

## COURSEWORK 2

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

---

# Mathematics for Machine Learning

---

*Author:*

Arnold Cheung (CID: 01184493)

Date: October 24, 2019

# 1 Differentiation

(a) First, expand the completed square form:

$$\begin{aligned} & (\mathbf{x} - \mathbf{c})^T \mathbf{C} (\mathbf{x} - \mathbf{c}) + c_0 \\ &= \mathbf{x}^T \mathbf{C} \mathbf{x} - \mathbf{x}^T \mathbf{C} \mathbf{c} - \mathbf{c}^T \mathbf{C} \mathbf{x} + \mathbf{c}^T \mathbf{C} \mathbf{c} + c_0 \\ &= \mathbf{x}^T \mathbf{C} \mathbf{x} - 2\mathbf{c}^T \mathbf{C} \mathbf{x} + \mathbf{c}^T \mathbf{C} \mathbf{c} + c_0 \end{aligned}$$

Then compare coefficients with  $f_1(\mathbf{x})$ , for  $\mathbf{C}$ :

$$\begin{aligned} \mathbf{x}^T \mathbf{C} \mathbf{x} &= \mathbf{x}^T \mathbf{B} \mathbf{x} - \mathbf{x}^T \mathbf{x} \\ \mathbf{x}^T \mathbf{C} \mathbf{x} &= \mathbf{x}^T (\mathbf{B} - \mathbf{I}) \mathbf{x} \\ \Rightarrow \mathbf{C} &= \mathbf{B} - \mathbf{I} \\ &= \begin{bmatrix} 4 & -2 \\ -2 & 4 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 4 \end{bmatrix} \\ &= \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix} \end{aligned}$$

For  $\mathbf{c}$ :

$$\begin{aligned} -2\mathbf{c}^T \mathbf{C} \mathbf{x} &= \mathbf{a}^T \mathbf{x} - \mathbf{b}^T \mathbf{x} \\ -2\mathbf{c}^T \mathbf{C} \mathbf{x} &= (\mathbf{a}^T - \mathbf{b}^T) \mathbf{x} \\ -2\mathbf{c}^T \mathbf{C} &= (\mathbf{a}^T - \mathbf{b}^T) \\ \mathbf{c}^T &= \left( \frac{\mathbf{b}^T - \mathbf{a}^T}{2} \right) \mathbf{C}^{-1} \\ \mathbf{c} &= \left( \left( \frac{\mathbf{b}^T - \mathbf{a}^T}{2} \right) \mathbf{C}^{-1} \right)^T \\ \mathbf{c} &= \left( \frac{1}{5} \left( \frac{\begin{bmatrix} -2 & 0 \end{bmatrix} - \begin{bmatrix} 0 & 1 \end{bmatrix}}{2} \right) \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix} \right)^T \\ \mathbf{c} &= \left( \frac{1}{5} \begin{bmatrix} -1 & 0 \end{bmatrix} \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix} \right)^T \\ \mathbf{c} &= \frac{1}{5} \begin{bmatrix} -3 \\ -2 \end{bmatrix} \end{aligned}$$

For  $c_0$ :

$$\begin{aligned} \mathbf{c}^T \mathbf{C} \mathbf{c} + c_0 &= 0 \\ c_0 &= -\mathbf{c}^T \mathbf{C} \mathbf{c} \\ c_0 &= -\frac{1}{25} \begin{bmatrix} -3 & -2 \end{bmatrix} \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix} \begin{bmatrix} -3 \\ -2 \end{bmatrix} \\ c_0 &= -\frac{3}{5} \end{aligned}$$

Therefore,  $\mathbf{x}^T \mathbf{B} \mathbf{x} - \mathbf{x}^T \mathbf{x} + 2\mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{x}$  can be written as:

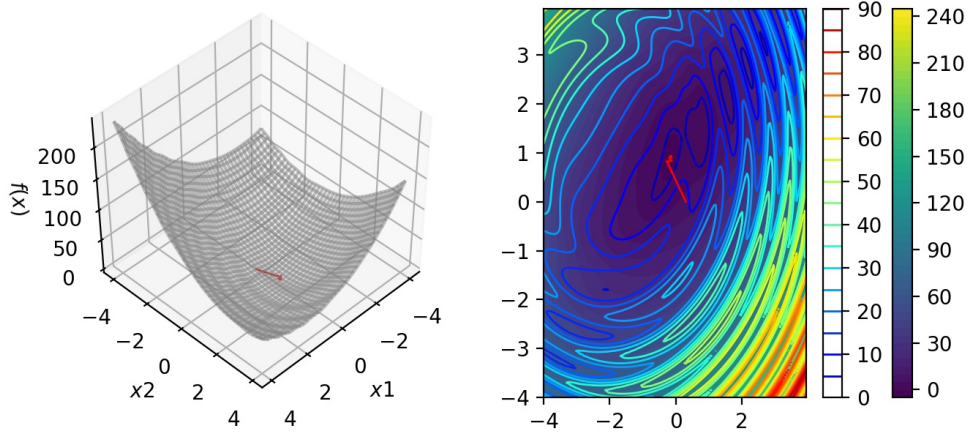
$$\left(\mathbf{x} - \frac{1}{5} \begin{bmatrix} -3 \\ -2 \end{bmatrix}\right)^T \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix} \left(\mathbf{x} - \frac{1}{5} \begin{bmatrix} -3 \\ -2 \end{bmatrix}\right) - \frac{3}{5}$$

- (b)  $f_1$  has minimum point can be deduced from  $\mathbf{C}$  being positive definite. The minimum point is located at  $\mathbf{x} = \mathbf{c} = \begin{bmatrix} -0.6 \\ -0.4 \end{bmatrix}$
- (c) The gradient of each function is given by the first derivative with respect to  $\mathbf{x}$ , the respective derivatives are:

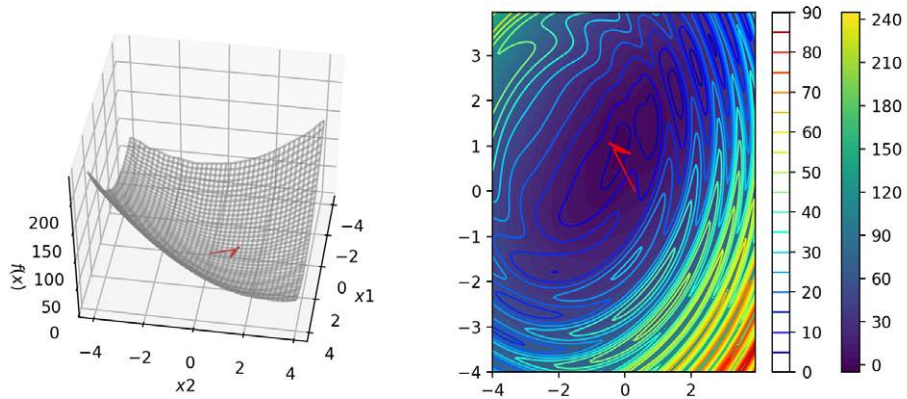
$$\begin{aligned} \frac{df_1}{d\mathbf{x}} &= 2\mathbf{x}^T \mathbf{B} - 2\mathbf{x}^T + \mathbf{a}^T + \mathbf{b}^T \\ \frac{df_2}{d\mathbf{x}} &= -2 \sin[(\mathbf{x} - \mathbf{b})^T (\mathbf{x} - \mathbf{b})] (\mathbf{x} - \mathbf{b})^T + 2(\mathbf{x} - \mathbf{a}) \mathbf{B} \\ \frac{df_3}{d\mathbf{x}} &= 2 \exp[(\mathbf{x} - \mathbf{a})^T (\mathbf{x} - \mathbf{a})] (\mathbf{x} - \mathbf{a})^T + 2 \exp[(\mathbf{x} - \mathbf{b})^T \mathbf{B} (\mathbf{x} - \mathbf{b})] (\mathbf{x} - \mathbf{b})^T \mathbf{B} \\ &\quad + \frac{1}{10} \text{tr}\left(\left(\frac{1}{100} \mathbf{I} + \mathbf{x} \mathbf{x}^T\right)^{-1} \frac{d\mathbf{x} \mathbf{x}^T}{d\mathbf{x}}\right) \end{aligned}$$

The code implementations can be found in the coding\_ansewrs.py file.

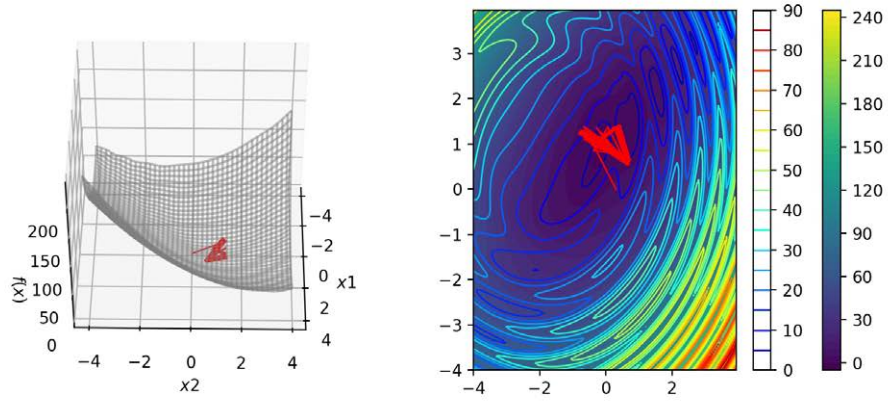
- (d) The continuous colour represents the value of the function at a given point, and the discrete contours represent the magnitude of the gradient of the function at a given point. All of the gradient descents below have a starting position of  $\begin{bmatrix} 0.3 \\ 0 \end{bmatrix}$  with 50 iterations, shown by the red line on both plots:



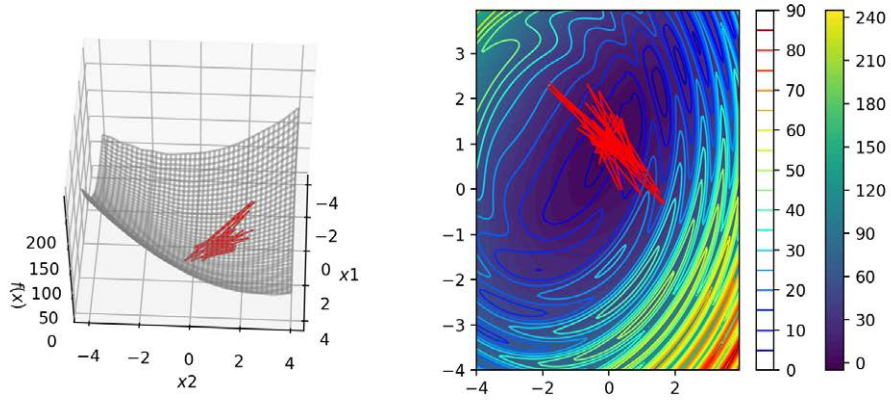
**Figure 1:** The 3D and contour plot of  $f_2$ , step size = 0.07.



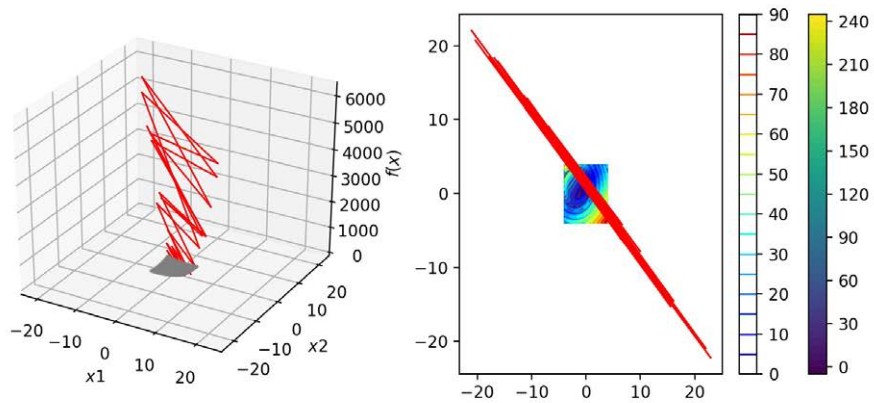
**Figure 2:** The 3D and contour plot of  $f_2$ , step size = 0.1.



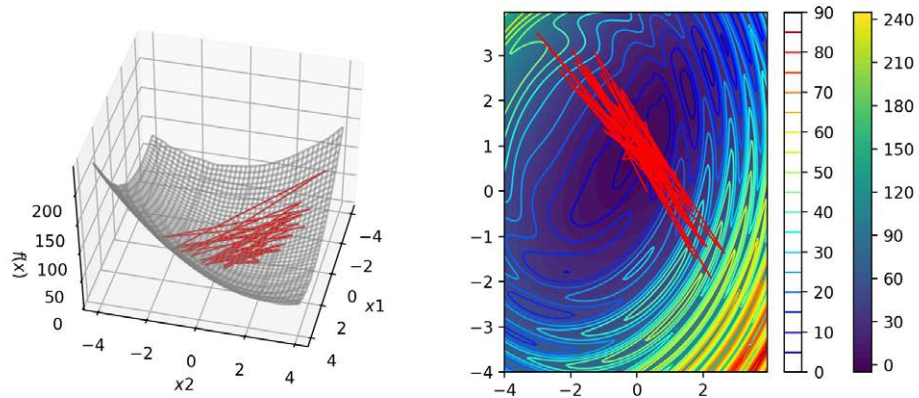
**Figure 3:** The 3D and contour plot of  $f_2$ , step size = 0.15.



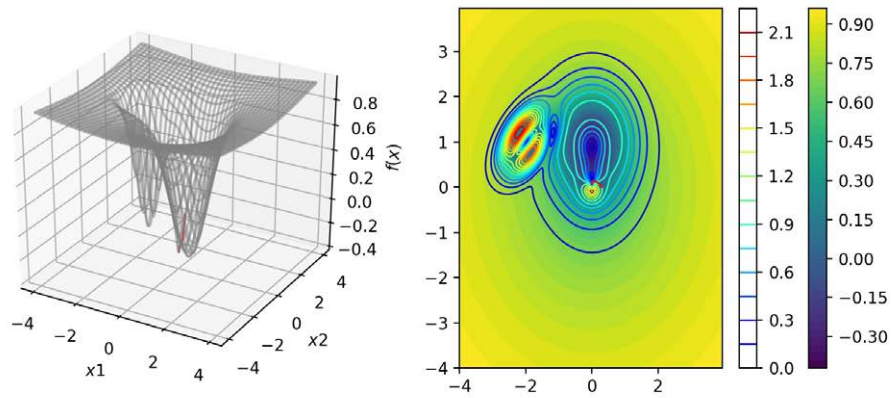
**Figure 4:** The 3D and contour plot of  $f_2$ , step size = 0.173.



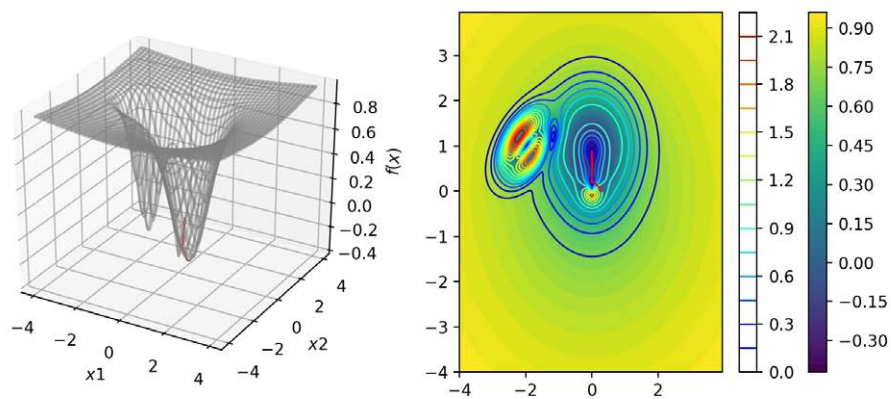
**Figure 5:** The 3D and contour plot of  $f_2$ , step size = 0.174.



**Figure 6:** The 3D and contour plot of  $f_2$ , step size = 0.1747.

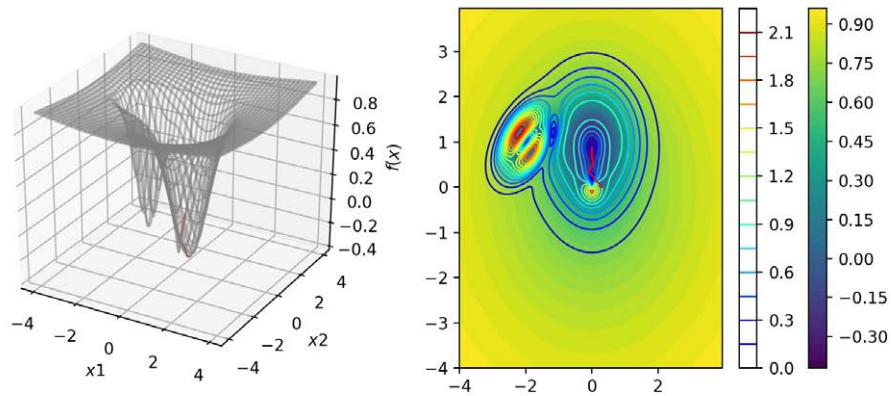


**Figure 7:** The 3D and contour plot of  $f_3$ , step size = 0.07.

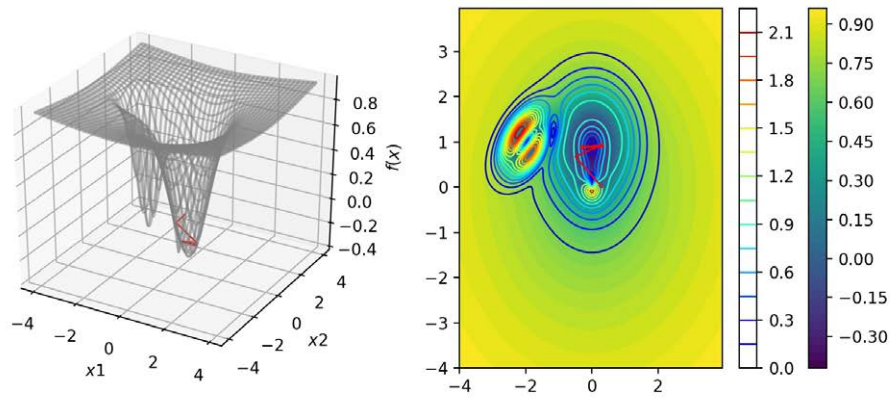


**Figure 8:** The 3D and contour plot of  $f_3$ , step size = 0.26.

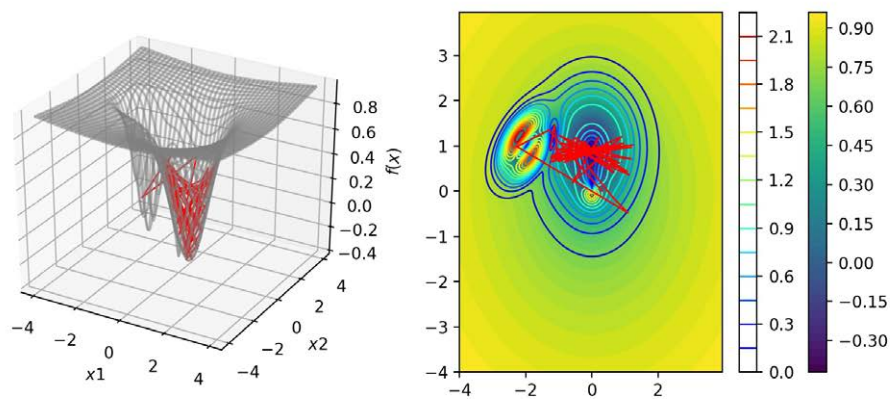




**Figure 9:** The 3D and contour plot of  $f_3$ , step size = 0.5.



**Figure 10:** The 3D and contour plot of  $f_3$ , step size = 1.0.



**Figure 11:** The 3D and contour plot of  $f_3$ , step size = 2.0.

- (e) For  $f_2$ , for step sizes small enough (smaller than 0.1), the gradient descent algorithm reliably locates the minimum around  $[-0.16, 0.93]$  of the function as shown in Figure 1. However, the algorithm begins to lose its convergence from step size around 0.1 and above, from Figure 2 to Figure 4, it can be seen that the output position of the algorithm fluctuates increasingly with an increasing step size (from 0.1 to 0.173). The result eventually diverges at step size = 0.174 as shown in Figure 5, where the algorithm is "climbing" out of the valley to a very high value, outside the range of the predefined axes. However, surprisingly when step size is increased to 0.1747 the algorithm is again contained within the predefined axes, as shown in Figure 6. This effect is likely due to the fact that when step size = 0.174, the algorithm lands on the right side of the periodic wave where the gradient points towards the minimum. And for step sizes beyond 1.0, the algorithm simply outputs values that increases exponentially.

For  $f_3$ , for step sizes small enough (smaller than 0.26), the gradient descent algorithm lands inside the local minimum of the small valley right next to the largest valley around  $[0,0]$  as shown in Figure 7. Step size = 0.26 is the first time it skips the small valley in to the lower minimum in the largest valley around  $[0, 0.89]$ . This remains true for increasing step sizes, although the lowest minimum cannot be reached for step sizes too big, such as 1.0 shown in Figure 10. For step sizes beyond 1.0, for example at 2.0 shown in Figure 11, although the algorithm is still contained within the big valley, it is increasingly unstable and even reaching positions in the second biggest valley occasionally.