

Numerical Optimization (COMP0120) Coursework 1

March 7, 2022

EXERCISE 1

We are given a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$.

$$f(x, y) = (y - \cos x)^2 + (y - x)^2$$

- (a) Calculate the gradient ∇f and the Hessian $\nabla^2 f$.

Sol: The gradient is

$$\begin{aligned}\nabla f &= \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} \\ &= \begin{pmatrix} 2(y - \cos x)(-1)(-\sin x) + 2(y - x)(-1) \\ 2(y - \cos x)(1) + 2(y - x) \end{pmatrix} \\ &= \begin{pmatrix} 2y \sin x - 2 \sin x \cos x + 2x - 2y \\ 4y - 2x - 2 \cos x \end{pmatrix}\end{aligned}$$

The Hessian is

$$\begin{aligned}\nabla^2 f &= \begin{pmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{pmatrix} \\ &= \begin{pmatrix} 2y \cos x - 2 \cos 2x + 2 & 2 \sin x - 2 \\ 2 \sin x - 2 & 4 \end{pmatrix}\end{aligned}$$

- (b) Find the minimiser x^* of the function f and describe how you obtained it. (Hint: You may compute it numerically but verify analytically that it satisfies the 1st order necessary condition.) Show that x^* is unique.

Sol: Observe that $f(x, y) \geq 0$ since the sum of squares is non-negative.

$$\begin{aligned}f(x, y)_{\min} &= (y - \cos x)^2 + (y - x)^2 = 0 \\ &\implies \begin{cases} (y - \cos x)^2 = 0 \\ (y - x)^2 = 0 \end{cases} \\ &\implies \begin{cases} y - \cos x = 0 \\ y - x = 0 \end{cases} \\ &\implies x = y = \cos x\end{aligned}$$

The solution plots in Figure 1. There is only one intersection between these two functions.

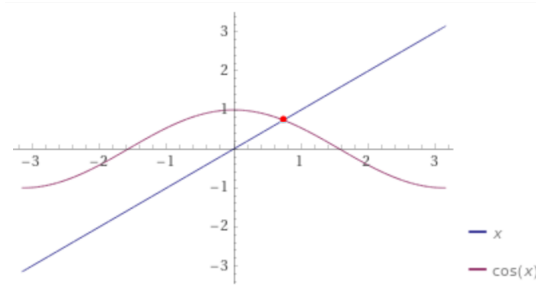


Figure 1: Plot for functions $y = x$ and $y = \cos x$

We cannot get the exact value solving the above equation, hence the solution for minimiser is approximated as $x = y = \cos x \approx 0.739$ and $x^* \approx (0.739, 0.739)^T$.

To verify x^* is the minimiser, we prove it satisfy 1st order necessary condition. f is continuously differentiable and substitute $x = y = \cos x$ into ∇f , check if $\nabla f = 0$.

$$\begin{aligned}\nabla f &= \begin{pmatrix} 2y \sin x - 2 \sin x \cos x + 2x - 2y \\ 4y - 2x - 2 \cos x \end{pmatrix} \\ &= \begin{pmatrix} 2 \cos x \sin x - 2 \sin x \cos x + 2x - 2x \\ 4x - 2x - 2x \end{pmatrix} \\ &= 0\end{aligned}$$

Proof of uniqueness: We want to prove x^* is unique by contradiction. Assume there exists one or more minimisers \tilde{x} ($\tilde{x} \neq x^*$). Solve $f(\tilde{x}, \tilde{y}) = 0$.

$$\begin{aligned}f(\tilde{x}, \tilde{y}) &= (\tilde{y} - \cos \tilde{x})^2 + (\tilde{y} - \tilde{x})^2 = 0 \\ \implies \tilde{x} &= \tilde{y} = \cos \tilde{x} \\ \implies \tilde{x} &= x^*\end{aligned}$$

Figure 1 shows a unique intersection

$\tilde{x} = x^*$ contradicts with $\tilde{x} \neq x^*$. Hence, x^* is unique.

- (c) Plot the function and its contours. Show that the level sets of f are bounded, Hint: To simplify the problem consider the level sets in the limit $|x|, |y|$ large.
Sol:

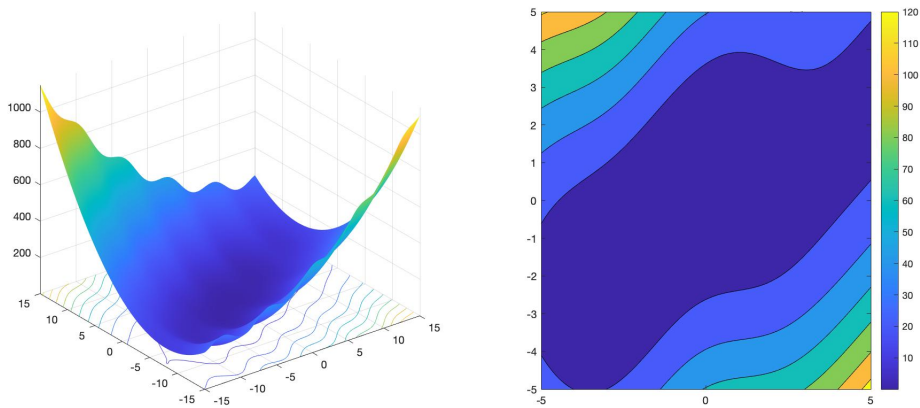


Figure 2: The function(left) and its contours(right)

The right plot in Figure 2 shows the contours of the function after scaling with factor $\alpha = 3$ using simple multiplication $\alpha X, \alpha Y$ and then doing a log transformation. We can find the minimiser x^* in this plot.

Show that the level sets of f are bounded.

When $|x|, |y| \rightarrow \infty$, $f(x, y) = (y - \cos x)^2 + (y - x)^2 \rightarrow y^2 + (y - x)^2$

Level sets of f is $\{(x, y) : y^2 + (y - x)^2 = C\}$, C is a constant.

Since $-(a^2 + b^2) \leq -2ab \leq a^2 + b^2$, $C = y^2 + (y - x)^2 = 2y^2 + x^2 - 2xy \leq 3y^2 + 2x^2$. Level sets of $3y^2 + 2x^2$ are bounded, hence level sets of f are also bounded.

- (d) Show that the gradient, ∇f , is locally Lipschitz continuous. (Hint: Recall integral form of Taylor's theorem and that Frobenious norm is an upper bound on the L_2 matrix norm.)

Sol: If the gradient, ∇f , is locally Lipschitz continuous, then $\exists L > 0$: $\|\nabla f(x) - \nabla f(\bar{x})\| \leq L\|x - \bar{x}\|$.

Using Taylor's theorem, $\nabla f(x + \delta) = \nabla f(x) + \int \nabla^2 f(x + t\delta) \delta dt$, then $\|\nabla f(x + \delta) - \nabla f(x)\| \leq \|\nabla^2 f(x + t\delta)\| \|\delta\|$.

$\|\nabla^2 f(x + t\delta)\| \leq \|\nabla^2 f(x + t\delta)\|_2 \leq \text{constant} \Rightarrow \|\nabla f(x) - \nabla f(\bar{x})\| \leq \text{constant} \times \|\delta\| = L\|x - \bar{x}\|$.

- (e) Show that the Hessian, $\nabla^2 f$, is locally Lipschitz continuous. (Hint: Young's inequality is a useful way of bounding products.)

EXERCISE 2

- (a) Apply steepest descent with an appropriate line search to the function f in **Ex.1** starting from $x_0 = (1, -1)^T$ and $x_0 = (-1, 0)^T$. Plot the iterates over the function contours. State your choice of a line search and any important parameters. What do you observe?

Sol: The starting points are respectively $x_0 = (1, -1)^T$ and $x_0 = (-1, 0)^T$. I applied the steepest descent with backtracking line search. The relevant parameters used are

```
alpha0 = 1; maxIter = 1e4; alpha_max = alpha0;
tol = 1e-6; ls0ptsSteepest.c1 = 1e-4; ls0ptsSteepest.rho = 0.1;
```

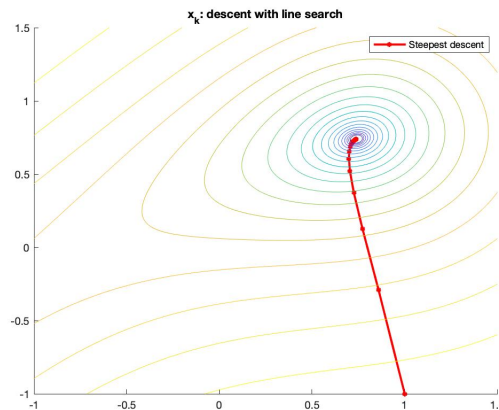


Figure 3: The plot of iterates over the function contours starting at $(1, -1)$

Figure 3 shows that the line goes from the stating point $(1, -1)$ into the plane flatly. It converges to the stationary point quickly and smoothly.

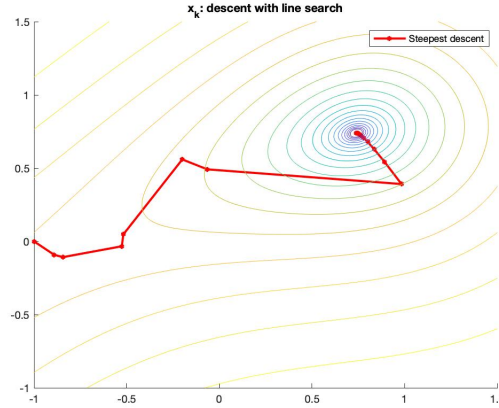


Figure 4: The plot of iterates over the function contours starting at $(-1, 0)$

Figure 4 shows that the line converge zigzag from the stating point $(-1, 0)$ to the stationary point. It rapidly falls off in the first few steps. When going closer to the convergence point, it uses more iterates to reach the minimum.

- (b) Investigate convergence of the steepest descent (a) iterates a posteriori and include one relevant error plot. What are the empirical convergence rates and how did you obtain them? Do they agree with the theoretical predictions? Paraphrase the relevant theoretical results.

Sol:

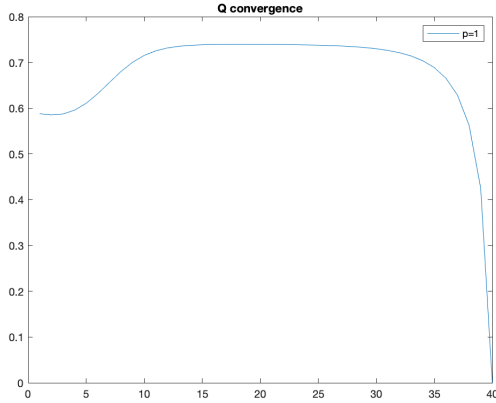


Figure 5: The Q convergence plot with the start point $(1, -1)$

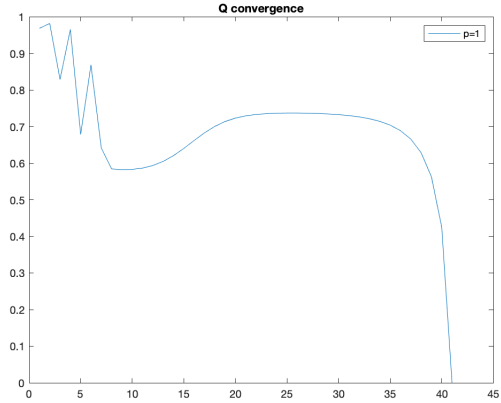


Figure 6: The Q convergence plot with the start point $(-1, 0)$

The convergence can be investigate by Q-linear convergence, which is $\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^p} < r$, $r \in (0, 1)$, k sufficiently large. The empirical convergence rate can be estimated from the error $e_k = \|x_k - x^*\|$. Theoretically, the convergence rate for steepest descent is linear convergence.

- (c) Apply Newton with an appropriate line search to the function f in **Ex.1** starting from $x_0 = (1, -1)^T$ and $x_0 = (-1, 0)^T$. Plot the iterates over the function contours. State your choice of a line search and any important parameters. What do you observe?

Sol: The starting points are respectively $x_0 = (1, -1)^T$ and $x_0 = (-1, 0)^T$. I applied Newton with backtracking line search. The relevant parameters used are

```
alpha0 = 1; maxIter = 1e4; alpha_max = alpha0;
tol = 1e-6; lsOptsNewton.c1 = 1e-4; lsOptsNewton.rho = 0.9;
```

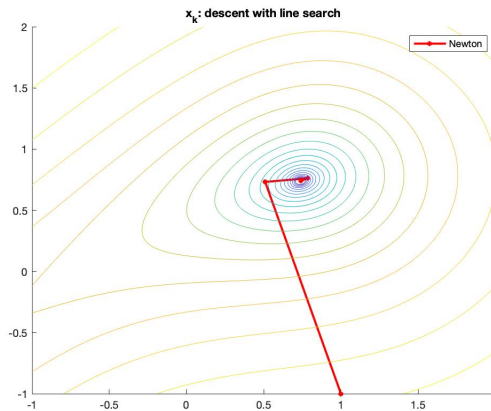


Figure 7: The plot of iterates over the function contours starting at $(1, -1)$

Figure 7 shows that the line goes from stating point $(1, -1)$ into the plane and converges rapidly.

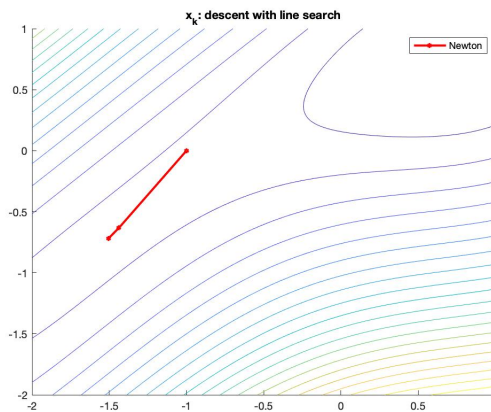


Figure 8: The plot of iterates over the function contours starting at $(-1, 0)$

Figure 8 shows that the line goes from stating point $(-1, 0)$ into the plane and stops after two iterations.

- (d) Investigate convergence of the Newton (c) iterates a posteriori and include one relevant error plot. What are the empirical convergence rates and how did you obtain them? Do

they agree with the theoretical predictions? Paraphrase the relevant theoretical results.
Sol:

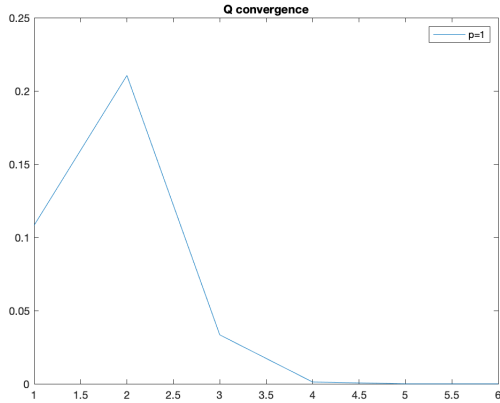


Figure 9: The Q convergence plot with the start point $(1, -1)$

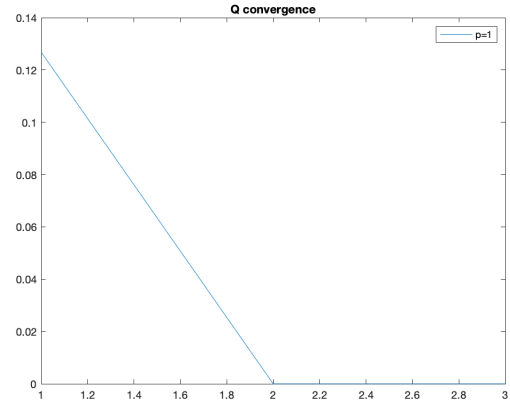


Figure 10: The Q convergence plot with the start point $(-1, 0)$

The empirical convergence rate can be estimated from the error $e_k = \|x_k - x^*\|$. Theoretically, the convergence rate for Newton is quadratic. The right plot above seems to disagree with the theoretical predictions since it has a linear trend.

- (e) Can global convergence of both methods be guaranteed or not and why? Paraphrase the relevant theoretical results.

Sol: Zoutendijk lemma says that for an iteration, $x_{k+1} = x_k + \alpha_k p_k$, $k = 0, 1, 2, \dots$. p_k is the descent direction and α_k satisfies the Wolfe conditions. Steepest descent is used with Wolfe conditions line search. f is non-convex and bounded below. ∇f is locally Lipschitz continuous. Thus global convergence of both methods satisfying Wolfe conditions will be guaranteed.

EXERCISE 3

- (a) Implement the dogleg trust region method for strictly convex functions (with s.p.d. Hessian). Your implementation should return the Cauchy point whenever the gradient and Newton steps are collinear. Include your implementation into the report. Highlight the part where you solve for the intersection point between the trust region and the dogleg path and provide a short narrative explanation.

```
pU = -(F.df(x_k)' * F.df(x_k)) / (F.df(x_k)' * F.d2f(x_k) * F.df(x_k));
pB = -inv(F.d2f(x_k)) * F.df(x_k);

if norm(pB) <= Delta
    % Use Newton, p = pB
    tau = 2;
elseif norm(pU) >= Delta
    % Use steepest descent, p = delta/norm(pU)*p_U;
    tau = Delta/norm(pU);
else
    pC = pB - pU;
    tau = (-pU.' * pC + sqrt((pU.' * pC)^2 - pC.' * pC * (pU.' * pU - Delta^2))) / (pC.' * pC);
    tau = tau + 1;
end

if (0 <= tau) && (tau <= 1)
    p = tau * pU;
elseif (1 <= tau) && (tau <= 2)
    p = pU + (tau - 1) * (pB - pU);
end
```

From slides, $p^B = -B^{-1}g_k$ $p^U = -\frac{g_k^T g_k}{g_k^T B g_k} g_k$

$$p(\tau) = \begin{cases} \tau p^U & , 0 \leq \tau \leq 1 \\ p^U + (\tau - 1)(p^B - p^U) & , 1 \leq \tau \leq 2 \end{cases}$$

(b) Apply the dogleg trust region method to minimise the Rosenbrock function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(x, y) = 100(y - x^2)^2 + (1 - x)^2$$

with two different starting points $x_0 = (0.5, 1)^T$ and $x_0 = (-1.5, 1)^T$. Plot the trajectories traced by the iterates over the function contours. State your choice of the stopping condition and any relevant parameters. What do you observe?

Sol: The starting points are respectively $x_0 = (0.5, 1)^T$ and $x_0 = (-1.5, 1)^T$. I applied the dogleg trust region method to minimise the Rosenbrock function. The relevant parameters used are

```
eta = 0.1; maxIter = 1000;
tol = 1e-6; Delta = 0.2;
```

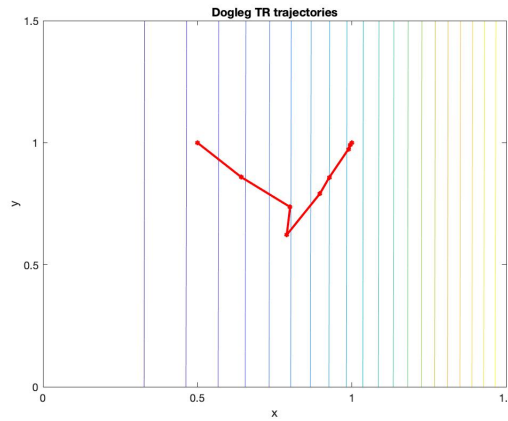


Figure 11: The trajectories traced by the iterates over the function contours starting at $(0.5, 1)$

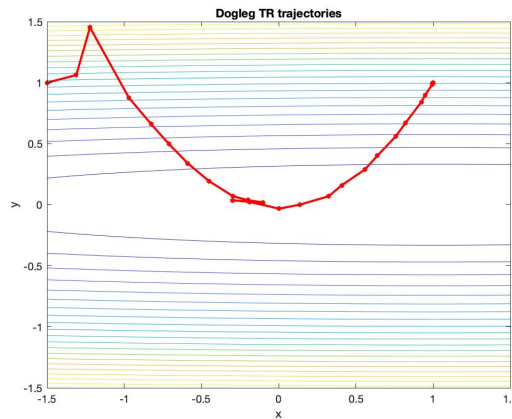


Figure 12: The trajectories traced by the iterates over the function contours starting at $(-1.5, 1)$

Figure 11 shows that the line goes from the starting point $(0.5, 1)$ into the plane. It rapidly converges to the stationary point.

Figure 12 shows that the line goes from the starting point $(-1.5, 1)$ into the plane. It takes many iterations to converge to the stationary point.

- (c) Investigate convergence of the dogleg iterates in **(b)** a posteriori and include one relevant error plot. What are the empirical convergence rates and how did you obtain them? Do they agree with the theoretical predictions? Paraphrase the relevant theoretical results.

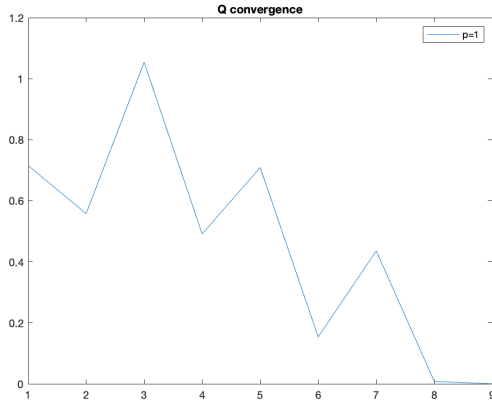


Figure 13: The Q convergence plot with the start point $(0.5, 1)$

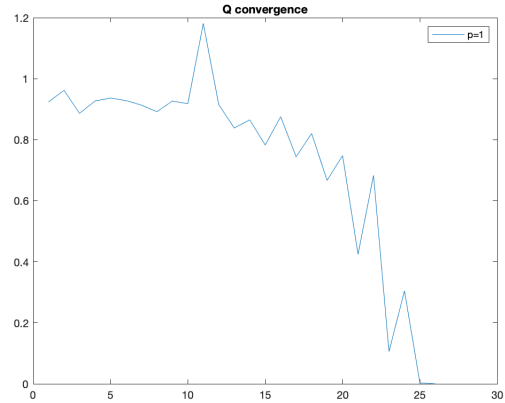


Figure 14: The Q convergence plot with the start point $(-1.5, 1)$

The empirical convergence rate can be estimated from the error $e_k = \|x_k - x^*\|$. The plots above show the error plots of the dogleg. Both of them converge within 25 iterations. The start point $(-1.5, 1)$ for Figure 14 uses more iterations to converge than $(0.5, 1)$ for Figure 13. This might be because the distance between $(-1.5, 1)$ and the minimum is larger than that between $(0.5, 1)$ and the minimum. Theoretically, the convergence rate for dogleg is superlinear local convergence. The right plot does not strictly follow linear convergence.

- (d) Can global convergence be expected or not, and why? Paraphrase the relevant theoretical results.

Sol: It is globally convergent. The dogleg algorithm gives p_k s.t. $m_k(p_k) \leq m_k(p_k^C)$, so it satisfies

$$m_k(0) - m_k(p) \geq c_1 \|g_k\| \min(\|\nabla_k\|, \frac{\|g_k\|}{\|B_k\|})$$

f is bounded below on the level set S and Lipschitz continuously differentiable in an open neighborhood of S , $\|p_k\| \leq \gamma \|\nabla_k\|$, $\gamma \geq 1$, then it has $\liminf \|g_k\| = 0$ and then $\lim \|g_k\| = 0$.

EXERCISE 4

- (a) Implement the Polar-Ribiere conjugate gradient method. (Hint: Modify the descentLineSearch.m template from tutorial 2.) Copy the relevant lines in your report.


```

df_k = F.df(x_k);
p_k = -df_k;
x_k1=x_k+alpha_k*p_k;
df_k1=F.df(x_k1);
B_k1 = (df_k1'*(df_k1 - df_k))/norm(df_k)^2;
p_k1 = -F.df(x_k1)+B_k1*p_k;

```

(b) Apply Polar-Ribiere conjugate gradient method to minimise the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(x, y) = x^2 + 4y^4 + 2y^2.$$

Try two initial points $x_0 = (-1, 2)^T$ and $x_0 = (-1, -0.25)^T$ and set the tolerance $tol = 1e^{-4}$. Plot the iterates over the function contours. State your choice of any relevant parameters.

Sol:

```

alpha0 = 1; tol = 1e-4; maxIter = 100;
ls0ptsCG_LS.c1 = 1e-4; ls0ptsCG_LS.c2 = 0.1;

```

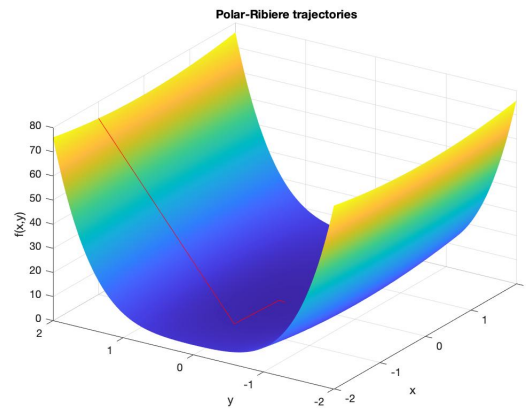


Figure 15: The trajectories traced by the iterates over the function contours starting at $(-1, 2)$

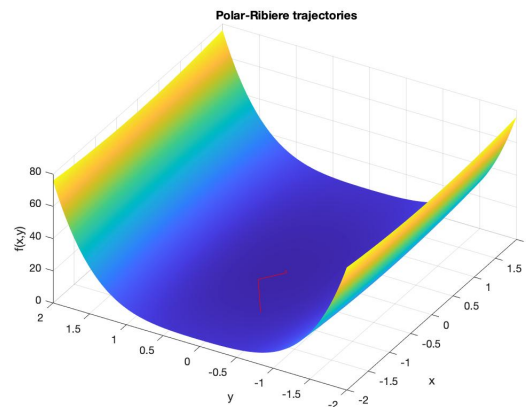


Figure 16: The trajectories traced by the iterates over the function contours starting at $(-1, -0.25)$

- (c) What is the main limitation of Polar-Ribiere method, do you observe it or any other problems in optimisation in (b), can you explain?

Sol: The main limitation of Polar-Ribiere is that it might not ensure global convergence. For example, if in some iteration, β_{k+1}^{PR} generated might be negative so that new p_k will not be the descent direction.

- (d) Can global convergence be guaranteed for this problem or not and why? If not, can you modify your Polar-Ribiere solver to guarantee global convergence and how? (Hint: What are the conditions under which global convergence can be proven for Polar-Ribiere?) Implement the modification, rerun your solver and compare its behaviour to the one in (b). Paraphrase the relevant theoretical results.

Sol: Global convergence is not guaranteed because β_{k+1}^{PR} might be negative and p_k will not be the descent direction. To solve this, we change the formular to $\beta_{k+1}^{newPR} = \max\{\beta_{k+1}^{PR}, 0\}$.

- (e) Investigate convergence of the Polar-Ribiere method in (b) a posteriori and include one relevant error plot. What are the empirical convergence rates and how did you obtain them? Do they agree with the theoretical predictions? Paraphrase the relevant theoretical results.

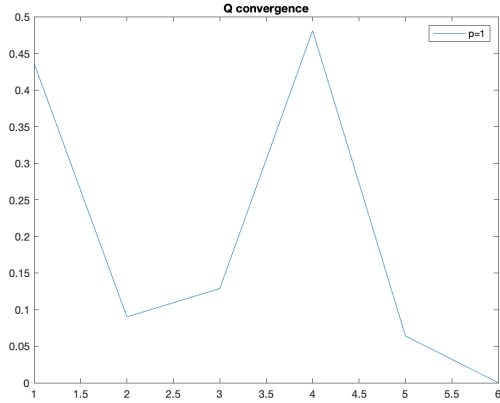


Figure 17: The Q convergence plot with the start point $(-1, 2)$

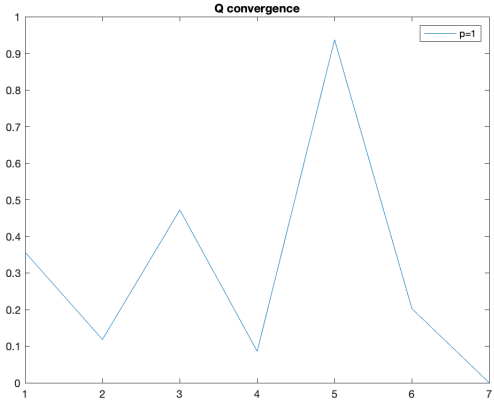


Figure 18: The Q convergence plot with the start point $(-1, -0.25)$

The empirical convergence rate can be estimated from the error $e_k = \|x_k - x^*\|$.

EXERCISE 5

- (a) Implement the BFGS method by modifying the descentLineSearch.m function from tutorial 2. Make your implementation efficient i.e. avoid explicitly forming the inverse Hessian matrix H_k . Copy the code lines implementing the update of H_k into your report and briefly explain what makes your implementation efficient.

```
H_k = @(x) eye(length(x));
p_k = -H_k(x_k)*F.df(x_k);
s_k = x_k - x_k0;
y_k = F.df(x_k) - F.df(x_k0);
rho = 1/(y_k'*s_k);
H_k = @(x) (eye(length(x)) - rho*s_k*y_k')*H_k(x_k0)*(eye(length(x)) - rho*y_k*s_k') + rho*s_k*(s_k');
```

(b) Minimise the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(x, y) = (x - 3y)^2 + x^4.$$

using BFGS implemented in (a) starting from $x_0 = (10, 10)^T$. Visualise the path traced by the iterates over the function contour. Which line search did you choose and why? State your choices of the relevant line search parameters.

Sol: The starting points are respectively $x_0 = (10, 10)^T$. I applied BFGS and the relevant parameters used are

```
maxIter = 300; tol = 1e-10; alpha0 = 1;
ls0pts_LS.c1 = 1e-4; ls0pts_LS.c2 = 0.5;
```

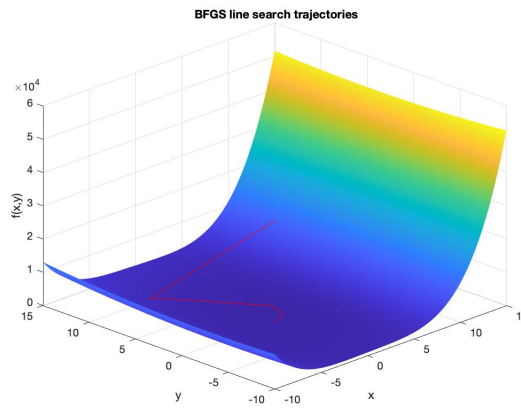


Figure 19: The trajectories traced by the iterates over the function contours starting at $x_0 = (10, 10)^T$

(c) Investigate convergence of the BFGS iterates in (b) a posteriori and include one relevant error plot. What is the empirical convergence rate and how did you obtain it? Does it agree with the theoretical prediction? Paraphrase the relevant theoretical result.

Sol:

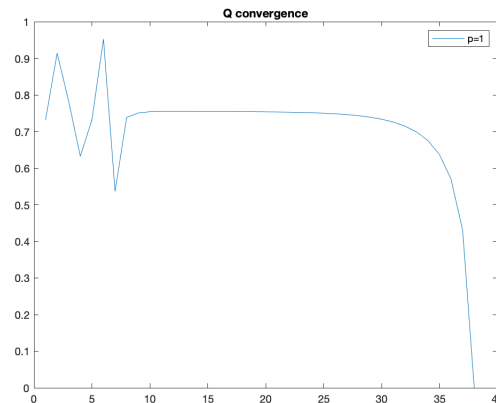


Figure 20: The Q convergence plot with the start point $x_0 = (10, 10)^T$

- (d) Can global convergence of BFGS and of what type be expected or not, and why? Paraphrase the relevant theoretical result.

Sol: The function is convex bounded below and Lipschitz continuously differentiable, then BFGS satisfying Wolfe line search condition is global convergence and it is superlinear convergence.

- (e) Investigate the quality of the respective (inverse) Hessian approximations computed by BFGS and SR-1. What do you observe? Do your observations agree with theoretical predictions? Paraphrase the relevant theoretical results.

Sol: (i)

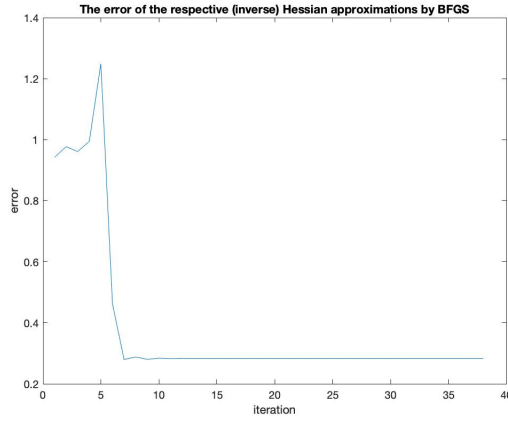


Figure 21: The error plot of the respective (inverse) Hessian approximations computed by BFGS ($\{\|I - H_k^{BFGS} \nabla^2 f(x_k)\|_2\}_{k \leq 0}$)

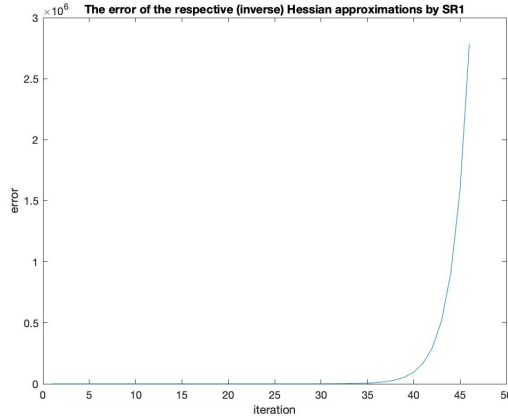


Figure 22: The error plot of the respective (inverse) Hessian approximations computed by SR1 ($\{\|I - (\nabla^2 f(x_k))^{-1} B_k^{SR1}\|_2\}_{k \leq 0}$)

(ii)

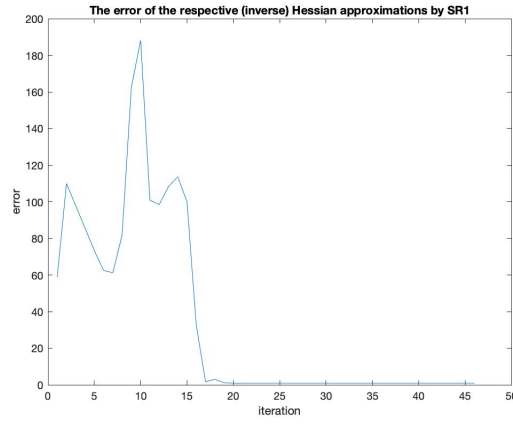


Figure 23: The error plot of the respective (inverse) Hessian approximations computed by SR1 ($\{\|B_k^{SR1} - \nabla^2 f(x_k)\|_2\}_{k \leq 0}$)

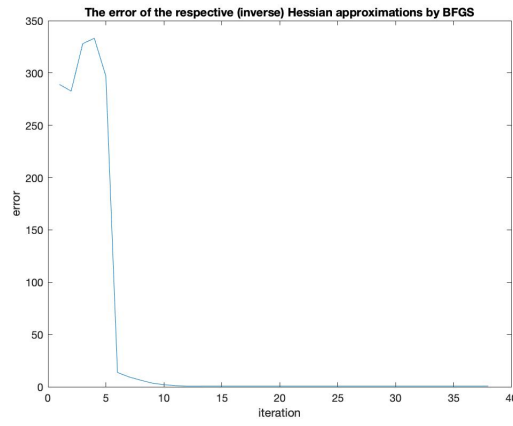


Figure 24: The error plot of the respective (inverse) Hessian approximations computed by BFGS ($\{\|(H_k^{BFGS})^{-1} - \nabla^2 f(x_k)\|_2\}_{k \leq 0}$)