CS/ECE/ME532 Period 17 Activity

Estimated Time: 25 min for P1, 15 min for P2, 25 min for P4

- 1. Alternative regularization formulas. This problem is about two alternative ways of solving the L_2 -regularized least squares problem.
 - a) Prove that for any $\lambda > 0$, the following matrix identity holds:

$$(\boldsymbol{A}^T\boldsymbol{A} + \lambda \boldsymbol{I})^{-1}\boldsymbol{A}^T = \boldsymbol{A}^T(\boldsymbol{A}\boldsymbol{A}^T + \lambda \boldsymbol{I})^{-1}$$

Hint: Start by considering the expression $\mathbf{A}^T \mathbf{A} \mathbf{A}^T + \lambda \mathbf{A}^T$ and factor it in two different ways (from the right or from the left).

- b) The identity proved in part a) shows that there are actually two equivalent formulas for the solution to the L_2 -regularized least squares problem. Suppose $\mathbf{A} \in \mathbb{R}^{8000 \times 100}$ and $\mathbf{y} \in \mathbb{R}^{8000}$, and use this identity to find \mathbf{w} that minimizes $||\mathbf{A}\mathbf{w} \mathbf{y}||_2^2 + \lambda ||\mathbf{w}||_2^2$ in two different ways. Which formula will compute more rapidly? Why? *Note:* The number of operations required for matrix inversion is proportional to the cube of the matrix dimension.
- c) A breast cancer gene database has approximately 8000 genes from 100 subjects. The label y_i is the disease state of the ith subject (+1 if no cancer, -1 if breast cancer). Suppose we build a linear classifier that combines the 8000 genes, say $\mathbf{g}_i, i = 1, 2, \ldots, 100$ to predict whether a subject has cancer $\hat{y}_i = \text{sign}\{\mathbf{g}_i^T \mathbf{w}\}$. Note that here \mathbf{g}_i and \mathbf{w} are 8000-by-1 vectors.
 - i. Write down the least squares problem for finding classifier weights \boldsymbol{w} given 100 labels. Does this problem have a unique solution?
 - ii. Write down a Tikhonov(ridge)-regression problem for finding the classifier weights given 100 labels. Does this problem have a unique solution? Which form of the identity in part a) leads to the most computationally efficient solution for the classifier weights?
- 2. The key idea behind proximal gradient descent is to reformulate the general regularized least-squares problem into a set of simpler scalar optimization problems. Consider the regularized least-squares problem

$$\min_{\boldsymbol{w}} ||\boldsymbol{z} - \boldsymbol{w}||_2^2 + \lambda r(\boldsymbol{w})$$

An upper bound and completing the square was used to simplify the generalized least-squares problem into this form. Let the i^{th} elements of z and w be z_i and w_i , respectively.

- a) Assume $r(\mathbf{w}) = ||\mathbf{w}||_2^2$. Write the regularized least-squares problem as a series of separable problems involving only w_i and z_i .
- b) Assume $r(\mathbf{w}) = ||\mathbf{w}||_1$. Write the regularized least-squares problem as a series of separable problems involving only w_i and z_i .
- 3. A script is available to compute a specified number of iterations of the proximal gradient descent algorithm for solving a Tikhonov-regularized least squares problem

$$\min_{m{w}} ||m{y} - m{X}m{w}||_2^2 + \lambda ||m{w}||_2^2$$

The provided script will get you started displaying the path taken by the weights in the proximal gradient descent iteration superimposed on a contour plot of the squared

error surface. Assume $\boldsymbol{y} = \begin{bmatrix} \sqrt{2} \\ 0 \\ 1 \\ 0 \end{bmatrix}$, the 4-by-2 $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$ has singular value decomposition $\boldsymbol{U} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$, $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix}$, and $\boldsymbol{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$. Complete

decomposition
$$\boldsymbol{U} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$
, $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix}$, and $\boldsymbol{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$. Complete

20 iterations of gradient descent in each case specified below.

Include the plots you generate below with your submission.

- a) What is the maximum value for the step size τ that will guarantee convergence?
- b) Start proximal gradient descent from the point $\boldsymbol{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$ using a step size of $\tau = 0.5$ and tuning parameter $\lambda = 0.5$. How do you explain the trajectory the weights take toward the optimum, e.g., why is it shaped this way? What direction does each iteration move in the regularization step?
- c) Repeat the previous case with $\lambda = 0.1$ What happens? How does λ affect each iteration and why?