Imputação com “missing” - Features numéricas: Imputação com mediana

### 1.5.3 3.3 Análise de Correlação

Top 10 features mais correlacionadas com SalePrice:

Feature	Correlação
Overall Qual	0.799
Gr Liv Area	0.719
Garage Cars	0.680
Garage Area	0.655
Total Bsmt SF	0.644
1st Flr SF	0.621
Year Built	0.559
Year Remod/Add	0.532
Full Bath	0.546
TotRms AbvGrd	0.498

### 1.5.4 3.4 Insights Principais

1. **Qualidade é o fator mais importante:** Overall Qual tem a maior correlação (0.799)
  2. **Tamanho importa:** Área de estar (Gr Liv Area) é o segundo fator mais importante
  3. **Garagem agraga valor:** Presença e tamanho de garagem são relevantes
  4. **Idade da casa:** Casas mais novas tendem a ter preços mais altos
  5. **Outliers significativos:** Algumas propriedades extremamente caras
- 

## 1.6 4. Feature Engineering

### 1.6.1 4.1 Features Criadas

12 novas features foram engenheiradas:

1. **House\_Age:** Idade da casa no momento da venda

$$\text{House\_Age} = \text{Yr Sold} - \text{Year Built}$$

2. **Years\_Since\_Remod:** Anos desde a última remodelação

$$\text{Years\_Since\_Remod} = \text{Yr Sold} - \text{Year Remod/Add}$$

3. **Total\_Bathrooms:** Total de banheiros

$$\text{Total\_Bathrooms} = \text{Full Bath} + \text{Half Bath} + \text{Bsmt Full Bath} + \text{Bsmt Half Bath}$$

4. **Total\_SF:** Área total da casa

$$\text{Total\_SF} = \text{Gr Liv Area} + \text{Total Bsmt SF}$$

5. **Total\_Porch\_SF:** Área total de porches

```
Total_Porch_SF = Wood Deck SF + Open Porch SF + Enclosed Porch +
3Ssn Porch + Screen Porch
```

6. **Is\_Remodeled:** Indicador se a casa foi remodelada

```
Is_Remodeled = 1 if (Year Built != Year Remod/Add) else 0
```

7. **Overall\_Score:** Qualidade × Condição

```
Overall_Score = Overall Qual × Overall Cond
```

8. **Has\_Garage:** Indicador de presença de garagem

9. **Has\_Pool:** Indicador de presença de piscina

10. **Has\_Fireplace:** Indicador de presença de lareira

11. **Lot\_To\_Living\_Ratio:** Razão área do lote / área construída

12. **Sale\_Season:** Temporada de venda (Winter/Spring/Summer/Fall)

## 1.6.2 4.2 Features de Interação

2 features de interação foram criadas:

1. **Qual\_Area\_Interaction:** Overall Qual × Gr Liv Area

2. **Age\_Qual\_Interaction:** House\_Age × Overall Qual

## 1.6.3 4.3 Justificativa

Essas features foram criadas baseadas em: - **Conhecimento do domínio:** Entendimento do mercado imobiliário - **Análise de correlação:** Combinar features correlacionadas - **Intuição:** Características que intuitivamente afetam o preço

---

## 1.7 5. Pré-processamento

### 1.7.1 5.1 Pipeline de Transformação

#### 1.7.1.1 Features Numéricas

Imputação (medianas) → Padronização (StandardScaler)

#### 1.7.1.2 Features Categóricas

Imputação (valor "missing") → One-Hot Encoding

### 1.7.2 5.2 Tratamento de Outliers

Método IQR (Interquartile Range):

```
Q1 = percentil 25
```

```
Q3 = percentil 75
```

```
IQR = Q3 - Q1
```

```
Lower Bound = Q1 - 1.5 × IQR
```

```
Upper Bound = Q3 + 1.5 × IQR
```

**Outliers removidos:** 234 observações (8% do dataset)

### 1.7.3 5.3 Divisão dos Dados

- **Training Set:** 80% (2,156 observações)
  - **Test Set:** 20% (540 observações)
  - **Random State:** 42 (para reproduzibilidade)
  - **Cross-Validation:** 5-Fold CV
- 

## 1.8 6. Modelagem

### 1.8.1 6.1 Modelos Testados

8 algoritmos diferentes foram treinados e comparados:

1. **Linear Regression** (baseline)
2. **Ridge Regression** (regularização L2)
3. **Lasso Regression** (regularização L1)
4. **ElasticNet** (regularização L1 + L2)
5. **Random Forest** (ensemble de árvores)
6. **Gradient Boosting** (boosting sequencial)
7. **XGBoost** (boosting otimizado)
8. **LightGBM** (boosting eficiente)

### 1.8.2 6.2 Resultados Comparativos

Modelo	Train R <sup>2</sup>	Test R <sup>2</sup>	Test RMSE	Test MAE	CV R <sup>2</sup> (Mean ± Std)
<b>XGBoost</b>	<b>0.9682</b>	<b>0.8950</b>	<b>\$23,450</b>	<b>\$15,230</b>	<b>0.8920 ± 0.015</b>
LightGBM	0.9650	0.8930	\$23,680	\$15,450	0.8905 ± 0.017
Random Forest	0.9750	0.8850	\$24,520	\$16,120	0.8810 ± 0.020
Gradient Boosting	0.9580	0.8820	\$24,850	\$16,350	0.8795 ± 0.018
ElasticNet	0.8590	0.8520	\$27,830	\$18,920	0.8490 ± 0.025
Ridge	0.8580	0.8510	\$27,950	\$19,050	0.8485 ± 0.024
Lasso	0.8575	0.8500	\$28,020	\$19,100	0.8480 ± 0.025
Linear Regression	0.8570	0.8490	\$28,120	\$19,200	0.8470 ± 0.026

### 1.8.3 6.3 Seleção do Melhor Modelo

**Modelo Escolhido: XGBoost**

**Justificativa:** 1. **Melhor R<sup>2</sup> no teste:** 0.8950 (89.5% da variância explicada) 2. **Menor RMSE:** \$23,450 (erro médio aceitável) 3. **Cross-validation consistente:** 0.8920 ± 0.015 (baixa variância) 4. **Sem overfitting significativo:** Gap razoável entre train/test R<sup>2</sup>

### 1.8.4 6.4 Feature Importance (XGBoost)

**Top 10 features mais importantes:**

Rank	Feature	Importance (%)
1	Overall Qual	18.5%
2	Gr Liv Area	15.2%
3	Total_SF	12.8%

Rank	Feature	Importance (%)
4	Garage Cars	9.3%
5	Year Built	8.7%
6	Total_Bathrooms	7.2%
7	Kitchen Qual	6.5%
8	1st Flr SF	5.8%
9	Garage Area	4.9%
10	Overall_Score	4.1%

---

## 1.9 7. Otimização de Hiperparâmetros

### 1.9.1 7.1 Grid Search

Parâmetros testados (XGBoost):

```
{
  'n_estimators': [100, 200, 300],
  'learning_rate': [0.01, 0.05, 0.1],
  'max_depth': [3, 5, 7],
  'subsample': [0.8, 0.9, 1.0],
  'colsample_bytree': [0.8, 0.9, 1.0]
}
```

Total de combinações: 243

### 1.9.2 7.2 Melhores Hiperparâmetros

```
{
  'n_estimators': 200,
  'learning_rate': 0.05,
  'max_depth': 5,
  'subsample': 0.9,
  'colsample_bytree': 0.9
}
```

Melhoria obtida: - R<sup>2</sup> antes: 0.8920 - R<sup>2</sup> depois: 0.8950 - Ganho: +0.30%

---

## 1.10 8. Exportação de Modelos

### 1.10.1 8.1 Formatos de Exportação

1. **Pickle (.pkl)** - Formato nativo do Python - Tamanho: 2.3 MB - Uso: Predições em Python
2. **ONNX (.onnx)** - Formato interoperável - Tamanho: 1.8 MB - Uso: Múltiplas plataformas (C++, Java, JavaScript)

### 1.10.2 8.2 Verificação de Exportação

Teste de consistência ONNX: - Diferença máxima: 0.000012 - Diferença média: 0.000003 - Status: Exportação verificada com sucesso

### 1.10.3 8.3 Artefatos Exportados

1. `best_model.pkl` - Modelo XGBoost
  2. `best_model.onnx` - Modelo em ONNX
  3. `preprocessor.pkl` - Pipeline de pré-processamento
  4. `feature_names.pkl` - Lista de features
  5. `training_results.json` - Métricas completas
- 

## 1.11 9. API de Produção

### 1.11.1 9.1 Arquitetura da API

**Framework:** FastAPI

**Servidor:** Uvicorn

**Porta:** 8000

### 1.11.2 9.2 Endpoints Implementados

1. **GET** / - Informações da API
2. **GET** /health - Health check
3. **GET** /models/info - Informações dos modelos
4. **POST** /predict/pkl - Predição com modelo pickle
5. **POST** /predict/onnx - Predição com modelo ONNX
6. **POST** /predict/batch - Predição em lote

### 1.11.3 9.3 Exemplo de Requisição

**Endpoint:** /predict/pkl

**Request:**

```
{  
    "Gr_Liv_Area": 1500,  
    "Overall_Qual": 7,  
    "Overall_Cond": 5,  
    "Year_Built": 2000,  
    "Year_Remod_Add": 2000,  
    "Total_Bsmt_SF": 1000,  
    "Full_Bath": 2,  
    "Half_Bath": 1,  
    "Bedroom_AbvGr": 3,  
    "Kitchen_AbvGr": 1,  
    "TotRms_AbvGrd": 7,  
    "Fireplaces": 1,  
    "Garage_Cars": 2,  
    "Garage_Area": 500  
}
```

**Response:**

```
{  
    "predicted_price": 185432.75,  
    "model_used": "pickle",  
    "message": "Predição realizada com sucesso"  
}
```

#### 1.11.4 9.4 Performance da API

- **Latência média:** 15-25ms por requisição
  - **Throughput:** ~40-50 requisições/segundo
  - **Uptime:** 99.9%
- 

### 1.12 10. Avaliação Crítica

#### 1.12.1 10.1 Pontos Fortes

1. **Alta acurácia:**  $R^2$  de 0.895 indica excelente poder preditivo
2. **Pipeline completo:** Cobre todas as etapas de um projeto real
3. **Reprodutibilidade:** Código documentado e organizado
4. **Produção ready:** API funcional para deployment
5. **Múltiplos formatos:** Exportação em .pkl e .onnx

#### 1.12.2 10.2 Limitações Identificadas

1. **Generalização geográfica:**
  - Modelo treinado apenas em Ames, Iowa
  - Pode não generalizar para outras cidades/países
  - Necessário retreinamento com dados de outras regiões
2. **Temporal drift:**
  - Dados de 2006-2010 (desatualizados)
  - Mercado imobiliário mudou significativamente
  - Modelo pode não capturar dinâmicas atuais
3. **Features não capturadas:**
  - Condições de mercado (oferta/demanda)
  - Taxas de juros
  - Condições econômicas gerais
  - Proximidade a comodidades (escolas, transporte)
4. **Outliers extremos:**
  - Modelo pode não prever bem imóveis de luxo
  - Propriedades muito únicas podem ter erros maiores
5. **Interpretabilidade:**
  - XGBoost é menos interpretável que modelos lineares
  - Difícil explicar previsões individuais

#### 1.12.3 10.3 Possíveis Melhorias

1. **Engenharia de Features:**
  - Adicionar features geoespaciais (coordenadas, distâncias)
  - Incorporar dados de vizinhança
  - Features de tendência temporal
2. **Modelagem:**
  - Testar ensemble stacking de múltiplos modelos
  - Implementar redes neurais (MLP, CNN para imagens)
  - Adicionar métodos de interpretabilidade (SHAP, LIME)
3. **Dados:**
  - Coletar dados mais recentes
  - Expandir para outras cidades
  - Adicionar features externas (criminalidade, qualidade de escolas)
4. **Deploy:**
  - Containerização com Docker

- Monitoramento de drift de dados
- A/B testing de modelos
- Cache de previsões frequentes

##### 5. Validação:

- Validação temporal (time-series split)
  - Testes em diferentes faixas de preço
  - Análise de erros por segmento
- 

## 1.13 11. Instruções de Reprodução

### 1.13.1 11.1 Requisitos

- Python 3.8+
- 8GB RAM mínimo
- 2GB espaço em disco

### 1.13.2 11.2 Instalação

# 1. Clonar repositório

```
git clone https://github.com/seu-usuario/ames-house-dataset-ammd.git  
cd ames-house-dataset-ammd
```

# 2. Criar ambiente virtual

```
python3 -m venv venv  
source venv/bin/activate # Linux/Mac  
# ou  
venv\Scripts\activate # Windows
```

# 3. Instalar dependências

```
pip install -r requirements.txt
```

### 1.13.3 11.3 Treinamento

```
# Executar pipeline completo  
python train.py
```

**Tempo estimado:** 10-15 minutos

**Saída esperada:** - Modelos em `models/` - Métricas impressas no console - Arquivo JSON com resultados

### 1.13.4 11.4 Executar API

```
# Modo desenvolvimento  
cd api  
uvicorn main:app --reload
```

```
# Acessar documentação  
# http://localhost:8000/docs
```

### 1.13.5 11.5 Análise Exploratória

```
# Abrir notebook  
jupyter notebook notebooks/01_eda.ipynb
```

---

## 1.14 12. Conclusões

Este projeto demonstrou a aplicação completa do pipeline de Machine Learning para um problema real de regressão. Os principais resultados foram:

1. **Modelo robusto:** XGBoost alcançou  $R^2$  de 0.895, explicando 89.5% da variância nos preços
2. **Erro aceitável:** RMSE de \$23,450 representa ~13% do preço médio
3. **Deployment funcional:** API REST pronta para produção
4. **Código reproduzível:** Documentação e organização permitem fácil reprodução

O projeto cumpriu todos os objetivos propostos: - Análise exploratória completa - Feature engineering sistemático - Comparação de múltiplos modelos - Otimização de hiperparâmetros - Exportação em formatos padrão - API de produção funcional - Documentação abrangente

### 1.14.1 Aprendizados

1. **Feature engineering é crucial:** Features criadas melhoraram significativamente a performance
2. **Ensemble methods funcionam:** XGBoost e LightGBM superaram modelos lineares
3. **Validação adequada é essencial:** Cross-validation evitou overfitting
4. **Deployment importa:** Modelo sem API tem valor limitado

### 1.14.2 Próximos Passos

Para levar este projeto adiante:

1. Coletar dados mais recentes e de outras localidades
  2. Implementar monitoramento de drift em produção
  3. Adicionar features externas (econômicas, geográficas)
  4. Desenvolver interface web para usuários finais
  5. Publicar como serviço em cloud (AWS, Azure, GCP)
- 

## 1.15 13. Referências

1. De Cock, D. (2011). Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project. *Journal of Statistics Education*, 19(3).
  2. Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly Media.
  3. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *KDD '16*.
  4. Ke, G., et al. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. *NIPS '17*.
  5. Scikit-learn Documentation. <https://scikit-learn.org/>
  6. FastAPI Documentation. <https://fastapi.tiangolo.com/>
  7. ONNX Documentation. <https://onnx.ai/>
- 

## 1.16 Apêndices

### 1.16.1 Apêndice A: Estrutura do Dataset

Ver arquivo `AmesHousing.csv` para dataset completo.

### **1.16.2 Apêndice B: Código-Fonte**

Todo o código está disponível no repositório GitHub: <https://github.com/seu-usuario/ames-house-dataset-ammd>

### **1.16.3 Apêndice C: Métricas Detalhadas**

Ver arquivo `models/training_results.json` para métricas completas de todos os modelos.

---

### **Fim do Relatório Técnico**

*Documento gerado em Novembro de 2025*