```python
import pandas as pd
import numpy as np
from sklearn import preprocessing
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
```

```python
%matplotlib inline
```

```python
df = pd.read_csv('WA_Fn-UseC_-Telco-Customer-Churn.csv')
```

```python
df.head
```

```
<bound method NDFrame.head of       customerID  gender  SeniorCitizen Partner Dependents  tenure  \
0     7590-VHVEG  Female              0     Yes         No       1
1     5575-GNVDE    Male              0      No         No      34
2     3668-QPYBK    Male              0      No         No       2
3     7795-CFOCW    Male              0      No         No      45
4     9237-HQITU  Female              0      No         No       2
...          ...     ...            ...     ...        ...     ...
7038  6840-RESVB    Male              0     Yes        Yes      24
7039  2234-XADUH  Female              0     Yes        Yes      72
7040  4801-JZAZL  Female              0     Yes        Yes      11
7041  8361-LTMKD    Male              1     Yes         No       4
7042  3186-AJIEK    Male              0      No         No      66

     PhoneService     MultipleLines InternetService OnlineSecurity ...  \
0              No  No phone service             DSL             No ...
1             Yes                No             DSL            Yes ...
2             Yes                No             DSL            Yes ...
3              No  No phone service             DSL            Yes ...
4             Yes                No     Fiber optic             No ...
...           ...               ...             ...            ... ...
7038          Yes               Yes             DSL            Yes ...
7039          Yes               Yes     Fiber optic             No ...
7040           No  No phone service             DSL            Yes ...
7041          Yes               Yes     Fiber optic             No ...
7042          Yes                No     Fiber optic            Yes ...

     DeviceProtection TechSupport StreamingTV StreamingMovies        Contract  \
0                  No          No          No              No  Month-to-month
1                 Yes          No          No              No        One year
2                  No          No          No              No  Month-to-month
3                 Yes         Yes          No              No        One year
4                  No          No          No              No  Month-to-month
...               ...         ...         ...             ...             ...
7038              Yes         Yes         Yes             Yes        One year
7039              Yes          No         Yes             Yes        One year
7040               No          No          No              No  Month-to-month
7041               No          No          No              No  Month-to-month
7042              Yes         Yes         Yes             Yes        Two year

     PaperlessBilling              PaymentMethod MonthlyCharges  TotalCharges  \
0                 Yes           Electronic check          29.85         29.85
1                  No               Mailed check          56.95        1889.5
2                 Yes               Mailed check          53.85        108.15
3                  No  Bank transfer (automatic)          42.30       1840.75
4                 Yes           Electronic check          70.70        151.65
...               ...                        ...            ...           ...
7038              Yes               Mailed check          84.80        1990.5
7039              Yes    Credit card (automatic)         103.20        7362.9
7040              Yes           Electronic check          29.60        346.45
7041              Yes               Mailed check          74.40         306.6
7042              Yes  Bank transfer (automatic)         105.65        6844.5

     Churn
0       No
1       No
2      Yes
3       No
4      Yes
...    ...
7038    No
7039    No
7040    No
7041   Yes
7042    No

[7043 rows x 21 columns]>
```

```python
data = pd.DataFrame(df)
```
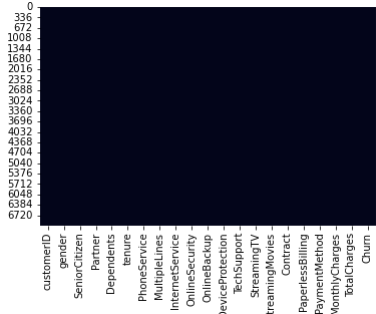
```python
data.head()
```

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMethod |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | No | No | No | No | Month-to-month | Yes | Electronic check |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | Yes | No | No | No | One year | No | Mailed check |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... | No | No | No | No | Month-to-month | Yes | Mailed check |
| 3 | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | ... | Yes | Yes | No | No | One year | No | Bank transfer (automatic) |
| 4 | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | ... | No | No | No | No | Month-to-month | Yes | Electronic check |

5 rows × 21 columns

```python
sns.heatmap(data.isnull(),cbar=False)
```

```
<AxesSubplot:>
```

```
In [8]:  df = df.drop(['customerID'], axis = 1)
         df.head()
```

Out[8]:

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMethod |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | Yes | No | No | No | No | Month-to-month | Yes | Electronic check |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | No | Yes | No | No | No | One year | No | Mailed check |
| 2 | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | Yes | No | No | No | No | Month-to-month | Yes | Mailed check |
| 3 | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | No | Yes | Yes | No | No | One year | No | Bank transfer (automatic) |
| 4 | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | No | No | No | No | No | Month-to-month | Yes | Electronic check |

```
In [9]:  df['TotalCharges'] = pd.to_numeric(df.TotalCharges, errors='coerce')
         df.isnull().sum()
```

Out[9]:
```
gender              0
SeniorCitizen       0
Partner             0
Dependents          0
tenure              0
PhoneService        0
MultipleLines       0
InternetService     0
OnlineSecurity      0
OnlineBackup        0
DeviceProtection    0
TechSupport         0
StreamingTV         0
StreamingMovies     0
Contract            0
PaperlessBilling    0
PaymentMethod       0
MonthlyCharges      0
TotalCharges       11
Churn               0
dtype: int64
```

```
In [10]:  df[np.isnan(df['TotalCharges'])]
```

Out[10]:

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMeth |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 488 | Female | 0 | Yes | Yes | 0 | No | No phone service | DSL | Yes | No | Yes | Yes | Yes | No | Two year | Yes | Bank trans (automat |
| 753 | Male | 0 | No | Yes | 0 | Yes | No | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | Two year | No | Mailed che |
| 936 | Female | 0 | Yes | Yes | 0 | Yes | No | DSL | Yes | Yes | No | No | Yes | Yes | Two year | No | Mailed che |
| 1082 | Male | 0 | Yes | Yes | 0 | Yes | Yes | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | Two year | No | Mailed che |
| 1340 | Female | 0 | Yes | Yes | 0 | No | No phone service | DSL | Yes | Yes | Yes | Yes | Yes | No | Two year | No | Credit ca (automat |
| 3331 | Male | 0 | Yes | Yes | 0 | Yes | No | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | Two year | No | Mailed che |
| 3826 | Male | 0 | Yes | Yes | 0 | Yes | Yes | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | Two year | No | Mailed che |
| 4380 | Female | 0 | Yes | Yes | 0 | Yes | No | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | Two year | No | Mailed che |
| 5218 | Male | 0 | Yes | Yes | 0 | Yes | No | No | No internet service | No internet service | No internet service | No internet service | No internet service | No internet service | One year | Yes | Mailed che |
| 6670 | Female | 0 | Yes | Yes | 0 | Yes | Yes | DSL | No | Yes | Yes | Yes | Yes | No | Two year | No | Mailed che |
| 6754 | Male | 0 | No | Yes | 0 | Yes | Yes | DSL | Yes | Yes | No | Yes | No | No | Two year | Yes | Bank trans (automat |

```
In [11]:  df[df['tenure'] == 0].index
```

Out[11]:  Int64Index([488, 753, 936, 1082, 1340, 3331, 3826, 4380, 5218, 6670, 6754], dtype='int64')

```
In [12]:  df.drop(labels=df[df['tenure'] == 0].index, axis=0, inplace=True)
          df[df['tenure'] == 0].index
```

Out[12]:  Int64Index([], dtype='int64')

```
In [13]:  df.fillna(df["TotalCharges"].mean())
```

Out[13]:

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMeth |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | Yes | No | No | No | No | Month-to-month | Yes | Electronic che |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | No | Yes | No | No | No | One year | No | Mailed che |
| 2 | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | Yes | No | No | No | No | Month-to-month | Yes | Mailed che |
| 3 | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | No | Yes | Yes | No | No | One year | No | Bank trans (automat |
| 4 | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | No | No | No | No | No | Month-to-month | Yes | Electronic che |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 7038 | Male | 0 | Yes | Yes | 24 | Yes | Yes | DSL | Yes | No | Yes | Yes | Yes | Yes | One year | Yes | Mailed che |
| 7039 | Female | 0 | Yes | Yes | 72 | Yes | Yes | Fiber optic | No | Yes | Yes | No | Yes | Yes | One year | Yes | Credit ca (automat |
| 7040 | Female | 0 | Yes | Yes | 11 | No | No phone service | DSL | Yes | No | No | No | No | No | Month-to-month | Yes | Electronic che |
| 7041 | Male | 1 | Yes | No | 4 | Yes | Yes | Fiber optic | No | No | No | No | No | No | Month-to-month | Yes | Mailed che |
| 7042 | Male | 0 | No | No | 66 | Yes | No | Fiber optic | Yes | No | Yes | Yes | Yes | Yes | Two year | Yes | Bank transf (automat |

7032 rows × 20 columns

```
In [14]:  df.isnull().sum()
```

```
Out[14]:    gender            0
            SeniorCitizen     0
            Partner           0
            Dependents        0
            tenure            0
            PhoneService      0
            MultipleLines     0
            InternetService   0
            OnlineSecurity    0
            OnlineBackup      0
            DeviceProtection  0
            TechSupport       0
            StreamingTV       0
            StreamingMovies   0
            Contract          0
            PaperlessBilling  0
            PaymentMethod     0
            MonthlyCharges    0
            TotalCharges      0
            Churn             0
            dtype: int64
```

```python
In [15]: df["SeniorCitizen"]= df["SeniorCitizen"].map({0: "No", 1: "Yes"})
         df.head()
```

Out[15]:

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMethod |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | No | Yes | No | 1 | No | No phone service | DSL | No | Yes | No | No | No | No | Month-to-month | Yes | Electronic check |
| 1 | Male | No | No | No | 34 | Yes | No | DSL | Yes | No | Yes | No | No | No | One year | No | Mailed check |
| 2 | Male | No | No | No | 2 | Yes | No | DSL | Yes | Yes | No | No | No | No | Month-to-month | Yes | Mailed check |
| 3 | Male | No | No | No | 45 | No | No phone service | DSL | Yes | No | Yes | Yes | No | No | One year | No | Bank transfer (automatic) |
| 4 | Female | No | No | No | 2 | Yes | No | Fiber optic | No | No | No | No | No | No | Month-to-month | Yes | Electronic check |

```python
In [16]: df["InternetService"].describe()
```

```
Out[16]: count         7032
         unique           3
         top     Fiber optic
         freq          3096
         Name: InternetService, dtype: object
```

```python
In [17]: df["InternetService"].describe(include=['object', 'bool'])
```

```
Out[17]: count         7032
         unique           3
         top     Fiber optic
         freq          3096
         Name: InternetService, dtype: object
```

```python
In [18]: numerical_cols = ['tenure', 'MonthlyCharges', 'TotalCharges']
         df[numerical_cols].describe()
```

Out[18]:

| | tenure | MonthlyCharges | TotalCharges |
|---|---|---|---|
| count | 7032.000000 | 7032.000000 | 7032.000000 |
| mean | 32.421786 | 64.798208 | 2283.300441 |
| std | 24.545260 | 30.085974 | 2266.771362 |
| min | 1.000000 | 18.250000 | 18.800000 |
| 25% | 9.000000 | 35.587500 | 401.450000 |
| 50% | 29.000000 | 70.350000 | 1397.475000 |
| 75% | 55.000000 | 89.862500 | 3794.737500 |
| max | 72.000000 | 118.750000 | 8684.800000 |

```python
In [19]: df["Churn"][df["Churn"]=="No"].groupby(by=df["gender"]).count()
```

```
Out[19]: gender
         Female    2544
         Male      2619
         Name: Churn, dtype: int64
```
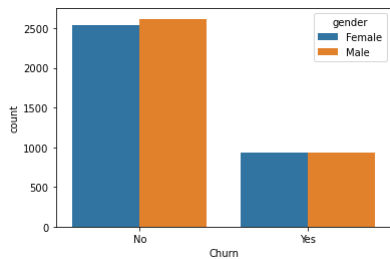
```python
In [20]: df["Churn"][df["Churn"]=="Yes"].groupby(by=df["gender"]).count()
```

```
Out[20]: gender
         Female    939
         Male      930
         Name: Churn, dtype: int64
```
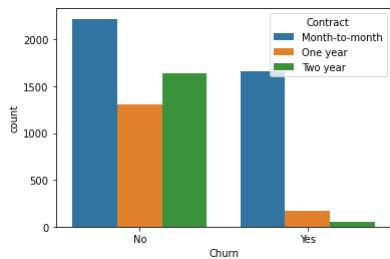
```python
In [21]: sns.countplot(x='Churn', data= df,hue='gender')
```
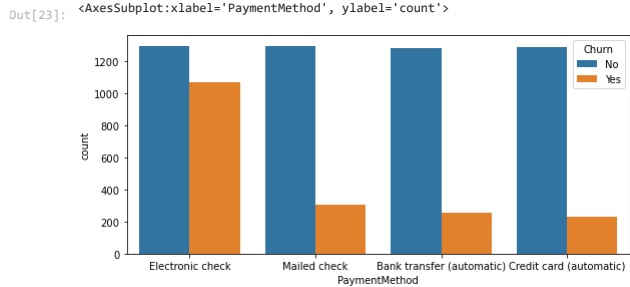
Out[21]: <AxesSubplot:xlabel='Churn', ylabel='count'>



```python
In [22]: sns.countplot(x='Churn', data= df,hue='Contract')
```

Out[22]: <AxesSubplot:xlabel='Churn', ylabel='count'>



```python
In [23]: plt.figure(figsize=(9,4))
         sns.countplot(x='PaymentMethod', data= df, hue='Churn')
```
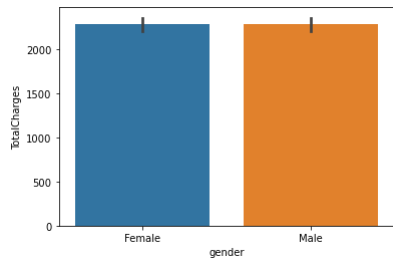
```
Out[23]: <AxesSubplot:xlabel='PaymentMethod', ylabel='count'>
```



```
In [24]: sns.barplot(x='gender',y='TotalCharges',data=df)
```

```
Out[24]: <AxesSubplot:xlabel='gender', ylabel='TotalCharges'>
```



```
In [25]: df["InternetService"].unique()
```

```
Out[25]: array(['DSL', 'Fiber optic', 'No'], dtype=object)
```

```
In [26]: df[df["gender"]=="Male"][["InternetService", "Churn"]].value_counts()
```
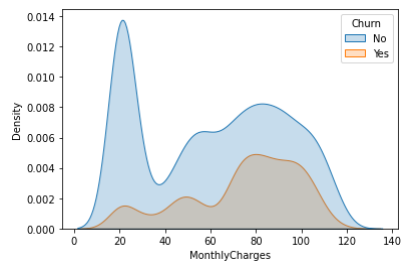
```
Out[26]: InternetService  Churn
         DSL              No       992
         Fiber optic      No       910
         No               No       717
         Fiber optic      Yes      633
         DSL              Yes      240
         No               Yes       57
         dtype: int64
```

```
In [27]: df[df["gender"]=="Female"][["InternetService", "Churn"]].value_counts()
```

```
Out[27]: InternetService  Churn
         DSL              No       965
         Fiber optic      No       889
         No               No       690
         Fiber optic      Yes      664
         DSL              Yes      219
         No               Yes       56
         dtype: int64
```

```
In [28]: sns.kdeplot(x='MonthlyCharges',data = df,hue= 'Churn',shade=True)
```

```
Out[28]: <AxesSubplot:xlabel='MonthlyCharges', ylabel='Density'>
```



```
In [29]: sns.set_context("paper",font_scale=1.1)
         ax = sns.kdeplot(df.MonthlyCharges[(df["Churn"] == 'No') ],
                          color="Red", shade = True);
         ax = sns.kdeplot(df.MonthlyCharges[(df["Churn"] == 'Yes') ],
                          ax =ax, color="Blue", shade= True);
         ax.legend(["Not Churn","Churn"],loc='upper right');
         ax.set_ylabel('Density');
         ax.set_xlabel('Monthly Charges');
         ax.set_title('Distribution of monthly charges by churn');
```
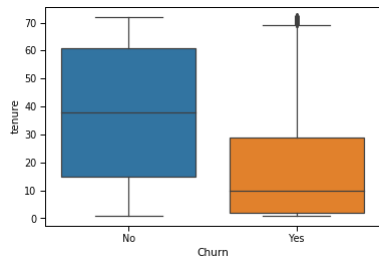


```
In [30]: sns.kdeplot(x='TotalCharges',data = df,hue= 'Churn',shade=True)
```

```
Out[30]: <AxesSubplot:xlabel='TotalCharges', ylabel='Density'>
```



```
In [31]: sns.boxplot(x='Churn',y='tenure',data=df)
```

```
In [32]: def object_to_int(dataframe_series):
             if dataframe_series.dtype=='object':
                 dataframe_series = LabelEncoder().fit_transform(dataframe_series)
             return dataframe_series
```

```
In [33]: df = df.apply(lambda x: object_to_int(x))
         df.head()
```

Out[33]:

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV | StreamingMovies | Contract | PaperlessBilling | PaymentMethod |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 2 |
| 1 | 1 | 0 | 0 | 0 | 34 | 1 | 0 | 0 | 2 | 0 | 2 | 0 | 0 | 0 | 1 | 0 | 3 |
| 2 | 1 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 3 |
| 3 | 1 | 0 | 0 | 0 | 45 | 0 | 1 | 0 | 2 | 0 | 2 | 2 | 0 | 0 | 1 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 |

```
In [34]: plt.figure(figsize=(14,7))
         df.corr()['Churn'].sort_values(ascending = False)
```

```
Out[34]: Churn               1.000000
         MonthlyCharges      0.192858
         PaperlessBilling    0.191454
         SeniorCitizen       0.150541
         PaymentMethod       0.107852
         MultipleLines       0.038043
         PhoneService        0.011691
         gender             -0.008545
         StreamingTV        -0.036303
         StreamingMovies    -0.038802
         InternetService    -0.047097
         Partner            -0.149982
         Dependents         -0.163128
         DeviceProtection   -0.177883
         OnlineBackup       -0.195290
         TotalCharges       -0.199484
         TechSupport        -0.282232
         OnlineSecurity     -0.289050
         tenure             -0.354049
         Contract           -0.396150
         Name: Churn, dtype: float64
         <Figure size 1008x504 with 0 Axes>
```

```
In [35]: X = df.drop(columns = ['Churn'])
         y = df['Churn'].values
```

```
In [36]: X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.30, random_state = 40, stratify=y)
```
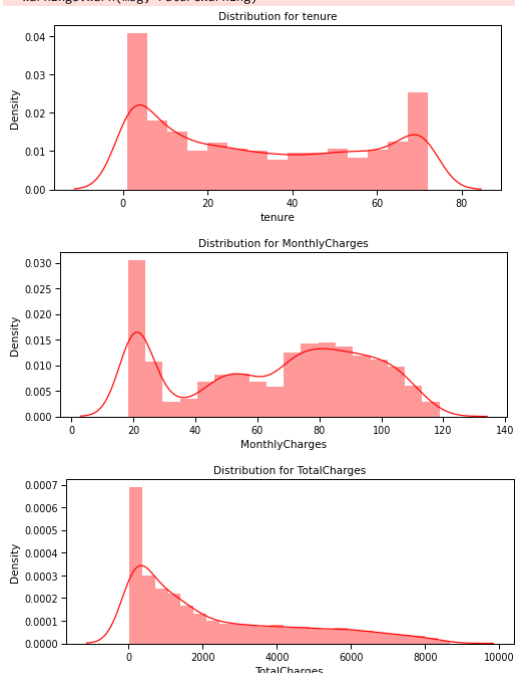
```
In [37]: def distplot(feature, frame, color='r'):
             plt.figure(figsize=(8,3))
             plt.title("Distribution for {}".format(feature))
             ax = sns.distplot(frame[feature], color= color)
```

```
In [38]: num_cols = ["tenure", 'MonthlyCharges', 'TotalCharges']
         for feat in num_cols: distplot(feat, df)
```
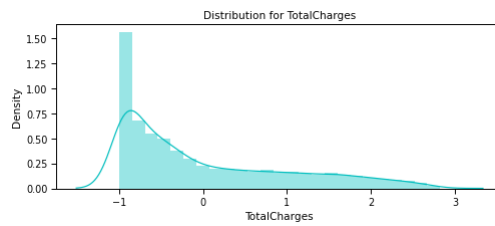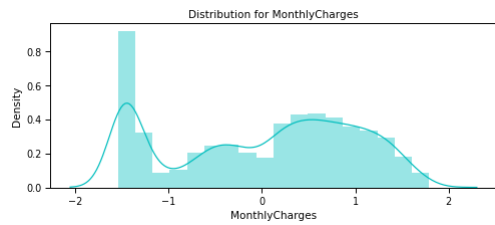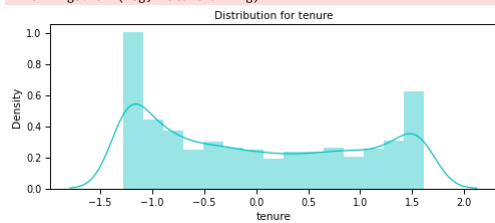
C:\Users\alokr\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
C:\Users\alokr\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
C:\Users\alokr\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)







```
In [39]: from sklearn.preprocessing import StandardScaler
```

```
In [40]: df_std = pd.DataFrame(StandardScaler().fit_transform(df[num_cols].astype('float64')),
                               columns=num_cols)
         for feat in numerical_cols: distplot(feat, df_std, color='c')
```

Distribution for tenure



Distribution for MonthlyCharges



Distribution for TotalCharges

In [41]:
```python
# Divide the columns into 3 categories, one ofor standardisation, one for label encoding and one for one hot encoding

cat_cols_ohe =['PaymentMethod', 'Contract', 'InternetService'] # those that need one-hot encoding
cat_cols_le = list(set(X_train.columns)- set(num_cols) - set(cat_cols_ohe)) #those that need label encoding
```

In [42]:
```python
scaler= StandardScaler()

X_train[num_cols] = scaler.fit_transform(X_train[num_cols])
X_test[num_cols] = scaler.transform(X_test[num_cols])
```

In [43]:
```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn import metrics
from sklearn.metrics import roc_curve
from sklearn.metrics import recall_score,confusion_matrix, precision_score, f1_score,accuracy_score, classification_report
```

In [44]:
```python
knn_model = KNeighborsClassifier(n_neighbors = 11)
knn_model.fit(X_train,y_train)
predicted_y = knn_model.predict(X_test)
accuracy_knn = knn_model.score(X_test,y_test)
print("KNN accuracy:",accuracy_knn)
```

KNN accuracy: 0.7753554502369668

In [45]:
```python
model_rf = RandomForestClassifier(n_estimators=500 , oob_score = True, n_jobs = -1,
                                  random_state =50, max_features = "auto",
                                  max_leaf_nodes = 30)
model_rf.fit(X_train, y_train)

# Make predictions
prediction_test = model_rf.predict(X_test)
print (metrics.accuracy_score(y_test, prediction_test))
```
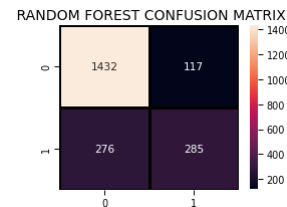
0.8137440758293839

In [46]:
```python
plt.figure(figsize=(4,3))
sns.heatmap(confusion_matrix(y_test, prediction_test),
            annot=True,fmt = "d",linecolor="k",linewidths=3)

plt.title(" RANDOM FOREST CONFUSION MATRIX",fontsize=14)
plt.show()
```



In [47]:
```python
lr_model = LogisticRegression()
lr_model.fit(X_train,y_train)
accuracy_lr = lr_model.score(X_test,y_test)
print("Logistic Regression accuracy is :",accuracy_lr)
```

Logistic Regression accuracy is : 0.8090047393364929

In [48]:
```python
dt_model = DecisionTreeClassifier()
dt_model.fit(X_train,y_train)
predictdt_y = dt_model.predict(X_test)
accuracy_dt = dt_model.score(X_test,y_test)
print("Decision Tree accuracy is :",accuracy_dt)
```

Decision Tree accuracy is : 0.737914691943128

In [ ]:

In [ ]: