



Ensemble Methods Project

M2 Computer Science MALIA

Guillaume Metzler

Institut de Communication (ICOM)

Université de Lyon, Université Lumière Lyon 2

Laboratoire ERIC UR 3083, Lyon, France

guillaume.metzler@univ-lyon2.fr

Groups of two persons are allowed to do this project.

The work shall be submitted at the following mail adress

guillaume.metzler@univ-lyon2.fr

before the deadline mentioned in your program (usually in March). More precisely, the submitted work shall contain a python or a notebook which contains the code that has used to produce the results and report (taking the form of a pdf file) where you will present the data and used methods to get them.

This report shall focus on the comparison of the results and the analysis of the later.

1 Datasets

You can find several data sets at the following address

[Download the data .](#)

you can also find a python code which gives you the opportunity to provide load and prepare them, such that they will be ready to be used.

[Prepare the Data](#)

For this work, you will have to use at least 20 different datasets.
Some indications on how to write the report are provided at the end of the document.

2 Work to perform

Methodology Your datasets can be divided into two groups according to a specific characteristic and maybe the procedure to used and/or the performance metrics you are going to use will be different.

For each of these groups of datasets, you will then need to compare the following algorithms :

- penalized logistic regression
- an ensemble method based on bagging
- a random forest method
- an ensemble method based on boosting
- an ensemble method based on stacking (learning the combination of voters)
- a gradient boosting method

What can be interesting also, is to test whether the combination is interesting. Thus, as a baseline, for instance, if you have decided to combine several linear SVM, to see the performance of a single linear SVM.

For the second group of data, it can also be interesting to test several sampling methods and/or cost-sensitive learning.

Code The code that will be submitted shall provide the same results as the one presented in the report!

3 A study of an approximation of kernel methods

The aim is to attempt to approximate kernel methods what we call Random Fourier Features, for the sake of simplicity, you can consider it as a boosting methods using cosine functions and using some landmarks

The landmarks allowed us to compute similarities with the weights of our training set. More precisely, if we consider a set of $l < m$ landmarks selected from our training set, the associated kernel is then of size $l \times l$ and our classifier was given by :

$$f(\mathbf{x}) = \sum_{k=1}^l \alpha_k y_k k(\mathbf{x}, \mathbf{x}_k).$$

To predict the label of a new data point, we only had to compute the similarity of the new point with the landmarks that had been chosen randomly.

We now propose to do two things :

- iteratively learn optimal landmarks,
- approximate our kernel by a sum of *cos*.

We will do all this in a *Gradient Boosting* context, and Algorithm 1 summarizes the procedure.

3.1 Understanding the algorithm

We focus on the different steps of the algorithm while keeping in mind the objective of this lab session.

3.2 Implementation

We now seek to implement this algorithm. You will then need to create a Python class (or functions) that can learn the model and perform prediction.

We should also keep in mind that our model depends on a hyperparameter λ that should ideally be tuned. Here, we will simply assign it a fixed value at the beginning, for example $\lambda = 1$.

For the implementation, you will need to solve optimization problems, which implies using certain solvers :

```
python from scipy import optimize from scipy.optimize import fminbound
```

In particular, we will use the function `optimize.fmin_l_bfgs_b` from the **scipy** library.

Algorithm 1: Gradient Boosting with RFF

Inputs : Dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$; Number of iterations T ;
Parameters γ and λ

Output $\text{sign}\left(H^0(\mathbf{x}) + \sum_{t=1}^T \alpha^t \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)\right)$
:

```
1:  $H^0 \leftarrow H^0(\mathbf{x}_i) = \frac{1}{2} \ln \frac{\sum_{j=1}^n (1+y_j)}{\sum_{j=1}^n (1-y_j)}$ 
2: for  $t = 1, \dots, T$  do
3:    $\forall i = 1, \dots, n, \quad w_i = \exp(-y_i H^{t-1}(\mathbf{x}_i))$ 
4:    $\forall i = 1, \dots, n, \quad \tilde{y}_i = y_i w_i$ 
5:   Draw  $\sim \mathcal{N}(0, 2\gamma)^d$ 
6:    $\mathbf{x}^t = \arg \min_{\mathbf{x}^t \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \exp\left(-\tilde{y}_i \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)\right)$ 
7:    $t = \arg \min_{\mathbf{x}^t \in \mathbb{R}^d} \lambda \|\mathbf{x}^t\|_2^2 + \frac{1}{n} \sum_{i=1}^n \exp\left(-\tilde{y}_i \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)\right)$ 
8:    $\alpha^t = \frac{1}{2} \ln \frac{\sum_{i=1}^n \left(1 + y_i \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)\right) w_i}{\sum_{i=1}^n \left(1 - y_i \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)\right) w_i}$ 
9:    $\forall i = 1, \dots, n, \quad H^t(\mathbf{x}_i) = H^{t-1}(\mathbf{x}_i) + \alpha^t \cos\left({}^t \cdot (\mathbf{x}_i - \mathbf{x}^t)\right)$ 
10: end for
```

Finally, for the implementation, you will need to compute a gradient and store all the parameter values of the model.

The experiments will be conducted on the previous datasets (see Section 1).

3.3 Expected work

Using the above code, try to explain how the methods work, by describing each step and the optimization problems that are solved and why it can be considered as a boosting approach.

To illustrate the method, you will conduct several experiments :

- you will show the frontier decision of the first iterations of the algorithms (the 3 first) on a 2D dimension dataset with few points
- you will then show the decision frontier of this method and compare it to those provided by a XGBoost (for instance) and a Kernel SVM on 2D dataset using the make_swiss_role and moon datasets with python.
- You will compare the performance and running time of the XGBoost, LGBM and this algorithm with respect to the number of iterations on one of the previous dataset (use graph for the representation) and use a kernel method as a baseline
- You will compare finally study the performances of different gradient based methods (Gradient Boosting, XGBoost, LGBM, this algorithm) on the datasets provided in

Section 1.

What follows is a suggestion on how to write your report.

Introduction

You will rapidly present the context of your project and the main target and difficulties that are linked to your datasets.

Methodology

It mainly consists in presenting in few words and equations the way the methods are working. Saying differently, the way you will tackle the problem and used concepts.

You first introduce the notations that will be used for the purpose.

You will then present the tools you will be using in the experimental part. We will start by talking about what we want to maximize, *i.e.* the performance measure by defining it before tackling the presentation of the algorithms (in a very short way).

For example, if you are testing an algorithm based on boosting, you briefly have to explain how it works. Don't hesitate to present the process in an abstract way, *i.e.* with mathematical notations and not just in words. Pseudo-code is also useful for summarizing the proposed approach.

If you are proposing several approaches for comparison purposes, you should take care to present the different approaches and justify why you are interested in them.

Note : it is not necessary to present all the algorithms used, but only those used to develop an "exotic" version.

Experiments

You will then draw up your experimental protocol, setting out the method(s) you have selected for the task in hand. This is generally composed of three parts.

Experimental Protocol

For each part, you will proceed as follows :

You will start by presenting the data you are using, usually in the form of a table that includes :

name of the data set, sample size, size of the data set and any other information you consider relevant.

Briefly present the experiments you will be running, the different algorithms you will be testing, the range of hyper-parameters you will be using, and how these hyper-parameters are optimized (cross-validation in k -folds, simple validation or do you choose to fix them). What are your training/validation/test sets ?

The information you provide in this section should enable the reader to reproduce the results you present in the next subsection.

Results

Here, you will present and analyze the results obtained, using graphs and/or tables. In addition to performance, you may also be interested in the speed of an algorithm.

The analysis should also highlight the advantages and disadvantages of the proposed methods.

This may involve using other performance measures/criteria to evaluate/comparison your algorithms.

As for the results, they can be presented in the form of a table similar to the one presented in Table 1.

Conclusion

The aim is to conclude your study about the studied methods. Review the proposed work and the main conclusions. It is also important to offer some perspectives on your work, based on the results obtained and the approach proposed. What method(s) not explored here could have been used to improve the results, and why do you think this is relevant given the experiments carried out ?

TABLE 1 – Mean test F1 over 10 iterations

Dataset	knn	svm_linear	svm_poly	svm_gauss
00.62% abalone20	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0
01.39% abalone17	13.7 \pm 6.3	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0
02.36% yeast6	44.0 \pm 14.6	0.0 \pm 0.0	38.9 \pm 18.0	0.0 \pm 0.0
03.31% wine4	5.0 \pm 5.0	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0
06.67% libras	81.1 \pm 12.7	71.0 \pm 10.0	88.3 \pm 7.9	75.6 \pm 11.0
10.23% pageblocks	84.9 \pm 2.2	75.7 \pm 2.9	82.3 \pm 2.5	78.6 \pm 2.5
10.98% yeast3	66.8 \pm 5.9	65.6 \pm 10.8	76.4 \pm 5.2	71.3 \pm 8.2
13.60% abalone8	22.7 \pm 1.3	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0
14.29% segmentation	85.1 \pm 3.3	45.8 \pm 3.7	64.1 \pm 4.9	50.0 \pm 4.6
22.73% hayes	14.8 \pm 10.2	0.0 \pm 0.0	39.4 \pm 19.9	0.0 \pm 0.0
23.52% vehicle	88.1 \pm 3.4	60.3 \pm 5.3	92.6 \pm 1.7	56.7 \pm 6.5
30.00% german	41.3 \pm 6.4	5.3 \pm 7.6	56.9 \pm 6.8	16.8 \pm 6.6
32.71% glass	70.2 \pm 4.0	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0
33.15% wine	86.5 \pm 4.3	89.6 \pm 5.9	90.1 \pm 5.8	86.2 \pm 6.1
34.90% pima	41.3 \pm 6.6	5.0 \pm 5.4	44.7 \pm 6.1	17.5 \pm 4.8
35.90% iono	82.5 \pm 4.1	88.9 \pm 4.0	92.0 \pm 3.4	93.3 \pm 3.4
37.50% autmpg	82.4 \pm 5.5	0.0 \pm 0.0	83.0 \pm 3.4	0.0 \pm 0.0
46.08% balance	95.3 \pm 1.9	94.5 \pm 1.2	99.9 \pm 0.3	97.3 \pm 0.8
Mean	55.9 \pm 5.4	33.4 \pm 3.2	52.7 \pm 4.8	35.7 \pm 3.0
Average Rank	2.00	3.72	1.72	2.56