

Nome: Michel Felipe de Sena Nassif

## Desafio Data Scientist - Qual a quantidade de bicicletas alugadas por hora em Montreal?

### Desafio:

*O desafio será desenvolver um modelo para prever a quantidade de bicicletas alugadas por hora ou por dia (fica a sua escolha).*

Neste desafio, utilizaremos os dados da **BIXI Montréal**. Também, os dados sobre as condições climáticas que foram obtidos através do Governo do Canadá.

#### I. Preparar e Explorar:

Os dados estão disponíveis publicamente em:

- <https://www.bixi.com/en/open-data>
- <http://www.climate.weather.gc.ca>

Mas também, neste [link](#), disponibilizamos os dados que foram coletados dessas bases públicas.

Nesta etapa, diferencie os tipos de variáveis e quais tratamentos fazer para elas. Para poder alcançar o objetivo do desafio, conecte os dados das duas fontes, e explore descritivamente a sua base criada.

#### II. Prever:

Nesta etapa, escolha algum modelo, implemente e avalie os resultados.

# Introdução

Estudando a base de dados de todas as bicicletas alugadas na cidade canadense de Montreal, no período de 15 de março a 15 de novembro de 2015 e 2016, assim como os dados climáticos da região, escolheu-se criar um modelo que juntasse as duas bases de informações e conseguisse prever o número de bicicletas alugadas na região por hora.

Através de uma primeira análise, foi possível concluir que se trata de um problema de modelagem supervisionada, no qual a partir de uma base de dados tratada deve-se ensinar o modelo a prever resultados de uma variável que dependa de outras totalmente independentes, presentes na nossa base de dados. Além disso, como, para o problema, a variável chave é o número de bicicletas alugadas por hora, a melhor estratégia para estimar esse valor é realizar um arquétipo de regressão que represente a dependência entre as variáveis.



Figura 1 – Exemplos de estratégias adotadas para criação de modelos.

Para a criação de um modelo preliminar, escolheu-se variáveis que levassem em conta aspectos temporais da rotina da cidade, como dia da semana ou hora e climáticos como a temperatura ou a umidade relativa do ar, de modo que as entradas escolhidas foram: ano, mês, data, hora, respectivo dia da semana - variando de 0 para segunda feira e 6 para domingo -, temperatura ( $^{\circ}\text{C}$ ), ponto de orvalho( $^{\circ}\text{C}$ ), umidade relativa do ar, velocidade do vento e visibilidade.

## Pré-processamento dos dados

Feita essa estratégia de resolução do problema, percebeu-se que se tratam de dados que precisaram ser manipulados. Os arquivos estavam separados por meses e para uma melhor manipulação escolheu-se agrupar esses arquivos em dois– um climático e um com todas as bicicletas alugadas. Para isso foi utilizado o seguinte código do arquivo *merging.py*, obtidos em um site<sup>1</sup>, capaz de agrupar vários arquivos .csv.

```
8 import pandas as pd
9 import os
10
11 files = [f for f in os.listdir('.') if os.path.isfile(f)]
12
13 merged = []
14
15 for f in files:
16     filename, ext = os.path.splitext(f)
17     if ext == '.csv':
18         read = pd.read_csv(f)
19         merged.append(read)
20
21 result = pd.concat(merged)
22
23 result.to_csv('merged.csv')
24
```

Figura 2 – Imagem do Código capaz de agrupar vários arquivos .csv.

No entanto, algumas variáveis precisavam ser transformadas e outras criadas, como por exemplo o dia da semana e o número de bicicletas alugadas por hora. Em um arquivo chamado *Preprocessamento.py* realizou-se essa etapa de transformação dos dados, assim como a integração entre as duas bases dados, para no final ter-se uma base de dados pronta para ser utilizada em um modelo.

Nessa etapa, como primeiro passo, importou-se as bibliotecas e as duas bases de dados: *datasetbike* referente aos aluguéis de bicicletas e o *datasetclimate* referente aos dados climáticos; criando-se em seguida a base de dados finais, representada por *datasetfinal*. Fonte de informação que apresenta os dados da lista *legenda*, como mostrado na figura 3.

Após importada as bases de dados, transferiu-se diretamente as informações de ano, mês, dia, hora, temperatura, ponto de orvalho, umidade relativa, velocidade do vento e visibilidade do *datasetclimate* para o *datasetfinal*. Em seguida, através da biblioteca *datetime*, utilizou-se a função *weekday()*, com os parâmetros de ano, mês e dia, para se achar o dia da semana correspondente e armazenando assim no banco de dados, também mostrada na figura 3.

```

8 # Importando as bibliotecas
9 import numpy as np
10 import pandas as pd
11 from datetime import date
12 import collections
13
14
15 # Importando os dados de aluguel de bicicleta e os climaticos
16 datasetbike = pd.read_csv('Bixi_time.csv', encoding = 'ISO-8859-1')
17 datasetclimate = pd.read_csv('Climate2.csv', encoding = 'ISO-8859-1')
18
19 #Criando o dataset final
20 legenda=['Ano', 'Mês', 'Dia', 'Hora', 'Dia da Semana', 'Temperatura(°C)',
21         'Ponto de Orvalho(°C)', 'Umidade Relativa', 'Velocidade do Vento(Km/h)',
22         'Visibilidade', 'Nº de Aluguéis por hora']
23 datasetfinal=pd.DataFrame(index= datasetclimate.index, columns=legenda)
24
25 #Importando as variaveis de tempo e clima para a base de dados
26 datasetfinal.iloc[:,0]=datasetclimate.iloc[:,2] #Ano
27 datasetfinal.iloc[:,1]=datasetclimate.iloc[:,3] #Mês
28 datasetfinal.iloc[:,2]=datasetclimate.iloc[:,4] #Dia
29 datasetfinal.iloc[:,5]=datasetclimate.iloc[:,7] #Temperatura
30 datasetfinal.iloc[:,6]=datasetclimate.iloc[:,9] #Ponto de Orvalho
31 datasetfinal.iloc[:,7]=datasetclimate.iloc[:,11] #Umidade Relativa
32 datasetfinal.iloc[:,8]=datasetclimate.iloc[:,15] #Velocidade do Vento
33 datasetfinal.iloc[:,9]=datasetclimate.iloc[:,17] #Visibilidade
24

```

Figura 3 – Passo 1 do pré-processamento de dados.

Por último, era necessário criar a variável de saída do número de bicicletas alugadas. Para isso criou-se uma lista auxiliar chamada *databike*, na qual todos os dados temporais de cada aluguel de bicicleta – do ano até a hora do aluguel - foram armazenados, para depois criar uma variável *counter* que junta todos os aluguéis que ocorreram na mesma hora. Vale notar, que para o problema foi escolhido levar em consideração a hora do aluguel da bicicleta e não sua devolução, uma vez que apenas 0,7% dos aluguéis usavam a bicicleta por mais de uma hora.

Tendo, o número de aluguéis que aconteceram na mesma hora, bastou-se apenas percorrer a base de dados do *datasetfinal*, já separada por dados de hora em hora, de modo que cada número de aluguéis por hora foi armazenado em sua determinada hora correspondente, criando assim a variável de saída no *datasetfinal*, convertendo-o depois para um arquivo .csv. Processo que pode ser inteiro visto na figura 4.

```

43 #Somando o nº de bicicletas alugadas em cada hora
44 databike = [" for x in range(len(datasetbike))"]
45 for i in range(len(datasetbike)):
46     databike[i]=datasetbike.iloc[i,1][0:13]
47 counter=collections.Counter(databike)
48
49
50 # Preenchendo no dataset a variável do nº de bicicletas alugadas por hora
51 valores = list(counter.values())
52 datas = list(counter.keys())
53 for i in range(len(datasetfinal)):
54     data=[datasetfinal.iloc[i,0], datasetfinal.iloc[i,1],datasetfinal.iloc[i,2],
55     datasetfinal.iloc[i,3]]
56     for k in range(len(datas)):
57         datacomp=[int(datas[k][0:4]),int(datas[k][5:7]),int(datas[k][8:10]),
58                 int(datas[k][11:13])]
59         ano=datacomp[0]
60         mes=datacomp[1]
61         dia=datacomp[2]
62         hora=datacomp[3]
63         if(ano==data[0] and mes ==int(data[1]) and dia ==int(data[2]) and hora==int(data[3])):
64             datasetfinal.iloc[i,10]=valores[k]
65
66 #Exportando para CSV:
67 datasetfinal.to_csv('DatasetFinal.csv', encoding='utf-8')

```

Figura 4 – Achando a variável de número de alugueis de bicicleta por hora.

## Criação e Aperfeiçoamento do Modelo

Com a base de dados disponível e todas as variáveis já manipuladas e criadas, passou-se para a etapa de criação do modelo, no arquivo chamado *CriaçãodoModelo.py*. Vale notar que foi necessário remover as linhas referentes aos períodos do começo de abril de 2015 e 2016 e final de novembro de 2015 e 2016, uma vez que na base de dados disponibilizada de alugueis de bicicleta não havia dados referentes a esse período.

Nesse arquivo, após importação do banco de dados, realizou-se um passo que teve como intuito localizar casos de dados faltando, adotando-se a estratégia de substituí-los pela média, no caso da matriz de variáveis independentes X e por 0 para a variável independente Y (uma vez que um NaN no número de alugueis de bicicletas afirma que não foi encontrada nenhuma locação de bicicleta no processo de contagem do pré-processamento de dados para aquela hora).

Realizado essa etapa, dividiu-se aleatoriamente a base de dados em uma parte referente ao treino o modelo (80% das informações) e outra para o teste (20% das informações), denominados respectivamente *training\_set* e *test\_set*. Processos que podem ser todos vistos na figura 5.

```

8 import numpy as np
9 import matplotlib.pyplot as plt
10 import pandas as pd
11 import statsmodels.formula.api as sm
12 from datetime import datetime
13
14 # Importando o Dataset
15 dataset = pd.read_csv('DatasetFinal.csv')
16 dataset['Horario'] = (dataset.apply(lambda row: datetime(
17     int(row['Ano']), int(row['Mês']), int(row['Dia']),
18     int(row['Hora'])), axis=1))
19
20 #Lidando com a possibilidade de nenhuma bicicleta ter sido alugada em um horário
21 dataset.iloc[:,11]=dataset.iloc[:,11].fillna(0)
22 np.argwhere(np.isnan(dataset.iloc[:,11]))
23
24 #Criando os vetores de entrada e saída
25 X = dataset.iloc[:, 1:11].values
26 y = dataset.iloc[:, 11].values
27
28 # Trabalhando com a possibilidade de falta de dados em X
29 from sklearn.preprocessing import Imputer
30 imputer = Imputer(missing_values = 'NaN', strategy = 'mean', axis = 0)
31 imputer = imputer.fit(X[:, 1:11])
32 X[:, 1:11] = imputer.transform(X[:, 1:11])
33
34
35 #Dividindo os dados em Training set e Test set
36 from sklearn.cross_validation import train_test_split
37 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)
38

```

Figura 5 – Etapa anterior a criação do modelo.

Realizado isso, criou-se um modelo inicial de regressão (*regressor*) com parâmetros genéricos para o nosso modelo utilizando o algoritmo de “Random Forest Regression”<sup>2</sup>, modelo robusto que já foi capaz de lidar muito bem com a nossa base de dados, como pode ser visto na tabela 2. Vale notar que foram testados outros processos de regressão como o SVR e o Polinomial, porém acabaram mostrando resultados inferiores pelo algoritmo escolhido para o problema.

Com o intuito de refinar o modelo, determinou-se as principais variáveis que afetam o número de alugueis de bicicletas realizados por hora. Para isso foi aplicado manualmente os passos do algoritmo de “Backward Elimination”<sup>3</sup>, no qual remove-se iterativamente todas as variáveis que apresentam um valor P maior que um nível de significância escolhido – no caso deste problema: 0,05. Com isso verificou-se que as variáveis do Ponto de Orvalho, dia do mês, ano não tinham muita correlação nosso vetor de saída, assim como a velocidade do vento que afeta o modelo, porém não o suficiente de acordo com os critérios estabelecidos (valor P de  $0,07 > 0,05$ ).

Resultando em um modelo reduzido as variáveis do mês, hora, dia da semana, temperatura, umidade relativa e visibilidade, que foram armazenadas em uma nova matriz denominada  $X_{opt}$ , utilizada novamente para a criação de um novo modelo regressão. Processos que são todos representados na figura 6.

```

41 # Ajustando o modelo de Random Forest Regression
42 from sklearn.ensemble import RandomForestRegressor
43 regressor = RandomForestRegressor()
44 regressor.fit(X_train, y_train)
45
46
47 #Olhando os valores de Adjusted R2 e P-Value
48 #Busca pelas variáveis de entrada mais importantes(P-Value<0.05)
49 X_opt= X_train[:,[0,1,2,3,4,5,6,7,8,9]]
50 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
51 regressor_OLS.summary()
52
53 #Remoção do Ponto de Orvalho
54 X_opt= X_train[:,[0,1,2,3,4,5,7,8,9]]
55 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
56 regressor_OLS.summary()
57
58 #Remoção do dia
59 X_opt= X_train[:,[0,1,3,4,5,7,8,9]]
60 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
61 regressor_OLS.summary()
62
63 #Remoção do ano
64 X_opt= X_train[:,[1,3,4,5,7,8,9]]
65 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
66 regressor_OLS.summary()
67 |
68 #Remoção da velocidade do vento
69 X_opt= X_train[:,[1,3,4,5,7,9]]
70 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
71 regressor_OLS.summary()
72 X_test=X_test[:,[1,3,4,5,7,9]]
73
74 # Ajustando o novo modelo
75 regressor.fit(X_opt, y_train)
76

```

Figura 6 – Criação do modelo e procura pelas variáveis mais importantes.

Por último, na figura 7, observa-se o código que procurou otimizar um pouco mais o modelo, olhando-se para seus hiperparâmetros, através da técnica de “Grid Search”<sup>4</sup>. Foram realizados três testes, um com parâmetros padrões e dois que percorreram uma rede de valores, representados pela lista *parameters*, resultados que podem ser vistos na tabela 2. Para então recriar o modelo de regressão novamente, dessa vez com os resultados dos melhores hiperparâmetros obtidos nos testes.

```

68 #Remoção da velocidade do vento
69 X_opt= X_train[:,[1,3,4,5,7,9]]
70 regressor_OLS=sm.OLS(endog=y_train, exog = X_opt).fit()
71 regressor_OLS.summary()
72 X_test=X_test[:,[1,3,4,5,7,9]]
73
74 # Ajustando o novo modelo
75 regressor.fit(X_opt, y_train)
76
77 #Otimizando os Hiperparâmetros do modelo
78 from sklearn.model_selection import GridSearchCV
79 parameters = [{}]#Teste 1 - Valores Padrões - Resultado:0.9158
80 parameters = [{ #Teste 2
81     'bootstrap': [True],
82     'max_depth': [None, 10, 50],
83     'max_features': [2, 3],
84     'min_samples_leaf': [1,3],
85     'min_samples_split': [2,3],
86     'n_estimators': [10,50, 100] #Resultado:0,9206
87 }]
88 parameters = [{ #Teste 3
89     'bootstrap': [True],
90     'max_depth': [40,50,60], #40
91     'max_features': [2,3,4],#4
92     'min_samples_leaf': [1],#1
93     'min_samples_split': [3,4],#3
94     'n_estimators': [100,200,500] #500
95 }]#Resultado:0,9266
96 grid_search = GridSearchCV(estimator = regressor, param_grid = parameters,
97                             cv =3, n_jobs=-1, )
98 grid_search = grid_search.fit(X_opt,y_train)
99 melhorresultado = grid_search.best_score_
100 melhoresparametros = grid_search.best_params_
101
102 regressor = RandomForestRegressor(n_estimators=500,bootstrap=True,max_depth=40,
103                                 max_features=4,min_samples_leaf=3)
104 regressor.fit(X_opt, y_train)

```

Figura 7 – Aplicação do Algoritmo de Grid Search

Para uma análise inicial do modelo, utilizou-se a função *.predict* do nosso regressor para o base de dados de teste *X\_test*. Foi realizada uma comparação direta de alguns valores entre os resultados previstos e o que de fato aconteceu com o vetor *y\_test*, dados disponíveis na tabela 1. Análise que apesar de simples, já consegue nos dar uma indicação da qualidade do modelo.



Tabela 1 – Comparação entre os resultados do modelo e a realidade no número de bicicletas alugadas por hora

Linha	Valor Real (y_test)	Previsão (y_pred)
0	74	62
500	1421	1163
912	447	652,5
1188	437	418,3
1573	88	188,7
2004	1282	1305
2063	1539	1579,1

A tabela 2 consegue fornecer uma análise mais detalhada dos modelos, nela foram utilizados os parâmetros de  $R^2$  Ajustado, que é capaz de nos mostrar a relação entre a qualidade do modelo e o número de variáveis utilizadas, penalizando valores mais altos de número de variáveis escolhidas; além do resultado obtido pelo parâmetro *.best\_score* da nossa função do Grid Search.

Tabela 2 – Índices de Performance para os modelos

Tipo de Modelo	Valor $R^2$ Ajustado	Nota Grid Search
Modelo Inicial com parâmetros padrões	0,799	0,9171
Modelo reduzido com parâmetros padrões	0,799	0,9158
Modelo reduzido e otimizado com Grid Search	0,799	0,9266

Por último, através do código da figura 8 foi realizada uma comparação visual através de dois gráficos que comparam os valores estimados com a realidade na variável de interesse do problema, representados na figura 9 e 10. Foram pegos pontos em intervalos de 100 e 100 linhas para a base de dados inteira na figura 9 e intervalos de 5 em 5 para o gráfico da figura 10, que pega o intervalo inteiro do mês de Junho de 2016. Em azul estão os valores previstos pelo modelo e em vermelho os valores reais.

```

110 # Gráfico - Análise Geral
111 eixoX=np.arange(0,10320,100)
112 eixoY=X[range(0,10320,100),:]
113 eixoY=regressor.predict(eixoY[:,[1,3,4,5,7,9]])
114 plt.scatter(eixoX, y[range(0,10320,100)], color = 'red')
115 plt.plot(eixoX, eixoY, color = 'blue')
116 plt.title('Comparação entre os resultados reais e o modelo - Dataset Inteiro')
117 plt.xlabel('Número da linha')
118 plt.ylabel('Nº de Bicicletas Alugadas')
119 plt.show()
120
121 # Gráfico2 - Análise do mês 6 de 2016
122 eixoX=np.arange(6288,7007,5)
123 eixoY=X[range(6288,7007,5),:]
124 eixoY=regressor.predict(eixoY[:,[1,3,4,5,7,9]])
125 plt.scatter(eixoX, y[range(6288,7007,5)], color = 'red')
126 plt.plot(eixoX, eixoY, color = 'blue')
127 plt.title('Comparação entre os resultados reais e o modelo - Junho de 2016')
128 plt.xlabel('Número da linha')
129 plt.ylabel('Nº de Bicicletas Alugadas')
130 plt.show()
131

```

Figura 8 – Código da montagem dos gráficos e previsão dos resultados

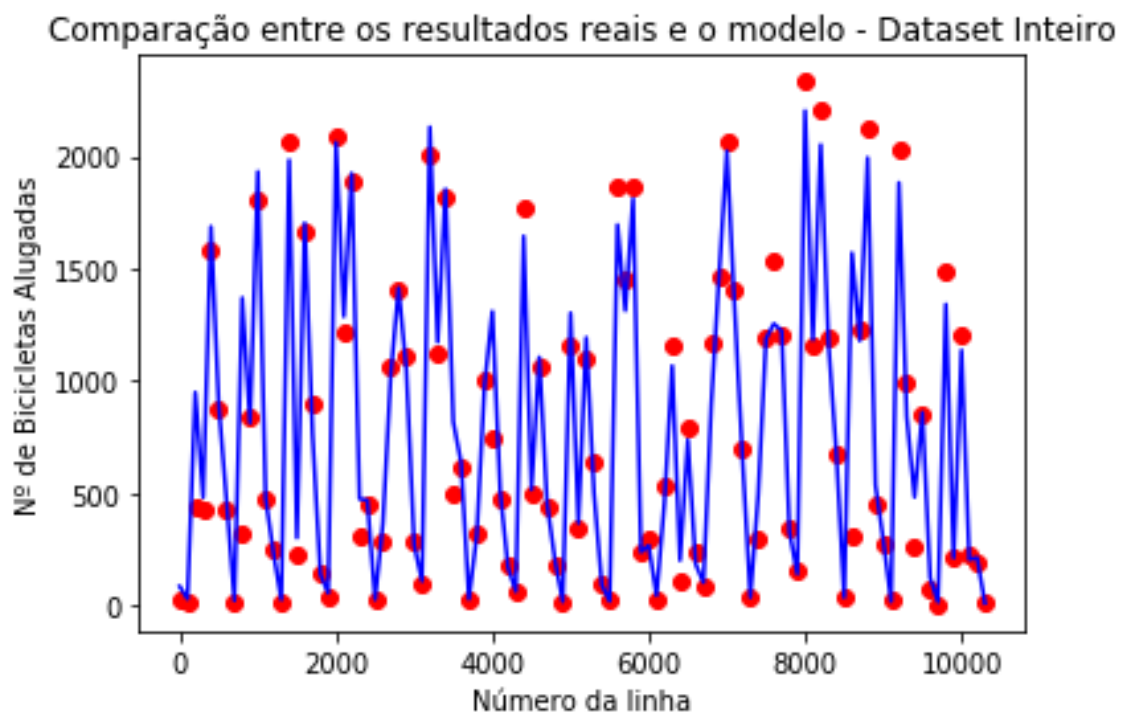


Figura 9 – Comparação entre a realidade em vermelho e o previsto em azul.

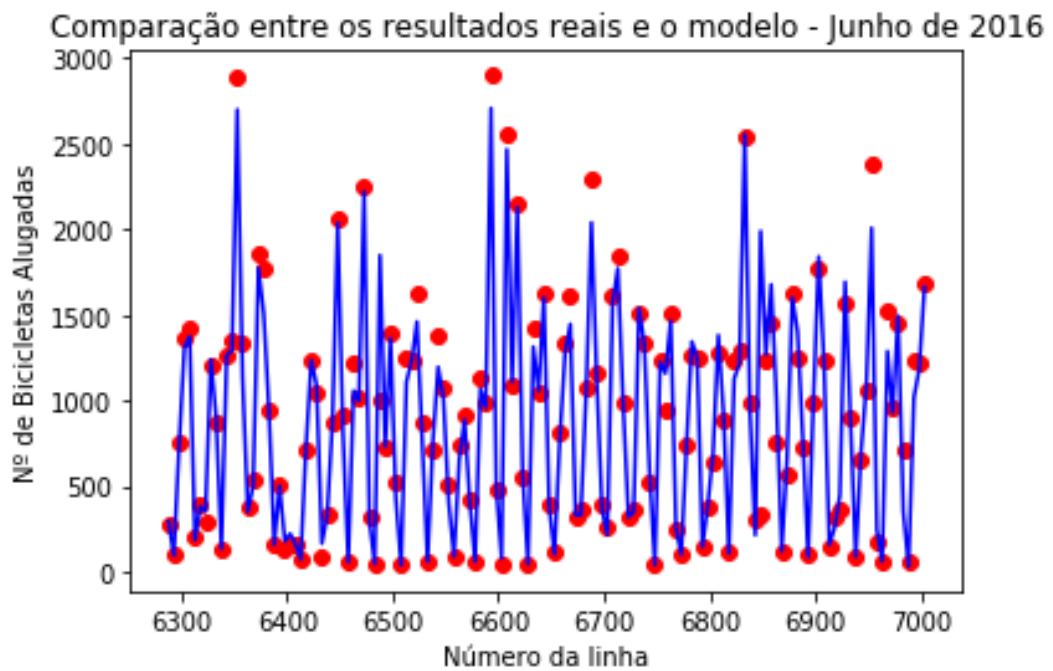


Figura 10 – Comparação para o mês de Junho de 2016 entre a realidade em vermelho e o previsto em azul.

## Conclusão e Análise dos Resultados

A partir dos resultados obtidos através das diferentes análises, é possível concluir que o modelo ainda consegue ser melhorado, através de uma melhor escolha ou manipulação das variáveis, uma vez que foram obtidos valores de  $R^2$  ajustado da ordem de 0,8.

Ou por uma pequena melhor predição dos valores da realidade, através de uma varredura mais profunda dos hiperparâmetros, já que foram obtidas notas no Grid Search da ordem de 90%. Contudo melhorias significativas nesse valor irá fazer com que o modelo sofra com “overfitting”, algo que ocorre com uma certa frequência para o algoritmo de “Random Forest” utilizado.

Porém, de maneira geral, o modelo de regressão utilizado foi capaz de se adequar bem aos valores da realidade, fornecendo respostas satisfatórias aos valores do número de alugueis realizados por hora na cidade de Montreal, sendo capaz também de pegar o comportamento cíclico em relação à madrugada e aos horários de pico do dia a dia no aluguel de bicicletas de uma cidade, como é possível de ser visualizado nos gráficos da figura 9 e 10.

# Referências

1. Leon Lovett. Merge multiple csv files.  
<https://www.linkedin.com/pulse/merge-multiple-csv-files-python-leon-lovett/>
2. Towardsdatascience.  
<https://towardsdatascience.com/hyperparameter-tuning-the-random-forest-in-python-using-scikit-learn-28d2aa77dd74>
3. Stat. Backward Elimination  
<https://www.stat.ubc.ca/~rollin/teach/643w04/lec/node42.html>
4. Scikit  
[http://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)