

IN4320 Machine Learning Assignment 2

February 28, 2018

When it comes to the programming part, this assignment can be put into one sentence: take two-class linear discriminant analysis [LDA, the classifier that assumes the class-conditional distributions to be Gaussian with the same covariance matrix and its parameters are estimated through maximum likelihood] and implement two different ways of performing semi-supervised learning for this classifier. You do need to do a bit more, however, . . . The second and more important part of the exercise is concerned with constructing/designing insightful experiments that illustrate the pros and cons of your methods.

When it comes to the implementation of your two semi-supervised approaches for LDA, you are certainly allowed to take any inspiration from other works, papers, web pages, etc., you are even allowed to implement existing methods. In any case, do provide proper references to where you got your inspiration from!

Now, let us make this challenging assignment a bit more concrete. Here are the more specific questions for you to answer and exercises for you to do.

Real

- a** Define and describe your two [really different?] ways of semi-supervised learning for LDA *on an algorithmic level*. Keep the descriptions for the two methods clearly separate. Before giving these descriptions, do note item **d**. The more different your two choices are, the easier it will be to solve those later exercises.
- b** Take the *MAGIC Gamma Telescope Data Set* from the UCI repository¹ and *first normalize all 10 features*² on the full data set once before all other experiments. Based on this normalized data set, make learning curves against the number of *unlabeled* samples for a total of 25 *labeled* samples in the training set. Check, at least, adding 0, 10, 20, 40, 80, 160, 320, and 640 unlabeled samples and see how the expected error rates change. Compare the two curve to the supervised error rates. Make sure

¹See <https://archive.ics.uci.edu/ml/machine-learning-databases/magic/magic04.data>. Note that the *last* column contains the class labels, which are encoded as g and h. The first 10 columns are the features.

²That is, make all 10 feature standard deviations equal to 1.

you repeat your experiments sufficiently often to get some nice and, possibly, smooth curves. Do you get significant changes in error rates?

- c** With the same preprocessed data set as in **b**, make the same type of plots, but now plot the *log*-likelihood³ [and not the error rate] versus the number of unlabeled data.

Imaginary

- d** Construct two artificial data sets. On the one data set, your first semi-supervised LDA should work well *in terms of the error rate* and improve over the regular supervised learner, but the second should give deteriorated performance on this same set: its performance should be worse than the supervised classifier. On the other data set, it should be the other way around: the second semi-supervised LDA should work better than the supervised learner and the first learner should fail to do so. Consider the setting in which you take few labeled samples and a large number of unlabeled samples [no need for learning curves]. Explain why the respective improvements and failures are expected.

My assessment: you should be able to keep your report within three pages.

³On a test set of course. Make sure that you test on sets of the same sizes or ensure in some other way that you can compare the performance between different amounts of added unlabeled data.