

# CSE 446: Machine Learning Winter 2018

## Assignment 3 w/o bonus questions

from  
Lukas Nies  
University of Washington

02/22/18

## Contents

<b>0</b>	<b>Policies</b>	<b>1</b>
0.1	List of Collaborators . . . . .	1
0.2	List of Acknowledgments . . . . .	1
0.3	Policies . . . . .	1
0.4	Note: Bonus not included! . . . . .	1
<b>1</b>	<b>Problem: Linear Regression on MNIST</b>	<b>2</b>
1.1	Closed Form Estimator . . . . .	2
1.2	Linear regression using gradient descent . . . . .	2
1.3	Linear Regression Using Stochastic Gradient Descent . . . . .	3
<b>2</b>	<b>Binary Classification with Logistic Regression</b>	<b>4</b>
<b>3</b>	<b>Multi-Class classification using Least Squares</b>	<b>6</b>
3.1	"One vs. all Classification" with Linear Regression . . . . .	6
<b>4</b>	<b>Probability and Maximum Likelihood Estimation</b>	<b>7</b>
4.1	Probability Review . . . . .	7
	Bibliography	8

## 0 Policies

### 0.1 List of Collaborators

My collaborator was Edith Heiter (discussed Problem 2 and 4). The development of the answers though was completely independent and individually.

### 0.2 List of Acknowledgments

None.

### 0.3 Policies

I have read and understood these policies.

### 0.4 Note: Bonus not included!

I will include some of the **bonus questions** in the same .pdf file as assignment 3 for better readability. Maybe not in this upload but in the upload until Monday 26th.

# 1 Problem: Linear Regression on MNIST

## 1.1 Closed Form Estimator

1. If one runs the Closed Form Estimator with  $\lambda = 0$  one encounters trying to invert a singular matrix ( $X^T X$ ) which is not possible per definition since the determinant is  $\det(X^T X) = 0$ . The matrix is therefore not invertible. To avoid this we introduce a regularization by adding the term  $\lambda \mathbb{1}_d$ . This is intuitively clear by considering the data itself: one digit consists of  $28 \times 28$  pixels where most pixels (at the edges and in the corners) don't carry any information about the digit itself. When calculating  $X^T X$  we get the same result: we have more "dimensions" than information for those "dimensions". In mathematical terms:  $X^T X$  is underdetermined.
2. For this part a grid search was implemented to search for different values of  $\lambda$  and the threshold to optimize the performance on the development set:
  - (a) The best result was found with  $\lambda = 0.02$  and a threshold of 0.5. The grid search ran for  $\lambda$  from 0.01 to 1 with steps of 0.01, the threshold ran from 0.1 to 1.0 in steps of 0.1.
  - (b) The average squared error using the parameters stated above is as follows:
    - Training error = 0.013045
    - Development error = 0.01420
    - Test error = 0.01626
  - (c) The misclassification error using the parameters stated above is as follows:
    - Training error = 0.93%
    - Development error = 1.08%
    - Test error = 1.76%
3. Samples with large values (far off the mean of the rest of the data points) have a strong influence on linear polynomial functions fitted through regression. This leads to large misclassification on most of the data points. A better model would be using a higher order polynomial to fit those samples more efficiently.

## 1.2 Linear regression using gradient descent

1. The proof is as follows:

$$\begin{aligned}
 \frac{\partial \mathcal{L}_\lambda}{\partial w} &= \frac{\partial}{\partial w} \left( \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (y_n - w^T x_n)^2 + \frac{\lambda}{2} \|w\|^2 \right) \\
 &= \frac{1}{N} \sum_{n=1}^N \left( -\frac{2x_n}{2} \right) (y_n - w^T x_n) + \left( \frac{2\lambda}{2} \mathbf{w} \right) \\
 &= -\frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n) x_n + \lambda \mathbf{w}
 \end{aligned}$$

2. We can rewrite this as a matrix expression:

$$\frac{\partial \mathcal{L}_w}{\partial w} = -\frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n) x_n + \lambda \mathbf{w} = -\frac{1}{N} X^T \cdot (Y - \hat{Y}) + \lambda \mathbf{w}$$

3. Stepsizes  $-10^{-2} \leq \eta \leq -10^{-1}$  worked well for this problem. For the error rate see figure 1. For generating the plots,  $\eta = \frac{1}{4} \times 10^{-1}$  and  $\eta = 10^{-2}$  were chosen. The lowest error I achieved is comparable to the closed form estimator, with

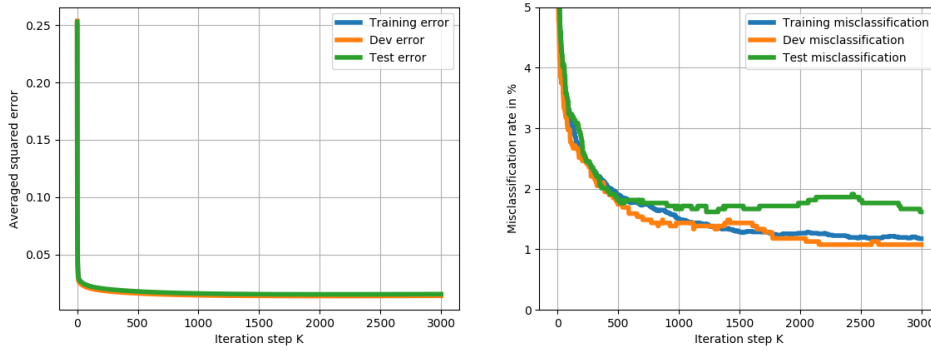


Figure 1: Plot of averaged squared errors (left) and misclassification loss in percent (right) for the gradient descent algorithm. For generating the plots,  $\eta = \frac{1}{4} \times 10^{-1}$  and  $\eta = 10^{-2}$  were chosen.

1.08% on the development set.

### 1.3 Linear Regression Using Stochastic Gradient Descent

1. The stochastic gradient descent diverges in this case with a learning rate for about  $\eta = -0.1$  at  $\lambda = 0.005$ . For too large values of  $\eta$  the algorithm might never find the global minimum and therefore the gradient gets larger and larger which leads to divergence.
2. Stepsizes  $\eta \leq -10^{-1}$  worked well for this problem. For the error rate see figure 2. For generating the plots, constant  $\eta = \frac{1}{4} \times 10^{-1}$  and  $\eta = 0.005$  were chosen. The lowest error I achieved is comparable to the closed form estimator, with 1.03% on the development set.

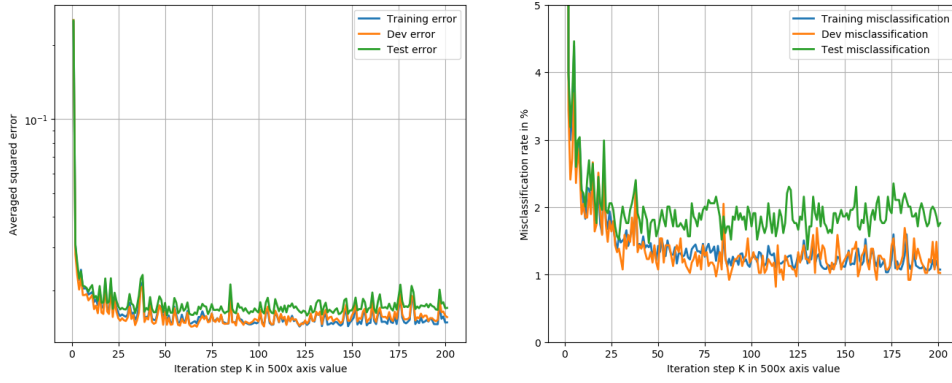


Figure 2: Plot of averaged squared errors (left, note the logarithmic vertical axis) and misclassification loss in percent (right) for the stochastic gradient descent algorithm. The horizontal axis shows the iteration steps for every 500th step. For generating the plots,  $\eta = \frac{1}{4} \times 10^{-1}$  and  $\eta = 0.005$  were chosen.

## 2 Binary Classification with Logistic Regression

1. For proofing this, we look at the cases  $y_n = 1$  and  $y_n = 0$  separately.

Case  $y_n = 1$ :

$$\begin{aligned}
 \frac{\partial \mathcal{L}_\lambda}{\partial w} &= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N \log p_w(y_n = 1 | x_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
 &= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(-wx_n)} + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
 &= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N [\log(1) - \log(1 + \exp(-wx_n))] + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
 &= -\frac{1}{N} \sum_{n=1}^N \left[ \frac{x_n \exp(-wx_n)}{1 + \exp(-wx_n)} \right] + \lambda \mathbf{w} \\
 &= -\frac{1}{N} \sum_{n=1}^N [x_n (\hat{y}_n^{-1} - 1) \hat{y}_n] + \lambda \mathbf{w} = -\frac{1}{N} \sum_{n=1}^N (1 - \hat{y}_n) x_n + \lambda \mathbf{w}
 \end{aligned}$$

Case  $y_n = 0$ :

$$\begin{aligned}
\frac{\partial \mathcal{L}_\lambda}{\partial w} &= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N \log p_w(y_n = 0 | x_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
&= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(+wx_n)} + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
&= \frac{\partial}{\partial w} \left( -\frac{1}{N} \sum_{n=1}^N [\log(1) - \log(1 + \exp(+wx_n))] + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\
&= -\frac{1}{N} \sum_{n=1}^N \left[ \frac{-x_n \exp(+wx_n)}{1 + \exp(+wx_n)} \right] + \lambda \mathbf{w} = -\frac{1}{N} \sum_{n=1}^N \left[ \frac{-x_n}{\exp(-wx_n) + 1} \right] + \lambda \mathbf{w} \\
&= -\frac{1}{N} \sum_{n=1}^N -x_n \hat{y}_n + \lambda \mathbf{w} = -\frac{1}{N} \sum_{n=1}^N (0 - \hat{y}_n) x_n + \lambda \mathbf{w}
\end{aligned}$$

2. We can rewrite this as a matrix expression:

$$\frac{\partial \mathcal{L}_\lambda}{\partial w} = -\frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n) x_n + \lambda \mathbf{w} = -\frac{1}{N} X^T \cdot (Y - \hat{Y}) + \lambda \mathbf{w}$$

3. Properties of logistic regression

- (a) Suppose the data is linear separable and  $\lambda = 0$ . In order to fit the data best, one would like to optimize the sigmoid function  $\frac{1}{1 + \exp(-wx)}$ . Since the data is linear separable all, data points with  $y_n = 0$  are left of  $x = 0$  and all points with  $y_n = 1$  are to the right. In this case, the optimal fit would be the Heaviside function (step function). In order to optimize the sigmoid function to approach the Heaviside function,  $w \rightarrow \inf$ . Hence, our weight vector would diverge.
- (b) If we suppose that  $d > n$  ( $\lambda = 0$ ) then the data matrix is sparse and several features will carry no information (equal 0). To fit a larger accumulation of features with value 0 the logistic function must approach, similar to previous question, the Heaviside function. Therefore the weight vector will diverge.
- (c) To avoid the divergence of the weight vector one can introduce regularization such that the algorithm stops early enough to give a good estimation without diverging too fast. If one does not consider this the algorithm might overfit strongly which influences the true error.

### **3 Multi-Class classification using Least Squares**

#### **3.1 "One vs. all Classification" with Linear Regression**



## 4 Probability and Maximum Likelihood Estimation

### 4.1 Probability Review

1. (a) Since the disease is quite rare it is likely that one does not have the disease even if one is tested positive. The amount of healthy persons in a group is larger than the number of sick persons and even if the test is highly accurate it's much more likely that one is healthy given the test is incorrect.
- (b) The probability of being tested positive ( $P$ ) given having the disease ( $\bar{H}$ ) is given by:

$$P(P|\bar{H}) = \frac{P(\bar{H}|P)P(P)}{P(\bar{H})} = 0.99. \quad (1)$$

By reversing this Bayesian theorem we can yield the probability having the disease given being tested positive:

$$P(P|\bar{H}) = \frac{P(\bar{H}|P)P(\bar{H})}{P(P)}, \quad (2)$$

where  $P(P)$  is the total probability being tested positive, which is the sum of the probabilities of being tested positive and being sick, or being tested positive and being healthy:

$$P(\bar{H}|P) = \frac{P(P|\bar{H})P(\bar{H})}{P(P)} = \frac{P(P|\bar{H})P(\bar{H})}{P(P|\bar{H})P(\bar{H}) + P(P|H)P(H)} \quad (3)$$

$$= \frac{0.99 \times 10^{-4}}{0.99 \times 10^{-4} + 0.01 \times (1 - 10^{-4})} = \frac{1}{102} = 0.98\% \quad (4)$$

2. We can rewrite the table as follows: It follows:

	S=0	S=1	
C=1	$P(S=0 \cap C=1)$	$P(S=1 \cap C=1)$	$P(C=1)$
C=0	$P(S=0 \cap C=0)$	$P(S=1 \cap C=0)$	$P(C=0)$
	$P(S=0)$	$P(S=1)$	1

	S=0	S=1	
C=1	$\frac{23}{154}$	$\frac{34}{154}$	$\frac{57}{154}$
C=0	$\frac{97}{154}$	$\frac{56}{154}$	$\frac{23}{154}$
	$\frac{64}{154}$	$\frac{90}{154}$	1

$$(a) \hat{p}(C = 1, S = 1) = P(C = 1 \cap S = 1) = \frac{43}{154} = 22.08\%$$

$$(b) \hat{p}(C = 1|S = 1) = \frac{P(C=1 \cap S=1)}{P(S=1)} = \frac{34}{90} = 37.78\%$$

$$(c) \hat{p}(C = 0|S = 0) = \frac{P(C=0 \cap S=0)}{P(S=0)} = \frac{41}{64} = 64.06\%$$

3. A hat over a parameter denotes an estimator of the parameter. The estimator estimates a probability value of an event based on a limited amount of samples whereas the actual probability of an event is given by looking at all samples.

## References