

## PROJETO: Prevendo Customer Churn em Operadoras de Telecom

Data: 09/03/2020

Problema de negócio: prever se um cliente vai cancelar seu plano SIM ou NÃO.

In [2]:

```
## Bibliotecas
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sea
from sklearn.feature_selection import VarianceThreshold
import warnings
warnings.filterwarnings('ignore')
```

In [3]:

```
# Alteração de parâmetros para visualizar dataset com grande número de colunas
pd.set_option("display.max_columns", 180)
```

### Carregamento e limpeza dos dados

In [4]:

```
# Carregando o dataset
df = pd.read_csv('C:/FCD/PythonSpark/Projeto_telecom/projeto4_telecom_treino.csv')
df.head(10)
```

Unnamed: 0	state	account_length	area_code	international_plan	voice_mail_plan	number_vmail_messages	total_day_minutes	total_day_calls	total_day_charge	total_eve_minutes	total_eve_calls	total_eve_charge	total_night_minutes	total_night_calls	total_night_charge	total_intl_calls	total_intl_charge	number_customer_service_calls	churn
0	1	KS	128	area_code_415	no	yes	25	265.1	0	0	0	0	0	0	0	0	0	0	0
1	2	OH	107	area_code_415	no	yes	26	161.6	0	0	0	0	0	0	0	0	0	0	0
2	3	NJ	137	area_code_415	no	no	0	243.4	0	0	0	0	0	0	0	0	0	0	0
3	4	OK	84	area_code_408	yes	no	0	299.4	0	0	0	0	0	0	0	0	0	0	0
4	5	OK	75	area_code_415	yes	no	0	166.7	0	0	0	0	0	0	0	0	0	0	0
5	6	AL	118	area_code_510	yes	no	0	223.4	0	0	0	0	0	0	0	0	0	0	0
6	7	MA	121	area_code_510	no	yes	24	218.2	0	0	0	0	0	0	0	0	0	0	0
7	8	MO	147	area_code_415	yes	no	0	157.0	0	0	0	0	0	0	0	0	0	0	0
8	9	LA	117	area_code_408	no	no	0	184.5	0	0	0	0	0	0	0	0	0	0	0
9	10	WV	141	area_code_415	yes	yes	37	256.6	0	0	0	0	0	0	0	0	0	0	0

In [5]:

```
# Verificando tipos dos atributos
df.dtypes
```

```
Unnamed: 0      int64
state           object
account_length  int64
area_code       object
international_plan  object
voice_mail_plan  object
number_vmail_messages  int64
total_day_minutes  float64
total_day_calls   int64
total_day_charge  float64
total_eve_minutes  float64
total_eve_calls   int64
total_eve_charge  float64
total_night_minutes  float64
total_night_calls  int64
total_night_charge float64
total_intl_calls   int64
total_intl_charge  float64
number_customer_service_calls  int64
churn           object
dtype: object
```

In [6]:

```
# Quantidade linhas x colunas
df.shape
```

```
(3333, 21)
```

In [7]:

```
# Verificação de valores null
df.isnull().values.any()
```

```
# Verificação valores vazios
df.value_counts()
```

```
False
```

In [8]:

```
# Verificando o balanceamento da variável target (churn), na qual, valor 'no' = cliente não desistiu e 'yes' = cliente desistiu da operadora;
churn = pd.DataFrame(df.churn.value_counts())
churn[["churns"]] = 100 * churn["churn"] / df.shape[0]
churn
```

```
# Constata-se que a variável TARGET está proporcionalmente desbalanceada e, portanto, é necessário o balanceamento dos dados para aplicar o algoritmo de machine learning.
```

```
churn
churn % churns
no      2050  65.009503
yes     483   14.491497
```

In [9]:

```
# Remoção da variável 'Unnamed: 0', pois trata-se de apenas um índice do dataset;
# Remoção da variável 'number_customer_service_calls', pois refere-se a um número de serviço do cliente (irrelevante para o modelo)
```

```
df.drop(["Unnamed: 0"], axis=1, inplace=True)
df.drop(["number_customer_service_calls"], axis=1, inplace=True)
```

In [10]:

```
df.head(10)
```

state	account_length	area_code	international_plan	voice_mail_plan	number_vmail_messages	total_day_minutes	total_day_calls	total_eve_minutes	total_eve_calls	total_night_minutes	total_night_calls	total_night_charge	total_intl_calls	total_intl_charge	churn
0	KS	128	area_code_415	no	yes	25	265.1	0	0	0	0	0	0	0	0
1	OH	107	area_code_415	no	yes	26	161.6	0	0	0	0	0	0	0	0
2	NJ	137	area_code_415	no	no	0	243.4	0	0	0	0	0	0	0	0
3	OK	84	area_code_408	yes	no	0	299.4	0	0	0	0	0	0	0	0
4	OK	75	area_code_415	yes	no	0	166.7	0	0	0	0	0	0	0	0
5	AL	118	area_code_510	yes	no	0	223.4	0	0	0	0	0	0	0	0
6	MA	121	area_code_510	no	yes	24	218.2	0	0	0	0	0	0	0	0
7	MO	147	area_code_415	yes	no	0	157.0	0	0	0	0	0	0	0	0
8	LA	117	area_code_408	no	no	0	184.5	0	0	0	0	0	0	0	0
9	WV	141	area_code_415	yes	yes	37	256.6	0	0	0	0	0	0	0	0

In [11]:

```
# Limpeza da variável 'area_code'
df.area_code = df.area_code.apply(lambda x: x.split('-')[2])
df.head()
```

state	account_length	area_code	international_plan	voice_mail_plan	number_vmail_messages	total_day_minutes	total_day_calls	total_eve_minutes	total_eve_calls	total_night_minutes	total_night_calls	total_night_charge	total_intl_calls	total_intl_charge	churn
0	KS	128	415	no	yes	25	265.1	0	0	0	0	0	0	0	0
1	OH	107	415	no	yes	26	161.6	0	0	0	0	0	0	0	0
2	NJ	137	415	no	no	0	243.4	0	0	0	0	0	0	0	0
3	OK	84	408	yes	no	0	299.4	0	0	0	0	0	0	0	0
4	OK	75	415	yes	no	0	166.7	0	0	0	0	0	0	0	0

In [12]:

```
# Convertendo variáveis categóricas para variáveis numéricas (0 e 1): Label Encoding
from sklearn.preprocessing import LabelEncoder
```

```
number = LabelEncoder()
df.international_plan = number.fit_transform(df.international_plan)
```

```
df.voice_mail_plan = number.fit_transform(df.voice_mail_plan)
```

```
df.state = number.fit_transform(df.state)
```

```
churn = number.fit_transform(df.churn)
```

```
df.head()
```

state	account_length	area_code	international_plan	voice_mail_plan	number_vmail_messages	total_day_minutes	total_day_calls	total_eve_minutes	total_eve_calls	total_night_minutes	total_night_calls	total_night_charge	total_intl_calls	total_intl_charge	churn
0	16	128	415	0	1	25	265.1	0	0	0	0	0	0	0	0
1	26	107	415	0	1	26	161.6	0	0	0	0	0	0	0	0
2	31	137	415	0	0	0	243.4	0	0	0	0	0	0	0	0
3	35	84	408	1	0	0	299.4	0	0	0	0	0	0	0	0
4	36	75	415	1	0	0	166.7	0	0	0	0	0	0	0	0

In [13]:

```
# Convertendo 'area_code' to numeric
df.area_code = pd.to_numeric(df['area_code'])
```

### Análise exploratória (Estatística descritiva)

In [14]:

```
# Análise de distribuição dos dados através de um histograma;
# Embarra a maior parte dos dados apresentarem distribuição normal, será necessário a padronização dos dados para o algoritmo Logistic Regression;
df.hist()
plt.show()
```



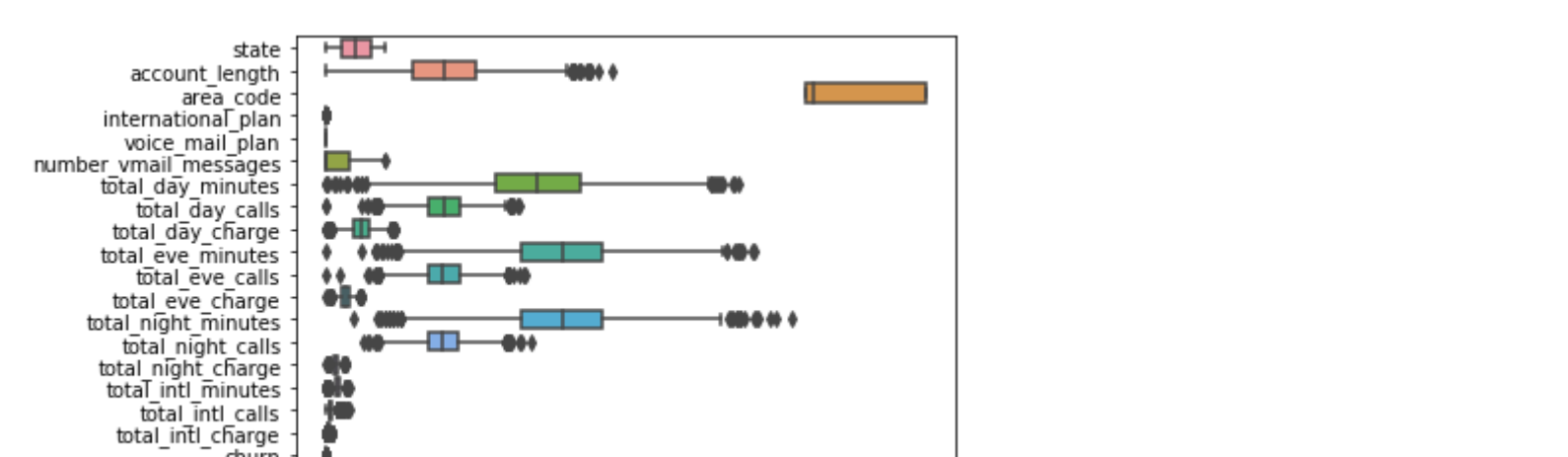
In [17]:

```
# Matriz de correlação com seaborn
corr_matrix = df.corr(method='pearson')
```

```
plt.figure(figsize=(12,10))
plt.title('Correlation Heatmap')
```

```
sea.heatmap(corr_matrix, square=True, cmap='coolwarm', ticklabels=True, yticklabels=True,
            annot=True, cbar=True, linewidths=5, linecolor='white', vmin=-1, vmax=1)
```

```
plt.show()
```



In [18]:

```
# Correlação com a variável target 'churn'
corr_target = abs(corr_matrix["churn"])
```

```
# Selecionando highly correlated features
```

```
# A correlação entre a variável target e as demais é baixa; apenas 4 variáveis mostraram uma correlação de apenas >=20 %.
```

```
relevant_features = corr_target[corr_target>0.2]
relevant_features
```

```
international_plan    0.259892
total_day_charge      0.285151
total_day_minutes     0.600600
Name: churn, dtype: float64
```

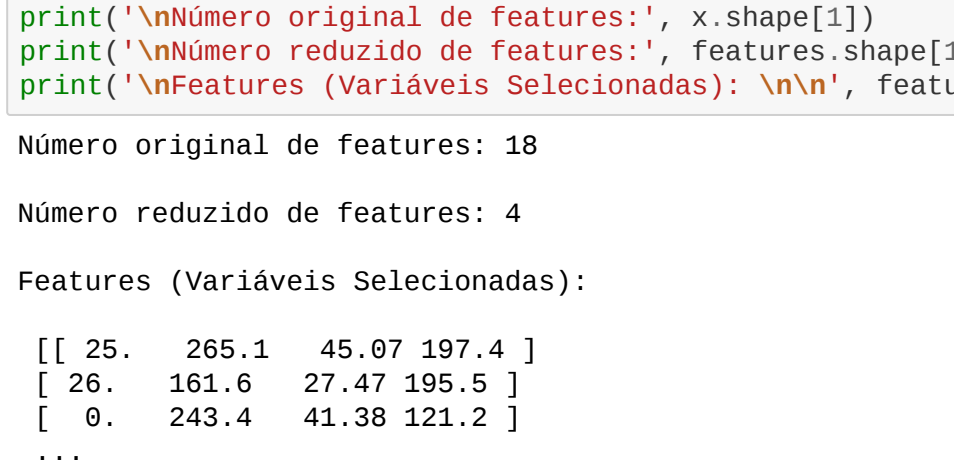
In [19]:

```
### Boxplot
sea.boxplot(data= df, orient='h')
```

```
# Verifica se diversos outliers presentes nas variáveis;
```

Out[19]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1906f9506a6>
```



### Pré-Processamento

#### Feature Selection - Seleção Univariada

In [20]:

```
# Extração de Variáveis com Testes Estatísticos Univariados
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
```

```
# Separando input e output
x= df.iloc[:,0:18]
y=df.iloc[:,18]
```

```
# Função para seleção de variáveis (f)
```

```
best_var = SelectKBest(score_func = chi2, k = 4)
```

```
# Executa a função de pontuação em (X, y) e obtém os recursos selecionados
```

```
fit = best_var.fit(x, y)
```

```
# Reduz X para os recursos selecionados
```

```
features = fit.transform(x)
```

```
# Resultados
```

```
print("\nNúmero original de features:", x.shape[1])
```

```
print("\nNúmero reduzido de features:", features.shape[1])
```

```
print("\nFeatures (Variáveis Selecionadas):", '\n\n', features)
```

```
Número original de features: 18
```

```
Número reduzido de features: 4
```

```
Features (Variáveis Selecionadas):
```

```
[ 25.  265.1  45.07 197.4 ]
```

```
[ 26.  161.6  27.47 195.5 ]
```

```
[ 0.  243.4  41.38 121.2 ]
```

```
...
```

```
[ 0.  189.8  39.74 288.8 ]
```

```
[ 0.  213.9  36.35 159.6 ]
```

```
[ 25.  234.4  39.85 265.9 ]
```

In [21]:

```
# As 4 variáveis mais relevantes conforme o teste univariado (qui-quadrado), foram:
```

```
number_vmail_messages, total_day_minutes, total_day_charges, total_eve_minutes
```

```
df_new = pd.DataFrame(data=features)
```

```
df_new.head()
```

ational_plan',
al_day_minutes',
ve_minutes',
ight_minutes',
al_intl_minutes',
0.01808938
0.06148953
0.05567697]

#### Feature Selection - Recursive Feature Elimination

In [22]:

```
# Import dos módulos
```

```
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression
```

```
import warnings
warnings.filterwarnings('ignore')
```

```
# Criação do modelo
```

```
modelo = LogisticRegression()
```

```
# RFE
```

```
rfe = RFE(modelo, 3)
```

```
fit = rfe.fit(x, y)
```

```
# Print dos resultados
```

```
print("Variáveis Preditoras:", x.columns[0:18])
```

```
print("Variáveis Selecionadas: %d" % fit.support_)
```

```
print("Ranking dos Atributos: %s" % fit.ranking_)
```

```
print("Número de Melhores Atributos: %d" % fit.n_features_)
```

```
Variáveis Preditoras: Index(['state', 'account_length', 'area_code', 'international_plan', 'voice_mail_plan', 'number_vmail_messages', 'total_day_minutes', 'total_day_calls', 'total_eve_minutes', 'total_eve_calls', 'total_night_minutes', 'total_night_calls', 'total_night_charge', 'total_intl_calls', 'total_intl_charge'], dtype='object')
```

```
Variáveis Selecionadas: False False False True True False False False False False False False
```

```
False False False True False
```

```
Ranking dos Atributos: [11 16 5 1 1 4 3 8 12 7 10 13 9 15 2 1 6]
```

```
Número de Melhores Atributos: 3
```

In [ ]:

```
# 4 Melhores variáveis conforme método RFE: 'international_plan', 'voice_mail_plan', 'total_intl_calls', 'total_eve_charge'
```

### Feature Selection - Ensemble

In [23]:

```
# Importância do Atributo com o Extra Trees Classifier
```

```
# Import dos Módulos
```

```
from sklearn.ensemble import ExtraTreesClassifier
```

```
# Criação do Modelo - Feature Selection
```

```
modelo = ExtraTreesClassifier()
```

```
modelo.fit(x, y)
```

```
# Print dos Resultados
```

```
print(x.columns[0:18])
```

```
print(modelo.feature_importances_)
```

```
df['state', 'account_length', 'area_code', 'international_plan', 'voice_mail_plan', 'number_vmail_messages', 'total_day_minutes', 'total_day_calls', 'total_eve_minutes', 'total_eve_calls', 'total_night_minutes', 'total_night_calls', 'total_night_charge', 'total_intl_calls', 'total_intl_charge']
```

```
dtype='object'
```

```
[0.0425953  0.04259807 0.02710827 0.07384055 0.01872  0.01890938
```

```
0.12242637 0.04776417 0.13252386 0.06236894 0.04316138 0.06148953
```

```
0.04592338 0.04526765 0.04754607 0.05229684 0.06058713 0.05567697]
```

In [24]:

```
# Ordenando o resultado do melhor para o menor
```

```
# As 3 variáveis mais importante neste método foram:
```

```
x = 'total_day_charges', 'total_day_minutes', 'international_plan'
```

```
df_importance = pd.DataFrame(data=df, columns=x)
```

```
df_importance.sort_values(by=x[0], ascending= False)
```

0	1	2	3
8	0.132524		
6	0.122426		
3	0.071841		
9	0.062369		
11	0.061410		
16	0.060587		
17	0.055677		
15	0.052297		
7	0.047764		
14	0.047546		
12	0.046923		
13	0.046268		
10	0.043161		
0	0.042595		
1	0.042398		
2	0.027108		
4	0.018720		
5	0.018009		

### PCA (Principal Component Analysis)