

Machine Learning Classification with Python for Direct Marketing *



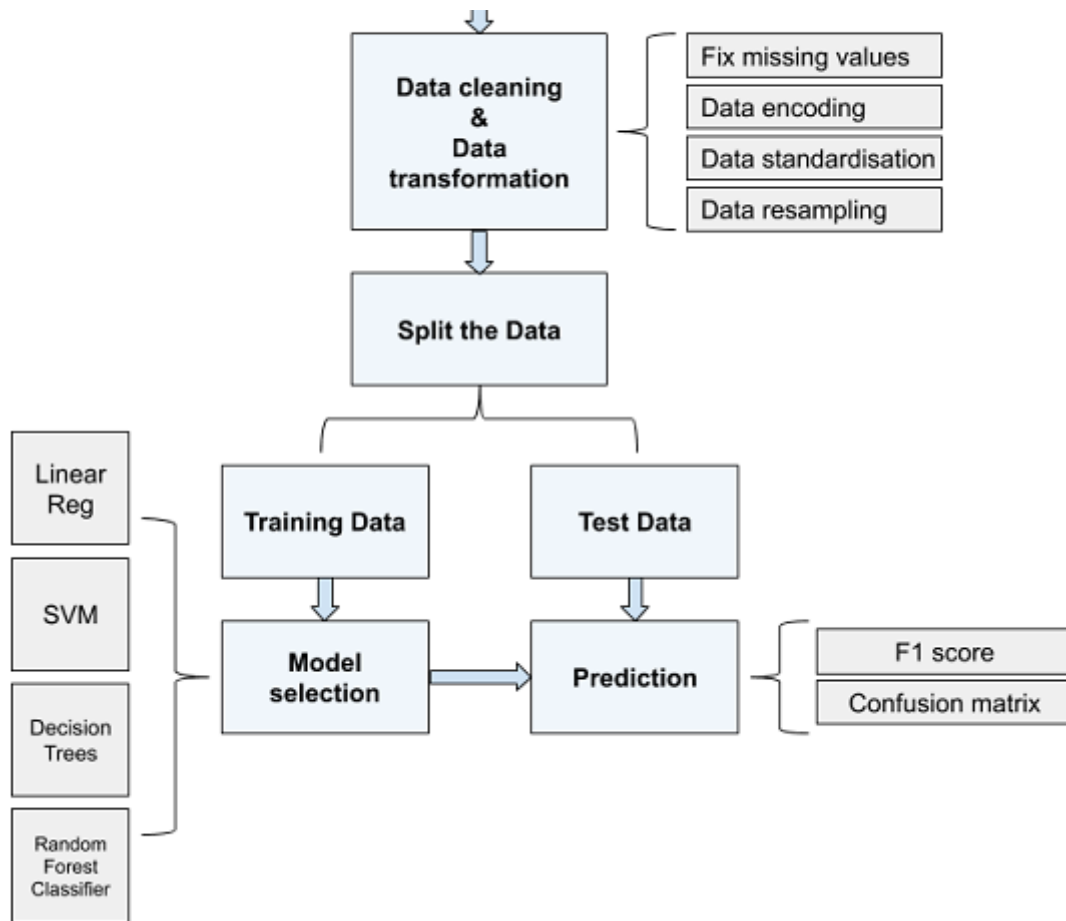
Sergey Medvedev

Jun 12 · 5 min read ★

How do you make business more time-efficient, slash costs and drive up sales? The question is timeless but not rhetorical. In the next few minutes of your reading time, I will apply a few classification algorithms to demonstrate how the use of the data analytic approach can contribute to that end. Together we'll create a predictive model that will help us customise the client databases we hand over to the telemarketing team so that they could concentrate resources on more promising clients first.

On course to that, we'll perform a number of actions on the dataset. First, we'll clean, encode, standardize and resample the data. Once the data is ready, we'll try four different classifiers on the training subset, make predictions and visualise them with a confusion matrix, and compute F1 score to elect the best model. These steps have been put together in the schema:





Project's schema

The dataset we'll be using here is not new to the town and you have probably come across it before. The data sample of 41,118 records was collected by a Portuguese bank between 2008 and 2013 and contains the results of a telemarketing campaign including customer's response to the bank's offer of a deposit contract (the binary target variable 'y'). That response is exactly what we are going to predict with the model. The dataset is available at Irvine's Machine Learning Repository of the University of California. So let's get started!

```

1 import pandas as pd
2 df = pd.read_csv('dataset_bank.csv', delimiter=';', decimal=',')
3 df.head()

```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

Outcome:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pdays	previous	poutcome	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m	nr.employed	y
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	mon	261	1	999	0	nonexistent	1.1	93.994	-36.4	4.857	5191	no
1	57	services	married	high school	unknown	no	no	telephone	may	mon	149	1	999	0	nonexistent	1.1	93.994	-36.4	4.857	5191	no
2	37	services	married	high school	no	yes	no	telephone	may	mon	226	1	999	0	nonexistent	1.1	93.994	-36.4	4.857	5191	no
3	40	admin.	married	basic.5y	no	no	no	telephone	may	mon	151	1	999	0	nonexistent	1.1	93.994	-36.4	4.857	5191	no
4	56	services	married	high school	no	no	yes	telephone	may	mon	307	1	999	0	nonexistent	1.1	93.994	-36.4	4.857	5191	no

Data Cleaning, Feature Selection, Feature Transformation

“I like to make a clean sweep of things”, Friedrich Nietzsche said. This is the data cleaning part!

With `df.isnull().sum()` query we make sure there are no missing values in the dataset (if there were any, `df.dropna(subset = ['feature_name'], inplace=True)` would drop them from the respected column).

In this example, I used Tableau Prep Builder data cleaning tool to trace and drop outliers, to make sure values in numerical features aren't strings, to rename some columns and to get rid of a few irrelevant ones, i.e. 'contact', 'month', 'day_of_week', 'duration', 'campaign', 'pdays', 'previous' and 'poutcome' (these columns describe a telephone call that has happened already, therefore they should not be used in our predictive model).

Next we shall transform the non-numerical labels of the categorical variables to numerical ones and convert them to integers. We do it like this:

```

1  from sklearn import preprocessing
2  num = preprocessing.LabelEncoder()
3
4  num.fit(["admin.", "blue-collar", "entrepreneur", "housemaid", "management",
5          "retired", "self-employed", "services", "student", "technician", "unemployed", "unknown"])
6  df['job'] = num.transform(df['job']).astype('int')
7
8  num.fit(["divorced", "married", "single", "unknown"])
9  df['marital'] = num.transform(df['marital']).astype('int')
10
11 num.fit(["basic.4y", "basic.6y", "basic.9y", "high.school", "illiterate", "professional.course", "university.degree"])
12 df['education'] = num.transform(df['education']).astype('int')
13
14 num.fit(["no", "yes", "unknown"])
15 df['housing_loan'] = num.transform(df['housing_loan']).astype('int')
16
17 num.fit(["no", "yes", "unknown"])
18 df['personal_loan'] = num.transform(df['personal_loan']).astype('int')
19
20 num.fit(["failure", "nonexistent", "success"])
21 df['poutcome'] = num.transform(df['poutcome']).astype('int')
22
23 num.fit(["yes", "no"])
24 df['y'] = num.transform(df['y']).astype('int')

```

It makes sense to run `df.dtypes` query just to ensure that the labels have turned to integers. We then apply `StandardScaler` from `sklearn.preprocessing` toolbox to standardize the numerical values of the other features that we expect to find use of in the model. The method standardizes features by removing the mean and scaling to unit variance:

```

1 from sklearn.preprocessing import StandardScaler
2 scaler=StandardScaler()
3
4 df['cons.price.idx'] = scaler.fit_transform(df[['cons.price.idx']].reshape(-1,1))
5 df['cons.conf.idx'] = scaler.fit_transform(df[['cons.conf.idx']].reshape(-1,1))
6 df['euribor3m'] = scaler.fit_transform(df[['euribor3m']].reshape(-1,1))

```

ML Classification with Python for Bank Telemarketing.py hosted with ❤️ by GitHub

view raw

Let's now rank the features of the dataset with recursive feature elimination (RFE) method and Random Forest Classifier algorithm as its estimator:

```

1 import numpy as np
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.feature_selection import RFE
4
5 X = np.asarray(df[['age', 'job', 'marital', 'education', 'housing_loan', 'personal_loan',
6                 'emp.var.rate', 'cons.price.idx', 'cons.conf.idx', 'euribor3m']])
7 y = np.asarray(df['y'])
8 rfc = RandomForestClassifier(n_estimators=40)
9 rfe = RFE(rfc, 6)
10 rfe_fit = rfe.fit(X, y)
11
12 print("Num Features: %s" % (rfe_fit.n_features_))
13 print("Selected Features: %s" % (rfe_fit.support_))
14 print("Feature Ranking: %s" % (rfe_fit.ranking_))

```

ML Classification with Python for Bank Telemarketing.py hosted with ❤️ by GitHub

view raw

Output:

```

Num Features: 6
Selected Features: [ True  True False  True False False  True False
 True  True]
Feature Ranking: [1 1 3 1 2 5 1 4 1 1]

```

At a later stage, when we'll be building a predictive model, we will make use of this feature ranking. We will play with the model trying to find an optimal combination of the highest-ranked features that would make prediction with satisfactory F1 score. Running ahead, that optimal combination will be:

```
1 X = np.asarray(df[['age', 'job', 'marital', 'education', 'housing_loan',
2                 'emp.var.rate', 'cons.conf.idx', 'euribor3m']])
3 y = np.asarray(df['y'])
```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

Class imbalance is the problem that often comes along with such classification cases as fraudulent credit card transactions or the results of online campaigns, for instance. After executing `df['y'].value_counts()` query we see that the two classes of the variable 'y' are not represented equally in our dataset also. After data cleaning there are 35584 records belonging to the class '0' and only 4517 records of the class '1' in the target variable 'y'. Prior to splitting the data into the training and testing samples, we should think of oversampling or undersampling the data.

To resample the data, let's apply `SMOTE` method for oversampling from

`imblearn.over_sampling` toolbox (for this step you may need to install `imblearn` package with Pip or Conda first):

```
1 from imblearn.over_sampling import SMOTE
2
3 sm=SMOTE(ratio='auto', kind='regular')
4 X_sampled,y_sampled=sm.fit_sample(X,y)
```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

It is as simple as that. Now the data is in balance with 35584 entries in each class:

```
1 Sampled_no = len(y_sampled[y_sampled==0])
2 Sampled_yes = len(y_sampled[y_sampled==1])
3 print([Sampled_no],[Sampled_yes])
```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

Output:

```
[35584] [35584]
```

Building a predictive model

Now that the data has been prepared, we are ready to train our model and make predictions. Let's first split the data into the training and testing sets:

```
1 from sklearn.model_selection import train_test_split
2 X_train,X_test,y_train,y_test = train_test_split(X_sampled,y_sampled,test_size=0.3,random_state=
```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

We will try four classification algorithms, i.e. Logistic Regression, Support Vector Machine, Decision Trees, and Random Forest, and then compute their F1 scores using a user-defined `scorer` function to choose the classifier with the highest score:

```
1 from sklearn.metrics import f1_score
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.svm import SVC
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.ensemble import RandomForestClassifier
6
7 lr = LogisticRegression(C=1, solver='lbfgs')
8 clf = SVC(kernel='rbf', gamma='auto')
9 dtree = DecisionTreeClassifier(criterion="entropy", max_depth=4)
10 rfc = RandomForestClassifier(n_estimators=40)
11
12 def scorer(i,j,k,l):
13     for every in (i,j,k,l):
14         every.fit(X_train,y_train)
15         print (every.__class__.__name__, 'F1 score = ', f1_score(y_test,every.predict(X_test)))
16 scorer (lr,clf,dtree,rfc)
```

Output:

```
LogisticRegression F1 score = 0.71245481432799659
SVC F1 score = 0.753674137005029
DecisionTreeClassifier F1 score = 0.7013983920155255
RandomForestClassifier F1 score = 0.923286257213907
```

F1 score is the weighted average of the precision and recall. You can read about how to interpret the precision and recall scores in my post [here](#).

Now, let's print a full classification report with the precision and recall for the Random Forest algorithm which has demonstrated the highest F1 score:

```
1 from sklearn.metrics import classification_report
2 yhat = rfc.predict(X_test)
3 print(classification_report(y_test,yhat))
```

ML Classification with Python for Bank Telemarketing.py hosted with ❤ by GitHub

[view raw](#)

Output:

	precision	recall	f1-score	support
0	0.91	0.93	0.92	10594
1	0.93	0.91	0.92	10757
micro avg	0.92	0.92	0.92	21351
macro avg	0.92	0.92	0.92	21351
weighted avg	0.92	0.92	0.92	21351

Finally, we can visualize the result with a confusion matrix:

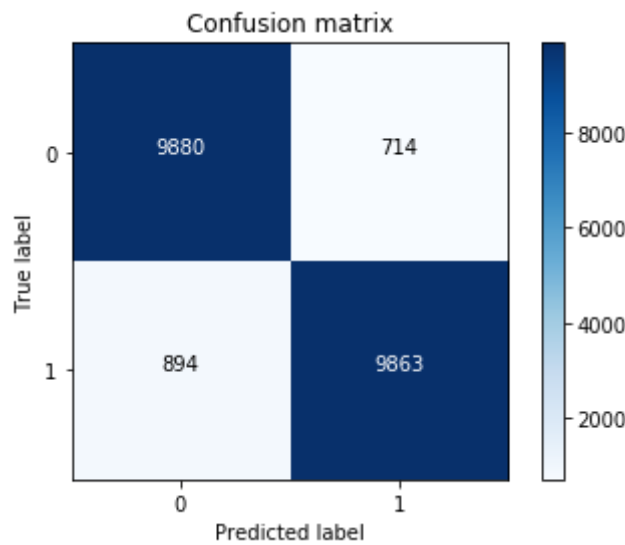
```
1 from sklearn.metrics import confusion_matrix
2 import itertools
3 import matplotlib.pyplot as plt
4
5 def plot_confusion_matrix(cm, classes,
6                           normalize=False,
7                           title='Confusion matrix',
8                           cmap=plt.cm.Blues):
9
10     plt.imshow(cm, interpolation='nearest', cmap=cmap)
11     plt.title(title)
12     plt.colorbar()
13     tick_marks = np.arange(len(classes))
14     plt.xticks(tick_marks, classes)
15     plt.yticks(tick_marks, classes)
16
17     fmt = '.2f' if normalize else 'd'
18     thresh = cm.max() / 2.
19     for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
```

```
20 plt.text(j, i, format(cm[i, j], fmt),
21         horizontalalignment="center",
22         color="white" if cm[i, j] > thresh else "black")
23
24 plt.ylabel('True label')
25 plt.xlabel('Predicted label')
26
27 plot_confusion_matrix(confusion_matrix(y_test, yhat), classes=['0','1'],normalize= False, titl
```

ML Classification with Python for Bank Telemarketing by hosted with ❤ by GitHub

[view raw](#)

Outcome:



Great! We've cleaned and transformed the data, selected the most relevant features, elected the best model and made a prediction with a decent score. Now we have a model that should help us customise the client databases we hand over to the telemarketing team so that they could center their efforts on those better positioned to react in the affirmative to the campaign first.

Thank you for reading !!

[Data Science](#)

[Data Analysis](#)

[Machine Learning](#)

[Classification](#)

[Python](#)

[About](#)

[Help](#)

[Legal](#)