

3rd Assignment - Optimization - Gradient Methods

Exercise 3.1.: $f(x, y) = x^2 + xy + y^2 + 5$

▷ Starting at $(1, 1)$, write 2 steps of conjugate gradient method for f .

1st step:

$$x_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \nabla f(x) = \begin{pmatrix} 2x + y \\ x + 2y \end{pmatrix} \quad \nabla f(x_0) = \begin{pmatrix} 3 \\ 3 \end{pmatrix}$$

$$z_1 = -\nabla f(x_0) = \begin{pmatrix} -3 \\ -3 \end{pmatrix}$$

Minimizing $f(x_0 + \alpha_1 z_1)$ with respect to α_1 , we get

$$\alpha_1 = - \frac{(z_1)^T \nabla f(x_0)}{(z_1)^T A (z_1)} \quad \text{where } A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \text{ the Hessian.}$$

$$\alpha_1 = - \frac{(-3 \ -3) \begin{pmatrix} 3 \\ 3 \end{pmatrix}}{(-3 \ -3) \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -3 \end{pmatrix}} = - \frac{-18}{(-9 \ -9) \begin{pmatrix} -3 \\ -3 \end{pmatrix}} = - \frac{-18}{54} = \frac{1}{3}$$

$$x_1 = x_0 + \alpha_1 z_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \frac{1}{3} \begin{pmatrix} -3 \\ -3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} //$$

- For the same starting point and function, do 2 steps for hypergradient descent method.

1st step:

$$\begin{pmatrix} x_0 \\ y_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \nabla f(x_0, y_0) = \begin{pmatrix} 3 \\ 3 \end{pmatrix} \quad \nabla f(x, y) = \begin{pmatrix} 2x+y \\ x+2y \end{pmatrix}$$

$$d_1 = -\nabla f(x_0, y_0) = \begin{pmatrix} -3 \\ -3 \end{pmatrix} \quad A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

$$(x_1, y_1) = (x_0, y_0) - \alpha_0 A^{-1} \nabla f(x_0, y_0)$$

$$A^{-1} = \frac{1}{4-1} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} = \begin{pmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{pmatrix}$$

$$\text{for } \alpha_0 = 0.01 = \frac{1}{100}$$

$$\begin{aligned} \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} &= \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \frac{1}{100} \begin{pmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{pmatrix} \begin{pmatrix} 3 \\ 3 \end{pmatrix} \\ &= \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \frac{1}{100} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} \end{aligned}$$

2nd step:

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} \quad \nabla f(x_1, y_1) = \begin{pmatrix} 297/100 \\ 297/100 \end{pmatrix}$$

$$\begin{aligned} (x_2, y_2) &= (x_1, y_1) - \alpha_1 A^{-1} \nabla f(x_1, y_1) \\ &= \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} - \alpha_1 \begin{pmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{pmatrix} \begin{pmatrix} 297/100 \\ 297/100 \end{pmatrix} = \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} - \alpha_1 \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} \end{aligned}$$

$$\text{To calculate } \alpha_1: \frac{\partial}{\partial \alpha_1} \underbrace{f\left[(x_1, y_1) - \alpha_1 A^{-1} \nabla f(x_1, y_1)\right]}_{g(\alpha_1)}$$

$$\Rightarrow \frac{\partial}{\partial \alpha_1} g(\alpha_1) = 0$$

$$\text{So, } g(\alpha_1) = 3 \left[\begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} - \alpha_1 \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} \right]^2 + 5$$

$$\frac{\partial}{\partial \alpha_1} g(\alpha_1) = -6 \left[\begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} - \alpha_1 \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} \right] \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} = 0 \Rightarrow \boxed{\alpha_1 = 1}$$

Therefore,

$$(x_2, y_2) = \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} - 1 \begin{pmatrix} 99/100 \\ 99/100 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} //$$

Exercise 3.2: Write down 1 step in the classical Newton method for the function $f(x, y) = (x+1)^2 + (y+3)^2 + 4$ starting at $(0, 0)$

$$\nabla f(x, y) = \begin{pmatrix} 2(x+1) \\ 2(y+3) \end{pmatrix} \quad H(x, y) = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \text{ the hessian matrix}$$

$$x_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \nabla f(0, 0) = \begin{pmatrix} 2 \\ 6 \end{pmatrix} \quad g_1 = -\nabla f(0, 0) = \begin{pmatrix} -2 \\ -6 \end{pmatrix}$$

$$x_1 = \overset{\text{Newton method}}{x_0 + [H(x_0)]^{-1} g_1} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \left[\begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \right]^{-1} \begin{pmatrix} -2 \\ -6 \end{pmatrix}$$

$$\hookrightarrow \left[\begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \right]^{-1} = \frac{1}{4-0} \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix}$$

$$\text{In general: } A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad A^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

$$x_1 = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} -2 \\ -6 \end{pmatrix} = \begin{pmatrix} -1 \\ -3 \end{pmatrix}$$

Exercise 3.3: $x^{(k+1)} = x^{(k)} - \frac{\alpha^{(k)}}{|I_k|} \sum_{i \in I_k} \nabla f_i(x^{(k)})$

size of minibatch equal to 2. Show that:

▷ $\sum_{i \in I_k} \nabla f_i(x)$ is a stochastic gradient

This expression is the sum of the individual functions $f_i(x)$ for i in the mini-batch I_k . Therefore, $\sum_{i \in I_k} \nabla f_i(x)$ is an estimate of the true gradient and it is "stochastic" because it is calculated based on a randomly selected subset of the data.

$|I_k| = 2$, I_k a set with two elements randomly selected from $\{1, 2, \dots, n\}$

$$g_k = \frac{1}{|I_k|} \sum_{i \in I_k} \nabla f_i(x) = \frac{1}{2} (\nabla f_{i_1}(x) + \nabla f_{i_2}(x))$$

$$\mathbb{E}(g_k) = \frac{1}{2} \mathbb{E}(\nabla f_{i_1}(x) + \nabla f_{i_2}(x)) \underset{\substack{\text{because are independent}}}{=} \frac{1}{2} \mathbb{E}(\nabla f_{i_1}(x)) + \frac{1}{2} \mathbb{E}(\nabla f_{i_2}(x)) \underset{\substack{\text{because are independent}}}{=} \frac{1}{2} \cdot 2 \mathbb{E}(\nabla f_{i_1}(x)) = \nabla F(x)$$

▷ The variance of this gradient is smaller than the variance of the standard stochastic gradient $\nabla f_{i_k}(x)$ (size of minibatch!)

Let's write down the variances for each gradient

$$\bullet \text{Var}(g_k) = \text{Var}\left(\frac{1}{2} (\nabla f_{i_1}(x) + \nabla f_{i_2}(x))\right)$$

$$= \frac{1}{4} \text{Var}(\nabla f_{i_1}(x)) + \frac{1}{4} \text{Var}(\nabla f_{i_2}(x)) + \frac{1}{2} \text{Cov}(\nabla f_{i_1}(x), \nabla f_{i_2}(x))$$

↳ because $\nabla f_{i_1}(x)$ and $\nabla f_{i_2}(x)$ are independent

$$\bullet \text{Var}(\nabla f_{i_k}(x))$$

The variance of g_k is the average of the variances of two individual gradients in the mini-batch and it is smaller than the variance of $\nabla f_{i_k}(x)$ when the mini-batch size is 2.

▷ Does your argument work for other sizes of the mini-batches?

Yes, because the main idea is that while mini-batch size increases, the variance of the stochastic gradient tends to decrease compared to the variance of the individual gradients in the mini-batch. This is because averaging over some samples tends to smooth out the noise associated with individual noise.

For example a mini-batch size m : $\bullet \text{Var}(g_k) = \frac{1}{m} \sum_{i \in I_k} \text{Var}(\nabla f_i(x))$
 $\bullet \text{Var}(\nabla f_{i_k}(x))$