

Assignment 4: Non-linear optimization

Assigned: October 19; Due: November 01.

Convergence rate of gradient descent [40 pts]

Consider an M -smooth function $f : \mathbb{R}^n \mapsto \mathbb{R}$. We consider the following optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

We apply the gradient descent method with constant step size $\alpha \leq 1/M$.

- (a) Show that the following inequality holds for each iteration k :

$$\|\nabla f(\mathbf{x}^k)\|^2 \leq \frac{2}{\alpha}(f(\mathbf{x}^k) - f(\mathbf{x}^{k+1}))$$

- (b) Show that the following inequality holds for each iteration k :

$$\min_{\ell=0, \dots, k-1} \|\nabla f(\mathbf{x}^\ell)\| \leq \sqrt{\frac{2}{\alpha k}(f(\mathbf{x}^0) - z^*)}$$

- (c) Bound the number of iterations required to reach a solution $\bar{\mathbf{x}}$ such that $\|\nabla f(\bar{\mathbf{x}})\| \leq \varepsilon$. Comment on the rate of convergence, compared to the results seen in class.

Newton's method [20 pts]

Define the following cubic functions

$$\begin{aligned} f(x) &= x^3 - 12x^2 + 10x + 3 \\ g(x) &= x^3 + 2x \end{aligned}$$

- Implement Newton's method to minimize f . Using a starting point of 10, plot the solution over 30 iterations. Repeat with a starting point of 2.5. Comment briefly on the performance of the algorithm.
- Implement Newton's method to minimize g . Using a starting point of 5, plot the solution over 30 iterations. Comment briefly on the performance of the algorithm.

Gradient descent vs. stochastic gradient descent [40 pts]

You have access to two datasets, each with $n = 1,200$ observations. Each dataset comprises two covariates stored in $1,200 \times 2$ matrices \mathbf{X}^A and \mathbf{X}^B , and $1,200 \times 1$ label vectors \mathbf{y}^A and \mathbf{y}^B . The goal is to build a *single* linear regression model for both datasets, with the following loss function:

$$\min_{\beta_1, \beta_2} \frac{1}{2n} \sum_{i=1}^n \min \left\{ (y_i^A - \beta_1 X_{i1}^A - \beta_2 X_{i2}^A)^2, (y_i^B - \beta_1 X_{i1}^B - \beta_2 X_{i2}^B)^2 \right\}$$

You have access to the following data files:

- dataA.csv: A matrix of size $1,200 \times 3$ that stores the first covariate $X_{11}^A, \dots, X_{n1}^A$ (first column), the second covariate $X_{12}^A, \dots, X_{n2}^A$ (second column), and the label y_1^A, \dots, y_n^A (third column).
- dataB.csv: A matrix of size $1,200 \times 3$ that stores the first covariate $X_{11}^B, \dots, X_{n1}^B$ (first column), the second covariate $X_{12}^B, \dots, X_{n2}^B$ (second column), and the label y_1^B, \dots, y_n^B (third column).

We will implement gradient descent and stochastic gradient descent for this problem. Throughout, use a starting point of $\beta^0 = (45, -15)$, 200 iterations, and a constant learning rate.

- Provide a contour plot of the loss function, for $\beta_1 \in [-25, 75]$ and $\beta_2 \in [-40, 40]$. What is the shape of the function? What are the implications for unconstrained optimization algorithms?
- Implement the gradient descent method with a learning rate of 0.01. Add the algorithm's trajectory on the contour plot from Question a. Repeat with a learning rate of 0.1, and provide a new trajectory plot. Comment briefly on the performance of the algorithm.
- Implement the stochastic gradient descent method with a learning rate of 0.01. Run the algorithm three times, with three random seeds. Report a contour plot of the cost function showing the three trajectories of your algorithm. Comment briefly.