

A real time approach targeted at buyer prediction, using behaviour classification in Tourism Web Analytics

Deep learning @ INAF – Pula 16 - 19 September 2019

Stefano Fadda - CRS4 - stefano@crs4.it

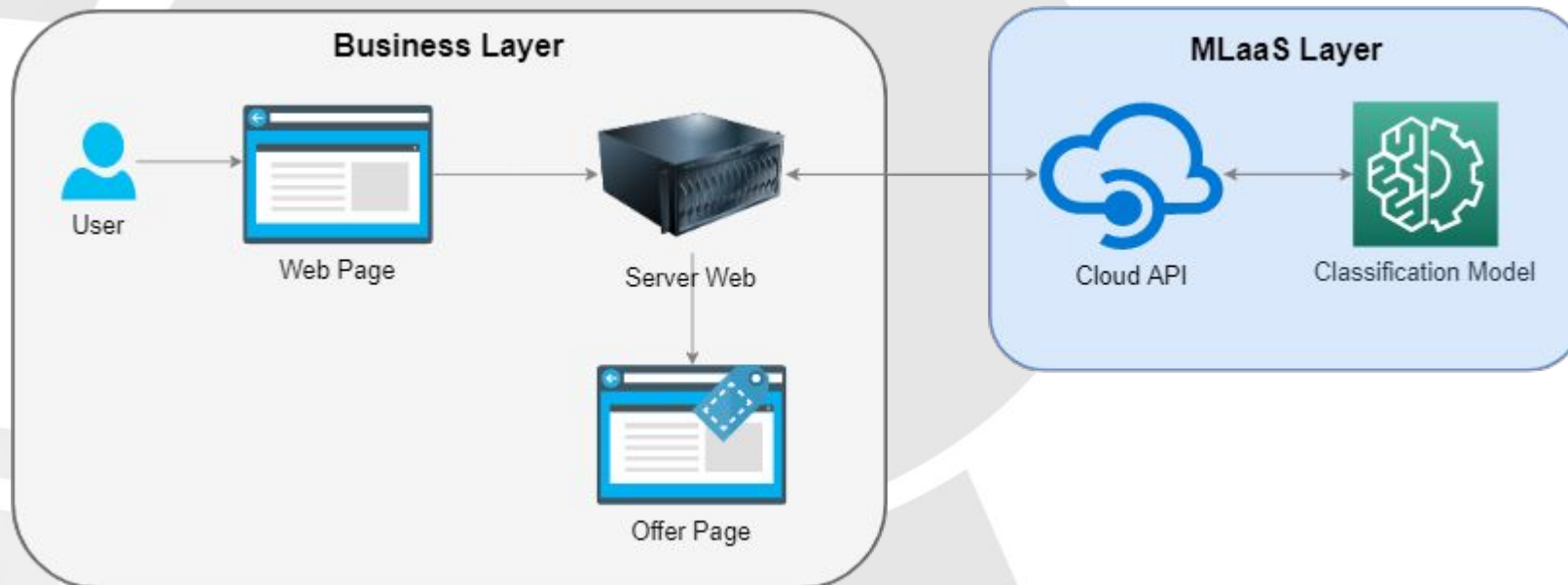
- context
- purchase prediction use case
- solution and methodology
- dataset
- models
- web application (MLaaS)
- conclusions



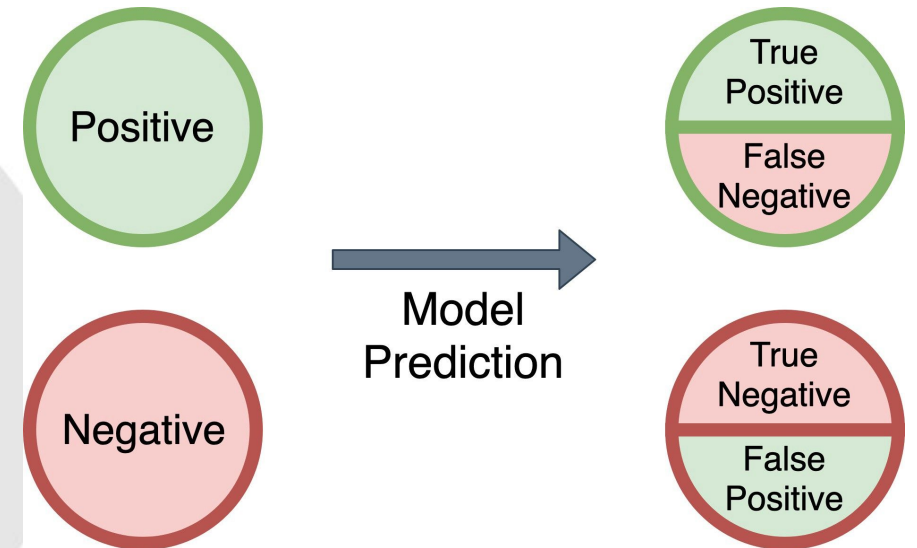
- statistical data from the partner agency about hotel/holidays structure
- user behaviour and satisfaction comprehension
- historical user navigation sessions to understand user behaviours



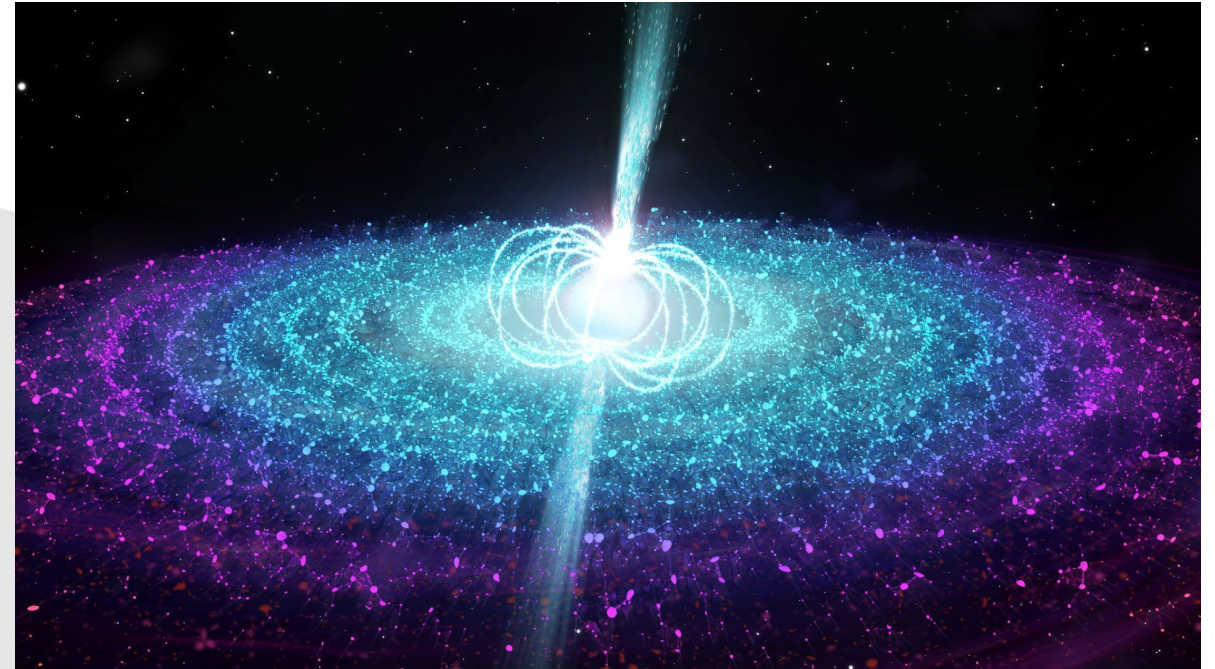
- automatic prediction of the user's purchase will
- possibility to perform the prediction in real time



- the problem can be addressed as a binary classification, so the data available can be cataloged as belonging to the category of who buys (positive) or not (negative)
- selection of all data about user navigations that are related to the two categories



- similar methodology has been applied to classification and identification of Pulsar
- a lot of false positives
- oversampling techniques to balance the dataset

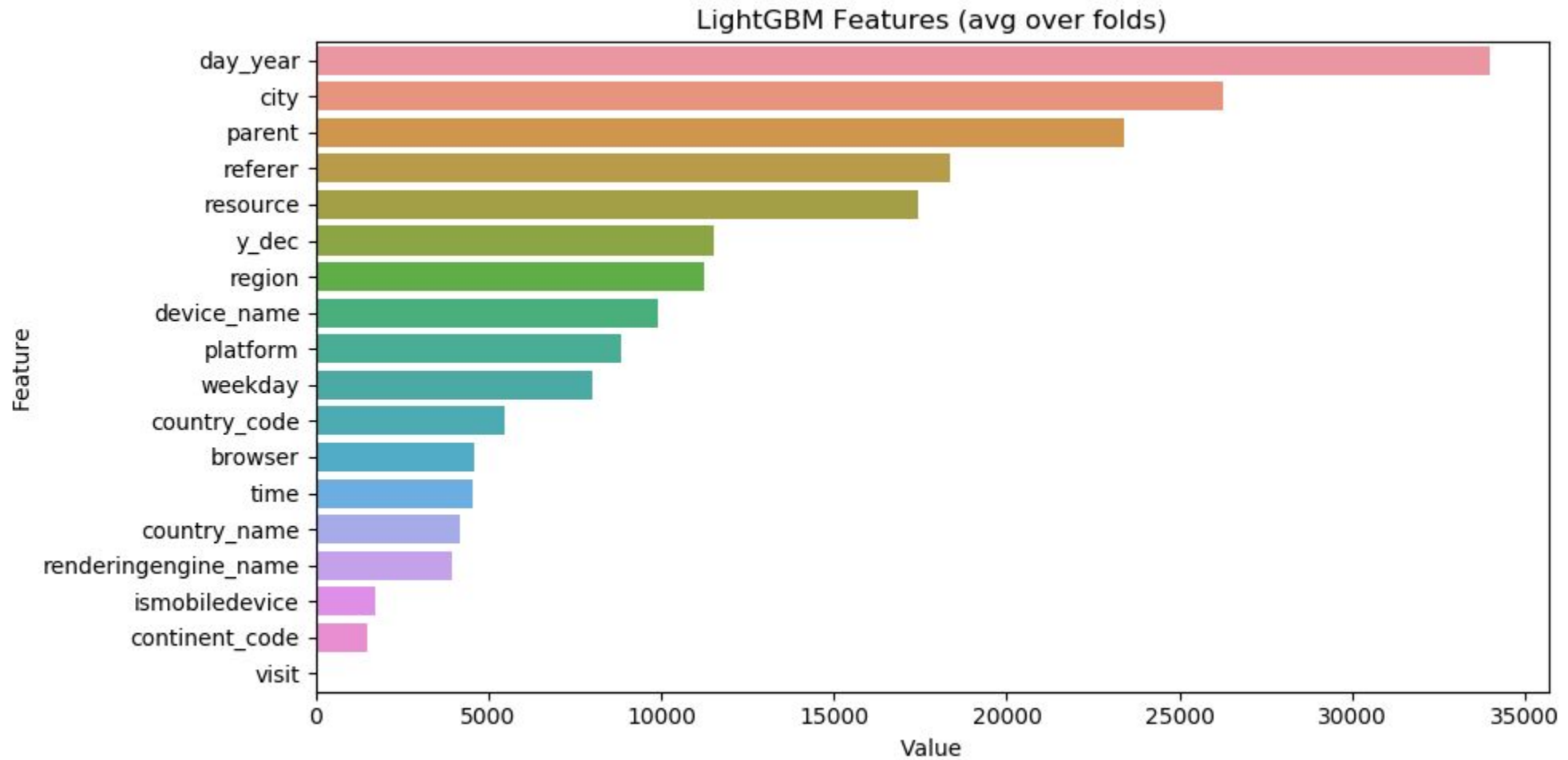


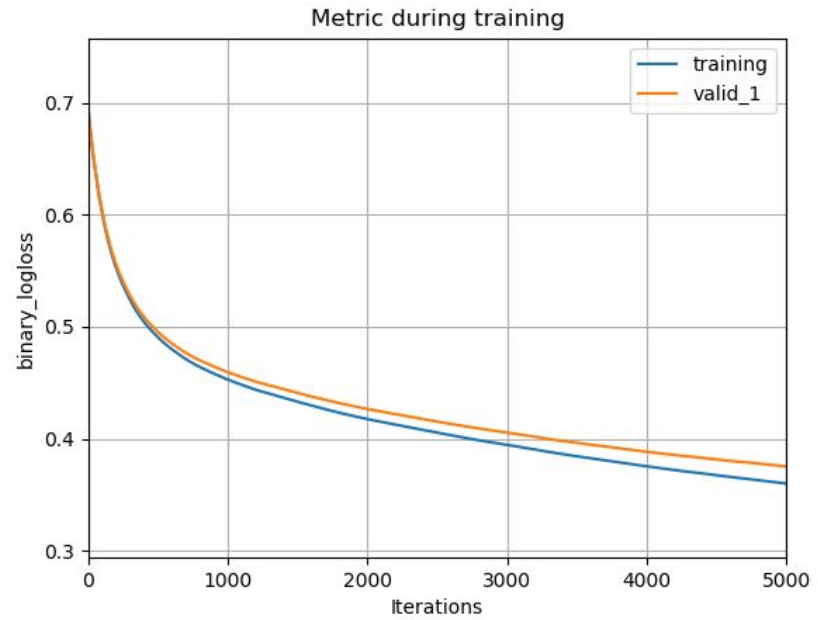
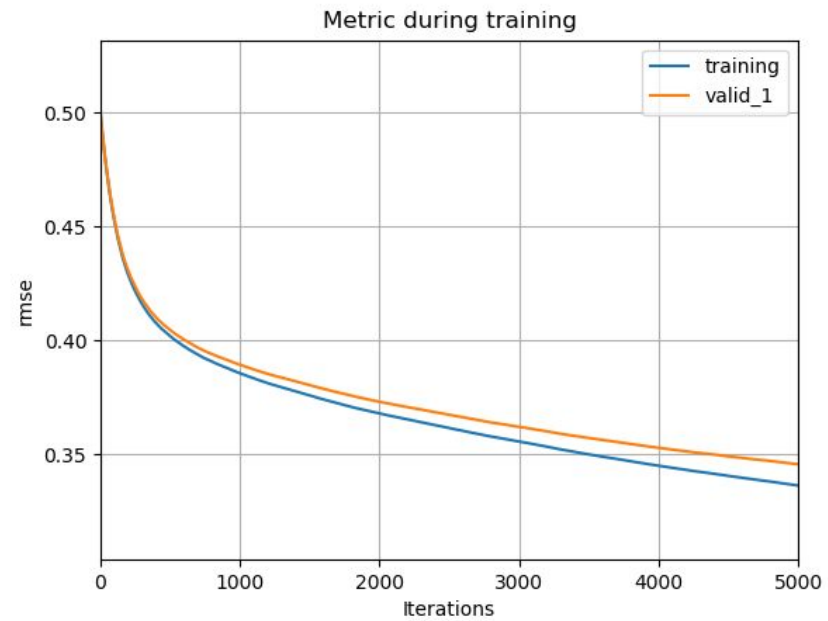
(Punia, Akhil, Ashish Sardana, and Monica Subashini. "Evaluating advanced machine learning techniques for pulsar detection from HTRU survey." 2017 International Conference on Intelligent Sustainable Systems (ICISS). IEEE, 2017.)

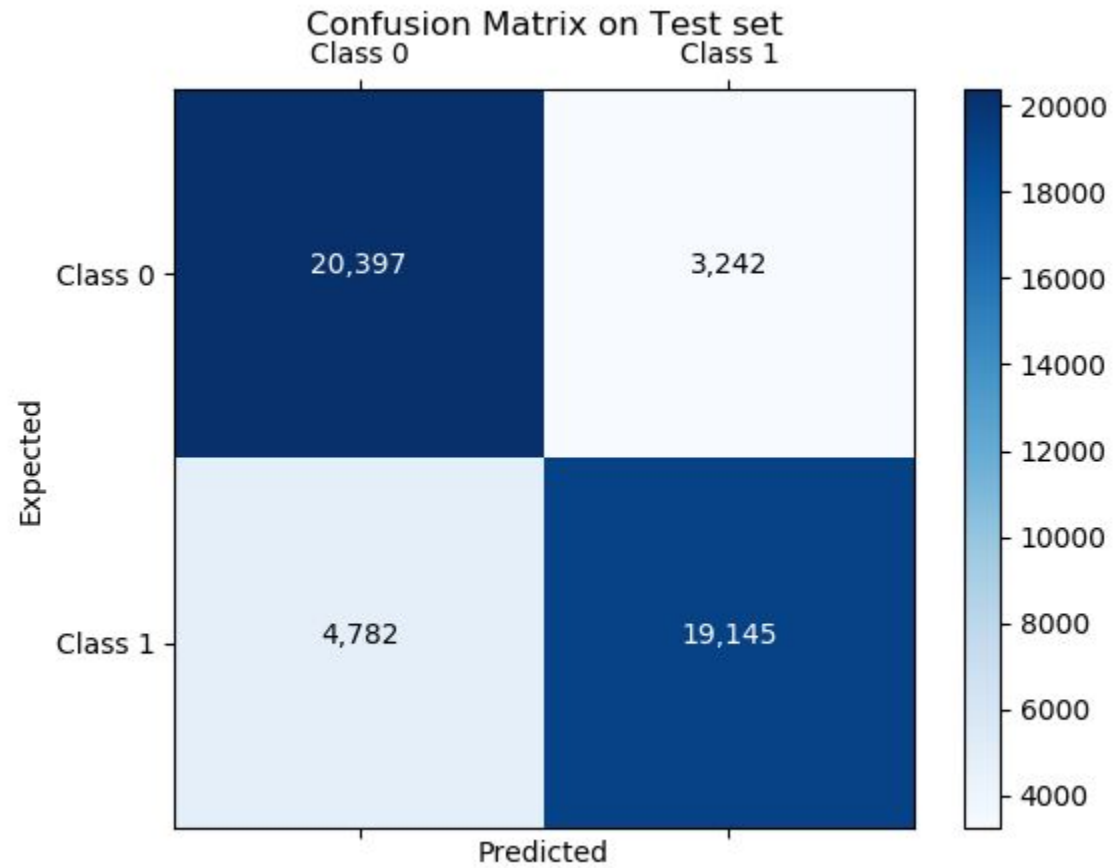
- positive user sessions represent 4% of the total
- dataset is extremely unbalanced on the negatives
- undersampling algorithm like RandomUnderSampler
- its obtained a balanced dataset of the two categories with the same weight
- dataset preprocessing phase to manage null data, duplicates and utility transformation

- analysis of the state-of-the-art of the main python libraries about binary classification: LightGBM, XGBoost, Scikit-learn
- RMSE is a measure of accuracy, to compare forecasting errors of different models for a particular dataset
- LightGBM has been identified as the faster in the training phase

	RMSE-train	RMSE-test	Training-time (sec.)
LightGBM	0.20426	0.22897	3.846010
XGBoost	0.23423	0.25058	20.176368
Scikit-learn	0.25721	0.26421	11.143738







- Web service as Restful API interface to execute the prediction on real-time navigation user data
- facilitate the integration with third party systems
- Web frontend to test the ML service, monitor the service and see the results

POST

none form-data x-www-form-urlencoded raw binary GraphQL BETA JSON (application/json)

```
1 [
2   {
3     "commit": 1,
4     "resource_id": 658566,
5     "session_id": 172660,
6     "visit": 0,
7     "resource": null,
8     "time": 32,
9     "created": "2018-08-09 10:14:42",
10    "referer": "https://www.google.it/",
11    "label": "ita",
12    "continent_code": "EU",
13    "country_code": "IT",
14    "country_name": "Italy",
15    "region": "10",
16    "city": "Morrovalle",
17    "postal_code": "63854",
18    "parent": "Firefox 61.0",
19    "platform_version": "6.1",
20    "renderingengine_name": "Gecko",
21    "ismobiledevice": null,
22    "device_name": "Windows Desktop"
23  }
24 ]
```

```
1 {
2   "status": 200,
3   "message": "true",
4   "body": {
5     "prediction": {
6       "commit": 0.0,
7       "predicted": 0.0,
8       "Accuracy": 0.8421052631578947,
9       "MSCLS": 0.1578947368421053,
10      "Precision": 0.0,
11      "F1_score": 0.0
12    }
13  }
14 }
```

The project is actually a work in progress, we evaluate the system with real live data in the near future

The approach and the technology is immediately applicable in other context