

YULAB tecnologia



Workshop Machine Learning utilizando Python

Guia didático

Imaginar, desenvolver, construir, aproximar, integrar, amadurecer e evoluir! Se você teve uma grande ideia, provavelmente passou por este ciclo, é a melhor forma de criar um sucesso. Com essas palavras em mente, construímos a Bemaker! Um makerspace que nasceu da vontade de impulsionar o movimento *Do It Yourself* (Faça você mesmo) para qualquer esfera da sociedade. Nossa missão é disponibilizar um espaço com ferramentas para que você construa seus sonhos! Venha fazer parte deste movimento!

Autores:

- Adelino Pinheiro Silva (adelino@yulab.com.br)
- Fabiano Calado (fabiano@yulab.com.br)
- Lucas Calado Alves Pereira (kalado@yulab.com.br)
- Thiago Amaral (thiago@yulab.com.br)



Contextualização

Este é um guia didático com objetivo de introduzir as técnicas de aprendizado de máquina e aplicá-las na prática. Este material é uma adaptação do guia de Helge Bjorland - [An Interactive Data Science Tutorial](#).

A ferramenta utilizada será o [Spyder](#) (que pode ser baixado no link) que roda em ambiente [Python](#). As demais ferramentas são os pacotes da biblioteca [SciPy](#).

O código completo para este workshop pode ser obtido [neste link](#).

1.1 Ambiente Spyder

O Spyder é uma interface de desenvolvimento integrada (IDE) para programação em Python que agrega algumas funcionalidades como a interpretação interativa.

A tela principal do Spyder pode ser visualizada na Figura 1.1.

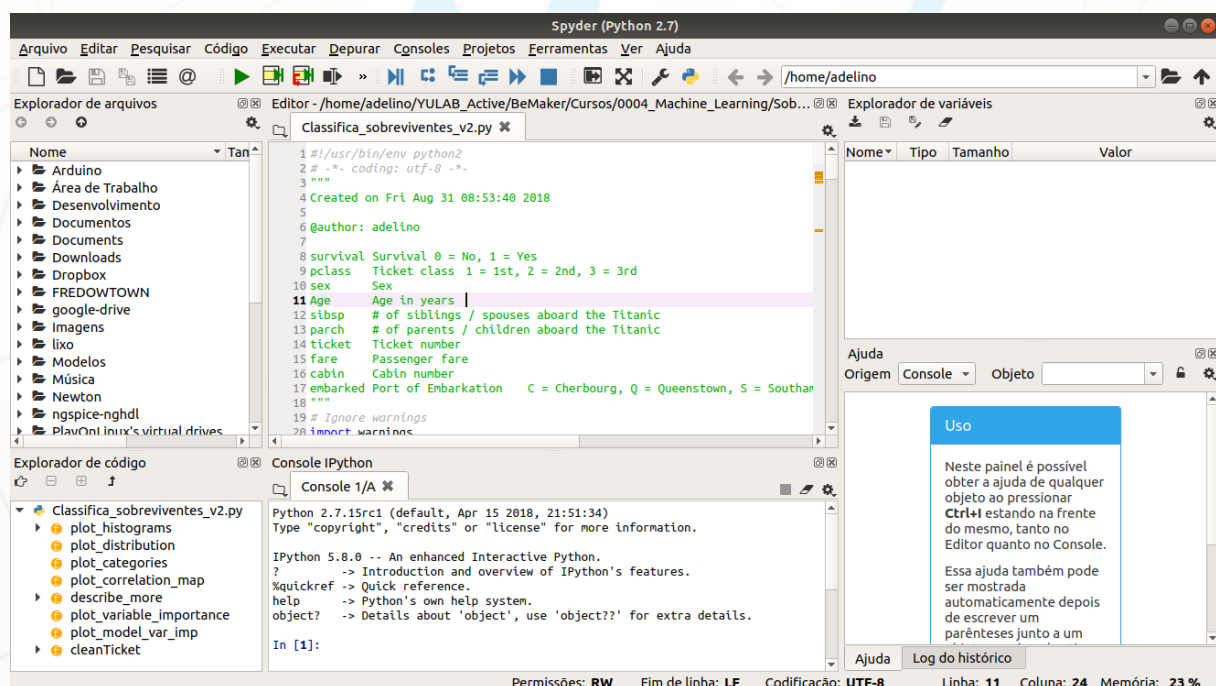


Figura 1.1: Tela de trabalho do Spyder como *layout* do Matlab.

Na tela do Spyder tem-se várias funcionalidades. As mais utilizadas são: a janela para o código, que é interpretado de forma interativa pelo Ipython; a lista de variáveis e o console (ou *prompt*) de comando.

Uma dica importante, para que os gráficos gerados apareçam em uma nova janela, e não no console siga o seguinte menu e altere a variável `Graphics backend` para `Automatic`.

Tools > preferences > IPython console > Graphics > Graphics backend > Backend: Automatic

Em seguida reinicie o *kernel* do Ipython com `Ctrl+..`



Colocando a mão na massa

Nesta etapa vamos colocar a mão na massa e programar no ambiente Spyder. O objetivo principal é utilizar as principais ferramentas de aprendizado de máquina sem se aprofundar na teoria.

2.1 Compreensão e Visualização dos Dados

2.1.1 Importando as bibliotecas

O primeiro passo é preparar o ambiente para manipular os dados. O bloco de código abaixo carrega as principais bibliotecas que serão utilizadas no desenvolvimento.

```
1 # Comando para ignorar mensagens de warnings
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 # Bibliotecas para lidar com tabelas e matrizes
6 import numpy as np
7 import pandas as pd
8
9 # Algoritmos de modelagem
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.naive_bayes import GaussianNB
14 from sklearn.svm import SVC, LinearSVC
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.ensemble import GradientBoostingClassifier
17 from sklearn.neural_network import MLPClassifier
18
19 # Algoritmos auxiliares para modelagem de dados
20 from sklearn.preprocessing import Imputer, Normalizer, scale
21 from sklearn.cross_validation import train_test_split, StratifiedKFold
22 from sklearn.feature_selection import RFECV
23
24 # Bibliotecas para visualizacao de dados
25 import matplotlib as mpl
26 import matplotlib.pyplot as plt
27 import matplotlib.pylab as pylab
28 import seaborn as sns
```

```

29 # Fechar todas janelas abertas
30 plt.close(" all ")
31
32 # Configuracoes de visualizacao
33 mpl.style.use( 'ggplot' )
34 sns.set_style( 'white' )
35 pylab.rcParams[ 'figure.figsize' ] = 8 , 6

```

2.1.2 Definindo funções auxiliares

A linguagem Python também permite fazer o encapsulamento de blocos de códigos. O tutorial original não explica o funcionamento das funções.

Basicamente são funções para visualização de dados.

```

1
2
3 def plot_histograms( df , variables , n_rows , n_cols ):
4     fig = plt.figure( figsize = ( 16 , 12 ) )
5     for i, var_name in enumerate( variables ):
6         ax=fig.add_subplot( n_rows , n_cols , i+1 )
7         df[ var_name ].hist( bins=10 , ax=ax )
8         ax.set_title( 'Skew: ' + str( round( float( df[ var_name ].skew() ) , ) ) ) #
9         ax.set_xticklabels( [] , visible=False )
10        ax.set_yticklabels( [] , visible=False )
11        fig.tight_layout() # Improves appearance a bit.
12        plt.show()
13
14 def plot_distribution( df , var , target , **kwargs ):
15     row = kwargs.get( 'row' , None )
16     col = kwargs.get( 'col' , None )
17     facet = sns.FacetGrid( df , hue=target , aspect=4 , row = row , col = col )
18     facet.map( sns.kdeplot , var , shade= True )
19     facet.set( xlim=( 0 , df[ var ].max() ) )
20     facet.add_legend()
21
22 def plot_categories( df , cat , target , **kwargs ):
23     row = kwargs.get( 'row' , None )
24     col = kwargs.get( 'col' , None )
25     facet = sns.FacetGrid( df , row = row , col = col )
26     facet.map( sns.barplot , cat , target )
27     facet.add_legend()
28
29 def plot_correlation_map( df ):
30     corr = titanic.corr()
31     _ , ax = plt.subplots( figsize =( 12 , 10 ) )
32     cmap = sns.diverging_palette( 220 , 10 , as_cmap = True )
33     _ = sns.heatmap(
34         corr ,
35         cmap = cmap ,
36         square=True ,
37         cbar_kws={ 'shrink' : .9 } ,
38         ax=ax ,
39         annot = True ,

```

```

40     annot.kws = { 'fontsize' : 12 }
41 )
42
43 def describe_more( df ):
44     var = [] ; l = [] ; t = []
45     for x in df:
46         var.append( x )
47         l.append( len( pd.value_counts( df[ x ] ) ) )
48         t.append( df[ x ].dtypes )
49     levels = pd.DataFrame( { 'Variable' : var , 'Levels' : l , 'Datatype' : t } )
50     levels.sort_values( by = 'Levels' , inplace = True )
51     return levels
52
53 def plot_variable_importance( X , y ):
54     tree = DecisionTreeClassifier( random_state = 99 )
55     tree.fit( X , y )
56     plot_model_var_imp( tree , X , y )
57
58 def plot_model_var_imp( model , X , y ):
59     imp = pd.DataFrame(
60         model.feature_importances_ ,
61         columns = [ 'Importance' ] ,
62         index = X.columns
63     )
64     imp = imp.sort_values( [ 'Importance' ] , ascending = True )
65     imp[ : 10 ].plot( kind = 'barh' )
66     print (model.score( X , y ))

```

2.1.3 Carregando dados

Está é uma etapa importante para qualquer tarefa de aprendizado de máquina. Os dados são toda informação disponível para superar os desafios.

Basicamente existem dois conjuntos de dados, os de treinamento e os de teste.

O conjunto de treinamento apresenta dos dados do problema, ou seja, as informações de entrada, e o resultado esperado. Este conjunto permite validar como o aprendizado de máquina esta calibrado.

Com os dados de teste verifica-se o poder de generalização do aprendizado de máquina.

O bloco de código a seguir é utilizado para carregar e fazer uma visualização prévia nos dados.

```

1 #get titanic & test csv files as a DataFrame
2 train = pd.read_csv("train.csv")
3 test   = pd.read_csv("test.csv")
4
5 full = train.append( test , ignore_index = True )
6 titanic = full[ :891 ]
7
8 del train , test
9

```



```
10 | print ( 'Datasets:' , 'full:' , full.shape , 'titanic:' , titanic.shape )
```

Para se familiarizar com os dados é possível acessar algumas linhas da tabela com o comando `titanic.head()`.

Basicamente os dados apresentam as seguintes informações:

1. **survival:** Indica se o passageiro foi um sobrevivente ou não, se o valor for 0 é um não sobrevivente e se o valor for 1 é um sobrevivente.
2. **pclass:** é a classe que a passagem dá acesso. Neste caso tem-se passagens de 1^a, 2^a e 3^a classe.
3. **Name:** Nome do passageiro.
4. **sex:** sexo do passageiro, e uma variável categórica, sendo `male` ou `female`.
5. **Age:** idade em anos.
6. **sibsp:** número de irmãos ou cônjuges a bordo do Titanic.
7. **parch:** número de pais ou filhos a bordo do Titanic.
8. **ticket:** número da passagem.
9. **faer:** valor da tarifa paga pelo passageiro.
10. **cabin:** Número da cabine.
11. **embarked:** Variável categoria que indica o porto em que o passageiro embarcou, sendo C para Cherbourg, Q para Queenstown e S para Southampton.

Uma variável numérica é aquela que apresenta valores inteiros ou números reais, enquanto a variável categórica assume apenas um número limitado, e geralmente fixo, de valores possíveis, como por exemplo o tipo sanguíneo ou o sexo.

2.1.4 Visualização de variáveis numéricas

A sequência de código a seguir permite verificar como as variáveis numéricas se distribuem.

```
1 | plot_correlation_map( titanic )
2 | plt.savefig( 'V3_Correlacao.png' )
3 |
4 | # Plot distributions of Age of passengers who survived or did not survive
5 | plot_distribution( titanic , var = 'Age' , target = 'Survived' , row = 'Sex' )
6 |
7 | # Exercise 1
8 | # Plot survival rate by Embarked
9 | plot_categories( titanic , cat = 'Embarked' , target = 'Survived' )
```

O desafio desta etapa é verificar as variáveis correlacionadas e entender como elas podem explicar o problema.

Exercício 1: Investigando as taxas de sobrevivência

Nesta etapas vamos o primeiro exercício, gerar gráficos que indique os índices de sobrevivência de acordo com as variáveis do problema.

Este é um exemplo de figura que pode ser gerada.

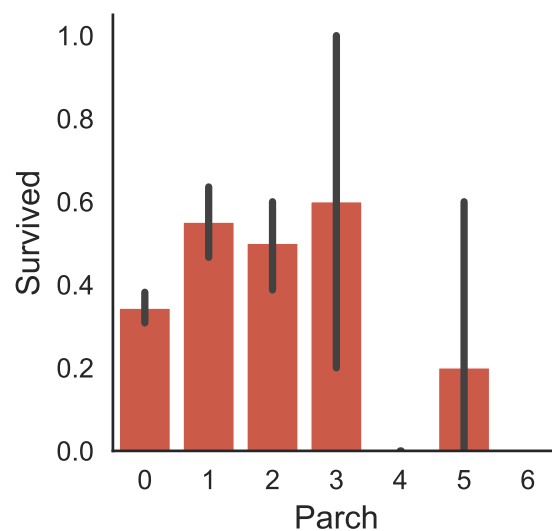


Figura 2.1: Taxa de sobrevivência de acordo com o numero de pais e filhos a bordo.

2.1.5 Visualização de variáveis categóricas

As variáveis Embarked, Pclass e Sex são tratadas nos dados como variáveis categóricas. Como os algoritmos recebem variáveis numéricas, faz-se necessário criar variáveis numéricas relacionadas às variáveis categóricas.

Desta forma, cada variável poderá assumir os valores 0 ou 1 se forem inseridas na em determinada categoria. Por exemplo, para o sexo, cria-se uma variável chamada male e o “valor” será codificado como uma variável binária (0 ou 1).

```
1 # Transform Sex into binary values 0 and 1
2 sex = pd.Series( np.where( full.Sex == 'male' , 1 , 0 ) , name = 'Sex' )
3
4 # Create a new variable for every unique value of Embarked
5 embarked = pd.get_dummies( full.Embarked , prefix='Embarked' )
6 embarked.head()
7
8 # Create a new variable for every unique value of Embarked
9 pclass = pd.get_dummies( full.Pclass , prefix='Pclass' )
10 pclass.head()
```


2.2 Preparação dos Dados

2.2.1 Variáveis ausentes

A maioria dos algoritmos de aprendizado de máquina exige que todas as variáveis tenham valores para usá-lo no treinamento do modelo. O método mais simples é preencher os valores ausentes com a média da variável em todas as observações no conjunto de treinamento.

```
1 # Create dataset
2 imputed = pd.DataFrame()
3 # Fill missing values of Age with the average of Age (mean)
4 imputed[ 'Age' ] = full.Age.fillna( full.Age.mean() )
5 # Fill missing values of Fare with the average of Fare (mean)
6 imputed[ 'Fare' ] = full.Fare.fillna( full.Fare.mean() )
7 imputed.head()
```

2.3 Engenharia de Características/Recursos - Criação de novas variáveis

2.3.1 Títulos dos nomes dos passageiros

Títulos refletem status social. No caso do Titanic, estes títulos podem prever probabilidade de sobrevivência?

Para verificar esta intuição podemos criar as variáveis referentes aos títulos.

```
1 title = pd.DataFrame()
2 # we extract the title from each name
3 title[ 'Title' ] = full[ 'Name' ].map( lambda name: name.split( ',' )[1].split( '.' )[0] )
4
5 # a map of more aggregated titles
6 Title_Dictionary = {
7     "Capt": "Officer",
8     "Col": "Officer",
9     "Major": "Officer",
10    "Jonkheer": "Royalty",
11    "Don": "Royalty",
12    "Sir": "Royalty",
13    "Dr": "Officer",
14    "Rev": "Officer",
15    "the Countess": "Royalty",
16    "Dona": "Royalty",
17    "Mme": "Mrs",
18    "Mlle": "Miss",
19    "Ms": "Mrs",
20    "Mr": "Mr",
```

```

21         "Mrs" :      "Mrs",
22         "Miss" :     "Miss",
23         "Master" :   "Master",
24         "Lady" :     "Royalty"
25     }
26
27 # we map each title
28 title[ 'Title' ] = title.Title.map( Title_Dictionary )
29 title = pd.get_dummies( title.Title )
30 #title = pd.concat( [ title , titles_dummies ] , axis = 1 )
31 title.head()

```

2.3.2 Informação na categoria das cabines

```

1 cabin = pd.DataFrame()
2
3 # replacing missing cabins with U (for Unknown)
4 cabin[ 'Cabin' ] = full.Cabin.fillna( 'U' )
5
6 # mapping each Cabin value with the cabin letter
7 cabin[ 'Cabin' ] = cabin[ 'Cabin' ].map( lambda c : c[0] )
8
9 # dummy encoding ...
10 cabin = pd.get_dummies( cabin[ 'Cabin' ] , prefix = 'Cabin' )
11
12 cabin.head()

```

2.3.3 O que é a classe do ticket?

```

1 # a function that extracts each prefix of the ticket, returns 'XXX' if no prefix (i.e
2 def cleanTicket( ticket ):
3     ticket = ticket.replace( '.' , '' )
4     ticket = ticket.replace( '/' , '' )
5     ticket = ticket.split()
6     ticket = map( lambda t : t.strip() , ticket )
7     ticket = list( filter( lambda t : not t.isdigit() , ticket ) )
8     if len( ticket ) > 0:
9         return ticket[0]
10    else:
11        return 'XXX'
12
13 ticket = pd.DataFrame()
14
15 # Extracting dummy variables from tickets:
16 ticket[ 'Ticket' ] = full[ 'Ticket' ].map( cleanTicket )
17 ticket = pd.get_dummies( ticket[ 'Ticket' ] , prefix = 'Ticket' )
18
19 ticket.shape
20 ticket.head()

```

2.3.4 Tamanho e tipo de família

```

1 family = pd.DataFrame()

```

```

2
3 # introducing a new feature : the size of families (including the passenger)
4 family[ 'FamilySize' ] = full[ 'Parch' ] + full[ 'SibSp' ] + 1
5
6 # introducing other features based on the family size
7 family[ 'Family_Single' ] = family[ 'FamilySize' ].map( lambda s : 1 if s == 1 else 0
8 family[ 'Family_Small' ] = family[ 'FamilySize' ].map( lambda s : 1 if 2 <= s <= 4 el
9 family[ 'Family_Large' ] = family[ 'FamilySize' ].map( lambda s : 1 if 5 <= s else 0
10
11 family.head()

```

2.4 Criando da base de dados do modelo final

Nesta etapa criam-se os conjuntos de treinamento e de teste para ter um conjunto de avaliação do modelo.

O conjunto de dados também é dividido por colunas em uma matriz (X) contendo os dados de entrada e um vetor (y) contendo o alvo (ou rótulos, sobreviveu ou não).

Selecione quais variáveis deseja incluir no conjunto de dados na lista abaixo:

- imputed;
- embarked;
- pclass;
- sex;
- family;
- cabin;
- ticket;
- title?

Inclua as variáveis que você gostaria de usar na função abaixo separadas por vírgula e, em seguida, execute o bloco de código abaixo.

```

1 # Select which features/variables to include in the dataset from the list below:
2 # imputed , embarked , pclass , sex , family , cabin , ticket
3
4 full_X = pd.concat( [ imputed , embarked , cabin , sex ] , axis=1 )
5 full_X.head()
6
7 # Create all datasets that are necessary to train , validate and test models
8 train_valid_X = full_X[ 0:891 ]
9 train_valid_y = titanic.Survived
10 test_X = full_X[ 891: ]
11 train_X , valid_X , train_y , valid_y = train_test_split( train_valid_X , train_valid_y
12 print (full_X.shape , train_X.shape , valid_X.shape , train_y.shape , valid_y.shape ,
13 plot_variable_importance(train_X , train_y)

```

2.5 Modelagem

A modelagem do aprendizado de máquina é basicamente a tipologia ou a técnica de classificação das variáveis de entrada para obter a variável de saída. Neste caso, se o passageiro sobreviveu ou não.

2.5.1 Seleção do Modelo, Treinamento e Avaliação

Utilize a o modelo que julgar mais adequado. Uma observação, no quadro de código a seguir retire o comentário do modelo que deseja entre as linhas 1 e 7. Caso utilize o modelo da linha 1 (*Random Florest*) retire o comentário da linha 15.

```
1 #model = RandomForestClassifier(n_estimators=100)
2 #model = SVC()
3 #model = GradientBoostingClassifier()
4 #model = KNeighborsClassifier(n_neighbors = 3)
5 #model = GaussianNB()
6 model = MLPClassifier(learning_rate = 'adaptive', hidden_layer_sizes = 200, max_iter =
7 #model = LogisticRegression()
8
9 model.fit( train_X , train_y )
10 print (model.score( train_X , train_y ) , model.score( valid_X , valid_y ))
11 #plot_model_var_imp(model, train_X , train_y)
12
13 #rfecv = RFECV(estimator = model, step = 1, cv = StratifiedKFold(2), scoring = 'accuracy')
14 rfecv = RFECV( estimator = model , step = 1 , cv = StratifiedKFold( train_y , 2 ) , scoring = 'accuracy')
15 #rfecv.fit( train_X , train_y )
16
17 test_Y = model.predict( test_X )
18 test_Y = np.asarray(test_Y, dtype=int)
19 passenger_id = full[891:].PassengerId
20 test = pd.DataFrame( { 'PassengerId': passenger_id , 'Survived': test_Y } )
21 test.shape
22 test.head()
23 test.to_csv( 'titanic_pred.csv' , index = False )
```

Esta é a etapa final da do processo. Vamos avaliar os resultados de acordo com a métrica do problema.

2.5.2 Métrica

A métrica para avaliar a pontuação é porcentagem de passageiros previstos corretamente. Este parâmetro é conhecido como precisão ou acurácia.

A precisão mede o quão bem a classificação binária identifica ou exclui corretamente uma condição. Ou seja, a precisão é a proporção de resultados verdadeiros (verdadeiros positivos e negativos) entre o número total de casos examinados.

2.5.3 Teste e submissão dos resultados

Para submeter o resultado no site do [kaggle](#) acesse o link e faça login.

Na tela do desafio do Titanic clique em Submit Predictions

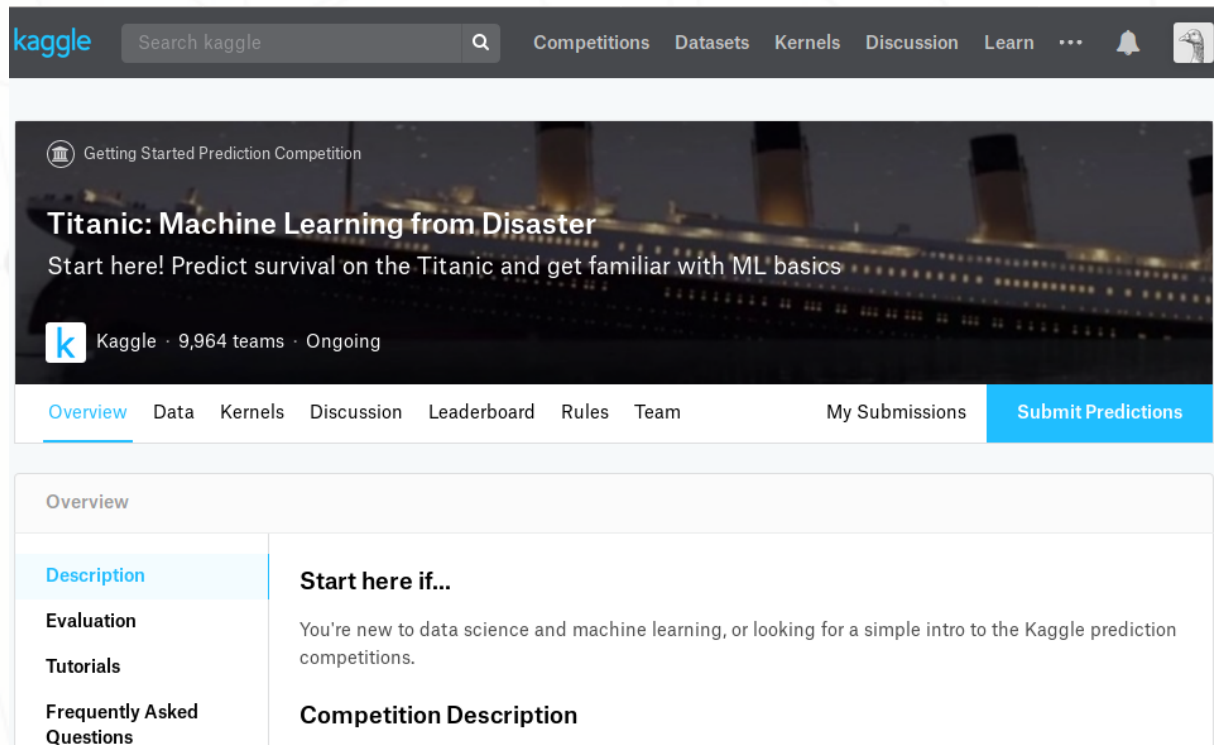


Figura 2.2: Tela de submissão dos resultados no Kaggle.

Considerações Finais

Conheça a BEMAKER nas redes sociais

www.facebook.com/bemakerbrasil

www.instagram.com/bemakerbr

www.linkedin.com/company/bemakerbr/

Bemaker - Eventos
Avenida Ivaí, 228, Dom Bosco
Belo Horizonte, MG

Yulab Tecnologia / Bemaker
Avenida Santa Matilde, 535, Dom Cabral
Belo Horizonte, MG

Este material está disponível no [GitLab](#).