
LAB 1: Binary classification and model selection

- This lab addresses binary classification and model selection on synthetic data.
- The aim of the lab is to play with the libraries and to get a practical grasp of what we have discussed in class.
- Follow the instructions below.

Goal:

This lab is divided in three parts depending of their level of complexity (**Beginner**, **Intermediate**, **Advanced**). Your goal is to complete entirely, at least, one of the three parts.

Setup instructions:

Running MATLAB

Download the file `regml2018_lab1.zip` from the syllabus on the course website (<http://lcs1.mit.edu/courses/regml/regml2018/#syllabus>), extract it and add all the sub-folders to the MATLAB path. This file includes all the code you need!

PART I: Beginner

Overture: Warm up

Run the file `gui_filter.m` and a GUI will start. Have a look at the various components. The objective of this section is to familiarize you with the interface.

- **I.A** Use the left panel to load a simulated classification dataset of type *Linear* (press the button `load data` to generate). Observe the generated data (the buttons `plot training` and `plot test` will allow you to toggle between training and test sets). See how the number of wrong labels affects the data set. You can also look at the other types of datasets, which are nonlinear.
- **I.B** From now, consider a *Linear* dataset, starting with 100 (resp. 1000) training (resp. test) samples, and some errors in the labels. Use the right panel to find a classifier for this dataset. For this, take a *regularized least squares* filter, associated to a *linear* kernel, and try various parameters. To choose the regularization parameter you can either choose KCV or set a fixed value. Press the button `run` to perform training and classification. Observe then the plot of the KCV error and the balance between training and test errors. Also have a look at the plot area on the left, where a separation function has appeared (again the buttons `plot training` and `plot test` allow you to switch between the two).

Allegro con brio: Analysis

Carry on the following experiments either using the GUI interface, when it is possible, or writing appropriate scripts.

- **I.C** Back in the shell, check the content of directory `./spectral_reg_toolbox`. There you will find, among others, the code for command `learn` (used for training), `patt_rec` (used for testing), `kcv` (used for model selection on the training set). For more information about the parameters and the usage of those scripts, type:

```
help learn
help patt_rec
help kcv
```

Finally, you may want to have a look at the content of directory `./dataset_scripts` and in particular to file `create_dataset.m`, that allows you to generate synthetic data of different kinds.

NOTE:

In the code we use a different notation from the one you have seen in the classes. In the *Regularized Least Squares* method (`rls.m`), the regularization parameter is τ instead of λ .

- **I.D** Generate noisy data of *Linear* type. Considering *linear-RLS*, observe how the training and test errors change as:
 - We change (increase or decrease) the regularization parameter τ .
 - The training set size grows (try various choices of n as long as your computer supports you!).
 - The amount of errors in the labels in the generated data grows.

Run training and testing for various choices of the suggested parameters.

- **I.E** Leaving all the other parameters fixed, choose an appropriate range for the regularization parameter, `tval=[t_min:t_step:t_max]`, and plot the training and the test errors for each τ . For doing this, you might need the function `learn` that you saw in I.C:

```
[alpha, err] = learn('lin', [], 'rls', t, X, Y, 'class');
```

In this context, you can have access to the training and test errors corresponding to the parameter τ with:

```
training_error = cell2mat(err);
Y_learnt = kernel('lin', [], Xt, X) * cell2mat(alpha);
test_error = learn_error(Y_learnt, Yt, 'class');
```

Use the KCV option to select by cross-validation the optimal regularization parameter, and see how it relates to your previous plot. If you could, would you use the test error to select the parameter?

- **I.F** Leaving all the other parameters fixed, choose an appropriate range for the number of points in the training set, `nval=[n_min: n_step:n_max]`, and plot the training and test errors. What do you observe as $n \rightarrow \infty$?

PART II: Intermediate

Crescendo: Advanced Analysis

- **II.A** Consider *gaussian-RLS* and perform parameter tuning in this case. This time, together with the regularization parameter τ , you'll have to choose an appropriate `sigma`, the kernel parameter.

- Try some (σ, τ) pairs and compare the obtained training error and test error.
- Fix τ and observe the effect of changing σ .
- Fix σ and observe the effect of changing τ .
- Do you notice (and if so, when) any overfitting/oversmoothing effects?

Try to confirm your results by exploring a range of parameters and plotting the training and test errors, as you did in I.E and I.F.

- **II.B** Consider *polynomial-RLS* and perform parameter tuning as in II.A. How does the choice of the kernel affect the learning behavior of the algorithm? In particular, compare the performances of the polynomial and Gaussian kernels on the *spiral* and *moons* datasets with respect to the number of examples in the training set (e.g. [10, 20, 50, 100, 1000]) and the value of the parameter regularization (e.g. [0.5, 0.1, 0.01, 0.001, 0.0001]).

PART III: Advanced

Finale: The Challenge

The challenge consists in a learning task using a real dataset, namely *USPS (United States Postal Service)*: This dataset contains a number of handwritten digits images. The problem is to train *the best classifier* that is able to discriminate between the digits `1` and `7`.

Once the classifiers are trained, they must be exported by means of the script `save_challenge_lab1` (to see how to use it, look at its code). The file `demo_lab1.m` contains a code snippet to perform a simple binary classification task by means of the previously presented scripts.

Submission: You should drop your results in a matrix file named `name-surname` to the link: <https://www.dropbox.com/request/j2z8yMoIjBxyfbRKaF6s> by the end of the challenge session.

The results will be presented during the next class. The score of your result will be based on the accuracy of the classifier on a *completely independently sampled* test set.

Deadline: 6:00 PM.