

# Bank Marketing Campaign Analysis

## What is a Term Deposit?

A **Term deposit** is a deposit that a bank or a financial institution offers with a fixed rate (often better than just opening deposit account) in which your money will be returned back at a specific maturity time.

## Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from plotly import tools
import plotly.plotly as py
import plotly.figure_factory as ff
import plotly.graph_objs as go
from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
init_notebook_mode(connected=True)

MAIN_PATH = '../input/'
df = pd.read_csv(MAIN_PATH + 'bank.csv')
term_deposits = df.copy()
# Have a grasp of how our data looks.
df.head()
```

Out[1]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	pou
0	59	admin.	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	ur
1	56	admin.	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	ur
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	ur
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	ur
4	54	admin.	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	ur

## Summary:

The mean age is approximately 41 years old (with a minimum of 18 years old and a maximum of 95 years old). The mean balance is 1,528. However, the standard deviation (std) is high, indicating significant distribution across the dataset.

As the data information suggests, it is advisable to drop the duration column since it is highly correlated with whether a potential client will buy a term deposit. Also, the duration is obtained after the call is made to the potential client, so if the target client has never received calls, this feature is not that useful. The reason why duration is highly correlated with opening a term deposit is that the more the bank talks to a target client, the higher the probability the target client will open a term deposit, as a higher duration implies higher interest (commitment) from the potential client.

Note: There aren't many insights we can gain from the descriptive dataset since most of the valuable information is located not in the "numeric" columns but in the "categorical columns".

In [2]: df.describe()

Out[2]:

	age	balance	day	duration	campaign	pdays	previous
count	11162.000000	11162.000000	11162.000000	11162.000000	11162.000000	11162.000000	11162.000000
mean	41.231948	1528.538524	15.658036	371.993818	2.508421	51.330407	0.832557
std	11.913369	3225.413326	8.420740	347.128386	2.722077	108.758282	2.292007
min	18.000000	-6847.000000	1.000000	2.000000	1.000000	-1.000000	0.000000
25%	32.000000	122.000000	8.000000	138.000000	1.000000	-1.000000	0.000000
50%	39.000000	550.000000	15.000000	255.000000	2.000000	-1.000000	0.000000
75%	49.000000	1708.000000	22.000000	496.000000	3.000000	20.750000	1.000000
max	95.000000	81204.000000	31.000000	3881.000000	63.000000	854.000000	58.000000

There are no missing values in the dataset. In the event of missing values, common strategies include filling them with the median, mean, or mode. While I generally prefer using the median, in this scenario, there is no need to fill any missing values.

```
In [3]: # No missing values.
```

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11162 entries, 0 to 11161
Data columns (total 17 columns):
age            11162 non-null int64
job             11162 non-null object
marital         11162 non-null object
education       11162 non-null object
default          11162 non-null object
balance          11162 non-null int64
housing          11162 non-null object
loan             11162 non-null object
contact          11162 non-null object
day              11162 non-null int64
month            11162 non-null object
duration         11162 non-null int64
campaign         11162 non-null int64
pdays            11162 non-null int64
previous         11162 non-null int64
poutcome         11162 non-null object
deposit           11162 non-null object
dtypes: int64(7), object(10)
memory usage: 1.4+ MB
```

```
In [4]: f, ax = plt.subplots(1,2, figsize=(16,8))
```

```
colors = ["#FA5858", "#64FE2E"]
labels = "Did not Open Term Suscriptions", "Opened Term Suscriptions"

plt.suptitle('Information on Term Suscriptions', fontsize=20)

df["deposit"].value_counts().plot.pie(explode=[0,0.25], autopct='%1.2f%%', ax=ax[0], shadow=True, colors=colors,
                                      labels=labels, fontsize=12, startangle=25)

# ax[0].set_title('State of Loan', fontsize=16)
ax[0].set_ylabel('% of Condition of Loans', fontsize=14)

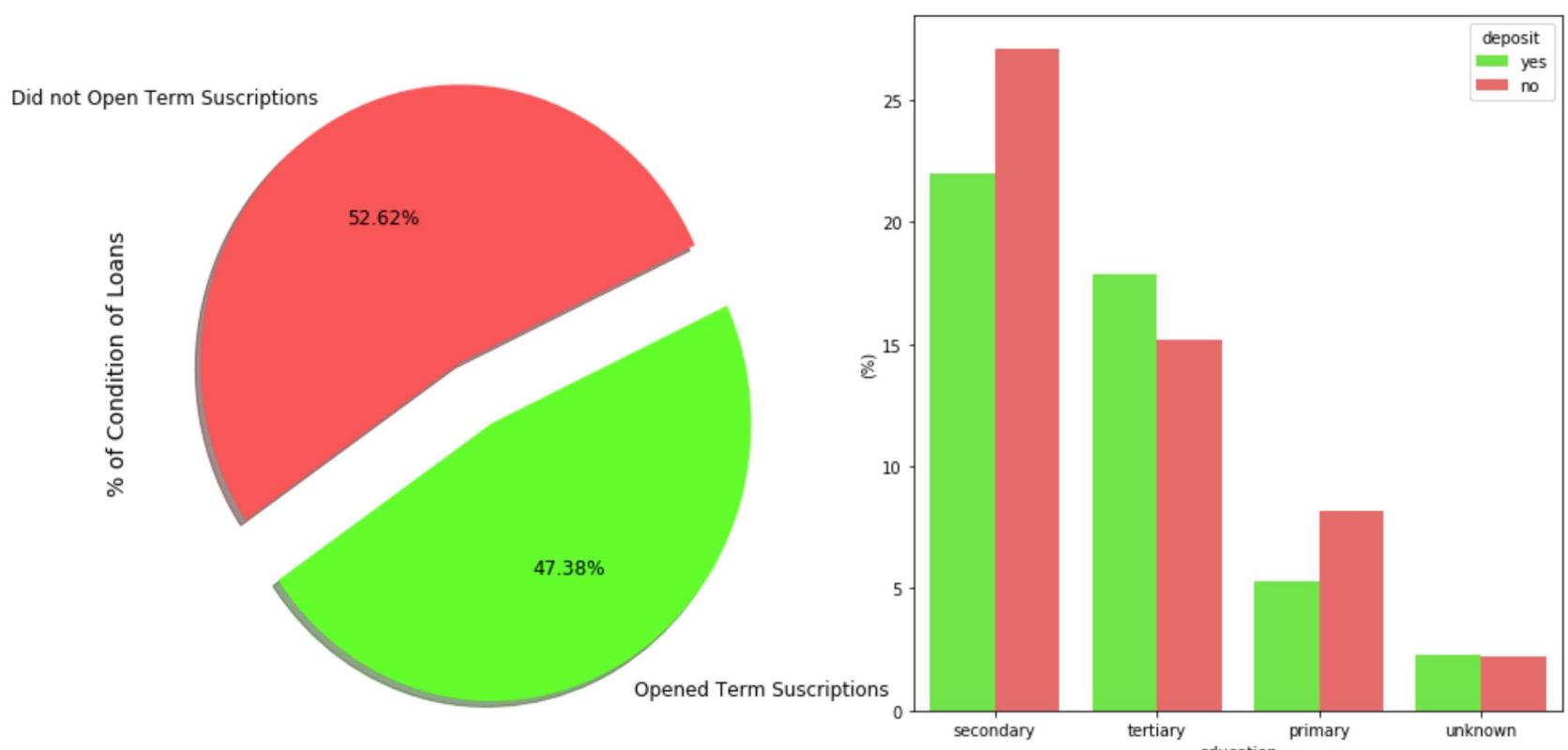
# sns.countplot('loan_condition', data=df, ax=ax[1], palette=colors)
# ax[1].set_title('Condition of Loans', fontsize=20)
# ax[1].set_xticklabels(['Good', 'Bad'], rotation='horizontal')
palette = ["#64FE2E", "#FA5858"]

sns.barplot(x="education", y="balance", hue="deposit", data=df, palette=palette, estimator=lambda x: len(x) /
ax[1].set(ylabel="(%)")
ax[1].set_xticklabels(df["education"].unique(), rotation=0, rotation_mode="anchor")
plt.show()
```

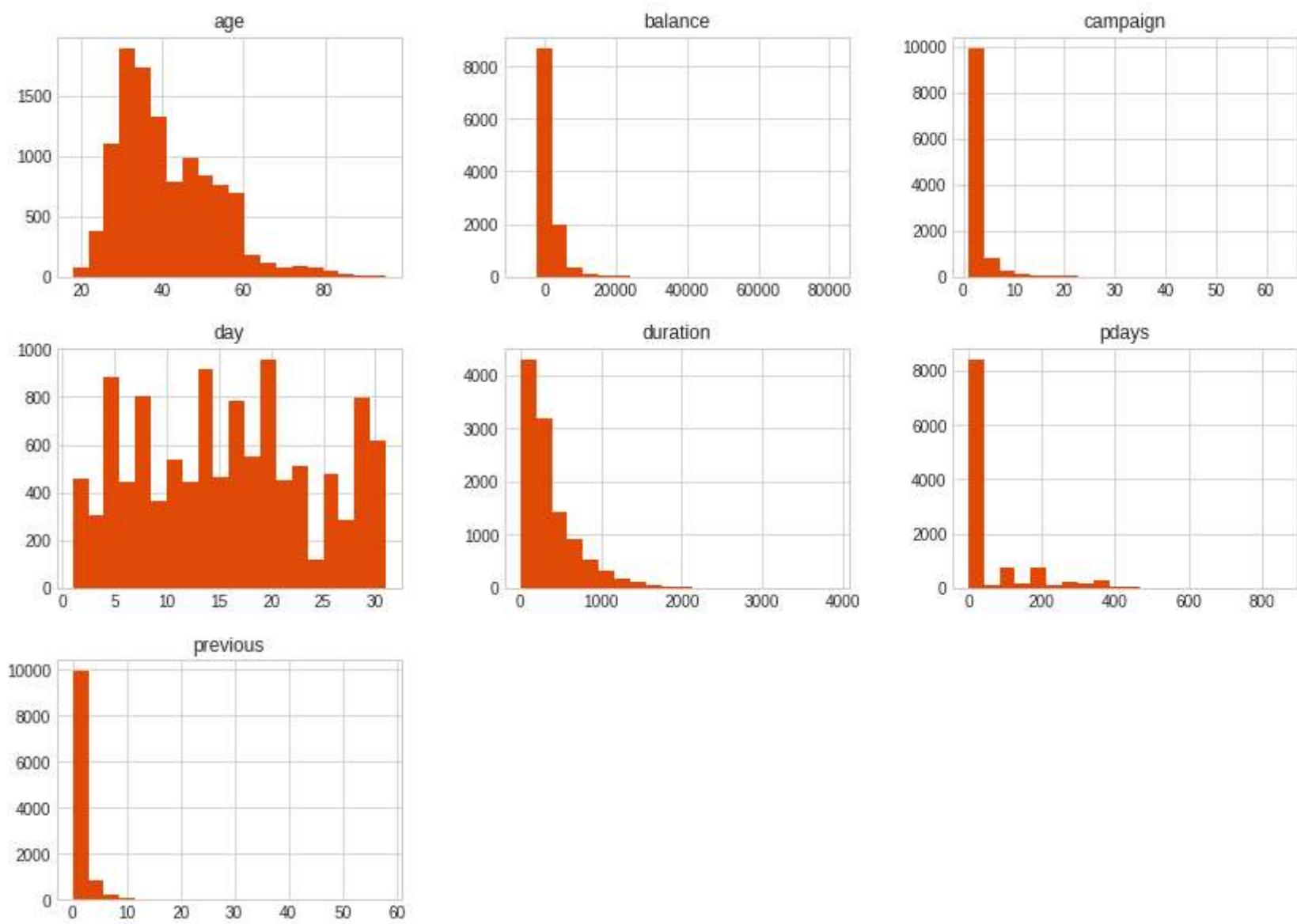
```
/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning:
```

```
Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.
```

Information on Term Suscriptions



```
In [5]: # Let's see how the numeric data is distributed.  
import matplotlib.pyplot as plt  
plt.style.use('seaborn-whitegrid')  
  
df.hist(bins=20, figsize=(14,10), color="#E14906")  
plt.show()
```



```
In [6]: df['deposit'].value_counts()
```

```
Out[6]: no      5873  
yes     5289  
Name: deposit, dtype: int64
```

```
In [7]: # plt.style.use('dark_background')
```

```
fig = plt.figure(figsize=(20,20))
ax1 = fig.add_subplot(221)
ax2 = fig.add_subplot(222)
ax3 = fig.add_subplot(212)

g = sns.boxplot(x="default", y="balance", hue="deposit",
                 data=df, palette="muted", ax=ax1)

g.set_title("Amount of Balance by Term Suscriptions")

# ax.set_xticklabels(df["default"].unique(), rotation=45, rotation_mode="anchor")

g1 = sns.boxplot(x="job", y="balance", hue="deposit",
                  data=df, palette="RdBu", ax=ax2)

g1.set_xticklabels(df["job"].unique(), rotation=90, rotation_mode="anchor")
g1.set_title("Type of Work by Term Suscriptions")

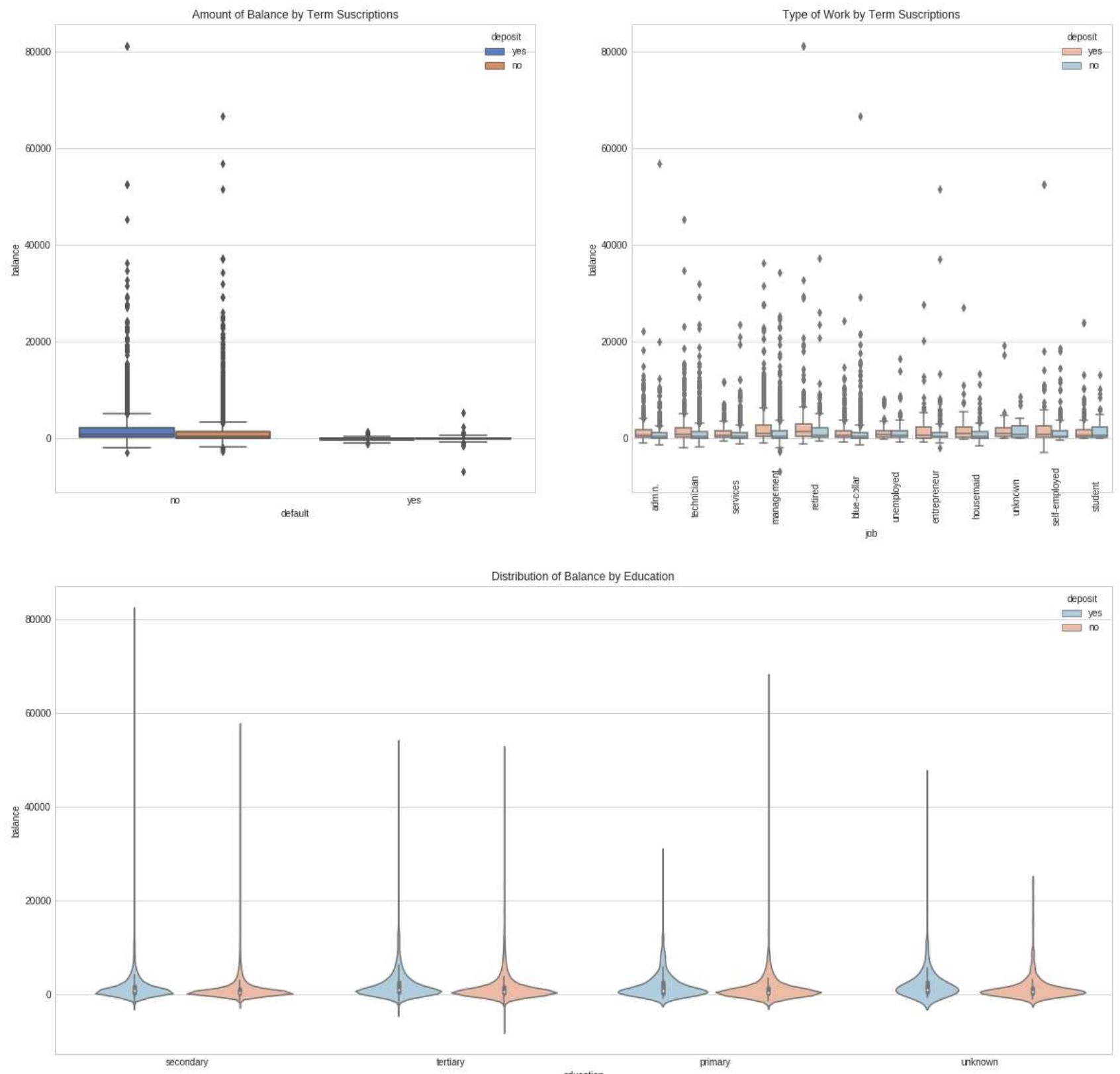
g2 = sns.violinplot(data=df, x="education", y="balance", hue="deposit", palette="RdBu_r")

g2.set_title("Distribution of Balance by Education")

plt.show()
```

```
/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning:
```

```
Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.
```



```
In [8]: df.head()
```

Out[8]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	pou
0	59	admin.	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	ur
1	56	admin.	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	ur
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	ur
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	ur
4	54	admin.	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	ur

## Analysis by Occupation:

Number of Occupations: Management is the most prevalent occupation in this dataset.

Age by Occupation: As expected, retired individuals have the highest median age, while students have the lowest.

Balance by Occupation: Management and retirees have the highest account balances.

```
In [9]: # Drop the Job Occupations that are "Unknown"  
df = df.drop(df.loc[df["job"] == "unknown"].index)
```

```
# Admin and management are basically the same let's put it under the same categorical value  
lst = [df]
```

```
for col in lst:  
    col.loc[col["job"] == "admin.", "job"] = "management"
```

```
In [10]: df.columns
```

```
Out[10]: Index(['age', 'job', 'marital', 'education', 'default', 'balance', 'housing',  
       'loan', 'contact', 'day', 'month', 'duration', 'campaign', 'pdays',  
       'previous', 'poutcome', 'deposit'],  
      dtype='object')
```

```
In [11]: import squarify
df = df.drop(df[df["balance"] == 0].index)
```

```
x = 0
y = 0
width = 100
height = 100

job_names = df['job'].value_counts().index
values = df['job'].value_counts().tolist()

normed = squarify.normalize_sizes(values, width, height)
rects = squarify.squarify(normed, x, y, width, height)

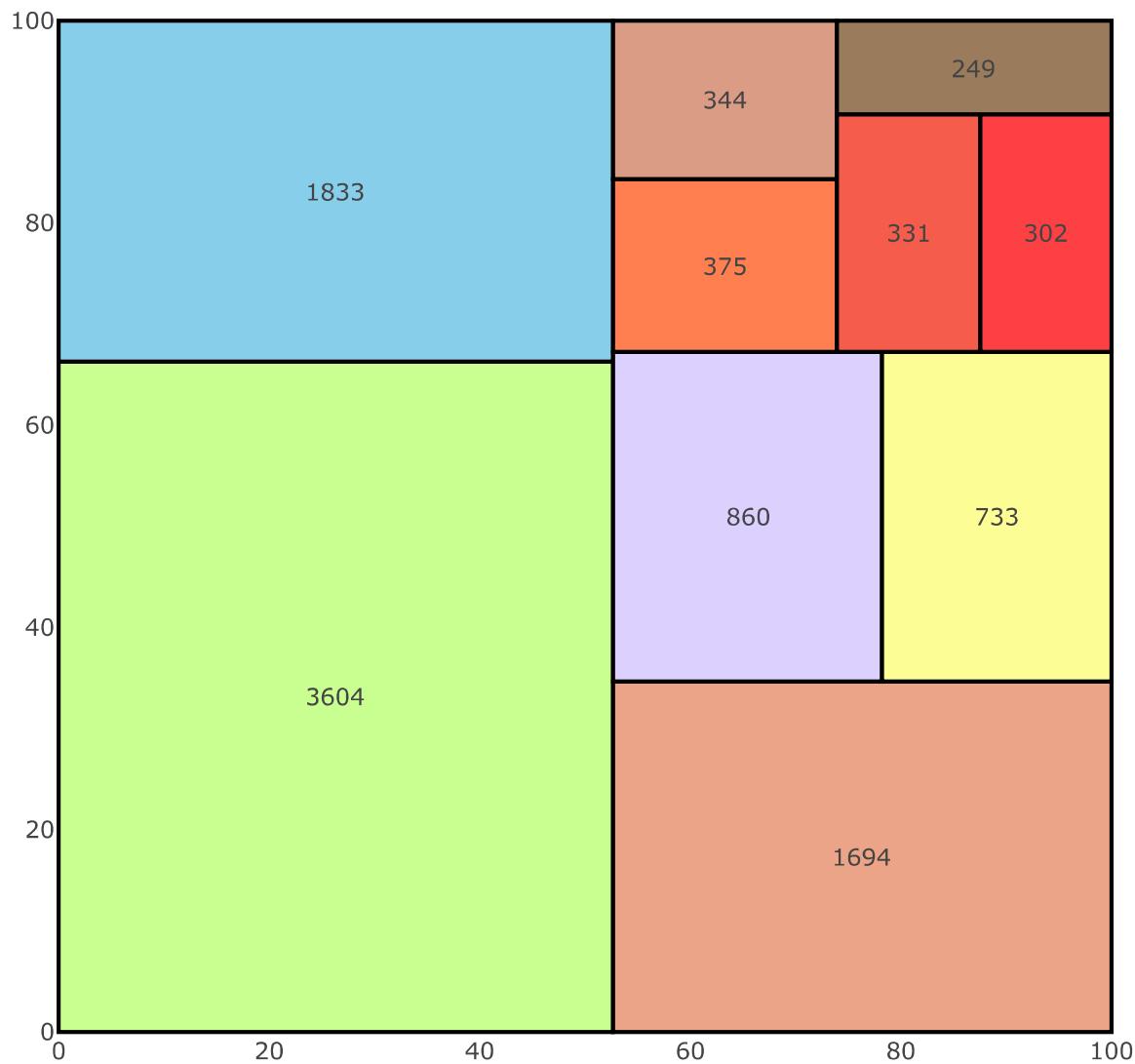
colors = ['rgb(200, 255, 144)', 'rgb(135, 206, 235)',
          'rgb(235, 164, 135)', 'rgb(220, 208, 255)',
          'rgb(253, 253, 150)', 'rgb(255, 127, 80)',
          'rgb(218, 156, 133)', 'rgb(245, 92, 76)',
          'rgb(252,64,68)', 'rgb(154,123,91)']

shapes = []
annotations = []
counter = 0

for r in rects:
    shapes.append(
        dict(
            type = 'rect',
            x0 = r['x'],
            y0 = r['y'],
            x1 = r['x'] + r['dx'],
            y1 = r['y'] + r['dy'],
            line = dict(width=2),
            fillcolor = colors[counter]
        )
    )
    annotations.append(
        dict(
            x = r['x']+(r['dx']/2),
            y = r['y']+(r['dy']/2),
            text = values[counter],
            showarrow = False
        )
    )
    counter = counter + 1
if counter >= len(colors):
    counter = 0

# For hover text
trace0 = go.Scatter(
    x = [ r['x']+(r['dx']/2) for r in rects],
    y = [ r['y']+(r['dy']/2) for r in rects],
    text = [ str(v) for v in job_names],
    mode='text',
)
layout = dict(
    title='Number of Occupations <br> <i>(From our Sample Population)</i>',
    height=700,
    width=700,
    xaxis=dict(showgrid=False,zeroline=False),
    yaxis=dict(showgrid=False,zeroline=False),
    shapes=shapes,
    annotations=annotations,
    hovermode='closest'
)
# With hovertext
figure = dict(data=[trace0], layout=layout)
iplot(figure, filename='squarify-treemap')
```

Number of Occupations  
*(From our Sample Population)*



***let's see which occupation tended to have more balance in their accounts***

```
In [12]: suscribed_df = df.loc[df["deposit"] == "yes"]

occupations = df["job"].unique().tolist()

# Get the balances by jobs
management = suscribed_df["age"].loc[suscribed_df["job"] == "management"].values
technician = suscribed_df["age"].loc[suscribed_df["job"] == "technician"].values
services = suscribed_df["age"].loc[suscribed_df["job"] == "services"].values
retired = suscribed_df["age"].loc[suscribed_df["job"] == "retired"].values
blue_collar = suscribed_df["age"].loc[suscribed_df["job"] == "blue-collar"].values
unemployed = suscribed_df["age"].loc[suscribed_df["job"] == "unemployed"].values
entrepreneur = suscribed_df["age"].loc[suscribed_df["job"] == "entrepreneur"].values
housemaid = suscribed_df["age"].loc[suscribed_df["job"] == "housemaid"].values
self_employed = suscribed_df["age"].loc[suscribed_df["job"] == "self-employed"].values
student = suscribed_df["age"].loc[suscribed_df["job"] == "student"].values

ages = [management, technician, services, retired, blue_collar, unemployed,
        entrepreneur, housemaid, self_employed, student]

colors = ['rgba(93, 164, 214, 0.5)', 'rgba(255, 144, 14, 0.5)',
          'rgba(44, 160, 101, 0.5)', 'rgba(255, 65, 54, 0.5)',
          'rgba(207, 114, 255, 0.5)', 'rgba(127, 96, 0, 0.5)',
          'rgba(229, 126, 56, 0.5)', 'rgba(229, 56, 56, 0.5)',
          'rgba(174, 229, 56, 0.5)', 'rgba(229, 56, 56, 0.5)']

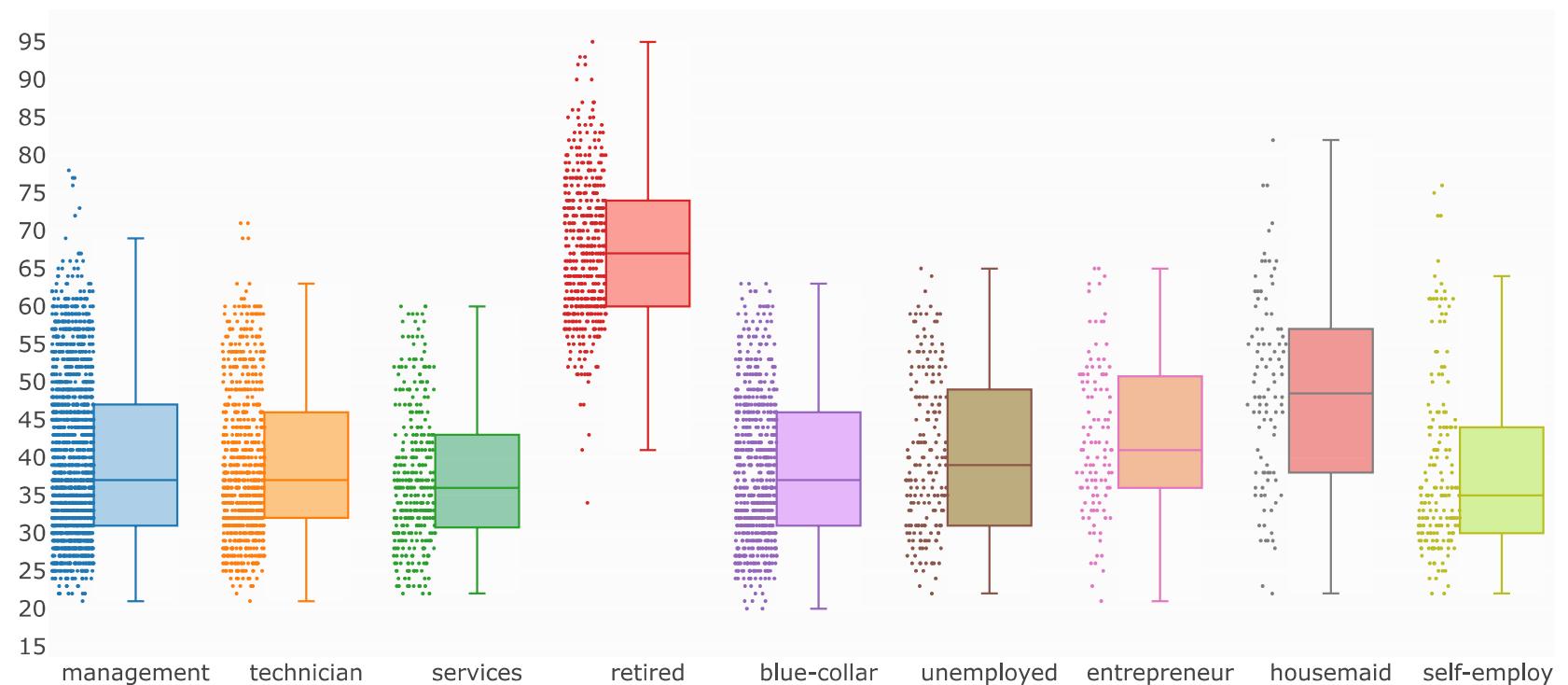
traces = []

for xd, yd, cls in zip(occupations, ages, colors):
    traces.append(go.Box(
        y=yd,
        name=xd,
        boxpoints='all',
        jitter=0.5,
        whiskerwidth=0.2,
        fillcolor=cls,
        marker=dict(
            size=2,
        ),
        line=dict(width=1),
    ))

layout = go.Layout(
    title='Distribution of Ages by Occupation',
    yaxis=dict(
        autorange=True,
        showgrid=True,
        zeroline=True,
        dtick=5,
        gridcolor='rgb(255, 255, 255)',
        gridwidth=1,
        zerolinecolor='rgb(255, 255, 255)',
        zerolinewidth=2,
    ),
    margin=dict(
        l=40,
        r=30,
        b=80,
        t=100,
    ),
    paper_bgcolor='rgb(224,255,246)',
    plot_bgcolor='rgb(251,251,251)',
    showlegend=False
)

fig = go.Figure(data=traces, layout=layout)
iplot(fig)
```

Distribution of Ages by Occupation



```
In [13]: # Balance Distribution

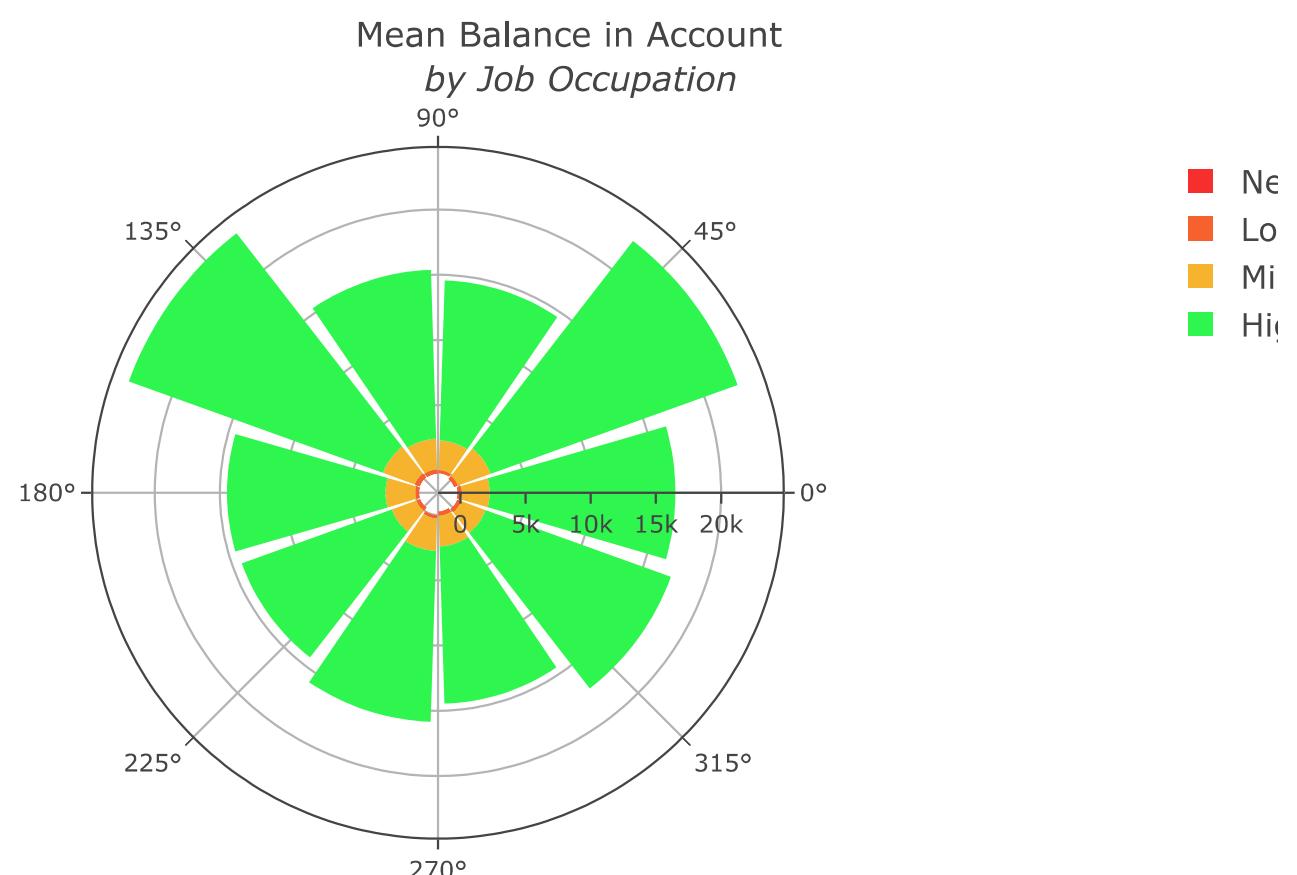
# Create a Balance Category
df["balance_status"] = np.nan
lst = [df]

for col in lst:
    col.loc[col["balance"] < 0, "balance_status"] = "negative"
    col.loc[(col["balance"] >= 0) & (col["balance"] <= 30000), "balance_status"] = "low"
    col.loc[(col["balance"] > 30000) & (col["balance"] <= 40000), "balance_status"] = "middle"
    col.loc[col["balance"] > 40000, "balance_status"] = "high"

# balance by balance_status
negative = df["balance"].loc[df["balance_status"] == "negative"].values.tolist()
low = df["balance"].loc[df["balance_status"] == "low"].values.tolist()
middle = df["balance"].loc[df["balance_status"] == "middle"].values.tolist()
high = df["balance"].loc[df["balance_status"] == "high"].values.tolist()

# Get the average by occupation in each balance category
job_balance = df.groupby(['job', 'balance_status'])['balance'].mean()

trace1 = go.Barpolar(
    r=[-199.0, -392.0, -209.0, -247.0, -233.0, -270.0, -271.0, 0, -276.0, -134.5],
    text=["blue-collar", "entrepreneur", "housemaid", "management", "retired", "self-employed",
          "services", "student", "technician", "unemployed"],
    name='Negative Balance',
    marker=dict(
        color='rgb(246, 46, 46)'
    )
)
trace2 = go.Barpolar(
    r=[319.5, 283.0, 212.0, 313.0, 409.0, 274.5, 308.5, 253.0, 316.0, 330.0],
    text=["blue-collar", "entrepreneur", "housemaid", "management", "retired", "self-employed",
          "services", "student", "technician", "unemployed"],
    name='Low Balance',
    marker=dict(
        color='rgb(246, 97, 46)'
    )
)
trace3 = go.Barpolar(
    r=[2128.5, 2686.0, 2290.0, 2366.0, 2579.0, 2293.5, 2005.5, 2488.0, 2362.0, 1976.0],
    text=["blue-collar", "entrepreneur", "housemaid", "management", "retired", "self-employed",
          "services", "student", "technician", "unemployed"],
    name='Middle Balance',
    marker=dict(
        color='rgb(246, 179, 46)'
    )
)
trace4 = go.Barpolar(
    r=[14247.5, 20138.5, 12278.5, 12956.0, 20723.0, 12159.0, 12223.0, 13107.0, 12063.0, 15107.5],
    text=["blue-collar", "entrepreneur", "housemaid", "management", "retired", "self-employed",
          "services", "student", "technician", "unemployed"],
    name='High Balance',
    marker=dict(
        color='rgb(46, 246, 78)'
    )
)
data = [trace1, trace2, trace3, trace4]
layout = go.Layout(
    title='Mean Balance in Account<br> <i> by Job Occupation</i>',
    font=dict(
        size=12
    ),
    legend=dict(
        font=dict(
            size=16
        )
    ),
    radialaxis=dict(
        ticksuffix='%'
    ),
    orientation=270
)
fig = go.Figure(data=data, layout=layout)
iplot(fig, filename='polar-area-chart')
```



## Marital Status

In this analysis, we didn't find any significant insights, except that most divorced individuals seem to have lower financial resources. No wonder, considering they have to split financial assets! Nevertheless, since no further insights have been found, we will proceed to clustering marital status with education status. Let's see if we can identify other groups of people in the sample population.

```
In [14]: df['marital'].value_counts()
```

```
Out[14]: married    5815
          single     3336
          divorced   1174
          Name: marital, dtype: int64
```

```
In [15]: df['marital'].unique()
```

```
Out[15]: array(['married', 'divorced', 'single'], dtype=object)
```

```
In [16]: df['marital'].value_counts().tolist()
```

```
Out[16]: [5815, 3336, 1174]
```

```
In [17]: vals = df['marital'].value_counts().tolist()
labels = ['married', 'divorced', 'single']

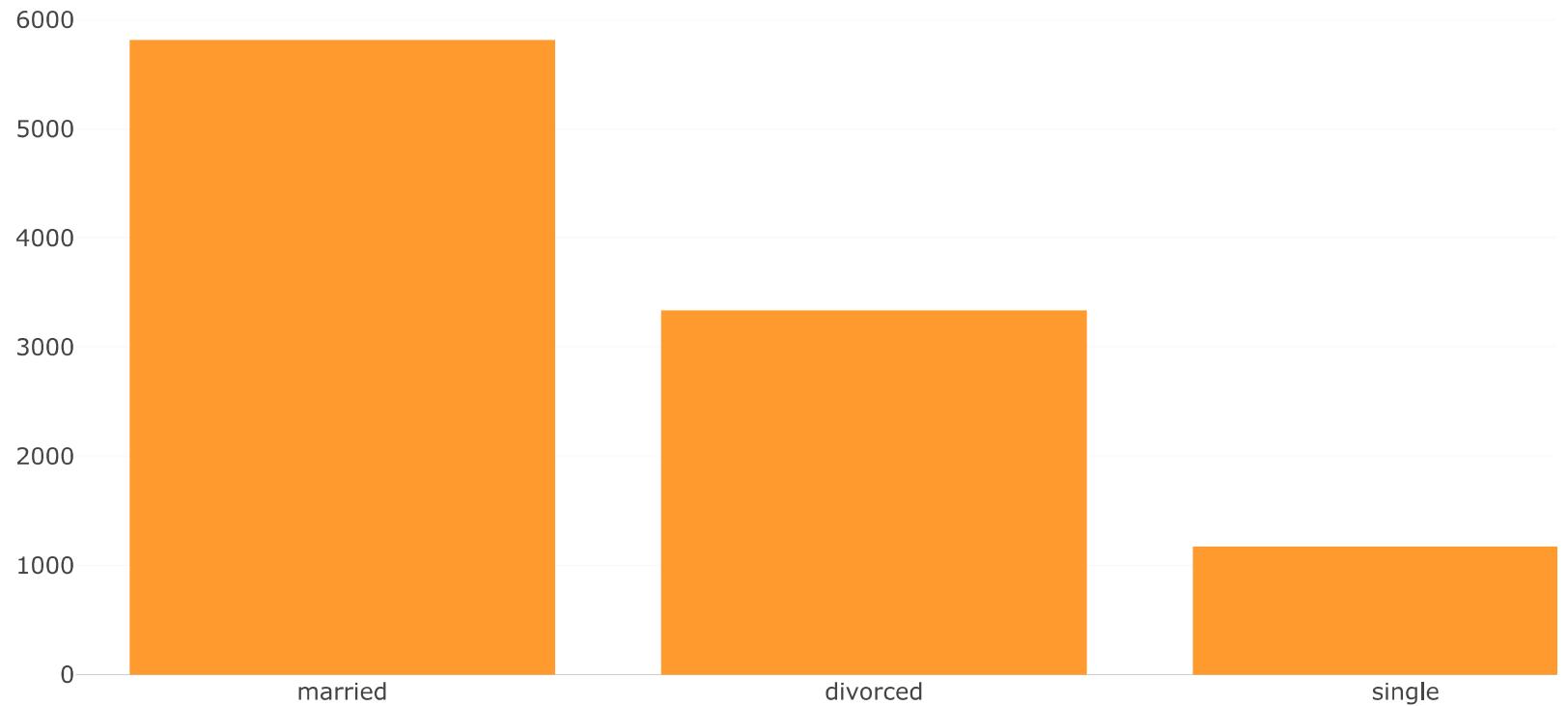
data = [go.Bar(
    x=labels,
    y=vals,
    marker=dict(
        color="#FE9A2E"
    )
)]

layout = go.Layout(
    title="Count by Marital Status",
)

fig = go.Figure(data=data, layout=layout)

iplot(fig, filename='basic-bar')
```

Count by Marital Status



```
In [18]: # Distribution of Balances by Marital status
single = df['balance'].loc[df['marital'] == 'single'].values
married = df['balance'].loc[df['marital'] == 'married'].values
divorced = df['balance'].loc[df['marital'] == 'divorced'].values

single_dist = go.Histogram(
    x=single,
    histnorm='density',
    name='single',
    marker=dict(
        color='#6E6E6E'
    )
)

married_dist = go.Histogram(
    x=married,
    histnorm='density',
    name='married',
    marker=dict(
        color='#2E9AFE'
    )
)

divorced_dist = go.Histogram(
    x=divorced,
    histnorm='density',
    name='divorced',
    marker=dict(
        color='#FA5858'
    )
)

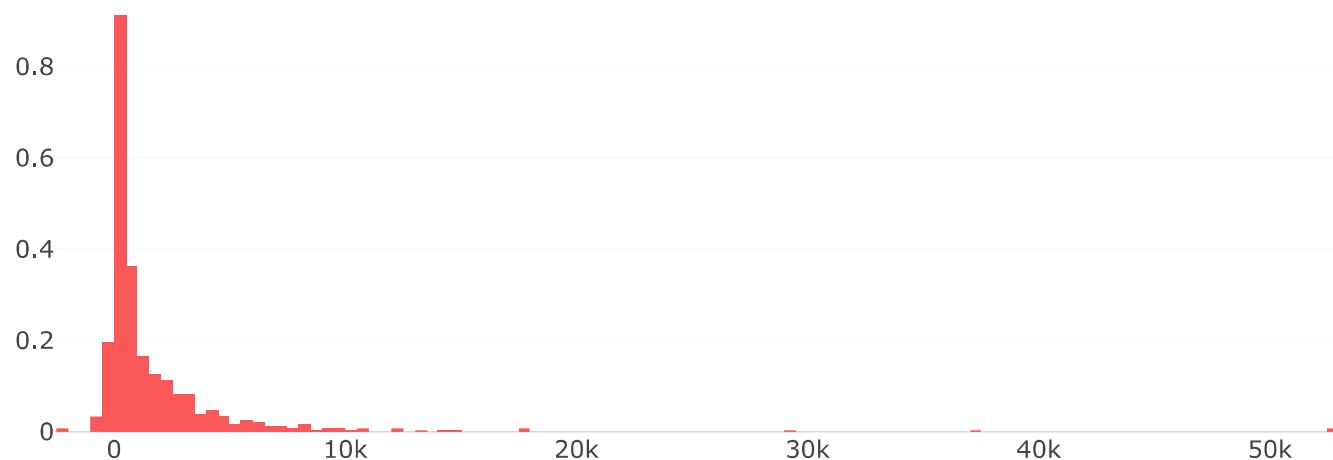
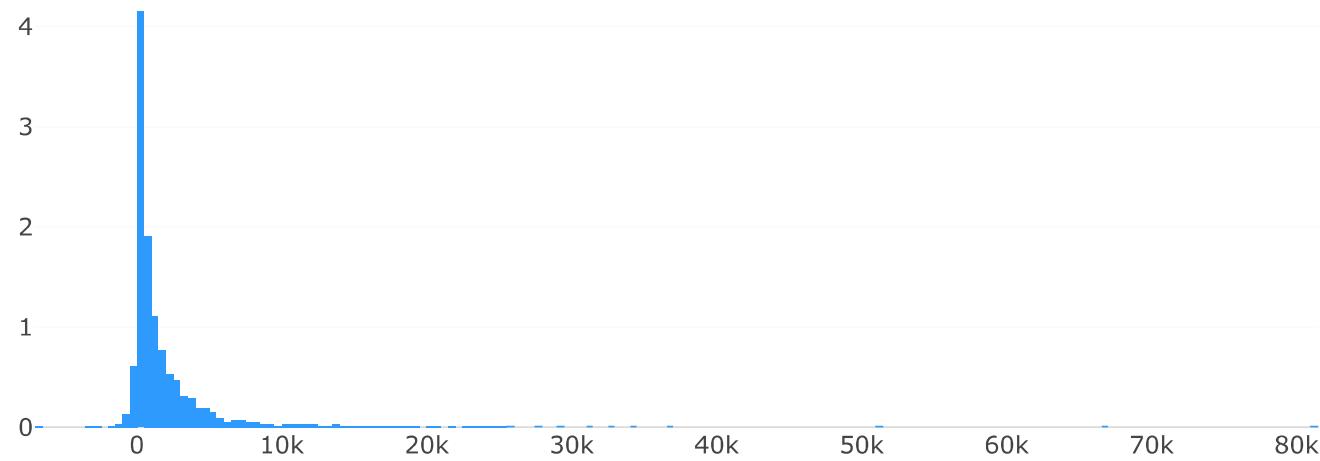
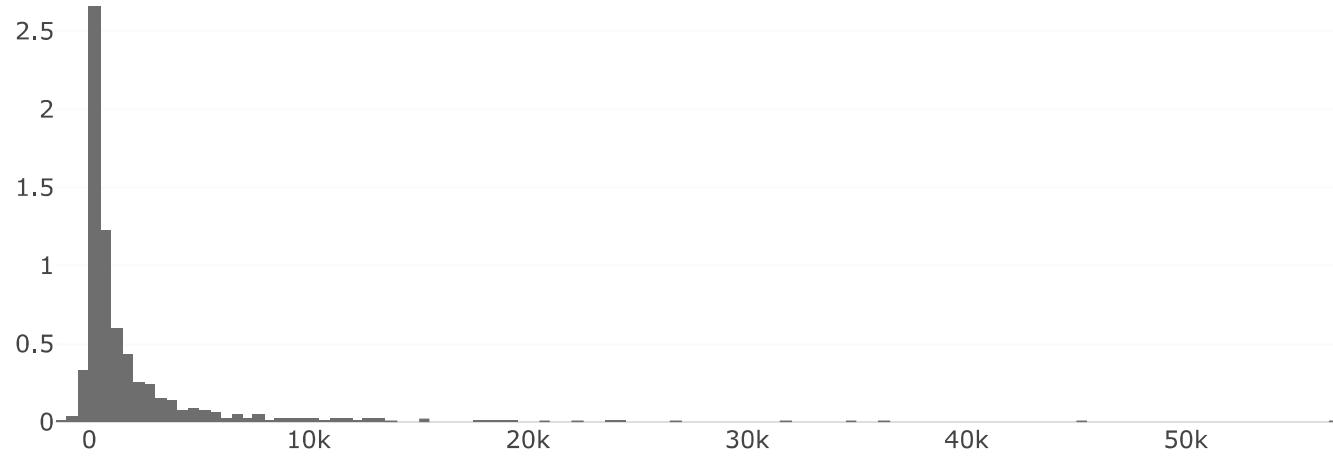
fig = tools.make_subplots(rows=3, print_grid=False)

fig.append_trace(single_dist, 1, 1)
fig.append_trace(married_dist, 2, 1)
fig.append_trace(divorced_dist, 3, 1)

fig['layout'].update(showlegend=False, title="Price Distributions by Marital Status",
                     height=1000, width=800)

iplot(fig, filename='custom-sized-subplot-with-subplot-titles')
```

### Price Distributions by Marital Status



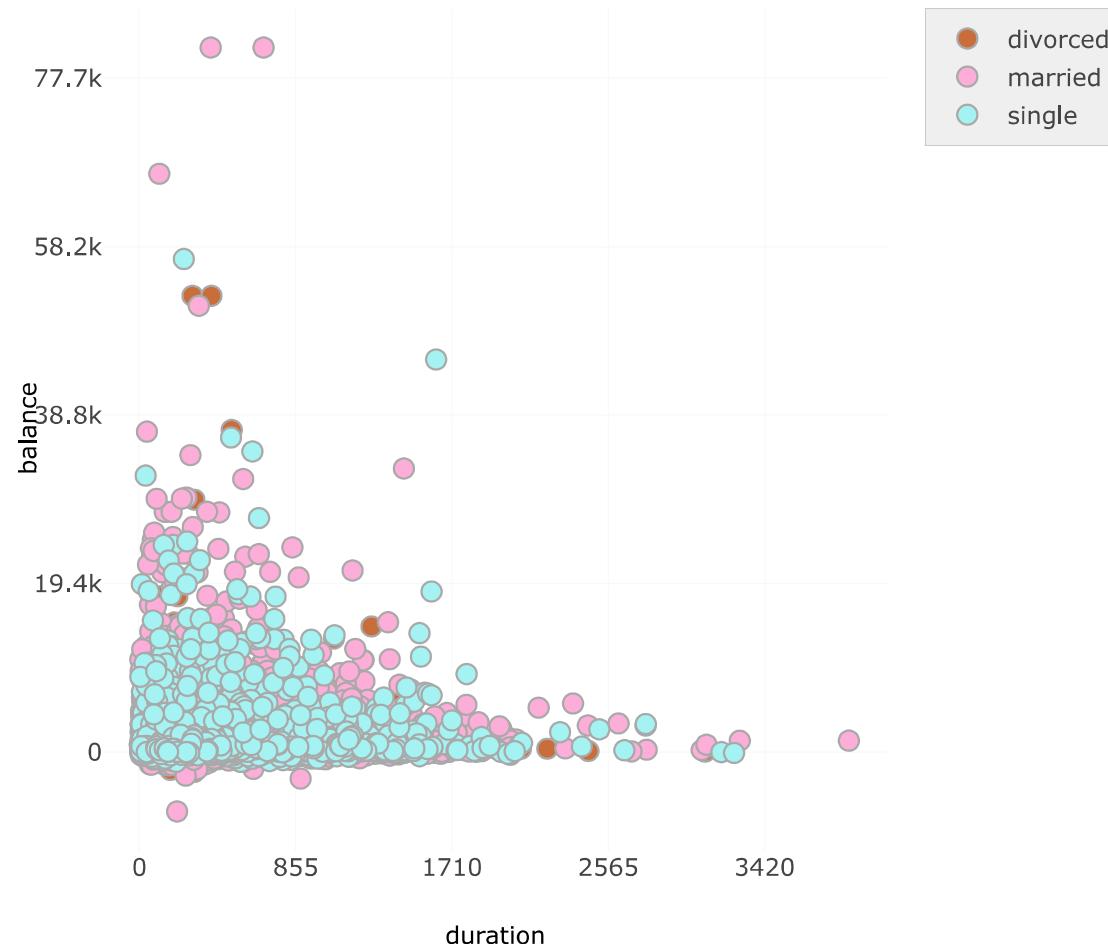
In [19]: df.head()

Out[19]:

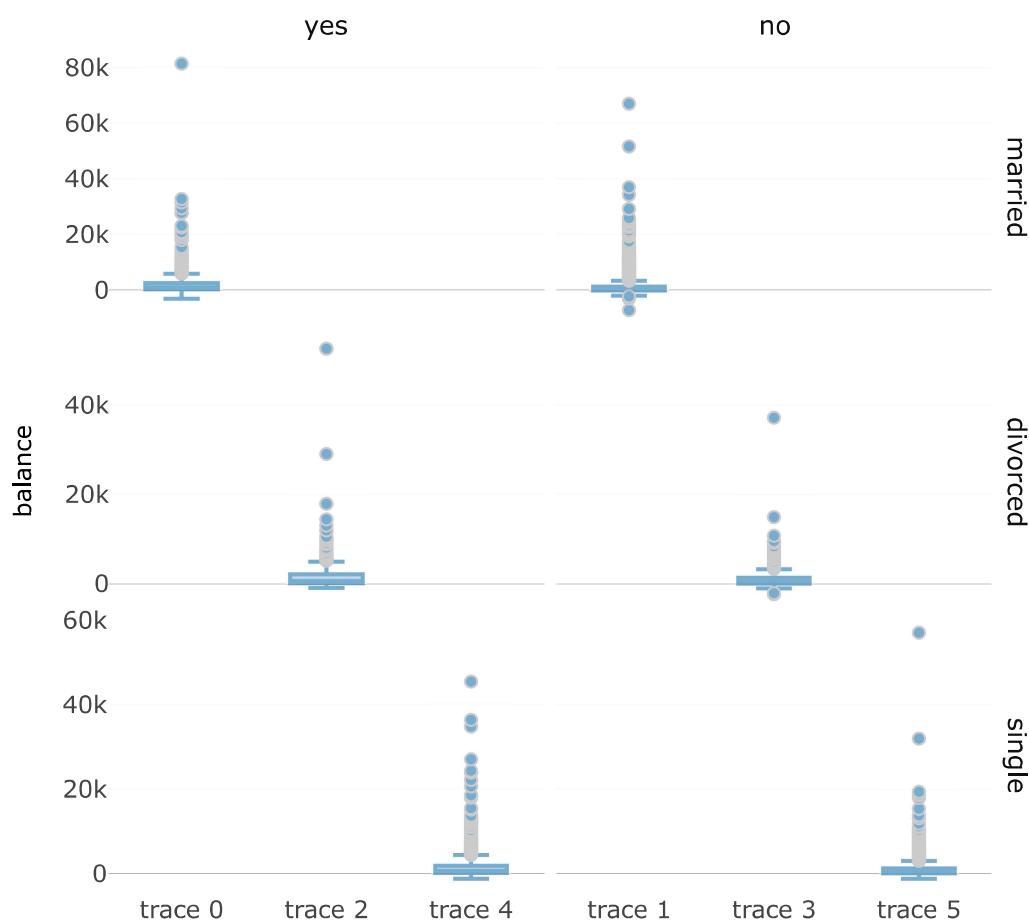
	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	f
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	

```
In [20]: # Notice how divorced have a considerably low amount of balance.
```

```
fig = ff.create_facet_grid(
    df,
    x='duration',
    y='balance',
    color_name='marital',
    show_boxes=False,
    marker={'size': 10, 'opacity': 1.0},
    colormap={'single': 'rgb(165, 242, 242)', 'married': 'rgb(253, 174, 216)', 'divorced': 'rgba(201, 109, 59, 0.5)'})
iplot(fig, filename='facet - custom colormap')
```



```
In [21]: # We have missed some important clients with some high balances.
# This shouldn't be happening.
fig = ff.create_facet_grid(
    df,
    y='balance',
    facet_row='marital',
    facet_col='deposit',
    trace_type='box',
)
iplot(fig, filename='facet - box traces')
```



```
In [22]: df.head()
```

Out[22]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	t
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	

## Clustering Marital Status and Education:

Marital Status: As discussed previously, divorce significantly affects an individual's balance.

Education: The level of education also significantly influences the balance of a prospect.

Loans: The presence of a previous loan significantly impacts the balance of the prospect.

```
In [23]:
```

```
df = df.drop(df.loc[df["education"] == "unknown"].index)
df['education'].unique()
```

Out[23]: array(['secondary', 'tertiary', 'primary'], dtype=object)

```
In [24]: df['marital/education'] = np.nan
lst = [df]

for col in lst:
    col.loc[(col['marital'] == 'single') & (df['education'] == 'primary'), 'marital/education'] = 'single/primary'
    col.loc[(col['marital'] == 'married') & (df['education'] == 'primary'), 'marital/education'] = 'married/primary'
    col.loc[(col['marital'] == 'divorced') & (df['education'] == 'primary'), 'marital/education'] = 'divorced/primary'
    col.loc[(col['marital'] == 'single') & (df['education'] == 'secondary'), 'marital/education'] = 'single/secondary'
    col.loc[(col['marital'] == 'married') & (df['education'] == 'secondary'), 'marital/education'] = 'married/secondary'
    col.loc[(col['marital'] == 'divorced') & (df['education'] == 'secondary'), 'marital/education'] = 'divorced/secondary'
    col.loc[(col['marital'] == 'single') & (df['education'] == 'tertiary'), 'marital/education'] = 'single/tertiary'
    col.loc[(col['marital'] == 'married') & (df['education'] == 'tertiary'), 'marital/education'] = 'married/tertiary'
    col.loc[(col['marital'] == 'divorced') & (df['education'] == 'tertiary'), 'marital/education'] = 'divorced/tertiary'

df.head()
```

Out[24]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	f
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	0
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	0
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	0
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	0
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	0

```
In [25]: pal = sns.cubehelix_palette(10, rot=-.25, light=.7)
g = sns.FacetGrid(df, row="marital/education", hue="marital/education", aspect=12, palette=pal)

g.map(sns.kdeplot, "balance", clip_on=False, shade=True, alpha=1, lw=1.5, bw=.2)
g.map(sns.kdeplot, "balance", clip_on=False, color="w", lw=1, bw=0)
g.map(plt.axhline, y=0, lw=2, clip_on=False)
```

Out[25]: <seaborn.axisgrid.FacetGrid at 0x793f15750b38>



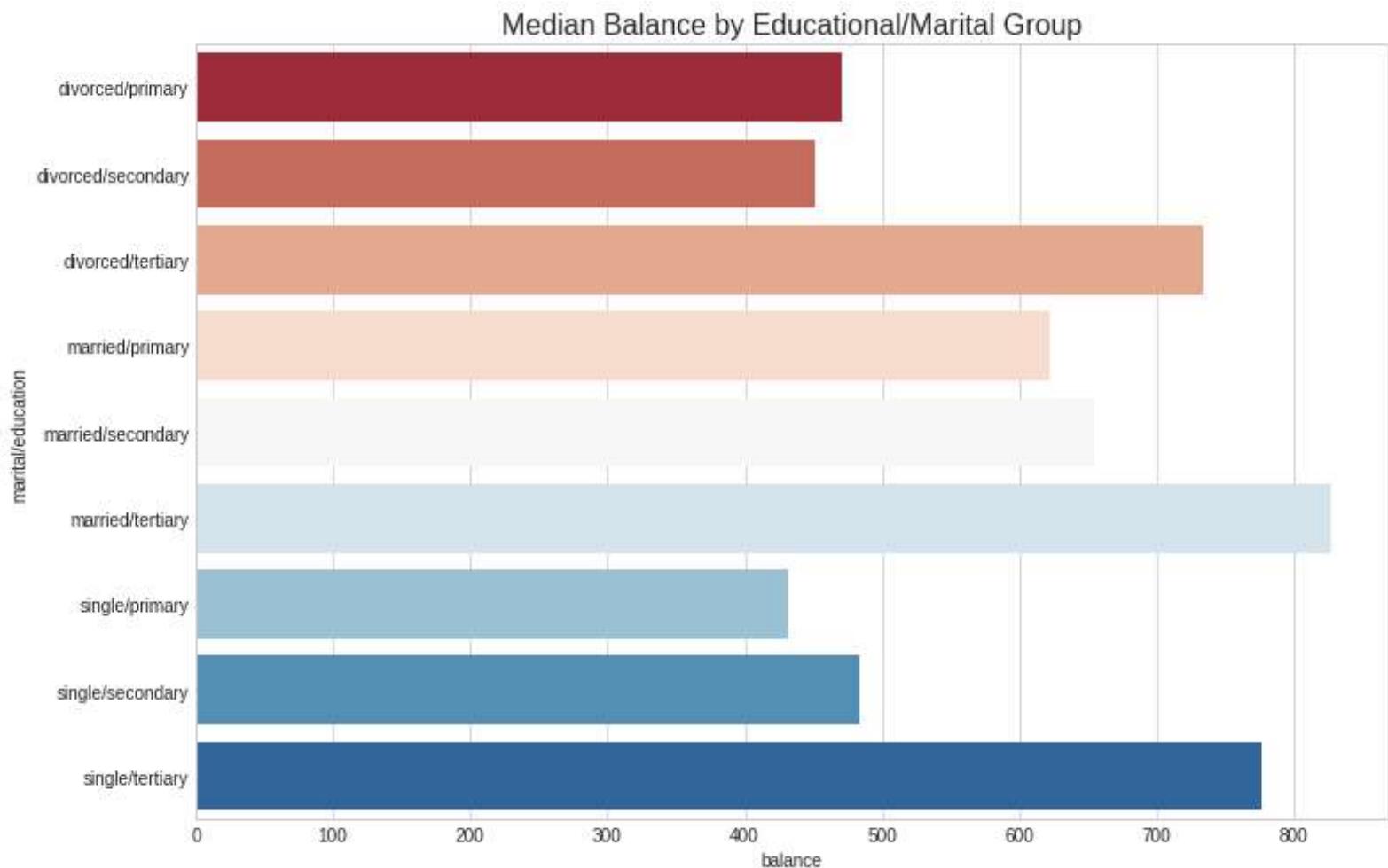
```
In [26]: education_groups = df.groupby(['marital/education'], as_index=False)[['balance']].median()

fig = plt.figure(figsize=(12,8))

sns.barplot(x="balance", y="marital/education", data=education_groups,
            label="Total", palette="RdBu")

plt.title('Median Balance by Educational/Marital Group', fontsize=16)
```

Out[26]: Text(0.5,1,'Median Balance by Educational/Marital Group')



```
In [27]: # Let's see the group who had Loans from the marital/education group
```

```
loan_balance = df.groupby(['marital/education', 'loan'], as_index=False)[['balance']].median()

no_loan = loan_balance[['balance']].loc[loan_balance['loan'] == 'no'].values
has_loan = loan_balance[['balance']].loc[loan_balance['loan'] == 'yes'].values

labels = loan_balance['marital/education'].unique().tolist()

trace0 = go.Scatter(
    x=no_loan,
    y=labels,
    mode='markers',
    name='No Loan',
    marker=dict(
        color='rgb(175,238,238)',
        line=dict(
            color='rgb(0,139,139)',
            width=1,
        ),
        symbol='circle',
        size=16,
    )
)
trace1 = go.Scatter(
    x=has_loan,
    y=labels,
    mode='markers',
    name='Has a Previous Loan',
    marker=dict(
        color='rgb(250,128,114)',
        line=dict(
            color='rgb(178,34,34)',
            width=1,
        ),
        symbol='circle',
        size=16,
    )
)
data = [trace0, trace1]
layout = go.Layout(
    title="The Impact of Loans to Married/Educational Clusters",
    xaxis=dict(
        showgrid=False,
        showline=True,
        linecolor='rgb(102, 102, 102)',
        titlefont=dict(
            color='rgb(204, 204, 204)'
        ),
        tickfont=dict(
            color='rgb(102, 102, 102)'
        ),
        showticklabels=False,
        dtick=10,
        ticks='outside',
        tickcolor='rgb(102, 102, 102)'
    ),
    margin=dict(
        l=140,
        r=40,
        b=50,
        t=80
    ),
    legend=dict(
        font=dict(
            size=10,
        ),
        yanchor='middle',
        xanchor='right',
    ),
    width=1000,
    height=800,
    paper_bgcolor='rgb(255,250,250)',
    plot_bgcolor='rgb(255,255,255)',
    hovermode='closest',
)
fig = go.Figure(data=data, layout=layout)
iplot(fig, filename='lowest-oecd-votes-cast')
```

## The Impact of Loans to Married/Educational Clusters



In [28]: `df.head()`

Out[28]:

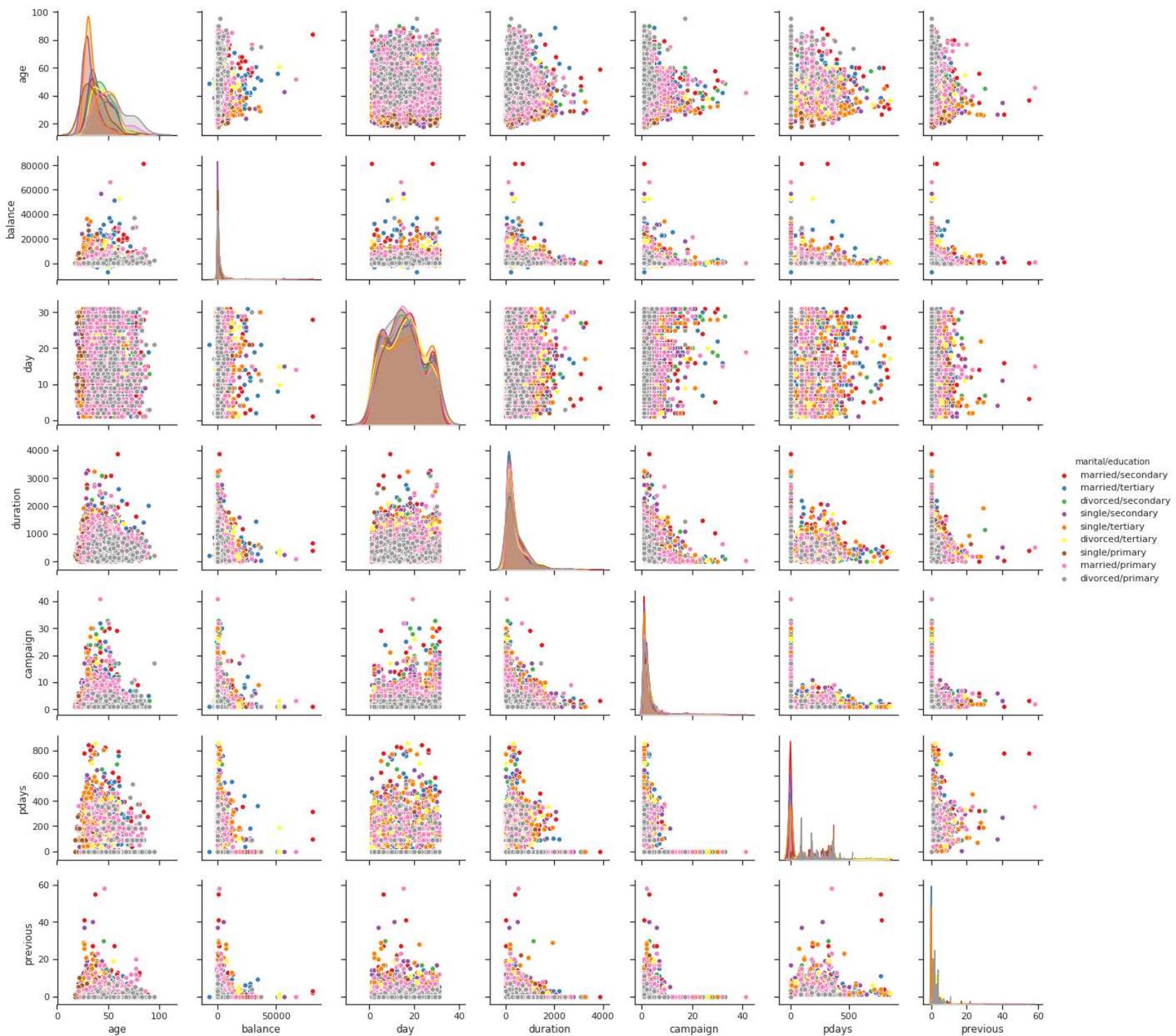
	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	target
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	no
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	no
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	no
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	no
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	no

```
In [29]: import seaborn as sns
sns.set(style="ticks")

sns.pairplot(df, hue="marital/education", palette="Set1")
plt.show()
```

/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning:

Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.



```
In [30]: df.head()
```

Out[30]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	t
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	

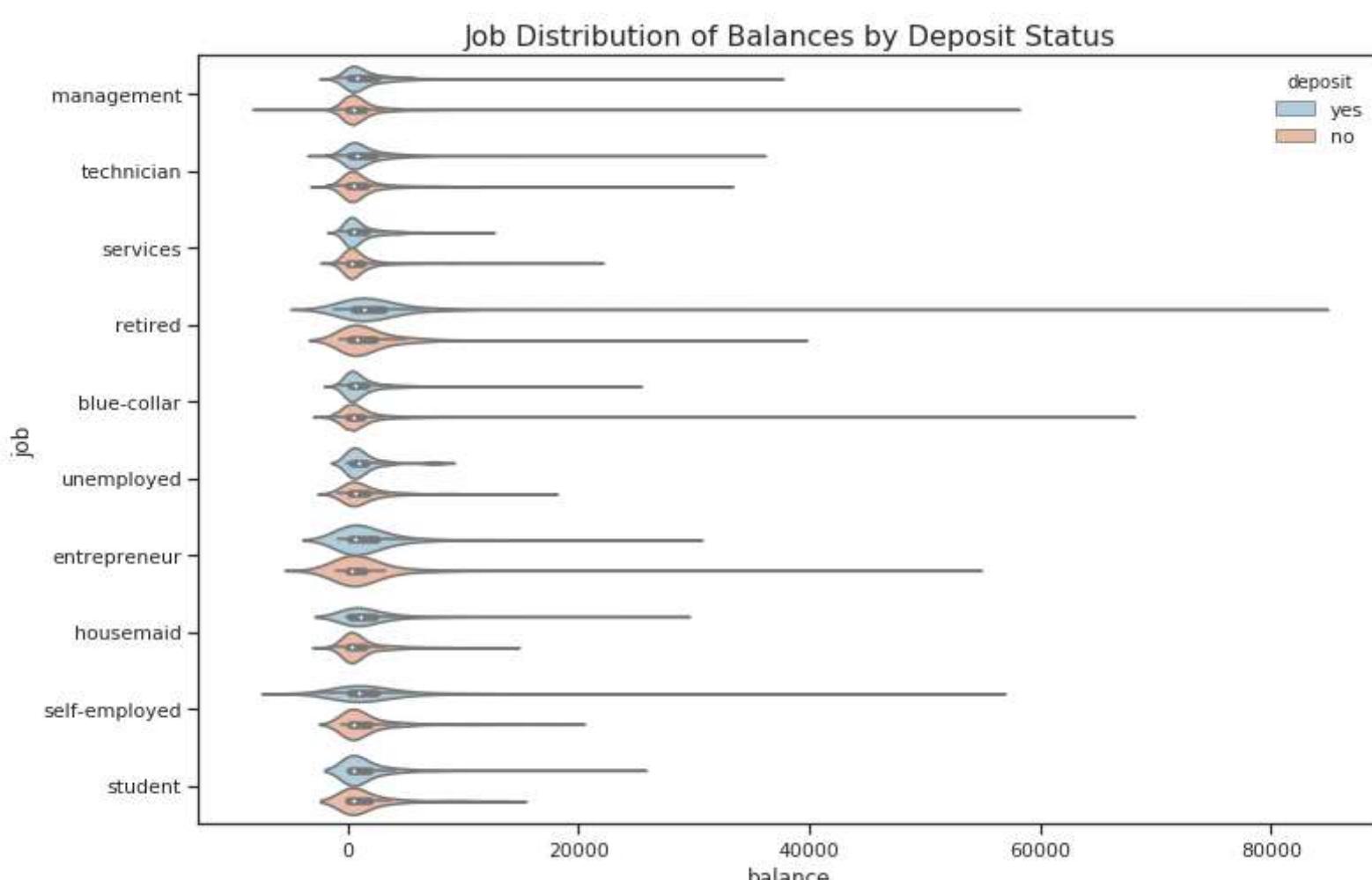
```
In [31]: fig = plt.figure(figsize=(12,8))

sns.violinplot(x="balance", y="job", hue="deposit", palette="RdBu_r",
                 data=df);

plt.title("Job Distribution of Balances by Deposit Status", fontsize=16)
plt.show()
```

/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning:

Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.



## Campaign Duration:

Campaign Duration: we observe a high correlation between duration and term deposits, indicating that the higher the duration, the more likely a client is to open a term deposit.

Average Campaign Duration: The average campaign duration is 374.76. Let's explore whether clients above this average were more likely to open a term deposit.

Duration Status: Individuals above the average duration were more likely to open a term deposit. Specifically, 78% of the group above average in duration opened term deposits, while only 32% of those below average opened term deposit accounts. This suggests that targeting individuals in the above-average category would be a promising strategy.

```
In [32]: df.drop(['marital/education', 'balance_status'], axis=1, inplace=True)
```

```
In [33]: df.head()
```

Out[33]:

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	t
0	59	management	married	secondary	no	2343	yes	no	unknown	5	may	1042	1	-1	0	
1	56	management	married	secondary	no	45	no	no	unknown	5	may	1467	1	-1	0	
2	41	technician	married	secondary	no	1270	yes	no	unknown	5	may	1389	1	-1	0	
3	55	services	married	secondary	no	2476	yes	no	unknown	5	may	579	1	-1	0	
4	54	management	married	tertiary	no	184	no	no	unknown	5	may	673	2	-1	0	

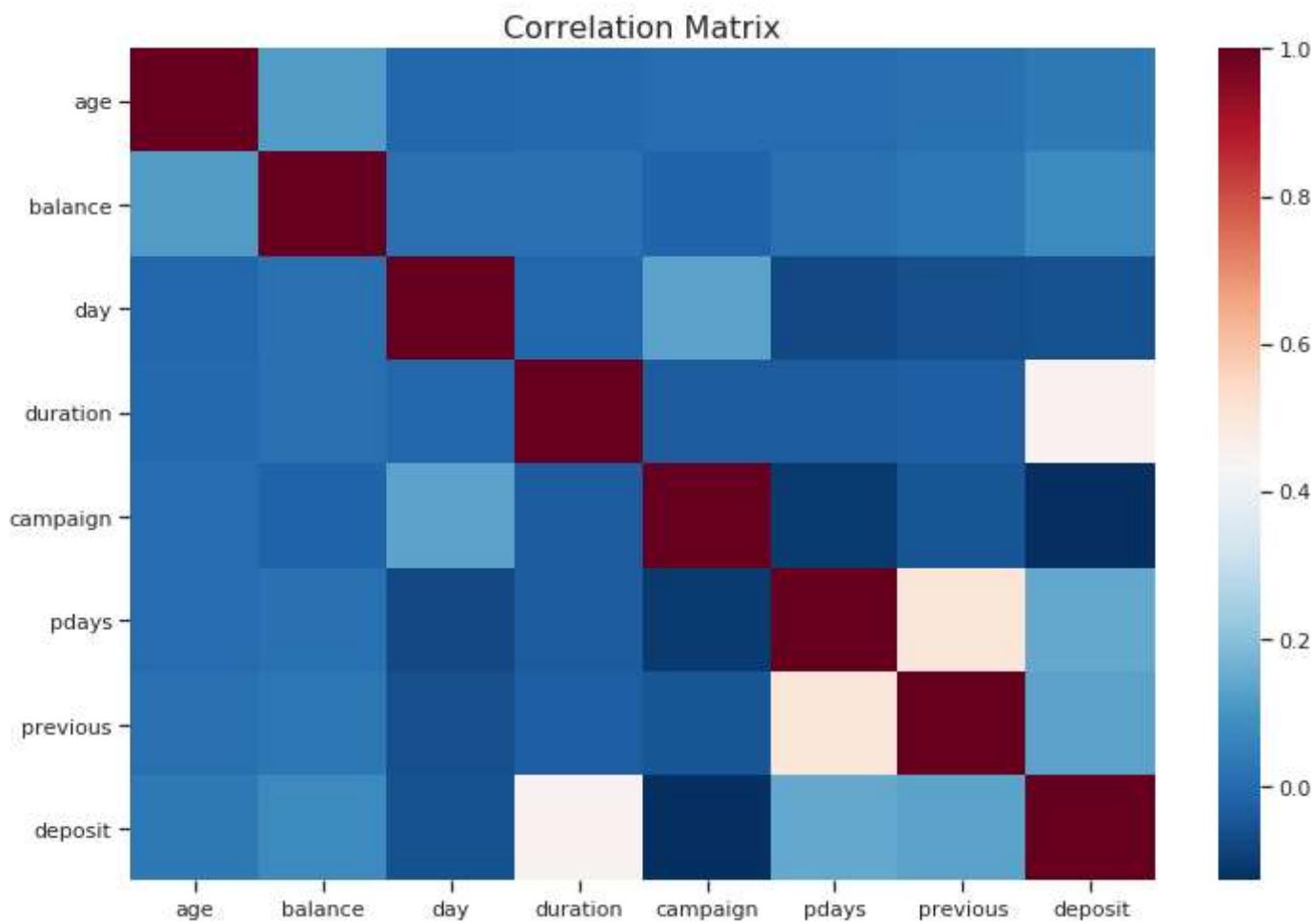
```
In [34]: # Let's drop marital/education and balance status
# Let's scale both numeric and categorical values
# Then let's use a correlation matrix
# With that we can determine if duration has influence on term deposits

from sklearn.preprocessing import StandardScaler, OneHotEncoder, LabelEncoder
fig = plt.figure(figsize=(12,8))
df['deposit'] = LabelEncoder().fit_transform(df['deposit'])

# Separate both dataframes into
numeric_df = df.select_dtypes(exclude="object")
# categorical_df = df.select_dtypes(include="object")

corr_numeric = numeric_df.corr()

sns.heatmap(corr_numeric, cbar=True, cmap="RdBu_r")
plt.title("Correlation Matrix", fontsize=16)
plt.show()
```



```
In [35]: sns.set(rc={'figure.figsize':(11.7,8.27)})
sns.set_style('whitegrid')
avg_duration = df['duration'].mean()

lst = [df]
df["duration_status"] = np.nan

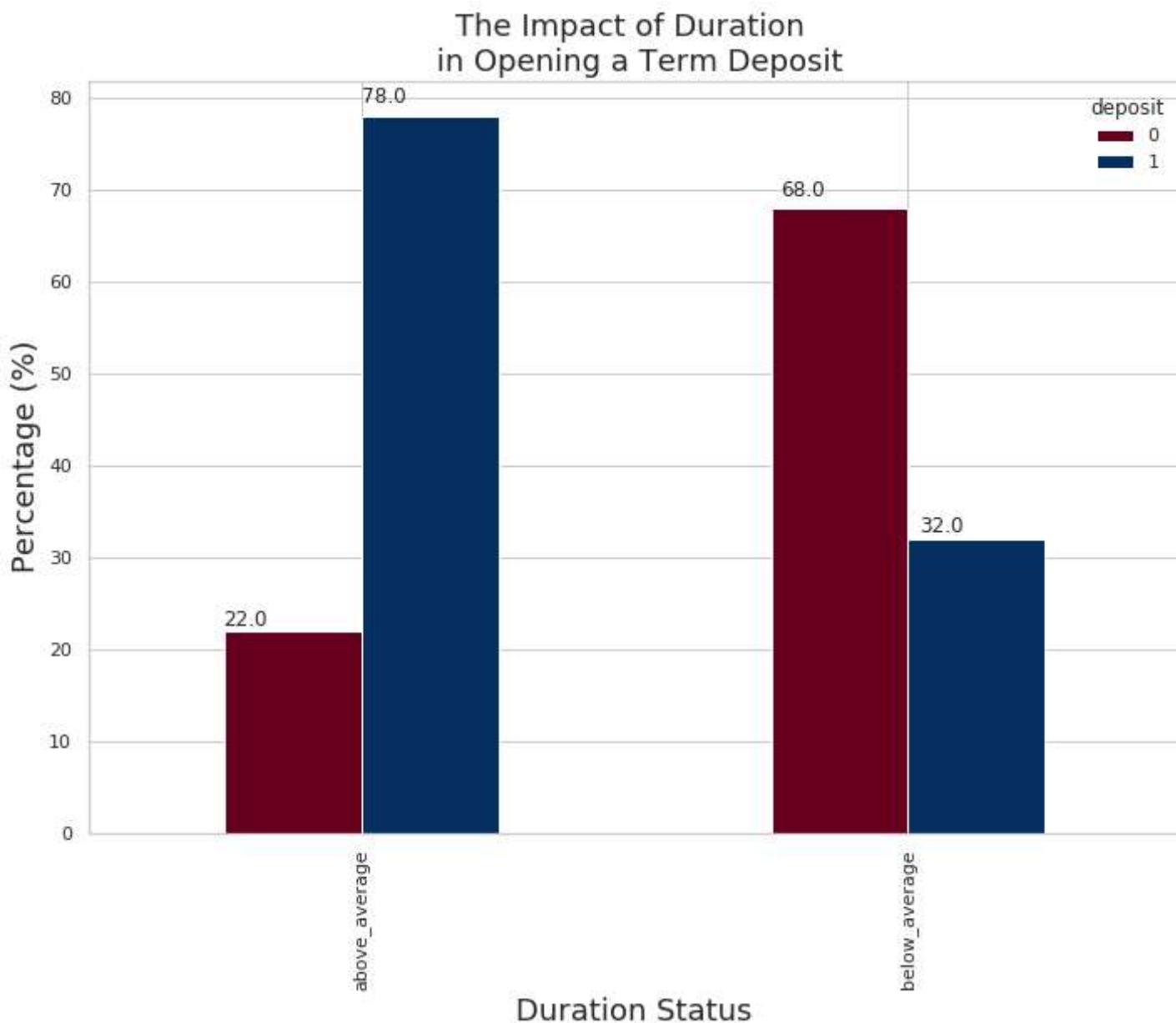
for col in lst:
    col.loc[col["duration"] < avg_duration, "duration_status"] = "below_average"
    col.loc[col["duration"] > avg_duration, "duration_status"] = "above_average"

pct_term = pd.crosstab(df['duration_status'], df['deposit']).apply(lambda r: round(r/r.sum(), 2) * 100, axis=1)

ax = pct_term.plot(kind='bar', stacked=False, cmap='RdBu')
plt.title("The Impact of Duration \n in Opening a Term Deposit", fontsize=18)
plt.xlabel("Duration Status", fontsize=18);
plt.ylabel("Percentage (%)", fontsize=18)

for p in ax.patches:
    ax.annotate(str(p.get_height()), (p.get_x() * 1.02, p.get_height() * 1.02))

plt.show()
```



## Classification Model:

```
In [36]: dep = term_deposits['deposit']
term_deposits.drop(labels=['deposit'], axis=1, inplace=True)
term_deposits.insert(0, 'deposit', dep)
term_deposits.head()
# housing has a -20% correlation with deposit let's see how it is distributed.
# 52 %
term_deposits["housing"].value_counts()/len(term_deposits)
```

Out[36]: no 0.526877  
yes 0.473123  
Name: housing, dtype: float64

```
In [37]: term_deposits["loan"].value_counts()/len(term_deposits)
```

Out[37]: no 0.869199  
yes 0.130801  
Name: loan, dtype: float64

## Stratified Sampling:

Stratified sampling is an important concept often overlooked when developing a model, whether for regression or classification. Remember that, to avoid overfitting our data, we must implement cross-validation. However, it's crucial to ensure that features with the greatest influence on our label (whether a potential client will open a term deposit or not) are equally distributed. What do I mean by this?

Personal Loans:

For instance, having a personal loan is a crucial feature in determining whether a potential client will open a term deposit. To confirm its substantial impact on the final output, you can check the correlation matrix above, where you'll see it has an -11% correlation with opening a deposit. What steps should we take before implementing stratified sampling in our train and test data?

We need to examine how our data is distributed.

After noticing that the column of loans contains 87% "no" (does not have personal loans) and 13% "yes" (have personal loans), We want to ensure that our training and test set maintains the same ratio of 87% "no" and 13% "yes."

```
In [38]: from sklearn.model_selection import StratifiedShuffleSplit
# Here we split the data into training and test sets and implement a stratified shuffle split.
stratified = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)

for train_set, test_set in stratified.split(term_deposits, term_deposits["loan"]):
    stratified_train = term_deposits.loc[train_set]
    stratified_test = term_deposits.loc[test_set]

    stratified_train["loan"].value_counts()/len(df)
    stratified_test["loan"].value_counts()/len(df)
```

```
Out[38]: no      0.196219
yes     0.029519
Name: loan, dtype: float64
```

```
In [39]: # Separate the Labels and the features.
train_data = stratified_train # Make a copy of the stratified training set.
test_data = stratified_test
train_data.shape
test_data.shape
train_data['deposit'].value_counts()
```

```
Out[39]: no      4697
yes     4232
Name: deposit, dtype: int64
```



In [40]:

```
from sklearn.base import BaseEstimator, TransformerMixin
from sklearn.utils import check_array
from sklearn.preprocessing import LabelEncoder
from scipy import sparse

class CategoricalEncoder(BaseEstimator, TransformerMixin):
    """Encode categorical features as a numeric array.

    The input to this transformer should be a matrix of integers or strings,
    denoting the values taken on by categorical (discrete) features.
    The features can be encoded using a one-hot aka one-of-K scheme
    (`encoding='onehot'`, the default) or converted to ordinal integers
    (`encoding='ordinal'`).

    This encoding is needed for feeding categorical data to many scikit-learn
    estimators, notably linear models and SVMs with the standard kernels.
    Read more in the :ref:`User Guide <preprocessing_categorical_features>`.

    Parameters
    -----
    encoding : str, 'onehot', 'onehot-dense' or 'ordinal'
        The type of encoding to use (default is 'onehot'):
        - 'onehot': encode the features using a one-hot aka one-of-K scheme
            (or also called 'dummy' encoding). This creates a binary column for
            each category and returns a sparse matrix.
        - 'onehot-dense': the same as 'onehot' but returns a dense array
            instead of a sparse matrix.
        - 'ordinal': encode the features as ordinal integers. This results in
            a single column of integers (0 to n_categories - 1) per feature.
    categories : 'auto' or a list of lists/arrays of values.
        Categories (unique values) per feature:
        - 'auto' : Determine categories automatically from the training data.
        - list : ``categories[i]`` holds the categories expected in the ith
            column. The passed categories are sorted before encoding the data
            (used categories can be found in the ``categories_`` attribute).
    dtype : number type, default np.float64
        Desired dtype of output.
    handle_unknown : 'error' (default) or 'ignore'
        Whether to raise an error or ignore if a unknown categorical feature is
        present during transform (default is to raise). When this is parameter
        is set to 'ignore' and an unknown category is encountered during
        transform, the resulting one-hot encoded columns for this feature
        will be all zeros.
        Ignoring unknown categories is not supported for
        ``encoding='ordinal'``.
    Attributes
    -----
    categories_ : list of arrays
        The categories of each feature determined during fitting. When
        categories were specified manually, this holds the sorted categories
        (in order corresponding with output of `transform`).
    Examples
    -----
    Given a dataset with three features and two samples, we let the encoder
    find the maximum value per feature and transform the data to a binary
    one-hot encoding.
    >>> from sklearn.preprocessing import CategoricalEncoder
    >>> enc = CategoricalEncoder(handle_unknown='ignore')
    >>> enc.fit([[0, 0, 3], [1, 1, 0], [0, 2, 1], [1, 0, 2]])
    ... # doctest: +ELLIPSIS
    CategoricalEncoder(categories='auto', dtype=<... 'numpy.float64'>,
                       encoding='onehot', handle_unknown='ignore')
    >>> enc.transform([[0, 1, 1], [1, 0, 4]]).toarray()
    array([[ 1.,  0.,  0.,  1.,  0.,  1.,  0.,  0.],
           [ 0.,  1.,  1.,  0.,  0.,  0.,  0.,  0.]])
    See also
    -----
    sklearn.preprocessing.OneHotEncoder : performs a one-hot encoding of
        integer ordinal features. The ``OneHotEncoder assumes`` that input
        features take on values in the range ``[0, max(feature)]`` instead of
        using the unique values.
    sklearn.feature_extraction.DictVectorizer : performs a one-hot encoding of
        dictionary items (also handles string-valued features).
    sklearn.feature_extraction.FeatureHasher : performs an approximate one-hot
        encoding of dictionary items or strings.
    """

    def __init__(self, encoding='onehot', categories='auto', dtype=np.float64,
                 handle_unknown='error'):
        self.encoding = encoding
        self.categories = categories
        self.dtype = dtype
        self.handle_unknown = handle_unknown

    def fit(self, X, y=None):
        """Fit the CategoricalEncoder to X.

        Parameters
        -----
        X : array-like, shape [n_samples, n_feature]
            The data to determine the categories of each feature.
        Returns
        -----
        self
        """

```

```

    if self.encoding not in ['onehot', 'onehot-dense', 'ordinal']:
        template = ("encoding should be either 'onehot', 'onehot-dense' "
                    "or 'ordinal', got %s")
        raise ValueError(template % self.handle_unknown)

    if self.handle_unknown not in ['error', 'ignore']:
        template = ("handle_unknown should be either 'error' or "
                    "'ignore', got %s")
        raise ValueError(template % self.handle_unknown)

    if self.encoding == 'ordinal' and self.handle_unknown == 'ignore':
        raise ValueError("handle_unknown='ignore' is not supported for"
                         " encoding='ordinal'")

X = check_array(X, dtype=np.object, accept_sparse='csc', copy=True)
n_samples, n_features = X.shape

self._label_encoders_ = [LabelEncoder() for _ in range(n_features)]

for i in range(n_features):
    le = self._label_encoders_[i]
    Xi = X[:, i]
    if self.categories == 'auto':
        le.fit(Xi)
    else:
        valid_mask = np.in1d(Xi, self.categories[i])
        if not np.all(valid_mask):
            if self.handle_unknown == 'error':
                diff = np.unique(Xi[~valid_mask])
                msg = ("Found unknown categories {0} in column {1}"
                       " during fit".format(diff, i))
                raise ValueError(msg)
            le.classes_ = np.array(np.sort(self.categories[i]))
    self.categories_ = [le.classes_ for le in self._label_encoders_]

return self

def transform(self, X):
    """Transform X using one-hot encoding.

    Parameters
    -----
    X : array-like, shape [n_samples, n_features]
        The data to encode.

    Returns
    -----
    X_out : sparse matrix or a 2-d array
        Transformed input.
    """
    X = check_array(X, accept_sparse='csc', dtype=np.object, copy=True)
    n_samples, n_features = X.shape
    X_int = np.zeros_like(X, dtype=np.int)
    X_mask = np.ones_like(X, dtype=np.bool)

    for i in range(n_features):
        valid_mask = np.in1d(X[:, i], self.categories_[i])

        if not np.all(valid_mask):
            if self.handle_unknown == 'error':
                diff = np.unique(X[~valid_mask, i])
                msg = ("Found unknown categories {0} in column {1}"
                       " during transform".format(diff, i))
                raise ValueError(msg)
            else:
                # Set the problematic rows to an acceptable value and
                # continue. The rows are marked `X_mask` and will be
                # removed later.
                X_mask[:, i] = valid_mask
                X[:, i][~valid_mask] = self.categories_[i][0]
                X_int[:, i] = self._label_encoders_[i].transform(X[:, i])

    if self.encoding == 'ordinal':
        return X_int.astype(self.dtype, copy=False)

    mask = X_mask.ravel()
    n_values = [cats.shape[0] for cats in self.categories_]
    n_values = np.array([0] + n_values)
    indices = np.cumsum(n_values)

    column_indices = (X_int + indices[:-1]).ravel()[mask]
    row_indices = np.repeat(np.arange(n_samples, dtype=np.int32),
                           n_features)[mask]
    data = np.ones(n_samples * n_features)[mask]

    out = sparse.csc_matrix((data, (row_indices, column_indices)),
                           shape=(n_samples, indices[-1]),
                           dtype=self.dtype).tocsr()

    if self.encoding == 'onehot-dense':
        return out.toarray()
    else:

```

```
    return out
```

```
In [41]: from sklearn.base import BaseEstimator, TransformerMixin

# A class to select numerical or categorical columns
# since Scikit-Learn doesn't handle DataFrames yet
class DataFrameSelector(BaseEstimator, TransformerMixin):
    def __init__(self, attribute_names):
        self.attribute_names = attribute_names
    def fit(self, X, y=None):
        return self
    def transform(self, X):
        return X[self.attribute_names]
```

```
In [42]: train_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 8929 entries, 9867 to 9672
Data columns (total 17 columns):
deposit      8929 non-null object
age          8929 non-null int64
job          8929 non-null object
marital       8929 non-null object
education    8929 non-null object
default       8929 non-null object
balance      8929 non-null int64
housing       8929 non-null object
loan          8929 non-null object
contact       8929 non-null object
day           8929 non-null int64
month         8929 non-null object
duration     8929 non-null int64
campaign     8929 non-null int64
pdays         8929 non-null int64
previous     8929 non-null int64
poutcome     8929 non-null object
dtypes: int64(7), object(10)
memory usage: 1.2+ MB
```

```
In [43]: from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

# Making pipelines
numerical_pipeline = Pipeline([
    ("select_numeric", DataFrameSelector(["age", "balance", "day", "campaign", "pdays", "previous", "duration"]),
     ("std_scaler", StandardScaler())),
])

categorical_pipeline = Pipeline([
    ("select_cat", DataFrameSelector(["job", "education", "marital", "default", "housing", "loan", "contact", "poutcome"])),
    ("cat_encoder", CategoricalEncoder(encoding='onehot-dense'))
])

from sklearn.pipeline import FeatureUnion
preprocess_pipeline = FeatureUnion(transformer_list=[
    ("numerical_pipeline", numerical_pipeline),
    ("categorical_pipeline", categorical_pipeline),
])
```

```
In [44]: X_train = preprocess_pipeline.fit_transform(train_data)
X_train
```

```
/opt/conda/lib/python3.6/site-packages/scikit-learn/preprocessing/data.py:645: DataConversionWarning:
  Data with input dtype int64 were all converted to float64 by StandardScaler.
/opt/conda/lib/python3.6/site-packages/scikit-learn/base.py:464: DataConversionWarning:
  Data with input dtype int64 were all converted to float64 by StandardScaler.
```

```
Out[44]: array([[ 1.14643868,  1.68761105,  1.69442818, ...,  0.        ,
   0.        ,  1.        ],
   [-0.86102339, -0.35066205, -0.5560058 , ...,  0.        ,
   0.        ,  1.        ],
   [-0.94466765, -0.20504785,  0.39154535, ...,  0.        ,
   0.        ,  1.        ],
   ...,
   [-0.86102339, -0.26889658, -1.02978138, ...,  0.        ,
   0.        ,  1.        ],
   [ 0.2263519 , -0.32166951,  0.50998924, ...,  0.        ,
   0.        ,  1.        ],
   [-0.61009063, -0.34740446,  1.69442818, ...,  1.        ,
   0.        ,  0.        ]])
```

```
In [45]: y_train = train_data['deposit']
y_test = test_data['deposit']
y_train.shape
```

```
Out[45]: (8929,)
```

```
In [46]: from sklearn.preprocessing import LabelEncoder

encode = LabelEncoder()
y_train = encode.fit_transform(y_train)
y_test = encode.fit_transform(y_test)
y_train_yes = (y_train == 1)
y_train
y_train_yes
```

```
Out[46]: array([False, False,  True, ...,  True,  True, False])
```

```
In [47]: some_instance = X_train[1250]
```

```
In [48]: # Time for Classification Models
import time

from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler, LabelEncoder

from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn import tree
from sklearn.neural_network import MLPClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.gaussian_process.kernels import RBF
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import GaussianNB

dict_classifiers = {
    "Logistic Regression": LogisticRegression(),
    "Nearest Neighbors": KNeighborsClassifier(),
    "Linear SVM": SVC(),
    "Gradient Boosting Classifier": GradientBoostingClassifier(),
    "Decision Tree": tree.DecisionTreeClassifier(),
    "Random Forest": RandomForestClassifier(n_estimators=18),
    "Neural Net": MLPClassifier(alpha=1),
    "Naive Bayes": GaussianNB()
}
```

```
In [49]: no_classifiers = len(dict_classifiers.keys())

def batch_classify(X_train, Y_train, verbose = True):
    df_results = pd.DataFrame(data=np.zeros(shape=(no_classifiers,3)), columns = ['classifier', 'train_score',
    count = 0
    for key, classifier in dict_classifiers.items():
        t_start = time.clock()
        classifier.fit(X_train, Y_train)
        t_end = time.clock()
        t_diff = t_end - t_start
        train_score = classifier.score(X_train, Y_train)
        df_results.loc[count, 'classifier'] = key
        df_results.loc[count, 'train_score'] = train_score
        df_results.loc[count, 'training_time'] = t_diff
        if verbose:
            print("trained {} in {:.2f} s".format(c=key, f=t_diff))
        count+=1
    return df_results
```

```
In [50]: df_results = batch_classify(X_train, y_train)
print(df_results.sort_values(by='train_score', ascending=False))

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning:
Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

trained Logistic Regression in 0.05 s
trained Nearest Neighbors in 0.42 s

/opt/conda/lib/python3.6/site-packages/sklearn/svm/base.py:196: FutureWarning:
The default value of gamma will change from 'auto' to 'scale' in version 0.22 to account better for unscaled
features. Set gamma explicitly to 'auto' or 'scale' to avoid this warning.

trained Linear SVM in 5.18 s
trained Gradient Boosting Classifier in 1.57 s
trained Decision Tree in 0.09 s
trained Random Forest in 0.20 s
trained Neural Net in 25.09 s
trained Naive Bayes in 0.04 s
   classifier  train_score  training_time
4      Decision Tree    1.000000     0.089817
5      Random Forest    0.998320     0.203567
1      Nearest Neighbors    0.863255     0.415755
3  Gradient Boosting Classifier    0.860343     1.568130
2      Linear SVM    0.852391     5.175151
6      Neural Net    0.850487    25.092250
0      Logistic Regression    0.830776     0.050623
7      Naive Bayes    0.721693     0.044644
```

## Avoiding Overfitting:

Brief Description of Overfitting:

Overfitting is an error in the modeling algorithm that considers random noise in the fitting process rather than focusing on the underlying pattern. This phenomenon becomes evident when a model demonstrates exceptional performance on the training set but performs poorly on the test set (unseen data for the model). Overfitting occurs when the model tailors itself too closely to the noise within the training data, failing to generalize well to new, unseen data. In the examples above, the Decision Tree Classifier and Random Forest classifiers are likely overfitting, as they yield nearly perfect accuracy scores (100% and 99%).

How to Avoid Overfitting:

The most effective method to prevent overfitting is to employ cross-validation. This involves splitting the training dataset into multiple subsets. For instance, if we choose a 3-fold cross-validation, two-thirds (66%) of the data will be used for training, and one-third (33%) will be used for testing. This testing process is repeated three times, iterating through different training and test sets. The primary objective of cross-validation is to capture the overall pattern of the data, promoting a more robust and generalized model.

```
In [51]: # Use Cross-validation.
from sklearn.model_selection import cross_val_score

# Logistic Regression
log_reg = LogisticRegression()
log_scores = cross_val_score(log_reg, X_train, y_train, cv=3)
log_reg_mean = log_scores.mean()

# SVC
svc_clf = SVC()
svc_scores = cross_val_score(svc_clf, X_train, y_train, cv=3)
svc_mean = svc_scores.mean()

# KNearestNeighbors
knn_clf = KNeighborsClassifier()
knn_scores = cross_val_score(knn_clf, X_train, y_train, cv=3)
knn_mean = knn_scores.mean()

# Decision Tree
tree_clf = tree.DecisionTreeClassifier()
tree_scores = cross_val_score(tree_clf, X_train, y_train, cv=3)
tree_mean = tree_scores.mean()

# Gradient Boosting Classifier
grad_clf = GradientBoostingClassifier()
grad_scores = cross_val_score(grad_clf, X_train, y_train, cv=3)
grad_mean = grad_scores.mean()

# Random Forest Classifier
rand_clf = RandomForestClassifier(n_estimators=18)
rand_scores = cross_val_score(rand_clf, X_train, y_train, cv=3)
rand_mean = rand_scores.mean()

# NeuralNet Classifier
neural_clf = MLPClassifier(alpha=1)
neural_scores = cross_val_score(neural_clf, X_train, y_train, cv=3)
neural_mean = neural_scores.mean()

# Naives Bayes
nav_clf = GaussianNB()
nav_scores = cross_val_score(nav_clf, X_train, y_train, cv=3)
nav_mean = nav_scores.mean()

# Create a Dataframe with the results.
d = {'Classifiers': ['Logistic Reg.', 'SVC', 'KNN', 'Dec Tree', 'Grad B CLF', 'Rand FC', 'Neural Classifier',
'Crossval Mean Scores': [log_reg_mean, svc_mean, knn_mean, tree_mean, grad_mean, rand_mean, neural_mean,
result_df = pd.DataFrame(data=d)

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning:
Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning:
Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning:
Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

/opt/conda/lib/python3.6/site-packages/sklearn/svm/base.py:196: FutureWarning:
The default value of gamma will change from 'auto' to 'scale' in version 0.22 to account better for unscaled
features. Set gamma explicitly to 'auto' or 'scale' to avoid this warning.

/opt/conda/lib/python3.6/site-packages/sklearn/svm/base.py:196: FutureWarning:
The default value of gamma will change from 'auto' to 'scale' in version 0.22 to account better for unscaled
features. Set gamma explicitly to 'auto' or 'scale' to avoid this warning.
```

```
In [52]: # All our models perform well but I will go with GradientBoosting.
result_df = result_df.sort_values(by=['Crossval Mean Scores'], ascending=False)
result_df
```

Out[52]:

Classifiers	Crossval Mean Scores
6 Neural Classifier	0.847801
7 Naives Bayes	0.847801
4 Grad B CLF	0.845224
1 SVC	0.840186
5 Rand FC	0.836936
0 Logistic Reg.	0.828425
2 KNN	0.804458
3 Dec Tree	0.785753

### Confusion Matrix:

		Prediction	
		0	1
Actual	0	TN	FP
	1	FN	TP

Insights of a Confusion Matrix:

The main purpose of a confusion matrix is to assess how well our model performs in classifying potential clients likely to subscribe to a term deposit. The confusion matrix comprises four terms: True Positives, False Positives, True Negatives, and False Negatives.

Positive/Negative: Types of Classes (labels) - ["No", "Yes"]

True/False: Correctly or Incorrectly classified by the model

True Negatives (Top-Left Square): The number of correct classifications of the "No" class or potential clients not willing to subscribe to a term deposit.

False Negatives (Top-Right Square): The number of incorrect classifications of the "No" class or potential clients not willing to subscribe to a term deposit.

False Positives (Bottom-Left Square): The number of incorrect classifications of the "Yes" class or potential clients willing to subscribe to a term deposit.

True Positives (Bottom-Right Square): The number of correct classifications of the "Yes" class or potential clients willing to subscribe to a term deposit.

```
In [53]: # Cross validate our Gradient Boosting Classifier
from sklearn.model_selection import cross_val_predict
y_train_pred = cross_val_predict(grad_clf, X_train, y_train, cv=3)
```

```
In [54]: from sklearn.metrics import accuracy_score
grad_clf.fit(X_train, y_train)
print ("Gradient Boost Classifier accuracy is %2.2f" % accuracy_score(y_train, y_train_pred))
```

Gradient Boost Classifier accuracy is 0.85

```
In [55]: from sklearn.metrics import confusion_matrix
# 4697: no's, 4232: yes
conf_matrix = confusion_matrix(y_train, y_train_pred)
f, ax = plt.subplots(figsize=(12, 8))
sns.heatmap(conf_matrix, annot=True, fmt="d", linewidths=.5, ax=ax)
plt.title("Confusion Matrix", fontsize=20)
plt.subplots_adjust(left=0.15, right=0.99, bottom=0.15, top=0.99)
ax.set_yticks(np.arange(conf_matrix.shape[0]) + 0.5, minor=False)
ax.set_xticklabels([""])
ax.set_yticklabels(['Refused T. Deposits', 'Accepted T. Deposits'], fontsize=16, rotation=360)
plt.show()
```



## Precision and Recall:

Recall:

Recall represents the total number of "Yes" labels in the dataset. It indicates how many "Yes" labels our model correctly identifies.

Precision:

Precision signifies how confident our model is when predicting that the actual label is a "Yes."

Recall-Precision Tradeoff:

There is a tradeoff between precision and recall—when precision increases, recall decreases, and vice versa. For example, if we raise precision from 30% to 60%, the model is selecting predictions it believes are 60% certain. In cases where the model rates an instance at 58% likelihood of being a potential client subscribing to a term deposit, it classifies it as a "No." However, that instance was actually a "Yes" (a potential client did subscribe to a term deposit). Hence, the higher the precision, the more likely the model is to miss instances that are actually a "Yes."

```
In [56]: # Let's find the scores for precision and recall.
from sklearn.metrics import precision_score, recall_score

print('Precision Score: ', precision_score(y_train, y_train_pred))

print('Recall Score: ', recall_score(y_train, y_train_pred))
```

Precision Score: 0.8244135732179458

Recall Score: 0.8553875236294896

```
In [57]: from sklearn.metrics import f1_score

f1_score(y_train, y_train_pred)
```

Out[57]: 0.8396149831845067

```
In [58]: y_scores = grad_clf.decision_function([some_instance])
y_scores
```

Out[58]: array([-3.88131785])

```
In [59]: # Increasing the threshold decreases the recall.  
threshold = 0  
y_some_digit_pred = (y_scores > threshold)
```

```
In [60]: y_scores = cross_val_predict(grad_clf, X_train, y_train, cv=3, method="decision_function")  
neural_y_scores = cross_val_predict(neural_clf, X_train, y_train, cv=3, method="predict_proba")  
naives_y_scores = cross_val_predict(nav_clf, X_train, y_train, cv=3, method="predict_proba")
```

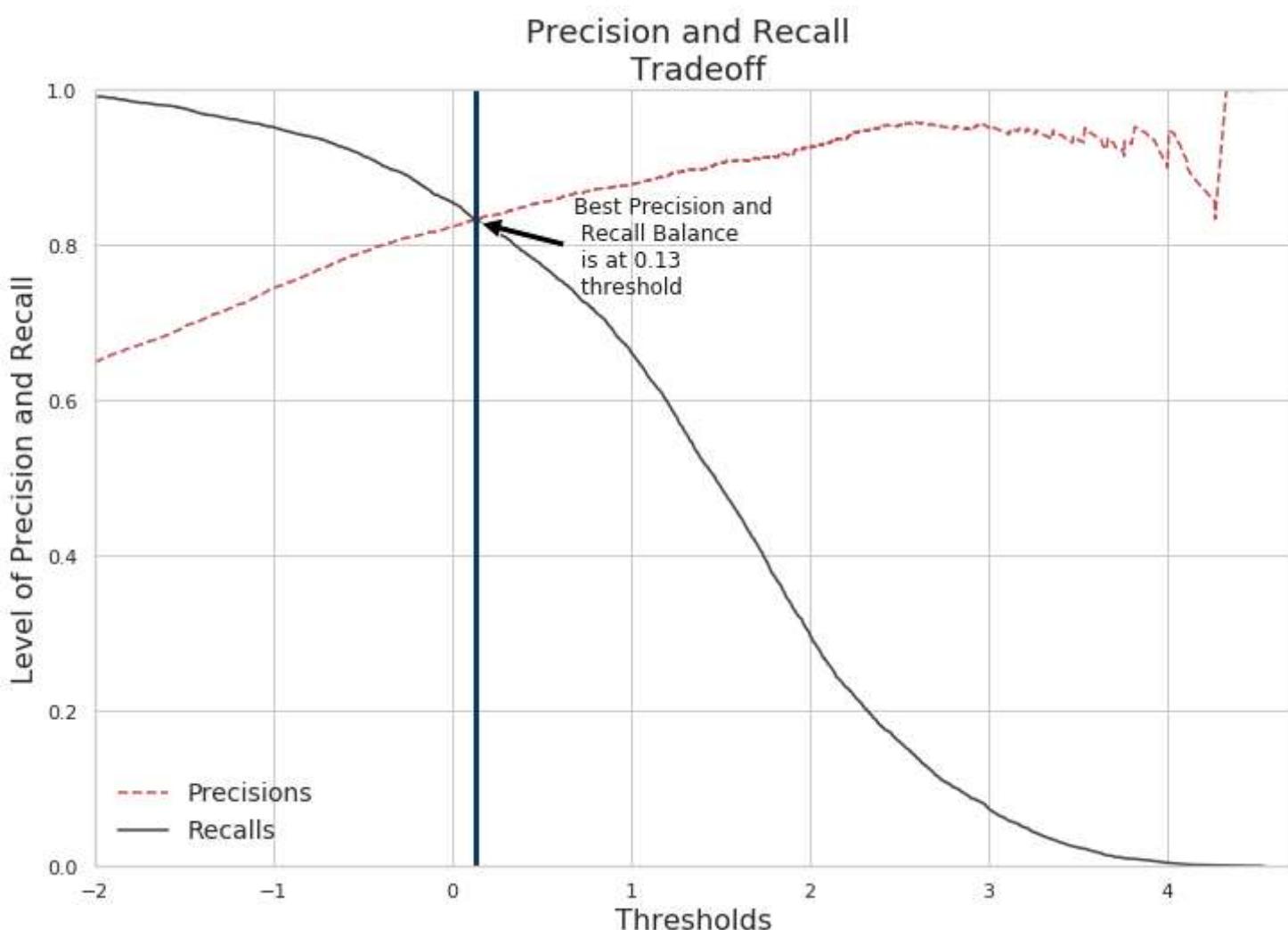
```
In [61]: # hack to work around issue #9589 introduced in Scikit-Learn 0.19.0  
if y_scores.ndim == 2:  
    y_scores = y_scores[:, 1]  
  
if neural_y_scores.ndim == 2:  
    neural_y_scores = neural_y_scores[:, 1]  
  
if naives_y_scores.ndim == 2:  
    naives_y_scores = naives_y_scores[:, 1]
```

```
In [62]: y_scores.shape
```

```
Out[62]: (8929,)
```

```
In [63]: # How can we decide which threshold to use? We want to return the scores instead of predictions with this code  
from sklearn.metrics import precision_recall_curve  
  
precisions, recalls, threshold = precision_recall_curve(y_train, y_scores)
```

```
In [64]: def precision_recall_curve(precisions, recalls, thresholds):  
    fig, ax = plt.subplots(figsize=(12,8))  
    plt.plot(thresholds, precisions[:-1], "r--", label="Precisions")  
    plt.plot(thresholds, recalls[:-1], "#424242", label="Recalls")  
    plt.title("Precision and Recall \n Tradeoff", fontsize=18)  
    plt.ylabel("Level of Precision and Recall", fontsize=16)  
    plt.xlabel("Thresholds", fontsize=16)  
    plt.legend(loc="best", fontsize=14)  
    plt.xlim([-2, 4.7])  
    plt.ylim([0, 1])  
    plt.axvline(x=0.13, linewidth=3, color="#0B3861")  
    plt.annotate('Best Precision and \n Recall Balance \n is at 0.13 \n threshold ', xy=(0.13, 0.83), xytext=(  
        textcoords="offset points",  
        arrowprops=dict(facecolor='black', shrink=0.05),  
        fontsize=12,  
        color='k')  
  
precision_recall_curve(precisions, recalls, threshold)  
plt.show()
```

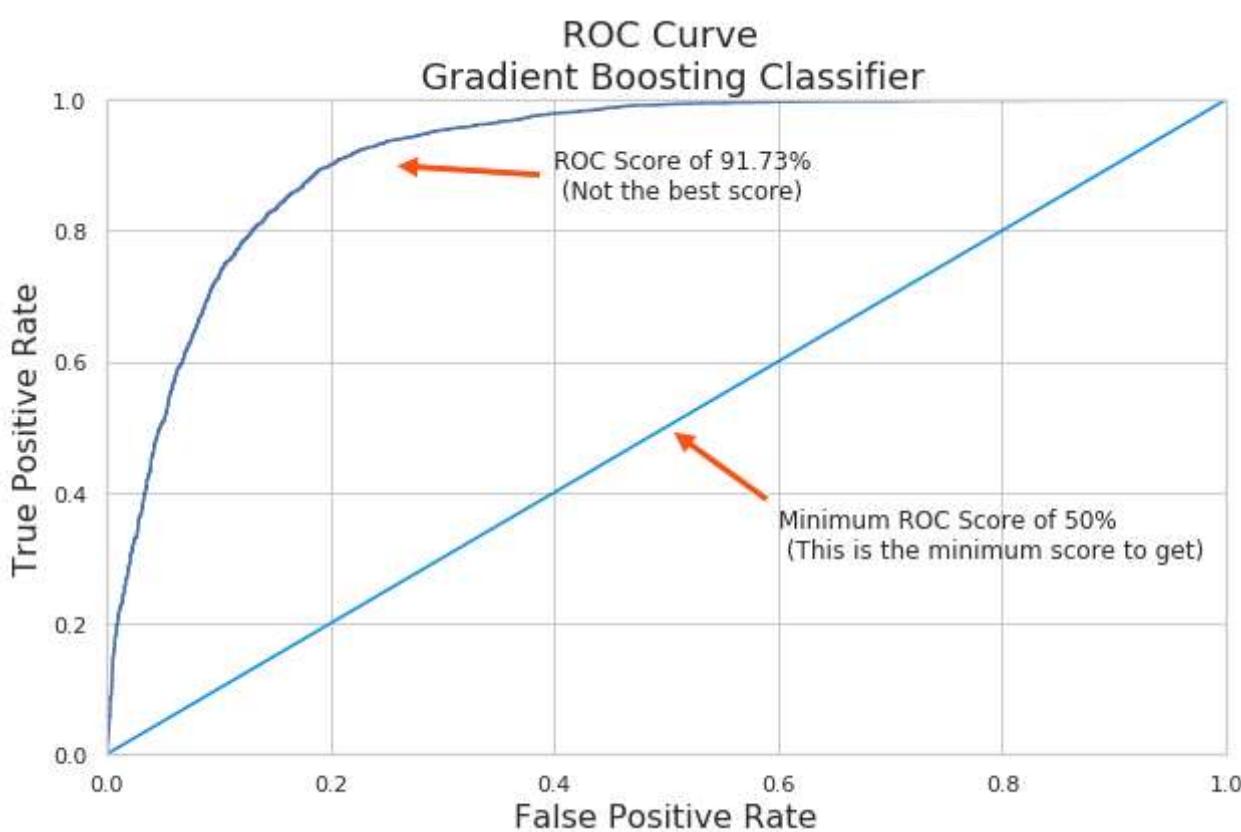


### ROC Curve (Receiver Operating Characteristic):

The ROC curve illustrates how effectively our classifier distinguishes between term deposit subscriptions (True Positives) and non-term deposit subscriptions. On the X-axis, we have False Positive Rates (Specificity), and on the Y-axis, we have the True Positive Rate (Sensitivity). As the curve moves, the threshold of classification changes, resulting in different values. The closer the curve is to the top-left corner, the better our model is at effectively separating both classes.

```
In [65]: from sklearn.metrics import roc_curve
# Gradient Boosting Classifier
# Neural Classifier
# Naives Bayes Classifier
grd_fpr, grd_tpr, threshold = roc_curve(y_train, y_scores)
neu_fpr, neu_tpr, neu_threshold = roc_curve(y_train, neural_y_scores)
nav_fpr, nav_tpr, nav_threshold = roc_curve(y_train, naives_y_scores)
```

```
In [66]: def graph_roc_curve(false_positive_rate, true_positive_rate, label=None):
    plt.figure(figsize=(10,6))
    plt.title('ROC Curve \n Gradient Boosting Classifier', fontsize=18)
    plt.plot(false_positive_rate, true_positive_rate, label=label)
    plt.plot([0, 1], [0, 1], '#0C8EE0')
    plt.axis([0, 1, 0, 1])
    plt.xlabel('False Positive Rate', fontsize=16)
    plt.ylabel('True Positive Rate', fontsize=16)
    plt.annotate('ROC Score of 91.73% \n (Not the best score)', xy=(0.25, 0.9), xytext=(0.4, 0.85),
                 arrowprops=dict(facecolor='#F75118', shrink=0.05),
                 )
    plt.annotate('Minimum ROC Score of 50% \n (This is the minimum score to get)', xy=(0.5, 0.5), xytext=(0.6,
                 arrowprops=dict(facecolor='#F75118', shrink=0.05),
                 )
graph_roc_curve(grd_fpr, grd_tpr, threshold)
plt.show()
```



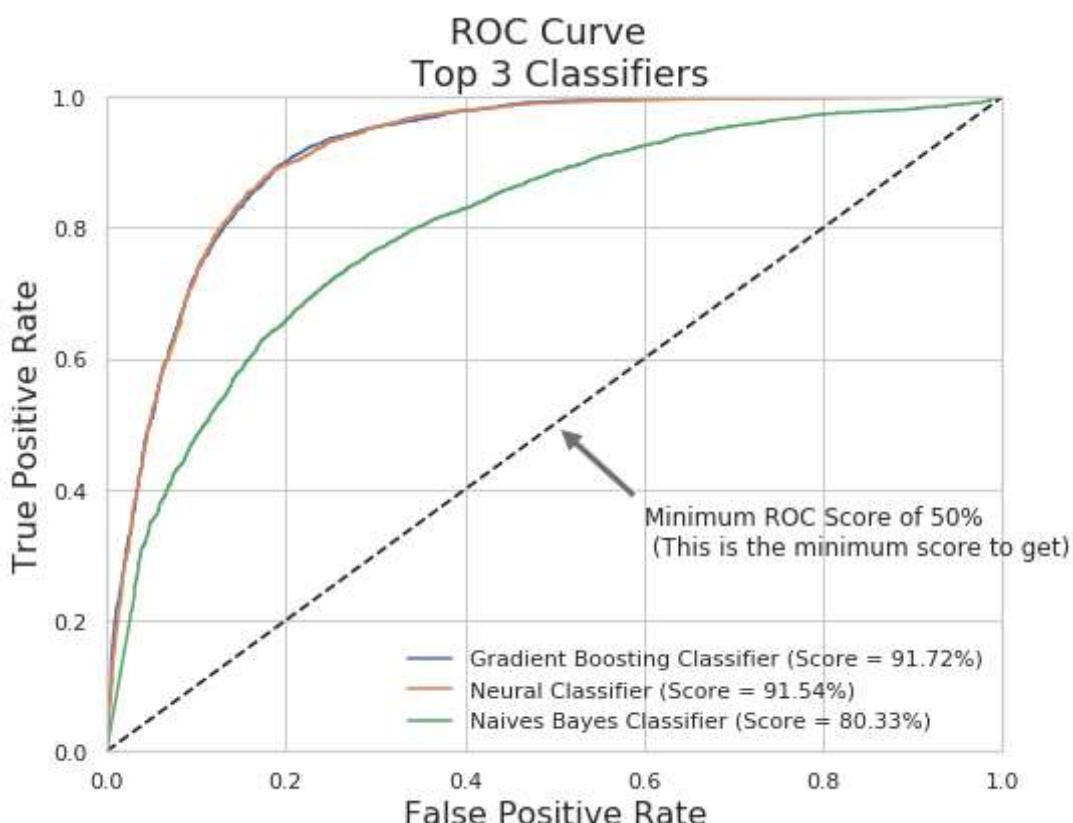
```
In [67]: from sklearn.metrics import roc_auc_score

print('Gradient Boost Classifier Score: ', roc_auc_score(y_train, y_scores))
print('Neural Classifier Score: ', roc_auc_score(y_train, neural_y_scores))
print('Naives Bayes Classifier: ', roc_auc_score(y_train, naives_y_scores))

Gradient Boost Classifier Score:  0.9173215628927768
Neural Classifier Score:  0.9162999408784837
Naives Bayes Classifier:  0.803363959942255
```

```
In [68]: def graph_roc_curve_multiple(grd_fpr, grd_tpr, neu_fpr, neu_tpr, nav_fpr, nav_tpr):
    plt.figure(figsize=(8,6))
    plt.title('ROC Curve \n Top 3 Classifiers', fontsize=18)
    plt.plot(grd_fpr, grd_tpr, label='Gradient Boosting Classifier (Score = 91.72%)')
    plt.plot(neu_fpr, neu_tpr, label='Neural Classifier (Score = 91.54%)')
    plt.plot(nav_fpr, nav_tpr, label='Naives Bayes Classifier (Score = 80.33%)')
    plt.plot([0, 1], [0, 1], 'k--')
    plt.axis([0, 1, 0, 1])
    plt.xlabel('False Positive Rate', fontsize=16)
    plt.ylabel('True Positive Rate', fontsize=16)
    plt.annotate('Minimum ROC Score of 50% \n (This is the minimum score to get)', xy=(0.5, 0.5), xytext=(0.6,
        arrowprops=dict(facecolor='#6E726D', shrink=0.05),
        )
    plt.legend()

graph_roc_curve_multiple(grd_fpr, grd_tpr, neu_fpr, neu_tpr, nav_fpr, nav_tpr)
plt.show()
```



```
In [69]: grad_clf.predict_proba([some_instance])
```

```
Out[69]: array([[0.97979311, 0.02020689]])
```

```
In [70]: # Let's see what does our classifier predict.
grad_clf.predict([some_instance])
```

```
Out[70]: array([0])
```

```
In [71]: y_train[1250]
```

```
Out[71]: 0
```

## Which Features Influence the Result of a Term Deposit Suscription?

### **DecisionTreeClassifier:**

The top three most important features for our classifier are:

Duration: Represents how long the conversation between the sales representative and the potential client took.

Contact: Indicates the number of contacts made to the potential client within the same marketing campaign.

Month: Refers to the month of the year.

```
In [72]: import numpy as np
import matplotlib.pyplot as plt
from sklearn import tree
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
plt.style.use('seaborn-white')

# Convert the columns into categorical variables
term_deposits['job'] = term_deposits['job'].astype('category').cat.codes
term_deposits['marital'] = term_deposits['marital'].astype('category').cat.codes
term_deposits['education'] = term_deposits['education'].astype('category').cat.codes
term_deposits['contact'] = term_deposits['contact'].astype('category').cat.codes
term_deposits['poutcome'] = term_deposits['poutcome'].astype('category').cat.codes
term_deposits['month'] = term_deposits['month'].astype('category').cat.codes
term_deposits['default'] = term_deposits['default'].astype('category').cat.codes
term_deposits['loan'] = term_deposits['loan'].astype('category').cat.codes
term_deposits['housing'] = term_deposits['housing'].astype('category').cat.codes

# Let's create new splittings as we did before, but now we have modified the data, so we need to do it once more
# Create train and test splits
target_name = 'deposit'
X = term_deposits.drop('deposit', axis=1)

label=term_deposits[target_name]

X_train, X_test, y_train, y_test = train_test_split(X,label,test_size=0.2, random_state=42, stratify=label)

# Build a classification task using 3 informative features
tree = tree.DecisionTreeClassifier(
    class_weight='balanced',
    min_weight_fraction_leaf = 0.01
)

tree = tree.fit(X_train, y_train)
importances = tree.feature_importances_
feature_names = term_deposits.drop('deposit', axis=1).columns
indices = np.argsort(importances)[::-1]

# Print the feature ranking
print("Feature ranking:")

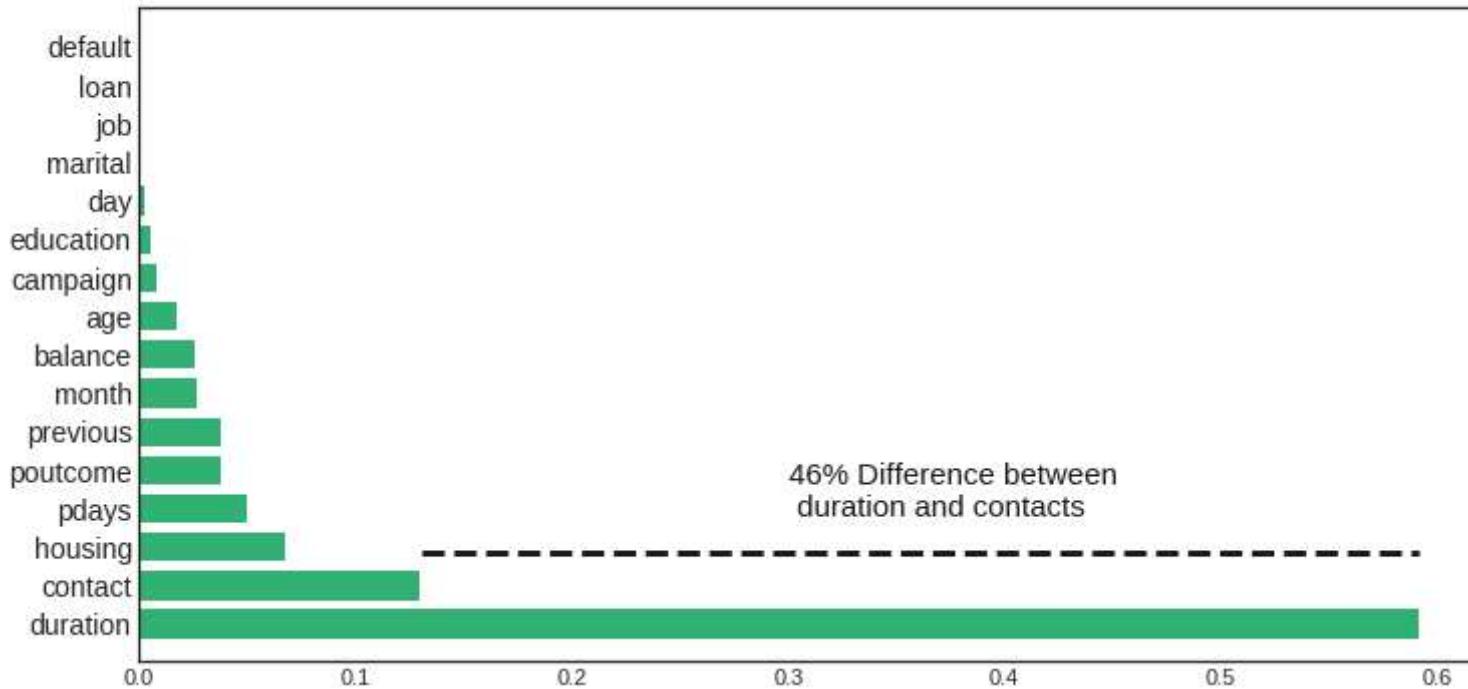
for f in range(X_train.shape[1]):
    print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))

# Plot the feature importances of the forest
def feature_importance_graph(indices, importances, feature_names):
    plt.figure(figsize=(12,6))
    plt.title("Determining Feature importances \n with DecisionTreeClassifier", fontsize=18)
    plt.barh(range(len(indices)), importances[indices], color='#31B173', align="center")
    plt.yticks(range(len(indices)), feature_names[indices], rotation='horizontal', fontsize=14)
    plt.ylim([-1, len(indices)])
    plt.axhline(y=1.85, xmin=0.21, xmax=0.952, color='k', linewidth=3, linestyle='--')
    plt.text(0.30, 2.8, '46% Difference between \n duration and contacts', color='k', fontsize=15)

feature_importance_graph(indices, importances, feature_names)
plt.show()

Feature ranking:
1. feature 11 (0.591310)
2. feature 8 (0.129966)
3. feature 6 (0.067020)
4. feature 13 (0.049923)
5. feature 15 (0.038138)
6. feature 14 (0.037830)
7. feature 10 (0.026646)
8. feature 5 (0.025842)
9. feature 0 (0.017757)
10. feature 12 (0.007889)
11. feature 3 (0.005280)
12. feature 9 (0.002200)
13. feature 2 (0.000147)
14. feature 1 (0.000050)
15. feature 7 (0.000000)
16. feature 4 (0.000000)
```

## Determining Feature importances with DecisionTreeClassifier



## GradientBoosting Classifier Wins

The Gradient Boosting classifier proves to be the best model for predicting whether a potential client will subscribe to a term deposit, achieving an accuracy of 84.6%.

```
In [73]: # Our three classifiers are grad_clf, nav_clf and neural_clf
from sklearn.ensemble import VotingClassifier
```

```
voting_clf = VotingClassifier(
    estimators=[('gbc', grad_clf), ('nav', nav_clf), ('neural', neural_clf)],
    voting='soft'
)
voting_clf.fit(X_train, y_train)
```

```
Out[73]: VotingClassifier(estimators=[('gbc', GradientBoostingClassifier(criterion='friedman_mse', init=None,
   learning_rate=0.1, loss='deviance', max_depth=3,
   max_features=None, max_leaf_nodes=None,
   min_impurity_decrease=0.0, min_impurity_split=None,
   min_samples_leaf=1,...=True, solver='adam', tol=0.0001,
   validation_fraction=0.1, verbose=False, warm_start=False)),
   flatten_transform=None, n_jobs=None, voting='soft', weights=None)
```

```
In [74]: from sklearn.metrics import accuracy_score

for clf in (grad_clf, nav_clf, neural_clf, voting_clf):
    clf.fit(X_train, y_train)
    predict = clf.predict(X_test)
    print(clf.__class__.__name__, accuracy_score(y_test, predict))

GradientBoostingClassifier 0.8463949843260188
GaussianNB 0.7514554411106136
MLPClassifier 0.6914464845499328
VotingClassifier 0.8226600985221675
```

## Conclusion

### What Actions should the Bank Consider?

#### Solutions for the Next Marketing Campaign:

Months of Marketing Activity: We observed that the month with the highest level of marketing activity was May, but it had the lowest effective rate for term deposit offers (-34.49%). For the next marketing campaign, it would be wise for the bank to focus on the months of March, September, October, and December. (December should be considered due to its lowest marketing activity, as there might be a reason for this.)

Seasonality: Potential clients tended to subscribe to term deposits during the fall and winter seasons. The next marketing campaign should focus its activities during these seasons.

Campaign Calls: Implementing a policy that limits calls to the same potential client to no more than three would save time and effort. Remember, excessive calls to the same potential client may lead to a higher likelihood of refusal to open a term deposit.

Age Category: Targeting potential clients in their 20s or younger and those in their 60s or older for the next campaign is recommended. The youngest category had a 60% chance of subscribing to a term deposit, while the eldest category had a 76% chance. Addressing these two categories could significantly increase the likelihood of more term deposit subscriptions.

Occupation: Students and retirees were the most likely to subscribe to a term deposit. Retired individuals, in particular, tend to have more term deposits, possibly seeking cash through interest payments. Students were another group that frequently subscribed to term deposits.

**House Loans and Balances:** Potential clients in the low and no balance categories were more likely to have a house loan. Individuals with average and high balances were less likely to have a house loan and more likely to open a term deposit. The next campaign should focus on individuals with average and high balances.

**Develop a Questionnaire during Calls:** Since the duration of the call positively correlates with a potential client opening a term deposit, introducing an engaging questionnaire during calls may increase conversation length. While this doesn't guarantee a term deposit subscription, it can enhance engagement, leading to a higher probability of subscription.

**Target Individuals with a Higher Duration (Above 375):** Targeting individuals above average in duration (above 375) is recommended, as there is a high likelihood that this group will open a term deposit account (78%). This would significantly boost the success rate of the next marketing campaign.

By combining these strategies and refining the target market, the next campaign is likely to be more effective than the current one.