Images haven't loaded yet. Please exit printing, wait for images to load, and try to print again.

Sep 29 · 9 min read

# Building A Logistic Regression in Python, Step by Step
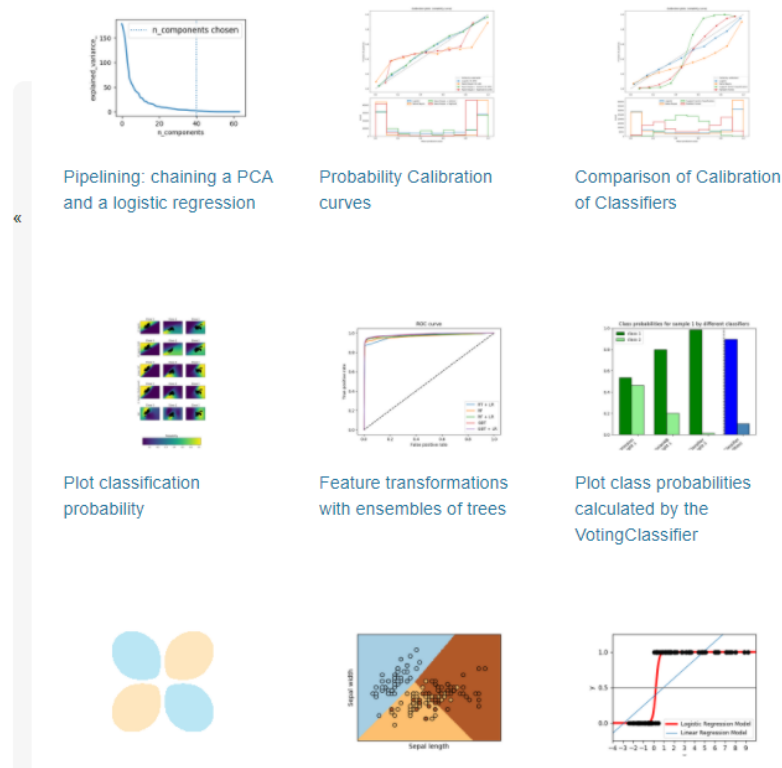


Pipelining: chaining a PCA and a logistic regression

Probability Calibration curves

Comparison of Calibration of Classifiers

Plot classification probability

Feature transformations with ensembles of trees

Plot class probabilities calculated by the VotingClassifier

Photo Credit: Scikit-Learn

( This article first appeared on *Datascience+* )

## Introduction

Logistic Regression is a Machine Learning classification algorithm that is used to predict the probability of a categorical dependent variable. In logistic regression, the dependent variable is a binary variable that contains data coded as 1 (yes, success, etc.) or 0 (no, failure, etc.). In other words, the logistic regression model predicts $P(Y=1)$ as a function of X.

## Logistic Regression Assumptions

- Binary logistic regression requires the dependent variable to be binary.

- For a binary regression, the factor level 1 of the dependent variable should represent the desired outcome.

- Only the meaningful variables should be included.

- The independent variables should be independent of each other. That is, the model should have little or no multicollinearity.

- The independent variables are linearly related to the log odds.

- Logistic regression requires quite large sample sizes.

Keeping the above assumptions in mind, let's look at our dataset.

## Data

The dataset comes from the UCI Machine Learning repository, and it is related to direct marketing campaigns (phone calls) of a Portuguese banking institution. The classification goal is to predict whether the client will subscribe (1/0) to a term deposit (variable y). The dataset can be downloaded from here.

```
import pandas as pd
import numpy as np
from sklearn import preprocessing
import matplotlib.pyplot as plt
plt.rc("font", size=14)
from sklearn.linear_model import LogisticRegression
from sklearn.cross_validation import train_test_split
import seaborn as sns
sns.set(style="white")
sns.set(style="whitegrid", color_codes=True)
```

The dataset provides the bank customers' information. It includes 41,188 records and 21 fields.

```
In [39]:  data = pd.read_csv('bank.csv', header=0)
          data = data.dropna()
          print(data.shape)
          print(list(data.columns))

          (41188, 21)
          ['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']
```

```
In [37]:  data.head()
```

Out[37]:

|   | age | job | marital | education | default | housing | loan | contact | month | day_of_week | ... | campaign | pdays | previous | poutcome | emp_var_rate |
|---|-----|-----|---------|-----------|---------|---------|------|---------|-------|-------------|-----|----------|-------|----------|----------|--------------|
| 0 | 44 | blue-collar | married | basic.4y | unknown | yes | no | cellular | aug | thu | ... | 1 | 999 | 0 | nonexistent | 1.4 |
| 1 | 53 | technician | married | unknown | no | no | no | cellular | nov | fri | ... | 1 | 999 | 0 | nonexistent | -0.1 |
| 2 | 28 | management | single | university.degree | no | yes | no | cellular | jun | thu | ... | 3 | 6 | 2 | success | -1.7 |
| 3 | 39 | services | married | high.school | no | no | no | cellular | apr | fri | ... | 2 | 999 | 0 | nonexistent | -1.8 |
| 4 | 55 | retired | married | basic.4y | no | yes | no | cellular | aug | fri | ... | 1 | 3 | 1 | success | -2.9 |

5 rows × 21 columns

Figure 1

**Input variables**

1. age (numeric)

2. job : type of job (categorical: "admin", "blue-collar", "entrepreneur", "housemaid", "management", "retired", "self-employed", "services", "student", "technician", "unemployed", "unknown")

3. marital : marital status (categorical: "divorced", "married", "single", "unknown")

4. education (categorical: "basic.4y", "basic.6y", "basic.9y", "high.school", "illiterate", "professional.course", "university.degree", "unknown")

5. default: has credit in default? (categorical: "no", "yes", "unknown")

6. housing: has housing loan? (categorical: "no", "yes", "unknown")

7. loan: has personal loan? (categorical: "no", "yes", "unknown")

8. contact: contact communication type (categorical: "cellular", "telephone")

9. month: last contact month of year (categorical: "jan", "feb", "mar", …, "nov", "dec")

10. day_of_week: last contact day of the week (categorical: "mon", "tue", "wed", "thu", "fri")

11. duration: last contact duration, in seconds (numeric). Important note: this attribute highly affects the output target (e.g., if duration=0 then y='no'). 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

12. campaign: number of contacts performed during this campaign and for this client (numeric, includes last contact)

13. 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)

14. previous: number of contacts performed before this campaign and for this client (numeric)

15. poutcome: outcome of the previous marketing campaign (categorical: "failure", "nonexistent", "success")

16. emp.var.rate: employment variation rate—(numeric)

17. cons.price.idx: consumer price index—(numeric)

18. cons.conf.idx: consumer confidence index—(numeric)

19. euribor3m: euribor 3 month rate—(numeric)

20. nr.employed: number of employees—(numeric)

**Predict variable (desired target):**

y—has the client subscribed a term deposit? (binary: "1", means "Yes", "0" means "No")

```
In [4]: data['education'].unique()

Out[4]: array(['basic.4y', 'unknown', 'university.degree', 'high.school',
                'basic.9y', 'professional.course', 'basic.6y', 'illiterate
```

Figure 2

Let us group "basic.4y", "basic.9y" and "basic.6y" together and call them "basic".

```
data['education']=np.where(data['education'] =='basic.9y',
'Basic', data['education'])
data['education']=np.where(data['education'] =='basic.6y',
'Basic', data['education'])
```

```python
data['education']=np.where(data['education'] =='basic.4y',
'Basic', data['education'])
```

After grouping, this is the columns:

```
In [6]: data['education'].unique()

Out[6]: array(['Basic', 'unknown', 'university.degree', 'high.school',
               'professional.course', 'illiterate'], dtype=object)
```

Figure 3

# Data exploration

```
In [7]: data['y'].value_counts()

Out[7]: 0    36548
        1     4640
        Name: y, dtype: int64
```

```python
In [17]: sns.countplot(x='y',data=data, palette='hls')
         plt.show()
         plt.savefig('count_plot')
```
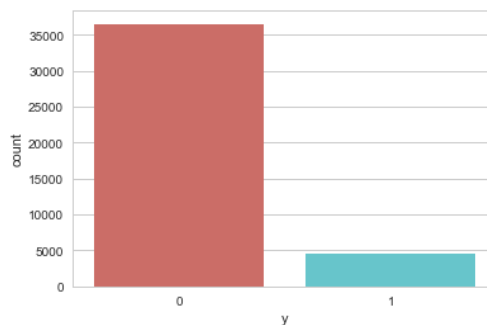
Figure 4

There are 36548 no's and 4640 yes's in the outcome variables.

Let's get a sense of the numbers across the two classes.

```
In [9]: data.groupby('y').mean()
```

| y | age | duration | campaign | pdays | previous | emp_var_rate | cons_price_idx | cons_conf_idx | euribor3m | nr_employed |
|---|-----|----------|----------|-------|----------|--------------|----------------|---------------|-----------|-------------|
| 0 | 39.911185 | 220.844807 | 2.633085 | 984.113878 | 0.132374 | 0.248875 | 93.603757 | -40.593097 | 3.811491 | 5176.166600 |
| 1 | 40.913147 | 553.191164 | 2.051724 | 792.035560 | 0.492672 | -1.233448 | 93.354386 | -39.789784 | 2.123135 | 5095.115991 |

Figure 5

*Observations*:

- The average age of customers who bought the term deposit is higher than that of the customers who didn't.

- The pdays (days since the customer was last contacted) is understandably lower for the customers who bought it. The lower the pdays, the better the memory of the last call and hence the better chances of a sale.

- Surprisingly, campaigns (number of contacts or calls made during the current campaign) are lower for customers who bought the term deposit.

We can calculate categorical means for other categorical variables such as education and marital status to get a more detailed sense of our data.

```
In [10]:  data.groupby('job').mean()
```
Out[10]:

| job | age | duration | campaign | pdays | previous | emp_var_rate | cons_price_idx | cons_conf_idx | euribor3m | nr_employed | y |
|---|---|---|---|---|---|---|---|---|---|---|---|
| admin. | 38.187296 | 254.312128 | 2.623489 | 954.319229 | 0.189023 | 0.015563 | 93.534054 | -40.245433 | 3.550274 | 5164.125350 | 0.129726 |
| blue-collar | 39.555760 | 264.542360 | 2.558461 | 985.160363 | 0.122542 | 0.248995 | 93.656656 | -41.375816 | 3.771996 | 5175.615150 | 0.068943 |
| entrepreneur | 41.723214 | 263.267857 | 2.535714 | 981.267170 | 0.138736 | 0.158723 | 93.605372 | -41.283654 | 3.791120 | 5176.313530 | 0.085165 |
| housemaid | 45.500000 | 250.454717 | 2.639623 | 960.579245 | 0.137736 | 0.433396 | 93.676576 | -39.495283 | 4.009645 | 5179.529623 | 0.100000 |
| management | 42.362859 | 257.058140 | 2.476060 | 962.647059 | 0.185021 | -0.012688 | 93.522755 | -40.489466 | 3.611316 | 5166.650513 | 0.112175 |
| retired | 62.027326 | 273.712209 | 2.476744 | 897.936047 | 0.327326 | -0.698314 | 93.430786 | -38.573081 | 2.770066 | 5122.262151 | 0.252326 |
| self-employed | 39.949331 | 264.142153 | 2.660802 | 976.621393 | 0.143561 | 0.094159 | 93.559982 | -40.488107 | 3.689376 | 5170.674384 | 0.104856 |
| services | 37.926430 | 258.398085 | 2.587805 | 979.974049 | 0.154951 | 0.175359 | 93.634659 | -41.290048 | 3.699187 | 5171.600126 | 0.081381 |
| student | 25.894857 | 283.683429 | 2.104000 | 840.217143 | 0.524571 | -1.408000 | 93.331613 | -40.187543 | 1.884224 | 5085.939086 | 0.314286 |
| technician | 38.507638 | 250.232241 | 2.577339 | 964.408127 | 0.153789 | 0.274566 | 93.561471 | -39.927569 | 3.820401 | 5175.648391 | 0.108260 |
| unemployed | 39.733728 | 249.451677 | 2.564103 | 935.316568 | 0.199211 | -0.111736 | 93.563781 | -40.007594 | 3.466583 | 5157.156509 | 0.142012 |
| unknown | 45.563636 | 239.675758 | 2.648485 | 938.727273 | 0.154545 | 0.357879 | 93.718942 | -38.797879 | 3.949033 | 5172.931818 | 0.112121 |

Figure 6

In [11]:  `data.groupby('marital').mean()`

Out[11]:

| | age | duration | campaign | pdays | previous | emp_var_rate | cons_price_idx | cons_conf_idx | euribor3m | nr_employed | y |
|---|---|---|---|---|---|---|---|---|---|---|---|
| marital | | | | | | | | | | | |
| divorced | 44.899393 | 253.790330 | 2.61340 | 968.639853 | 0.168690 | 0.163985 | 93.606563 | -40.707069 | 3.715603 | 5170.878643 | 0.103209 |
| married | 42.307165 | 257.438623 | 2.57281 | 967.247673 | 0.155608 | 0.183625 | 93.597367 | -40.270659 | 3.745832 | 5171.848772 | 0.101573 |
| single | 33.158714 | 261.524378 | 2.53380 | 949.909578 | 0.211359 | -0.167989 | 93.517300 | -40.918698 | 3.317447 | 5155.199265 | 0.140041 |
| unknown | 40.275000 | 312.725000 | 3.18750 | 937.100000 | 0.275000 | -0.221250 | 93.471250 | -40.820000 | 3.313038 | 5157.393750 | 0.150000 |

In [12]:  `data.groupby('education').mean()`

Out[12]:

| | age | duration | campaign | pdays | previous | emp_var_rate | cons_price_idx | cons_conf_idx | euribor3m | nr_employed | y |
|---|---|---|---|---|---|---|---|---|---|---|---|
| education | | | | | | | | | | | |
| Basic | 42.163910 | 263.043874 | 2.559498 | 974.877967 | 0.141053 | 0.191329 | 93.639933 | -40.927595 | 3.729654 | 5172.014113 | 0.087029 |
| high.school | 37.998213 | 260.886810 | 2.568576 | 964.358382 | 0.185917 | 0.032937 | 93.584857 | -40.940641 | 3.556157 | 5164.994735 | 0.108355 |
| illiterate | 48.500000 | 276.777778 | 2.277778 | 943.833333 | 0.111111 | -0.133333 | 93.317333 | -39.950000 | 3.516556 | 5171.777778 | 0.222222 |
| professional.course | 40.080107 | 252.533855 | 2.586115 | 960.765974 | 0.163075 | 0.173012 | 93.569864 | -40.124108 | 3.710457 | 5170.155979 | 0.113485 |
| university.degree | 38.879191 | 253.223373 | 2.563527 | 951.807692 | 0.192390 | -0.028090 | 93.493466 | -39.975805 | 3.529663 | 5163.226298 | 0.137245 |
| unknown | 43.481225 | 262.390526 | 2.596187 | 942.830734 | 0.226459 | 0.059099 | 93.658615 | -39.877816 | 3.571098 | 5159.549509 | 0.145003 |

Figure 7

# Visualizations

```
%matplotlib inline
pd.crosstab(data.job,data.y).plot(kind='bar')
plt.title('Purchase Frequency for Job Title')
plt.xlabel('Job')
plt.ylabel('Frequency of Purchase')
plt.savefig('purchase_fre_job')
```
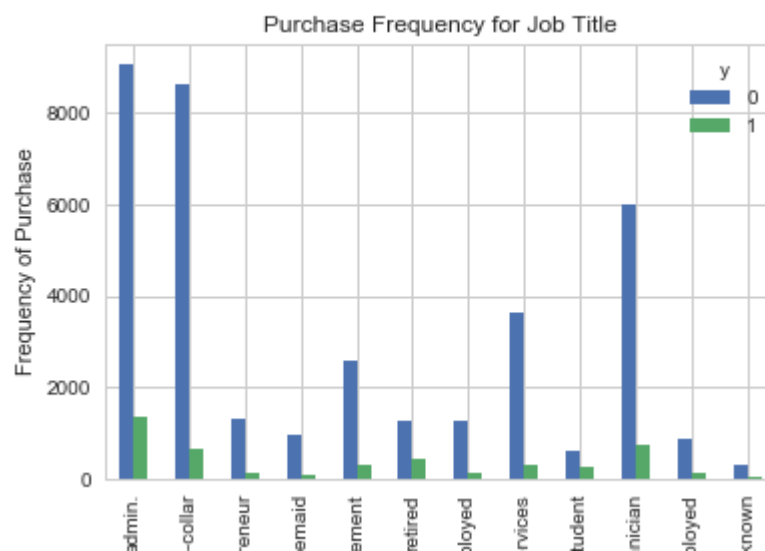


Figure 8

The frequency of purchase of the deposit depends a great deal on the job title. Thus, the job title can be a good predictor of the outcome variable.

```
table=pd.crosstab(data.marital,data.y)
table.div(table.sum(1).astype(float),
axis=0).plot(kind='bar', stacked=True)
plt.title('Stacked Bar Chart of Marital Status vs Purchase')
plt.xlabel('Marital Status')
plt.ylabel('Proportion of Customers')
plt.savefig('mariral_vs_pur_stack')
```
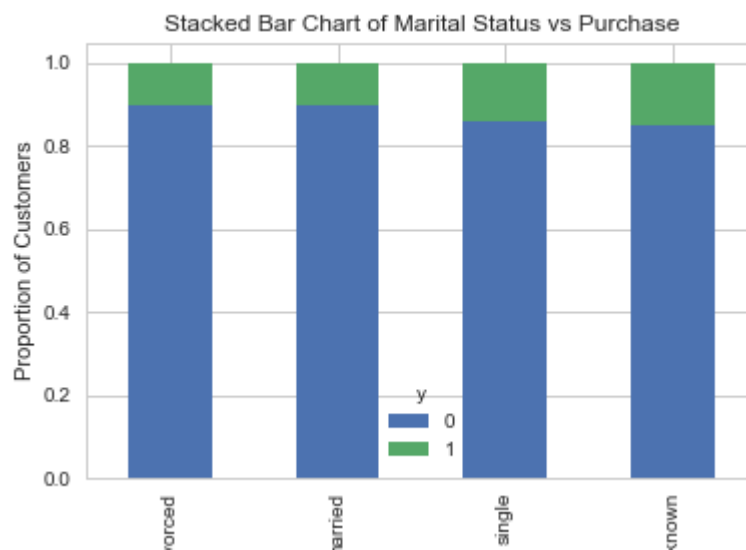


Figure 9

The marital status does not seem a strong predictor for the outcome variable.

```
table=pd.crosstab(data.education,data.y)
table.div(table.sum(1).astype(float),
axis=0).plot(kind='bar', stacked=True)
plt.title('Stacked Bar Chart of Education vs Purchase')
plt.xlabel('Education')
plt.ylabel('Proportion of Customers')
plt.savefig('edu_vs_pur_stack')
```
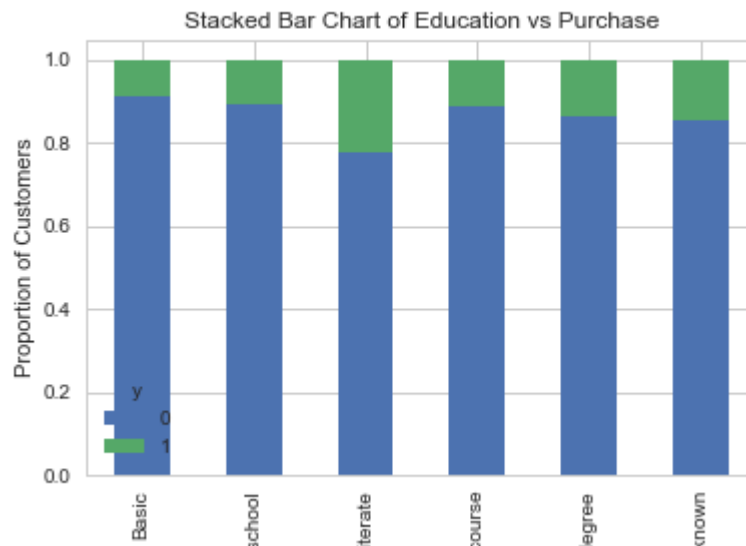
Figure 10

Education seems a good predictor of the outcome variable.

```
pd.crosstab(data.day_of_week,data.y).plot(kind='bar')
plt.title('Purchase Frequency for Day of Week')
plt.xlabel('Day of Week')
plt.ylabel('Frequency of Purchase')
plt.savefig('pur_dayofweek_bar')
```
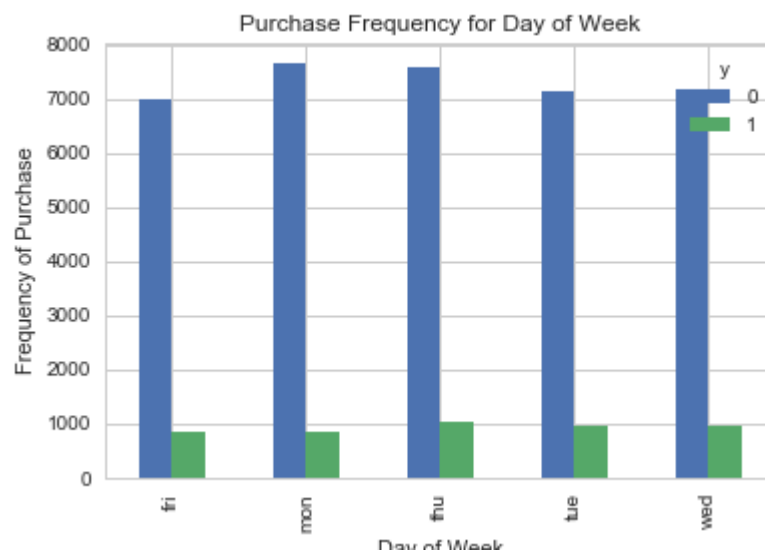


Figure 11

Day of week may not be a good predictor of the outcome.

```
pd.crosstab(data.month,data.y).plot(kind='bar')
plt.title('Purchase Frequency for Month')
plt.xlabel('Month')
plt.ylabel('Frequency of Purchase')
plt.savefig('pur_fre_month_bar')
```
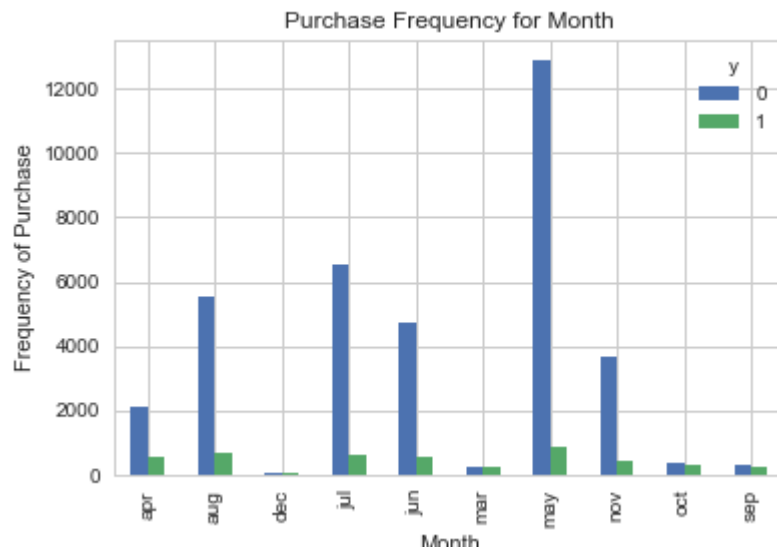


Figure 12

Month might be a good predictor of the outcome variable.

```
data.age.hist()
plt.title('Histogram of Age')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.savefig('hist_age')
```
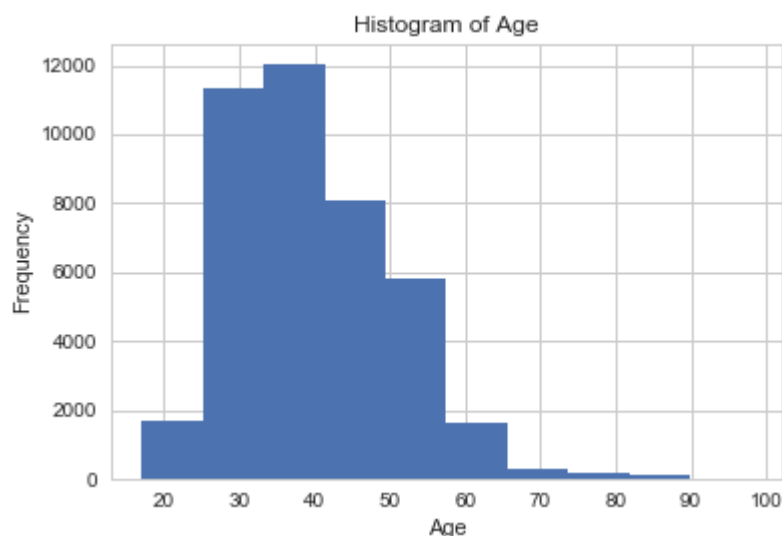
Figure 13

Most of the customers of the bank in this dataset are in the age range of 30–40.

```
pd.crosstab(data.poutcome,data.y).plot(kind='bar')
plt.title('Purchase Frequency for Poutcome')
plt.xlabel('Poutcome')
plt.ylabel('Frequency of Purchase')
plt.savefig('pur_fre_pout_bar')
```
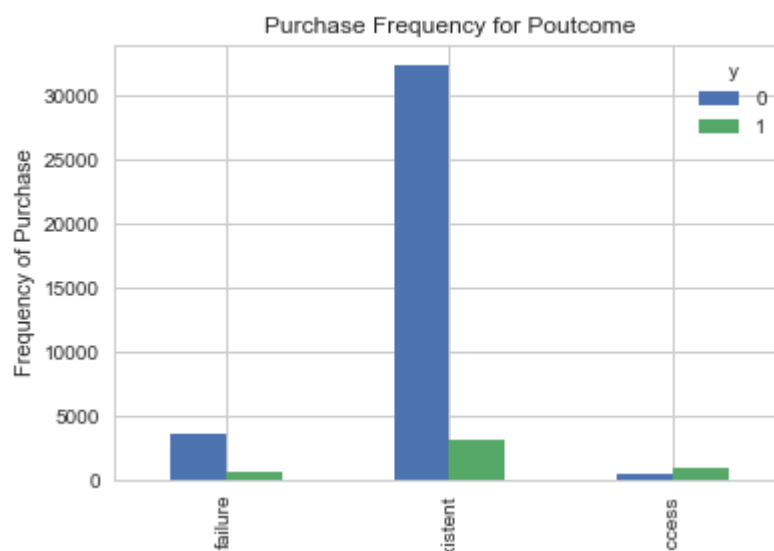


Figure 14

Poutcome seems to be a good predictor of the outcome variable.

# Create dummy variables

That is variables with only two values, zero and one.

```
cat_vars=
['job','marital','education','default','housing','loan','con
tact','month','day_of_week','poutcome']
for var in cat_vars:
    cat_list='var'+'_'+var
    cat_list = pd.get_dummies(data[var], prefix=var)
    data1=data.join(cat_list)
    data=data1
```

```
cat_vars=
['job','marital','education','default','housing','loan','con
tact','month','day_of_week','poutcome']
data_vars=data.columns.values.tolist()
to_keep=[i for i in data_vars if i not in cat_vars]
```

Our final data columns will be:

```
data_final=data[to_keep]
data_final.columns.values
```

```
array(['age', 'duration', 'campaign', 'pdays', 'previous', 'emp_var_rate',
       'cons_price_idx', 'cons_conf_idx', 'euribor3m', 'nr_employed', 'y',
       'job_admin.', 'job_blue-collar', 'job_entrepreneur',
       'job_housemaid', 'job_management', 'job_retired',
       'job_self-employed', 'job_services', 'job_student',
       'job_technician', 'job_unemployed', 'job_unknown',
       'marital_divorced', 'marital_married', 'marital_single',
       'marital_unknown', 'education_Basic', 'education_high.school',
       'education_illiterate', 'education_professional.course',
       'education_university.degree', 'education_unknown', 'default_no',
       'default_unknown', 'default_yes', 'housing_no', 'housing_unknown',
       'housing_yes', 'loan_no', 'loan_unknown', 'loan_yes',
       'contact_cellular', 'contact_telephone', 'month_apr', 'month_aug',
       'month_dec', 'month_jul', 'month_jun', 'month_mar', 'month_may',
       'month_nov', 'month_oct', 'month_sep', 'day_of_week_fri',
       'day_of_week_mon', 'day_of_week_thu', 'day_of_week_tue',
       'day_of_week_wed', 'poutcome_failure', 'poutcome_nonexistent',
       'poutcome_success'], dtype=object)
```

Figure 15

```
data_final_vars=data_final.columns.values.tolist()
y=['y']
X=[i for i in data_final_vars if i not in y]
```

# Feature Selection

Recursive Feature Elimination (RFE) is based on the idea to repeatedly construct a model and choose either the best or worst performing feature, setting the feature aside and then repeating the process with the rest of the features. This process is applied until all features in the dataset are exhausted. The goal of RFE is to select features by recursively considering smaller and smaller sets of features.

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


logreg = LogisticRegression()


rfe = RFE(logreg, 18)
rfe = rfe.fit(data_final[X], data_final[y] )
print(rfe.support_)
print(rfe.ranking_)
```

```
[False False False False  True False False False  True False False  True
 False False False  True False  True False False False False False False
 False False False False False False False False  True False False False
 False False False False False False  True  True  True False False False
  True  True  True False False False  True False False  True  True  True
  True]
[35 33 12 40  1 13 17 16  1 27 11  1 24 39 42  1 31  1  1 19 21 41  2  3  4
 43  6  7 38  8 10 15  1 14 44 36 29 37 20 30 28 23  1  1  1 18 22 25  1  1
  1 32  5  9  1 34 26  1  1  1  1]
```

Figure 16

The RFE has helped us select the following features: "previous", "euribor3m", "job_blue-collar", "job_retired", "job_services", "job_student", "default_no", "month_aug", "month_dec", "month_jul", "month_nov", "month_oct", "month_sep", "day_of_week_fri", "day_of_week_wed", "poutcome_failure", "poutcome_nonexistent", "poutcome_success".

```
cols=["previous", "euribor3m", "job_blue-collar",
"job_retired", "job_services", "job_student", "default_no",
      "month_aug", "month_dec", "month_jul", "month_nov",
"month_oct", "month_sep", "day_of_week_fri",
"day_of_week_wed",
      "poutcome_failure", "poutcome_nonexistent",
"poutcome_success"]
X=data_final[cols]
y=data_final['y']
```

## Implementing the model

```
import statsmodels.api as sm
logit_model=sm.Logit(y,X)
result=logit_model.fit()
print(result.summary())
```

```
                        Logit Regression Results
==============================================================================
Dep. Variable:                      y   No. Observations:                41188
Model:                          Logit   Df Residuals:                    41170
Method:                           MLE   Df Model:                           17
Date:                Sat, 18 Nov 2017   Pseudo R-squ.:                  0.1844
Time:                        02:47:55   Log-Likelihood:                -11826.
converged:                       True   LL-Null:                       -14499.
                                        LLR p-value:                     0.000
==============================================================================
                         coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
previous               0.2385      0.051      4.642      0.000       0.138       0.339
euribor3m             -0.4981      0.012    -40.386      0.000      -0.522      -0.474
job_blue-collar       -0.3222      0.049     -6.549      0.000      -0.419      -0.226
job_retired            0.3821      0.069      5.552      0.000       0.247       0.517
job_services          -0.2423      0.065     -3.701      0.000      -0.371      -0.114
job_student            0.3540      0.086      4.107      0.000       0.185       0.523
default_no             0.3312      0.056      5.943      0.000       0.222       0.440
month_aug              0.4272      0.055      7.770      0.000       0.319       0.535
month_dec              0.8061      0.163      4.948      0.000       0.487       1.125
month_jul              0.7319      0.056     13.094      0.000       0.622       0.841
month_nov              0.2706      0.064      4.249      0.000       0.146       0.395
month_oct              0.8043      0.087      9.258      0.000       0.634       0.975
month_sep              0.5906      0.096      6.160      0.000       0.403       0.778
day_of_week_fri       -0.0044      0.046     -0.097      0.923      -0.094       0.085
day_of_week_wed        0.1226      0.044      2.771      0.006       0.036       0.209
poutcome_failure      -1.8438      0.100    -18.412      0.000      -2.040      -1.647
poutcome_nonexistent  -1.1344      0.070    -16.253      0.000      -1.271      -0.998
poutcome_success       0.0912      0.114      0.803      0.422      -0.131       0.314
==============================================================================
```

Figure 17

The p-values for most of the variables are smaller than 0.05, therefore, most of them are significant to the model.

## Logistic Regression Model Fitting

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.3, random_state=0)
from sklearn.linear_model import LogisticRegression
from sklearn import metrics
logreg = LogisticRegression()
logreg.fit(X_train, y_train)
```

*LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1, penalty='l2', random_state=None, solver='liblinear', tol=0.0001, verbose=0, warm_start=False*

**Predicting the test set results and calculating the accuracy**

```
y_pred = logreg.predict(X_test)
print('Accuracy of logistic regression classifier on test
set: {:.2f}'.format(logreg.score(X_test, y_test)))
```

*Accuracy of logistic regression classifier on test set: 0.90*

## Cross Validation

Cross validation attempts to avoid overfitting while still producing a prediction for each observation dataset. We are using 10-fold Cross-Validation to train our Logistic Regression model.

```
from sklearn import model_selection
from sklearn.model_selection import cross_val_score
kfold = model_selection.KFold(n_splits=10, random_state=7)
modelCV = LogisticRegression()
scoring = 'accuracy'
results = model_selection.cross_val_score(modelCV, X_train,
y_train, cv=kfold, scoring=scoring)
print("10-fold cross validation average accuracy: %.3f" %
(results.mean()))
```

*10-fold cross validation average accuracy: 0.897*

The average accuracy remains very close to the Logistic Regression model accuracy; hence, we can conclude that our model generalizes well.

## Confusion Matrix

```
from sklearn.metrics import confusion_matrix
confusion_matrix = confusion_matrix(y_test, y_pred)
print(confusion_matrix)
```

*[[10872 109]*
*[ 1122 254]]*

The result is telling us that we have 10872+254 correct predictions and 1122+109 incorrect predictions.

# Compute precision, recall, F-measure and support

To quote from Scikit Learn:

The precision is the ratio tp / (tp + fp) where tp is the number of true positives and fp the number of false positives. The precision is intuitively the ability of the classifier to not label a sample as positive if it is negative.

The recall is the ratio tp / (tp + fn) where tp is the number of true positives and fn the number of false negatives. The recall is intuitively the ability of the classifier to find all the positive samples.

The F-beta score can be interpreted as a weighted harmonic mean of the precision and recall, where an F-beta score reaches its best value at 1 and worst score at 0.

The F-beta score weights the recall more than the precision by a factor of beta. beta = 1.0 means recall and precision are equally important.

The support is the number of occurrences of each class in y_test.

```
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.91 | 0.99 | 0.95 | 10981 |
| 1 | 0.70 | 0.18 | 0.29 | 1376 |
| avg / total | 0.88 | 0.90 | 0.87 | 12357 |

Figure 18

**Interpretation**: Of the entire test set, 88% of the promoted term deposit were the term deposit that the customers liked. Of the entire test set, 90% of the customer's preferred term deposits that were promoted.

# ROC Curve

```
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve
logit_roc_auc = roc_auc_score(y_test,
logreg.predict(X_test))
fpr, tpr, thresholds = roc_curve(y_test,
logreg.predict_proba(X_test)[:,1])
plt.figure()
plt.plot(fpr, tpr, label='Logistic Regression (area =
%0.2f)' % logit_roc_auc)
plt.plot([0, 1], [0, 1],'r--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic')
plt.legend(loc="lower right")
plt.savefig('Log_ROC')
plt.show()
```
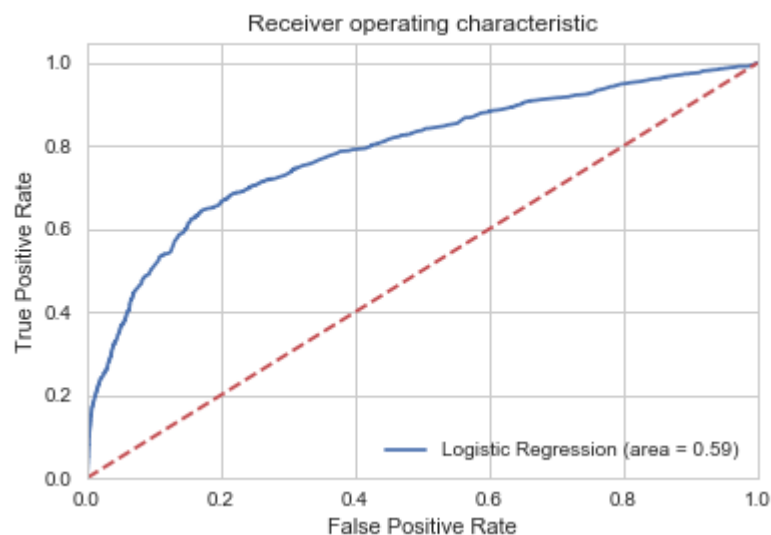


Figure 19

The receiver operating characteristic (ROC) curve is another common tool used with binary classifiers. The dotted line represents the ROC curve of a purely random classifier; a good classifier stays as far away from that line as possible (toward the top-left corner).

The Jupyter notebook used to make this post is available here. I would be pleased to receive feedback or questions on any of the above.

Reference: Learning Predictive Analytics with Python book