

---

# ART.T458 Machine Learning

## Midterm Assignment

### (Final report assigned by Shimosaka)

---

Masamichi Shimosaka  
Department of Computer Science  
Tokyo Institute of Technology

#### Instruction

- Choose the problems from your preference and solve the chosen problems.
- Use Tex / MS word.
- Submit the A4 sized / letter sized report (pdf format file) by **6th 17:00 August 2018 via OCW-i**.
- The maximum points you can earn in each problem is shown in each section title.
- Let  $s_i \geq 0$  be the score you earned in  $i$ -th problem, the final score used in the evaluation  $q$  is processed by

$$q = \min(30, \sum_i s_i).$$

- You can earn 2 pts (that is independent of  $q$ ), if you show constructive criticism to my part of this course.
- You do not need to include the source code of your implementation in your report. You are welcome to show the URL of your code via github or some other publicly available repository services.
- You can use Japanese as well as English.

#### Problem 1 (10 pts)

We consider a binary classification with a linear logistic regression. Let  $\mathbf{x} \in \mathbb{R}^d$ , and  $\mathbf{w} \in \mathbb{R}^d$  be an  $d$ -dimensional input vector, and a parameter of the model, respectively. The classifier is represented by  $f(\mathbf{x}) = 2\llbracket \mathbf{w}^\top \mathbf{x} \rrbracket - 1$ , where  $\llbracket c \rrbracket$  denotes an indicator function that returns 1 if  $c$  is true, otherwise returns 0. With the supervised dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , we consider an optimization problem for the logistic regression. The optimization problem can be written as

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$
$$J(\mathbf{w}) := \sum_{i=1}^n (\ln(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i))) + \lambda \mathbf{w}^\top \mathbf{w}.$$

With some artificial dataset (see Toy Dataset section), we consider to implement this optimization methods in the following way.

1. Implement batch steepest gradient method<sup>1</sup>.
2. Implement Newton based method
3. Compare the performance of the above two optimization methods by showing  $J(\mathbf{w}^{(t)})$  w.r.t.  $t$ , where  $\mathbf{w}^{(t)}$  represents the parameter at  $t$ -th iteration.

---

<sup>1</sup>We assume here the learning rate is constant for simplicity. Consider the upper bound of Lipsitz constant of the gradient of this objective.

## Problem 2 (10 pts + optional 10 pts)

We consider *lasso*, where the square loss, and the L1 regularization are employed for linear regression. In this problem, we employ proximal gradient method (PG). So as to make the discussion simple, we use the following objective:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left( (\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A} (\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1 \right).$$

Implement PG for lasso and show the results in a couple of conditions. In this question, use the same learning rate  $\eta_t = L^{-1}$ , where  $L$  depicts the Lipsitz constant of the gradient of the objective, which is derived from the Hessian matrix  $2\mathbf{A}$  (i.e. use the maximum eigen value of  $2\mathbf{A}$  as the inverse of the learning rate:  $\eta_t^{-1}$ ).

1. Show the result of PG in terms of  $\|\mathbf{w}^{(t)} - \hat{\mathbf{w}}\|$  of the objective function w.r.t. the number of iteration. Use semi log plot. Use the following condition:

$$\mathbf{A} = \begin{pmatrix} 3 & 0.5 \\ 0.5 & 1 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \quad 2).$$

To verify the property of L1 regularization, run the experiment with  $\lambda = 2, 4, 6$ . Recall that the numerical result can be found in the slide used in the lecture. You can also verify the result by using `cvx (matlab) / cvxopt (python)`.

2. (Option: You can earn additional 5 pts) Implement the accelerated proximal gradient method (APG) and compare the performance of APG with that of PG in the above condition.
3. (Option: You can earn additional 5 pts) Implement AdaGrad optimization and compare the performance of it with PG, and APG. Use the following condition:

$$\mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \quad 2), \lambda = 0.89.$$

Recall that the numerical result can be found in the slide used in the lecture.

## Problem 3 (20 pts + optional 20 pts)

We consider the dual of the support vector machine (L2-regularized hinge loss based binary classifier). The original optimization problem of this classification can be represented as

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \left( \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i) + \lambda \|\mathbf{w}\|_2^2 \right), \quad (1)$$

where  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\mathbf{w} \in \mathbb{R}^d$ ,  $y_i \in \{\pm 1\}$ , and  $\lambda > 0$  denotes  $i$ -th input variable, the parameter vector, and the label for  $i$ -th input data, and a coefficient of the regularization term, respectively.

1. Verify that the dual Lagrange function of this optimization can be written as

$$\begin{aligned} & \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\operatorname{maximize}} && -\frac{1}{4\lambda} \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} + \boldsymbol{\alpha}^\top \mathbf{1} \\ & \text{subject to} && \mathbf{0} \leq \boldsymbol{\alpha} \leq \mathbf{1} \end{aligned}, \quad (2)$$

where  $\mathbf{1}$  and  $\mathbf{0}$  denote a  $n$  dimensional vector whose elements are all 1 and 0, respectively, and  $\mathbf{K} \in \mathbb{R}^{n \times n}$  denotes a symmetric square matrix, and its  $i$ -th row and  $j$ -th element can be represented by  $y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$ .

2. From the KKT condition, verify the optimal weight parameter  $\mathbf{w}$  given by  $\boldsymbol{\alpha}$  can be written as

$$\hat{\mathbf{w}} = \frac{1}{2\lambda} \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (3)$$

3. Implement minimization for the negative dual Lagrange function using projected gradient (for simplicity). For the validity, show the score of the dual Lagrange function (use (2)), and sum of hinge loss function and the regularization, w.r.t. number of iteration (use (1))

and (3)), respectively. Confirm that the duality gap reaches 0 for the convergence. Here we assume projected gradient just computes

$$\alpha^{(t)} = P_{[0,1]^n} \left( \alpha^{(t-1)} - \eta_t \left( \frac{1}{2\lambda} K \alpha - \mathbf{1} \right) \right),$$

where  $\eta_t$  represents learning rate at  $t$ -th iteration, and  $P_{[0,1]^n}$  depicts a projection operator that each of the input cast into  $[0, 1]$ .

4. (Option: You can earn additional 20 pts) Investigate optimization algorithm known as dual coordinate descent technique used in LibLinear and sequential minimal optimization (SMO) also used in LibSVM, and summarize the basic concept of the two algorithms. Then run the three algorithms including projected gradient, SMO, and dual coordinate descent. Finally compare the performance in terms of the convergence speed.

#### Problem 4 (15 pts + optional 20 pts)

We consider a binary classification problem, where the hinge loss function, and L1 square regularization are leveraged. Let  $\mathbf{x} \in \mathbb{R}^d$  be an input, and  $\mathbf{w} \in \mathbb{R}^d$  be a parameter of the model, respectively. We consider a binary discriminant function with linear regression as  $f(\mathbf{x}) = 2\llbracket \mathbf{w}^\top \mathbf{x} \geq 0 \rrbracket - 1$ , where  $\llbracket c \rrbracket$  returns 1 when  $c$  is true, otherwise it returns 0, and With a supervised dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , the learning problem can be formalized as

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \left( \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i) + \lambda \|\mathbf{w}\|_1 \right), \quad (4)$$

where  $y_i \in \{\pm 1\}$  denotes a binary label for  $i$ -th training example.

1. Derive a linear program from (4) (with auxiliary variable  $\xi_i \geq \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i) \geq 0$ , and  $e_i \geq |w_i| \geq 0$ ). Recall that linear program can be written as

$$\begin{array}{ll} \text{minimize}_{\mathbf{z}} & \mathbf{c}^\top \mathbf{z} \\ \text{subject to} & \mathbf{A} \mathbf{z} \leq \mathbf{b} \end{array}$$

(See *LpBoost* that deals with LP from L1-regularized hinge loss model)

2. By using some artificial dataset (see Toy Dataset section), implement this problem via cvx (in matlab) / cvxopt (in python) (just for reference) and a (batch) proximal subgradient method. Then confirm that the parameter properly converges. The specification of the dataset should be described in the report.
3. (Option: You can earn additional 20 pts) Derive the dual Lagrange function of this linear program, then investigate *LpBoost* for automated featurization by using the dual problem, and the process how to find relevant features in a boosting way. Implement this dual program by a column generation technique. For simplicity, use Dataset II described in this document and use decision stump as single weak learner, a candidate of the feature. Show the classification accuracy as well as the number of features used in the model and the number of support vectors.

#### Problem 5 (15 pts)

We consider a binary nonlinear classification problem by using some nonlinear transformation and linear classifier. Let  $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{\pm 1\}$  be  $i$ -th training example in  $n$  sized dataset, i.e. we have a supervised dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , and  $\phi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{\tilde{d}}$  be a feature map from the original  $d$  dimensional feature space to  $\tilde{d}$  dimensional feature space. For example, use Gaussian kernel as the feature map such as

$$\{\phi(\mathbf{x})\}_k = \exp(-\alpha \|\mathbf{x} - \mathbf{x}_{z_k}\|_2^2),$$

where  $\{\phi(\mathbf{x})\}_k$  represents a  $k$ -th feature map among  $\tilde{d}$  dimensional feature maps,  $z_k$  represents an index of the training and  $\alpha > 0$  is a hyper parameter of the kernel. It should be noted that  $z_k$  is chosen by the randomly permuted sequence of indices. Solve the following problems using Dataset I (see Toy Dataset section).

1. Declare machine learning problem within linear binary classification. Use Gaussian kernel for handling nonlinear separation. Specify the loss function, regularization, respectively.
2. Implement, and run the ml problem as you declare. Compare the performance of the model by modifying the size of  $n$ ,  $\tilde{d} < n$ , and  $\alpha$ .

### Problem 6 (10 pts + optional 30 pts)

We consider a matrix optimization problem with the following objective:

$$\operatorname{argmin}_{\mathbf{Z} \in \mathbb{R}^{m \times n}} \left( \sum_{i,j \notin Q} |A_{i,j} - Z_{i,j}|^2 + \lambda \|\mathbf{Z}\|_* \right)$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  represents data matrix (user vs. movie rating data in recommendation systems), and also contains null value on  $Q$  denotes.  $\|\cdot\|_*$  represents a nuclear norm of the matrix,  $\lambda > 0$  denotes a hyper parameter for the regularization. In this optimization problem,  $\mathbf{Z}$  corresponds to the recovered data of the incomplete data matrix  $\mathbf{A}$ . In the scenario of the recommendation systems, the inferred  $\mathbf{Z}$  at location  $Q$  corresponds to the inferred score of movie rating. In this setting, we consider a nuclear

1. Describe the definition of the nuclear norm of a matrix by investigating it from www.
2. Define the proximal operation with the nuclear norm.
3. (Option: You can earn additional 15 pts) With some dataset (see Toy Dataset section), implement proximal gradient method for this machine learning problem and shows the recovered data  $\mathbf{Z}$  by using surface plotting.
4. (Option: You can earn additional 15 pts) Implement non-negative matrix factorization as alternative approach to recover  $\mathbf{A}$  and compare the performance by choosing the hyper parameters. Discuss the advantages and disadvantages of the two methods you implement here.

### Problem 7 (15 pts)

In this problem, we consider recent advances in adaptive gradient based techniques.

1. Pick up at least 3 algorithms on adaptive gradient methods other than AdaGrad, Adam methods explained in the lecture. It should be noted that you need reference (including authors, titles, and the conference information).
2. Summarize the characteristics of the three algorithms in comparison with AdaGrad, and Adam.

### Problem 8 (15 pts)

We consider an extension of regularization design in terms of structured sparsity.

1. Investigate structured sparsity and summarize the application with structured sparsity.
2. Investigate and summarize the definition of machine learning with structured sparsity such as overlapped group lasso. Recall that group lasso is the simplest machine learning problem on structured sparsity.
3. Investigate alternating direction multiplier methods (ADMM) for solving the overlapped group lasso and summarize the procedure for overlapped group lasso.

### Problem 9 (10 pts)

Consider to use publicly available dataset for your practice of machine learning.

1. Specify the dataset used here. You need to declare the URL.
2. Declare the machine learning problem (such as regression, and classification)

3. Declare the machine learning algorithms you employ here. Pick up at least three algorithms and explain the reason(s) why you choose the three algorithms among so many candidates.
4. Declare the software libraries you employ here<sup>2</sup>.
5. Report the result of the machine learning among three algorithms you employ here. It should be noted that you have to take care of the way of performance evaluation.

## Toy Datasets

Use the following datasets for some problems, if necessary. It might be better to give seed to a random number generator in the initialization step. Though the following codes are described in matlab, MS hopes students do not have difficulty in generating the similar dataset in python or other programming language implementation.

### Dataset I

```
n = 100;
x = 3 * (rand(n, 2) - 0.5);
radius = [x(:, 1).^2 + x(:, 2).^2];
y = (radius > 0.7 + 0.1 * randn(n, 1)) & (radius < 2.2 + 0.1 * randn(n, 1));
y = 2 * y - 1;
```

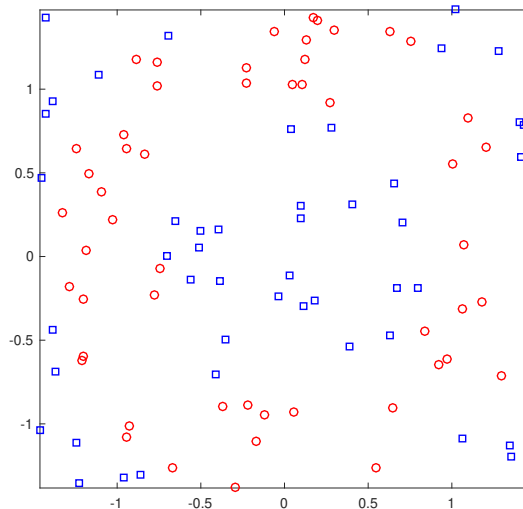


Figure 1: Dataset I

### Dataset II

```
n = 40;
omega = randn(1, 1);
noise = 0.8 * randn(n, 1);

x = randn(n, 2);
y = 2 * (omega * x(:, 1) + x(:, 2) + noise > 0) - 1;
```

---

<sup>2</sup>You can use any types of the machine learning libraries here.

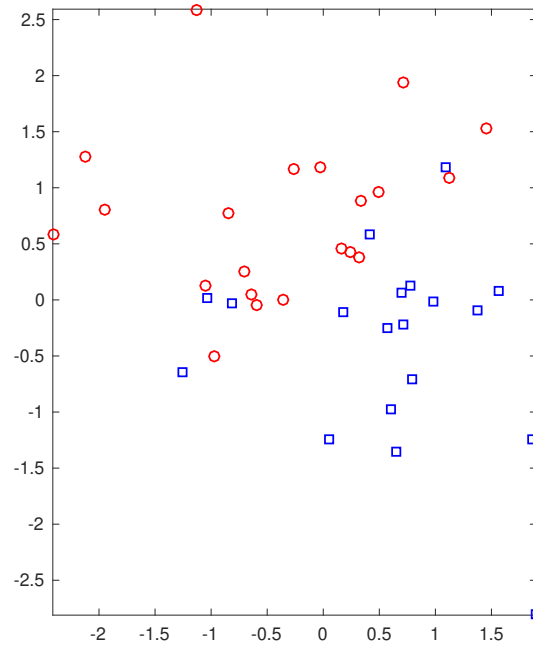


Figure 2: Dataset II

### Dataset III

```

m = 20;
n = 40;

r = 2;

A = rand(m, r) * rand(r, n);

ninc = 100;

Q = randperm(m * n, ninc);

A(Q) = NaN;

```

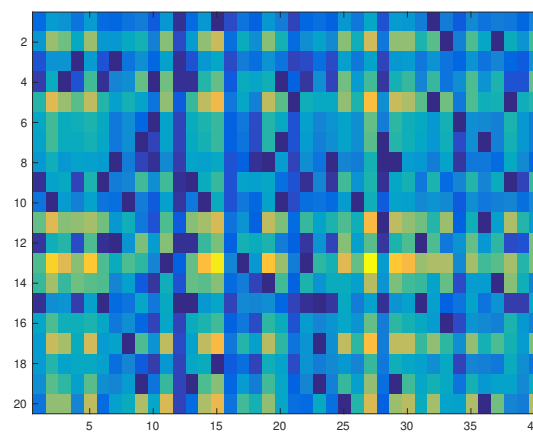


Figure 3: Dataset III