

Modern Optimization Techniques – Group 01

Exercise Sheet 05

Submitted by: Muhammad Inaam Ashraf (Matrikel-Nr: 307524)

Semester 2 MSc. Data Analytics

Question 1: Exact Newton Method

1. Let us (theoretically) optimize the following function using a Newton Descent approach:

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = \exp(x^2 + y^2) = e^{x^2+y^2}$$

(a) Compute the gradient $\nabla f(x, y)$ and the Hessian $\nabla^2 f(x, y)$!

$$\nabla f(x, y) = \begin{bmatrix} \nabla_x f(x, y) \\ \nabla_y f(x, y) \end{bmatrix} = \begin{bmatrix} e^{x^2+y^2} 2x \\ e^{x^2+y^2} 2y \end{bmatrix}$$

$$\nabla^2 f(x, y) = \begin{bmatrix} \nabla_x(\nabla_x f(x, y)) & \nabla_y(\nabla_x f(x, y)) \\ \nabla_x(\nabla_y f(x, y)) & \nabla_y(\nabla_y f(x, y)) \end{bmatrix} = \begin{bmatrix} e^{x^2+y^2} 2(2x^2 + 1) & e^{x^2+y^2} 4xy \\ e^{x^2+y^2} 4xy & e^{x^2+y^2} 2(2y^2 + 1) \end{bmatrix}$$

(b) Compute $\nabla^2 f(x, y)^{-1}$ using Cramer's Rule:

$$\nabla^2 f(x, y)^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$ad - bc = (e^{x^2+y^2} 4x^2 + e^{x^2+y^2} 2)(e^{x^2+y^2} 4y^2 + e^{x^2+y^2} 2) - (e^{x^2+y^2} 4xy)(e^{x^2+y^2} 4xy)$$

$$ad - bc = ((e^{x^2+y^2})^2 4x^2 4y^2 + (e^{x^2+y^2})^2 8x^2 + (e^{x^2+y^2})^2 8y^2 + (e^{x^2+y^2})^2 4) - ((e^{x^2+y^2})^2 4x^2 4y^2)$$

$$ad - bc = (e^{x^2+y^2})^2 8x^2 + (e^{x^2+y^2})^2 8y^2 + (e^{x^2+y^2})^2 4$$

$$\nabla^2 f(x, y)^{-1} = \frac{1}{(e^{x^2+y^2})^2 8x^2 + (e^{x^2+y^2})^2 8y^2 + (e^{x^2+y^2})^2 4} \begin{bmatrix} e^{x^2+y^2} 4y^2 + e^{x^2+y^2} 2 & -e^{x^2+y^2} 4xy \\ -e^{x^2+y^2} 4xy & e^{x^2+y^2} 4x^2 + e^{x^2+y^2} 2 \end{bmatrix}$$

$$\nabla^2 f(x, y)^{-1} = \frac{1}{e^{x^2+y^2} 8x^2 + e^{x^2+y^2} 8y^2 + e^{x^2+y^2} 4} \begin{bmatrix} 4y^2 + 2 & -4xy \\ -4xy & 4x^2 + 2 \end{bmatrix}$$

(c) Compute the update step of the Newton Algorithm:

$$\Delta_{x,y} = -\nabla^2 f(x,y)^{-1} \nabla f(x,y) = -\frac{1}{e^{x^2+y^2} 8x^2 + e^{x^2+y^2} 8y^2 + e^{x^2+y^2} 4} \begin{bmatrix} 4y^2 + 2 & -4xy \\ -4xy & 4x^2 + 2 \end{bmatrix} \begin{bmatrix} e^{x^2+y^2} 2x \\ e^{x^2+y^2} 2y \end{bmatrix}$$

$$\Delta_{x,y} = -\frac{1}{e^{x^2+y^2} 8x^2 + e^{x^2+y^2} 8y^2 + e^{x^2+y^2} 4} \begin{bmatrix} e^{x^2+y^2} 8xy^2 + e^{x^2+y^2} 4x - e^{x^2+y^2} 8xy^2 \\ -e^{x^2+y^2} 8x^2y + e^{x^2+y^2} 8x^2y + e^{x^2+y^2} 4y \end{bmatrix}$$

$$\Delta_{x,y} = -\frac{1}{8x^2 + 8y^2 + 4} \begin{bmatrix} 4x \\ 4y \end{bmatrix} = \begin{bmatrix} -\frac{x}{2x^2 + 2y^2 + 1} \\ -\frac{y}{2x^2 + 2y^2 + 1} \end{bmatrix}$$

Question 2: Quasi-Newton Method: BFGS

Learn a logistic regression using the BFGS Quasi-Newton Method from the lecture, for the following data:

$$X = \begin{bmatrix} 1 & 1 & 5 & 2 & 3 \\ 1 & 2 & 1 & 1 & 2 \\ 1 & 3 & 9 & -2 & 6 \\ 1 & 1 & 4 & 0 & -1 \\ 1 & 2 & 2 & 1 & 3 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Start with a random initialization of β_0 and show the convergence of the algorithm! Also start with a diagonal Hessian in the first iteration and $\mu_0 = 0.001$

Note: Show the working for two iterations without coding the method.

We have

$$\mathcal{L}(X, \beta, Y) = -\sum_{i=1}^m y_i \log(\sigma(x_i \beta)) + (1 - y_i) \log(1 - \sigma(x_i \beta))$$

$$\nabla_{\beta} \mathcal{L}(X, \beta, Y) = -X^T (Y - \hat{Y})$$

Initial Hessian inverse and $\beta^{(0)}$ are given by

$$A^{(0)} = I, \quad \beta^{(0)} = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

We have

$$\Delta \beta^{(0)} = -A^{(0)} \nabla_{\beta^{(0)}} \mathcal{L}(X, \beta^{(0)}, Y)$$

Computing \hat{Y}

$$X\beta^{(0)} = \begin{bmatrix} 1 & 1 & 5 & 2 & 3 \\ 1 & 2 & 1 & 1 & 2 \\ 1 & 3 & 9 & -2 & 6 \\ 1 & 1 & 4 & 0 & -1 \\ 1 & 2 & 2 & 1 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ -2 \\ -4 \\ 2 \end{bmatrix}, \text{ and } \hat{Y} = \frac{1}{1 + e^{-X\beta^{(0)}}} = \begin{bmatrix} 0.269 \\ 0.881 \\ 0.119 \\ 0.018 \\ 0.881 \end{bmatrix}$$

Now

$$\nabla_{\beta^{(0)}} \mathcal{L}(X, \beta^{(0)}, Y) = -X^T(Y - \hat{Y}) = -\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 1 & 2 \\ 5 & 1 & 9 & 4 & 2 \\ 2 & 1 & -2 & 0 & 1 \\ 3 & 2 & 6 & -1 & 3 \end{bmatrix} \left(\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.269 \\ 0.881 \\ 0.119 \\ 0.018 \\ 0.881 \end{bmatrix} \right) = \begin{bmatrix} -0.832 \\ -0.832 \\ -2.868 \\ -1.939 \\ -2.092 \end{bmatrix}$$

$$\Delta\beta^{(0)} = -A^{(0)} \nabla_{\beta^{(0)}} \mathcal{L}(X, \beta^{(0)}, Y) = -\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} -0.832 \\ -0.832 \\ -2.868 \\ -1.939 \\ -2.092 \end{bmatrix} = \begin{bmatrix} 0.832 \\ 0.832 \\ 2.868 \\ 1.939 \\ 2.092 \end{bmatrix}$$

So

$$\beta^{(1)} = \beta^{(0)} + \mu_0 \Delta\beta^{(0)} = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} + 0.001 \begin{bmatrix} 0.832 \\ 0.832 \\ 2.868 \\ 1.939 \\ 2.092 \end{bmatrix} = \begin{bmatrix} \mathbf{1.0008} \\ \mathbf{.0008} \\ \mathbf{-.9971} \\ \mathbf{.0019} \\ \mathbf{1.0021} \end{bmatrix}$$

$$s^{(1)} = \beta^{(1)} - \beta^{(0)} = \begin{bmatrix} 1.0008 \\ .0008 \\ -.9971 \\ .0019 \\ 1.0021 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{.0008} \\ \mathbf{.0008} \\ \mathbf{.0029} \\ \mathbf{.0019} \\ \mathbf{.0021} \end{bmatrix}$$

$$g^{(1)} = \nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y) - \nabla_{\beta^{(0)}} \mathcal{L}(X, \beta^{(0)}, Y)$$

$$\nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y) = -X^T(Y - \hat{Y})$$

Computing \hat{Y}

$$X\beta^{(1)} = \begin{bmatrix} 1 & 1 & 5 & 2 & 3 \\ 1 & 2 & 1 & 1 & 2 \\ 1 & 3 & 9 & -2 & 6 \\ 1 & 1 & 4 & 0 & -1 \\ 1 & 2 & 2 & 1 & 3 \end{bmatrix} \begin{bmatrix} 1.0008 \\ .0008 \\ -.9971 \\ .0019 \\ 1.0021 \end{bmatrix} = \begin{bmatrix} -.974 \\ 2.011 \\ -1.962 \\ -3.989 \\ 2.016 \end{bmatrix}, \text{ and } \hat{Y} = \frac{1}{1 + e^{-X\beta^{(1)}}} = \begin{bmatrix} 0.274 \\ 0.882 \\ 0.123 \\ 0.018 \\ 0.883 \end{bmatrix}$$

$$\nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y) = -X^T(Y - \hat{Y}) = -\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 1 & 2 \\ 5 & 1 & 9 & 4 & 2 \\ 2 & 1 & -2 & 0 & 1 \\ 3 & 2 & 6 & -1 & 3 \end{bmatrix} \left(\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.274 \\ 0.882 \\ 0.123 \\ 0.018 \\ 0.883 \end{bmatrix} \right) = \begin{bmatrix} -0.82 \\ -0.81 \\ -2.8 \\ -1.93 \\ -2.04 \end{bmatrix}$$

$$g^{(1)} = \nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y) - \nabla_{\beta^{(0)}} \mathcal{L}(X, \beta^{(0)}, Y) = \begin{bmatrix} -0.82 \\ -0.81 \\ -2.8 \\ -1.93 \\ -2.04 \end{bmatrix} - \begin{bmatrix} -0.832 \\ -0.832 \\ -2.868 \\ -1.939 \\ -2.092 \end{bmatrix} = \begin{bmatrix} \mathbf{0.012} \\ \mathbf{0.023} \\ \mathbf{0.067} \\ \mathbf{0.005} \\ \mathbf{0.047} \end{bmatrix}$$

Finally,

$$A^{(1)} = A^{(0)} + \frac{(s^{(1)} - A^{(0)}g^{(1)})s^{(1)T} + s^{(1)}(s^{(1)} - A^{(0)}g^{(1)})^T}{s^{(1)T}g^{(1)}} - \frac{(s^{(1)} - A^{(0)}g^{(1)})^T g^{(1)}}{(s^{(1)T}g^{(1)})^2} s^{(1)} s^{(1)T}$$

Now,

$$s^{(1)} - A^{(0)}g^{(1)} = \begin{bmatrix} .0008 \\ .0008 \\ .0029 \\ .0019 \\ .0021 \end{bmatrix} - \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.012 \\ 0.023 \\ 0.067 \\ 0.005 \\ 0.047 \end{bmatrix} = \begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix}$$

$$A^{(1)} = A^{(0)} + \frac{\begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix} \begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} + \begin{bmatrix} .0008 \\ .0008 \\ .0029 \\ .0019 \\ .0021 \end{bmatrix} \begin{bmatrix} -.011 & -.022 & -.065 & -.003 & -.045 \end{bmatrix}}{\begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} \begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix}} - \frac{\begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix} \begin{bmatrix} -.011 & -.022 & -.065 & -.003 & -.045 \end{bmatrix}}{\left(\begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} \begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix} \right)^2}$$

$$A^{(1)} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} + \begin{bmatrix} -.057 & -.085 & -.261 & -.075 & -.185 \\ -.085 & -.113 & -.356 & -.139 & -.254 \\ -.261 & -.356 & -.118 & -.406 & -.796 \\ -.075 & -.139 & -.406 & -.038 & -.283 \\ -.185 & -.254 & -.796 & -.283 & -.567 \end{bmatrix} - \begin{bmatrix} -.045 & -.045 & -.155 & -.105 & -.113 \\ -.045 & -.045 & -.155 & -.105 & -.113 \\ -.155 & -.155 & -.535 & -.362 & -.390 \\ -.105 & -.105 & -.362 & -.245 & -.264 \\ -.113 & -.113 & -.390 & -.264 & -.285 \end{bmatrix}$$

$$A^{(1)} = \begin{bmatrix} .99 & -.04 & -.11 & .03 & -.07 \\ -.04 & .93 & -.2 & -.03 & -.14 \\ -.11 & -.2 & .42 & -.04 & -.41 \\ .03 & -.03 & -.04 & 1.21 & -.02 \\ -.07 & -.14 & -.41 & -.02 & .72 \end{bmatrix}$$

Second iteration:

We have

$$\Delta\beta^{(1)} = -A^{(1)}\nabla_{\beta^{(1)}}\mathcal{L}(X, \beta^{(1)}, Y)$$

$$\Delta\beta^{(1)} = - \begin{bmatrix} .99 & -.04 & -.11 & .03 & -.07 \\ -.04 & .93 & -.2 & -.03 & -.14 \\ -.11 & -.2 & .42 & -.04 & -.41 \\ .03 & -.03 & -.04 & 1.21 & -.02 \\ -.07 & -.14 & -.41 & -.02 & .72 \end{bmatrix} \begin{bmatrix} -0.82 \\ -0.81 \\ -2.8 \\ -1.93 \\ -2.04 \end{bmatrix} = \begin{bmatrix} 0.39 \\ -0.20 \\ 0 \\ 2.17 \\ 0.12 \end{bmatrix}$$

So

$$\beta^{(2)} = \beta^{(1)} + \mu_0 \Delta \beta^{(1)} = \begin{bmatrix} 1.0008 \\ .0008 \\ -.9971 \\ .0019 \\ 1.0021 \end{bmatrix} + 0.001 \begin{bmatrix} 0.39 \\ -0.20 \\ 0 \\ 2.17 \\ 0.12 \end{bmatrix} = \begin{bmatrix} \mathbf{1.0012} \\ \mathbf{.0006} \\ \mathbf{-.9971} \\ \mathbf{4.1042} \\ \mathbf{1.0022} \end{bmatrix}$$

$$s^{(2)} = \beta^{(2)} - \beta^{(1)} = \begin{bmatrix} 1.0012 \\ .0006 \\ -.9971 \\ 4.1042 \\ 1.0022 \end{bmatrix} - \begin{bmatrix} 1.0008 \\ .0008 \\ -.9971 \\ .0019 \\ 1.0021 \end{bmatrix} = \begin{bmatrix} \mathbf{.0008} \\ \mathbf{.0008} \\ \mathbf{.0029} \\ \mathbf{.0019} \\ \mathbf{.0021} \end{bmatrix}$$

$$g^{(2)} = \nabla_{\beta^{(2)}} \mathcal{L}(X, \beta^{(2)}, Y) - \nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y)$$

$$\nabla_{\beta^{(2)}} \mathcal{L}(X, \beta^{(2)}, Y) = -X^T (Y - \hat{Y})$$

Computing \hat{Y}

$$X\beta^{(2)} = \begin{bmatrix} 1 & 1 & 5 & 2 & 3 \\ 1 & 2 & 1 & 1 & 2 \\ 1 & 3 & 9 & -2 & 6 \\ 1 & 1 & 4 & 0 & -1 \\ 1 & 2 & 2 & 1 & 3 \end{bmatrix} \begin{bmatrix} 1.0012 \\ .0006 \\ -.9971 \\ 4.1042 \\ 1.0022 \end{bmatrix} = \begin{bmatrix} -.969 \\ 2.014 \\ -1.966 \\ -3.989 \\ 2.019 \end{bmatrix}, \text{ and } \hat{Y} = \frac{1}{1 + e^{-X\beta^{(2)}}} = \begin{bmatrix} 0.275 \\ 0.882 \\ 0.123 \\ 0.018 \\ 0.883 \end{bmatrix}$$

$$\nabla_{\beta^{(2)}} \mathcal{L}(X, \beta^{(2)}, Y) = -X^T (Y - \hat{Y}) = - \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 1 & 2 \\ 5 & 1 & 9 & 4 & 2 \\ 2 & 1 & -2 & 0 & 1 \\ 3 & 2 & 6 & -1 & 3 \end{bmatrix} \left(\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ -1 \end{bmatrix} - \begin{bmatrix} 0.275 \\ 0.882 \\ 0.123 \\ 0.018 \\ 0.883 \end{bmatrix} \right) = \begin{bmatrix} -0.8189 \\ -0.808 \\ -2.798 \\ -1.93 \\ -2.043 \end{bmatrix}$$

$$g^{(2)} = \nabla_{\beta^{(2)}} \mathcal{L}(X, \beta^{(2)}, Y) - \nabla_{\beta^{(1)}} \mathcal{L}(X, \beta^{(1)}, Y) = \begin{bmatrix} -0.8189 \\ -0.808 \\ -2.798 \\ -1.93 \\ -2.043 \end{bmatrix} - \begin{bmatrix} -0.82 \\ -0.81 \\ -2.8 \\ -1.93 \\ -2.04 \end{bmatrix} = \begin{bmatrix} \mathbf{0.001} \\ \mathbf{0.0008} \\ \mathbf{0.002} \\ \mathbf{0.0033} \\ \mathbf{0.0018} \end{bmatrix}$$

Finally,

$$A^{(2)} = A^{(1)} + \frac{(s^{(2)} - A^{(1)} g^{(2)}) s^{(2)T} + s^{(2)} (s^{(2)} - A^{(1)} g^{(2)})^T}{s^{(2)T} g^{(2)}} - \frac{(s^{(2)} - A^{(1)} g^{(2)})^T g^{(2)}}{(s^{(2)T} g^{(2)})^2} s^{(2)} s^{(2)T}$$

Now,

$$s^{(2)} - A^{(1)} g^{(2)} = \begin{bmatrix} .0008 \\ .0008 \\ .0029 \\ .0019 \\ .0021 \end{bmatrix} - \begin{bmatrix} .99 & -.04 & -.11 & .03 & -.07 \\ -.04 & .93 & -.2 & -.03 & -.14 \\ -.11 & -.2 & .42 & -.04 & -.41 \\ .03 & -.03 & -.04 & 1.21 & -.02 \\ -.07 & -.14 & -.41 & -.02 & .72 \end{bmatrix} \begin{bmatrix} 0.001 \\ 0.0008 \\ 0.002 \\ 0.0033 \\ 0.0018 \end{bmatrix} = \begin{bmatrix} -.0004 \\ -.0001 \\ .0003 \\ -.0017 \\ -.00009 \end{bmatrix}$$

$$\begin{aligned}
 A^{(1)} = A^{(0)} + & \begin{bmatrix} -.0004 \\ -.0001 \\ .0003 \\ -.0017 \\ -.00009 \end{bmatrix} \begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} + \begin{bmatrix} .0008 \\ .0008 \\ .0029 \\ .0019 \\ .0021 \end{bmatrix} \begin{bmatrix} -.0004 & -.0001 & .0003 & -.0017 & -.00009 \end{bmatrix} \\
 & \begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} \begin{bmatrix} 0.001 \\ 0.0008 \\ 0.002 \\ 0.0033 \\ 0.0018 \end{bmatrix} \\
 & - \begin{bmatrix} -.0004 & -.0001 & .0003 & -.0017 & -.00009 \end{bmatrix} \begin{bmatrix} -.011 \\ -.022 \\ -.065 \\ -.003 \\ -.045 \end{bmatrix} \begin{bmatrix} .0008 \\ .0008 \\ .0029 \\ .0019 \\ .0021 \end{bmatrix} \\
 & - \begin{bmatrix} .0008 & .0008 & .0029 & .0019 & .0021 \end{bmatrix} \begin{bmatrix} 0.001 \\ 0.0008 \\ 0.002 \\ 0.0033 \\ 0.0018 \end{bmatrix})^2
 \end{aligned}$$

$$A^{(2)} = \begin{bmatrix} .96 & -.04 & -.09 & .09 & -.07 \\ -.04 & .94 & -.21 & -.066 & -.142 \\ -.09 & -.21 & .418 & .043 & -.4 \\ -.09 & -.066 & .043 & .704 & -.046 \\ -.07 & -.14 & -.40 & -.046 & .716 \end{bmatrix}$$