

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:

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

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 i th subject (+1 if no cancer, -1 if breast cancer). Suppose we build a linear classifier that combines the 8000 genes, say $\mathbf{g}_i, i = 1, 2, \dots, 100$ 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.
- Write down the least squares problem for finding classifier weights \mathbf{w} given 100 labels. Does this problem have a unique solution?
 - 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?

SOLUTION:

- a) Factoring the expression from the hint in two different ways, we have:

$$\mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}) = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}) \mathbf{A}^T$$

The terms in parentheses are always invertible, so we can multiply on the right by $(\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1}$ and on the left by $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1}$ to we obtain the desired identity.

- b) The two equivalent formulas for the solution to the L_2 -regularized least squares problem are:

$$\hat{\mathbf{w}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$$

The bottleneck in computation is the matrix inversion. Since $\mathbf{A} \in \mathbb{R}^{8000 \times 100}$, this means $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}) \in \mathbb{R}^{100 \times 100}$, whereas $(\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}) \in \mathbb{R}^{8000 \times 8000}$. So the second formula is on the order of 80^3 slower to compute.

c) Let $\mathbf{A} = \begin{bmatrix} \mathbf{g}_1^T \\ \mathbf{g}_2^T \\ \dots \\ \mathbf{g}_{100}^T \end{bmatrix}$ and $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{100} \end{bmatrix}$.

- i. $\min_{\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|_2^2$. Since $\mathbf{A} \in \mathbb{R}^{100 \times 8000}$, \mathbf{A} has rank at most 100 which is much less than the number of weights \mathbf{w} . No unique solution.
- ii. $\min_{\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$ always has a unique solution given by $\mathbf{w} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$. However, this requires inverting an 8000 by 8000 matrix, so it is much more efficient to use the above identity and obtain the solution as $\mathbf{w} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$ as this expression only requires inverting a 100 by 100 matrix.

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

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

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 \mathbf{z} and \mathbf{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 .

SOLUTION:

- a) $\|\mathbf{w}\|_2^2 = \sum_{i=1}^n w_i^2$, so the regularized least-squares problem is

$$\min_{\mathbf{w}} \sum_{i=1}^n (z_i - w_i)^2 + \lambda w_i^2$$

which is equivalent to

$$\min_{w_i} (z_i - w_i)^2 + \lambda w_i^2, i = 1, 2, \dots, n$$

b) In this case $\|\mathbf{w}\|_1 = \sum_{i=1}^n |w_i|$ so the regularized least-squares problem is

$$\min_{\mathbf{w}} \sum_{i=1}^n (z_i - w_i)^2 + \lambda |w_i|$$

which is equivalent to

$$\min_{w_i} (z_i - w_i)^2 + \lambda |w_i|, i = 1, 2, \dots, n$$

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

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

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 $\mathbf{y} = \begin{bmatrix} \sqrt{2} \\ 0 \\ 1 \\ 0 \end{bmatrix}$, the 4-by-2 $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ has singular value

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

20 iterations of gradient descent in each case specified below.

Include the plots you generate below with your submission.

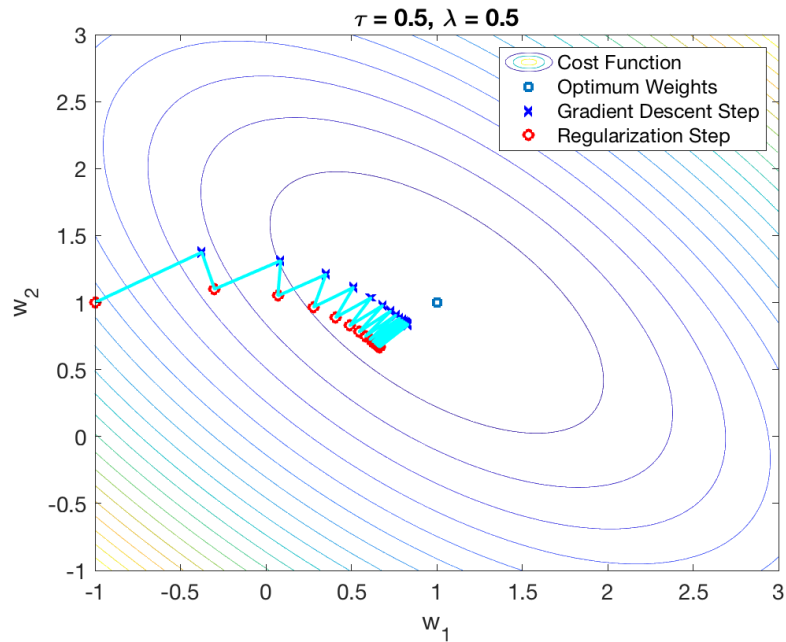
a) What is the maximum value for the step size τ that will guarantee convergence?

b) Start proximal gradient descent from the point $\mathbf{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$ 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?

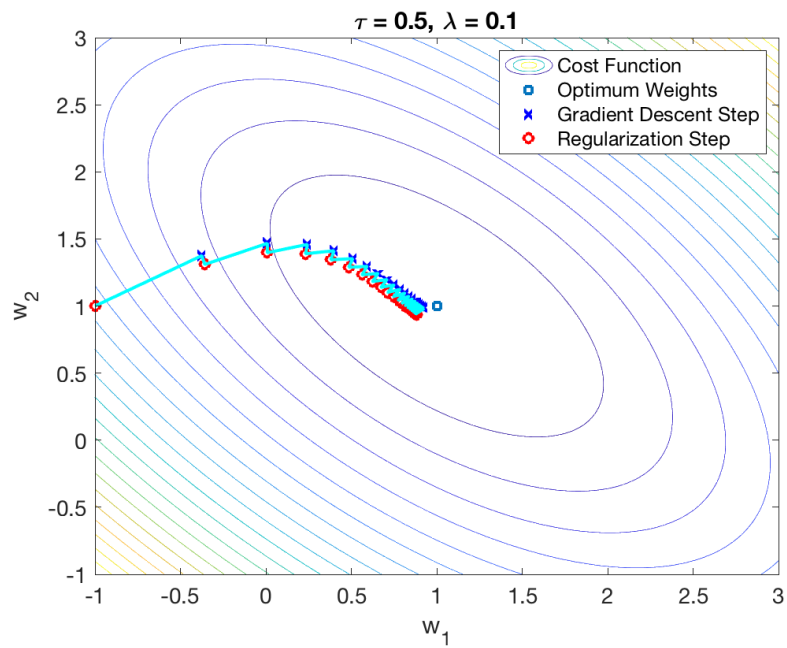
SOLUTION:

- a) The proximal gradient descent algorithm is derived assuming $\tau < 1/\|\mathbf{X}\|_{op}^2$ which implies $\tau < 1$ in this problem.



b)

Tikhonov regularization penalizes the length of the solution, so application of the regularizer forces the most recent iteration, marked by the blue x, toward the origin. Then the gradient at that point is used to determine the direction of the next update, that update is shrunk towards the origin, etc. Note that the solution does not converge to the bottom of the least-squares cost function because the regularizer pulls it towards the origin (smaller norm).



c)

Decreasing λ places less emphasis on the regularizer and results in less shrinkage toward the origin at each step. Conversely, if we increased λ then we'd be weighting the norm of the solution more heavily and there would be greater shrinkage of each iteration towards the origin.