

Statistical Machine Learning (W4400)

Spring 2014

<http://stat.columbia.edu/~porbanz/teaching/W4400/>

Peter Orbanz

porbanz@stat.columbia.edu

Lu Meng

lumeng@stat.columbia.edu

Homework 2

Due: 4 March 2014

Homework submission: We will collect your homework **at the beginning of class** on the due date. If you cannot attend class that day, you can leave your solution in my postbox in the Department of Statistics, 10th floor SSW, at any time before then.

Problem 0 (Linear Classification)

Consider a perceptron classifier in \mathbb{R}^2 , given by a hyperplane with orthogonal vector v_H and offset c , and two points \mathbf{x}_1 and \mathbf{x}_2 . Suppose that

$$v_H := \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} \quad c := \frac{1}{2\sqrt{2}} \quad \mathbf{x}_1 := \begin{pmatrix} -3 \\ 0 \end{pmatrix} \quad \mathbf{x}_2 := \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} .$$

1. Compute the classification result for \mathbf{x}_1 and \mathbf{x}_2 .
2. If the classifier is given by the same pair $(v_H, -c)$, but was trained as an SVM with margin 1, do the results change? Please explain your answer.
3. Which cost function does the perceptron cost function approximate, and why do we approximate it?

Problem 1 (Perceptron)

Our Perceptron implementation will use the homogeneous coordinate representation of vectors, i.e. vectors $\mathbf{x} \in \mathbb{R}^d$ are represented by $(\mathbf{x}, 1) = (x_1, \dots, x_d, 1)^t \in \mathbb{R}^{d+1}$. Recall that this representation turns affine hyperplanes

$$\{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{v}_H, \mathbf{x} \rangle - c = 0\}$$

with normal vector $\mathbf{v}_H = (v_1, \dots, v_d)^t$ into hyperplanes through the origin of the form

$$\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} \mid (\mathbf{v}_H, -c)^t (\mathbf{x}, 1) = 0\} .$$

To work with a Perceptron, we need a two-class set of linearly separable data. For the purposes of testing the algorithm, we will random-generate this data synthetically. A sample set of n data values in d -dimensional space will be represented as a $n \times (d+1)$ -matrix S , so each row is the (homogeneous coordinate) representation of a single data point. This matrix contains the data of both classes, so in addition, we need class labels. Class labels are represented as a vector \mathbf{y} of length n , with entries in $\{-1, +1\}$.

This problem requires linearly separable data. You can download the function `fakedata.R` from the class homepage, which generates such data synthetically: If \mathbf{z} is a the vector $(\mathbf{v}_H, -c)$ and n the size of the training data set, `fakedata(z, n)` returns the matrix S and the class label vector \mathbf{y} . (Alternatively, you can also write your own implementation of `fakedata`, if you are so inclined—any function that generates linearly separable, Euclidean data is fine.)

Homework questions:

1. To evaluate a Perceptron solution (a hyperplane classifier trained by a Perceptron algorithm), write a function `classify(S, z)` (where $\mathbf{z} = (v, -c)$ is the vector defining the hyperplane) and return class label

vector y as described above. S is a sample data set and v the Perceptron weight vector (which will be returned by the Perceptron training algorithm).

2. Implement the Perceptron training algorithm (see slide 38) with $\alpha(k) = \frac{1}{k}$, where k is the number of the current iteration. The training algorithm should be an R function `perceptrain(S,y)` which returns a list containing z and `Z_history`, where z is again the normal vector z_H of the hyperplane. `Z_history` is a matrix containing the history of the normal vector throughout the training run, i.e. a $(\text{number of iterations}) \times (d + 1)$ matrix, the rows of which are the interim results for z .

A recommended sanity check is to train your `perceptrain` function on a `fakedata` sample, and to make sure that the classifier returned by your Perceptron algorithm correctly classifies its own training data.

3. Generate a new 3D random vector z , run `fakedata(z,100)`, and train your Perceptron on this set. Re-run `fakedata` with the same z to produce a test set, and check whether it is classified correctly.
4. Convert the data and the vectors into their corresponding 2D representation to make it more suitable for plotting.

As a solution, please let us have:

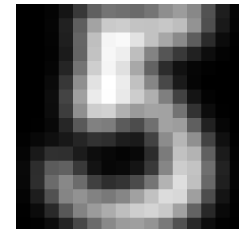
- A printout of **all** your program code.
- One plot showing the test data set and the classifier hyperplane, and one showing the training data and the trajectory of the algorithm by visualizing `Z_history`.

Problem 2 (SVM)

In this problem, we will apply a support vector machine to classify hand-written digits. You do not have to implement the SVM algorithm: The R library `e1071` provides an implementation, see

<http://cran.r-project.org/web/packages/e1071/index.html>

Download the digit data set from the course website. The zip archive contains two files: Both files are text files. The file `uspsdata.txt` contains a matrix with one data point (= vector of length 256) per row. The 256-vector in each row represents a 16×16 image of a handwritten number. The file `uspsc1.txt` contains the corresponding class labels. The data contains two classes—the digits 5 and 6—so the class labels are stored as -1 and +1, respectively. The image on the right shows the first row, re-arranged as a 16×16 matrix and plotted as a gray scale image.



- Randomly select about 20% of the data and set it aside as a test set.
- Train a linear SVM with soft margin. Cross-validate the margin parameter.
- Train an SVM with soft margin and RBF kernel. You will have to cross-validate both the soft-margin parameter and the kernel bandwidth.
- After you have selected parameter values for both algorithms, train each one with the parameter value you have chosen. Then compute the misclassification rate (the proportion of misclassified data points) on the test set.

Homework questions:

1. Plot the cross-validation estimates of the misclassification rate. Please plot the rate as
 - (a) a function of the margin parameter in the linear case.
 - (b) a function of the margin parameter and the kernel bandwidth in the non-linear case.

2. Report the test set estimates of the misclassification rates for both cases, with the parameter values you have selected, and compare the two results. Is a linear SVM a good choice for this data, or should we use a non-linear one?