

Nanodegree Engenheiro de Machine Learning

Projeto final

Márcio Alexandre Costa
13 de maio de 2019

I. Definição

Visão geral do projeto

Construir um modelo que possa prever o valor de venda das casas com base na localização e característica da casa ideal.

Descrição do problema

O objetivo é criar um modelo de Aprendizagem de Máquinas para prever o valor de venda de casa em Ames / Iowa nos Estados Unidos.

Com base na descrição da casa ideal, como localização, quantidade de quartos, banheiros, tamanho da casa e muitas outras descrições (79 descrições estão disponíveis no dataset deste projeto) que o comprador pode usar para descrever e ser usado pelo modelo para prever o valor de venda destas casas.

O modelo poderá servir para balizar o sonho da casa ideal verso a valor que se deseja ou pode pagar por esta "casa dos sonhos". Sabemos que nós, seres humanos somos muito ruins em prever e pessoalmente que acredito que um algoritmo bem treinado e constantemente atualizado pode ajudar e muito com a previsão do valor de compra e/ou venda de casas.

Segue 2 estudos acadêmicos que se utiliza de modelos de Aprendizagem de Máquinas para resolver problemas de predição de valor de venda de casas.

De Cook, Dean. "Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project." Journal of Statistics Education, vol. 19, no. 3, 2011 -

<https://amstat.tandfonline.com/doi/abs/10.1080/10691898.2011.11889627#aHR0cHM6Ly9hbXN0YXQudGFuZGZvbmxpbmUuY29tL2RvaS9wZGYvMTAuMTA4MC8xMDY5MTg5OC4yMDEuLjExODg5NjI3P25lZWRYBY2Nlc3M9dHJ1ZUBAQDA=> .

De Hujia Yu, Jiafu Wu. "Real Estate Price Prediction with Regression and Classification." CS 229 Autumn 2016 Project Final Report - http://cs229.stanford.edu/proj2016/report/WuYu_HousingPrice_report.pdf

A ideia deste projeto é utilizar algoritmo de Regressão Linear que melhor prever o valor de vende das casas e utilizá-lo no dia-a-dia para ajudar potenciais clientes e corretores a selecionar as melhores casas a ser visitadas com base no valor que se despõe a pagar e as características da casa almejada.

Métricas

O modelo será avaliado pelo RMSE (Root Mean Squared Error) entre o logaritmo do valor previsto e o logaritmo do valor de venda observado.

$$\text{RMSE}_{fo} = \left[\sum_{i=1}^N (z_{fi} - z_{oi})^2 / N \right]^{1/2}$$

Σ = summation ("add up")

$(z_{fi} - z_{oi})^2$ = differences, squared

N = sample size.

O RMSE é a medida de erro mais comumente usada para aferir a qualidade do ajuste de um modelo é a chamada RAIZ DO ERRO MÉDIO QUADRÁTICO.

Ela é a raiz do erro médio quadrático da diferença entre a predição e o valor real.

O RMSE é independente de escala e um ótimo candidato para representar explicitamente o erro de modelos lineares.

II. Análise

Exploração dos dados

O dataset deste projeto está disponibilizado no Kaggle, no link <https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data>.

O conjunto de dados contém 79 características e 1.460 entradas, sendo 43 características com valores categóricas (não numéricas) e 36 características com valores numéricas. O conjunto de dados contém 19 características com valores ausentes (não preenchidos).

Segue abaixo uma pequena amostra do dataset que será trabalhado. Inicialmente o dataset continha 1.460 registros e 79 variáveis sobre a casa.

```
[ ] df_train = pd.read_csv('gdrive/My Drive/Colab Notebooks/train.csv')
df_train.head()
```

| | Id | MSSubClass | MSZoning | LotFrontage | LotArea | Street | Alley | LotShape | LandContour | Utilities | ... | PoolArea | PoolQC | Fence | MiscFeature | MiscVal | MoSold | YrSold | SaleType | SaleCondition | SalePrice |
|---|----|------------|----------|-------------|---------|--------|-------|----------|-------------|-----------|-----|----------|--------|-------|-------------|---------|--------|--------|----------|---------------|-----------|
| 0 | 1 | 60 | RL | 65.0 | 8450 | Pave | NaN | Reg | Lvl | AllPub | ... | 0 | NaN | NaN | NaN | 0 | 2 | 2006 | WD | Normal | 208500 |
| 1 | 2 | 20 | RL | 80.0 | 9600 | Pave | NaN | Reg | Lvl | AllPub | ... | 0 | NaN | NaN | NaN | 0 | 5 | 2007 | WD | Normal | 181500 |
| 2 | 3 | 60 | RL | 68.0 | 11250 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN | NaN | NaN | 0 | 9 | 2006 | WD | Normal | 223500 |
| 3 | 4 | 70 | RL | 60.0 | 9550 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN | NaN | NaN | 0 | 2 | 2006 | WD | Abnorml | 140000 |
| 4 | 5 | 60 | RL | 84.0 | 14260 | Pave | NaN | IR1 | Lvl | AllPub | ... | 0 | NaN | NaN | NaN | 0 | 12 | 2008 | WD | Normal | 250000 |

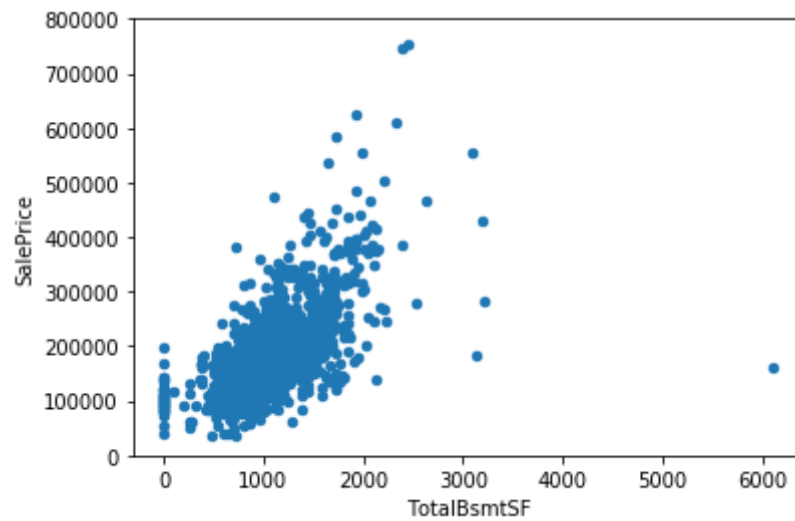
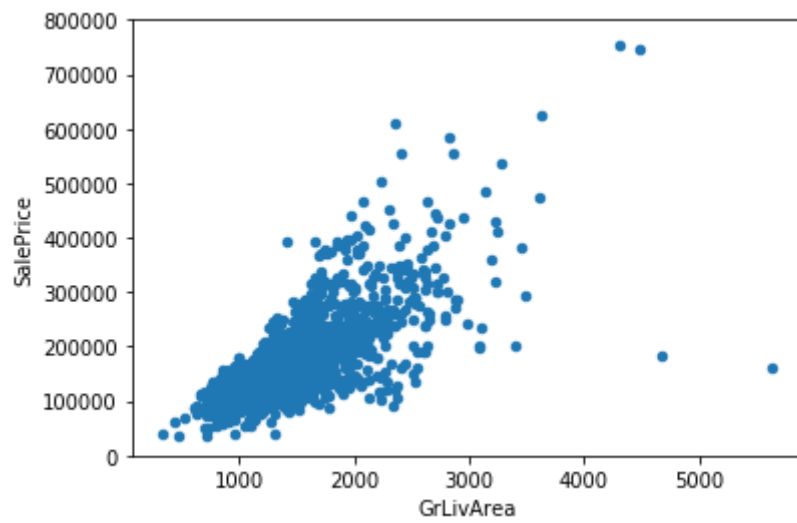
5 rows x 81 columns

Os dados estatísticos da variável alvo – SalePrice

```
[6] df_train['SalePrice'].describe()
```

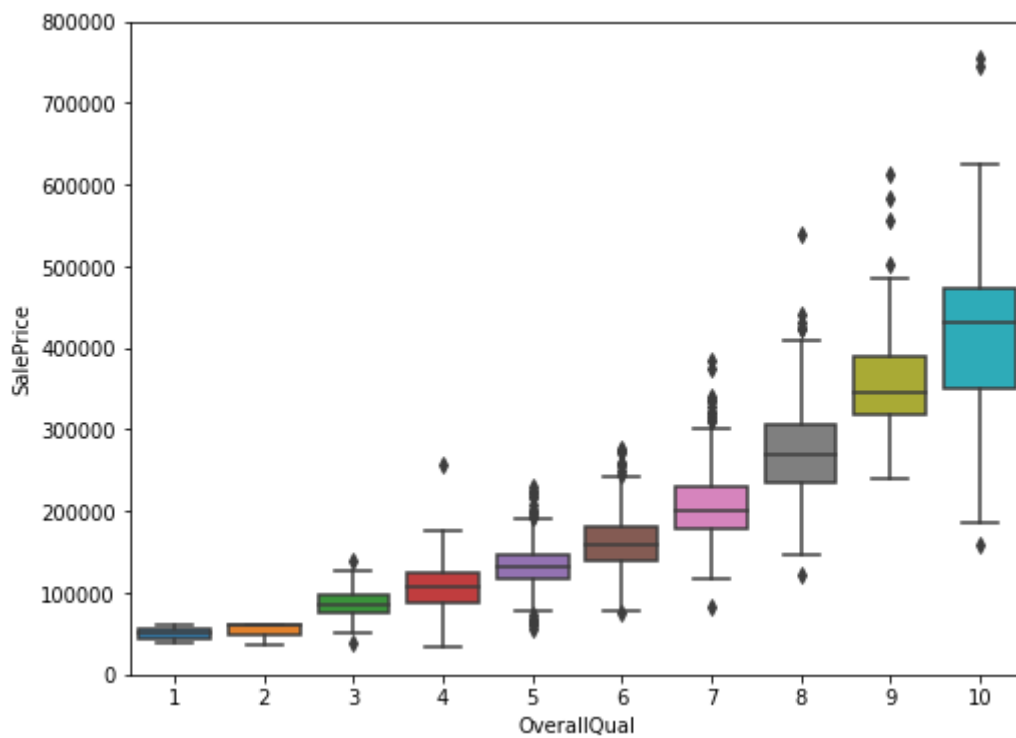
```
count    1460.000000
mean     180921.195890
std       79442.502883
min       34900.000000
25%      129975.000000
50%      163000.000000
75%      214000.000000
max       755000.000000
Name: SalePrice, dtype: float64
```

Dois gráficos entre as duas variáveis numéricas com a variável alvo – SalePrice.



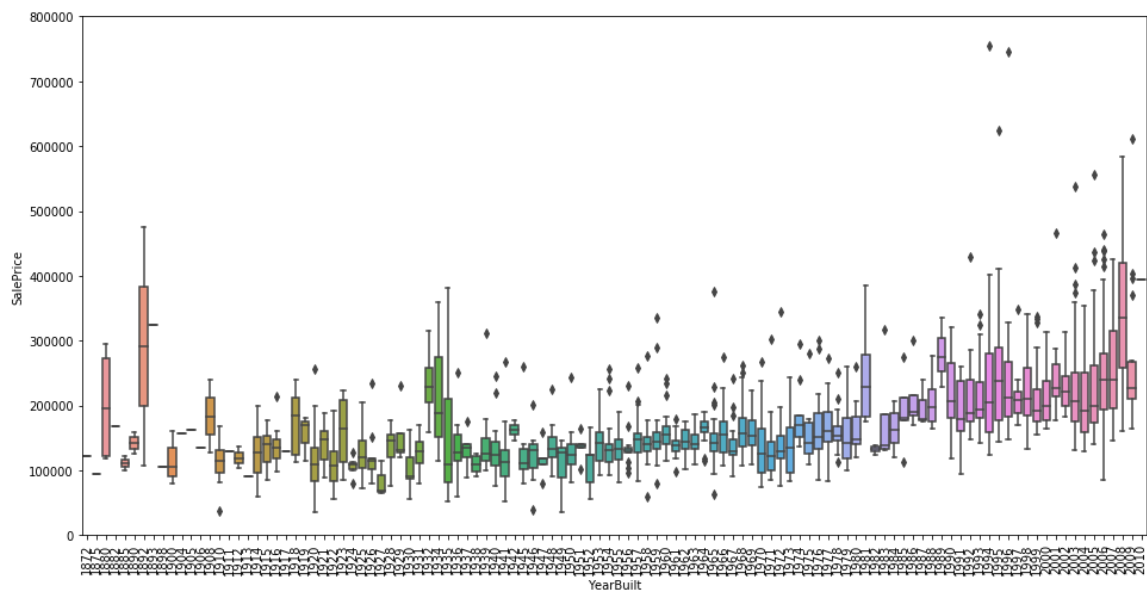
Ambas as variáveis GrLivArea e TotalBsmtSF tem ótima linearidade com a variável alvo – SalePrice.

Verificando a linearidade entre a variável alvo – SalePrice com duas variáveis categóricas.



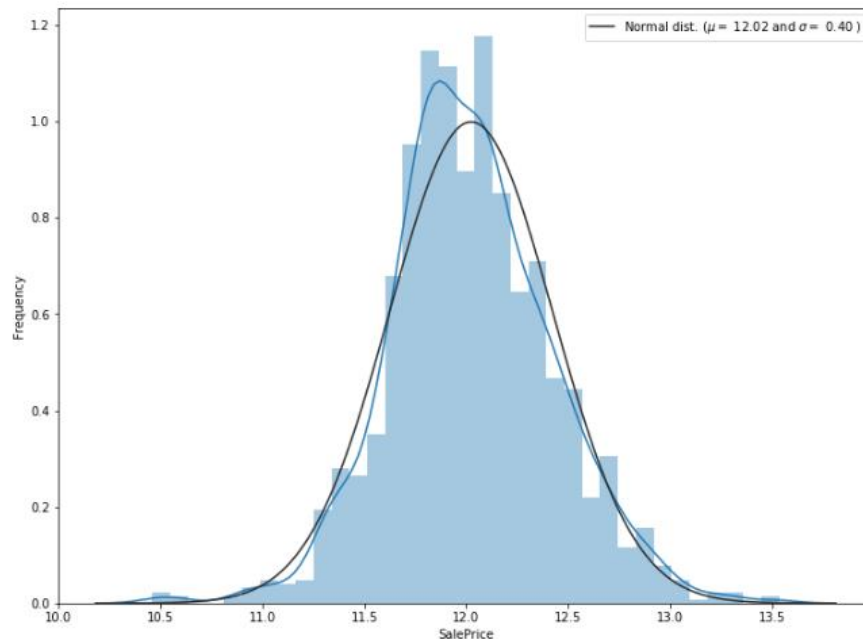
OverallQual e SalePrice possui relação, como é possível ver acima.

Agora a variável YearBuilt não mostra nenhuma relação com a variável SalePrice.

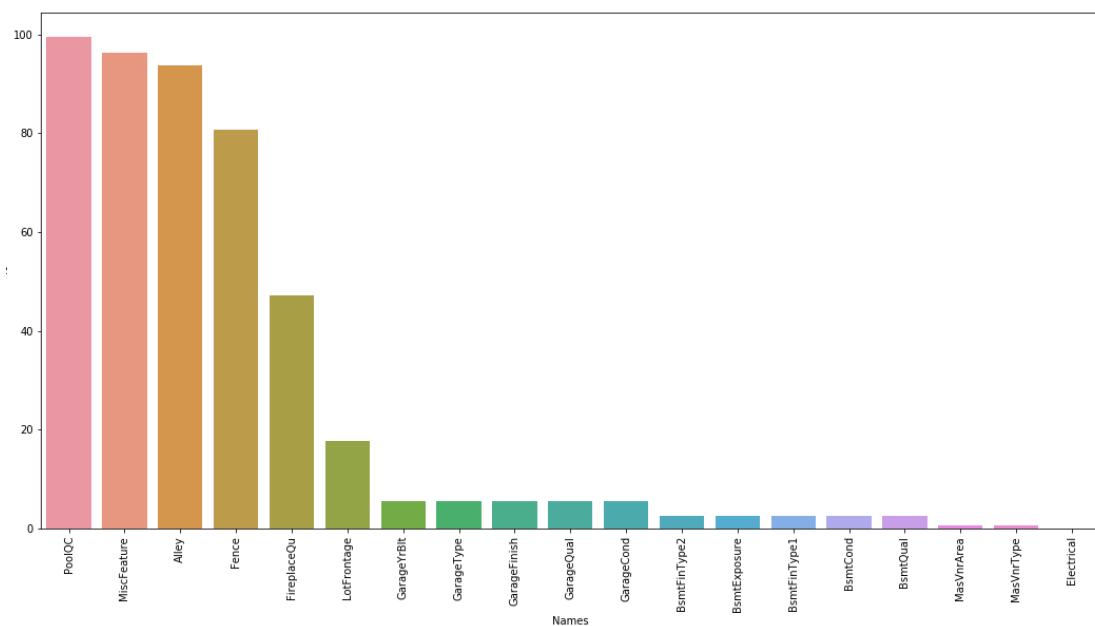


Visualização exploratória

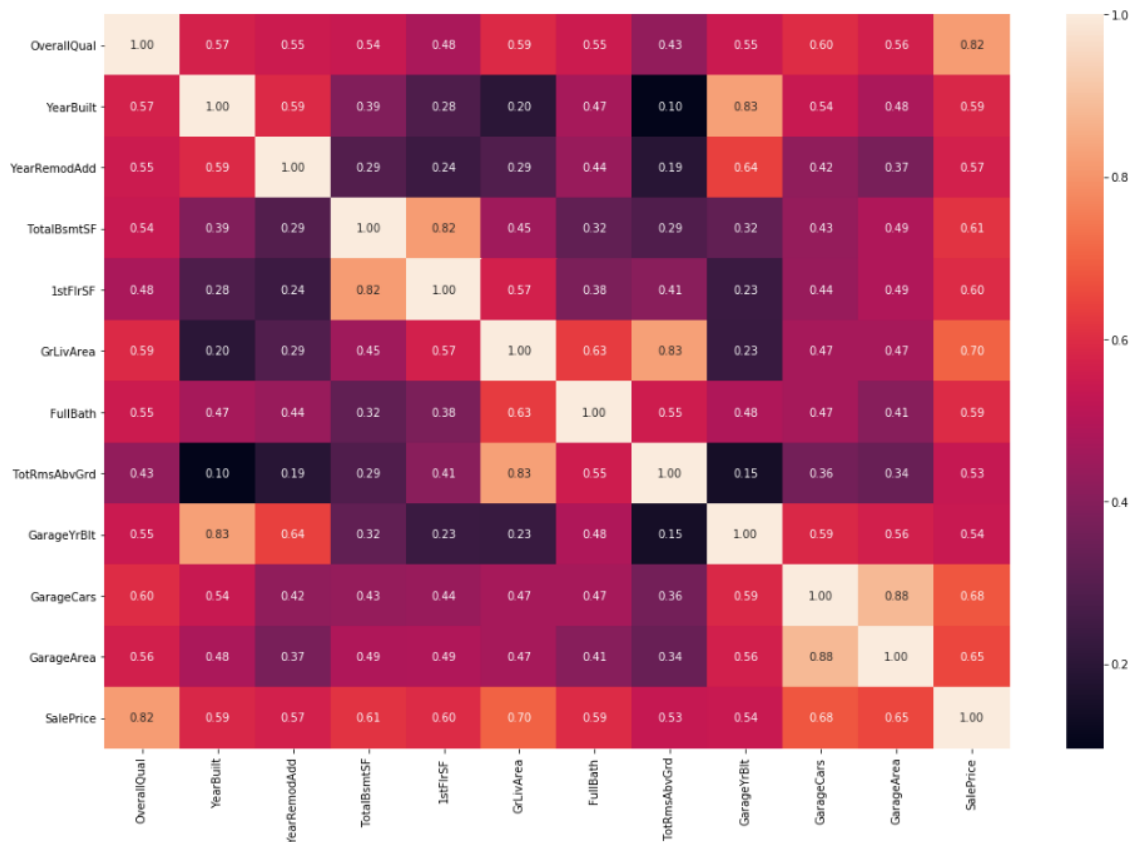
O gráfico da variável SalePrice do dataset (variável alvo, o modelo usará as características da casa para chegar no valor de venda da casa) após transformação logarítmo.



Estas são as variáveis (características da casa) do dataset que tem dados não preenchidos e com as porcentagens de dados ausentes em cada uma destas variáveis. As variáveis PoolQC, MiscFeatures e Alley tem mais 80% de dados ausentes.



Abaixo um gráfico onde mostra as correlações das características da casa que tenham correlação acima de 0.5 (50%).



Se observamos a última coluna que tem a correlação entre a SalePrice e as demais variáveis, podemos ver que as maiores correlação são:

- OverallQual – 0.82
- GrLivArea – 0.70
- GarageCars – 0.68
- GarageArea – 0.65
- TotalBsmtSF – 0.61
- 1stFlrSF – 0.60

Agora, observando a correlação entre as variáveis features, temos forte correlação entre:

- GarageArea e GarageCars – 0.88
- YearBuilt e GarageYrBlt – 0.83
- TotRmsAbvGrd e GrLivArea – 0.83
- TotalBsmtSF e 1stFlrSF – 0.82

Algoritmos e técnicas

Utilizados algoritmos de Regressão da aprendizagem de máquina supervisionada. Os algoritmos de regressão trabalhados nesta proposta são com uma pequena introdução sobre cada um destes algoritmos:

- **Lasso**

A técnica Lasso (Least Absolute Shrinkage and Selection Operator) é uma técnica de Regressão de encolhimento frequentemente usado com os dados tem multicolinearidade (quando as variáveis independentes são fortemente correlacionadas entre si).

Uma regressão com diversos coeficientes regressores torna o modelo como um todo muito mais complexo e pode tirar características de interpretabilidade. Uma forma de eliminar esse problema, que pode absorver o ruído dos dados e causar o sobreajustar (overfitting), esses métodos fazem a retenção de um subconjunto de coeficientes regressores o que não somente reduz a complexidade do modelo e a forma que o mesmo é calculado e construído, como reduz o erro e de quebra minimiza qualquer possibilidade do modelo ter sobreajustar.

A técnica Lasso usa mecanismo de penalização dos coeficientes com um alto grau de correlação entre si.

Equação

$$\min_{\alpha, \beta} \frac{1}{N} \sum_{i=1}^N f(x_i, y_i, \alpha, \beta) \text{ subject to } \|\beta\|_1 \leq t \quad [2]$$

[2] Fonte: [https://en.wikipedia.org/wiki/Lasso_\(statistics\)](https://en.wikipedia.org/wiki/Lasso_(statistics))

- **Simple e/ou Multivariate Linear Regression**

Usado para estimar valores contínuos. Assume uma relação linear entre a variável dependente (neste projeto a SalePrice) cujo dados são valores contínuos com as variáveis independentes (neste projeto as características da casa). Podendo (y) ser calculado a partir de valores de entrada (x)

Simple Linear Regression: $y = b_0 + b_1x$

Multivariate Linear Regression: $y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$, onde n = número de variáveis independentes.

Característica:

- Dependem de uma relação linear entre as variáveis independentes com a variável dependente.
- Sensíveis a Outliers.

- **Support Vector Regression (SVR)**

O SVR é uma técnica mais avançada, ela se diferencia de uma técnica de regressão linear simples tentando minimizar o erro e individualizar o hiperplano com a maximização das margens através dos vetores de suporte (linha tracejada).

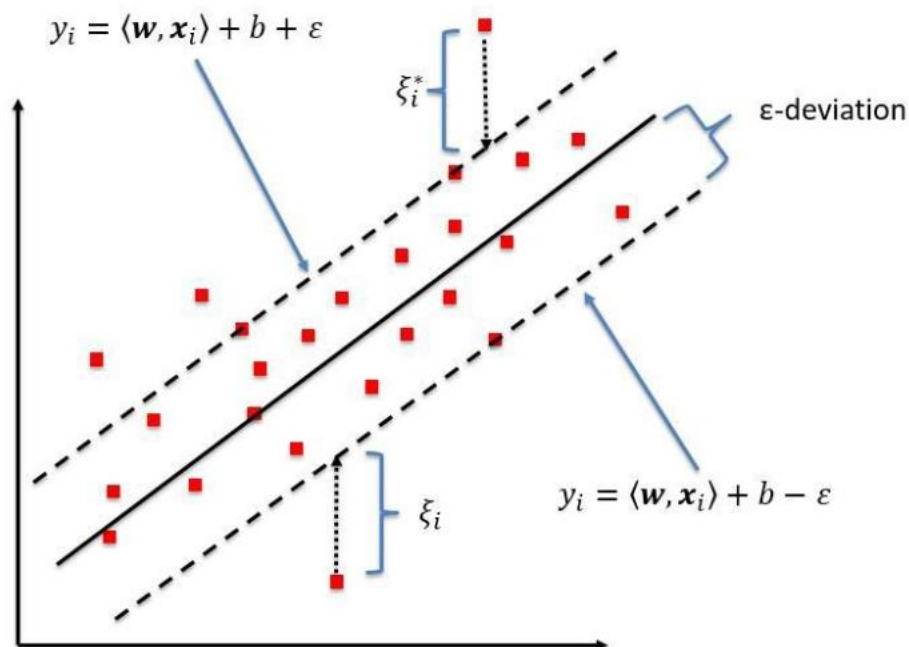


Fig 3: https://www.researchgate.net/figure/Schematic-of-the-one-dimensional-support-vector-regression-SVR-model-Only-the-points_fig5_320916953

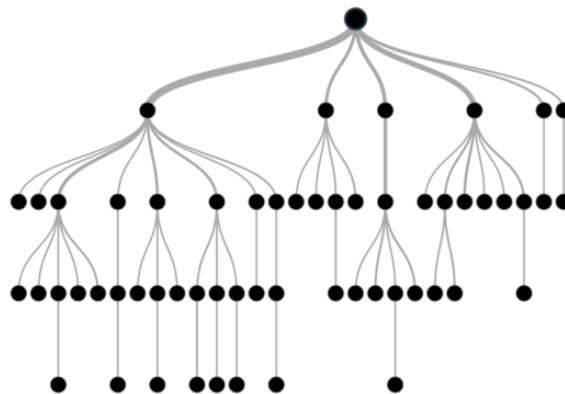
É técnica que se utiliza de Kernel transformando os dados em um espaço de maior dimensionalidade sendo possível a separação linear.

Característica:

- Não é muito influenciado por ruídos nos dados
- Mais fácil de usar do que as redes neurais
- É lento
- A implementação do algoritmo é complexa e fechada (black box).

- **Decision Tree Regression**

Árvore de decisão é um tipo de algoritmo de aprendizagem supervisionada (com uma variável alvo pré-definida), muito utilizada em problemas de classificação. Ele funciona para ambas as variáveis categóricas e contínuas de entrada e de saída. Na árvore de decisão, dividimos a população ou amostra em dois ou mais conjuntos homogêneos (ou sub-populações) com base nos divisores/diferenciadores mais significativos das variáveis de entrada.



Característica:

- Modelos preditivos de alta precisão;
- Estabilidade;
- Útil em exploração de dados;
- Facilidade de interpretação;
- Mapeiam muito bem relações não-lineares;
- Podem ser adaptados para resolver vários tipos de problema (classificação ou regressão).
- Poder facilmente se sobreajustar (overfitting)

Mais informações em português - <https://www.vooo.pro/insights/um-tutorial-completo-sobre-a-modelagem-baseada-em-tree-arvore-do-zero-em-r-python/>

- **Random Forest Regression.**

Random Forest Regression é um método de aprendizagem de máquina versátil e capaz de executar tarefas de regressão e de classificação. É um tipo de método de aprendizado de 'ensemble', onde um grupo de modelos fracos são combinados para formar um modelo mais forte.

Similar em princípio com o Decision Tree. Na floresta aleatória, crescemos múltiplas árvores ao invés de uma única árvore no modelo do CART (Classification And Regression Tree).

Característica:

- Poder de lidar com dados em grandes volumes e com muitas dimensões;
- O modelo produz o grau de importância das variáveis;
- Eficaz para estimar os dados faltantes.

Com relação a aplicação destes algoritmos, na primeira abordagem, os parâmetros utilizados nestes algoritmos são parâmetros padrões, sem nenhuma customização, para atingir um melhor resultado.

Após a definição de qual é o melhor algoritmo, irei usar o GridSearch, para chegar nos melhores parâmetros deste algoritmo.

Benchmark

Iremos criar um modelo com o algoritmo 'Naive' de Regressão para utilizar como modelo de benchmark comparando com os algoritmos de regressão mais avançados.

O resultado do algoritmo DummyRegressor que foi utilizado como algoritmo de regressão 'Naive' teve o RMSE de:

```
RMSE on Training set : 0.13071262536971964  
RMSE on Test set : 0.1269396054397447
```

III. Metodologia

Pré-processamento de dados

Após entendimentos dos dados, foram feitos os seguintes pré-processamentos:

- Tratamento dos dados ausentes dos dados.

MiscFeature = Preenchido com None (No Misc Feature).

Alley = Preenchido com None (No Alley access).

Fence = Preenchido com None (No Fence).

FireplaceQu = Preenchido com None (No FireplaceQu).

LotFrontage = Preenchido com o valor mediano do bairro.

GarageType | GarageFinish | GarageQual | GarageCond = Todas foram preenchidas com None (No Garage).

GarageYrBlt | GarageArea | GarageCars = Todas foram preenchidas com 0.

BsmtFinSF1, 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'BsmtFullBath',

'BsmtHalfBath' = Preenchido com 0.

'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2' =

Preenchido com None (No basement).

MasVnrType = Preenchido com None (masonry veneer).

MasVnrArea = Preenchido com 0 (masonry veneer).

MSZoning = Preenchido com RL.

Utilities = Eliminado do dataset (Todos os dados no dataset tem o mesmo valor – AllPub).

Functional = Preenchido com "Typ" (significa típico).

Electrical = Preenchido com o valor mais frequente 'SBrkr'.

KitchenQual = Preenchido com o valor mais frequente 'TA'.

Exterior1st = Preenchido com o valor mais frequente.

Exterior2nd = Preenchido com o valor mais frequente.

SaleType = Preenchido com o valor mais frequente 'WD'.

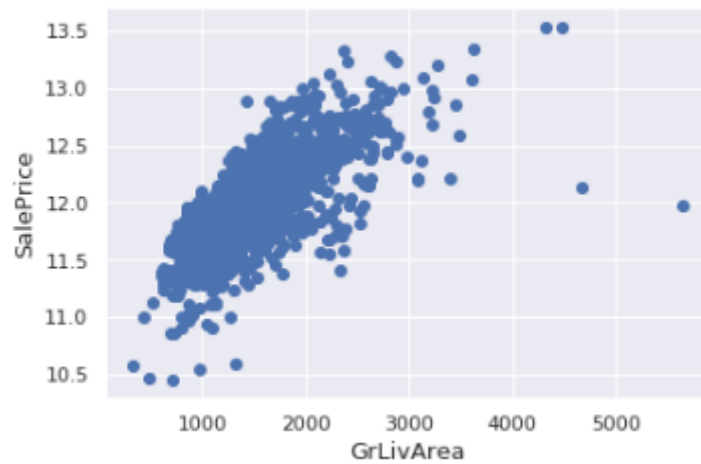
MSSubClass = Preenchido com None (No building class).

PoolQC = Eliminado do dataset.

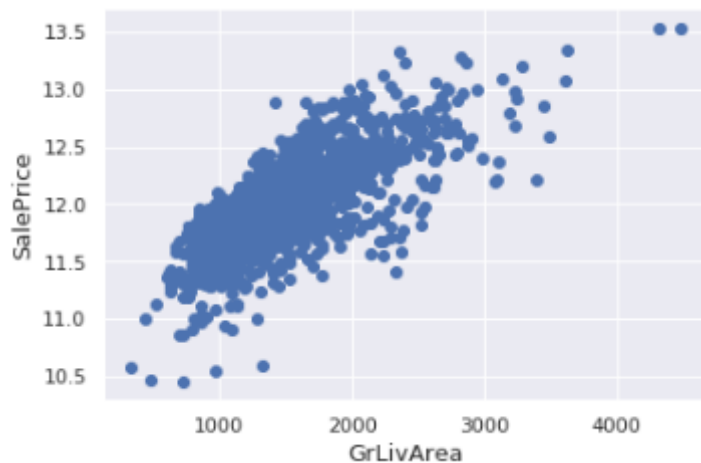
- Eliminação de Outliers observando a variável SalePrice.

Correlacionando os valores de SalePrice com GrLiveArea que são 2 variáveis numérica com maior correlação.

Antes



Depois de eliminado os 2 registros onde GrLiveArea maior que 4000 e SalePrice maior que 11.5.



- Normalização dos dados numéricos utilizando o MinMaxScaler da biblioteca Scikit Learn.
- Convertendo algumas variáveis pseudo numéricas para variáveis categóricas.

O dataset passou a ter a dimensão de 1458 registros e 333 variáveis.

- Por último, convertendo as variáveis categóricas em variáveis 'dummy', com valor 0 ou 1.

```
[27] df_train.head()
```

| | LotFrontage | LotArea | OverallQual | YearBuilt | YearRemodAdd | MasVnrArea | BsmtFinSF1 | BsmtFinSF2 | BsmtFinSF3 | TotalBsmtSF | ... | SaleType_ConvLw | SaleType_New | SaleType_Oth | SaleType_M0 | SaleCondition_Abnorml | SaleCondition_AdjAcnd | SaleCondition_Alloca | SaleCo |
|---|-------------|----------|-------------|-----------|--------------|------------|------------|------------|------------|-------------|-----|-----------------|--------------|--------------|-------------|-----------------------|-----------------------|----------------------|--------|
| 0 | 0.150606 | 0.033420 | 0.666667 | 0.949275 | 0.883333 | 0.12250 | 0.322669 | 0.0 | 0.064212 | 0.266999 | ... | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |
| 1 | 0.202055 | 0.030795 | 0.555556 | 0.753623 | 0.433333 | 0.00000 | 0.440584 | 0.0 | 0.121575 | 0.393637 | ... | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |
| 2 | 0.180959 | 0.040507 | 0.666667 | 0.934783 | 0.866667 | 0.10125 | 0.222121 | 0.0 | 0.185788 | 0.286962 | ... | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |
| 3 | 0.133562 | 0.038561 | 0.666667 | 0.311594 | 0.333333 | 0.00000 | 0.098720 | 0.0 | 0.231164 | 0.236808 | ... | 0 | 0 | 0 | 1 | 1 | 0 | 0 | |
| 4 | 0.215753 | 0.060576 | 0.777778 | 0.927536 | 0.833333 | 0.21875 | 0.299360 | 0.0 | 0.209760 | 0.357143 | ... | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |

5 rows x 333 columns

Implementação

Foi criado uma função 'train_predict' para servi de pipeline dos modelos.

```
[70] # Criação do pipeline de treino dos modelos
def train_predict(learner, X_train, y_train, X_test, y_test):
    """
    inputs:
    - learner: the learning algorithm to be trained and predicted on
    - X_train: features training set
    - y_train: income training set
    - X_test: features testing set
    - y_test: income testing set
    """
    results = {}

    # Medindo o tempo de criação do modelo
    start = time()
    learner = learner.fit(X_train, y_train)
    end = time()

    # Calculando tempo de treino
    results['train_time'] = end - start

    # Medindo o tempo de prever
    start = time()
    predictions_train = learner.predict(X_train)
    predictions_test = learner.predict(X_test)
    end = time()

    # Calculando o tempo de prever
    results['pred_time'] = end - start

    # RMSE de Treino
    results['rmse_train'] = rmse_cv_train(learner).mean()

    # RMSE de Teste
    results['rmse_test'] = rmse_cv_test(learner).mean()

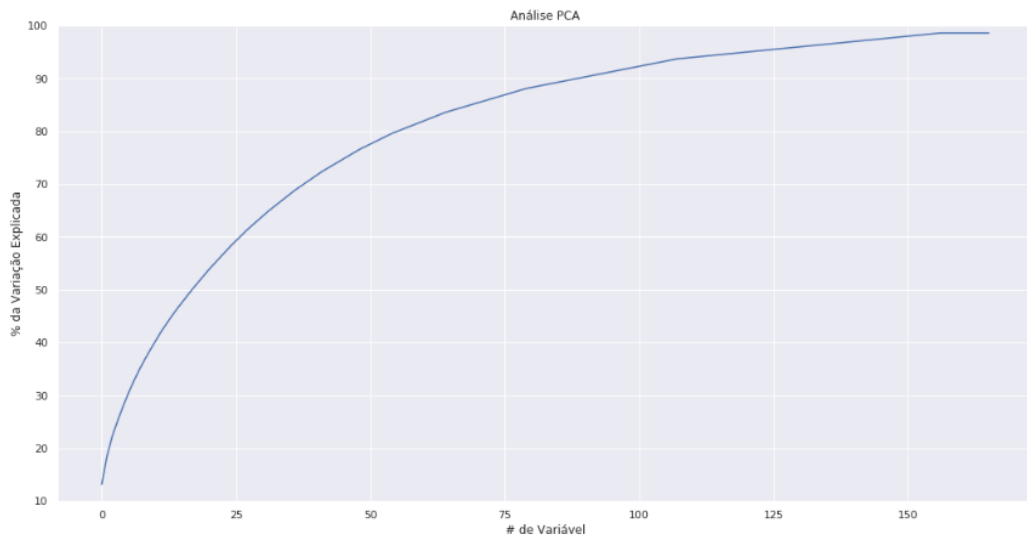
    # Sucesso
    print("Modelo {} treinou {} e para prever com o RMSE de treino de {} e o RMSE de teste {} samples.".format(learner.__class__.__name__, results['train_time'], results['pred_time'], results['rmse_train'], results['rmse_test']))

    # Return dos resultados
    return results
```

A função acima tem como objetivo rodar todos os algoritmos definidos e medir o tempo de criação (fit), o tempo de predição, além de gerar o RMSE do dataset de treino e o de teste.

Apenas para o rodar o algoritmo "Linear Regression Simple", com um única variável, foi preciso rodar fora do pipeline (função 'train_predict'), devido a necessidade de extrair apenas um variável do dataset de treino e teste, que foi escolhida a variável 'GrLivArea' por ter uma correlação mais forte com a variável SalePrice.

Os modelos foram testados utilizando o dataset todo e com o dataset reduzido utilizando a técnica de PCA (Principal Component Analysis). Para explicar 85% do preço de venda de casa, são necessários utilizar no mínimo 75 características.



Refinamento

Este são os resultados utilizando o dataset completo com as 333 variáveis.

```
[73] df_desempenho = pd.DataFrame(desempenho)
df_desempenho
```

| | DummyRegressor | Lasso | LinearRegression | SVR | DecisionTreeRegressor | RandomForestRegressor |
|-------------------|----------------|----------|------------------|----------|-----------------------|-----------------------|
| pred_time | 0.001091 | 0.003893 | 4.860640e-03 | 0.026185 | 0.004561 | 0.016517 |
| rmse_test | 0.126940 | 0.126940 | 1.936900e+10 | 0.085973 | 0.076008 | 0.074165 |
| rmse_train | 0.130713 | 0.130713 | 1.174931e+09 | 0.074627 | 0.067612 | 0.075746 |
| train_time | 0.000462 | 0.006949 | 4.780412e-02 | 0.022225 | 0.043612 | 0.410411 |

Este são os resultados utilizando o dataset reduzido com PCA, com somente as 75 principais variáveis.

```
[77] df_desempenho_pca = pd.DataFrame(desempenho_pca)
df_desempenho_pca
```

| | DummyRegressor | Lasso | LinearRegression | SVR | DecisionTreeRegressor | RandomForestRegressor |
|-------------------|----------------|----------|------------------|----------|-----------------------|-----------------------|
| pred_time | 0.000148 | 0.002028 | 2.270699e-03 | 0.009251 | 0.004313 | 0.014241 |
| rmse_test | 0.126940 | 0.126940 | 1.936900e+10 | 0.085973 | 0.076008 | 0.074165 |
| rmse_train | 0.130713 | 0.130713 | 1.174931e+09 | 0.074627 | 0.067612 | 0.075746 |
| train_time | 0.001254 | 0.003347 | 7.127047e-03 | 0.008519 | 0.067811 | 0.930622 |

Comparando o tempo e o RMSE dos algoritmos com seus parâmetros padrões entre o dataset completo (df_desempenho) e o dataset reduzido (df_desempenho_pca com PCA), concluo que não há necessidade de utilização do PCA, como pode ver abaixo. [1]

```
[78] resultado = df_desempenho - df_desempenho_pca
      resultado
```

| | DummyRegressor | Lasso | LinearRegression | SVR | DecisionTreeRegressor | RandomForestRegressor |
|------------|----------------|----------|------------------|----------|-----------------------|-----------------------|
| pred_time | 0.000021 | 0.001795 | 0.001669 | 0.015383 | 0.001900 | 0.001714 |
| rmse_test | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| rmse_train | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| train_time | -0.000359 | 0.004143 | 0.049516 | 0.012399 | -0.021146 | -0.516016 |

Os valores 0.000000 significa que não houve diferença da métrica RMSE utilizando os algoritmos entre o dataset completo e o dataset reduzido.

[1] O jupyter notebook foi computado utilizando o serviço da Google Colab - <https://colab.research.google.com>. Pode ser que rodando este notebook em notebooks, possam obter outros resultados e outra conclusão.

Utilizado o GridSearch para definir os melhores parâmetros do **algoritmo que mais desempenho no dataset de teste.**

IV. Resultados

Modelo de avaliação e validação

Os parâmetros padrões utilizado para cada algoritmo foram:

```
# Inicialize os modelos
clf_A = DummyRegressor()
clf_B = linear_model.Lasso(alpha=0.1, random_state = 40)
clf_C = LinearRegression()
clf_D = SVR(gamma='scale', C=1.0, epsilon=0.2)
clf_E = DecisionTreeRegressor(random_state=40)
clf_F = RandomForestRegressor(max_depth=2, n_estimators=100, random_state=40)
```

O resultado dos desempenhos dos algoritmos com seus parâmetros padrões foram:


```
[73] df_desempenho = pd.DataFrame(desempenho)
df_desempenho
```

| | DummyRegressor | Lasso | LinearRegression | SVR | DecisionTreeRegressor | RandomForestRegressor |
|------------|----------------|----------|------------------|----------|-----------------------|-----------------------|
| pred_time | 0.000161 | 0.003680 | 4.600763e-03 | 0.024902 | 0.004576 | 0.016091 |
| rmse_test | 0.126940 | 0.126940 | 1.936900e+10 | 0.085973 | 0.076008 | <u>0.074165</u> |
| rmse_train | 0.130713 | 0.130713 | 1.174931e+09 | 0.074627 | 0.067612 | 0.075746 |
| train_time | 0.000899 | 0.006837 | 5.665088e-02 | 0.021402 | 0.043694 | 0.422792 |

Os parâmetros utilizados para melhorar o algoritmo Random Forest Regression foram:

```
[84] parameters = {
    'n_estimators' : [100, 150, 300],
    'max_depth' : [3, 5, 10],
    'min_samples_split' : [3, 5, 10]
}
```

Estes parâmetros foram automatizados utilizando com GridSearch para melhorar os resultados RMSE do Random Forest Regression, o RMSE ficou:

Unoptimized model

Score on Test set : 0.6495593292427022

RMSE on Test set : 0.07416484852022152

Optimized Model

Score on Test set : 0.859131653368277

RMSE on Test set : 0.053563052292514804

Uma redução de ~28% do RMSE.

Justificativa

Comparando com o RMSE do benchmark, o qual foi utilizado o 'DummyRegressor' que se tinha os seguintes resultados:

RMSE on Training set : 0.13071262536971964

RMSE on Test set : 0.1269396054397447

Com os resultados obtidos com Random Forest Regression

RMSE on Test set : 0.053563052292514804

Temos uma melhora de ~58% em redução do erro calculado pelo RMSE.

Recomendamos a utilização do modelo de aprendizagem de máquina com o algoritmo Random Forest Regression para prever o valor de venda de casa em Ames / Iowa nos Estados Unidos com base nas 79 características da casa.

V. Conclusão

Forma livre de visualização

```
[78] resultado = df_desempenho - df_desempenho_pca
      resultado
```

| | DummyRegressor | Lasso | LinearRegression | SVR | DecisionTreeRegressor | RandomForestRegressor |
|------------|----------------|----------|------------------|----------|-----------------------|-----------------------|
| pred_time | 0.000085 | 0.001439 | -0.003647 | 0.016572 | 0.002398 | 0.001805 |
| rmse_test | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| rmse_train | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| train_time | -0.000666 | 0.004936 | 0.041267 | 0.012957 | -0.024661 | -0.514149 |

Na visualização acima é uma comparação de desempenho (tempo) e eficiência (RMSE), como pode ver, com 75 características utilizando o PCA, pode se checar ao mesmo nível de eficiência quando utilizados as 333 características, os valores 0.000000 significa que não houve diferença da métrica RMSE utilizando os algoritmos entre o dataset completo e o dataset reduzido.

Para treinar os algoritmos DummyRegressor (Benchmark), DecisionRegressor e RandomForestRegressor, o tempo de ficou ligeiramente maior. E somente o Linear Regression que o tempo de prever maior com o dataset reduzido (PCA), lembrando que o LinearRegression foi o pior algoritmo para este dataset.

Reflexão

O processo seguindo neste projeto pode ser resumido assim:

- Definição do problema.
- Coleta de dados sobre o problema.
- Entendimento dos dados.
- Definição do benchmark para avaliação do modelo proposto posteriormente e qual seria a métrica de comparação.
- Pré processamento dos dados, trabalhar com os ausentes, outliers, correlação e até a redução de características.
- Pré Treinar o modelo e avaliar o resultado.
- Comparar o resultado e definir o melhor algoritmo de regressão para resolver este problema.
- Trabalhar com os melhores parâmetros do algoritmo escolhido.
- Comparação final.

Melhorias

Entendo que este projeto poderia ser enriquecido o treino de outros algoritmos de Regressão, tais como: Ridge, Elastic-Net, SGD (Stochastic Gradiente Descent) e outros.

O uso de modelo de Classificação onde transformaria a variável alvo em categórica, ou seja, preço de venda entre 100K-150K, 150K-200K, 200K-250K.... e utilizar modelos de classificação ao invés de Regressão poderia ser verificado. Como já realizado no estudo de Hujia Yu.[3]

[3] De Hujia Yu, Jiafu Wu. "Real Estate Price Prediction with Regression and Classification." CS 229 Autumn 2016 Project Final Report - http://cs229.stanford.edu/proj2016/report/WuYu_HousingPrice_report.pdf