

# 10-701 Introduction to Machine Learning

## Homework 3, **version 1.3**

*Due Oct 30, 11:59 am*

---

### Rules:

1. Homework submission is done via CMU Autolab system. Please package your writeup and code into a zip or tar file, *e.g.*, let `submit.zip` contain `writeup.pdf` and your code. Submit the package to <https://autolab.cs.cmu.edu/courses/10701-f15>.
  2. Like conference websites, repeated submission is allowed. Please feel free to refine your answers since we will only grade the latest version. Submitting incomplete solutions early will be helpful in preventing last minute panic as well.
  3. Autolab may allow submission after the deadline, note however it is because of the late day policy. Please see course website for policy on late submission.
  4. We recommend that you typeset your homework using appropriate software such as L<sup>A</sup>T<sub>E</sub>X. If you are writing please make sure your homework is cleanly written up and legible. The TAs will not invest undue effort to decrypt bad handwriting.
  5. You are allowed to collaborate on the homework, but you should write up your own solution and code. Please indicate your collaborators in your submission.
-

# 1 Neural Networks (50 Points) (Zhiting)

## 1.1 Neural network for regression

Figure.1 shows a two-layer neural network which learns a function  $f : X \rightarrow Y$  where  $X = (X_1, X_2) \in \mathbb{R}^2$ . The weights  $\mathbf{w} = \{w_1, \dots, w_6\}$  can be arbitrary. There are two possible choices for the function implemented by each unit in this network:

- S: signed sigmoid function  $S(a) = \text{sign}[\sigma(a) - 0.5] = \text{sign}[\frac{1}{1+\exp\{-a\}} - 0.5]$
- L: linear function  $L(a) = ca$

where in both cases  $a = \sum_i w_i X_i$ .

1. Assign proper activation functions (S or L) to each unit in Figure.1 so this neural network simulates a linear regression:  $Y = \beta_1 X_1 + \beta_2 X_2$ .
2. Assign proper activation functions (S or L) for each unit in Figure.1 so this neural network simulates a binary logistic regression classifier:  $Y = \arg \max_y P(Y = y|X)$ , where  $P(Y = 1|X) = \frac{\exp(\beta_1 X_1 + \beta_2 X_2)}{1 + \exp(\beta_1 X_1 + \beta_2 X_2)}$ , and  $P(Y = -1|X) = \frac{1}{1 + \exp(\beta_1 X_1 + \beta_2 X_2)}$ . Derive  $\beta_1$  and  $\beta_2$  in terms of  $w_1, \dots, w_6$ .
3. Assign proper activation functions (S or L) to each unit in Figure.1 so this neural network simulates a boosting classifier which combines two logistic regression classifiers,  $f_1 : X \rightarrow Y_1$  and  $f_2 : X \rightarrow Y_2$ , to produce its final prediction:  $Y = \text{sign}[\alpha_1 Y_1 + \alpha_2 Y_2]$ . Use the same distribution in problem 1.1.2 for  $f_1$  and  $f_2$ . Derive  $\alpha_1$  and  $\alpha_2$  in terms of  $w_1, \dots, w_6$ .

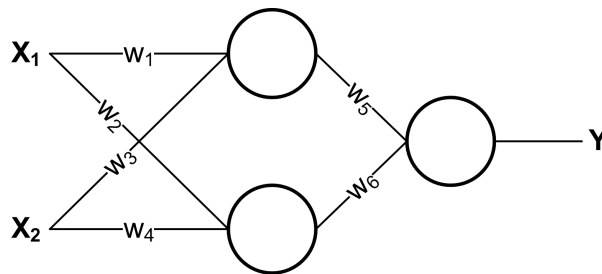


Figure 1: A two-layer neural network.

## 1.2 Convolutional neural networks

1. Count the total number of parameters in LeNet (pp.46, slides of Lecture.8). How many parameters in all of the convolutional layers? How many parameters in all of the fully-connected layers?

Note:

- (a) The filter size of each convolutional and pooling(subsampling) layer:  
C1:  $5 \times 5$  (i.e., each unit of C1 has a  $5 \times 5$  receptive field in its preceding layer);  
S2:  $2 \times 2$ ;  
C3:  $5 \times 5$ ;  
S4:  $2 \times 2$ ;
- (b) Fully-connected layers in LeNet include C5, F6, and OUTPUT

2. In a convolutional layer the units are organized into planes, each of which is called a feature map. The units within a feature map (indexed  $q$ ) have different inputs, but all share a common weight vector,  $\mathbf{w}^{(q)}$ . A convolutional network is usually trained through backpropagation. Let  $J^{(q)}$  be the number of units in the  $q$ th feature map,  $z_j^{(q)}$  the activation of the  $j$ th unit,  $x_{ji}^{(q)}$  the  $i$ th input for the  $j$ th unit,  $w_i^{(q)}$  the  $i$ th element of  $\mathbf{w}^{(q)}$ ,  $L$  the training loss. Derive the gradient of  $w_i^{(q)}$ .

### 1.3 Gradient vanishing/explosion

In this problem we will study the difficulty of back-propagation in training deep neural networks. For simplicity, we consider the simplest deep neural network: one with just a single neuron in each layer, where the output of the neuron in the  $j$ th layer is  $z_j = \sigma(a_j) = \sigma(w_j z_{j-1} + b_j)$ . Here  $\sigma$  is some activation function whose derivative on  $x$  is  $\sigma'(x)$ . Let  $m$  be the number of layers in the neural network,  $L$  the training loss.

1. Derive the derivative of  $L$  w.r.t.  $b_1$  (the bias of the neuron in the first layer).
2. Assume the activation function is the usual sigmoid function  $\sigma(x) = 1/(1 + \exp\{-x\})$ . The weights  $\mathbf{w}$  are initialized to be  $|w_j| < 1$  ( $j = 1, \dots, m$ ).
  - (a) Explain why the above gradient ( $\partial L / \partial b_1$ ) tends to vanish ( $\rightarrow 0$ ) when  $m$  is large.
  - (b) Even if  $|w|$  is large, the above gradient would also tend to vanish, rather than explode ( $\rightarrow \infty$ ). Explain why. (A rigorous proof is not required.)
3. One of the approaches to (partially) address the gradient vanishing/explosion problem is to use the rectified linear (ReLU) activation function instead of the sigmoid. The ReL activation function is  $\sigma(x) = \max\{0, x\}$ . Explain why ReL can alleviate the gradient vanishing problem as faced by sigmoid.
4. A second approach to (partially) address the gradient vanishing/explosion problem is layer-wise pre-training. Restricted Boltzmann machine (RBM) is one of the widely-used models for layer-wise pre-training. Figure 2 shows an example of RBM which includes  $K$  hidden units  $\mathbf{h}$ , and  $J$  input units  $\mathbf{v}$ . Let us define the joint distribution as the following general form:

$$P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp \left( \sum_i \theta_i \phi_i(\mathbf{v}, \mathbf{h}) \right), \quad (1)$$

where  $Z = \sum_{\mathbf{v}, \mathbf{h}} \exp(\sum_i \theta_i \phi_i(\mathbf{v}, \mathbf{h}))$  is the normalization term;  $\phi_i(\mathbf{v}, \mathbf{h})$  are some features;  $\theta_i$  are the parameters corresponding to the weights in the RBM. Consider the simplest learning algorithm, gradient descent. Show that

$$\frac{\partial \log P(\mathbf{v})}{\partial \theta_i} = \sum_{\mathbf{h}} \phi_i(\mathbf{v}, \mathbf{h}) P(\mathbf{h}|\mathbf{v}) - \sum_{\mathbf{v}, \mathbf{h}} \phi_i(\mathbf{v}, \mathbf{h}) P(\mathbf{v}, \mathbf{h}). \quad (2)$$

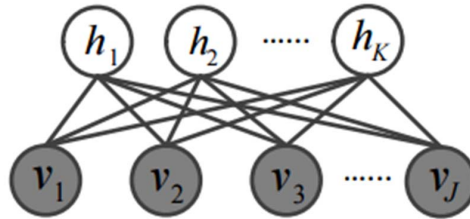


Figure 2: A restricted Boltzmann machine.

## 2 Support Vector Machines (50 Points) (Yuntian)

### 2.1 Support Vector Regression (25 Points)

We now extend support vector machines (SVM) to regression problems. Recall that in regression problems, we have  $n$  data points  $(x_i, y_i)_{i=1}^n$  where  $x_i \in \mathbb{R}^m$  and  $y_i \in \mathbb{R}$ . Given a function class  $\mathcal{F}$  (e.g. linear or quadratic functions), we want to fit a function  $f \in \mathcal{F}$  on the training set:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} C \sum_{i=1}^n l(f(x_i), y_i) + R(f) \quad (3)$$

where  $l(\cdot, \cdot)$  is the loss function,  $R(f)$  is the regularization term,  $C$  controls the regularization strength. The first part tries to fit data, and the second part penalizes complex  $f$  to avoid over-fitting.

In the support vector regression (SVR) framework, we consider linear function class  $\mathcal{F} = \{x \rightarrow w^T x\}$  (we do not consider interception term for simplicity). We use  $\ell_2$ -regularizer  $R(f) = \frac{1}{2} \|w\|_2^2$  for  $f(x) = w^T x$ . For the loss function  $l$ , similar to the hinge-loss function in SVM classification, we employ an  $\epsilon$ -insensitive error function

$$l_\epsilon(f(x), y) = \begin{cases} 0 & \text{if } |f(x) - y| < \epsilon \\ |f(x) - y| - \epsilon & \text{otherwise.} \end{cases} \quad (4)$$

Then we get the following optimization problem:

$$w^* = \operatorname{argmin}_w C \sum_{i=1}^n l_\epsilon(w^T x_i, y_i) + \frac{1}{2} \|w\|_2^2. \quad (5)$$

1. Write down the dual problem of SVR. (Hint: follow the derivations for SVM)
2. Write down the KKT conditions, and explain what are the “support vectors”.
3. Derive a kernelized version of SVR. For a test point  $x$ , write down the prediction rule.
4. Give one reason why do we usually solve the dual problem of SVR and SVM instead of the primal.
5. Implement SVR on a **1-D toy dataset**. Each line of the dataset contains a training instance  $(x_i, y_i)$  (separated by a tab). For this problem, you need to
  - Use RBF kernel  $k(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|_2^2}{2h^2})$ , and take  $h = 0.5$ ,  $C = 4$ ,  $\epsilon = 0.1$ .
  - Plot the prediction curve for  $x \in [0, 1]$  and show the support vectors versus other training points in the training dataset.

(Hint: You are allowed to use optimization toolkits such as CVX or Matlab’s inbuilt function quadprog to solve the dual problem.)

### 2.2 Support Kernel Machines (20 Points)

In SVM, the kernel function can be viewed as a similarity measure between data points. In some classification scenarios, features may come from different sources or modalities, e.g. in some tasks the data may contain both image features and text features. In that case, since these are different representations, they have different measures of similarity corresponding to different kernels. In such a case, we want to learn a combination of kernels instead of using a single kernel. There is significant amount of work in combining kernels, here we adapt the notations in [1].

We begin by considering a linear case of Support Kernel Machine (SKM). Suppose the data points  $x_i \in \mathcal{X} = \mathbb{R}^k$ . We also assume we are given a decomposition of  $\mathbb{R}^k = \mathbb{R}^{k_1} \times \cdots \times \mathbb{R}^{k_m}$ , so that each data

point  $x_i$  can be decomposed into  $m$  block components, i.e.  $x_i = (x_{1i}, \dots, x_{mi})$  where each  $x_{ji}$  is in general a vector. In real tasks, each block may correspond to a certain kind of representation, e.g.  $x_{1i}$  may correspond to image features and  $x_{2i}$  may be text features.

Our goal is to find a linear classifier of the form  $y = \text{sign}(w^T x + b)$  where  $w$  has the same block decomposition  $w = (w_1, \dots, w_m) \in \mathbb{R}^{k_1 + \dots + k_m}$ . Recall that in linear SVM the objective is:

$$\begin{aligned} & \underset{\substack{w \in \mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_m} \\ \xi_i \geq 0, b \in \mathbb{R}}}{\text{minimize}} & \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \xi_i \end{aligned} \quad (6)$$

$$\text{subject to} \quad y_i \left( \sum_j w_j^T x_{ji} + b \right) \geq 1 - \xi_i, \quad \forall i \in \{1, \dots, n\} \quad (7)$$

In SKM, we encourage the sparsity of the vector  $w$  at the level of blocks. The primal problem for the SKM is defined as:

$$\begin{aligned} & \underset{\substack{w \in \mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_m} \\ \xi_i \geq 0, b \in \mathbb{R}}}{\text{minimize}} & \frac{1}{2} \left( \sum_{j=1}^m d_j \|w_j\|_2 \right)^2 + C \sum_{i=1}^n \xi_i \end{aligned} \quad (8)$$

$$\text{subject to} \quad y_i \left( \sum_j w_j^T x_{ji} + b \right) \geq 1 - \xi_i, \quad \forall i \in \{1, \dots, n\} \quad (9)$$

where  $d_j > 0$  can be seen as constant.

1. By introducing dual variables  $\alpha_i \geq 0$  and  $\beta_i \geq 0$ , we get the Lagrangian function

$$\mathcal{L} = \frac{1}{2} \left( \sum_{j=1}^m d_j \|w_j\|_2 \right)^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i \left( y_i \left( \sum_{j=1}^m w_j^T x_{ji} + b \right) - 1 + \xi_i \right) - \sum_{i=1}^n \beta_i \xi_i \quad (10)$$

Denote  $\gamma = \sum_{j=1}^m d_j \|w_j\|_2$ .

- (a) Show that at the **minimum** of the Lagrangian function, i.e.  $w = \text{argmin}_w \mathcal{L}$  for this and the following questions,

$$\|w_j\|_2 d_j \gamma = w_j^T \sum_{i=1}^n \alpha_i y_i x_{ji}, \quad \forall j \in \{1, \dots, m\} \quad (11)$$

- (b) Show that  $\left\| \sum_{i=1}^n \alpha_i y_i x_{ji} \right\|_2 \leq d_j \gamma, \forall j \in \{1, \dots, m\}$ .

Note: for (b) you can get full credit if you only consider  $w_j \neq 0$ , but you can get 5 extra points if you include  $w_j = 0$  case in your proof. A hint is that  $\mathcal{L}$  is not differentiable w.r.t.  $w_j$  if  $w_j = 0$ , and you may refer to [this link](#) for how to deal with that case by using  $\partial \|x\|_2 = \{g : \|g\|_2 \leq 1\}$  if  $x = 0$ .

- (c) Show that

- if  $\left\| \sum_i \alpha_i y_i x_{ji} \right\|_2 < d_j \gamma$ , then  $w_j = 0$ ,
- if  $\left\| \sum_i \alpha_i y_i x_{ji} \right\|_2 = d_j \gamma$ , then  $\exists \eta_j > 0$ , such that  $w_j = \eta_j \sum_i \alpha_i y_i x_{ji}$ .

2. Recall from homework 2 that  $\ell_1$  norm can encourage sparsity. Explain the effect of the regularization term  $\frac{1}{2} \left( \sum_{j=1}^m d_j \|w_j\|_2 \right)^2$ .
3. Now we extend the above analysis to a kernelized version. Assume that we have a mapping  $\phi : \mathcal{X} \rightarrow \mathbb{R}^f$  which is generally a non-linear function. We assume that  $\phi(x)$  has  $m$  block components  $\phi(x) = (\phi_1(x), \dots, \phi_m(x))$ , and we also assume  $w$  has the same decomposition  $w = (w_1, \dots, w_m)$ . Show that at the **minimum** of the Lagrangian function,  $\exists \eta_j \geq 0$  such that  $w_j = \eta_j \sum_{i=1}^n \alpha_i y_i \phi_j(x_i)$ .

## 2.3 SVM Error Analysis (5 Points)

In this problem, we want to analyze the error of SVM classification. Assume that we have  $n$  data points  $(x_i, y_i)_{i=1}^n$  where  $x_i \in \mathbb{R}^m$  and  $y_i = 1, \dots, K$ . Assume that we train an SVM classifier  $f_{(x_1, y_1), \dots, (x_n, y_n)}$  on these  $n$  data points.

For a randomly drawn test data point  $(x_{n+1}, y_{n+1})$ , the prediction is  $y_{n+1}^{\text{pred}} = f_{(x_1, y_1), \dots, (x_n, y_n)}(x_{n+1})$ . We assume that the  $n$  training data points and the test data point  $(x_{n+1}, y_{n+1})$  are drawn i.i.d from some unknown underlying distribution. The expected error rate is defined as:

$$\text{err} = \mathbb{E}_{(x_1, y_1), \dots, (x_n, y_n)} \mathbb{E}_{(x_{n+1}, y_{n+1})} [\mathbf{1} \{f_{(x_1, y_1), \dots, (x_n, y_n)}(x_{n+1}) \neq y_{n+1}\}] \quad (12)$$

where the indicator function  $\mathbf{1} \{A\} = 1$  if  $A$  is true, otherwise 0.

1. Show the the expected error rate is equal to the expectation of leave-one-out cross validation error for  $n + 1$  data points.
2. In the lecture, we have the statement that “the leave-one-out cross-validation error does not depend on the dimensionality of the feature space but only on the number of support vectors”. Show that this statement is true by explaining why

$$\text{err}_{\text{loocv}} \leq \frac{n_s}{n + 1} \quad (13)$$

where  $\text{err}_{\text{loocv}}$  is the leave-one-out cross-validation error for training set  $(x_i, y_i)_{i=1}^{n+1}$ ,  $n_s$  is the number of support vectors.

## References

- [1] Francis R Bach, Gert R.G. Lanckriet, and Michael I Jordan. Multiple kernel learning, conic duality, and the smo algorithm. In *ICML*, 2004. [4](#)