

HT 2019

PROBLEM SHEET 3

Integration in Higher Dimensions, Convergence, 1-Dimensional Root-Finding

3.1. The 2-dimensional Trapezium rule, with  $m$  strips in each dimension, is

$$\int_{[a,b] \times [c,d]} f(x,y) d(x,y) \approx \sum_{i=0}^n \sum_{j=0}^n w_{ij} f\left(a+i\left(\frac{b-a}{n}\right), c+j\left(\frac{d-c}{n}\right)\right)$$

where  $\underline{W} = (w_{ij})$  is an  $(n+1) \times (n+1)$  matrix of scalars depending on  $m, a, b, c, d$ . We can derive it by iterating the Trapezium rule first considering  $g(x) = \int_c^d f(x,y) dy$  :

$$\int_{[a,b] \times [c,d]} f(x,y) d(x,y) = \int_a^b \int_c^d f(x,y) dy dx \approx \text{(the Composite Trapezium Rule with } n \text{ strips)}$$

$$\approx \frac{b-a}{2n} \left[ \int_c^d f(a,y) dy + 2 \sum_{i=1}^{n-1} \int_c^d f\left(a+i\left(\frac{b-a}{n}\right), y\right) dy + \int_c^d f(b,y) dy \right] \approx$$

(By applying the Composite Trapezium Rule with  $n$  strips for each integral we obtain:)

$$\begin{aligned} & \approx \frac{b-a}{2n} \cdot \frac{d-c}{2n} \left[ \left( f(a,c) + 2 \sum_{j=1}^{n-1} f\left(a, c+j\left(\frac{d-c}{n}\right)\right) + f(a,d) \right) + \left( 2 \sum_{i=1}^{n-1} f\left(a+i\left(\frac{b-a}{n}\right), c\right) + \right. \right. \\ & + 4 \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} f\left(a+i\left(\frac{b-a}{n}\right), c+j\left(\frac{d-c}{n}\right)\right) + 2 \sum_{i=1}^{n-1} f\left(a+i\left(\frac{b-a}{n}\right), d\right) \left. \right) + \left( f(b,c) + \right. \\ & \left. \left. + 2 \sum_{j=1}^{n-1} f\left(b, c+j\left(\frac{d-c}{n}\right)\right) + f(b,d) \right) \right] = \frac{(b-a)(d-c)}{4n^2} \sum_{i=0}^n \sum_{j=0}^n w_{ij} f\left(a+i\left(\frac{b-a}{n}\right), c+j\left(\frac{d-c}{n}\right)\right), \end{aligned}$$

where

$$w_{ij} = \begin{cases} 1, & \text{if } (i,j) \in \{(0,0), (0,n), (n,0), (n,n)\} \\ 2, & \text{if } (i \in \{0,n\} \text{ and } j \notin \{0,n\}) \text{ OR } (i \notin \{0,n\} \text{ and } j \in \{0,n\}) \\ 4, & \text{if } i \notin \{0,n\} \text{ and } j \notin \{0,n\} \end{cases}$$

$$W = \begin{pmatrix} 1 & 2 & 2 & 2 & \dots & 2 & 2 & 1 \\ 2 & 4 & 4 & 4 & \dots & 4 & 4 & 2 \\ 2 & 4 & 4 & 4 & \dots & 4 & 4 & 2 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 2 & 4 & 4 & 4 & \dots & 4 & 4 & 2 \\ 1 & 2 & 2 & 2 & \dots & 2 & 2 & 1 \end{pmatrix},$$

In 3 dimensions, I would expect:

$$w_{ijk} = \begin{cases} 1, & \text{if } i \in \{0, n\} \text{ and } j \in \{0, n\} \text{ and } k \in \{0, n\} \\ 2, & \text{if exactly two of the indices } \in \{0, n\} \\ 4, & \text{if exactly one of the indices } \in \{0, n\} \\ 8, & \text{if } i \notin \{0, n\} \text{ and } j \notin \{0, n\} \text{ and } k \notin \{0, n\} \end{cases}$$

3.2. Let's start with the Monte Carlo estimate using  $N$  random variables (i.i.d.):

$$MC_N[f, R] = A(R) \frac{1}{N} \sum_i f(\underline{x}_i)$$

Then, by adding one random sample, we get:

$$\begin{aligned} MC_{(N+1)}[f, R] &= A(R) \frac{1}{N+1} \sum_{i=1}^{N+1} f(\underline{x}_i) = A(R) \frac{1}{N+1} \left( \sum_{i=1}^N f(\underline{x}_i) + f(\underline{x}_{N+1}) \right) = \\ &= A(R) \cdot \frac{1}{N+1} \cdot \frac{N}{A(R)} MC_N[f, R] + A(R) \cdot \frac{1}{N+1} f(\underline{x}_{N+1}) \end{aligned}$$

$$\text{So, we get that } MC_{(N+1)}[f, R] = \frac{N}{N+1} MC_N[f, R] + A(R) \cdot \frac{1}{N+1} f(\underline{x}_{N+1}).$$

The error of the Monte Carlo estimate using  $N$  r.v.s. is

$$e_N = \text{err}(MC_N)[f, R] = \sqrt{\frac{V}{N}}, \text{ where } V \text{ is a constant.}$$

We will prove that:

$$\textcircled{A} \quad \frac{e_{N+1}}{e_N} \xrightarrow{N \rightarrow \infty} 1 : \lim_{N \rightarrow \infty} \frac{e_{N+1}}{e_N} = \lim_{N \rightarrow \infty} \frac{\sqrt{\frac{V}{N+1}}}{\sqrt{\frac{V}{N}}} = \lim_{N \rightarrow \infty} \sqrt{\frac{N}{N+1}} = 1$$

$$\textcircled{B} \quad \frac{|e_{N+2} - e_{N+1}|}{|e_{N+1} - e_N|} \xrightarrow{N \rightarrow \infty} 1 :$$

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{|e_{N+2} - e_{N+1}|}{|e_{N+1} - e_N|} &= \lim_{N \rightarrow \infty} \frac{\left| \sqrt{\frac{V}{N+2}} - \sqrt{\frac{V}{N+1}} \right|}{\left| \sqrt{\frac{V}{N+1}} - \sqrt{\frac{V}{N}} \right|} = \lim_{N \rightarrow \infty} \frac{\frac{1}{\sqrt{N+1}} - \frac{1}{\sqrt{N+2}}}{\frac{1}{\sqrt{N}} - \frac{1}{\sqrt{N+1}}} = \\ &= \lim_{N \rightarrow \infty} \frac{\sqrt{N+2} - \sqrt{N+1}}{\sqrt{N+1} \cdot \sqrt{N+2}} \cdot \frac{\sqrt{N} \cdot \sqrt{N+1}}{\sqrt{N+1} - \sqrt{N}} = \lim_{N \rightarrow \infty} \sqrt{\frac{N}{N+2}} \cdot \frac{(N+2)(N+1)}{\sqrt{N+2} + \sqrt{N+1}} \cdot \frac{\sqrt{N+1} + \sqrt{N}}{(N+1) - N} = \\ &= \lim_{N \rightarrow \infty} \sqrt{\frac{N}{N+2}} \cdot \frac{\sqrt{N+1} + \sqrt{N}}{\sqrt{N+2} + \sqrt{N+1}} = \lim_{N \rightarrow \infty} \sqrt{\frac{1}{1 + \frac{2}{N}}} \cdot \frac{\sqrt{1 + \frac{1}{N}} + 1}{\sqrt{1 + \frac{2}{N}} + \sqrt{1 + \frac{1}{N}}} = 1 \cdot \frac{2}{2} = 1 \end{aligned}$$

From  $\textcircled{A}$  and  $\textcircled{B}$  we deduce that the iterative algorithm we use for the Monte Carlo estimate has logarithmic convergence.

3.3. We want to approximate  $\int_R f(\underline{x}) d\underline{x}$  by partitioning  $R$  into  $k$  predetermined disjoint sub-regions  $R_1, R_2, \dots, R_k$  and using Monte Carlo integration on each sub-region:

$$\int_R f(\underline{x}) d\underline{x} = \sum_i \int_{R_i} f(\underline{x}) d\underline{x} \approx \sum_i MC_{N_i}[f, R_i] \quad (*)$$

We will find the variance of  $(*)$ :

$$\begin{aligned} \text{Var} \left[ \sum_{i=1}^k MC_{N_i}[f, R_i] \right] & \stackrel{\text{independent r.v.s.}}{=} \sum_{i=1}^k \text{Var} [MC_{N_i}[f, R_i]] = \sum_{i=1}^k \text{Var} [A(R_i) \cdot \frac{1}{N_i} \cdot \sum_{j=1}^{N_i} f(\underline{x}_{ij})] = \\ & = \sum_{i=1}^k A^2(R_i) \cdot \frac{1}{N_i^2} \text{Var} \left[ \sum_{j=1}^{N_i} f(\underline{x}_{ij}) \right] \end{aligned}$$

Now, we'll calculate the general term  $\text{Var} \left[ \sum_{j=1}^{N_i} f(\underline{x}_{ij}) \right]$  for all  $i \in \{1, 2, \dots, k\}$ :

$$\text{Var} \left[ \sum_{j=1}^{N_i} f(\underline{x}_{ij}) \right] \stackrel{\text{independent r.v.s.}}{=} \sum_{j=1}^{N_i} \text{Var} [f(\underline{x}_{ij})] \stackrel{\text{identically distributed}}{=} N_i \text{Var}[f(\underline{x}_i)]$$

Going back, we obtain:

$$\text{Var} \left[ \sum_{i=1}^k MC_{N_i}[f, R_i] \right] = \sum_{i=1}^k A^2(R_i) \cdot \frac{1}{N_i^2} \cdot N_i \cdot \text{Var}[f(\underline{x}_i)] = \sum_{i=1}^k \frac{A^2(R_i)}{N_i} \text{Var}[f(\underline{x}_i)].$$

Now, we want to find the optimal allocation of samples between sub-regions. That means that we want to minimize  $\text{Var} \left[ \sum_{i=1}^k MC_{N_i}[f, R_i] \right] = \sum_{i=1}^k \frac{A^2(R_i)}{N_i} \text{Var}[f(\underline{x}_i)]$ . As  $A(R_i)$  and  $\text{Var}[f(\underline{x}_i)]$  are considered in these case given constants, we wish to solve the constrained optimization problem: ( $a_i = A^2(R_i) \text{Var}[f(\underline{x}_i)]$ )

$$\underset{N_1, N_2, \dots, N_k}{\text{minimize}} \sum_{i=1}^k \frac{a_i}{N_i} \text{ subject to } \sum_{i=1}^k N_i = N, \forall i. N_i \geq 0.$$

We form the Lagrangian:

$$\Lambda(1, \underline{x}) = f(\underline{x}) - 1g(\underline{x}), \text{ where } \underline{x} = \begin{pmatrix} N_1 \\ \vdots \\ N_k \end{pmatrix}, f(\underline{x}) = \sum_{i=1}^k \frac{a_i}{N_i}, g(\underline{x}) = \sum_{i=1}^k N_i - N$$

And we want to find its stationary points:

- $g(\underline{x}) = 0 \Rightarrow \sum_{i=1}^k N_i - N = 0$
- $\frac{\partial f}{\partial N_i} = 1 \frac{\partial g}{\partial N_i} \Rightarrow -\frac{a_i}{N_i^2} = 1 \Rightarrow$  We will use  $\alpha = -1$ . Then, we have

$$\frac{a_i}{N_i^2} = \alpha \Rightarrow N_i = \sqrt{\frac{a_i}{\alpha}}$$

Putting this back in  $g(\underline{x}) = 0$ , we get

$$\sum_{i=1}^k \sqrt{\frac{a_i}{\alpha}} - N = 0 \Rightarrow \frac{\sum_{i=1}^k \sqrt{|a_i|}}{\sqrt{|\alpha|}} = N \Rightarrow |\alpha| = \left( \frac{\sum_{i=1}^k \sqrt{|a_i|}}{N} \right)^2 \Rightarrow N_i = \frac{N \sqrt{|a_i|}}{\sum_{i=1}^k \sqrt{|a_i|}}$$

By replacing  $a_i$  by its value, we get

$$N_i = \frac{N \sqrt{|A^2(R_i) \text{Var}[f(x_i)]|}}{\sum_{j=1}^K \sqrt{|A^2(R_j) \text{Var}[f(x_j)]|}} = \frac{A(R_i) N \sqrt{|\text{Var}[f(x_i)]|}}{A(R_i) \sum_{j=1}^K \sqrt{|\text{Var}[f(x_j)]|}}$$

$$N_i = \frac{N \sqrt{|\text{Var}[f(x_i)]|}}{\sum_{j=1}^K \sqrt{|\text{Var}[f(x_j)]|}}$$

For computing the Monte Carlo estimate, we need to calculate:

$$\sum_{i=1}^K \text{MC}_{N_i} [f, R_i] = \sum_{i=1}^K A(R_i) \frac{1}{N_i} \sum_{j=1}^{N_i} f(\underline{x}_{ij}) = \sum_{i=1}^K A(R_i) \cdot \frac{\sum_{j=1}^{N_i} \sqrt{|\text{Var}[f(x_{ij})]|}}{N_i \sqrt{|\text{Var}[f(x_i)]|}} \cdot \sum_{j=1}^{N_i} f(\underline{x}_{ij})$$

We are given:

- $R_1, R_2, \dots, R_K$  sub-regions of  $R$
- $\text{Var}[f(x_i)]$ , for all  $i \in \{1, 2, \dots, K\}$
- $f$ , so we can calculate  $f(\underline{x}_{ij})$ , where  $\underline{x}_{i1}, \underline{x}_{i2}, \dots, \underline{x}_{iN_i}$  uniformly distributed on  $R_i$ ; independent and
- $N$  (we choose an  $N$  that is big enough to provide a good approximation for  $\int_R f(\underline{x}) d\underline{x}$ ).

The algorithm is:

1. Calculate  $A(R_1), A(R_2), \dots, A(R_K)$
2. Calculate  $\sqrt{|\text{Var}[f(x_i)]|}$  for all  $i \in \{1, 2, \dots, K\}$
3. Calculate every  $N_i = \left\lceil \frac{N \sqrt{|\text{Var}[f(x_i)]|}}{\sum_{j=1}^K \sqrt{|\text{Var}[f(x_j)]|}} \right\rceil$ , for  $i \in \{1, 2, \dots, K\}$  (we want integers here)
4. Calculate every  $f(\underline{x}_{ij})$ , for all  $i \in \{1, 2, \dots, K\}$  and  $j \in \{1, 2, \dots, N_i\}$
5.  $\text{MC}_{N_i} [f, R_i]$  will be equal to  $\sum_{i=1}^K A(R_i) \cdot \frac{1}{N_i} \cdot \sum_{j=1}^{N_i} f(\underline{x}_{ij})$ .

If we consider  $\text{Var}[f(x_i)]$  unknown, we need to add iteratively (starting from two samples for each region) according to the current estimates on each region. We add more samples to the regions which need to be evaluated more.

3.5.  $L > 0$ , constant

$$x_{n+1} = x_n + L e^{-x_n} - 1$$

(a) In general, for  $f: \mathbb{R} \rightarrow \mathbb{R}$ , the iterative step of Newton's method for finding the root  $x^*$  of  $f$  is:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

For  $f(x) = \frac{1}{L} e^x - 1$ , we have

$$f'(x) = f''(x) = \frac{1}{L} e^x, \text{ so the iterative step is:}$$

$$x_{n+1} = x_n - \frac{\frac{1}{L} e^{x_n} - 1}{\frac{1}{L} e^{x_n}} = x_n - 1 + L e^{-x_n}, \text{ so we have the function } f.$$

Now, let's calculate its root  $x^*$ :

$$f(x^*) = 0 \Rightarrow \frac{1}{L} e^{x^*} - 1 = 0 \Rightarrow e^{x^*} = L \Rightarrow x^* = \ln L$$

(b) Let  $i = (x^* - c, x^* + c)$  and  $f' = f''$  are strictly increasing

$$A(c) = \frac{\max_{\beta \in i} |f''(\beta)|}{\min_{\alpha \in i} |f'(\alpha)|} \stackrel{L > 0}{=} \frac{f''(x^* + c)}{f'(x^* - c)} = \frac{\frac{1}{L} e^{x^* + c}}{\frac{1}{L} e^{x^* - c}} = e^{x^* + c - x^* + c} \Rightarrow$$

$$\Rightarrow A(c) = e^{2c}, \text{ with } c > 0$$

If we start with  $x_0 \in (\ln L - c, \ln L + c)$ , with  $\frac{c A(c)}{2} < 1$ , then Newton's method converges at least quadratically to  $x^*$ . So, if we take for example  $c = \frac{1}{2}$ , we then have

$$\frac{c A(c)}{2} = \frac{\frac{1}{2} \cdot e}{2} = \frac{1}{4} e < 1 \text{ (True), so we can start with } x_0 \in (\ln L - \frac{1}{2}, \ln L + \frac{1}{2}).$$

(c) If  $L=2$  and  $x_0=100$ , the error, which will initially be  $e_0 = 100 - \ln 2$  will decrease by approximately 1 at each step, as  $x_n - x_{n+1} = L e^{-x_n} - 1$ , and as we start with a big positive value of  $x_0$ ,  $L e^{-x_n}$  will be very small compared to 1, so we decrease the error by 1, therefore the approximate number of iterations needed is  $> 100$  (actually we need 104 steps, so our estimate is pretty close).

If  $L=2$  and  $x_0=-100$ , we first see that  $x_1$  will be  $[2e^{100}-101]$ , which is very big and positive, so, by the same method as at the previous case, we will need approximately  $2e^{100}-100$  steps to get close to  $x^* = \ln 2$ .

(d) We now do not have a way to know the answer, which is  $\ln L$  for a constant  $L > 0$ , so we need to choose  $x_0$  in a different way to ensure that our method converges to the correct result.

First, we'll start with an interval in which we are sure that  $x^* = \ln L$  is.

For  $L > 1$ , we have  $\ln L < L$ , as the function  $g(x) = x - \ln x$  has  $g'(x) = 1 - \frac{1}{x} > 0$ , for  $x > 1$ , so  $g$  is strictly increasing on  $(1, +\infty)$   $\Rightarrow g(x) > g(1)$  for  $x > 1 \Rightarrow g(x) > 0$  as  $g(1) = 1 > 0$ .  
Also, obviously  $\ln L > -\frac{1}{L}$  for  $L > 1$  as  $\ln L > 0$ .

For  $L \in (0, 1)$  we have  $\ln L > -\frac{1}{L}$ , as we can say that  $L = \frac{1}{x}$ , with  $x > 1$  and we get to  $\ln \frac{1}{x} > -x$ , or  $-\ln x > -x$ , or  $\ln x < x$ , which we already proved above.

Also, obviously  $\ln L < L$  for  $L \in (0, 1)$  as  $\ln L < 0$ .

So, we are sure that  $\ln L \in (-\frac{1}{L}, L)$ , values that we can calculate.

$$(a_0, b_0) = \left(-\frac{1}{L}, L\right)$$

Now, we'll use the "interval bisection" method until we get to an interval  $(a, b)$ , in which we know that we have the root, such that  $\frac{b-a}{4} < \frac{1}{8}$ , or  $b-a < \frac{1}{2}$ . This because we can write  $(a, b) = (m-2c, m+2c)$ , with  $m = \frac{a+b}{2}$  and  $c = \frac{b-a}{2}$  and we want that

$$B(c) = \frac{\max_{p \in [m-2c, m+2c]} |f''(p)|}{\min_{p \in [m-2c, m+2c]} |f'(p)|} = e^{4c}, \text{ where } i = (m-2c, m+2c) \text{ to have } \frac{cB(c)}{2} < 1, \text{ or } ce^{4c} < 2, \text{ and this happens}$$

for all  $c \in (0, \frac{1}{4})$ . So, we run the method until  $b-a < \frac{1}{2}$  and then initialize  $x_0 = m$ , so

$x_0 = \frac{a+b}{2}$ , so  $x^* \in (x_0 - c, x_0 + c)$ , where  $c = \frac{b-a}{2} < \frac{1}{4}$ , therefore we can run now Newton's method, which will converge quadratically to the root and we do this until

$|x_n - x_{n-1}| < 10^{-10} |x_n|$ , therefore the absolute error will be  $10^{-10} \ln 2$ , so the approximate root will be  $\ln 2 \pm 10^{-10} \ln 2$  ( $x_n$ ).

3.6] Let  $f: [a, b] \rightarrow \mathbb{R}$  convex, differentiable, and strictly decreasing, with  $f(b) < 0$ .

Let  $x_0 \in [a, b]$  with  $f(x_0) \geq 0$

First of all, we will prove that Newton's method:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \Rightarrow x_{n+1} - x_n = -\frac{f(x_n)}{f'(x_n)} \quad (*)$$

produces, in our case, an increasing sequence  $(x_0, x_1, \dots)$ . To do that, we will use the following property of a convex function:

P: A differentiable function of one variable is convex on an interval if and only if the function lies above all of its tangents:

$$f(x) \geq f(y) + f'(y)(x-y)$$

for all  $x$  and  $y$  in the interval.

For this to work, we will consider the domain of  $f$  to be  $\mathbb{R}$ , but  $f(x) = f(a)$  for  $x \in (-\infty, a)$  and  $f(x) = f(b)$  for  $x \in (b, +\infty)$ , so that  $f$  is still convex on  $\mathbb{R}$ , it is differentiable and strictly decreasing (on  $[a, b]$ ).

The problem comes if we get to a value  $x_n \notin [a, b]$  and that would imply  $f'(x_n) = 0$ , so we would need to stop iterating Newton's method. Therefore, we will prove by induction on  $n$  that  $x_n \in [a, b]$ , and later that  $x_n \in [x_0, x^*]$ , where  $x^*$  is the (unique) root of  $f$ :

P(0):  $x_0 \in [a, b]$ , obviously

$f(a)f(b) < 0$  and  $f$  is strictly dec.

Now, we'll assume that  $x_n \in [a, b]$  and we'll prove that  $x_{n+1} \in [a, b]$

$$x_n, x_{n+1} \in \mathbb{R} \stackrel{P}{\Rightarrow} f(x_{n+1}) \geq f(x_n) + f'(x_n)(x_{n+1} - x_n) \stackrel{*}{=} f(x_n) - f'(x_n) \frac{f(x_n)}{f'(x_n)} = 0$$

So,  $f(x_{n+1}) \geq 0 \Rightarrow x_{n+1} \in (-\infty, x^*]$

From P(n), we deduce that  $f(x_n) \geq 0$ , too, so

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \stackrel{\substack{\geq 0 \\ f'(x_n) < 0}}{\geq} x_n \Rightarrow x_{n+1} \in [x_n, +\infty) \quad \left| \begin{array}{l} \uparrow \\ \Rightarrow x_{n+1} \in [x_n, x^*] \text{ and by following} \end{array} \right.$$

the induction we get  $x_{n+1} \in [x_0, x^*] \subset [a, b] \Rightarrow x_{n+1} \in [a, b]$

So, we obtained that:

•  $(x_n)_{n \geq 0}$  is an increasing sequence | (Weinsthaas)

•  $(x_n)_{n \geq 0}$  is bounded above by  $x^*$  |  $\Rightarrow (x_n)_{n \geq 0}$  is convergent to  $l \in [x_0, x^*]$ .

By going back to Newton's iterative step and by applying  $\lim_{n \rightarrow \infty}$  to both terms, we get:

$$\lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} x_n - \frac{f(x_n)}{f'(x_n)} \stackrel{(f \text{ is continuous and } f' \neq 0 \text{ as } l \in [a, b])}{\Rightarrow} l = l - \frac{f(l)}{f'(l)} \Rightarrow \frac{f(l)}{f'(l)} = 0 \Rightarrow f(l) = 0 \Rightarrow l = x^*$$

Therefore,  $(x_n)_{n \geq 0}$  is increasing and converges to  $x^*$ .

Now, we consider

$f(x) = e^{\lambda(x-1)} - x$ , with  $\lambda > 1 \Rightarrow f$  is convex (difference of two convex functions)

and we want the lower root of  $f$ .

$$(b) \frac{df}{dx} = \lambda e^{\lambda(x-1)} - 1$$

$$\frac{d^2f}{dx^2} = \lambda^2 e^{\lambda(x-1)}$$

$$\frac{df}{dx}(x^+) = 0 \Rightarrow \lambda e^{\lambda(x^+-1)} - 1 = 0 \Rightarrow e^{\lambda(x^+-1)} = \frac{1}{\lambda} \Rightarrow x^+-1 = -\frac{\ln \lambda}{\lambda} \Rightarrow x^+ = 1 - \frac{\ln \lambda}{\lambda}$$

We have  $f(0) = e^{-\lambda} > 0$

$$f(x^+) = e^{1-\frac{\ln \lambda}{\lambda}-1} - 1 + \frac{\ln \lambda}{\lambda} = e^{-\ln \lambda} - 1 + \frac{\ln \lambda}{\lambda} = \frac{1}{\lambda} - 1 + \frac{\ln \lambda}{\lambda} = \frac{1+\ln \lambda - 1}{\lambda} < 0$$

because for  $x > 1$ , the function  $g(x) = 1 + \ln x - x$  has  $g'(x) = \frac{1}{x} - 1 < 0 \Rightarrow g$  is strictly decreasing, so  $g(x) < g(1) = 0$ , so for  $\lambda > 1$  we have  $1 + \ln \lambda - 1 < 0 \Rightarrow f(x^+) < 0$

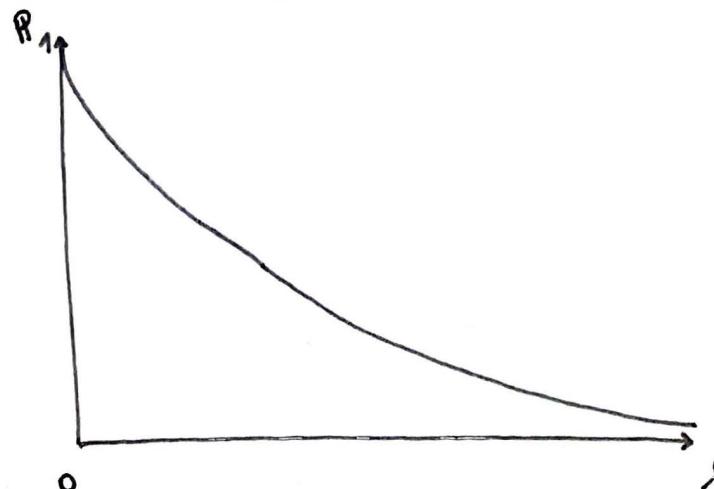
So,  $f(0)f(x^+) < 0$

$f$  is continuous and strictly decreasing  $\left| \begin{array}{l} \Rightarrow (\exists!) x^* \in (0, x^+) \text{ such that } f(x^*) = 0 \\ \text{exists and it's unique} \end{array} \right.$

The conditions from (a) apply here too, as we have:

$f: [0, x^+] \rightarrow \mathbb{R}$ : convex, differentiable, strictly decreasing (on this interval, easy to prove), with  $f(0) \geq 0$  and  $f(x^+) < 0$ .

(c) To find the extinction probability, we find the root of  $f$  starting from  $(0, x^+)$ , we first need to reduce the interval by using "interval bisection" until we get to an interval which is small enough to ensure quadratic convergence. (same method as in 3.5)



As  $\lambda \rightarrow \infty$ ,  $P(\text{extinction}) \rightarrow 0$ .

3.7.  $x^*$  is a root of  $f(x)=0$ , with  $\frac{df}{dx}(x^*) \neq 0$  and  $\frac{d^2f}{dx^2}(x^*) \neq 0$

$$e_n = x_n - x^*$$

(a) We want to first show that

$$e_{n+1} = e_n e_{n-1} \left( \frac{\frac{f(x_n)}{e_n} - \frac{f(x_{n-1})}{e_{n-1}}}{\frac{e_n - e_{n-1}}{e_n - e_{n-1}}} \right) \left( \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \right) \quad (\text{since } e_n - e_{n-1} = x_n - x^* - x_{n-1} + x^* = x_n - x_{n-1})$$

$$e_{n+1} = e_n e_{n-1} \cdot \frac{e_{n-1} f(x_n) - e_n f(x_{n-1})}{e_n f(x_n) - f(x_{n-1})}$$

$$e_{n+1} = \frac{e_{n-1} f(x_n) - e_n f(x_{n-1})}{f(x_n) - f(x_{n-1})} = \frac{(x_{n-1} - x^*) f(x_n) - (x_n - x^*) f(x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$x_{n+1} - x^* = \frac{x_{n-1} f(x_n) - x_n f(x_{n-1}) - x^* (f(x_n) - f(x_{n-1}))}{f(x_n) - f(x_{n-1})}$$

$$x_{n+1} = \frac{x_{n-1} f(x_n) - x_n f(x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$\text{Secant method iteration: } x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})} = \frac{x_{n-1} f(x_n) - x_n f(x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$x_n f(x_n) - x_n f(x_{n-1}) - x_n f(x_n) + x_{n-1} f(x_n) = x_{n-1} f(x_n) - x_n f(x_{n-1}) \quad \text{YES!}$$

Now, we'll show that

$$e_n e_{n-1} \left( \frac{\frac{f(x_n)}{e_n} - \frac{f(x_{n-1})}{e_{n-1}}}{\frac{e_n - e_{n-1}}{e_n - e_{n-1}}} \right) \left( \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \right) = \frac{1}{2} e_n e_{n-1} \frac{d^2 f}{dx^2}(\beta) / \frac{df}{dx}(\alpha), \text{ for } \alpha, \beta \in I,$$

where  $I$  contains  $x_{n-1}, x_n$  and  $x^*$

From Lemma 5.6 we have

$$\frac{f(x_n)}{e_n} - \frac{f(x_{n-1})}{e_{n-1}} = \frac{f(x^* + e_n)}{e_n} - \frac{f(x^* + e_{n-1})}{e_{n-1}} = \frac{(e_n - e_{n-1})}{2} \frac{d^2 f}{dx^2}(\beta), \text{ with } \beta \in I$$

From Taylor's Theorem, we have

$$f(x_n) = f(x_{n-1}) + (x_n - x_{n-1}) \frac{df}{dx}(\alpha), \text{ with } \alpha \in I \Rightarrow \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} = \frac{df}{dx}(\alpha)$$

From these two, we get:

$$e_{n+1} \left( \frac{\frac{e_n e_{n-1}}{2} \frac{d^2 f}{dx^2}(\beta)}{e_n e_{n-1}} \right) \left( \frac{1}{\frac{df}{dx}(x)} \right) = \frac{1}{2} e_n e_{n-1} \frac{d^2 f}{dx^2}(\beta) / \frac{df}{dx}(\alpha), \text{ with } \alpha, \beta \in I.$$

Therefore, we proved that

$$e_{n+1} = e_n e_{n-1} \left( \frac{\frac{f(x_n) - f(x_{n-1})}{e_n - e_{n-1}}}{\frac{f(x_n) - f(x_{n-1})}{e_n - e_{n-1}}} \right) \left( \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \right) = \frac{1}{2} e_n e_{n-1} \frac{\frac{d^2 f}{dx^2}(\beta)}{\frac{df}{dx}(\alpha)}, \text{ for } \alpha, \beta \in I, \text{ where}$$

$I$  is an interval that contains  $x_{n-1}, x_n$  and  $x^*$ .

(b) If  $x_0$  and  $x_1$  lie in  $(x^* - c, x^* + c)$  and let

$$A(c) = \frac{\max_{i \in I} \left| \frac{d^2 f}{dx^2}(i) \right|}{\min_{j \in I} \left| \frac{df}{dx}(j) \right|} \quad \text{with} \quad \frac{c A(c)}{2} < 1$$

$$\text{We have } \frac{\frac{d^2 f}{dx^2}(\beta)}{\frac{df}{dx}(\alpha)} \leq A(c) \quad (\forall \alpha, \beta \in (x^* - c, x^* + c) \text{ that also contains } x_{n-1}, x_n)$$

$$e_{n+1} = \frac{1}{2} e_n e_{n-1} \frac{\frac{d^2 f}{dx^2}(\beta)}{\frac{df}{dx}(\alpha)} \leq \frac{1}{2} e_n e_{n-1} A(c) \quad (*)$$

If we have  $x_0, x_1 \in (x^* - c, x^* + c)$  and  $\frac{c A(c)}{2} < 1$ , we'll prove by induction that

$$|e_n| \leq p^n (|e_0| + |e_1|) \text{ where } p = \frac{c A(c)}{2} < 1$$

~~$$\text{Base case: } |e_0| \leq p^0 |e_0| = |e_0| \quad (\text{True}) ; |e_1| \leq$$~~

~~$$\text{Inductive step: iH} = |e_n| \leq p^n |e_0|$$~~

$$|e_{n+1}| \leq \frac{1}{2} |e_n| |e_{n-1}| A(c) \Rightarrow (\text{as } |e_n| \leq |e_{n-1}| \leq \dots \leq |e_2| \leq |e_1|, |e_0| < c)$$

$$\Rightarrow |e_{n+1}| \leq \frac{1}{2} c A(c) |e_n| \leq p |e_n| \stackrel{iH}{\leq} p^{n+1} |e_0|$$

$$x_0, x_1 \in (x^* - c, x^* + c)$$

If we have  $x_0, x_1 \in (x^*-c, x^*+c)$  and  $\frac{cA(c)}{2} < 1$ , we'll prove by strong induction that  $|e_n| \leq p^{\lfloor \frac{n}{2} \rfloor} \max\{|e_0|, |e_1|\}$ , where  $p = \frac{cA(c)}{2} < 1$

Base cases:

$$P(0): |e_0| \leq p^0 \max\{|e_0|, |e_1|\} \quad (\text{YES})$$

$$P(1): |e_1| \leq p^0 \max\{|e_0|, |e_1|\} \quad (\text{YES})$$

Inductive step:

We assume that  $P(n-1)$  and  $P(n)$  are true and we'll show that  $P(n+1)$  is true (IH)

$$P(n-1): |e_{n-1}| \leq p^{\lfloor \frac{n-1}{2} \rfloor} \max\{|e_0|, |e_1|\} \stackrel{p < 1}{\downarrow} \leq \max\{|e_0|, |e_1|\} \stackrel{x_0, x_1 \in (x^*-c, x^*+c) \Rightarrow |e_0|, |e_1| < c}{\downarrow} \leq c$$

$$P(n): |e_n| \leq p^{\lfloor \frac{n}{2} \rfloor} \max\{|e_0|, |e_1|\} \leq \max\{|e_0|, |e_1|\} \quad (**)$$

I  $m=2k, k \in \mathbb{Z} \Rightarrow$

$$\Rightarrow |e_{n+1}| = |e_{2k+1}| \stackrel{(*)}{\leq} \frac{A(c)}{2} |e_{2k}| |e_{2k-1}| \stackrel{(**)}{\leq} \frac{A(c)}{2} c |e_{2k-1}| = p |e_{2k-1}| \stackrel{(IH)}{\leq} p \cdot p^{k-1} \max\{|e_0|, |e_1|\} = p^k \max\{|e_0|, |e_1|\} = p^{\lfloor \frac{n+1}{2} \rfloor} \max\{|e_0|, |e_1|\}$$

II  $m=2k+1, k \in \mathbb{Z} \Rightarrow$

$$\Rightarrow |e_{n+1}| = |e_{2k+2}| \stackrel{(*)}{\leq} \frac{A(c)}{2} |e_{2k+1}| |e_{2k}| \stackrel{(**)}{\leq} \frac{A(c)}{2} c |e_{2k}| = p |e_{2k}| \stackrel{(IH)}{\leq} p \cdot p^k \max\{|e_0|, |e_1|\} = p^{\lfloor \frac{n+1}{2} \rfloor} \max\{|e_0|, |e_1|\}$$

So, we showed that  $|e_n| \leq p^{\lfloor \frac{n}{2} \rfloor} \max\{|e_0|, |e_1|\} \Rightarrow e_n \xrightarrow[(p < 1)]{n \rightarrow \infty} 0 \Rightarrow$  the Secant method converges!

~~(c)  $\frac{|e_{n+1}|}{|e_n| |e_{n-1}|} = \frac{|e_{n+1}|}{|e_n| |e_{n-1}|} \leq \frac{A(c)}{2} \stackrel{c \rightarrow 0}{\rightarrow} \frac{1}{2} \frac{\left| \frac{d^2 f}{dx^2}(x^*) \right|}{\left| \frac{df}{dx}(x^*) \right|} = K > 0 \Rightarrow \frac{|e_{n+1}|}{|e_n| |e_{n-1}|} \rightarrow C, \text{ for some } C < K$~~  (the terms are always positive, so  $\frac{|e_{n+1}|}{|e_n| |e_{n-1}|}$  gets bounded by  $K$ , therefore the limit lies)

The fact that it converges to  $x^*$  comes from  $e_n \xrightarrow{n \rightarrow \infty} 0 \Rightarrow x_n - x^* \xrightarrow{n \rightarrow \infty} 0 \Rightarrow x_n \xrightarrow{n \rightarrow \infty} x^*$ .

~~(c)  $\frac{|e_{n+1}|}{|e_n| |e_{n-1}|} = \frac{|e_{n+1}|}{|e_n| |e_{n-1}|} \leq \frac{A(c)}{2} \stackrel{c \rightarrow 0}{\rightarrow} \frac{1}{2} \frac{\left| \frac{d^2 f}{dx^2}(x^*) \right|}{\left| \frac{df}{dx}(x^*) \right|} = C > 0 \Rightarrow$~~

$$\Rightarrow \frac{|e_{n+1}|}{|e_n| |e_{n-1}|} \rightarrow C, \text{ for some } C > 0.$$

(d) Now, let's find the order of convergence of the Secant method:

We know that the Secant method converges, so we can say that:

$$\frac{|e_{n+1}|}{|e_n|^q} = a_m \quad (\forall m \geq 0, a_m \in \mathbb{R} \text{ constant as } n \rightarrow \infty)$$

$$\text{Then, } |e_{n+1}| = a_m |e_n|^q = a_m (a_{m-1} |e_{m-1}|^q)^q = a_m a_{m-1}^q |e_{m-1}|^{q^2}$$

This leads to:

$$\frac{|e_{n+1}|}{|e_n| |e_{n-1}|} = \frac{a_m a_{m-1}^q |e_{m-1}|^{q^2}}{a_{m-1} |e_{m-1}|^q |e_{m-1}|} = a_m a_{m-1}^{q-1} |e_{m-1}|^{q^2-q-1} \xrightarrow{n \rightarrow \infty} C \Rightarrow$$

$$\Rightarrow |e_{m-1}|^{q^2-q-1} \xrightarrow{n \rightarrow \infty} \frac{C}{a_m a_{m-1}^{q-1}} \Rightarrow q^2 - q - 1 = 0 \Rightarrow q_{1/2} = \frac{1 \pm \sqrt{5}}{2} \left| \begin{array}{l} \text{constant} > 0 \\ q \geq 1 \end{array} \right. \Rightarrow$$

$$\Rightarrow q = \frac{1 + \sqrt{5}}{2} \Rightarrow q \approx 1.618, \text{ the order of convergence for the Secant method.}$$

**3.4.**  $x_{k+2} = \frac{9}{4} x_{k+1} - \frac{1}{2} x_k$

(i) First, we solve the auxiliary equation:

$$\begin{aligned} 1^2 &= \frac{9}{4} 1 - \frac{1}{2} & \Delta &= \frac{49}{16} \\ 1^2 - \frac{9}{4} 1 + \frac{1}{2} &= 0 \Rightarrow \lambda_{1/2} = \frac{\frac{9}{4} \pm \frac{7}{4}}{2} \Rightarrow \lambda_1 = 2; \lambda_2 = \frac{1}{4} \end{aligned}$$

Then, our solutions are:

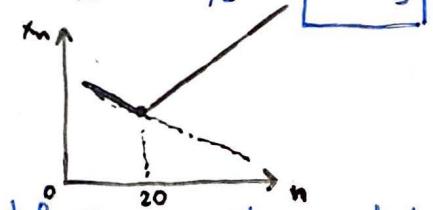
$$x_k = A_1 \lambda_1^k + A_2 \lambda_2^k = A_1 \cdot 2^k + A_2 \cdot \left(\frac{1}{4}\right)^k, \text{ for } A_1 \text{ and } A_2 \text{ to be determined.}$$

By using  $x_1 = \frac{1}{3}$  and  $x_2 = \frac{1}{12}$ , we get

$$\begin{cases} A_1 \cdot 2 + A_2 \cdot \frac{1}{4} = \frac{1}{3} \\ A_1 \cdot 4 + A_2 \cdot \frac{1}{16} = \frac{1}{12} \end{cases} \Rightarrow \begin{cases} 4A_1 + \frac{1}{2}A_2 = \frac{2}{3} \\ 4A_1 + \frac{1}{16}A_2 = \frac{1}{12} \end{cases} \Rightarrow \left(\frac{1}{2} - \frac{1}{16}\right)A_2 = \frac{2}{3} - \frac{1}{12} \Rightarrow \frac{7}{16}A_2 = \frac{7}{12} \Rightarrow A_2 = \frac{4}{3}$$

$$2A_1 = \frac{1}{3} - \frac{1}{4}A_2 = \frac{1}{3} - \frac{1}{3} = 0 \Rightarrow A_1 = 0$$

$$\text{Therefore, we get } x_k = \frac{4}{3} \left(\frac{1}{4}\right)^k, k \in \mathbb{N}_+$$



(ii) By implementing the recurrence relation, we observe that until  $x_{20}$  we obtain predictable answers and from there up to  $x_{100}$ , which will be  $\approx 1.88 \cdot 10^{-12}$ , the value of  $x_k$  increases, although our formula would make it decrease. That is because at step  $k=20$ , when we tried to find  $x_{21}$ , the relative error was  $\frac{x_{k+1}-x_k}{x_k}$ , so, because of roundoff errors,  $x_{k+1}$  increased instead. 12.