


```
[36]: # Balançamento
target_valores = dados.valores[:, 18]
target_valores = target_valores.astype('int64')
balanecedor = SMOTETK()
escala_pronto, target_balancedo = balanecedor.fit_resample(escala_pronto, target_valores)
```

Analisando multicolinearidade de atributos

```
In [37]: # Construindo matriz de correlação - Tabela
corr = pd.DataFrame(escala_pronto, columns = dados.columns[0:18]).corr(method = 'pearson') corr

Out[37]:
```

	state	account_length	international_plan	voice_mail_plan	number_vmail_messages	total_day_minutes
state	1.000000	-0.019438	0.054757	-0.039700	-0.037747	0.079211
account_length	-0.019438	1.000000	0.023696	0.015760	0.006080	-0.003967
international_plan	0.054757	0.023696	1.000000	0.044761	0.036932	0.017680
voice_mail_plan	-0.039700	0.015760	0.044761	1.000000	0.963915	-0.106416
number_vmail_messages	-0.037747	0.003080	0.036932	0.963915	1.000000	-0.102342
total_day_minutes	0.079211	-0.003967	0.017680	-0.106416	-0.102342	1.000000
total_day_calls	-0.026135	0.033740	-0.040010	-0.024705	-0.022082	0.043593
total_day_charge	0.079209	-0.003965	0.017667	-0.106409	-0.102336	1.000000
total_eve_charge	0.020349	-0.015566	0.017661	-0.036517	-0.041751	0.152181
total_eve_calls	0.010816	0.024822	-0.032027	-0.002669	-0.005646	0.015130
total_eve_charge	0.020348	-0.015564	0.017664	-0.036502	-0.041732	0.152157
total_night_minutes	0.002450	-0.044033	-0.072469	-0.050754	-0.044793	0.001479
total_night_calls	-0.005604	-0.030704	0.028123	0.014580	0.003747	0.037132
total_night_charge	0.002430	-0.044239	-0.072464	-0.050779	-0.044846	0.001496
total_intl_minutes	0.059447	-0.000243	0.169371	0.018973	0.019840	-0.002287
total_intl_calls	-0.054175	0.024250	-0.059637	0.044337	0.046478	0.024478
total_intl_charge	0.059592	-0.000255	0.169124	0.018959	0.019819	-0.002034
number_customer_service_calls	0.054514	-0.003767	-0.088424	-0.016801	-0.010680	-0.018222

Para facilitar a visualização, vou gerar um mapa de calor a partir da tabela acima de correlação. Dessa forma, poderemos observar de forma mais clara possíveis atributos colineares.

```
In [38]: # Construindo matriz de correlação - Gráfico

# Definindo área de plotagem
fig, ax = plt.subplots(figsize = (20, 15))

# construção do gráfico
sns.heatmap(corr, linewidths = 1, annot = True, cmap = 'mako')
ax.set_title('Matriz de Correlação', fontsize = 30)

# Mostrar gráfico de correlação
plt.show()
```



Se observarmos além das variáveis correspondentes iguais (correlação = 1 para as variáveis analisadas com elas mesmas), observamos alguns outros coeficientes de correlação muito fortes, que podem afetar o algoritmo.

A fim de identificar e retirar os de nosso modelo, vamos utilizar a técnica VIF (Variance Inflation Factor). Neste método, cada variável é selecionada e sofre um processo de regressão em relação às demais variáveis. Os fatores são resultados destas regressões e, se eles forem altos, significam forte correlação. Fatores acima de 5 indicam alta multicolinearidade e as variáveis podem ser prejudiciais ao modelo.

Abaixo vamos calcular o VIF de cada variável do dataset e verificar quais superam 5, retirando-as do dataset para realização do treinamento do modelo.

```
In [ ]: # Salvando o plot
fig.savefig("plots/correlation-matrix.png")

In [39]: # Construção de um dataframe com os dados já normalizados e padronizados, unindo com os nomes das colun
as
vif_df = pd.DataFrame(escala_pronto, columns = dados.columns[0:18])

# Acrescentando constante na matriz de variáveis
X = add_constant(vif_df)

# Imprimindo VIFs:
pd.Series([variance_inflation_factor(X.values, i) for i in range(X.shape[1])], index=X.columns)
```

Acima observamos que os VIFs das variáveis 'total_night_minutes', 'total_night_charge', 'total_intl_minutes' e 'total_intl_charge' superaram 5, então, serão retiradas do dataset.

```
In [40]: # Retirando variáveis colineares
escala_pronto = pd.DataFrame(escala_pronto).drop(columns = [11, 13, 14, 16]).values

Machine Learning

In [41]: # Divisão do dataset em dados de treino e dados de teste
x_treino, x_teste, y_treino, y_teste = train_test_split(escala_pronto, target_balancedo, \
                                                    test_size = 0.3)
```

Modelo 1 - Logistic Regression com Train-Test-Split comum e sem ajuste de hiper-parâmetros

```
In [42]: # Instanciando Modelo 1 - Logistic Regression
modelo1 = LogisticRegression()

In [43]: # Treinamento do Modelo 1
modelo1.fit(x_treino, y_treino)

Out[43]: LogisticRegression()
```

```
In [44]: # Fazendo as previsões com o Modelo 1
previsoes1 = modelo1.predict(x_teste)
```

Avaliando o Modelo 1

```
In [45]: # Relatório do Modelo 1
print("Relatório de Classificação:\n", classification_report(y_teste, previsoes1, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes1) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes1)))
```

	precision	recall	f1-score	support
0	0.7988	0.7538	0.7756	853
1	0.7680	0.8110	0.7889	857
accuracy				0.7825
macro avg	0.7834	0.7824	0.7823	1710
weighted avg	0.7833	0.7825	0.7823	1710

Acurácia: 78.25%
AUC: 0.78

Explicações e Conclusões - Modelo 1

A acurácia ficou acima de 70%, o que é aceitável para o modelo e problema de negócio.

O recall também mostrou boas porcentagens, mostrando que o balançamento da variável target funcionou e o modelo aprendeu sem vies para uma das opções de saída.

O AUC também ficou satisfatório, tendo em vista que indica a precisão das previsões e quanto maior, melhor.

No entanto, é válido tentarmos realizar algumas melhorias no algoritmo, como o método Cross Validation.

Modelo 1 - Logistic Regression com Cross Validation e sem ajuste de hiper-parâmetros

```
In [46]: # Cross Validation com Modelo 1
resultado = cross_val_score(modelo1, escala_pronto, target_balancedo, cv = KFold(5, True))

In [47]: # Observando acurácia de cada fold
for accuracy in range(5):
    print("Acurácia fold %d: %.2f%%" % ((accuracy.numerator + 1), resultado.take(accuracy) * 100))
```

Acurácia fold 1: 77.99%
Acurácia fold 2: 76.23%
Acurácia fold 3: 77.37%
Acurácia fold 4: 79.47%
Acurácia fold 5: 78.23%

```
In [48]: # Imprimindo média de acurácia - modelo 2
print("Acurácia Final: %.2f%%" % (resultado.mean() * 100))

Acurácia Final: 78.16%
```

Observamos que o cross validation não funcionou da forma esperada, tendo em vista que a acurácia teve uma leve queda.

Salvando Modelo 1

```
In [49]: arquivo1 = 'modelos/modelo_classificador_lr.sav'
pickle.dump(modelo1, open(arquivo1, 'wb'))
```

Modelo 2 - Ajustando Hiper-Parâmetros

O Grid Search Parameter tuning faz combinações dos parâmetros do algoritmo (neste caso, Regressão Logística). Os parâmetros são fornecidos a função que cria uma tabela com os melhores valores para os parâmetros testados.

```
In [50]: # Definindo os valores e parâmetros que serão testados
valores_grid = ['penalty' in ['l1', 'l2', 'l1', 'l2'], 'C' in [0.001, 0.01, 0.1, 1, 10, 100, 1000], \
               'class_weight' in ['dict', 'balanced', 'None'], 'solver' in ['lbfgs', 'liblinear']]

In [51]: # Criando o grid e treinando o modelo com todos os parâmetros apresentados
grid = GridSearchCV(estimator = modelo1, param_grid = valores_grid,
                    scoring = 'roc_auc', cv = 5)
```

```
Out[51]: GridSearchCV(estimator=LogisticRegression(),
                    param_grid={'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000],
                    'class_weight': ['dict', 'balanced', 'None'],
                    'penalty': ['l1', 'l2'],
                    'solver': ['lbfgs', 'liblinear']})
```

```
In [52]: # Analisando resultados
print("Acurácia: %.3f" % (grid.best_score_ * 100))
print("Melhores Parâmetros do Modelo:\n", grid.best_params_)
```

Acurácia: 78.33%
Melhores Parâmetros do Modelo:
'C': 0.01, 'class_weight': 'balanced', 'penalty': 'l1', 'solver': 'liblinear'

Explicações e Conclusões - Modelo 2

Entre valores testados para os parâmetros escolhidos (regra de penalização, força de regularização, pesos das classes e algoritmo interno de otimização), o método mostrou melhores resultados para C = 0.1 e penalidade = l1, peso = balanced e algoritmo de otimização = liblinear (bom resultado com dataset pequeno). Os valores default para esses parâmetros são, respectivamente, 1, l2, None e lbfgs.

Mesmo com aparente pequena melhora na acurácia, em um grande conjunto de dados, esses pequenos percentuais podem representar muitos valores previstos de forma correta.

Nesse sentido, vou salvar um novo modelo com os valores do resultado acima nos respectivos parâmetros.

Construindo e Treinando Modelo 2 (novo modelo com ajuste de hiper-parâmetros)

```
In [53]: # Instanciando modelo com ajuste de hiper-parâmetros
# (utilizando random_state para "fidelizar" os dados)
modelo2 = LogisticRegression(C=0.1, penalty = 'l1', solver = 'liblinear', class_weight = 'balanced',
                             random_state = 10)
```

```
In [54]: # Treinamento do Modelo 2
modelo2.fit(x_treino, y_treino)
```

```
Out[54]: LogisticRegression(C=0.1, class_weight='balanced', penalty='l1',
                             random_state=10, solver='liblinear')
```

Salvando Modelo 2

```
In [55]: arquivo2 = 'modelos/modelo_classificador_lr_gridsearchcv.sav'
pickle.dump(modelo2, open(arquivo2, 'wb'))
```

Conforme observado anteriormente, apesar de resultados adequados, pode ser que a regressão logística não seja o melhor algoritmo para este caso de classificação binária.

Neste sentido, acredito ser válido testar a performance de outros algoritmos como LDA, KNN, Naive Bayes, CART e SVM.

Modelo 3 - Linear Discriminant Analysis

```
In [56]: # Instanciando Modelo 3 - Linear Discriminant Analysis
modelo3 = LinearDiscriminantAnalysis()
```

```
In [57]: # Treinamento do Modelo 3
modelo3.fit(x_treino, y_treino)
```

```
Out[57]: LinearDiscriminantAnalysis()
```

```
In [58]: # Fazendo previsões com o Modelo 3
previsoes3 = modelo3.predict(x_teste)
```

Avaliando Modelo 3

```
In [59]: # Relatório do Modelo 3
print("Relatório de Classificação:\n", classification_report(y_teste, previsoes3, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes3) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes3)))
```

	precision	recall	f1-score	support
0	0.8020	0.7597	0.7803	853
1	0.7727	0.8133	0.7925	857
accuracy				0.7865
macro avg	0.7874	0.7865	0.7864	1710
weighted avg	0.7873	0.7865	0.7864	1710

Acurácia: 78.65%
AUC: 0.79

Explicações e Conclusões - Modelo 3

Aparentemente não houve melhora considerável, tendo em vista que acurácia e AUC permaneceram praticamente iguais.

Salvando Modelo 3

```
In [60]: arquivo3 = 'modelos/modelo_classificador_lda.sav'
pickle.dump(modelo3, open(arquivo3, 'wb'))
```

Modelo 4 - KNN - K-Nearest Neighbors

```
In [61]: # Instanciando o Modelo 4 - KNN
modelo4 = KNeighborsClassifier()
```

```
In [62]: # Treinamento do Modelo 4
modelo4.fit(x_treino, y_treino)
```

```
Out[62]: KNeighborsClassifier()
```

```
In [63]: # Fazendo previsões com o Modelo 4
previsoes4 = modelo4.predict(x_teste)
```

Avaliando Modelo 4

```
In [64]: # Relatório do Modelo 4
print("Relatório de Classificação:\n", classification_report(y_teste, previsoes4, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes4) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes4)))
```

	precision	recall	f1-score	support
0	0.9843	0.8089	0.8880	853
1	0.8385	0.9872	0.9068	857
accuracy				0.8982
macro avg	0.9114	0.8980	0.8974	1710
weighted avg	0.9112	0.8982	0.8974	1710

Acurácia: 89.82%
AUC: 0.90

Explicações e Conclusões - Modelo 4

Tanto acurácia como AUC melhoraram de forma expressiva.

Apesar disso, parece que o modelo aprendeu de forma levemente desbalanceada sobre as classes 0 e 1 (recall 0 = 80%, recall 1 = 99%).

Salvando Modelo 4

```
In [65]: arquivo4 = 'modelos/modelo_classificador_knn.sav'
pickle.dump(modelo4, open(arquivo4, 'wb'))
```

Modelo 5 - Naive Bayes

```
In [66]: # Instanciando Modelo 5 - Naive Bayes
modelo5 = GaussianNB()
```

```
In [67]: # Treinamento do Modelo 5
modelo5.fit(x_treino, y_treino)
```

```
Out[67]: GaussianNB()
```

```
In [68]: # Fazendo previsões com o Modelo 5
previsoes5 = modelo5.predict(x_teste)
```

Avaliando Modelo 5

```
In [69]: # Relatório de Classificação:\n", classification_report(y_teste, previsoes5, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes5) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes5)))
```

	precision	recall	f1-score	support
0	0.8375	0.7796	0.8075	853
1	0.7948	0.8495	0.8212	857
accuracy				0.8146
macro avg	0.8161	0.8145	0.8144	1710
weighted avg	0.8161	0.8146	0.8144	1710

Acurácia: 81.46%
AUC: 0.81

Explicações e Conclusões - Modelo 5

Apesar de aprender de forma um pouco mais equilibrada sobre as duas classes (observar recall), neste algoritmo observamos a queda da acurácia e da AUC.

Salvando Modelo 5

```
In [70]: arquivo5 = 'modelos/modelo_classificador_nb.sav'
pickle.dump(modelo5, open(arquivo5, 'wb'))
```

Modelo 6 - CART - Classification and Regression Trees

```
In [71]: # Instanciando Modelo 6 - CART
modelo6 = DecisionTreeClassifier()
```

```
In [72]: # Treinamento do Modelo 6
modelo6.fit(x_treino, y_treino)
```

```
Out[72]: DecisionTreeClassifier()
```

```
In [73]: # Fazendo previsões com o Modelo 6
previsoes6 = modelo6.predict(x_teste)
```

Avaliando Modelo 6

```
In [74]: print("Relatório de Classificação:\n", classification_report(y_teste, previsoes6, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes6) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes6)))
```

	precision	recall	f1-score	support
0	0.9034	0.8664	0.8925	853
1	0.8743	0.9253	0.8991	857
accuracy				0.8959
macro avg	0.8973	0.8958	0.8958	1710
weighted avg	0.8973	0.8959	0.8958	1710

Acurácia: 89.59%
AUC: 0.90

Explicações e Conclusões - Modelo 6

Utilizando o CART, observamos aumento considerável na acurácia e na AUC, no entanto, ainda encontramos certo desequilíbrio no treinamento das duas classes.

Salvando Modelo 6

```
In [75]: arquivo6 = 'modelos/modelo_classificador_cart.sav'
pickle.dump(modelo6, open(arquivo6, 'wb'))
```

Modelo 7 - SVM - Support Vector Machines

```
In [76]: # Instanciando Modelo 7 - SVM
modelo7 = SVC(probability = True)
```

```
In [77]: # Treinamento do Modelo 7
modelo7.fit(x_treino, y_treino)
```

```
Out[77]: SVC(probability=True)
```

```
In [78]: # Fazendo previsões com o Modelo 7
previsoes7 = modelo7.predict(x_teste)
```

Avaliando Modelo 7

```
In [79]: print("Relatório de Classificação:\n", classification_report(y_teste, previsoes7, digits=4))
print("Acurácia: %.2f%%" % (accuracy_score(y_teste, previsoes7) * 100))
print("AUC: %.2f" % (roc_auc_score(y_teste, previsoes7)))
```

	precision	recall	f1-score	support
0	0.9034	0.8992	0.9013	853
1	0.9001	0.9043	0.9022	857
accuracy				0.9018
macro avg	0.9018	0.9017	0.9018	1710
weighted avg	0.9018	0.9018	0.9018	1710

Acurácia: 90.18%
AUC: 0.90

Explicações e conclusões - Modelo 7

Paralelo com o caso anterior (diferenciando-se apenas na questão do balançamento das classes), este modelo pareceu ser bastante performático, obtendo excelente acurácia e AUC.

O treinamento para a classe 1 ficou ligeiramente superior se comparado à classe 0, porém, foi o melhor que obtive até o momento.

Salvando Modelo 7

```
In [80]: arquivo7 = 'modelos/modelo_classificador_svm.sav'
pickle.dump(modelo7, open(arquivo7, 'wb'))
```

Modelo escolhido: Modelo 7 - Support Vector Machines (SVM)

Motivo: ótima acurácia e AUC, além de ter aprendido de forma equilibrada sobre as duas classes.

Previsões com o modelo mais performático, coleta das probabilidades e salvando o trabalho

```
In [81]: # Baixando dados de teste
teste = pd.read_csv('dataset/projeto4_telecom_teste.csv', index_col = 0)
```

Ajustando dados de teste para igual estrutura de dados de treino

Criar um pipeline padrão de transformação para os dados de teste

- 1) transformar coluna state para numérica (churn rate)
- 2) retirar as colunas area_code, total_night_minutes, total_night_charge, total_intl_minutes e total_intl_charge
- 3) transformar colunas international_plan, voice_mail_plan em numérica e churn (0 e 1)
- 4) Normalizar dados
- 5) padronizar dados

```
In [82]: # Função 1 para executar o pipeline
# Preparando dataset completo
def ajuste_dataset_completo(dataset_teste):
    global teste
    # retirar a coluna area_code
    dataset_teste = dataset_teste.drop(['area_code', 'total_night_minutes', 'total_night_charge', 'total_intl_minutes', ], axis = 1)
    # substituindo os estados por seus churn rates
    for i in list(churn_rate.state):
        if i in list(churn_rate.state):
            dataset_teste.state.values[dataset_teste.state == i] = churn_rate.state.values[churn_rate.state == i]
    # salvando colunas churn, international_plan e voice_mail_plan em numéricas (0 e 1)
    dataset_teste.churn = [1 if item == 'yes' else 0 for item in dataset_teste.churn]
    dataset_teste.international_plan = [1 if item == 'yes' else 0 for item in dataset_teste.international_plan]
    dataset_teste.voice_mail_plan = [1 if item == 'yes' else 0 for item in dataset_teste.voice_mail_plan]
    # salvando novo dataset
    teste = dataset_teste
```

```
In [83]: # Aplicando pipeline parte 1 ao dataset de teste
ajuste_dataset_teste(teste)
```

```
In [84]: # Separando preditoras de target
x = teste.values[:, 0:14]
y = teste.values[:, 14]
```

```
In [85]: # Função 2 para executar o pipeline
# Normalizando e padronizando preditoras
def ajuste_preditoras(preditoras):
    global x
    # normalizando dados
    scaler = MinMaxScaler(feature_range = (0, 1))
    preditoras = scaler.fit_transform(preditoras)
    # padronizando dados
    escala2 = StandardScaler().fit(preditoras)
    preditoras = escala2.transform(preditoras)
    # salvando novo dataset
    x = preditoras
```

```
In [86]: # Aplicando pipeline parte 1 ao dataset de teste
ajuste_preditoras(x)
```

```
In [87]: # Baixando modelo
modelo = pickle.load(open(arquivo7, 'rb'))
```

```
In [88]: # Gerando previsões com o modelo escolhido
previsoes = modelo.predict(x)
```

```
In [89]: # Coletando e salvando as probabilidades calculadas pelo modelo
probabilidades = modelo.predict_proba(x)
```

Salvando os resultados

```
In [93]: # Previsões
previsoes = pd.DataFrame({'ID': np.arange(1, 1668), 'churn': previsoes})
previsoes = previs
```