

# Análise Comparativa: *Lung Cancer*

---



Bárbara da S. Oliveira, Carlos A. de S. Monteiro,  
Carlos A. M. de Pinho, Hugo S. Sousa,  
Larissa A. Barbosa, Luiz H. M. de Souza,  
Matheus A. Constancio e Rafael G. da Silva  
7 de novembro de 2022

***Data Sillers***

# CONTEÚDO

1. Introdução
2. Pré-processamento
3. Metodologia e Experimento
4. Conclusão

# Introdução

---

# ANÁLISE COMPARATIVA

- O objetivo dessa etapa é comparar modelos para encontrar um que possa ser melhor utilizado para a detecção de câncer pulmonar
- Para isso, precisamos antes realizar a preparação e o pré-processamento dos dados

# DADOS ORIGINAIS

Variável	Tipo
GENDER	Nominal
AGE	Discreta
SMOKING	Nominal
YELLOW_FINGERS	Nominal
ANXIETY	Nominal
PEER_PRESSURE	Nominal
CHRONIC DISEASE	Nominal
FATIGUE	Nominal
ALLERGY	Nominal
WHEEZING	Nominal
ALCOHOL CONSUMING	Nominal
COUGHING	Nominal
SHORTNESS OF BREATH	Nominal
SWALLOWING DIFFICULTY	Nominal
CHEST PAIN	Nominal
LUNG_CANCER	Nominal

**Tabela 1:** Variáveis e seus tipos

# **Pré-processamento**

---

# TRATAMENTO DE DADOS

```
1 # Criar Pipeline para cada tipo de variavel
2
3 # Para variaveis nominais
4 nominal_preprocessor = Pipeline([
5     # Dados discrepantes
6     ('missing', SimpleImputer(strategy='most_frequent')),
7     # Dados faltantes
8     ('encoder', OneHotEncoder(sparse=False)),
9     # Normalizacao
10    ('normalization', MinMaxScaler()),
11 ])
12
13 # Para variaveis discretas
14 discrete_preprocessor = Pipeline([
15     ('missing', SimpleImputer(strategy='mean')),
16     ('normalization', MinMaxScaler())
17 ])
```

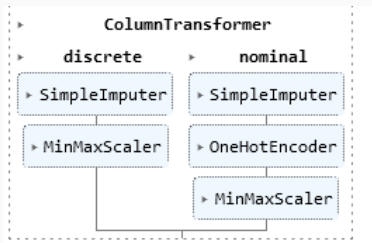
# TRATAMENTO DE DADOS

- ***SimpleImputer***: Substitui qualquer valor faltante com um valor estático (média, moda, mediana), neste caso, a média
- ***OneHotEncoder***: Encoda as características dos nossos dados em um *array* numérico, onde cada valor representará um dos dados
- ***MinMaxScaler***: Normaliza os dados escalando eles para ficarem entre uma determinada faixa de valores



# TRATAMENTO DE DADOS

```
1 # Agora podemos utilizar os pre processadores
2 # Criados acima para transformar os dados
3 preprocessing = ColumnTransformer([
4     ("discrete", discrete_preprocessor, discrete_columns),
5     ("nominal", nominal_preprocessor, nominal_columns)
6 ])
```



**Figura 1:** Transformação das variáveis

# **Metodologia e Experimento**

---

## ESCOLHA DO MODELO

Iremos analisar quatro modelos de aprendizado de máquina para solucionar nosso problema:

- ***Logistic Regression*** (LR)
- ***K-Nearest-Neighbors*** (KNN)
- ***Support Vector Machine*** (SVM)
- ***Naive Baiyes*** (NB)

Para encontrarmos, dentro desses modelos, a melhor configuração viável para o nosso problema, utilizaremos a classe *GridSearch* do *SickitLearn*.

# CONFIGURAÇÃO DO EXPERIMENTO

```
1 # Nome do modelo, Chamada do metodo e Parametros de Teste
2 models = [
3     ("LR", LogisticRegression(solver='saga', max_iter=1000),
4     {"penalty": ['none', 'l1', 'l2']} ),
5
6     ("KNN", KNeighborsClassifier(metric='euclidean'),
7     {"n_neighbors": np.arange(1, 31, 2), 'weights': ["uniform", "
8         distance"]}),
9
10    ("SVM", SVC(max_iter=10000),
11    {'C':[1, 10, 100, 1000], 'gamma':[1, 0.1, 0.001, 0.0001], '
12        kernel' ['linear', 'rbf']})),
13 ]
```

## CRITÉRIOS DE AVALIAÇÃO

Os modelos serão comparados através dos seguintes parâmetros:

- ***accuracy***: Proporção entre os dados que foram corretamente previstos (como positivos ou negativos) com o total de dados observados
- ***precision***: Proporção entre dados corretamente previstos como positivos e o total de observações positivas
- ***recall***: Proporção entre dados corretamente previstos como positivos com o total de observações. Em outras palavras, esse parâmetro ajuda a identificar a sensibilidade do nosso modelo
- ***f1***: Média ponderada entre *precision* e *recall*, portanto levando em conta tanto falsos positivos quanto falsos negativos

# MÉTODO DE VALIDAÇÃO

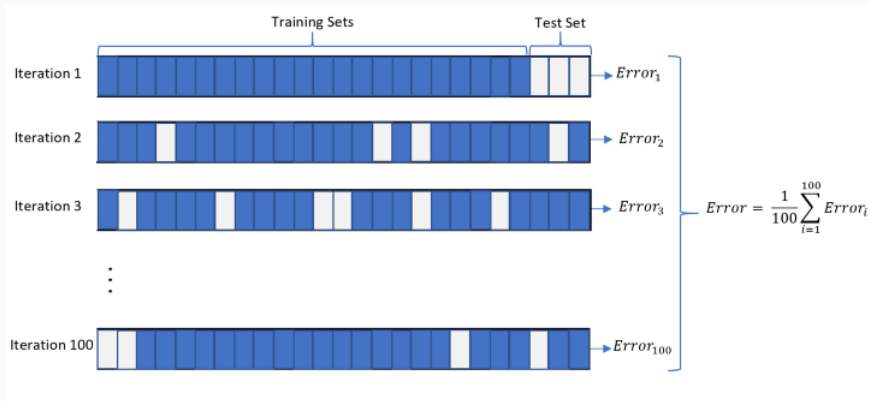
Iremos separar o conjunto de dados em dois conjuntos: **teste** e **treino** para realizar a validação cruzada.

```
1 X = df.drop(columns=[target_column], axis=1)
2 y = (
3     df[[target_column]]
4     .replace({"YES": 1, "NO": 0})
5     .to_numpy()
6     .ravel()
7 )
8
9 #Separa em conjuntos de teste e treino
10 cv = ShuffleSplit(n_splits=30, train_size=0.8, random_state=42)
```

# TESTE DO MODELO

- Para a verificação do modelo, foi utilizado o *Cross Validation* ou Validação Cruzada, onde há uma divisão entre os dados para treinamento
- Mais especificamente, foi utilizado o método de Monte Carlo, onde a divisão feita entre os dados é arbitrária, nesse caso, **80%** foram utilizados para o treinamento do modelo
- Além disso, fazemos essa separação 30 vezes, para que o algoritmo nunca realize o treino nos mesmos valores previamente selecionados

# MONTE CARLO



**Figura 2:** Representação do Método de Monte Carlo



# MÉTODO DE VALIDAÇÃO

```
1 results = {}
2 for model_name, model, model_params in models:
3     model_gs = GridSearchCV(model, model_params, scoring='
4         accuracy')
5     approach = Pipeline([
6         ("preprocessing", preprocessing),
7         ("model", model_gs)
8     ])
9     model_results = cross_validate(
10         approach,
11         X=X,
12         y=y,
13         # Critérios de avaliação
14         scoring=['accuracy', 'f1', 'precision', 'recall'],
15         cv=cv,
16         n_jobs=-1,
17         return_train_score=False
18     )
```

# Conclusão

---

# RESULTADOS

score	KNN	LR	NB	SVM
fit_time	$0.704 \pm 0.141$	$0.351 \pm 0.054$	$0.053 \pm 0.008$	$1.314 \pm 1.039$
score_time	$0.014 \pm 0.007$	$0.011 \pm 0.004$	$0.013 \pm 0.007$	$0.014 \pm 0.009$
test_accuracy	$0.904 \pm 0.035$	$0.930 \pm 0.034$	$0.912 \pm 0.033$	$0.923 \pm 0.033$
test_f1	$0.946 \pm 0.021$	$0.960 \pm 0.020$	$0.949 \pm 0.020$	$0.956 \pm 0.020$
test_precision	$0.933 \pm 0.036$	$0.948 \pm 0.034$	$0.962 \pm 0.033$	$0.948 \pm 0.033$
test_recall	$0.960 \pm 0.021$	$0.974 \pm 0.018$	$0.937 \pm 0.023$	$0.966 \pm 0.019$

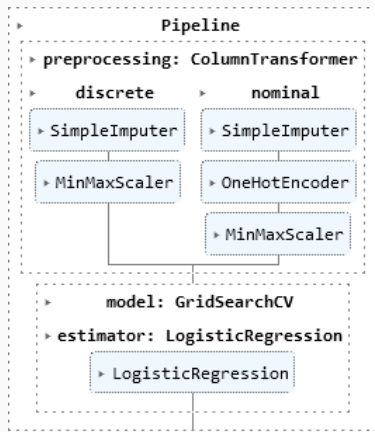
**Tabela 2:** Resultados obtidos

## PERSISTÊNCIA DO MODELO

Com isso, definimos que o melhor modelo é o de **Regressão Logística**, portanto podemos obter os melhores parâmetros desse modelo e salvar esse modelo em disco para utilizar na próxima fase da análise.

```
1 #Obtem o modelo e os parametros ganhadores
2 model_name, model, model_params = [foo for foo in models if foo
    [0] == winner][0]
3
4 #Obtem o melhor modelo denovo pelo GridSearchCV
5 model_gs = GridSearchCV(model, model_params, scoring='accuracy')
6 approach = Pipeline([
7     ("preprocessing", preprocessing),
8     ("model", model_gs)
9 ])
10 approach.fit(X, y) #Seleciona o approach
```

# PERSISTÊNCIA DO MODELO



**Figura 3:** Modelo Escolhido

## PRÓXIMOS PASSOS

Com o modelo selecionado e treinado, agora podemos:

- Construir um formulário para obter dados de usuários
- Aplicar o modelo escolhido nesses dados para prever a possibilidade de câncer pulmonar
- Realizar o *deploy* desse formulário online (utilizando o **Streamlit** e **StreamlitShare**)