

Responsible AI

Project Description

A large bank has asked us to evaluate the marketing algorithms they use for retail banking. Their sophisticated phone marketing algorithm predicts whether a certain person will subscribe to a term deposit or not. Based on that assessment, the bank then optimises its phone calling strategy. With this algorithm, the bank has been successful in predicting which clients are more likely to subscribe to their term deposits.

Management is now interested in finding out how a classification model can lead to certain decision-making processes.

Data Dictionary

Input variables:

bank client data:

- age (numeric)
- job : type of job (categorical: 'admin', 'blue-collar','entrepreneur','housemaid','management','retired','self-employed','services','student','technician','unemployed','unknown')
- marital : marital status (categorical: 'divorced','married','single','unknown'; note: 'divorced' means divorced or widowed)
- education (categorical: 'basic.4y','basic.6y','basic.9y','high.school','illiterate','professional.course','university.degree','unknown')
- default: has credit in default? (categorical: 'no','yes','unknown')
- housing: has housing loan? (categorical: 'no','yes','unknown')
- loan: has personal loan? (categorical: 'no','yes','unknown')

related with the last contact of the current campaign:

- contact: contact communication type (categorical: 'cellular','telephone')
- month: last contact month of year (categorical: 'jan', 'feb', 'mar', ..., 'nov', 'dec')
- day_of_week: last contact day of the week (categorical: 'mon','tue','wed','thu','fri')
- duration: last contact duration, in seconds (numeric). Important note this attribute highly affects the output target (e.g. if duration=0 then y= 'no'). Yet, the duration is not known before a call is performed. Also, after the end of the call y is obviously known. Thus, this input should only be included for benchmark purposes and should be discarded if the intention is to have a realistic predictive model.
- other attributes:
- campaign: number of contacts performed during this campaign and for this client (numeric, includes last contact)
- pdays: number of days that passed by after the client was last contacted from a previous campaign (numeric, 999 means client was not previously contacted)
- previous: number of contacts performed before this campaign and for this client (numeric)
- poutcome: outcome of the previous marketing campaign (categorical: 'failure','nonexistent','success')

social and economic context attributes

- emp.var.rate: employment variation rate - quarterly indicator (numeric)
- cons.price.idx: consumer price index - monthly indicator (numeric)
- cons.conf.idx: consumer confidence index - monthly indicator (numeric)
- euribor3m: euribor 3 month rate - daily indicator (numeric)
- nr.employed: number of employees - quarterly indicator (numeric)

Output variable (desired target):

21 - y - has the client subscribed a term deposit? (binary: 'yes','no')

Summary

Number of employees and age has influence on who will subscribe to term deposit

Import Libraries

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

import sklearn

import shap
import statsmodels.api as sm

import datetime
from datetime import datetime, timedelta

import scipy.stats

import pandas_profiling
from pandas_profiling import ProfileReport

#import graphviz

# Import xgboost as xgb
# from xgboost import XGBClassifier, XGBRegressor
# from xgboost import to_graphviz, plot_importance

# from sklearn.experimental import enable_hist_gradient_boosting
# from sklearn.linear_model import ElasticNet, Lasso, LinearRegression, LogisticRegression, Ridge
# from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor, ExtraTreesClassifier, ExtraTree
# from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor, HistGradientBoostingClassifier

%matplotlib inline
# Set the default autosave frequency in seconds
autosave 60
sns.set_style('dark')
sns.set(font_scale=1.2)

plt.rc('axes', labelsize=14)
plt.rc('text', fontstyle='normal')
plt.rc('ytick', labelsize=12)

# from sklearn.pipeline import Pipeline
# from sklearn.model_selection import RepeatedStratifiedKFold
# from sklearn.feature_selection import RFECV, SelectKBest, f_classif, f_regression, chi2

from sklearn.inspection import import_permutation_importance
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.tree import export_graphviz, plot_tree
from sklearn.metrics import confusion_matrix, ClassificationReport, mean_absolute_error, mean_squared_error
from sklearn.metrics import plot_confusion_matrix, plot_precision_recall_curve, plot_roc_curve, accuracy_score
from sklearn.metrics import auc, f1_score, precision_score, recall_score, roc_auc_score

# from tpot import TPOTClassifier, TPOTRegressor
# from imblearn.under_sampling import RandomUnderSampler
# from imblearn.over_sampling import RandomOverSampler
# from imblearn.over_sampling import SMOTE

import warnings
warnings.filterwarnings('ignore')

# Import pickle
# from pickle import dump, load

# Use Folium library to plot values on a map.
# Import folium

# Use Feature-Engine library

# Import feature engine missing data imputers as ndi
# from feature_engine.outlier_removers import Winsorizer
# from feature_engine import categorical_encoders as ce

# from pycaret.classification import *
# from pycaret.clustering import *
# from pycaret.regression import *

pd.set_option('display.max_columns',None)
#pd.set_option('display.max_rows',100)
pd.set_option('display.width', 1000)
pd.set_option('display.float_format', '{:1.2f}'.format)

random.seed(0)
np.random.seed(0)
np.set_printoptions(suppress=True)
```

Autosaving every 60 seconds

Exploratory Data Analysis

```
In [2]: df = pd.read_csv("bank-additional-full.csv", sep=';')

In [3]: df

Out[3]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pdays
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	mon	149	1	999
1	57	services	married	high.school	unknown	no	no	telephone	may	mon	261	1	999
2	37	services	married	high.school	no	yes	no	telephone	may	mon	226	1	999
3	40	admin.	married	basic.6y	no	no	no	telephone	may	mon	151	1	999
4	56	services	married	high.school	no	no	yes	telephone	may	mon	307	1	999
...
41183	73	retired	married	professional.course	no	yes	no	cellular	nov	fri	334	1	999
41184	46	blue-collar	married	professional.course	no	no	no	cellular	nov	fri	383	1	999
41185	56	retired	married	university.degree	no	yes	no	cellular	nov	fri	189	2	999
41186	44	technician	married	professional.course	no	no	no	cellular	nov	fri	442	1	999
41187	74	retired	married	professional.course	no	yes	no	cellular	nov	fri	239	3	999

41188 rows x 14 columns

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41188 entries, 0 to 41187
Data columns (total 14 columns):
#   column                Non-Null Count  Dtype
---  ---                ---
0   age                   41188 non-null    int64
1   job                   41188 non-null    object
2   marital               41188 non-null    object
3   education             41188 non-null    object
4   default               41188 non-null    object
5   housing               41188 non-null    object
6   loan                  41188 non-null    object
7   contact               41188 non-null    object
8   month                 41188 non-null    object
9   day_of_week           41188 non-null    object
10  duration              41188 non-null    int64
11  campaign              41188 non-null    int64
12  pdays                 41188 non-null    int64
13  previous              41188 non-null    int64
14  emp.var.rate          41188 non-null    float64
15  cons.price.idx         41188 non-null    float64
16  cons.conf.idx         41188 non-null    float64
17  euribor3m             41188 non-null    float64
18  nr.employed           41188 non-null    float64
19  y                     41188 non-null    object
dtypes: float64(5), int64(5), object(11)
memory usage: 6.4 MB
```

```
In [5]: df.describe(include='all')

Out[5]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pdays
count	41188.00	41188	41188	41188	41188	41188	41188	41188	41188	41188	41188.00	41188.00	41188.00
unique	nan	12	4	8	3	3	3	2	10	5	nan	nan	nan
top	nan	admin.	married	university.degree	no	yes	no	cellular	may	thu	nan	nan	nan
freq	nan	10422	24928	12168	32588	21576	33950	26144	13769	8623	nan	nan	nan
std	40.02	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	258.29	2.57	962.48
mean	10.42	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	259.28	2.77	186.91
min	17.00	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	0.00	1.00	0.00
50%	32.00	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	102.00	1.00	999.00
75%	38.00	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	180.00	2.00	999.00
max	75.00	47.00	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	319.00	3.00	999.00
95%	47.00	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	4918.00	56.00	999.00

```
In [6]: df.shape

Out[6]: (41188, 14)
```

```
In [7]: df.columns

Out[7]: Index(['age', 'job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'month', 'day_of_week', 'duration', 'campaign', 'pdays', 'previous', 'emp.var.rate', 'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr.employed', 'y'], dtype='object')
```

Groupby Function

```
In [8]: df.groupby('by=y').mean()

Out[8]:
```

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m	nr.employed
y										
no	39.91	220.84	2.63	984.11	0.13	0.25	93.60	-40.59	3.81	5176.17
yes	40.91	553.19	2.05	792.04	0.49	-1.23	93.35	-39.79	2.12	5095.12

```
In [9]: df.groupby('by=y').count()

Out[9]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pdays	previous	pout
no	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548	36548
yes	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640	4640

```
In [10]: df.y.value_counts()

Out[10]: no      36548
yes       4640
Name: y, dtype: int64
```

Pandas Profiling Reports

```
In [11]: profile = ProfileReport(df=df, title='Bank Marketing Report', minimal=True)
```

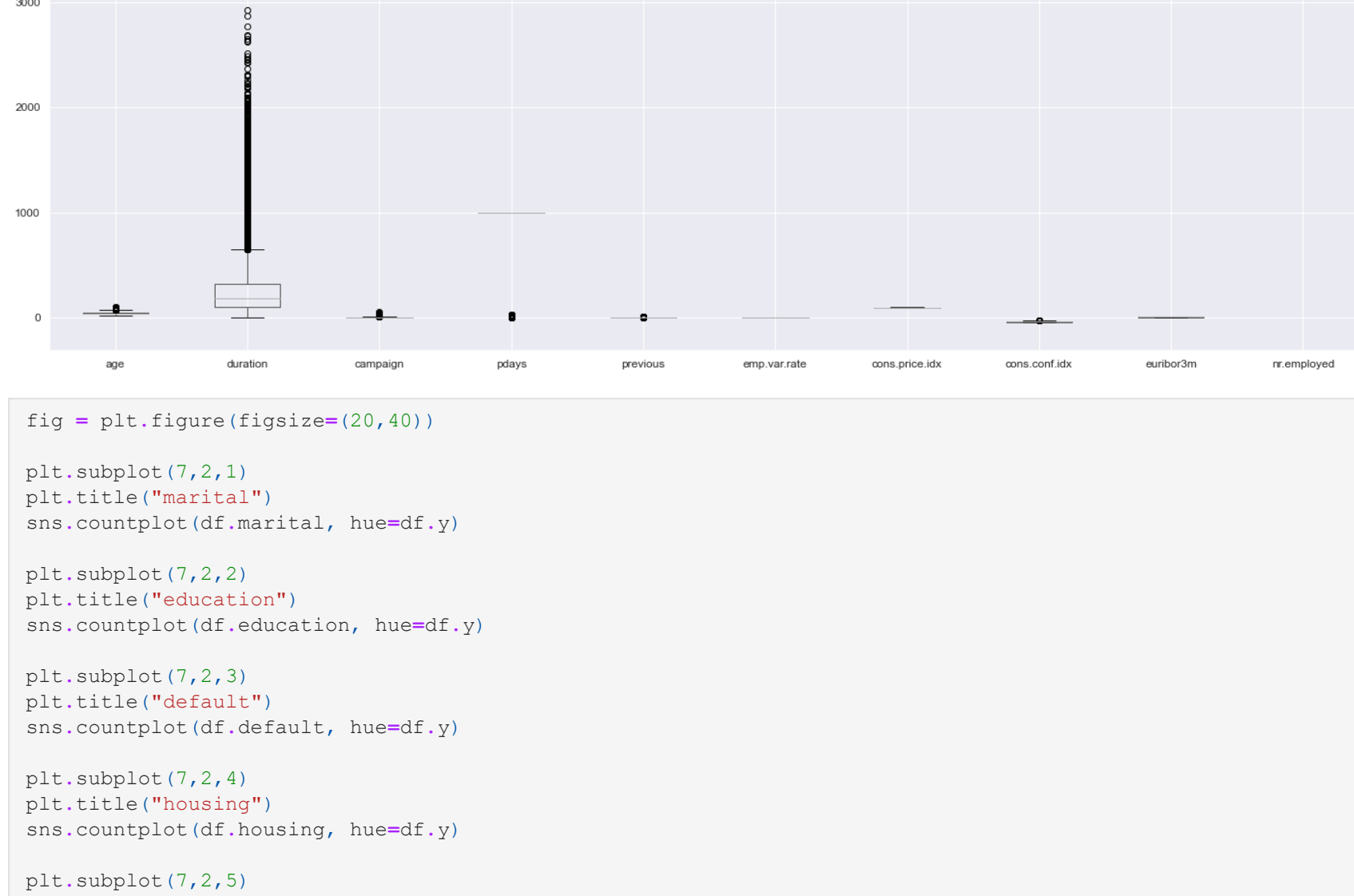
```
In [12]: profile.to_notebook_iframe()
```

Bank Marketing Report

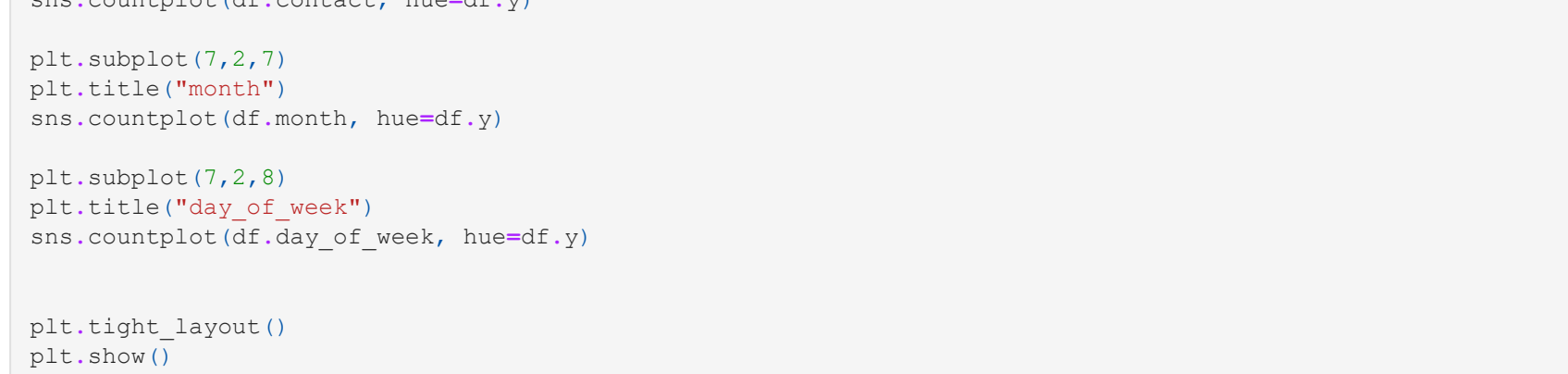
Overview

Variables

Overview



Variables



```
In [13]: profile.to_file("your_report.html")
```

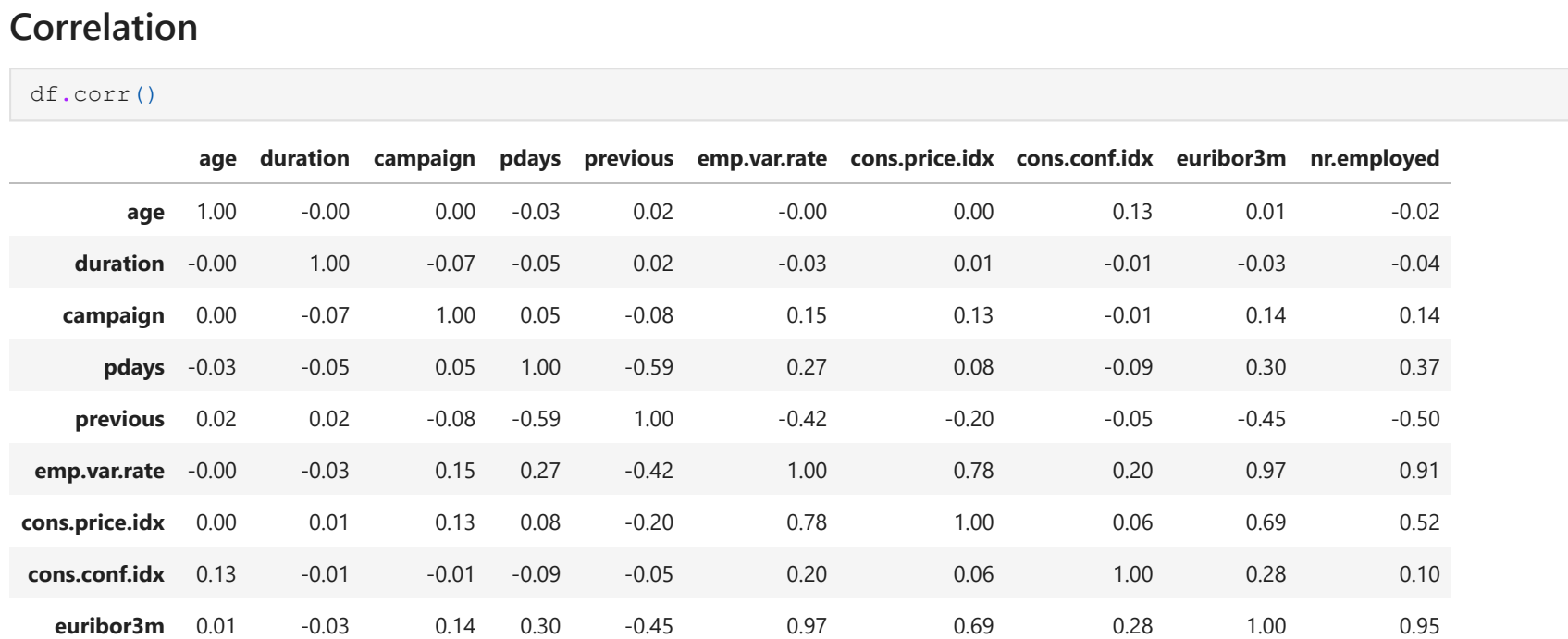
Data Visualization

Univariate Data Exploration

```
In [14]: df.hist(bins=50, figsize=(20,15))
plt.suptitle('Feature Distribution', x=0.5, y=1.02, ha='center', fontsize=20)
plt.tight_layout()
plt.show()
```



```
In [15]: df.boxplot(figsize=(20,10))
plt.suptitle('BoxPlot', x=0.5, y=1.02, ha='center', fontsize=20)
plt.tight_layout()
plt.show()
```



```
In [16]: fig = plt.figure(figsize=(20,40))
plt.subplot(7,2,1)
plt.title('marital')
sns.countplot(df.marital, hue=df.y)

plt.subplot(7,2,2)
plt.title('education')
sns.countplot(df.education, hue=df.y)

plt.subplot(7,2,3)
plt.title('default')
sns.countplot(df.default, hue=df.y)

plt.subplot(7,2,4)
plt.title('housing')
sns.countplot(df.housing, hue=df.y)

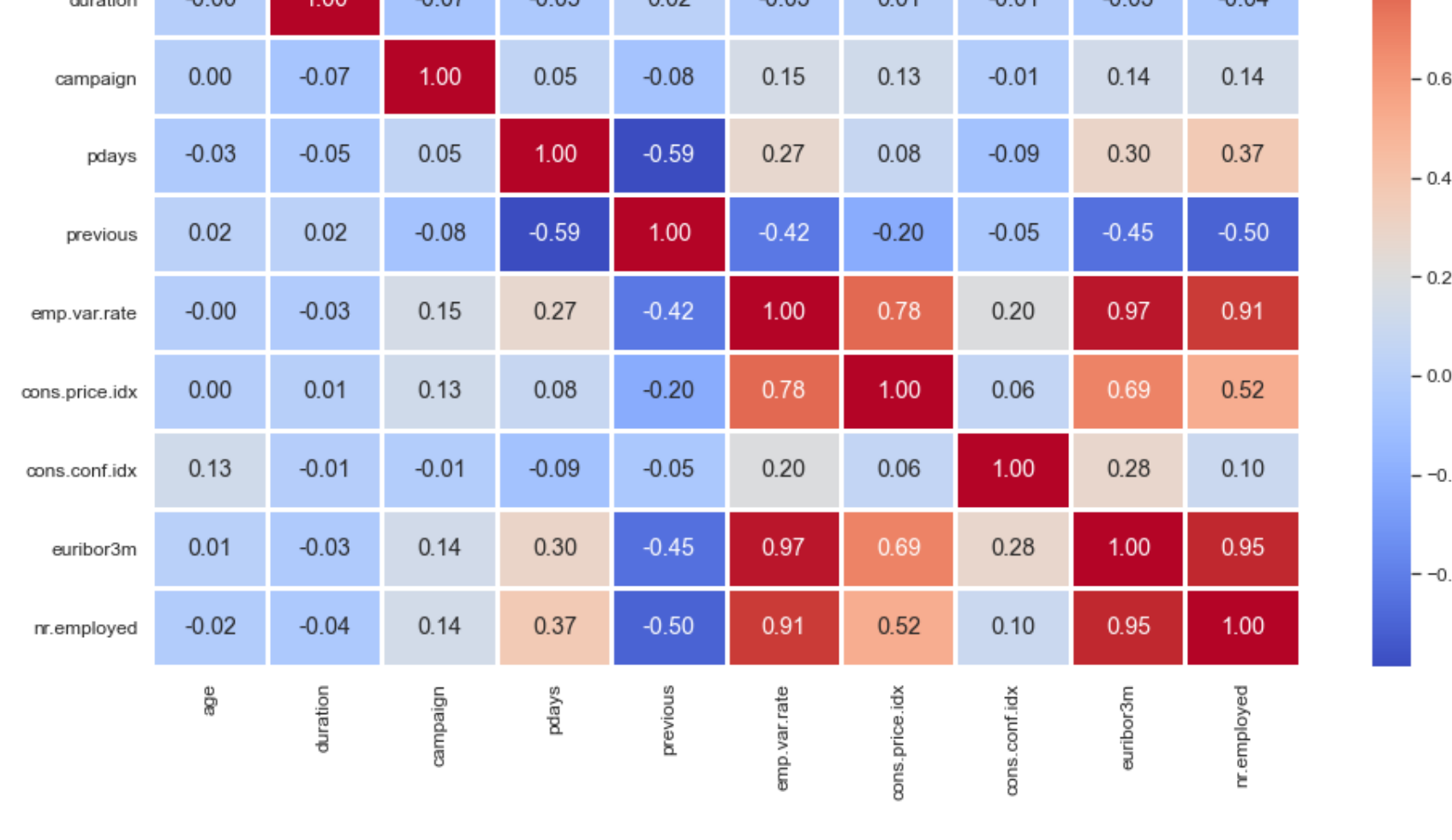
plt.subplot(7,2,5)
plt.title('loan')
sns.countplot(df.loan, hue=df.y)

plt.subplot(7,2,6)
plt.title('contact')
sns.countplot(df.contact, hue=df.y)

plt.subplot(7,2,7)
plt.title('month')
sns.countplot(df.month, hue=df.y)

plt.subplot(7,2,8)
plt.title('day of week')
sns.countplot(df.day_of_week, hue=df.y)

plt.tight_layout()
plt.show()
```



```
In [17]: df.corr()

Out[17]:
```

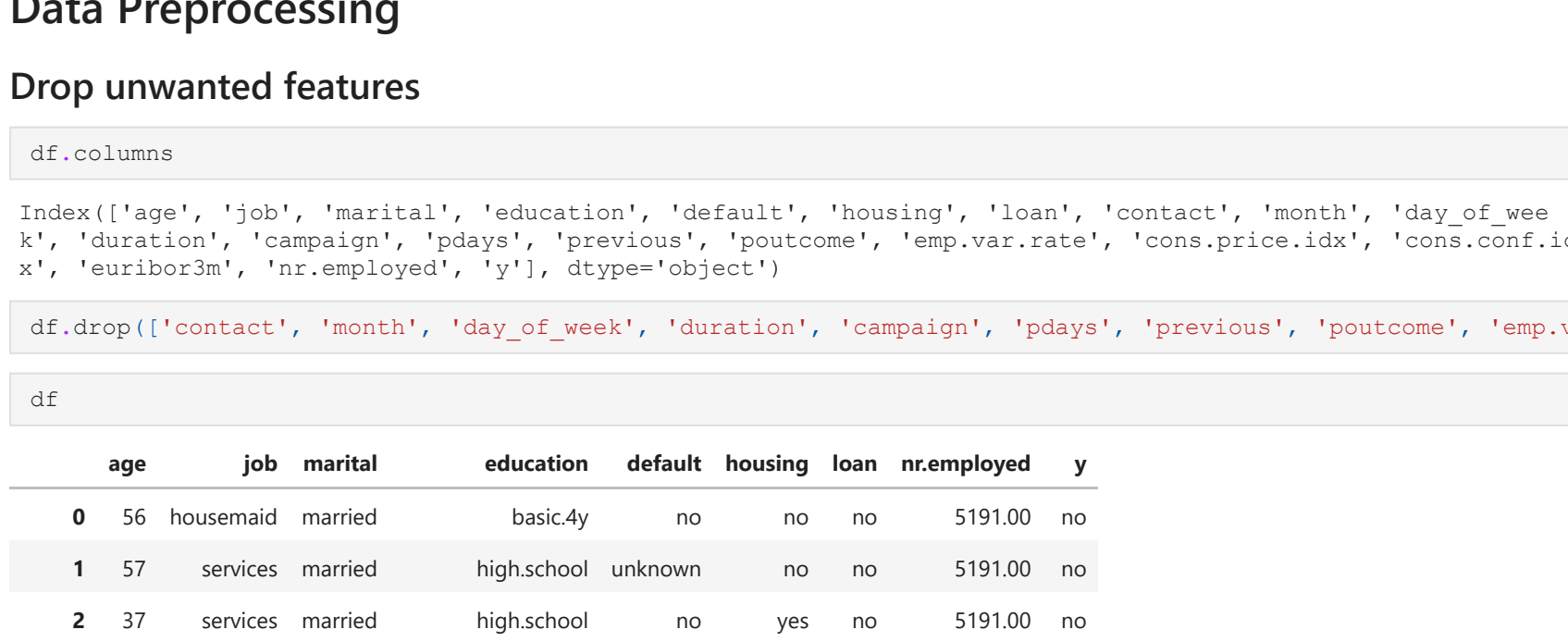
	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m	nr.employed
age	1.00	-0.00	0.00	-0.03	0.02	-0.00	0.00	0.13	0.01	-0.02
duration	-0.00	1.00	-0.07	-0.05	0.02	-0.03	0.01	-0.01	-0.03	-0.04
campaign	0.00	-0.07	1.00	0.05	-0.08	0.15	0.13	-0.01	0.14	0.14
pdays	0.03	-0.05	0.05	1.00	-0.59	0.27	0.08	-0.09	0.30	0.37
previous	0.02	0.02	-0.08	-0.59	1.00	-0.42	-0.20	-0.05	-0.45	-0.50
emp.var.rate	-0.00	-0.03	0.15	0.27	-0.42	1.00	0.78	0.20	0.67	0.91
cons.price.idx	0.00	0.01	0.13	0.08	-0.20	0.78	1.00	0.06	0.69	0.52
cons.conf.idx	0.13	-0.01	-0.01	-0.09	-0.05	0.20	0.06	1.00	0.28	0.10
euribor3m	0.01	-0.03	0.14	0.30	-0.45	0.97	0.69	0.28	1.00	0.95
nr.employed	-0.02	-0.04	0.14	0.37	-0.50	0.91	0.52	0.10	0.95	1.00

```
In [18]: plt.figure(figsize=(16,9))
sns.heatmap(df.corr(), cmap='coolwarm', annot=True, fmt='.2f', linewidth=2)
plt.title('Correlation Heatmap')
plt.show()
```



Pairplots

```
In [19]: sns.pairplot(df.sample(500))
plt.suptitle('Pairplots of features', x=0.5, y=1.02, ha='center', fontsize=20)
plt.show()
```



Data Preprocessing

Drop unwanted features

```
In [20]: df.columns

Out[20]: Index(['age', 'job', 'marital', 'education', 'default', 'housing', 'loan', 'nr.employed', 'y'], dtype='object')
```

```
In [21]: df.drop(['contact', 'month', 'day_of_week', 'duration', 'campaign', 'pdays', 'previous', 'poutcome', 'emp.va'], axis=1, inplace=True)

In [22]: df
```

```
Out[22]:
```

	age	job	marital	education	default	housing	loan	nr.employed	y
0	56	housemaid	married	basic.4y	no	no	no	5191.00	no
1	57	services	married	high.school	unknown	no	no	5191.00	no
2	37	services	married	high.school	no	yes	no	5191.00	no
3	40	admin.	married	basic.6y	no	no	no	5191.00	no
4	56	services	married	high.school	no	no	yes	5191.00	no
...
41183	73	retired	married	professional.course	no	yes	no	4963.60	yes
41184	46	blue-collar	married	professional.course	no	no	no	4963.60	no
41185	56	retired	married	university.degree	no	yes	no	4963.60	no
41186	44	technician	married	professional.course	no	no	no	4963.60	yes
41187	74	retired	married	professional.course	no	yes	no	4963.60	no

23638 rows x 9 columns

Treat Missing Values

```
In [23]: df.isnull().sum()

Out[23]: age      0
job      0
marital   0
education 0
default   0
housing   0
loan      0
nr.employed 0
y         0
dtype: int64
```

Treat Duplicate Values

```
In [24]: df.duplicated(keep='first').sum()

Out[24]: 17550
```

```
In [25]: df[df.duplicated(keep=False)] #Check duplicate values

Out[25]:
```

	age	job	marital	education	default	housing	loan	nr.employed	y
0	56	housemaid	married	basic.4y	no	no	no	5191.00	no
2	37	services	married	high.school	no	yes	no	5191.00	no
6	59	admin.	married	professional.course	no	no	no	5191.00	no
7	41	blue-collar	married	unknown	unknown	no	no	5191.00	no
9	25	services	single	high.school	no	yes	no	5191.00	no
...
41163	35	technician	divorced	basic.4y	no	yes	no	4963.60	yes
41172	31	admin.	single	university.degree	no	yes	no	4963.60	yes
41173	62	retired	married	university.degree	no	yes	no	4963.60	yes
41174	62	retired	married	university.degree	no	yes	no	4963.60	yes
41181	37	admin.	married	university.degree	no	yes	no	4963.60	yes

24701 rows x 9 columns

```
In [26]: df.drop_duplicates(inplace=True)

In [27]: df
```

```
Out[27]:
```

	age	job	marital	education	default	housing	loan	nr.employed	y
0	56	housemaid	married	basic.4y	no	no	no	5191.00	no
1	57	services	married	high.school	unknown	no	no	5191.00	no
2	37	services	married	high.school	no	yes	no	5191.00	no
3	40	admin.	married	basic.6y	no	no	no	5191.00	no
4	56	services	married	high.school	no	no	yes	5191.00	no
...
41183	73	retired	married	professional.course	no	yes	no	4963.60	yes
23634	46	blue-collar	married	professional.course	no	no	no	4963.60	no
23635	56	retired	married	university.degree	no	yes	no	4963.60	no
23636	44	technician							

