

CS/ECE/ME532 Activity 17

Estimated Time: 25 min for P1, 15 min for P2, 25 min for P3

- 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?

- 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 .

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?

1a.) $(\underline{A}^T \underline{A} + \lambda \underline{I})^{-1} \underline{A}^T = \underline{A}^T (\underline{A} \underline{A}^T + \lambda \underline{I})^{-1}$ if $\lambda > 0$.

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

$$\rightarrow (\underline{A}^T \underline{A} + \lambda \underline{I})^{-1} \underline{A}^T (\underline{A} \underline{A}^T + \lambda \underline{I}) = \underline{A}^T \quad \underline{\text{LHS}}$$

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

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

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

$$\rightarrow \underline{I} \underline{A}^T = \underline{A}^T \quad \checkmark$$

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

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

$$\rightarrow (\underline{A}^T \underline{A} \underline{A}^T + \underline{A}^T \lambda) (\underline{A} \underline{A}^T + \lambda \underline{I})^{-1}$$

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

$$\rightarrow \underline{A}^T \underline{I} = \underline{A}^T \quad \checkmark$$

1b.) $A, 8000 \times 100$, $y: 8000 \times 1$

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

$\underline{A}^T \underline{A}$ is 100×100 $\underline{A} \underline{A}^T$ is 8000×8000

This expression will invert much faster.

1c.) i.) $\underline{w}_{\min} = (\underline{G}^T \underline{G})^{-1} \underline{G}^T \underline{y}$

$$\min_{\underline{w}} \|\underline{G}\underline{w} - \underline{y}\|_2^2$$
$$\begin{matrix} & 100 \\ 8000 & \left[\underline{G} \right] \left[\underline{w} \right] = \left[\underline{y} \right] \\ 8000 \times 100 & 100 \times 1 & 8000 \times 1 \end{matrix}$$

If $\text{rank } \underline{G} = \text{rank}(\underline{G}; \underline{y})$

& $\text{rank } \underline{G} = \dim(\underline{w}) = 100$, then there is a unique solution.

ii) $\min_{\underline{w}} \|\underline{G}\underline{w} - \underline{y}\|_2^2 + \lambda \|\underline{w}\|_2^2$

$$\rightarrow \underline{w}_{\min} = \underline{G}^T (\underline{G} \underline{G}^T + \lambda \underline{I})^{-1} \underline{y} = (\underline{G}^T \underline{G} + \lambda \underline{I})^{-1} \underline{G}^T \underline{y}$$

$\underline{G}^T \underline{G}$ is 100×100

will invert much faster.

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

$$f(\underline{w}) = \|\underline{z} - \underline{w}\|_2^2 + \lambda r(\underline{w})$$

$$a.) \text{ If } r(\underline{w}) = \|\underline{w}\|_2^2$$

$$\rightarrow r(\underline{w}) = \sum_{i=1}^m w_i^2$$

$$\rightarrow \min_{\underline{w}} \|\underline{z} - \underline{w}\|_2^2 + \lambda \sum_{i=1}^m w_i^2 \rightarrow \min_{w_i} \sum_i \sqrt{z_i^2 + w_i^2} + \lambda w_i^2$$

$$b.) \min_{w_i} \sum_i \sqrt{z_i^2 + w_i^2} + \lambda |w_i|$$

$$3a.) \tau \leq 1/\|\underline{A}\|_{op} = 1/\delta_1 = 1.$$

3b.)

Each regularization step

pulls the trajectory at the direction perpendicular to the surface.

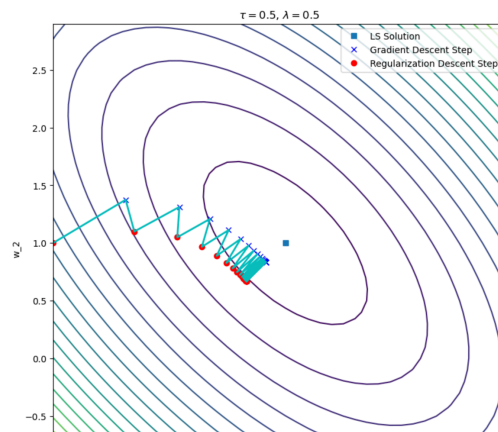
```

w_init = np.array([[-1],[1]])
lam = 0.5;
it = 20
tau = 0.5
W,Z = prxgraddescent_l2(X,y,tau,lam,w_init,it)

# Concatenate gradient and regularization steps to display trajectory
G = np.zeros((2,0))
for i in range(it):
    G = np.hstack((G,np.hstack([W[:,i],Z[:,i+1]])))

plt.figure(figsize=(9,9))
plt.contour(W1,w2,20)
plt.plot(w_ls[0],w_ls[1],"s", label="LS Solution")
plt.plot(Z[0,i],Z[1,i],"bx",linewidth=2, label="Gradient Descent Step")
plt.plot(W[0,i],W[1,i],"ro",linewidth=2, label="Regularization Descent Step")
plt.plot(G[0,i],G[1,i],"-c",linewidth=2)
plt.legend()
plt.xlabel('w_1')
plt.ylabel('w_2')
plt.title('$\\tau = 0.5, \lambda = 0.5$')

```



3c.)

As $\lambda \downarrow$,
the trajectory is pulled
"less hard" in the \perp
direction by the regularizing
term.

This allows us to get
closer to LS solution
ultimately.

