

Homework #5

CS 260: Machine Learning Algorithms

Prof. Ameet Talwalkar

Due: 11/12/15, 8am

Please abide by the [Academic Integrity Policy](#)

1 Bias-Variance Tradeoff

Consider a dataset with n data points (\mathbf{x}_i, y_i) , $\mathbf{x}_i \in \mathbb{R}^{p \times 1}$, drawn from the following linear model:

$$y = \mathbf{x}^\top \boldsymbol{\beta}^* + \varepsilon,$$

where ε is a Gaussian noise and the star sign is used to differentiate the true parameter from the estimators that will be introduced later. Consider the L_2 regularized linear regression as follows:

$$\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right\},$$

where $\lambda \geq 0$ is the regularization parameter. Let $X \in \mathbb{R}^{n \times p}$ denote the matrix obtained by stacking \mathbf{x}_i^\top in each row. [Properties of an affine transformation](#) of a gaussian random variable will be useful throughout this problem.

- Find the closed form solution for $\hat{\boldsymbol{\beta}}_\lambda$ and its distribution.
- Calculate the bias term $\mathbb{E}[\mathbf{x}^\top \hat{\boldsymbol{\beta}}_\lambda] - \mathbf{x}^\top \boldsymbol{\beta}^*$ as a function of λ and some *fixed* test point \mathbf{x} .
- Calculate the variance term $\mathbb{E} \left[\left(\mathbf{x}^\top \hat{\boldsymbol{\beta}}_\lambda - \mathbb{E}[\mathbf{x}^\top \hat{\boldsymbol{\beta}}_\lambda] \right)^2 \right]$ as a function of λ and some *fixed* test point \mathbf{x} .
- Use the results from parts (b) and (c) and the bias-variance theorem to analyze the impact of λ in the squared error. Specifically, which term dominates when λ is small or large?

2 Kernelized Perceptron

Given a set of training samples $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ where $y \in \{-1, 1\}$, the Perceptron algorithm learns a weight vector \mathbf{w} by iterating through all training samples. For each \mathbf{x}_i , if the prediction is incorrect, we update \mathbf{w} by $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$. Now we would like to *kernelize* the Perceptron algorithm. Assume we map \mathbf{x} to $\varphi(\mathbf{x})$ through a nonlinear feature mapping φ , and we want to learn a new weight vector \mathbf{w} that makes prediction by $y = \text{sign}(\mathbf{w}^\top \varphi(\mathbf{x}))$. Further assume that we initial the algorithm with $\mathbf{w} = 0$.

- Show that \mathbf{w} is always a linear combination of feature vectors, i.e. $\mathbf{w} = \sum_{i=1}^N \alpha_i \varphi(\mathbf{x}_i)$.
- Show that while the update rule for \mathbf{w} for a kernelized Perceptron does depend on the explicit feature mapping $\varphi(\mathbf{x})$, the prediction can be re-expressed and thus depends only on the inner products between nonlinear transformed features.
- Show that we do not need to explicitly store \mathbf{w} at training or test time. Instead, we can implicitly use it by maintaining all the α_i . Please give the outline of the algorithm that would allow us to not store \mathbf{w} . You should indicate how α_i is initialized, when to update α_i , and how it is updated.

3 Kernels

Mercer's theorem implies that a bivariate function $k(\cdot, \cdot)$ is a positive definite kernel function iff, for any N and any $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, the corresponding kernel matrix K is positive semidefinite, where $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. Recall that a matrix $A \in \mathbb{R}^{n \times n}$ is positive semidefinite if all of its eigenvalues are non-negative, or equivalently, if $\mathbf{x}^\top A \mathbf{x} \geq 0$ for arbitrary vector $\mathbf{x} \in \mathbb{R}^{n \times 1}$.

Suppose $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are positive definite kernel functions with corresponding kernel matrices K_1 and K_2 . Use Mercer's theorem to show that the following kernel functions are positive definite.

- a. $K_3 = a_1 K_1 + a_2 K_2$, for $a_1, a_2 \geq 0$.
- b. K_4 defined by $k_4(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}')$ where $f(\cdot)$ is an arbitrary real valued function.
- c. K_5 defined by $k_5(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$.

4 Soft Margin Hyperplanes

The function of the slack variables used in the optimization problem for soft margin hyperplanes has the form: $\xi \mapsto \sum_{i=1}^n \xi_i$. Instead, we could use $\xi \mapsto \sum_{i=1}^n \xi_i^p$, with $p > 1$.

- a. Give the dual formulation of the problem in this general case.
- b. How does this more general formulation ($p > 1$) compare to the standard setting ($p = 1$) discussed in lecture? Is the general formulation more or less complex? Justify your answer.

5 Programming

In this problem, you will experiment with SVMs on a real-world dataset. You will implement a linear SVM (i.e., an SVM using the original features. You will also use a widely used SVM toolbox called LibSVM to experiment with kernel SVMs.

Dataset: We have provided the *Splice Dataset* from UCI's machine learning data repository.¹ The provided binary classification dataset has 60 input features, and the training and test sets contain 1,000 and 2,175 samples, respectively (the files are called `splice_train.mat` and `splice_test.mat`).

5.1 Data preprocessing

Preprocess the training and test data by

- a. computing the mean of each dimension and subtracting it from each dimension
- b. dividing each dimension by its standard deviation

Notice that the mean and standard deviation should be estimated from the *training data* and then applied to both datasets. Explain why this is the case. Also, report the mean and the standard deviation of the third and 10th features on the test data.

¹<https://archive.ics.uci.edu/ml/datasets/Molecular+Biology+%28Splice-junction+Gene+Sequences%29>.

5.2 Implement linear SVM

Please fill in the Matlab functions `trainsvm` in `trainsvm.m` and `testsvm` in `testsvm.m`.

The input of `trainsvm` contain training feature vectors and labels, as well as the tradeoff parameter C . The output of `trainsvm` contain the SVM parameters (weight vector and bias). In your implementation, you need to solve SVM in its primal form

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i \\ & \xi_i \geq 0, \forall i \end{aligned}$$

Please use the `quadprog` function in Matlab to solve the above quadratic problem.

For `testsvm`, the input contains testing feature vectors and labels, as well as SVM parameters. The output contains the test accuracy.

5.3 Cross validation for linear SVM

Use 5-fold cross validation to select the optimal C for your implementation of linear SVM.

- Report the cross-validation accuracy (averaged accuracy over each validation set) and average training time (averaged over each training subset) on different C taken from $\{4^{-6}, 4^{-5}, \dots, 4, 4^2\}$. How does the value of C affect the cross validation accuracy and average training time? Explain your observation.
- Which C do you choose based on the cross validation results?
- For the selected C , report the test accuracy.

5.4 Use linear SVM in LibSVM

LibSVM is widely used toolbox for SVMs, and it has a Matlab interface. Download LibSVM from <http://www.csie.ntu.edu.tw/~cjlin/libsvm/> and install it according to the README file (make sure to use the Matlab interface provided in the LibSVM toolbox). For each C from $\{4^{-6}, 4^{-5}, \dots, 4, 4^2\}$, apply 5-fold cross validation (use `-v` option in LibSVM) and report the cross validation accuracy and average training time.

- Is the cross validation accuracy the same as that in 5.3? Note that LibSVM solves linear SVM in dual form while your implementation does it in primal form.
- How does LibSVM compare with your implementation in terms of training time?

5.5 Use kernel SVM in LibSVM

LibSVM supports a number of kernel types. Here you need to experiment with the polynomial kernel and RBF (Radial Basis Function) kernel.

- Polynomial kernel.** Please tune C and `degree` in the kernel. For each combination of (C, degree) , where $C \in \{4^{-3}, 4^{-4}, \dots, 4^6, 4^7\}$ and `degree` $\in \{1, 2, 3\}$, report the 5-fold cross validation accuracy and average training time.
- RBF kernel.** Please tune C and `gamma` in the kernel. For each combination of (C, gamma) , where $C \in \{4^{-3}, 4^{-4}, \dots, 4^6, 4^7\}$ and `gamma` $\in \{4^{-7}, 4^{-6}, \dots, 4^{-1}, 4^{-2}\}$, report the 5-fold cross validation accuracy and average training time.

Based on the cross validation results of Polynomial and RBF kernel, which kernel type and kernel parameters will you choose? Report the corresponding test accuracy for the configuration with the highest cross validation accuracy.

Submission

Please provide the following as part of your submission:

- Provide your answers to problems 1-4, 5.1, and 5.3-5.5 in hardcopy. The papers need to be stapled and submitted at the beginning of class on the due date.
- Please put all of your code in a single folder named `[lastname]_[firstname]_hw5`, and submit a single `.zip` file containing this folder called `[lastname]_[firstname]_hw5.zip`. The only acceptable languages are MATLAB and Octave.
- You MUST include the main function called `CS260_hw5.m` in your root folder in your zip file. After running this main file, your program should be able to generate all of the results needed for this programming assignment, either as plots or console outputs. You can have multiple files (i.e your sub-functions), however, we require that once we unzip your folder and execute your main file, your program should execute correctly. Please double-check your program before submitting.
- Please submit this zip file by the due date by following the instructions on [this form](#).
- You are encouraged to collaborate, but collaboration must be limited to discussion only and you need to write down / implement your solution on your own. You also need to list with whom you have discussed the HW problems.