

CS/ECE/ME532 Activity 20

Estimated time: 15 min for P1, 20 min for P2, 15 min for P3

1. An exponential loss function $f(w)$ is defined as

$$f(w) = \begin{cases} e^{-2(w-1)}, & w < 1 \\ e^{w-1}, & w \geq 1 \end{cases}$$

- a) Is $f(w)$ convex? Why? *Hint:* Graph the function.
- b) Is $f(w)$ differentiable everywhere? If not, where not?
- c) The “differential set” $\partial f(\mathbf{w})$ is the set of subgradients $\mathbf{v} \in \partial f(\mathbf{w})$ for which $f(\mathbf{u}) \geq f(\mathbf{w}) + (\mathbf{u} - \mathbf{w})^T \mathbf{v}$. Find the differential set for $f(w)$ as a function of w .
2. We are trying to predict whether a certain chemical reaction will take place as a function of our experimental conditions: temperature, pressure, concentration of catalyst, and several other factors. For each experiment $i = 1, \dots, m$ we record the experimental conditions in the vector $\mathbf{x}_i \in \mathbb{R}^n$ and the outcome in the scalar $b_i \in \{-1, 1\}$ (+1 if the reaction occurred and -1 if it did not). We will train our linear classifier to minimize hinge loss. Namely, we solve:

$$\underset{\mathbf{w}}{\text{minimize}} \quad \sum_{i=1}^m (1 - b_i \mathbf{x}_i^T \mathbf{w})_+ \quad \text{where } (u)_+ = \max(0, u) \text{ is the hinge loss operator}$$

- a) Derive a gradient descent method for solving this problem. Explicitly give the computations required at each step. *Note:* you may ignore points where the function is non-differentiable.
- b) Explain what happens to the algorithm if you land at a \mathbf{w}^k that classifies all the points perfectly, and by a substantial margin.
3. You have four training samples $y_1 = 1, \mathbf{x}_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$, $y_2 = 2, \mathbf{x}_2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$, $y_3 = -1, \mathbf{x}_3 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$, and $y_4 = -2, \mathbf{x}_4 = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$. Use cyclic stochastic gradient descent to find the first two updates for the LASSO problem

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

assuming a step size of $\tau = 1$ and $\mathbf{w}^{(0)} = \mathbf{0}$. Also indicate the data used for the first six updates.