

CS/ECE/ME532 Activity 18

Estimated time: 15 mins for P1, 20 mins for P2, 15 mins for P3, 20 mins for P4

1. 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. You recall from the previous period that the least-squares problem for finding classifier weights has no unique solution.
Your hypothesis is that a relatively small number of the 8000 genes are predictive of the cancer state. Identify a regularization strategy consistent with this hypothesis and justify your choice.
2. Consider the least-squares problem $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ where $\mathbf{y} = 4$ and $\mathbf{X} = \begin{bmatrix} 2 & 1 \end{bmatrix}$.
 - a) Does this problem have a unique solution? Why or why not?
 - b) Sketch the contours of the cost function $f(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ in the $w_1 - w_2$ plane.
 - c) Now consider the LASSO $\min_{\mathbf{w}} \|\mathbf{w}\|_1$ subject to $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 < 1$. Find the solution using the following steps
 - i. Repeat your sketch from part b).
 - ii. Add a sketch of $\|\mathbf{w}\|_1 = c$
 - iii. Find the \mathbf{w} that satisfies $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 = 1$ with the minimum possible value of $\|\mathbf{w}\|_1$.
 - d) Use your insight from the previous part to sketch the set of solutions to the problem $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$ for $0 < \lambda < \infty$.
3. The script provided has a function that will compute a specified number of iterations of the proximal gradient descent algorithm for solving the ℓ_1 -regularized least-squares problem

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

The 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 for the cost function defined in problem 2. part b) starting from $\mathbf{w}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$. The script assumes $\lambda = 4$ and $\tau = 1/4$.

Include the plots you generate below with your submission.

a) How many iterations does it take for the algorithm to converge to the solution? What is the converged value for \mathbf{w} ?

b) Change to $\lambda = 2$. How many iterations does it take for the algorithm to converge to the solution? What is the converged value for \mathbf{w} ?

c) Explain what happens to the weights in the regularization step.

4. Use the proximal gradient algorithm to solve $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + 4\|\mathbf{w}\|_1$ for the parameters defined in problem 2.

a) What is the maximum value for the step size in the negative gradient direction, τ ?

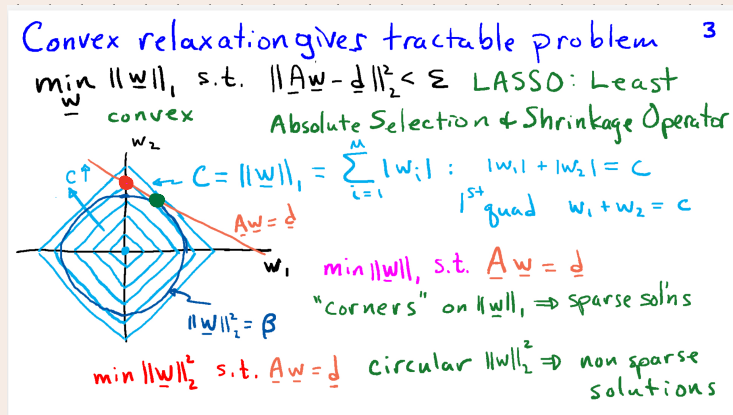
b) Suppose $\tau = 0.1$ and you start at $\mathbf{w}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$. Calculate the first two complete iterations of the proximal gradient algorithm and depict $\mathbf{w}^{(0)}, \mathbf{z}^{(1)}, \mathbf{w}^{(1)}, \mathbf{z}^{(2)}$ and $\mathbf{w}^{(2)}$ on a sketch of the cost function identical to the one you created in problem 2.b).

532 Activity 18 DEVIN BRESSER

1. 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. You recall from the previous period that the least-squares problem for finding classifier weights has no unique solution.

Your hypothesis is that a relatively small number of the 8000 genes are predictive of the cancer state. Identify a regularization strategy consistent with this hypothesis and justify your choice.

A good regularization strategy to eliminate features is l1-regularization aka LASSO regularization. This is because of the geometry of the unit norm “ball” for the l1-norm, which is more of a 45-degree rotated square. The lowest l1-norms lie at the “corners” of the unit diamond (along axes), and thus, l1-norm minimization forces some dimensions (features) to zero.



- smallest ℓ_1 norm is along corners, where some dimensions are 0

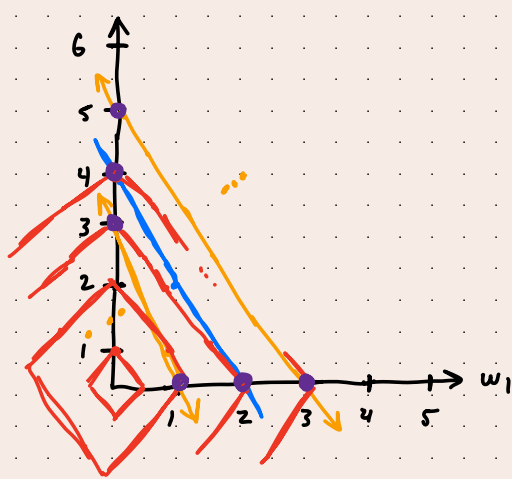
2. Consider the least-squares problem $\min_w \|y - Xw\|_2^2$ where $y = 4$ and $X = \begin{bmatrix} 2 & 1 \end{bmatrix}$.

- Does this problem have a unique solution? Why or why not?
- Sketch the contours of the cost function $f(w) = \|y - Xw\|_2^2$ in the $w_1 - w_2$ plane.
- Now consider the LASSO $\min_w \|w\|_1$ subject to $\|y - Xw\|_2^2 < 1$. Find the solution using the following steps:
 - Repeat your sketch from part b).
 - Add a sketch of $\|w\|_1 = c$.
 - Find the w that satisfies $\|y - Xw\|_2^2 = 1$ with the minimum possible value of $\|w\|_1$.
- Use your insight from the previous part to sketch the set of solutions to the problem $\min_w \|y - Xw\|_2^2 + \lambda \|w\|_1$ for $0 < \lambda < \infty$.

a.) $\text{rank}(X) = 1 < \dim(w) = 2$. No unique solution.

$$\begin{matrix} X & w & = & y \\ \begin{matrix} 1 \times 2 & 2 \times 1 & 1 \times 1 \end{matrix} & & & \end{matrix}$$

$$b.) f(w) = \|y - Xw\|_2^2 = \|4 - [2 \ 1]w\|_2^2 = (4 - 2w_1 - w_2)^2$$



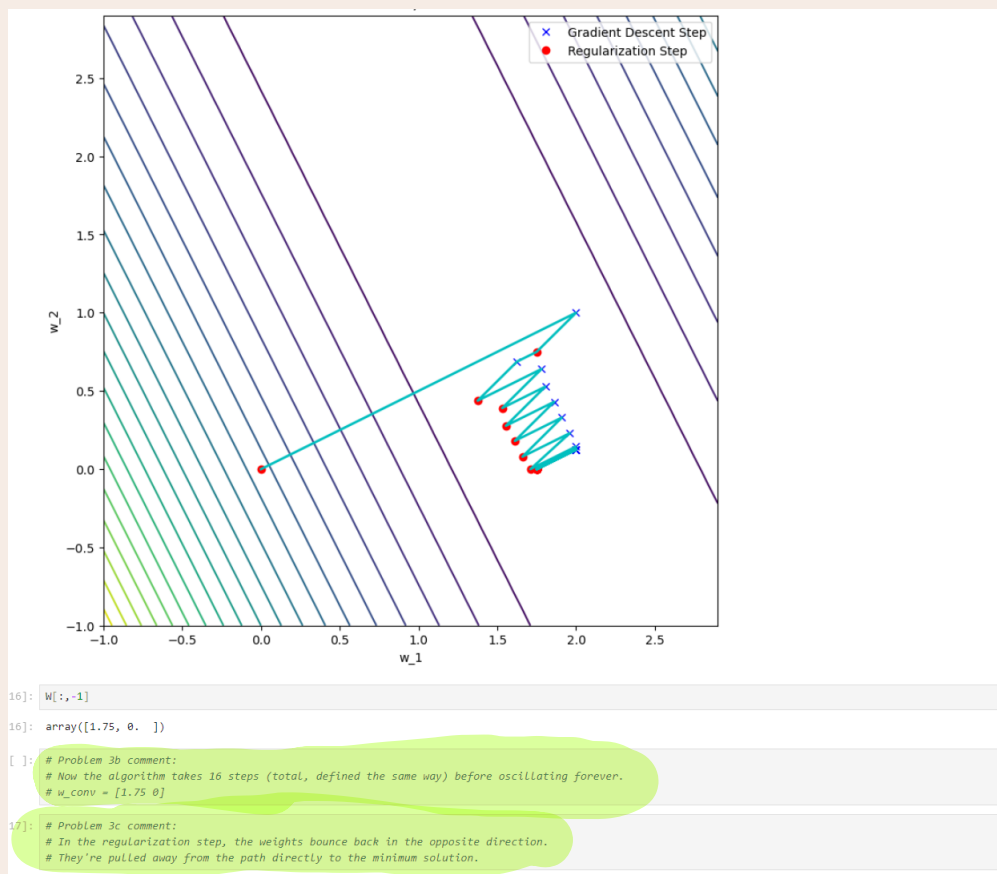
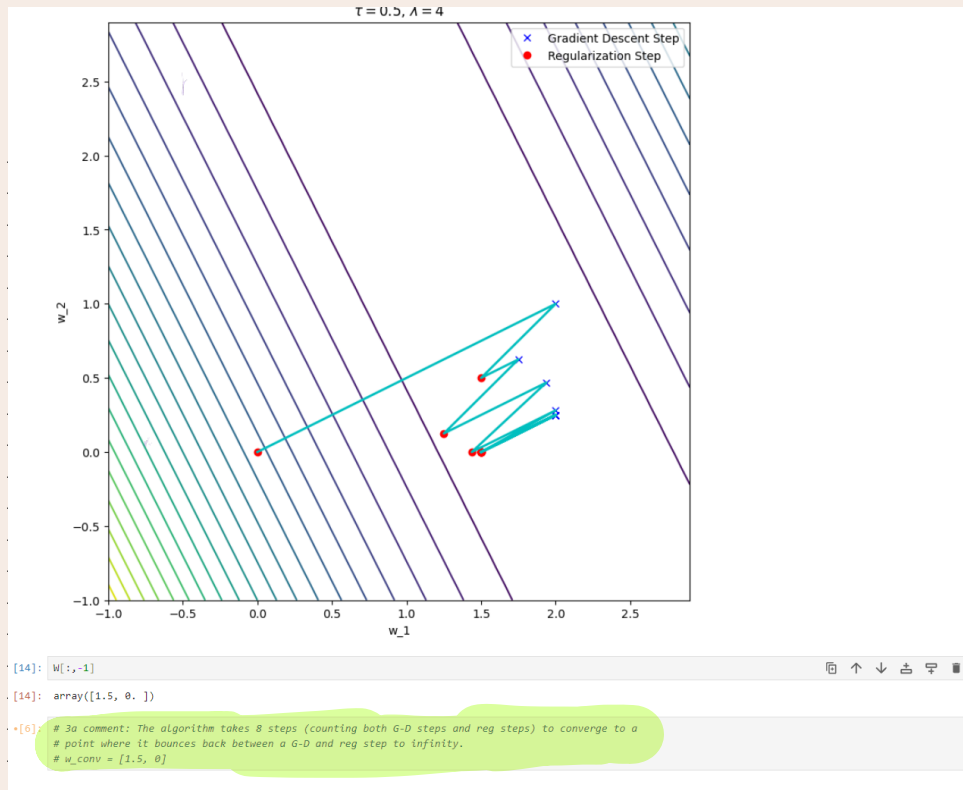
$$(4 - 2w_1 - w_2)^2 = c \rightarrow c = 0, 1, 2, 3, \dots$$

c.) Added L_1 norm rectangles: (red)

As we can see the minimum L_1 - norm is achieved along the w_1 & w_2 axes

d.) These are the purple dots • They will be spaced apart further & further with a squared relationship.

P3



$$X = \begin{bmatrix} 2 & 1 \end{bmatrix}, \quad y = 4$$

4.)

a.) calculator gives $S_1 = 2.236$

$$\rightarrow \tau < \frac{1}{\|A\|_{op}^2} = \frac{1}{(2.236)^2} = 0.2$$

b.) $\tau = 0.1, \quad \underline{w}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\underline{z}^{(0)} = \underline{w}^{(0)} - 2\tau \underline{X}^T (\underline{X} \underline{w}^{(0)} - y) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} (-4) = \begin{bmatrix} 1.6 \\ 0.8 \end{bmatrix}$$

$$\underline{w}^{(1)} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{z}^{(0)} - \underline{w}\|_2^2 + 0.4 \|\underline{w}\|_1$$

$$\rightarrow \begin{cases} w_1^{(1)} = (|z_1^{(0)}| - \frac{0.4}{2})_+ \operatorname{sign}(z_1^{(0)}) = (1.6 - 0.8)_+ \cdot 1 = 0.8 \\ w_2^{(1)} = (|z_2^{(0)}| - \frac{0.4}{2})_+ \operatorname{sign}(z_2^{(0)}) = (0.8 - 0.8)_+ \cdot 1 = 0 \end{cases} \rightarrow \underline{w}_1 = \begin{bmatrix} 0.8 \\ 0 \end{bmatrix}$$

$$\rightarrow \underline{z}^{(1)} = \underline{w}^{(1)} - 2\tau \underline{X}^T (\underline{X} \underline{w}^{(1)} - y) = \begin{bmatrix} 0.8 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} ([2 \ 1] \begin{bmatrix} 0.8 \\ 0 \end{bmatrix} - 4)$$

$$= \begin{bmatrix} 0.8 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} (-2.4) = \begin{bmatrix} 1.76 \\ 0.48 \end{bmatrix}$$

$$\rightarrow \begin{cases} w_1^{(2)} = (|z_1^{(1)}| - \frac{0.4}{2})_+ \operatorname{sign}(z_1^{(1)}) = (1.76 - 0.8)_+ \cdot 1 = 0.96 \\ w_2^{(2)} = (|z_2^{(1)}| - \frac{0.4}{2})_+ \operatorname{sign}(z_2^{(1)}) = (0.48 - 0.8)_+ \cdot 1 = 0 \end{cases} \rightarrow \underline{w}^{(2)} = \begin{bmatrix} 0.96 \\ 0 \end{bmatrix}$$

$$\rightarrow \underline{z}^{(2)} = \underline{w}^{(2)} - 2\tau \underline{X}^T (\underline{X} \underline{w}^{(2)} - y) = \begin{bmatrix} 0.96 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} ([2 \ 1] \begin{bmatrix} 0.96 \\ 0 \end{bmatrix} - 4)$$

$$= \begin{bmatrix} 0.96 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} (-2.08) = \begin{bmatrix} 1.792 \\ 0.416 \end{bmatrix}$$

1) Initialize $\underline{w}_0 = 0, \quad 0 < \tau < 1/\|A\|_{op}^2$
 2) LS-GD $\underline{z}^{(k)} = \underline{w}^{(k)} - \tau A^T (A \underline{w}^{(k)} - \underline{d})$ (grad. descent from current estimate of \underline{w} to $\underline{z}^{(k)}$ defines $\underline{z}^{(k)}$)
 3) Regularize $\underline{w}^{(k+1)} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \lambda \tau r(\underline{w})$ (Solve the proximal (Lasso) regularized problem. Find a \underline{w} that has been regularized that it comes to $\underline{z}^{(k)}$)
 4) Check if converged if $\|\underline{w}^{(k+1)} - \underline{w}^{(k)}\|_2^2 < \epsilon$, Stop
 If iteration does not result in significant change to \underline{w} , else, to 2.)

