

183.605

Machine Learning for Visual Computing

Assignment 1

Group 12:
Hanna Huber (0925230)
Lena Trautmann (1526567)
Elisabeth Wetzler ()

November 18, 2016

- Upload a zip-file with the required programs. You can choose the programming language.
- Add a PDF document with answers to all of the questions of the assignment (particularly all required plots) and description and discussion of results.

1 Assignment 1

1.1 Part 1: Binary classification and the perceptron

1.1.1 Reading data

Tasks:

- Read the data using functions of your programming language resp. simulation software.

First, we wrote a new .csv-file with our .py-script to delete all extra blanks. Afterwards we could read in the data using `dlmread()`.

- Plot the input vectors in \mathbb{R}^2 and visualize corresponding target values (e.g. by using color).

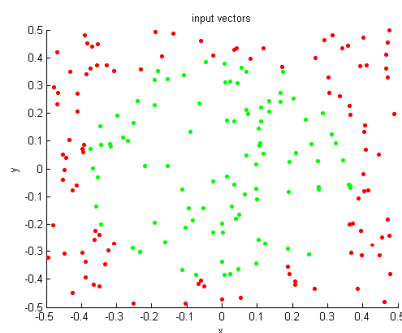


Figure 1: Plot of the input vectors with the target value visualized by colour.

Figure 1 shows the input vectors.

- Use the feature transformation $(x_1, x_2) \rightarrow (x_1^2, x_2^2)$ and plot the data in the new feature space. The data should now be linearly separable.

Figure 2 shows the transformed input vectors.

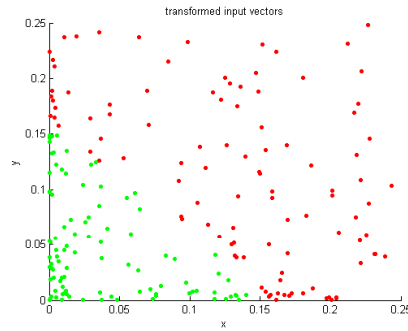


Figure 2: Plot of the transformed input vectors with the target value visualized by colour.

1.1.2 Perceptron training algorithm

Tasks:

- Implement both functions. Use homogeneous coordinates and the corresponding augmented weight vector $\mathbf{w} \in \mathbb{R}^3$.

The functions `percTrain` and `perc` are implemented in the files `percTrain.m` and `perc.m`, respectively.

- Plot the data and decision boundary in \mathbb{R}^2 , both in the original data space (see e.g. Figure 3) and in the feature space of basis functions, each at three different stages of the training: after the first iteration, after approximately half of the required iterations, and after convergence.

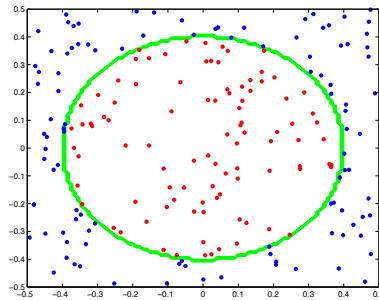


Figure 3: Plot of the decision boundary in the original data space found by the perceptron (green curve) together with labelled data points.

Figures 4, 5, 6, 7 show perceptron learning after the first iteration, half of the iterations needed and after the final iteration. They show that the algorithm converges much faster using online learning (see Figures 4 and 5). Figures 5 and 7 illustrate how the transformation into the feature space of basis functions makes the data linearly separable. The non-linear decision boundary in the original data space is obtained by applying the inverse transformation to the linear decision boundary in feature space.

- Initialize $\mathbf{w} = \mathbf{0}$. What is the influence of the learning rate?

The weight vector is initialized as $\mathbf{w} = \mathbf{0}$. This way, the learning rate γ merely scales the weight vector $\mathbf{w}^{(j)} = \gamma \sum_{i \in M} x_i t_i$, where j denotes the current iteration and M the set of data points that have been misclassified and used to update \mathbf{w} up to this point. As $(\gamma_1 \sum_{i \in M} x_i t_i)^T (x_k t_k) \leq 0 \Leftrightarrow (\gamma_2 \sum_{i \in M} x_i t_i)^T (x_k t_k) \leq 0$ for any $\gamma_1, \gamma_2 > 0$, the classification of the k th data point and thus the learning behaviour of the perceptron are not influenced by the learning rate.

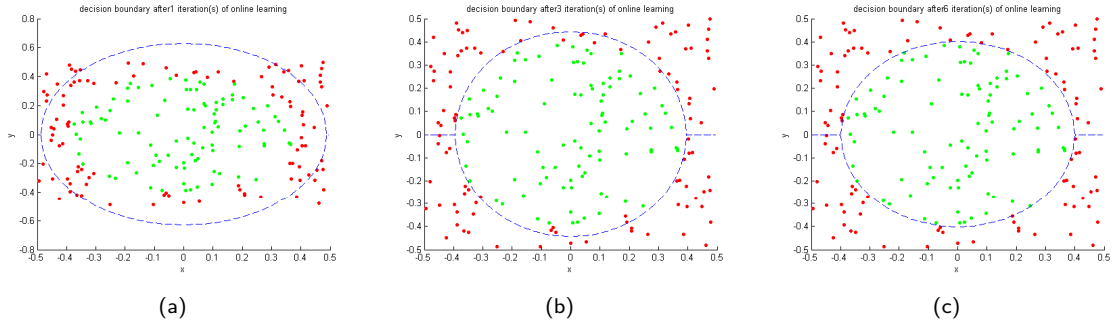


Figure 4: Perceptron decision boundary in the original data space at iterations #1, #3 and #6 of online learning.

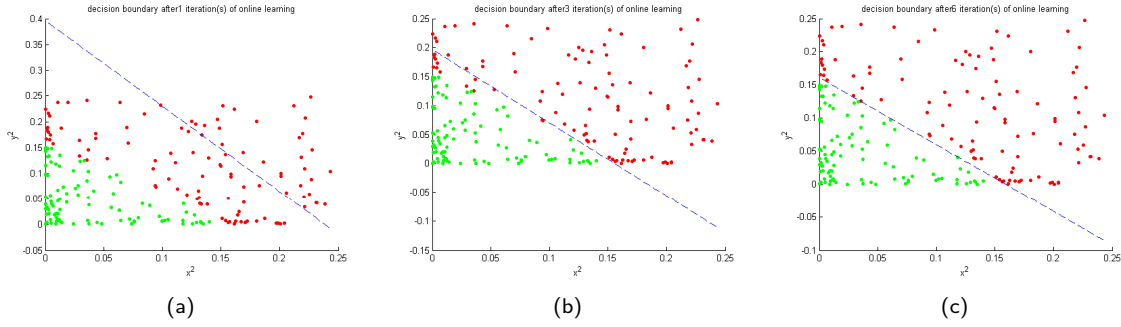


Figure 5: Perceptron decision boundary in the feature space of basis functions at iterations #1, #3 and #6 of online learning.

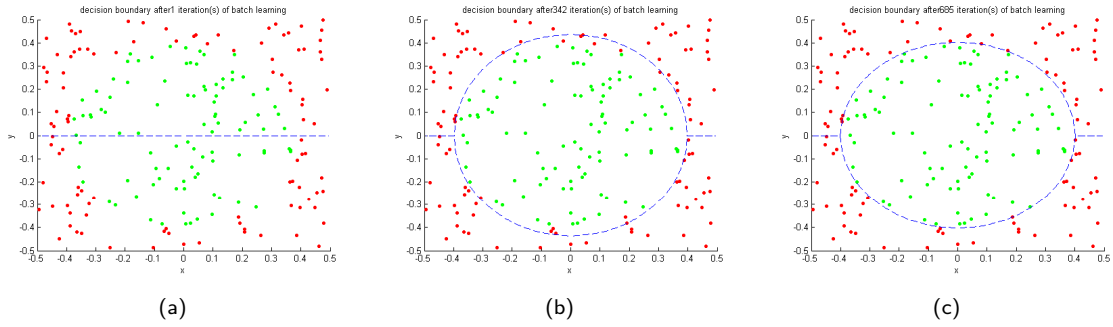


Figure 6: Perceptron decision boundary in the original data space at iterations #1, #342 and #685 of batch learning.

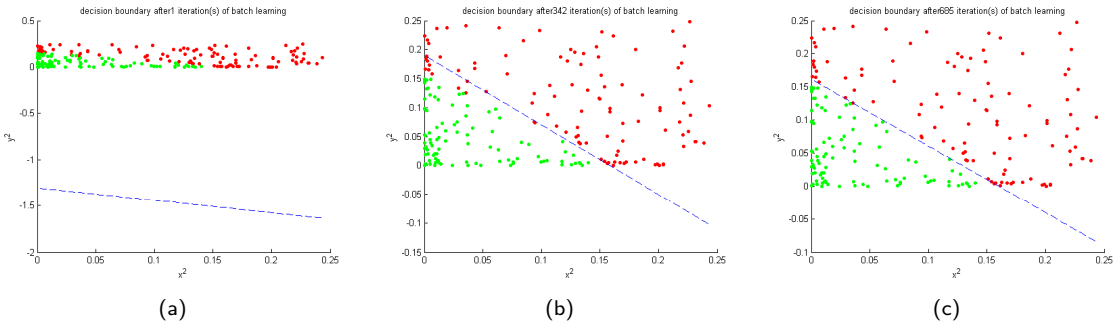


Figure 7: Perceptron decision boundary in the feature space of basis functions at iterations #1, #342 and #685 of batch learning.

1.2 Part 2: Linear basis function models for regression

1.2.1 Experimental setup

At the end of the setup we have:

Table 1: Weight vector derived for the LMS-rule (online and closed form)

	ϕ_0	ϕ_1	ϕ_2
wOnline	-2.4164	-6.7862	0.7742
wClosed	-2.4193	-6.7593	0.7645

- xtrain and ttrain: the training data
- phi: the transformation function
- xtrain_phi: the transformation of xtrain

1.2.2 Optimization: LMS-learning rule vs. closed form

Tasks:

- What is the resulting weight vector when using the LMS-rule?

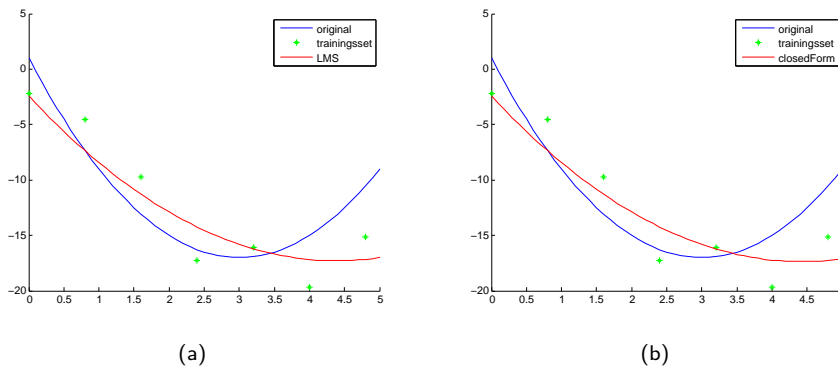


Figure 8: Optimization using the LMS-rule with online learning (a) and in closed form (b)

The values of the resulting weight vector is summarized in Table 1 and the corresponding curve plotted in Figure 8(a).

- How can you determine the optimal \mathbf{w}^* in closed form? Compare \mathbf{w}^* with the outcome of the LMS-rule training.

The closed form is calculated by updating the weight vector, determined by all N points, in one step. The resulting values are also summarized in Table 1 and the corresponding curve plotted in Figure 8(b).

- What is the influence of γ ? Which value for γ represents a good tradeoff between number of iterations and convergence?

The learning rate γ controls the influence of the single data points. Using high learning rates during online learning, the regression reacts quickly and moves towards the current data point. This is disadvantageous if the current data point is an outlier. Consequently, if the learning rate is too high, the number of iterations until convergence increases as well. The same problem occurs, if the rate is too low and the regression reacts too slowly. For this assignment, we observed that a learning rate of $\gamma = 0.0001$ is a suitable choice.

Also, the difference between online and batch learning decreases with decreasing values of γ . This is due to the fact that the weight vector \mathbf{w} only changes very little with each data point during online learning. As a result, the update factor for subsequent datapoints is also only slightly changed. This way, the difference between adding the update factors subsequently or at the same time is less remarkable.

1.2.3 Model-complexity and model-selection

Determine \mathbf{w}^* in closed form for 2000 different training sets, in which only the t_i are varied according to $\mathcal{N}(\mu = y_i, \sigma = 16)$, while the x_i remain unchanged.

Tasks¹:

¹In all tasks \mathcal{E} refers to the expected value with respect to the random variable \mathbf{w}^* , i.e. $\mathcal{E} \equiv \mathcal{E}_{\mathbf{w}^*}$

- Select a fixed x' , which is not an observation of the training set, but lies between two observations (e.g. $x' = 2$)
- Estimate the *mean squared error*

$$mse = \mathcal{E}(f(x') - f_{\mathbf{w}^*}(x'))^2, \quad (1)$$

i.e., the mean of the squared residuals of the models prediction $f_{\mathbf{w}^*}(x')$ from the true function value $f(x')$ for all $0 \leq d \leq 8$ ($d = 0$ corresponds to a constant function) using at least 2000 trials.

- Estimate by the same way the quantities $bias^2 = (f(x') - \mathcal{E}f_{\mathbf{w}^*}(x'))^2$ and $var = (f_{\mathbf{w}^*}(x') - \mathcal{E}f_{\mathbf{w}^*}(x'))^2$.
- Plot mse , $bias^2$ and var against d together in one plot. What is the relation of the quantities?
- (optional) Generate the above plots only for $d = 8$, but minimize instead of $E(\mathbf{w})$ the regularized error function

$$E_{\lambda}(\mathbf{w}) = \sum_{i=1}^N (t_i - \mathbf{w}^T \Phi(x_i))^2 + \lambda \|\mathbf{w}\|^2, \quad (2)$$

i.e. $\mathbf{w}^* = \arg \min_{\mathbf{w}} E_{\lambda}(\mathbf{w})$. Plot the quantities against λ instead of d . Hint: The minimum can be obtained in closed form (see lecture slides).