

# INF582: DATA SCIENCE - LEARNING FROM DATA

ÉCOLE POLYTECHNIQUE

## AXA Data Challenge - Assignment

Data Science and Mining (DaSciM) Team

February 5, 2016

### 1 Description of the Assignment

Whether in a contact centre or bank branch environment, workforce managers face the constant challenge of balancing the priorities of service levels and labour costs. In the case where the demand (inbound calls, outbound calls, emails, web chats, etc.) is greater than supply (the agents themselves), the price, in the form of reduced service levels, falling customer satisfaction and poor agent morale, rises. On the other side, when supply is greater than demand, service levels tend to improve, but at the cost of idle and unproductive agents. The key to optimising the bottom line performance of a contact centre is to find a harmonious balance between supply and demand. This bottom line performance is directly impacted by the direct costs of hiring and employing your agents, but it is also influenced by client satisfaction, agent morale and other factors. Taking all the above into consideration, the basis of any good staffing plan is an accurate workload forecast. An accurate forecast gives us the opportunity to predict workload in order to get the right number of staff in place to handle it.

The specific project constitutes an AXA data challenge, where its purpose is to apply data mining and machine learning techniques for the development of an *inbound call forecasting system*. The forecasting system should be able to predict the number of incoming calls for the AXA call center in France, on a per “half-hour” time slot basis. The prediction is for three (3) days ahead in time. More specifically, based on the history of the incoming calls up to a specific time stamp (you cannot use data/features that corresponds to future time slots), the proposed model should be able to predict the number of the calls, received three (3) days later. In this way, the problem can be seen as a regression problem and the goal is to design a model that achieves to predict the incoming calls of the AXA call center in France with high accuracy. A detailed description of the dataset that will be used for the training of your proposed models is given in Section 2. The specific dataset includes telephony data retrieved from AXA call centers, and corresponds to the period spanning the calendar years 2011 and 2012. Last but not least, the final evaluation of your model will be given by using a *Leaderboard platform* (more details about the *Leaderboard platform* will be given soon).

#### 1.1 Data Challenge Awards Ceremony

As the data are provided by AXA Assistance and this data challenge forms part of the activities of the *AXA-X DaSciS chair* - after the evaluation of your submissions there will be a reception organized by the chair and prizes will be awarded to the best submissions. You will be informed on the details in due time.

## 2 Dataset Description

In this section, we present the structure of the training dataset (`train_2011_2012.csv`<sup>1</sup>) that will be used for the training of your model. As mentioned previously, the training dataset includes telephony data derived from AXA call centers, and correspond to the calendar years 2011 and 2012. Figure 1 shows how the training dataset has been derived. Each one of the rows of the dataset corresponds to the number of incoming calls for each different combination of values for the following attributes: DATE (time stamps in half hour slots), SPLIT\_COD, ACD\_COD, ASS\_ASSIGNMENT. Please consider that some combinations may not be present on the dataset. For a detailed description of the attributes refer to the `field_description.xlsx`<sup>2</sup> file.

The *objective of your work* here is to train a model (or a set of models) able to predict the number of incoming calls (CSPL\_RECEIVED\_CALLS) for three days after the current/given date for each different set of values of the attributes: Date(time stamps in half hour slots), ASS\_ASSIGNMENT. For the purposes of this project please limit the analysis of your model solely on the calls made within the territory of France. Nevertheless, it's up to you how to use the entire dataset if this improves your models.

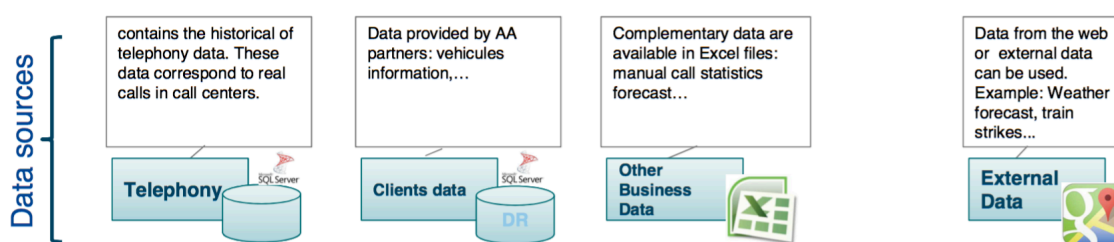


Figure 1: Data Sources of the training dataset.

Concerning external data, that can be used to improve the predictive performance of your model, we supply as well weather data<sup>3</sup> for years 2011 (`meteo_2011.csv`) and 2012 (`meteo_2012.csv`), respectively. Each row of the weather dataset offers information about the weather of a city, at a specific date-time. Figure 2 gives an example of the first five rows of the dataset. The features represented by the columns of the weather data files are presented below:

- date,
- dept\_nb (Department number),
- city,
- temperature,
- wind\_dir (wind direction),
- precip (precipitation),
- pressure\_hPa (pressure in hPa).

As you can easily observe, the dataset has some missing values for some attributes (NULL). In the preprocessing task, you should take care of a number of similar cases. In the case of features that take numerical values, one approach could be to replace the missing values with the mean value of this feature. Some other features may not be useful at the prediction task. It would be helpful to explore the dataset and try to deal with such cases. Additionally, some of the features take values that correspond to a string (e.g., the TPER\_TEAM feature takes values Jours and Nuit). In such cases, we can create two new features (i.e., add two new columns to the data matrix) which correspond to the two possible strings. Thus, if the TPER\_TEAM feature takes the value Jours, the feature that corresponds to Jours will become equal to 1, while the feature that corresponds to Nuit will become equal to 0.

<sup>1</sup>Training dataset: [http://moodle.lix.polytechnique.fr/data\\_challenge/train\\_2011\\_2012.7z](http://moodle.lix.polytechnique.fr/data_challenge/train_2011_2012.7z).

<sup>2</sup>Dataset description: [http://moodle.lix.polytechnique.fr/data\\_challenge/field\\_description.zip](http://moodle.lix.polytechnique.fr/data_challenge/field_description.zip).

<sup>3</sup>Weather data: [http://moodle.lix.polytechnique.fr/data\\_challenge/meteo.7z](http://moodle.lix.polytechnique.fr/data_challenge/meteo.7z).

```
>>> train_data = pd.read_csv('meteo_2012.csv',header=0)
>>> print "Size of the dataset: ", train_data.shape
Size of the dataset: (6299857, 8)
>>> print train_data.head()
  2012-01-01 00:00  44      Nantes-Atlantique  12.9  22.7  OSO  0  1019.8
0  2012-01-01 00:00  59                Dunkerque  12.2  32.4   SO  0  1010.4
1  2012-01-01 00:00  59            Lille-Lesquin  12.1  27.0  OSO  0  1012.1
2  2012-01-01 00:00  60                Creil    12.3  22.7  SSO  0  1015.2
3  2012-01-01 00:00  60          Beauvais-Tille  12.1  22.7  OSO  0  1014.6
4  2012-01-01 00:00  61  Alençon - Valframbert  12.5  19.1   SO  0  1017.0
>>> █
```

**Figure 2:** First five rows of the weather dataset for 2012 (using the `pandas` module).

As part of the preprocessing step, you can also apply feature selection techniques to keep a subset with the most informative features or dimensionality reduction methods (e.g., Linear Discriminant Analysis). It is also possible to create new features that do not exist in the dataset, but can be useful in the forecasting task. Thus, you can create a new feature (i.e., add a new column to the data matrix) to represent this information (this is known as feature engineering or generation).

### 3 Summary of the Pipeline

The pipeline that will be followed in the project is similar to the one followed in the labs. In the following, we briefly describe each part of the pipeline.

- *Data pre-processing:* After loading the data, a preprocessing task should be done to transform the data into an appropriate format. In the previous section, we discussed some of these points.
- *Feature engineering - Dimensionality reduction:* The next step involves the feature engineering task, i.e., how to select a subset of the features that will be used in the learning task (feature selection) or how to create new features from the already existing ones (see also previous section). Moreover, it is possible to apply dimensionality reduction techniques in order to improve the performance of the algorithms.
- *Learning algorithm:* The next step of the pipeline involves the selection of the appropriate learning (i.e., regression) algorithm for the problem. At this point, you can test the performance of a number of different algorithms and choose the best one. Additionally, you can follow an ensemble learning approach, combining many regression algorithms.
- *Evaluation:* In Section 4, we describe in detail how the evaluating will be performed.

### 4 Evaluation

You will build your model based on the training data contained in the `train_2011_2012.csv` file. To do this, you can apply *cross-validation* techniques<sup>4</sup>. The goal of cross-validation is to define a dataset to test the model in the training phase, in order to limit problems like overfitting and have an insight on how the model will generalize to an independent dataset (i.e., an unknown dataset, like the test dataset that will be used to assess your model).

In  $k$ -fold cross-validation, the original sample is randomly partitioned into  $k$  equal size subsamples. On the  $k$  subsamples, a single subsample is retained as the validation data for testing the model, and the remaining  $k - 1$  subsamples are used as training data. The cross-validation process is then repeated  $k$  times (the folds), with each of the  $k$  subsamples used exactly once as the validation (i.e., test) data. The  $k$  results from the folds can then be averaged (or otherwise combined) to produce a single estimation (average accuracy of the model).

<sup>4</sup>Wikipedia's lemma for *Cross-validation*: [http://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](http://en.wikipedia.org/wiki/Cross-validation_(statistics)).

Of course, having a good model that achieves good accuracy under cross validation does not guarantee that the same accuracy will be also achieved for the test data. Thus, the final evaluation of your model will be done on the test dataset contained in the `submission.txt`<sup>5</sup> file. So, after having a model that performs well under cross-validation, you should train the model using the whole training dataset and test it on the test dataset as described below.

### Submission file

For the final evaluation of your model, you have to predict the number of calls that will be received at a number of different combination of values of the following attributes: `DATE` (corresponding to half hour slots), `ASS_ASSIGNMENT`. More specifically, get the predicted number of calls for each instance (row) contained in the `submission.txt` file. Each row of the `submission.txt` file corresponds to a different combination `DATE` (corresponding to half hour slots) and `ASS_ASSIGNMENT` (Table 1 presents a snapshot of the `submission.txt` file). In the `submission.txt` example file, all the prediction values are set equal to zero. You must replace those values with your predicted ones. Do not change the format of the file (fields separated by tab). The final evaluation of your model will be made based on the mean square error (MSE) metric.

DATE	ASS_ASSIGNMENT	Prediction
2012-01-03 00:00:00.000	CAT	0
2012-01-03 00:00:00.000	Tlphonie	0
2012-01-03 00:00:00.000	Tech. Inter	0
2012-01-03 00:00:00.000	Tech. Axa	0
2012-01-03 00:00:00.000	Services	0

**Table 1:** First five rows of the `submission.txt` example file

The data corresponding to the required dates (the listed dates in the `submission.txt` file) are omitted from the dataset. Moreover the data on a 2-day window a priori to those dates listed in the `submission.txt` file, are also omitted to ensure that you will not use them for the predictions of your submission.

## 5 Useful Python Libraries

In this section, we briefly discuss some useful tools that can be useful in the project and you are encouraged to use.

- For the preprocessing task which also involves some initial data exploration, you may use the `pandas` Python library for data analysis<sup>6</sup>.
- A very powerful machine learning library in Python is `scikit-learn`<sup>7</sup>. It can be used in the preprocessing step (e.g., for feature selection) and in the calls forecasting task (a plethora of regression algorithms have been implemented in `scikit-learn`). Recall that we have already used the `scikit-learn` in the labs.
- Finally, you are always encouraged to propose and develop your own learning algorithms or use the ones developed in the labs.

## 6 Details about the Submission of the Project

Your final evaluation for the project will be based on both the accuracy (according to the mean square error (MSE) metric) of the proposed model as well as on your total approach to the problem.

**As part of the project, you have to submit the following:**

<sup>5</sup>Testing dataset: [http://moodle.lix.polytechnique.fr/data\\_challenge/submission.txt](http://moodle.lix.polytechnique.fr/data_challenge/submission.txt)

<sup>6</sup><http://pandas.pydata.org/>.

<sup>7</sup><http://scikit-learn.org/>.

1. Your **final submission file** (`submission.txt`), which contains the estimated number of calls.
2. A **2-5 pages report**, in which you should describe the approach and the methods that you used in the project. Since this is a real data science task, we are interested to know how you dealt with each part of the pipeline, e.g., if you have created new features and why, which algorithms did you use for the calls forecasting task and why, their performance (accuracy and training time), approaches that finally didn't work but is interesting to present them, and in general, whatever you think that is interesting to report). Also, in the report, please provide the names of the team members and the identifier of your team (e.g., INF582).
3. A directory with the code of your implementation.
4. Create a `.zip` file with the `team_identifier.zip` (the identifier of your team), containing the code and the report and **submit it to moodle (one submission per team)**.
5. **Deadline: Friday, March 4, 23:59.**