

# Coursework 1

## Mathematics for Machine Learning (CO-496)

### Instructions

This coursework has both writing and coding components. The coding part must be done in python. The code will be run on the standard CSG installation and will be tested on the labTS <sup>1</sup> system. On starting the course you will be given a gitlab repository. Every commit you make to this repository can be marked automatically.

You are not permitted to use any symbolic manipulation libraries (e.g. `sympy`) or automatic differentiation tools (e.g. `tensorflow`, `pytorch`) for your submitted code (though, of course, you may find these useful for checking your answers). You should not need to import anything other than `numpy` for the submitted code for this assignment.

The writing assignment requires plots, which you can create using any method of your choice. You should not submit the code used to create these plots.

No aspect of your submission may be hand-drawn. You are strongly encouraged to use  $\text{\LaTeX}$  to create the written component.

You are required to submit the following:

- A `.pdf` file for your written answers, submitted through CATe
- A `.py` file which implements all the methods for the coding exercises, submitted through gitlab.

---

<sup>1</sup><https://teaching.doc.ic.ac.uk/labts>

# 1 Differentiation

In this question, we define the following constants:

$$\mathbf{B} = \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}, \quad \mathbf{a} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} 0 \\ -1 \end{pmatrix}.$$

We define also the following functions, which are all  $\mathbb{R}^2 \rightarrow \mathbb{R}$

$$f_1(\mathbf{x}) = \mathbf{x}^T \mathbf{x} + \mathbf{x}^T \mathbf{B} \mathbf{x} - \mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{x}$$

$$f_2(\mathbf{x}) = \sin((\mathbf{x} - \mathbf{a})^T (\mathbf{x} - \mathbf{a})) + (\mathbf{x} - \mathbf{b})^T \mathbf{B} (\mathbf{x} - \mathbf{b})$$

$$f_3(\mathbf{x}) = 1 - (\exp(-(x - a)^T (x - a)) + \exp(-(\mathbf{x} - \mathbf{b})^T \mathbf{B} (\mathbf{x} - \mathbf{b}))) - \frac{1}{10} \log \left| \frac{1}{100} \mathbf{I} + \mathbf{x} \mathbf{x}^T \right|$$

a) Write  $f_1(\mathbf{x})$  in the completed square form  $(\mathbf{x} - \mathbf{c})^T \mathbf{C} (\mathbf{x} - \mathbf{c}) + c_0$ .

$$\mathbf{c} = \frac{1}{2}(\mathbf{I} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b}) = \begin{pmatrix} \frac{1}{6} \\ \frac{1}{6} \end{pmatrix}$$

$$\mathbf{C} = \mathbf{I} + \mathbf{B} = \begin{pmatrix} 4 & -1 \\ -1 & 4 \end{pmatrix}$$

$$c_0 = -\frac{1}{4}(\mathbf{a} - \mathbf{b})^T (\mathbf{I} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b}) = -\frac{1}{6}$$

3 marks for

- correct  $\mathbf{C}$
- correct  $\mathbf{c}$
- correct  $c_0$

lose a mark for any incorrect notation

b) Explain how you can tell that  $f_1$  has a minimum point. State the minimum value of  $f_1$  and find the input which achieves this minimum.

The eigenvalues of  $\mathbf{I} + \mathbf{B}$  are 5 and 3, which are both positive. Hence the quadratic form is positive definite. The minimum value is  $-\frac{1}{6}$ , which occurs at the point  $(\frac{1}{6}, \frac{1}{6})$ . Or an argument using the Hessian (evals 10 and 6) is fine too. NB need to state which matrix the evals are for

2 marks for

- the appropriate eigenvalues and accompanying explanation
- the minimum value stated and the input which achieves this value

c) Write a python function `grad_f1(x)` that return the gradient for  $f_1$ .

4 marks

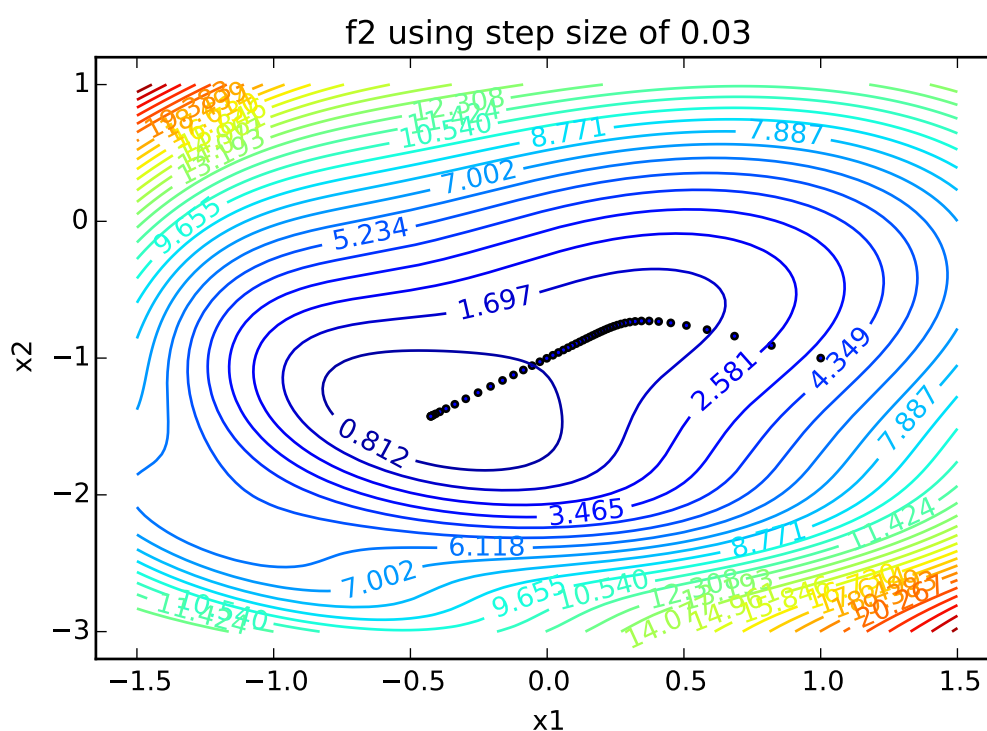
d) Write a python function `grad_f2(x)` that return the gradient for  $f_2$ .

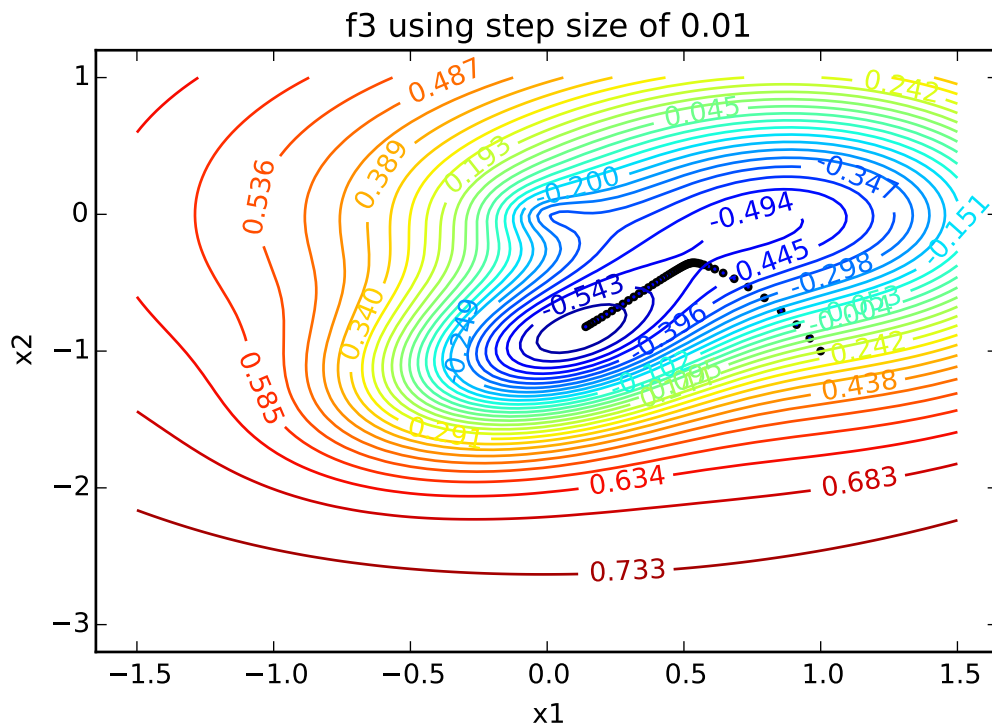
6 marks

- e) Optional: Write a python function `grad_f3(x)` that return the gradient for  $f_3$ . If you don't do this question you still need the gradient for the next part, but you can use an automatic differentiation package. `autograd` is recommended as a good lightweight option, but you could also use e.g. `pytorch`, `tensorflow`, `mxnet` or `chainer`.

There are no marks for this question, but you can still submit the code to labs to verify your answer.

- f) Use your gradients to implement a gradient descent algorithm with 50 iterations to find a local minimum for both  $f_2$  and  $f_3$ . Show the steps of your algorithm on a contour plot of the function. Start from the point  $(1, -1)$  and state the step size you used. Produce separate contour plots for the two functions, using first component of  $\mathbf{x}$  on the  $x$  axis and the second on the  $y$ .





4 marks for:

- sensible scales, plots labeled, contours visible (use at least 20)
- correct functions shown
- gradient descent for first plot with 50 points, converging to minimum
- gradient descent for second plot, converging to one of the minima

g) For the two functions  $f_2$  and  $f_3$ , discuss the qualitative differences you observe when performing gradient descent with step sizes varying between 0.1 and 1, again starting the point  $(1, -1)$ . Briefly describe also what happens in the two cases with grossly mis-specified step-sizes (i.e. greater than 1), with a reason to explain the difference in behaviour.

$f_2$ :

- There is only one minimum, so unless it diverges GD always converges to this point
- For large stepsizes GD diverges

$f_3$

- There are three minima. GD converges to different ones at around 0.3 (we don't get the third unless we start from a different point).
- For large step sizes we observe oscillations and not divergence to infinity.

To explain the difference: the gradients of  $f_2$  blow up as  $|\mathbf{x}| \rightarrow \infty$ , so divergence, whereas the gradients of  $f_3$  tend to zero far away from the stationary points. This is because the exponential terms are of negative definite forms, so gradients  $\rightarrow 0$  very quickly, and the logdet term gradient looks like  $(\mathbf{I} + \mathbf{x}\mathbf{x}^\top)^{-1}\mathbf{x}$  which tends to zero like  $\frac{\mathbf{x}}{r^2}$ .

6 marks for:

- observation about number of minima for  $f_2$
- observation about large stepsizes for  $f_2$
- observation about number of minima for  $f_3$
- observation about different behaviour for different moderate stepsizes for  $f_3$
- observation about large stepsizes for  $f_3$
- explanation for the difference in behaviour