

# STAT 6910-001 – Principles of ML – Homework #4

Due: 5:00 PM 10/25/19

## 1. Conceptual questions (8 pts).

- (a) (4 pts) Linear classifiers. Discuss the differences between LDA, Logistic Regression, separating hyperplanes, and the Optimal Soft-Margin Hyperplane. In what situations would you use each of the classifiers? *Hint:* You may find the discussion in sections 4.5.2 and 12.2 in the ESL book or sections 9.1, 9.2, and 9.5 in the ISL book helpful.
- (b) (4 pts) Describe how you would apply the SVM to the multiclass case.

## 2. Kernels (22 pts).

- (a) (4 pts) To what feature map  $\Phi$  does the kernel

$$k(\mathbf{u}, \mathbf{v}) = (\langle \mathbf{u}, \mathbf{v} \rangle + 1)^3$$

correspond? Assume the inputs have an arbitrary dimension  $d$  and the inner product is the dot product.

- (b) (18 pts) Let  $k_1, k_2$  be symmetric, positive-definite kernels over  $\mathbb{R}^D \times \mathbb{R}^D$ , let  $a \in \mathbb{R}^+$  be a positive real number, let  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  be a real-valued function, and let  $p : \mathbb{R} \rightarrow \mathbb{R}$  be a polynomial with positive coefficients. For each of the functions  $k$  below, state whether it is necessarily a positive-definite kernel. If you think it is, prove it. If you think it is not, give a counterexample.
  - i.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
  - ii.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) - k_2(\mathbf{x}, \mathbf{z})$
  - iii.  $k(\mathbf{x}, \mathbf{z}) = ak_1(\mathbf{x}, \mathbf{z})$
  - iv.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z})k_2(\mathbf{x}, \mathbf{z})$
  - v.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z})$
  - vi.  $k(\mathbf{x}, \mathbf{z}) = p(k_1(\mathbf{x}, \mathbf{z}))$

- 3. **Alternative OSM hyperplane (10 pts).** An alternative way to extend the max-margin hyperplane to nonseparable data is to solve the following quadratic program (another name for an optimization problem with a quadratic objective function):

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

The only difference with respect to the OSM hyperplane is that we are now squaring the slack variables. This assigns a stronger penalty to data points that violate the margin.

- (a) (4 pts) Which loss is associated with the above quadratic program? In other words, show that learning a hyperplane by the above optimization problem is equivalent to ERM with a certain loss. **Hint:** you should not use the Lagrangian or KKT conditions to do this. Instead, look at the notes where we first talked about the OSM. In those notes, we argued that the OSM optimization problem is equivalent to regularized ERM with the hinge loss. You should follow a similar procedure.
- (b) (4 pts) Argue that the second set of constraints can be dropped without changing the solution.

- (c) (2 pts) Identify an advantage and a disadvantage of this loss compared to the hinge loss.

4. **Subgradient methods for the optimal soft margin hyperplane (25 pts).**

In this problem you will implement the subgradient and stochastic subgradient methods for minimizing the convex but nondifferentiable function

$$J(b, \mathbf{w}) = \frac{1}{n} \sum_{i=1}^n L(y_i, \mathbf{w}^T \mathbf{x}_i + b) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where  $L(y, t) = \max(0, 1 - yt)$  is the hinge loss. As we saw in class, this corresponds to the optimal soft margin hyperplane.

- (a) (5 pts) Determine  $J_i(b, \mathbf{w})$  such that

$$J(b, \mathbf{w}) = \sum_{i=1}^n J_i(b, \mathbf{w}).$$

Determine a subgradient  $\mathbf{u}_i$  of each  $J_i$  with respect to the variable  $\boldsymbol{\theta} = [b \quad \mathbf{w}]^T$ . A subgradient of  $J$  is then  $\sum_i \mathbf{u}_i$ .

*Note:* Recall that if  $f(\mathbf{z}) = g(h(\mathbf{z}))$  where  $g : \mathbb{R} \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ , and both  $g$  and  $h$  are differentiable, then the chain rule gives

$$\nabla f(\mathbf{z}) = \nabla h(\mathbf{z}) \cdot g'(h(\mathbf{z})).$$

If  $g$  is convex and  $h$  is differentiable, then the same formula gives a subgradient of  $f$  at  $\mathbf{z}$  where  $g'(h(\mathbf{z}))$  is replaced by a subgradient of  $g$  at  $h(\mathbf{z})$ .

- (b) (6 pts) Download the file `nuclear.mat` from Canvas. The variables  $x$  and  $y$  contain training data for a binary classification problem. The variables correspond to the total energy and tail energy of waveforms produced by a nuclear particle detector. The classes correspond to neutrons and gamma rays. Neutrons have a slightly larger tail energy for the same total energy relative to gamma rays, which allows the two particle types to be distinguished. This is a somewhat large data set ( $n = 20,000$ ), and subgradient methods are appropriate given their scalability. Implement the subgradient method for minimizing  $J$  and apply it to the nuclear data. Submit two figures: One showing the data and the learned line, the other showing  $J$  as a function of iteration number. Also report the estimated hyperplane parameters and the minimum achieved value of the objective function.

*Comments:*

- Use  $\lambda = 0.001$ . Since this is a linear problem in a low dimension, we don't need much regularization.
  - Use a step-size of  $\alpha_j = 100/j$ , where  $j$  is the iteration number.
  - To compute the subgradient of  $J$ , write a subroutine to find the subgradient of  $J_i$  and then sum those results.
  - Since the objective will not be monotone decreasing, determining a good stopping rule can be tricky. For this problem, just look at the graph of the objective function and "eyeball it" to decide when the algorithm has converged.
  - Debugging goes faster if you start out with just a subsample of the data.
- (c) (6 pts) Now implement the stochastic subgradient method with a minibatch size of  $m = 1$ . Be sure to cycle through all data points before starting a new loop through the data. Report/hand in the same items as for part (b).

*More comments:*

- Use the same  $\lambda$ , stopping strategy, and  $\alpha_j$  as in part (b). Here  $j$  indexes the number of times you have cycled (randomly) through the data; i.e., it is the number of epochs you have trained for.

- Your plot of  $J$  versus iteration number will have roughly  $n$  times as many points as in part (b) since you have  $n$  updates for every one update of the full subgradient method.
- (d) (3 pts) Comment on the empirical rate of convergence of the stochastic subgradient method relative to the subgradient method. Explain your finding.
- (e) (5 pts) Submit your code.

## 5. Classification (15 pts).

- (a) (2 pts) Download a dataset for classification from the UCI Machine Learning Repository (you can Google it). Report the number of classes and the number of features in your chosen dataset.
- (b) (3 pts) Apply logistic regression to this dataset. Calculate the test error based on  $k$ -fold cross-validation (you may select  $k$ ) and report the training and test error. You may use any built-in methods in your language of choice.
- (c) (8 pts) Apply the SVM to this dataset using 2 different kernels: linear and Gaussian. Describe the tuning parameters in each case and set these parameters by cross-validation. Report your final test error, training error, and tuning parameters after cross-validation. You may use any built-in methods in your language of choice.
- (d) (2 pts) Which of the three methods you applied (logistic regression, SVM with linear kernel, SVM with Gaussian kernel) worked best? Do your results make sense?

*Comments:*

- The linear kernel for SVM is equivalent to applying the optimal soft-margin hyperplane.
- You will probably want to choose a dataset that does not have any categorical features.
- Depending on which code you use, you may want to choose data with a binary classification problem.
- The SVM can take a while to train so you may want to consider sample size when choosing a dataset.

## 6. Kernel Ridge Regression (20 pts). This problem examines the similarities and differences between KRR with and without offset.

- (a) (3 pts) In KRR with offset, give a formula for the offset  $b$  using the kernel.
- (b) (2 pts) This problem shows you how to compute the KRR with offset solution without a bunch of loops. The kernel matrix between all pairs of the training data is the  $n \times n$  matrix  $K = [k(\mathbf{x}_i, \mathbf{x}_j)]$ . The “centered” training kernel matrix  $\tilde{K}$  is the  $n \times n$  matrix whose entries are  $\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}_j) \rangle$  where  $\tilde{\Phi}(\mathbf{x}) := \Phi(\mathbf{x}) - \frac{1}{n} \sum_{\ell=1}^n \Phi(\mathbf{x}_\ell)$  and  $\Phi$  is a feature map corresponding to the kernel  $k$ . This is the matrix that we get in KRR with offset. We can calculate  $\tilde{K}$  using the following formula:

$$\tilde{K} = K - KO - OK + OKO, \quad (1)$$

where  $O$  is a square matrix with all entries equal to  $1/n$ . Prove that this relationship is true.

- (c) (4 pts) Now consider a test data set  $\mathbf{x}'_1, \dots, \mathbf{x}'_m$  and let  $K'$  be the  $n \times m$  kernel matrix between pairs of training and test points with entries  $k(\mathbf{x}_i, \mathbf{x}'_j)$ , and let  $\tilde{K}'$  be the  $n \times m$  centered version of  $K'$  whose entries are  $\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}'_j) \rangle$ . Determine a formula analogous to equation (1) for relating  $\tilde{K}'$  to  $K'$ . Also determine a formula for computing the predicted outputs on all the test points. Your formula should yield a column vector of length  $m$ .
- (d) Let's revisit the body fat data we saw earlier. We saw that a linear fit was reasonable, but now let's try a nonlinear fit. Use the first 150 examples for training, and the remainder for testing.
  - i. (3 pts) Implement kernel ridge regression *with* offset. Report the mean squared error (MSE) of the training data AND the test data. Report the offset  $b$  using your formula you derived earlier.

- ii. (4 pts) Implement kernel ridge regression *without* offset. Report the MSE of the training data AND the test data. Comment on any differences in performance with respect to KRR with offset (the previous part).
- iii. (4 pts) Please submit your code.

Additional comments:

- Use the Gaussian kernel with parameters  $\sigma = 15$  and  $\lambda = 0.003$ .
- The `dist2` function is useful for computing squared distances between a bunch of points.
- Parts (b) and (c) will save you from having to write a bunch of loops. Commands like `ones` and `repmat` in Matlab, or `np.ones` and `np.tile` in Python may also be useful here.