

Class Project 1

Submission Deadline: Oct 31st, 2016

Introduction

In this project, you will learn to use the concepts we have seen in the lectures and practiced in the labs on a real-world dataset, start to finish. You will do exploratory data analysis to understand your dataset and your features, do feature processing and engineering to clean your dataset and extract more meaningful information, implement and use machine learning methods on real data, analyze your model and generate predictions using those methods and report your findings.

Grading. Project 1 counts 10% to your final grade in the course. Within Project 1, we will grade your code, report and competition score (each counting one third). (Project 2 will count 30%).

Logistics

Group formation. For Project 1, we will *randomly* assign you to a group of 3 students (two in exceptional cases). We are trying to maximize fairness, and hope you keep in mind that a good data science team combines a diverse set of skills, and greatly benefits from inter-disciplinary backgrounds.

After completion of this first project, we will decide on the group assignment method for the second project based on our experience and your feedback for the first one.

Deliverables at a glance. (More details and grading criteria further down)

- **Code.** In Python. *No external libraries allowed!* For this first project, we want you to implement and use the methods we have seen in class. (External libraries will be allowed in Project 2).
- **Written Report.** You will write a maximum 2 page PDF report on your findings, using LaTeX.
- **Competitive Part.** To give you immediate feedback and a fair ranking, we use the competition platform kaggle.com to score your predictions. You can submit whenever and almost as many times as you like, up until the final submission deadline.

The Dataset. For this course, we are providing you with our own online competition based on one of the most popular machine learning challenges recently - finding the Higgs boson - using original data from CERN.

Step 1 - Getting Started

Create an account using your epfl.ch email and head over to the competition arena

<http://inclass.kaggle.com/epfml-project-1>

Download the training dataset, available in .csv format. To load the data, use the same code we used during the labs. You can find an example of a csv loading function in our provided code from labs 1 and 2.

Step 2 - Implement ML Methods

We want you to implement and use the methods we have seen in class and in the labs. You will need to provide working implementations of the functions in Table 1. If you have not finished them during the labs, you should start by implementing the first ones to have a working toolbox before diving in the dataset.

Function	Details
<code>least_squares_GD(y, tx, gamma, max_iters)</code>	Linear regression using gradient descent
<code>least_squares_SGD(y, tx, gamma, max_iters)</code>	Linear regression using stochastic gradient descent
<code>least_squares(y, tx)</code>	Least squares regression using normal equations
<code>ridge_regression(y, tx, lambda_)</code>	Ridge regression using normal equations
<code>logistic_regression(y, tx, gamma, max_iters)</code>	Logistic regression using gradient descent or SGD
<code>reg_logistic_regression(y, tx, lambda_, gamma, max_iters)</code>	Regularized logistic regression using gradient descent or SGD

Table 1: List of functions to implement. In the above method signatures, for iterative methods, `gamma` is the step-size and `max_iters` is the number of steps to run. `lambda_` is always the regularization parameter. (Note that we have used the trailing underscore because `lambda` is a reserved word in Python with a different meaning).

You should take care of the following

- All code should be easily readable and commented.
- Note that we might automatically call your provided methods and evaluate for correct implementation

Here are some best practices of scientific computing as a reference: <http://arxiv.org/pdf/1210.0530>.

Step 3 - Submitting your Predictions

Once you have a working model (using the above methods or a modified one), you can send your predictions to the Kaggle competition to see how your model is doing against the other teams. You can submit whenever and as many times as you like, until the deadline.

Your predictions must be in `.csv` format, see `sample-submission.csv`. You must use the same datapoint ids as in the test set `test.csv`. To generate `.csv` output from Python, use our provided helper functions in `helpers.py` (see project 1 folder on github).

After a submission, Kaggle will compute your score on the test set, and will show you your score and ranking in the leaderboard.

This is useful to see how you compare against other teams, but you should not consider this score as the only evaluation of your model. Always estimate your test error by using a local validation set, or local cross-validation! This is important to avoid overfitting the test set online. Also, it allows you to make experiments faster, and save uploading bandwidth :). You are only allowed a maximum of 5 submissions to Kaggle per day.

Improving your predictions. While the above described method implementations must be part of your code submission, you can now implement additional modifications of these basic methods above. You can construct better features for the task, or perform better data preprocessing for this particular dataset, or even implement an additional modification of one of the above mentioned ML methods. Note that it is not allowed to use external libraries, code or data in this project. (It will be allowed in Project 2).

Step 4 - Final Submission of Your Project

Your final submission to the CMT system (a standard system used for scientific conferences) must consist of the following:

- **Report:** Your 2 page report as .pdf
- **Code:** The complete executable and documented Python code, as one .zip file.
Rules for the code part:
 - *Reproducibility:* In your submission, you must provide a script `run.py` which produces *exactly* the same .csv predictions which you used in your best submission to the competition on Kaggle.
 - *Documentation:* Your ML system must be clearly described in your PDF report and also well-documented in the code itself. A clear ReadMe file must be provided. The documentation must also include all data preparation, feature generation as well as cross-validation steps that you have used.
 - In addition to your customized system, don't forget that your code submission must still also include the 6 *basic method implementations* as described above in step 2.
 - *No use of external ML libraries* is allowed in Project 1. (It will be allowed in Project 2).
 - No external datasets allowed.

Submission URL: <https://cmt3.research.microsoft.com/EPFML2016>

Submission deadline: Oct 31st, 2016 (midnight)

Step 5 - Profit

... ;)

Physics Background

The Higgs boson is an elementary particle in the Standard Model of physics which explains why other particles have mass. Its discovery at the Large Hadron Collider at CERN was announced in March 2013. In this project, you will apply machine learning techniques to actual CERN particle accelerator data to recreate the process of “discovering” the Higgs particle. For some background, physicists at CERN smash protons into one another at high speeds to generate even smaller particles as by-products of the collisions. Rarely, these collisions can produce a Higgs boson. Since the Higgs boson decays rapidly into other particles, scientists don't observe it directly, but rather measure its “decay signature”, or the products that result from its decay process. Since many decay signatures look similar, it is our job to estimate the likelihood that a given event's signature was the result of a Higgs boson (signal) or some other process/particle (background). In practice, this means that you will be given a vector of features representing the decay signature of a collision event, and asked to predict whether this event was signal (a Higgs boson) or background (something else). To do this, you will use the binary classification techniques we have discussed in the lectures.

If you're interested in more background on this dataset, we point you to the longer description here: https://higgsml.lal.in2p3.fr/files/2014/04/documentation_v1.8.pdf.

Note that understanding the physics background is not necessary to perform well in this machine learning challenge as part of the course.

Appendix

Grading Criteria

- **Competitive Part (counts one third).** The final rank of your team in the (private) leaderboard will be translated linearly to a scale from 4 to 6.
- **Code (counts one third).** In Python. *No external libraries allowed!* For this first project, we want you to implement and use the methods we have seen in class. The code will be graded by two TAs independently, according to the criteria described above in Step 4.
- **Written Report (counts one third).** You will write a maximum 2 page PDF report on your findings, using LaTeX. The code will be graded by two TAs independently, and we will provide you feedback. The main criteria will be if you were able to correctly use, implement and describe the 6 baseline methods mentioned in Step 2 above. This counts half for the written report. In addition, we will grade you on the scientific contribution you made additionally, to improve your predictions. For this part, the criteria are
 - scientific novelty
 - creativity
 - reproducibility
 - solid comparison baselines supporting your claims
 - writeup quality

As usual, your code and report will be automatically checked for plagiarism.

Guidelines for Machine Learning Projects

Now that you have implemented few basic methods, you should use this toolbox on the dataset. Here are a few things that you might want to try.

Exploratory data analysis You should learn about your dataset - figure out which features are continuous, which ones are categorical, check if there are obvious relationships between the features, take a look at the distribution of each feature, and so on. Check https://en.wikipedia.org/wiki/Exploratory_data_analysis.

Feature processing Cleaning your dataset by removing useless features and values, combining others, finding better representations of the features to feed your model, scaling the features, and so on. Check this article on feature engineering: <http://machinelearningmastery.com/discover-feature-engineering-how-to-engineer-features-and-how-to-get-good-at-it/>.

Determining whether a method overfits or underfits You should be able to diagnose whether your model is over- or underfitting the data and take actions to fix the problems with your model. Recommended reading: *Advice on applying machine learning methods* by Andrew Ng: <http://cs229.stanford.edu/materials/ML-advice.pdf>.

Applying methods and visualizing Beyond simply applying the models we have seen, it helps to try to understand what the ML model is doing. Try to find out which datapoints are wrongly classified and, if possible, why this is the case. Then use this information to improve your model. Check Peter Domingo's *Useful things to know about machine learning*: <http://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf>

Accurately estimate how well your method is doing By applying cross-validation and estimating the test error.

Report Guidelines

In addition to finding a good model for the data, you will need to explain your methodology in a report. For the first project, this will help you getting used to writing, and prepare you for the more extensive Project 2.

Clearly describe your used methods, state your conclusions and argue that the results you obtained make (or do not make) sense, and the reasons behind it. Keep the report short and to the point, with a strict limit of 2 pages (Project 2 will allow 4 pages). References are allowed to be put on a extra third page.

To get started more easily with writing the report, we provide you a LaTeX template here

github.com/epfml/ML_course/tree/master/projects/project1/latex-example-paper

The file also contains some more helpful information on how to write a scientific report or paper. We will also help you learn it during the exercise session and office hours if you ask us.

For more guidelines on what makes a good report, see the grading criteria above. In particular, don't forget to take care about

- *Reproducibility*: Not only in the code, but also in the report, do include complete details about each algorithm you tried, e.g. what lambda values you used for ridge regression? How exactly did you do that feature transformation? how many folds did you use for cross-validation? etc...
- *Baselines*: Give clear experimental evidence: When you added this new combined feature, or changed the regularization, by how much did that increase or decrease the test error? It is crucial to always report such obtained differences in the evaluation metrics, and to include several properly implemented baseline algorithms as a comparison to your approach.

Some additional resources on LaTeX:

- <http://www.maths.tcd.ie/~dwilkins/LaTeXPrimer/> - tutorial on LaTeX
- <http://www.stdout.org/~winston/latex/latexsheet-a4.pdf> - cheat sheet collecting most of all useful commands in LaTeX
- <http://mirror.switch.ch/ftp/mirror/tex/info/first-latex-doc/first-latex-doc.pdf> - example how to create a document with Latex
- <http://en.wikibooks.org/wiki/LaTeX> - detailed tutorial on LaTeX

Producing figures for LaTeX in Python

There are some good visualization tools in Python. “matplotlib” is probably the single most used Python package for 2D-graphics. The relevant tutorials are as follow:

- Matplotlib tutorial: <http://www.labri.fr/perso/nrougier/teaching/matplotlib/>
- Matplotlib tutorial: <https://sites.google.com/site/scigraphs/tutorial>
- Matplotlib Tutorial: http://jakevdp.github.io/plt_tutorial/

Regarding other useful Python data visualization libraries, please refer to this blog for more information.