

Learning from Data – Week 3

Assignment 3: Support Vector Machines and Feature Analysis

General remarks

This assignment is meant to get your acquainted with Support Vector Machines. You will also have to explore feature contribution quite a bit further, and do some feature analysis. Additionally, you're encouraged to experiment a little further with report writing, by making your own choices (how do you want to report on the experiments? How about feature description? And feature analysis? Experiment!)

For the practical parts, you are asked to train and test a few models using SVM, and produce what you think the best settings are for the data we're using. You are also asked to understand what role different features are playing.

Please note: for this assignment you will work in pairs. If possibly, it might be a good idea to work in different pairs than the previous paired assignment, but you don't necessarily have to change. In the report, please include a section at the end where you clarify who did what.

What you have to hand in, as a pair:

- code with your SVM model (see Exercise 3.1). You have to assume that we will run it like this:
`LFDassignment3_SVM_yournames.py <trainset> <testset>`
where `trainset` and `testset` have the same format. You are given the `trainset`, but we're holding out the test set. Please, make sure that your script is taking arguments.
- research report (based on the template) that describes what you've done and also includes or directly incorporates answers/discussion to all questions in this document. Please, make sure you submit a pdf file.

Deadline: 30 September 2019, 23:59.

Data

For this assignment, we will be using the review data we used for Assignment 1 and Assignment 2, and which you find on Nestor in this assignment's materials.

Exercise 3.1 – Support Vector Machines

You will be working on the *binary sentiment classification* with the dataset from Assignment 2, which is uploaded on Nestor again for this assignment. We are again withholding the test set, and we will evaluate your model on it.

3.1.1 Default settings

Run a support vector machine with a linear kernel (`cls = svm.SVC(kernel='linear', C=1.0)`) on the binary sentiment classification. Use default settings and report results using either cross-validation or a portion of the dataset as development set. Just make sure you document what you have done.

3.1.2 Setting C

Try to change the value of the C parameter, and see if you can get better results (either with x-validation or on a development set, just use the same settings as above). Please, in the report include details of what you observe by changing C, and also explain what the C parameter is used for.

3.1.3 Using a non-linear kernel

Do the same as above, but use a radial basis function (rbf) kernel. You will also have to specify a *gamma* parameter, in addition to C (e.g. `cls = svm.SVC(kernel='rbf', gamma=0.7, C=1.0)`) You can check information on the *gamma* (and C) parameter here: scikit-learn.org/stable/auto_examples/svm/plot_rbf_parameters.html. Experiment with changing these values, and report what you observe. It has been claimed that rbf kernels are not great for text classification, and are normally outperformed by linear kernels — do you observe the same?

You are free to use grid search to find the best parameters, but bear in mind that it might be computationally challenging.

3.1.4 Implementation differences

In scikit there are two SVM implementations: `svm.SVC` and `svm.LinearSVC`. By checking the documentation online, and by testing both implementations for your experiments, please write up a short paragraph in your report where you explain what the differences between the two implementations are, and any usage recommendation you might have (when to use one implementation or the other).

3.1.5 Best SVM model

Based on your observations and results on x-validated data or on your development set, make a final decision on your kernel and parameters, and report it in the document. Send in

your best SVM model, like you did for Assignment 2, calling it: `LFDassignment3_SVM_yournames.py`. You are expected to experiment with adding features, too. Possibly n-grams might help? Or part-of-speech information? You have to build larger, more informed systems, so this is a good place to start experimenting. We will run it on held-out data, and send back the results to you. Please, make sure that your script takes arguments, as you did for Assignment 2 – it will be run like this:

```
LFDassignment3_SVM_yournames.py <trainset> <testset>.
```

3.1.6 Feature contribution

As mentioned in the previous section (Section 3.1.5), you are encouraged to experiment with additional features. For this portion of the exercise, you are also asked to more closely *explore* which features actually contribute the most. You can use the attribute `svm.coef_` to do this. Try to explore how to use it, and what information it gives you (weight? direction?), and comment on the results you obtain with your model. (Please note that this attribute only works when using a linear kernel.)