# 1. Credit card applications

Commercial banks receive *a lot* of applications for credit cards. Many of them get rejected for many reasons, like high loan balances, low income levels, or too many inquiries on an individual's credit report, for example. Manually analyzing these applications is mundane, error-prone, and time-consuming (and time is money!). Luckily, this task can be automated with the power of machine learning and pretty much every commercial bank does so nowadays. In this notebook, we will build an automatic credit card approval predictor using machine learning techniques, just like the real banks do.



We'll use the Credit Card Approval dataset (http://archive.ics.uci.edu/ml/datasets/credit+approval) from the UCI Machine Learning Repository. The structure of this notebook is as follows:

- First, we will start off by loading and viewing the dataset.
- We will see that the dataset has a mixture of both numerical and non-numerical features, that it contains values from different ranges, plus that it contains a number of missing entries.
- We will have to preprocess the dataset to ensure the machine learning model we choose can make good predictions.
- After our data is in good shape, we will do some exploratory data analysis to build our intuitions.
- Finally, we will build a machine learning model that can predict if an individual's application for a credit card will be accepted.

In [160]:
```python
# Import pandas
# ... YOUR CODE FOR TASK 1 ...
import pandas as pd

# Load dataset
cc_apps = pd.read_csv('datasets/cc_approvals.data', header=None)

# Inspect data
# ... YOUR CODE FOR TASK 1 ...
display(cc_apps.head())
```

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|-----|-----|---|
| 0 | b | 30.83 | 0.000 | u | g | w | v | 1.25 | t | t | 1 | f | g | 00202 | 0 | + |
| 1 | a | 58.67 | 4.460 | u | g | q | h | 3.04 | t | t | 6 | f | g | 00043 | 560 | + |
| 2 | a | 24.50 | 0.500 | u | g | q | h | 1.50 | t | f | 0 | f | g | 00280 | 824 | + |
| 3 | b | 27.83 | 1.540 | u | g | w | v | 3.75 | t | t | 5 | t | g | 00100 | 3 | + |
| 4 | b | 20.17 | 5.625 | u | g | w | v | 1.71 | t | f | 0 | f | s | 00120 | 0 | + |

# 2. Inspecting the applications

The output may appear a bit confusing at its first sight, but let's try to figure out the most important features of a credit card application. The features of this dataset have been anonymized to protect the privacy, but [this blog (http://rstudio-pubs-static.s3.amazonaws.com/73039_9946de135c0a49daa7a0a9eda4a67a72.html)](http://rstudio-pubs-static.s3.amazonaws.com/73039_9946de135c0a49daa7a0a9eda4a67a72.html) gives us a pretty good overview of the probable features. The probable features in a typical credit card application are `Gender`, `Age`, `Debt`, `Married`, `BankCustomer`, `EducationLevel`, `Ethnicity`, `YearsEmployed`, `PriorDefault`, `Employed`, `CreditScore`, `DriversLicense`, `Citizen`, `ZipCode`, `Income` and finally the `ApprovalStatus`. This gives us a pretty good starting point, and we can map these features with respect to the columns in the output.

As we can see from our first glance at the data, the dataset has a mixture of numerical and non-numerical features. This can be fixed with some preprocessing, but before we do that, let's learn about the dataset a bit more to see if there are other dataset issues that need to be fixed.

```
In [162]:   # Print summary statistics
            cc_apps_description = cc_apps.describe()
            print(cc_apps_description)

            print('\n')

            # Print DataFrame information
            cc_apps_info = cc_apps.info()
            print(cc_apps_info)

            print('\n')

            # Inspect missing values in the dataset
            # ... YOUR CODE FOR TASK 2 ...
            display(cc_apps.tail(17))
```

```
                2           7          10             14
count  690.000000  690.000000  690.00000     690.000000
mean     4.758725    2.223406    2.40000    1017.385507
std      4.978163    3.346513    4.86294    5210.102598
min      0.000000    0.000000    0.00000       0.000000
25%      1.000000    0.165000    0.00000       0.000000
50%      2.750000    1.000000    0.00000       5.000000
75%      7.207500    2.625000    3.00000     395.500000
max     28.000000   28.500000   67.00000  100000.000000


<class 'pandas.core.frame.DataFrame'>
RangeIndex: 690 entries, 0 to 689
Data columns (total 16 columns):
0     690 non-null object
1     690 non-null object
2     690 non-null float64
3     690 non-null object
4     690 non-null object
5     690 non-null object
6     690 non-null object
7     690 non-null float64
8     690 non-null object
9     690 non-null object
10    690 non-null int64
11    690 non-null object
12    690 non-null object
13    690 non-null object
14    690 non-null int64
15    690 non-null object
dtypes: float64(2), int64(2), object(12)
memory usage: 86.3+ KB
None
```

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 673 | ? | 29.50 | 2.000 | y | p | e | h | 2.000 | f | f | 0 | f | g | 00256 | 17 | - |
| 674 | a | 37.33 | 2.500 | u | g | i | h | 0.210 | f | f | 0 | f | g | 00260 | 246 | - |
| 675 | a | 41.58 | 1.040 | u | g | aa | v | 0.665 | f | f | 0 | f | g | 00240 | 237 | - |
| 676 | a | 30.58 | 10.665 | u | g | q | h | 0.085 | f | t | 12 | t | g | 00129 | 3 | - |
| 677 | b | 19.42 | 7.250 | u | g | m | v | 0.040 | f | t | 1 | f | g | 00100 | 1 | - |
| 678 | a | 17.92 | 10.210 | u | g | ff | ff | 0.000 | f | f | 0 | f | g | 00000 | 50 | - |
| 679 | a | 20.08 | 1.250 | u | g | c | v | 0.000 | f | f | 0 | f | g | 00000 | 0 | - |
| 680 | b | 19.50 | 0.290 | u | g | k | v | 0.290 | f | f | 0 | f | g | 00280 | 364 | - |
| 681 | b | 27.83 | 1.000 | y | p | d | h | 3.000 | f | f | 0 | f | g | 00176 | 537 | - |
| 682 | b | 17.08 | 3.290 | u | g | i | v | 0.335 | f | f | 0 | t | g | 00140 | 2 | - |
| 683 | b | 36.42 | 0.750 | y | p | d | v | 0.585 | f | f | 0 | f | g | 00240 | 3 | - |
| 684 | b | 40.58 | 3.290 | u | g | m | v | 3.500 | f | f | 0 | t | s | 00400 | 0 | - |
| 685 | b | 21.08 | 10.085 | y | p | e | h | 1.250 | f | f | 0 | f | g | 00260 | 0 | - |
| 686 | a | 22.67 | 0.750 | u | g | c | v | 2.000 | f | t | 2 | t | g | 00200 | 394 | - |
| 687 | a | 25.25 | 13.500 | y | p | ff | ff | 2.000 | f | t | 1 | t | g | 00200 | 1 | - |
| 688 | b | 17.92 | 0.205 | u | g | aa | v | 0.040 | f | f | 0 | f | g | 00280 | 750 | - |
| 689 | b | 35.00 | 3.375 | u | g | c | h | 8.290 | f | f | 0 | t | g | 00000 | 0 | - |

# 3. Splitting the dataset into train and test sets

Now, we will split our data into train set and test set to prepare our data for two different phases of machine learning modeling: training and testing. Ideally, no information from the test data should be used to preprocess the training data or should be used to direct the training process of a machine learning model. Hence, we first split the data and then preprocess it.

Also, features like `DriversLicense` and `ZipCode` are not as important as the other features in the dataset for predicting credit card approvals. To get a better sense, we can measure their statistical correlation (https://realpython.com/numpy-scipy-pandas-correlation-python/) to the labels of the dataset. But this is out of scope for this project. We should drop them to design our machine learning model with the best set of features. In Data Science literature, this is often referred to as *feature selection*.

```python
In [164]:  # Import train_test_split
           # ... YOUR CODE FOR TASK 3 ...
           from sklearn.model_selection import train_test_split

           # Drop the features 11 and 13
           cc_apps.drop(columns=[11, 13], inplace=True)

           # Split into train and test sets
           cc_apps_train, cc_apps_test = train_test_split(cc_apps, test_size=0.33, random_
           display(cc_apps_train.head())
           display(cc_apps_test.head())
```

|     | 0 | 1     | 2     | 3 | 4 | 5 | 6  | 7    | 8 | 9 | 10 | 12 | 14   | 15 |
|-----|---|-------|-------|---|---|---|----|------|---|---|----|----|------|----|
| 382 | a | 24.33 | 2.500 | y | p | i | bb | 4.50 | f | f | 0  | g  | 456  | -  |
| 137 | b | 33.58 | 2.750 | u | g | m | v  | 4.25 | t | t | 6  | g  | 0    | +  |
| 346 | ? | 32.25 | 1.500 | u | g | c | v  | 0.25 | f | f | 0  | g  | 122  | -  |
| 326 | b | 30.17 | 1.085 | y | p | c | v  | 0.04 | f | f | 0  | g  | 179  | -  |
| 33  | a | 36.75 | 5.125 | u | g | e | v  | 5.00 | t | f | 0  | g  | 4000 | +  |

|     | 0 | 1     | 2   | 3 | 4 | 5  | 6  | 7   | 8 | 9 | 10 | 12 | 14  | 15 |
|-----|---|-------|-----|---|---|----|----|-----|---|---|----|----|-----|----|
| 286 | a | ?     | 1.5 | u | g | ff | ff | 0.0 | f | t | 2  | g  | 105 | -  |
| 511 | a | 46.00 | 4.0 | u | g | j  | j  | 0.0 | t | f | 0  | g  | 960 | +  |
| 257 | b | 20.00 | 0.0 | u | g | d  | v  | 0.5 | f | f | 0  | g  | 0   | -  |
| 336 | b | 47.33 | 6.5 | u | g | c  | v  | 1.0 | f | f | 0  | g  | 228 | -  |
| 318 | b | 19.17 | 0.0 | y | p | m  | bb | 0.0 | f | f | 0  | s  | 1   | +  |

# 4. Handling the missing values (part i)

Now we've split our data, we can handle some of the issues we identified when inspecting the DataFrame, including:

- Our dataset contains both numeric and non-numeric data (specifically data that are of `float64`, `int64` and `object` types). Specifically, the features 2, 7, 10 and 14 contain numeric values (of types float64, float64, int64 and int64 respectively) and all the other features contain non-numeric values.
- The dataset also contains values from several ranges. Some features have a value range of 0 - 28, some have a range of 2 - 67, and some have a range of 1017 - 100000. Apart from these, we can get useful statistical information (like `mean`, `max`, and `min`) about the features that have numerical values.
- Finally, the dataset has missing values, which we'll take care of in this task. The missing values in the dataset are labeled with '?', which can be seen in the last cell's output of the second task.

Now, let's temporarily replace these missing value question marks with NaN.

```python
In [166]: # Import numpy
          # ... YOUR CODE FOR TASK 4 ...
          import numpy as np

          # Replace the '?'s with NaN in the train and test sets
          cc_apps_train = cc_apps_train.replace('?', np.nan)
          cc_apps_test = cc_apps_test.replace('?', np.nan)
```

# 5. Handling the missing values (part ii)

We replaced all the question marks with NaNs. This is going to help us in the next missing value treatment that we are going to perform.

An important question that gets raised here is *why are we giving so much importance to missing values*? Can't they be just ignored? Ignoring missing values can affect the performance of a machine learning model heavily. While ignoring the missing values our machine learning model may miss out on information about the dataset that may be useful for its training. Then, there are many models which cannot handle missing values implicitly such as Linear Discriminant Analysis (LDA).

So, to avoid this problem, we are going to impute the missing values with a strategy called mean imputation.

In [168]:
```python
# Impute the missing values with mean imputation

cc_apps_train[[2, 7, 10, 14]].fillna( cc_apps_train.mean(), inplace=True)
cc_apps_test[[2, 7, 10, 14]].fillna(cc_apps_train.mean(), inplace=True)

# Count the number of NaNs in the datasets and print the counts to verify
# ... YOUR CODE FOR TASK 5 ...
print(cc_apps_train.isna().sum())
print(cc_apps_test.isna().sum())
```

```
0      8
1      5
2      0
3      6
4      6
5      7
6      7
7      0
8      0
9      0
10     0
12     0
14     0
15     0
dtype: int64
0      4
1      7
2      0
3      0
4      0
5      2
6      2
7      0
8      0
9      0
10     0
12     0
14     0
15     0
dtype: int64
```

# 6. Handling the missing values (part iii)

We have successfully taken care of the missing values present in the numeric columns. There are still some missing values to be imputed for columns 0, 1, 3, 4, 5, 6 and 13. All of these columns contain non-numeric data and this is why the mean imputation strategy would not work here. This needs a different treatment.

We are going to impute these missing values with the most frequent values as present in the respective columns. This is good practice (https://www.datacamp.com/community/tutorials/categorical-data) when it comes to imputing missing values for categorical data in general.

```python
In [170]:  # Iterate over each column of cc_apps_train
           for col in cc_apps_train.columns :
               # Check if the column is of object type
               if cc_apps_train[col].dtype == 'object':
                   # Impute with the most frequent value
                   cc_apps_train = cc_apps_train.fillna( cc_apps_train[col].value_counts(
                   cc_apps_test = cc_apps_test.fillna(cc_apps_test[col].value_counts().in

           # Count the number of NaNs in the dataset and print the counts to verify
           # ... YOUR CODE FOR TASK 6 ...
           print(cc_apps_train.isna().sum())
           print(cc_apps_train.isna().sum())
```

```
0      0
1      0
2      0
3      0
4      0
5      0
6      0
7      0
8      0
9      0
10     0
12     0
14     0
15     0
dtype: int64
0      0
1      0
2      0
3      0
4      0
5      0
6      0
7      0
8      0
9      0
10     0
12     0
14     0
15     0
dtype: int64
```

# 7. Preprocessing the data (part i)

The missing values are now successfully handled.

There is still some minor but essential data preprocessing needed before we proceed towards building our machine learning model. We are going to divide these remaining preprocessing steps into two main tasks:

1. Convert the non-numeric data into numeric.
2. Scale the feature values to a uniform range.

First, we will be converting all the non-numeric values into numeric ones. We do this because not only it results in a faster computation but also many machine learning models (like XGBoost) (and especially the ones developed using scikit-learn) require the data to be in a strictly numeric format. We will do this by using the `get_dummies()` method from pandas.

```python
In [172]:  # Convert the categorical features in the train and test sets independently
           cc_apps_train = pd.get_dummies(cc_apps_train[[0, 1, 3, 4, 5, 6]])
           cc_apps_test = pd.get_dummies(cc_apps_test[[0, 1, 3, 4, 5, 6]])

           # Reindex the columns of the test set aligning with the train set
           cc_apps_test = cc_apps_test.reindex(columns=cc_apps_train.columns, fill_value=
```

# 8. Preprocessing the data (part ii)

Now, we are only left with one final preprocessing step of scaling before we can fit a machine learning model to the data.

Now, let's try to understand what these scaled values mean in the real world. Let's use `CreditScore` as an example. The credit score of a person is their creditworthiness based on their credit history. The higher this number, the more financially trustworthy a person is considered to be. So, a `CreditScore` of 1 is the highest since we're rescaling all the values to the range of 0-1.

```python
In [174]:  # Import MinMaxScaler
           # ... YOUR CODE FOR TASK 8 ...
           from sklearn.preprocessing import MinMaxScaler

           # Segregate features and labels into separate variables
           X_train, y_train = cc_apps_train.iloc[ : , :-1].values, cc_apps_train.iloc[ :
           X_test, y_test = cc_apps_test.iloc[ : , :-1].values, cc_apps_test.iloc[ : ,  [

           # Instantiate MinMaxScaler and use it to rescale X_train and X_test
           scaler = MinMaxScaler(feature_range=(0, 1))
           rescaledX_train = scaler.fit_transform(X_train)
           rescaledX_test = scaler.transform(X_test)
```

# 9. Fitting a logistic regression model to the train set

Essentially, predicting if a credit card application will be approved or not is a classification (https://en.wikipedia.org/wiki/Statistical_classification) task. According to UCI, our dataset contains more instances that correspond to "Denied" status than instances corresponding to "Approved" status. Specifically, out of 690 instances, there are 383 (55.5%) applications that got denied and 307 (44.5%) applications that got approved.

This gives us a benchmark. A good machine learning model should be able to accurately predict the status of the applications with respect to these statistics.

Which model should we pick? A question to ask is: *are the features that affect the credit card approval decision process correlated with each other?* Although we can measure correlation, that is outside the scope of this notebook, so we'll rely on our intuition that they indeed are

correlated for now. Because of this correlation, we'll take advantage of the fact that generalized linear models perform well in these cases. Let's start our machine learning modeling with a Logistic Regression model (a generalized linear model).

In [176]:
```python
# Import LogisticRegression
from sklearn.linear_model import LogisticRegression

# Instantiate a LogisticRegression classifier with default parameter values
logreg = LogisticRegression()

# Fit logreg to the train set
logreg.fit(rescaledX_train,y_train)
```

Out[176]:
```
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)
```

# 10. Making predictions and evaluating performance

But how well does our model perform?

We will now evaluate our model on the test set with respect to classification accuracy (https://developers.google.com/machine-learning/crash-course/classification/accuracy). But we will also take a look the model's confusion matrix (http://www.dataschool.io/simple-guide-to-confusion-matrix-terminology/). In the case of predicting credit card applications, it is important to see if our machine learning model is equally capable of predicting approved and denied status, in line with the frequency of these labels in our original dataset. If our model is not performing well in this aspect, then it might end up approving the application that should have been approved. The confusion matrix helps us to view our model's performance from these aspects.

In [178]:
```python
# Import confusion_matrix
# ... YOUR CODE FOR TASK 10 ...
from sklearn.metrics import confusion_matrix
# Use logreg to predict instances from the test set and store it
y_pred = logreg.predict(rescaledX_test)

# Get the accuracy score of logreg model and print it
print("Accuracy of logistic regression classifier: ", logreg.score(rescaledX_t

# Print the confusion matrix of the logreg model
# ... YOUR CODE FOR TASK 10 ...
confusion_matrix(y_test, y_pred)
```

```
Accuracy of logistic regression classifier:  0.9956140350877193
```

Out[178]:
```
array([[227,   0],
       [  1,   0]])
```

# 11. Grid searching and making the model perform better

Our model was pretty good! In fact it was able to yield an accuracy score of 100%.

For the confusion matrix, the first element of the of the first row of the confusion matrix denotes the true negatives meaning the number of negative instances (denied applications) predicted by the model correctly. And the last element of the second row of the confusion matrix denotes the true positives meaning the number of positive instances (approved applications) predicted by the model correctly.

But if we hadn't got a perfect score what's to be done?. We can perform a grid search (https://machinelearningmastery.com/how-to-tune-algorithm-parameters-with-scikit-learn/) of the model parameters to improve the model's ability to predict credit card approvals.

scikit-learn's implementation of logistic regression (http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html) consists of different hyperparameters but we will grid search over the following two:

- tol
- max_iter

```
In [180]:  # Import GridSearchCV
           # ... YOUR CODE FOR TASK 11 ...
           from sklearn.model_selection import GridSearchCV
           # Define the grid of values for tol and max_iter
           tol = [0.01, 0.001, 0.0001]
           max_iter = [100, 150, 200]

           # Create a dictionary where tol and max_iter are keys and the lists of their va
           param_grid = {'tol':tol, 'max_iter':max_iter}
```

# 12. Finding the best performing model

We have defined the grid of hyperparameter values and converted them into a single dictionary format which `GridSearchCV()` expects as one of its parameters. Now, we will begin the grid search to see which values perform best.

We will instantiate `GridSearchCV()` with our earlier `logreg` model with all the data we have. We will also instruct `GridSearchCV()` to perform a cross-validation (https://www.dataschool.io/machine-learning-with-scikit-learn/) of five folds.

We'll end the notebook by storing the best-achieved score and the respective best parameters.

While building this credit card predictor, we tackled some of the most widely-known preprocessing steps such as **scaling**, **label encoding**, and **missing value imputation**. We finished with some **machine learning** to predict if a person's application for a credit card would get approved or not given some information about that person.

In [182]:
```python
# Instantiate GridSearchCV with the required parameters
grid_model = GridSearchCV(estimator=logreg, param_grid=param_grid, cv=5)

# Fit grid_model to the data
grid_model_result = grid_model.fit(rescaledX_train, y_train)

# Summarize resultsy
best_score, best_params = grid_model_result.best_score_, grid_model_result.bes
print("Best: %f using %s" % (best_score, best_params))

# Extract the best model and evaluate it on the test set
best_model = grid_model_result.best_estimator_
print("Accuracy of logistic regression classifier: ", best_model.score(rescale
```

```
Best: 0.984848 using {'max_iter': 100, 'tol': 0.01}
Accuracy of logistic regression classifier:  0.9956140350877193
```