# Accuracy vs Efficiency of Numerical Methods
## How to program a Calculator

Jake Darby

**Abstract**

This document will discuss and analyse various numerical methods for computing functions commonly found on calculators. The aim of this paper is to, for each set of functions, compare and contrast several algorithms in regards to their effiency and accuracy.

# 1 Introduction

For many thousands of years all calculations that a Human might want performing had to be done by hand. For simple calculations such as addition, subtraction and multiplication this was not such an issue, but as society evolved we wanted to know the answer to increasingly hard questions. The greeks saught to find a value for $\pi$, and ended up with the bounds that $\frac{223}{71} < \pi < \frac{22}{7}$, which while sufficient for their needs is not sufficient for ours in the modern era.

Nowadays we have computers and calculators to calculate functions such as square roots or the trigonometric functions. These devices however must be told how to correctly calculate these values, those methods used will be looked at here. Specifically, for most purposes it does not matter if some small error is present in the results, as long as the value found is close enough to the actual value.

## 1.1 Code and Computers used

During this project I will be discussing the implementation of various algorithms. I will be implementing these algorithms in the C programming language, using the C11 standard.

I chose the C programming language to implement my algorithms in because, once it compiles to binary machine code, the programs produced tend to be very efficient. This is partly due to the low-level of C programming, having relatively close control over direct CPU actions; however this does come at the cost of losing higher functionality that many other programming languages offer. A second reason for the efficiency is due to C's long history, originally being developed int 1969-1970, which has lead to several very efficient compilers.

I will be implementing most programs using C's built in primitive types, typically `int`, `unsigned int` and `double`. On a computer an `int` is an integer that can represent both positive and negative bits using twos compliment, this gives an `int` using $n$ bits a minimum value of $-2^n$ and a maximum value of $2^n - 1$. Typically a computer will store an `int` as 32 bits, though some computers may use more or less bits. An `unsigned int` is very similar to an `int`, but does not represent negative values, and thus an `unsigned int` of $n$ bits has a minimum value of $0$ and a maximum value of $2^{n+1} - 1$.

If an integer of a specific number of bits is needed then the header `stdint.h` may be used which defines `int_N` and `uint_N` which repsectively represent `int` of N bits and `unsigned int` of N bits; The typical values of N are 8, 16, 32 and 64.

In C a `double` is a floating point representation of a real value, that typically follows the IEEE 754 standard for double-precision binary floating points. This standard has:

- 1 bit for the sign of the number, $s$

- 11 bits for the exponent, $e$

- 52 bits for the significand, $b = b_0 b_1 b_2 \ldots b_{51}$

- A value that is understood to be:

$$(-1)^s \left(1 + \sum_{i=1}^{52} b_{52-i} 2^{-i}\right) \times 2^{e-1023}$$

This gives a `double` value a precision of around 15-17 significant decimal digits. While this is good for most applications, there are applications where we may want even more precision than this. To solve this I will be implementing certain algorithms using the GNU Multiple Precision Arithmetic Library (reffered to as GMP) as well as GNU MPFR Library(reffered to as GMP), wich was built upon GMP to correct and optimise the original. These libraries allow the use of arbitrary precision real values, given enough memory space, as well as integers longer than C's standard integer types can hold.

The downside is that the increased precision does increase computation time for these calculations.

I will be compiling and testing all of my code on a benchmark machine running a light version of Ubuntu ¡VERSION¿, using the GNU C Compiler. The specifications of the machine, that may impact perrformance are:

- An Intel i5-4690K processor running at 4GHz. This processor uses a 64 bit architecture.

- 8Gb of DDR3 RAM

- A modern chipset on the motherboard

TODO: Expand the introduction and make it more eloquent

# 2 General Definitions and Theorems

This section will discuss those definitions and theorems that will be applicable and refferenced later in the document.

## 2.1 Methods

In this document we will look at various functions, such as root functions, trigonometric functions, among others. Despite the variety of functions being analysed there are several methods that are useful for more than one function, we will discuss these methods below.

### 2.1.1 Newton-Raphson Method

The Newton-Raphson Method is named after Sir Isaac Newton and Joseph Raphson. It is a method that takes a continuously differentiable function $f$ and it's derivative $f'$, as well as an initial guess $x_0$, to create successively more accurate solutions to $x$ where $f(x) = 0$.

The specific definition of the Newton-Raphson method that I will be using in this document is below:

**Definition 2.1.1.1.** Given $f \in \mathcal{C}(\mathbb{R})$, $f'$ being the derivative of $f$, and $x_0 \in \mathbb{R}$; then we define:

$$x_{n+1} := x_n - \frac{f(x)}{f'(x)} \forall n \in \mathbb{N}$$

The hope is that, if we start with a good initial guess, that the method will converge to some $x \in \mathbb{R} : f(x) = 0$.

We derive the above method in the following way. If we have an approximation $x_n \in \mathbb{R} :$ $f(x_n) \approx 0$ we consider the tangent to $f(x)$ at $x_n$, which is given by $y = f'(x_n)(x-x_n)+f(x_n)$. If we set $y = 0$ and solve for $x$ we get that

$$x = x_n - \frac{f(x)}{f'(x)}$$

The iterative formula follows by letting $x_{n+1} := x$ in the above.

TODO: Re-arrange into a better order
TODO: Refferences
TODO: Possibly include discussions of pros and cons

### 2.1.2 Taylor Series Expansion

The Taylor Series formulation was created by Brook Taylor in 1715, based off of the work of Scottish mathematician James Gregory. The Taylor Series describes a method of representing a given function by a polynomial, where any finite number of terms will give an approximation to the function itself.

**Definition 2.1.2.1.** Given $f : \mathbb{R} \to \mathbb{R}$ which is infinitely differentiable at $a \in \mathbb{R}$, we define the Taylor Series of $f$ at $a$ to be:

$$\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!}(x-a)^n$$

We can then, from our definition of a Taylor Series, define a polynomial that will approximate our function $f$ at $x \in \mathcal{I} \subset \mathbb{R}$

**Definition 2.1.2.2.** Given $f : \mathbb{R} \to \mathbb{R}$ which has a Taylor Series of $\sum_{n=0}^{\infty} c_n x^n$, we define the Taylor Polynomial of degree $N \in \mathbb{N}$ to be

$$p_N(x) := \sum_{n=0}^{N} c_n x^n = c_0 + c_1 x + c_2 x^2 + \cdots + c_N x^N$$

Some examples of simple Taylor Series are:

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n \qquad \qquad \forall x \in (-1, 1)$$

TODO: Eloquate this section
TODO: Discuss when a Taylor Series is not equal to it's function
TODO: Describe intervals for which Taylor Series converge
TODO: Create more examples of Taylor Series that are not used later in the document
TODO: Possibly prove Taylor's Theorem

### 2.1.3 CORDIC

TODO: Describe the CORDIC Algorithm in general
TODO: Let other sections define specific changes to CORDIC

## 2.2 Errors

Errors are very useful in this document for discussing convergence and measuring the accuracy of a particular method. There are two distinct types of error that we are interested in:

**Definition 2.2.1.** If we have a value $v$ and it's approximation $\tilde{v}$, then the absolute error is

$$\epsilon := |v - \tilde{v}|$$

The absolute error will be useful in guaranteeing a certain level of accuracy that a given implementation of a method will give. Uses of absolute error in the document will use $\epsilon$ as their absolute error variable.

**Definition 2.2.2.** If we have a value $v$ and it's approximation $\tilde{v}$, then the relative error is

$$\eta := \frac{\epsilon}{|v|} = \left| 1 - \frac{\tilde{v}}{v} \right|$$

The relative error will also be useful in guaranteeing a certain level of accuracy that a given implementation of a method will give, particularly when the absolute error varies with the size of the input. Uses of relative error in the document will use $\eta$ as their absolute error variable.

As the previous two errors are hard or impossible to accurately estimate, in a practical manner, then we want an error estimate that we can calculate in an algorithm. Thus we define the iteration error, as the absolute difference of consecutive iterations.

**Definition 2.2.3.** If we have the sequence $(x_n)_{n \in \mathbb{N}}$, then the iteration error is defined as:

$$\delta_n := |x_n - x_{n-1}|$$

For our practical purposes we can note that it is almost impossible to calculate $\epsilon_n$ exactly, while $\delta_n$ is easy to calculate and tends to be a good enough approximation of $\epsilon_n$; particularly if the convergence is rapid.

## 2.3 Convergence

definition
In this section we consider what it means for a sequence to converge to a limit value, and some results that will be useful in future chapters.

**Definition 2.3.1.** A sequence $(x_n \in \mathbb{R} : n \in \mathbb{N})$ converges to $x$ uniformly if

$$\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : \epsilon_n := |x - x_n| < \tau \forall n \in [N, \infty) \cap \mathbb{Z}$$

*Remark* 2.3.1.1. We will typically use the notation that $\lim_{n \to \infty} |x_n - x| = 0$, to denote that $(x_n : n \in \mathbb{N})$ converges to $x$.

**Theorem 2.3.1.** $(x_n \in \mathbb{R} : n \in \mathbb{N})$ *converges to $x$ uniformly if and only if*

$$\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : |x_n - x_m| < \tau \forall m, n \in [N, \infty) \cap \mathbb{Z}$$

*Proof.* For $\implies$ :
Suppose that $(x_n : n \in \mathbb{N})$ converges to $x$ uniformly. Then $\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : |x_n - x| < \tau \forall n \in [N, \infty) \cap \mathbb{Z}$.
Thus suppose $N \in \mathbb{N}$ is such that $|x_n - x| < \frac{\tau}{2} \forall n \in [N, \infty) \cap \mathbb{Z}$.
Then if $n, m \geq N$ we see that

$$|x_n - x_m| \leq |x_n - x| + |x_m - x| \leq \tau$$

For $\impliedby$:
Omitted for brevity. $\qquad \square$

**Definition 2.3.2.** If $(x_n \in \mathbb{R} : n \in \mathbb{N})$ is a sequence that converges to $x$, then it is said to converge:

- Linearly if $\lambda \in \mathbb{R}^+$ and
$$\lim_{n \to \infty} \frac{|x_{n+1} - x|}{|x_n - x|} = \lambda$$

- Quadratically if $\lambda \in \mathbb{R}^+$ and
$$\lim_{n \to \infty} \frac{|x_{n+1} - x|}{|x_n - x|^2} = \lambda$$

- Order $\alpha \in \mathbb{R}_0^+$ if $\lambda \in \mathbb{R}^+$ and

$$\lim_{n \to \infty} \frac{|x_{n+1} - x|}{|x_n - x|^\alpha} = \lambda$$

For the following proof we will have $\epsilon_n := |x^* - x_n|$.

**Theorem 2.3.2.** *Let $f$ be a twice differentiable function, $x^*$ be a solution to $f(x) = 0$ and $(x_n : n \in \mathbb{N})$ be a sequence produced by the Newton-Raphson Method from some initial point $x_0$. If the following are satisfied, then $(x_n : n \in \mathbb{N}_0)$ converges quadratically to $x^*$:*

$\mathbf{NR_1:}$ $f'(x) \neq 0 \forall x \in I := [x^* - r, x^* + r]$, where $r \in [|x^* - x_0|, \infty)$

$\mathbf{NR_2:}$ $f''(x)$ is continuous $\forall x \in I$

$\mathbf{NR_3:}$ $M|\epsilon_0| < 1$ where $M := \sup\left\{\left|\frac{f''(x)}{f'(x)}\right| : x \in I\right\}$

*Proof.* By Taylor's Theorem with Lagrange Remainders we have that

$$0 = f(x^*) = f(x_n) + (x^* - x_n)f'(x_n) + \frac{1}{2}(x^* - x_n)^2 f''(y_n)$$

where $0 < |x^* - y_n| < |x^* - x_n|$.

Then we get the following derivation:

$$f(x_n) + (x^* - x_n)f'(x_n) = -\frac{1}{2}(x^* - x_n)^2 f''(y_n)$$

$$\implies \left(\frac{f(x_n)}{f'(x_n)} - x_n\right) + x^* = -\frac{1}{2}\frac{f''(y_n)}{f'(x_n)}(x^* - x_n)^2 \qquad \text{as NR}_3 \implies f'(x_n) \neq 0$$

$$\implies x^* - x_{n+1} = -\frac{1}{2}\frac{f''(y_n)}{f'(x_n)}(x^* - x_n)^2$$

$$\implies \epsilon_{n+1} = \frac{1}{2}\left|\frac{f''(y_n)}{f'(x_n)}\right|\epsilon_n^2 \qquad \text{by taking absolute values}$$

As $\text{NR}_2$ holds then $M$ exists and is positive, and thereforewe have:

$$\epsilon_n \leq M\epsilon_{n-1}^2 \leq M^{2^n - 1}\epsilon_0^{2^n}$$

We now aim to show that we have convergence, i.e. $\lim_{n \to \infty} x_n = x^*$; to do this it suffices to show that $\lim_{n \to \infty} \epsilon_n = 0$.

Consider the sequence $(z_n := M^{2^n - 1}\epsilon_0^{2^n} : n \in \mathbb{N}_0)$. We know that $0 \leq \epsilon_n \leq z_n \forall n \in \mathbb{N}_0$, so it then follows that if $\lim_{n \to \infty} z_n = 0$, then $\lim_{n \to \infty} \epsilon_n = 0$ by the Squeeze Theorem **??**.

Now as $M\epsilon_0 < 1$ by $\text{NR}_3$, then we see that:

$$\begin{aligned} \lim_{n \to \infty} z_n &= \lim_{n \to \infty} (M\epsilon_0)^{2^n - 1}\epsilon_0 \\ &= \epsilon_0 \lim_{n \to \infty} (M\epsilon_0)^{2^n - 1} \\ &= \epsilon_0 \cdot 0 \qquad \text{because } \lim_{n \to \infty} r^{m_n} = 0 \text{where } |r| < 1, m_n \geq n \forall n \in \mathbb{N} \\ &= 0 \end{aligned}$$

Now to show that this sequence converges quadratically we see that $\epsilon_{n+1} = \frac{1}{2} \left| \frac{f''(y_n)}{f'(x_n)} \right| \epsilon_n^2$, and therefore $\frac{\epsilon_{n+1}}{\epsilon_n^2} = \frac{1}{2} \left| \frac{f''(y_n)}{f'(x_n)} \right|$.

Because $|x^* - y_n| < |x^* - x_n|$ and $\lim_{n\to\infty} x_n = x^*$, then it follows that $\lim_{n\to\infty} y_n = x^*$. Therefore we see that

$$\lim_{n\to\infty} \frac{\epsilon_{n+1}}{\epsilon_n} = \frac{1}{2} \left| \frac{f''(x^*)}{f'(x^*)} \right| \in \mathbb{R}^+$$

Hence as the above limit exists and is positive then the sequence is quadratically convergent.

$\square$

## 2.4 Efficiency and Accuracy Metrics

# 3 Division and Multiplication

## 3.1 Introduction

Though the idea of Division and Multiplication can seem fairly simple, particularly from an abstract pure mathematical point of view, these calculations can be computationally difficult. This section will show a few algorithms designed to calculated these values, and discuss their implementation and efficiency.

## 3.2 Multiplication

### 3.2.1 Basic Methods

### 3.2.2 Advanced Methods

### 3.2.3 Analysis

## 3.3 Division

### 3.3.1 Basic Methods

### 3.3.2 Advanced Methods

### 3.3.3 Analysis

TODO: Section to be filled out, no work currently done on this section

# 4 Trigonometric Functions

This section will focus on trigonometric functions, which are commonly used cyclic functions. These functions have been studied for hundreds of years, and can be challenging to calculate. We will discuss several methods of calculating them below before comparing methods.
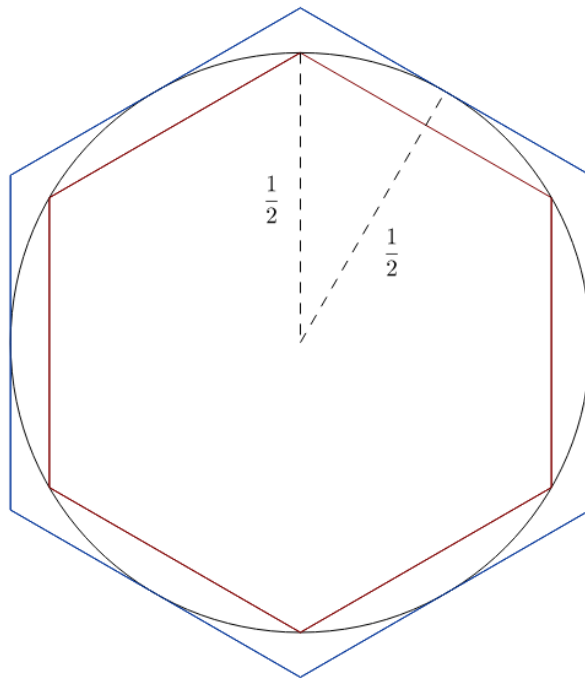TODO: Extend and Eloquate introduction

## 4.1 Calculating $\pi$

Several of the methods in this section require that we already know the value of $\pi$, for example when we are applying several trig identities. Here we will briefly discuss several methods for calculating the value of $\pi$, so that we may use this value in later subsections.

The first method to consider is the method used by ancient mathematicians, such as the Greeks and Chinese. We know that if the radius of the circle is $\frac{1}{2}$, then the circumference of the circle is $\pi$, and the value is between the perimeters of the inner and outer polygon perimeters. The internal perimeter is $p_n = n \sin(\frac{\pi}{n})$ and the external perimeter is $P_n = n \tan(\frac{\pi}{n})$.

Figure 4.1.1: Ancient method of calculating $\pi$



As we know the values of $\tan(\frac{\pi}{6})$ and $\sin(\frac{\pi}{6})$, then we can calculate $P_6$ and $p_6$. It has be shown that $P_{2n} = \frac{2p_n P_n}{p_n + P_n}$ and $p_{2n} = \sqrt{p_n P_{2n}}$, which allows us to create an iterative method to approximate $\pi$, by taking the mid-point of the successive polygon perimeters.

Other common historical methods for approximating $\pi$ are to use infinite series. One such method uses the series expansion of $\tan^{-1}$, which is discussed in detail below, where $\tan^{-1}(1) = \frac{\pi}{4}$. This gives the following approximation using $N$ terms:

$$\pi = 4 \sum_{n=0}^{N} \frac{(-1)^n}{2n+1} = \sum_{n=0}^{N} \frac{8}{(4n+1)(4n+3)} \tag{4.1.1}$$

10

This sequence converges very slowly, with sublinear convergence, to the correct value. More modern methods have typically revolved around finding more rapidly converging infinite series, examples include Ramanujan's series:

$$\frac{1}{\pi} = \frac{2\sqrt{2}}{9801} \sum_{n=0}^{\infty} \frac{(4n)!(1103 + 26390n)}{(k!)^n 396^{4n}} \qquad (4.1.2)$$

or the Chudnovsky algorithm:

$$\frac{1}{\pi} = 12 \sum_{n=0}^{\infty} \frac{(-1)^n (6n)!(13591409 + 545140134n)}{(3n)!(n!)^3 640320^{3n+\frac{3}{2}}} \qquad (4.1.3)$$
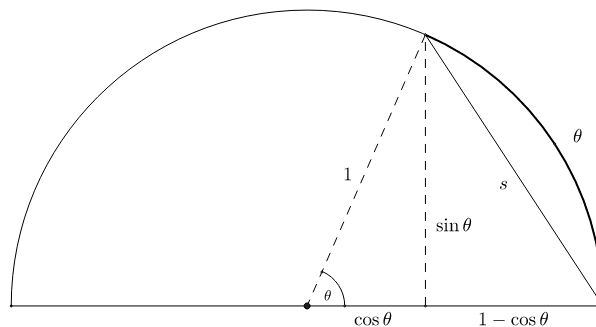
This final series is extremely rapidly convergant to the value of $\frac{1}{\pi}$, for example just the first term gives $\pi$ accurate to 13 decimal places while we can get $\pi$ accurate to 1000 decimal places with summing just 71 terms. Compared to Equation **??** which takes the summation of 500 terms to acheive the same 1000 digits of accuracy.

To get large degrees of accuracy for $\pi$ is extremely computer intensive and using the `mpfr` requires the number of bits of precision and number of terms to be set. This makes calculating $\pi$ to a large number of decimal places, for example 1000000, computationally infeasible on a regular home computer. Therefore for our purposes we will use the precalculated value of $\pi$ to 1000000 decimal places as listed on `http://www.exploratorium.edu/pi/pi_archive/Pi10-6.html`

## 4.2   Geometric Method

The first method I will be discussing is a method based on geometric properties that are derived on a circle, and we will start by considering values of $\cos$ in the range $[0, \frac{\pi}{2}]$. To do this we will consider the following figure of the unit circle:
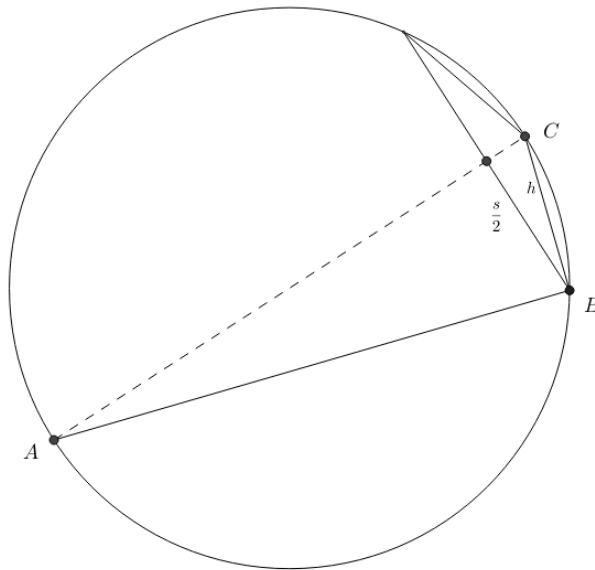
Figure 4.2.1: Diagram showing angles to be dealt with



Here theta will be given in radians, and we can note that the labelled arc has length $\theta$ due the formula for the circumference of a circle. By using the following derivation we can find a formula for $\theta$ in terms of $s$:

$$\begin{aligned}
s^2 &= \sin^2 \theta + (1 - \cos \theta)^2 \\
&= (\sin^2 \theta + \cos^2 \theta) + 1 - 2\cos \theta \\
&= 2 - 2\cos \theta \qquad \qquad \text{By using } \sin^2 \theta + \cos^2 \theta = 1 \\
\cos \theta &= 1 - \frac{s^2}{2}
\end{aligned}$$

We will now consider a second diagram which will allow us to calculate an approximate value of $s$.

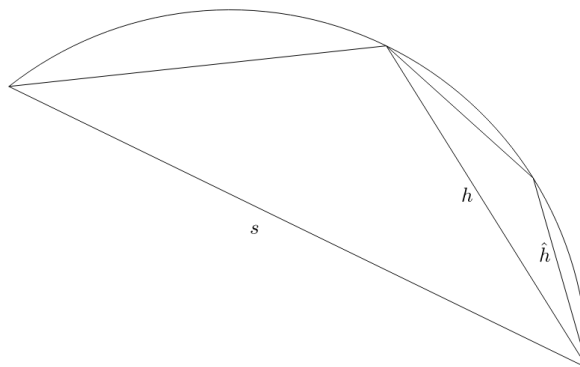Figure 4.2.2: Diagram detailing how to calculate $s$



We will first note that by an elementary geometry result we can know that the angle $ABC$ is a right-angle; also we can consider that $h$ is an approximation of $\frac{\theta}{2}$, which will become relevant later. Now because $AC$ is a diameter of our circle then it's length is 2 and thus, by utilising Pythagarus' Theorem, we get that the length of $AB$ is $\sqrt{AC^2 - BC^2} = \sqrt{4 - h^2}$.

From here we consider the area of triangle $ABC$, which can be calculated as $\frac{1}{2} \cdot h \cdot \sqrt{4 - h^2}$ and as $\frac{1}{2} \cdot 2 \cdot \frac{s}{2}$; by equating these two, squaring both sides and re-arranging we get that $s^2 = h^2(4 - h^2)$. Now we have the basis for a method that will allow us to calculate $\cos \theta$.

To complete our method we will consider introducing a new line that is to $h$ what $h$ is to $s$ as shown in the diagram below:

Figure 4.2.3: Detailing the recursive steps



It is easy to see that if we repeat the steps above we get that $h^2 = \hat{h}^2(4 - \hat{h}^2)$, and it also follows that $\hat{h} \approx \frac{\theta}{4}$. Using this we can take an initial guess of $h_0 := \frac{\theta}{2^k}$, for some $k \in \mathbb{N}$, and

then calculate $h_{n+1}^2 = h_n^2(4 - h_n^2)$ where $n \in [0, k] \cap \mathbb{Z}$; finally we calculate $\cos \theta = 1 - \frac{h_k^2}{2}$, giving the following algorithm:

Algorithm 4.2.1: Geometric calculation of $\cos$

```
1   geometric_cos (θ ∈ [0, π/2], k ∈ ℕ)
2       h₀ := θ/2ᵏ
3       n := 0
4       while n < K:
5           h²ₙ₊₁ := h²ₙ · (4 - h²ₙ)
6           n ↦ n + 1
7       return 1 - h²ₖ/2
```
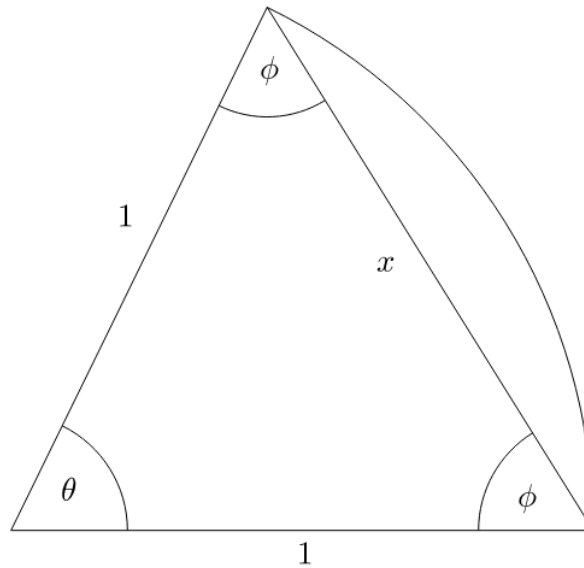
Now we can use the above pseudocode to calculate any trigonometric function value by using various trigonometric identities. First we suppose $\theta \in \mathbb{R}$, then we can repeatedly apply the identity $\cos \theta = \cos(\theta \pm 2\pi)$ to either add or subtrack $2\pi$ until we have a value $\theta' \in [0, 2pi)$. Once we have this value we can utilise the following assignment to calculate $\cos \theta$:

$$\cos \theta = \begin{cases} \cos \theta' & : \quad \theta' \in [0, \frac{\pi}{2}] \\ -\cos(\pi - \theta') & : \quad \theta' \in [\frac{\pi}{2}, \pi] \\ -\cos(\theta' - \pi) & : \quad \theta' \in [\pi, \frac{3\pi}{2}] \\ \cos(2\pi - \theta') & : \quad \theta' \in [\frac{3\pi}{2}, 2\pi) \end{cases}$$

Using Algorithm 4.2.1 we can also easily calculate both $\sin \theta$ and $\tan \theta$, by further use of trigonometric identities. In particular we note that $\sin \theta = \cos(\theta - \frac{\pi}{2}$ and $\tan \theta = \frac{\sin \theta}{\cos \theta}$. Hence we can now calculate the trigonometric function value of any angle.

We now wish to analyse the error of our approximation for $\cos$, as the other methods have errors that are derivative of the error for approximating $\cos$. Now Figure 4.2.4 shows an arc of a circle which creates chord $x$, with this we will be able to calculate the exact length of the chord and thus work on the error of our approximations.

Figure 4.2.4: Diagram to find actual arc approximation

To start we will note that $\phi = \frac{\pi - \theta}{2} = \frac{\pi}{2} - \frac{\theta}{2}$, and then by using the Sine Law we get

$$\frac{x}{\sin\theta} = \frac{1}{\sin\phi} \implies x = \frac{\sin\theta}{\sin\phi}$$

Now we can recall the double angle formula for $\sin$, which gives $\sin\theta = 2\sin\frac{\theta}{2}\cos\frac{\theta}{2}$, and also $\sin\phi = \cos\frac{\theta}{2}$. This allows us to see that

$$x = \frac{2\sin\frac{\theta}{2}\cos\frac{\theta}{2}}{\cos\frac{\theta}{2}} = 2\sin\frac{\theta}{2}$$

Therefore we see that $h_n$ is approximating the chord length associated with angle $\theta 2^{n-k}$, and thus $\epsilon_n = |h_n - 2\sin(\theta 2^{n-k-1})|$. Now as $h_0 = \theta 2^{-k} \approx 2\sin(\theta 2^{-k-1})$ then if follows that $\exists\phi$ such that $h_0 = 2\sin(\phi 2 - k - 1)$, from this we can see that $\phi = 2^{k+1}\sin^{-1}(\theta 2^{-k-1})$. We will uses these facts to prove a couple of propositions.

**Proposition 4.2.1.** $h_n = 2\sin(\phi 2^{n-k-1}) \forall n \in [0, k] \cap \mathbb{Z}$ *where* $\phi := 2^{k+1}\sin^{-1}(\theta 2^{-k-1})$.

*Proof.* Proceed by induction on $n \in [0, k] \cap \mathbb{Z}$.

**H**$(n)$: $h_n = 2\sin(\phi 2^{n-k-1})$

**H**$(0)$:
$$2\sin(\phi 2^{-k-1}) = 2\sin(\sin^{-1}(\phi 2^{-k-1}))$$
$$= \theta 2^{-k}$$
$$= h_0 \qquad\qquad \text{by definition of } h_0$$

**H**$(n) \implies$ **H**$(n+1)$:

$$h_{n+1} = h_n\sqrt{4 - h_n^2}$$
$$= 2\sin(\phi 2^{n-k-1})\sqrt{4 - 4\sin^2(\phi 2^{n-k-1})} \qquad\qquad \text{by H}(n)$$
$$= 4\sin(\phi 2^{n-k-1})\cos(\phi 2^{n-k-1})$$
$$= 2\sin(\phi 2^{n-k}) \qquad\qquad \text{by the use of double angle formulas}$$

$\square$

**Proposition 4.2.2.** $h_n > 2\sin(\theta 2^{n-k-1}) \forall n \in [0, k] \cap \mathbb{Z}$

*Proof.* We start by considering the expansion of the exact value of $h_n$.

$$h_n = 2\sin(\phi 2^{n-k-1})$$
$$= 2\sin(2^{n-k-1}(2^{k+1}\sin^{-1}(\theta 2^{-k-1})))$$
$$= 2\sin(2^n\sin^{-1}(\theta 2^{-k-1}))$$
$$= 2\sin(\theta 2^{n-k-1} + \tfrac{1}{6}\theta^3 2^{n-3k-3} + \mathcal{O}(2^{-5k})) \quad \text{Detailed in section ??}$$

Now as we know that $n \leq k$, then it follows that $\theta 2^{n-k-1} \leq \frac{1}{2}\theta$.

Also as $\theta \leq \frac{\pi}{2}$ we know that $\theta 2^{n-k-1} \leq \frac{\pi}{4}$.

We can also show that $\frac{1}{6}\theta^3 2^{n-3k-3} + \mathcal{O}(2^{-5k}) \leq \frac{\pi}{4}$, though the proof is ommited here for brevity; therefore we see that $\phi 2^{n-k-1} \leq \frac{\pi}{2}$, and obviously that $\phi 2^{n-k-1} > \theta 2^{n-k-1}$.

Hence, as $\sin$ is an increasing function in the range $[0, \frac{\pi}{2}]$, we conclude that

$$h_n = 2\sin(\phi 2^{n-k-1}) > 2\sin(\theta 2^{n-k-1})$$

. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\square$

With these two propositions we can now consider the error of our approximation of $\cos$. First we will prove the following proposition regarding the error of the approximation of $s$:

**Proposition 4.2.3.** *If $\epsilon_n := |h_n - 2\sin(\theta 2^{n-k-1})| \forall n \in [0,k] \cap \mathbb{Z}$, then $\epsilon_k < 2^k \epsilon_0)$.*

*Proof.* $\epsilon_n = h_n - 2\sin(\theta 2^{n-k-1})$ as $h_n > 2\sin(\theta 2^{n-k-1})$ by Proposition 4.2.2.

Now we see that:

$$\epsilon_{n+1} = h_{n+1} - 2\sin(\theta 2^{n-k})$$
$$= h_n\sqrt{4 - h_n^2} - 4\sin(\theta 2^{n-k-1})\cos(\theta 2^{n-k-1})$$

If we consider the equation $\alpha\beta - \gamma\delta = (\alpha - \gamma) + \alpha(\beta - 1) - \gamma(\delta - 1)$ and apply it to our current formula we get:

$$\epsilon_{n+1} = (h_n - 2\sin(\theta 2^{n-k-1})) + h_n(\sqrt{4 - h_n^2} - 1) - 2\sin(\theta 2^{n-k-1})(2\cos(\theta 2^{n-k-1}) - 1)$$
$$= \epsilon_n + h_n(\sqrt{4 - h_n^2} - 1) - 2\sin(\theta 2^{n-k-1})(2\cos(\theta 2^{n-k-1}) - 1)$$
$$= 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2) - 2\sin(\theta 2^{n-k-1})(2\cos(\theta 2^{n-k-1}) - 2)$$
$$= 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2) + 2\sin(\theta 2^{n-k-1})(2 - 2\cos(\theta 2^{n-k-1}))$$
$$< 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2\cos(\theta 2^{n-k-1}))$$
$$< 2\epsilon_n + h_n(\sqrt{4 - 4\sin^2(\theta 2^{n-k-1})} - 2\cos(\theta 2^{n-k-1}))$$
$$= 2\epsilon_n + h_n(2\cos(\theta 2^{n-k-1}) - 2\cos(\theta 2^{n-k-1}))$$
$$= 2\epsilon_n$$

The inequalities in the above derivation arrise from the fact that $h_n > 2\sin(\theta 2^{n-k-1})$ by Proposition 4.2.2.

Hence as we now know that $\epsilon_{n+1} < 2\epsilon_n$, we then see that $\epsilon_n < 2^n \epsilon_0$. Therefore we prove our statement that

$$\epsilon_k < 2^k \epsilon_0$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\square$

Obviously $\epsilon_k = |h_k - s|$, and we can now use this to find the error of our final answer. First we will start by letting $\mathcal{C} := 1 - \frac{1}{2}h_k^2$ and note that analytically $cos\theta = 1 - \frac{1}{2}s^2$. Therefore we will now consider $\epsilon_{\mathcal{C}} = |\mathcal{C} - \cos(\theta)|$:

$$\epsilon_{\mathcal{C}} = |1 - \frac{h_k^2}{2} - 1 + \frac{s^2}{2}|$$

$$= \frac{1}{2}|h_k^2 - s^2|$$

$$= \frac{1}{2}|h_k h_k - 2\sin(\frac{\theta}{2})2\sin(\frac{\theta}{2})|$$

$$= \frac{1}{2}(h_k h_k - 2\sin(\frac{\theta}{2})2\sin(\frac{\theta}{2})) \qquad \text{as } 2\sin(\frac{\theta}{2}) < h_k$$

$$= \frac{1}{2}(2\epsilon_k + h_k(h_k - 2) - 2\sin(\frac{\theta}{2})(2\sin(\frac{\theta}{2}) - 2))$$

$$< \frac{1}{2}(2\epsilon_k + h_k(h_k - 2\sin(\frac{\theta}{2})))$$

$$= \frac{1}{2}(2 + h_k)\epsilon_k$$

$$= \frac{1}{2}(2 + 2\sin(\frac{\phi}{2}))\epsilon_k$$

$$= (1 + \sin(\frac{\phi}{2}))\epsilon_k$$

$$\leq 2\epsilon_k$$

As $\epsilon_{\mathcal{C}} \leq 2\epsilon_k$, then by Proposition 4.2.3 we see that $\epsilon_{\mathcal{C}} < 2^{k+1}\epsilon_0$. Now to consider $\epsilon_0$ we first observe that $\epsilon_0 = \theta 2^{-k} - 2\sin\theta 2^{-k-1}$, and therefore we can conclude that:

$$\epsilon_{\mathcal{C}} < 2\theta - 2^{k+2}\sin(\theta 2^{-k-1})$$

If we then wish to calculate $\cos\theta$ accurate to $N$ decimal places then we are looking to find $k \in \mathbb{N}$ such that

$$2\theta - 2^{k+2}\sin(\theta 2^{-k-1}) < 10^{-N} \implies 2^{k+2}\sin(\theta 2^{-k-1}) > 2\theta - 10^{-N}$$

For an example of the above in action we will be taking $\theta = 0.5$. The table below shows the minimum $k \in \mathbb{N}$ to guarantee $N$ digits of accuracy in the result:

| $N$ | $k$ |
|---|---|
| 5 | 6 |
| 10 | 14 |
| 50 | 80 |
| 100 | 163 |
| 1000 | 1658 |

As can be seen the value of $k$ required to acheive $N$ digits of accuracy increases roughly linearly when $\theta = 0.5$. Testing for other values of $\theta$ reveals them to have similar required values for $k$, at least within the same order of each other.

Another consideration for Algorithm 4.2.1 is that we could "run it in reverse" to attain an algorithm for the inverse cosine function. To start take line 7 which is $\mathcal{C} = 1 - \frac{1}{2}h_k^2$, which can be re-arranged to give $h_k^2 = 2 - 2\mathcal{C}$, where we know $\mathcal{C}$ as our initial value.

Line 5 is a little more difficult, but by re-arranging we see that $h_n^4 - 4h_n^2 + h_{n+1}^2 = 0$, which can be solved via the quadratic formula to give $h_n^2 = 2 \pm \sqrt{4 - h_{n+1}^2}$. Now we can make

the observation that if $x \in \mathbb{R}_0^+$, then $\cos^{-1}(-x) = \pi - \cos^{-1}(x)$ and so we can restric our algorithm to only consider $x \in [0,1]$. With this we know that $\theta \in [0, \frac{\pi}{2}]$, and thus $h_k \leq \sqrt{2}$. Therefore as $h_{n+1} > h_n \forall n \in [0, k-1] \cap \mathbb{Z}$ we see that $h_n^2 \leq 2 \forall n \in [0, k] \cap \mathbb{Z}$. This allows us to ascertain that to reverse Line 5 we perform $h_n^2 = 2 - \sqrt{4 - h_{n+1}^2}$.

Finally line 2 is reversed by returning the value $2^k h_0$; therefore we get the following algorithm for $\cos^{-1}(x)$ where $x \in [0,1]$:

Algorithm 4.2.2: Geometric calculation of $\cos^{-1}$

```
1   geometric_aCos (x ∈ [0,1], k ∈ ℕ)
2       h_k := 2 - 2x
3       n := k - 1
4       while  n ≥ 0:
5           h_n^2 := 2 - √(4 - h_{n+1}^2)
6           n ↦ n - 1
7       return  2^k h_0
```

Similar to the regular trigonometric functions we can use trigonometric identities to calculate the inverse trigonometric functions from $\cos^{-1}$. To start we recall that $\cos^{-1}(-x) = -\cos(x)$ where $x \in [0,1]$, then we can use the identities that $\sin^{-1}(x) = \frac{\pi}{2} - \cos^{-1}(x)$ and $\tan^{-1}(x) = \sin^{-1}(\frac{x}{\sqrt{x^2+1}})$.

If we suppose that all operations in the method are accurately computed then Algorithm 4.2.2 is a computation with high accuracy. This is because there is no initial guess, such as in Algorithm 4.2.1, and so the only introduction of error is assuming that $2^k h_0 \approx \theta$. However as we discuss in detail in Section **??**, calculating square roots is not a simple task and thus will introduce error to the method in general; therefore the accuracy of the method is roughly as accurate as our method of calculating square roots.

## 4.3  Taylor Series

If we consider our definition of a McClaurin Series from Section **??**, we can use this to approximate our Trigonometric Functions. Consider first $\cos\theta$, for which we know that $\frac{d}{d\theta}\cos\theta = -\sin\theta$; it then follows that $\frac{d^2}{d\theta^2}\cos\theta = -\cos\theta$, $\frac{d^3}{d\theta^3}\cos\theta = \sin\theta$ and $\frac{d^4}{d\theta^4}\cos\theta = \cos\theta$.

If we let $f(x) = \cos x$ and use the known values $\cos(0) = 1$ and $\sin(0) = 0$, then we see that:

$$f^{(n)}(0) = \begin{cases} 1 & : \ 4 \mid n \\ 0 & : \ 4 \mid n-1 \\ -1 & : \ 4 \mid n-2 \\ 0 & : \ 4 \mid n-3 \end{cases}$$

By simplifying this by ommitting the $0$ coefficient terms we get the following series:

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} x^{2n} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \cdots \tag{4.3.1}$$

By using similar working we can get that the series associated with $\sin x$):

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \cdots \tag{4.3.2}$$

Before we go any further we need to consider when Equations 4.3.1 and 4.3.2 converge to their respective functions. To do this we will use the ratio test for series as defined in **??**, using Equation 4.3.1 we see that

$$
\begin{aligned}
L_{\mathcal{C}} &= \lim_{n \to \infty} \left| \frac{a_{n+1}}{a_n} \right| \\
&= \lim_{n \to \infty} \left| \frac{\frac{(-1)^{n+1}}{(2n+2)!} x^{2n+2}}{\frac{(-1)^n}{(2n)!} x^{2n}} \right| \\
&= \frac{(2n)!}{(2n+2)!} |x|^2 \\
&= \frac{1}{(2n+2)(2n+1)} |x|^2
\end{aligned}
$$

Now it is easy to see that, $L_{\mathcal{C}} = 0$ for all values of $x$ as the fractional component decreases as $n$ increases and $|x|^2$ is a constant. Therefore we can conclude that Equation 4.3.1 converges to $\cos(x)$ for all values of $x$. We can use a very similar deduction to show that Equation 4.3.2 converges to $\sin(x)$ for all values of $x$.

The above means that $\cos$ and $\sin$ can be approximated using Taylor Polynomials, in particular for a given $N \in \mathbb{N}$:

$$\cos x \approx \sum_{n=0}^{N} \frac{(-1)^n}{(2n)!} x^{2n} \quad \text{and} \quad \sin x \approx \sum_{n=0}^{N} \frac{(-1)^n}{(2n+1)!} x^{2n+1}$$

This allows us to create the following two methods for computing $\cos x$ and $\sin x$:

Algorithm 4.3.1: Taylor computation of $\cos$ and $\sin$

```
1    taylor_cos (x ∈ ℝ, N ∈ ℕ)
2        C := 0
3        n := 0
4        while n < N:
5            C ↦ C + (−1)ⁿ · 1/(2n)! x^2n
6            n ↦ n + 1
7        return C
8
9    taylor_sin (x ∈ ℝ, N ∈ ℕ)
10       S := 0
11       n := 0
12       while n < N:
13           S ↦ S + (−1)ⁿ · 1/(2n+1)! x^2n+1
14           n ↦ n + 1
15       return S
```

As these two methods are obviously very similar and the fact that $\sin(x) = \cos(x - \frac{\pi}{2})$, we will continue by examining only the taylor method for approximating $\cos$. We will assume that

any calculations for $\sin$ are transformed into a problem of finding a $\cos$ value.

It should be noted that this $\cos$ algorithnm is particularly inefficient to calculate on a computer implementation; this is primarily due to the way in which the update of $\mathcal{C}$ is calculated each loop.

In each loop we are calculating $x^{2n}$, which has a naieve complexity of $\mathcal{O}(2n)$. However what we are actually calculating $x^{2(n-1)} \cdot x^2$ and thus if we store the values of $x^{2(n-1)}$ and $x^2$, the complexity of this step drops to $\mathcal{O}(1)$. Similarly we are also calculating $\frac{1}{(2n)!}$ in each loop which, by the same logic, is $\frac{1}{2(n-1)!} \cdot \frac{1}{(2n)(2n-1)}$, and we can use the same storage and update method as for $x^{2n}$.

As another step towards optimizing the algorithm we can start with an initial value of $\mathcal{C} = 1$, and then perform two updates of $\mathcal{C}$ each loop until we reach or surpass $N$. This saves calculating $(-1)^n$ each loop, by explicitly performing two different calculations. Implementing all of the above gives us the following two updated methods:

Algorithm 4.3.2: Taylor computation of $\cos$ optimised

```
1    taylor_cos (x ∈ ℝ, N ∈ ℕ)
2        𝒞 := 1
3        x₂ := x²
4        a := 1
5        b := 1
6        n := 1
7        while  n < N :
8            a ↦ a · 1/((2n-1)(2n))
9            b ↦ b · x₂
10           𝒞 ↦ 𝒞 - a · b
11           a ↦ a · 1/((2n+1)(2n+2))
12           b ↦ b · x₂
13           𝒞 ↦ 𝒞 + a · b
14           n ↦ n + 2
15       return  𝒞
```

As the next term of the polynomial is known definitively then we can see that it is very easy to calculate the error of our approximation. We see that

$$\begin{aligned}
\epsilon_N &= |\cos(x) - \text{taylor\_cos(x, N)}| \\
&= \mathcal{O}(|x|^{N'+1}) \qquad\qquad \text{where } N' \text{ is the smallest} \\
&\qquad\qquad\qquad\qquad\qquad \text{odd integer such that } N' \geq N \\
&\leq \frac{1}{(2(N'+1))!}|x|^{N'+1} \\
&\leq \frac{1}{(2(N+1))!}|x|^{N+1}
\end{aligned}$$

If we place bounds on the value of $\cos$ calculated as in Section 4.2, then we know that $|x| \leq \frac{\pi}{2}$, and thus we get the following bound for the error of our approximation:

$$\epsilon_N \le \frac{\pi^{N'+1}}{2^{N'+1}(2(N'+1))!}$$

Thus if we find $N \in \mathbb{N}$ such that $\frac{\pi^{N+1}}{2^{N+1}(2(N+1)!)} < \tau \in \mathbb{R}+$ then we know that $\epsilon_N < \tau$. If we consider $\tau = 10^k$, then we can find $N \in \mathbb{N}$ such that our approximation is accurate to $k$ decimal places. Below is a table which details some values of $k$ and the corresponding minimum $N$ to guarantee $k$ decimal places of accuracy:

| $k$ | $N$ |
| --- | --- |
| 5 | 4 |
| 10 | 7 |
| 50 | 21 |
| 100 | 36 |
| 1000 | 233 |

Now for $\tan x$ we can either calculate both $\sin x$ and $\cos x$ using taylor_cos(x,N) and divide the resulting value, or we can calculate $\tan x$ directly using a Taylor expansion.

In calculating the McClaurin series for $\tan x$ we start by letting $\tan x = \sum_{n=0}^{\infty} a_n x^n$, and then noting that as $\tan x$ is an odd series then it's McClaurin series only contains non-zero coefficients for odd powers of $x$; therefore we get that $\tan x = \sum_{n=0}^{\infty} a_{2n+1} x^{2n+1} = a_1 x + a_3 x^3 + a_5 x^5 + \cdots$.

Next we consider that $\frac{d}{dx} \tan x = 1 + \tan^2 x$, and knowing the McClaurin series form of $\tan x$ we get the following:

$$\sum_{n=0}^{\infty} (2n+1)a_{2n+1} x^{2n} = 1 + (\sum_{n=0}^{\infty} a_{2n+1} x^{2n+1})^2$$
$$= 1 + a_1^2 x^2 + (2a_1 a_3)x^4 + (2a_1 a_5 + a_3^2)x^6 + \cdots$$

Considering the co-efficients of powers on the right hand side of the above equation we see that $2a_1 a_3 = a_1 a_3 + a_{3a1} = a_1 a_{4-1} + a_3 a_{4-3}$ and $2a_1 a_5 + a_3^2 = a_1 a_5 + a_3 a_3 + a_5 a_1 = a_1 a_{6-1} + a_3 a_{6-3} + a_5 a_{6-5}$. This indicates that our general form for the co-efficient of $2n$ on the right hand side is $\sum_{k=1}^{n} a_{2k-1} a_{2n-2k+1}$, and thus returning to our equation we get

$$a_1 + \sum_{n=1}^{\infty} (2n+1)a_{2n+1} x^{2n} = 1 + \sum_{n=1}^{\infty} (\sum_{k=1}^{n} a_{2k-1} a_{2n-2k+1})x^{2n}$$

Using this we conclude that $a_1 = 1$ and $a_{2n+1} = \frac{1}{2n+1} \sum_{k=1}^{n} a_{2k-1} a_{2n-2k+1} \forall n \in \mathbb{N}$. We can note immediately that the calculation of any previous co-efficients will provide no help in calculating later co-efficients and so the entire sum must be calculated each loop, while also storing each co-efficient already calcualted.

This means that the complexity to calculate co-efficient $a_{2n+1}$ is $\mathcal{O}(n)$ and will be the $n^{\text{th}}$ such calculation, making the complexity of calculating $n$ co-efficients to be $\mathcal{O}(n^2)$. Comparing this to the taylor_cos method we see that to calculate up to $n$ co-efficients of both $\cos$ and $\sin$ has complexity $\mathcal{O}(n)$. Therfore it is more efficient to calculate $\tan$ by calculating both $\cos$ and $\sin$ using Algorithm **??**, and performing division than directly using Taylor Polynomial

approximation.

We would also like to be able to calculate the inverse trigonometric functions using this method, which means we need to find our McClaurin series of the inverse trigonometric functions. The simplest of these is $\tan^{-1}$, where we start by recalling that $\frac{d}{dx}\tan^{-1}x = \frac{1}{1+x^2}$ and then by intergrating both sides we get:

$$\tan^{-1}x = \int \frac{1}{1+x^2}dx$$

$$= \int (1-(-x^2))^{-}1dx$$

$$= \int \sum_{n=0}^{\infty}(-x^2)^n dx \qquad \text{by Equation ??}$$

$$= \int \sum_{n=0}^{\infty}(-1)^n x^{2n} dx$$

$$= c + \sum_{n=0}^{\infty}\frac{(-1)^n}{2n+1}x^{2n+1}$$

As $\tan^{-1}(0) = 0$ then we see that $c = 0$ and thus gives us the following formula for $\tan^{-1}$:

$$\tan^{-1}x = \sum_{n=0}^{\infty}\frac{(-1)^n}{2n+1}x^{2n+1}$$

Now due to the restrictions from Equation **??** the above is only valid for $x \in [-1,1]$, but we know that the domain of $\tan^{-1}$ is $x \in \mathbb{R}$. To fix this we will first recognise that $\tan^{-1}(-x) = -\tan^{-1}(x)$, so we can restric our problem to $x \in \mathbb{R}_0^+$. Now if we take the double angle formula for $\tan$:

$$\tan(\alpha + \beta) = \frac{\tan(\alpha) + \tan(\beta)}{1 - \tan(\alpha)\tan(\beta)}$$

By substituting $\alpha = \tan^{-1}(x)$ and $\beta = \tan^{-1}(x)$ into the above then we get

$$\tan^{-1}(x) + \tan^{-1}(y) = \tan^{-1}\left(\frac{x+y}{1-xy}\right)$$

Using this, suppose we are looking for $\tan^{-1}(z)$ where $z \in (1,\infty)$ and let $y = 1$, then $\tan^{-1}(y) = \frac{\pi}{4}$. We can then re-arrange the equation $z = \frac{x+1}{1-x}$ to get $x = \frac{z-1}{z+1}$; finally as $z > 1$, then $0 < x < 1$. This allows us to calculate:

$$\tan^{-1}(z) = \frac{\pi}{4} + \tan^{-1}\left(\frac{z-1}{z+1}\right)$$

In the above the calculated value is in the range $[0,1]$ and so it is valid to use a Taylor polynomial using our McClaurin series above. This gives the following method

Algorithm 4.3.3: Taylor Method for $\tan^{-1}$

| | |
|---|---|
| 1 | $\mathrm{taylor\_aTan}\,(x \in [0,1], N \in \mathbb{N})$ |
| 2 | $\quad \mathcal{T} := 0$ |
| 3 | $\quad x_2 := x^2$ |

```
4          y := x
5          n := 0
6          while  n < N :
7                  T ↦ T + (1/(2n+1))y
8                  y ↦ y · x_2
9                  T ↦ T − (1/(2n+2))y
10                 y ↦ y · x_2
11                 n ↦ n + 2
12         return  T
```

Similar to Algorithm **??** the error of Algorithm 4.3.3 is easy to calculate. We see that

$$\epsilon_N = |\tan^{-1}(x) - \text{taylor\_aTan}(x, N)|$$

$$\leq \frac{1}{2N+3}|x|^{2N+3}$$

$$\leq \frac{1}{2N+3} \qquad\qquad \text{as } x \leq 1$$

The next function we will consider is $\sin^{-1}$, which starts it's derivation in much the same way as $\tan^{-1}$. First we start by recalling that $\frac{d}{dx}\sin^{-1}(x) = (1-x^2)^{-\frac{1}{2}}$, then by taking integrals of both sides we get the following derivation:

$$\sin^{-1}(x) = \int (1-x^2)^{-\frac{1}{2}}dx$$

$$= \int \sum_{n=0}^{\infty} \binom{-\frac{1}{2}}{n}(-x^2)^n$$

$$= c + \sum_{n=0}^{\infty}(-1)^n \left(\prod_{k=1}^{n}\frac{-\frac{1}{2}-k+1}{k}\right)\frac{x^{2n+1}}{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{(-1)^n}{n!(2n+1)}\left(\prod_{k=1}^{n}\frac{1}{2}-k\right)x^{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{(-1)^{2n}}{n!(2n+1)}\left(\prod_{k=1}^{n}\frac{2k-1}{2}\right)x^{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{1}{n!(2n+1)2^n}\left(\prod_{k=1}^{n}2k-1\right)x^{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{1}{n!(2n+1)2^n}(1\times 3\times 5\times \cdots \times(2n-1))x^{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{1}{n!(2n+1)2^n}\times\frac{1\times 2\times 3\times\cdots\times(2n)}{2\times 4\times\cdots\times(2n)}x^{2n+1}$$

$$= c + \sum_{n=0}^{\infty}\frac{(2n)!}{(n!)^2(2n+1)4^n}x^{2n+1}$$

As $\sin^{-1}(0) = 0$ then we see that $c = 0$. Because the above is valid for $x \in (-1, 1)$, and we know the values of $\sin^{-1}(-1)$ and $\sin^{-1}(1)$, then we can have the following method for evaluating $\sin^{-1}$:

```
1    taylor_aSin (x ∈ [-1, 1], N ∈ ℕ)
2        if  x = 1:
3            return  π/2
4        if  x = -1:
5            return  -π/2
6        S := x
7        x₂ := x²
8        y := x
9        a := 1
10       b := 1
11       c := 1
12       n := 1
13       while  n < N:
14           a ↦ 2n · (2n - 1) · a
15           b ↦ n² · b
16           c ↦ 4 · c
17           y ↦ x₂ · y
18           S ↦ S + a/(b·c·(2n+1)) · y
19           n ↦ n + 1
20       return  S
```

The error for this method is similar to the $\tan^{-1}$ method, in that $\epsilon_N \leq \frac{(2(N+1))!}{((N+1)!)(2N+1)4^{N+1}}$. Finally we note that $\cos^{-1}(x) = \frac{\pi}{2} - \sin^{-1}(x)$, and thus can be calculated from a value calculated with Algorithm 4.3.4.

## 4.4 CORDIC

CORDIC is an algorithm that stands for COordinate Rotation DIgital Computer and can be used to calculate many functions, including Trigonometric Values. The CORDIC algorithm works by utilising Matrix Rotations of unit vectors. This algorithm is less accurate than some other methods but has the advantage of being able to be implemented for fixed point real numbers in efficient ways using only addition and bitshifting.

CORDIC works by taking an initial guess of $\mathbf{x}_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ which can be rotated through an anti-clockwise angle of $\gamma$ by the matrix

$$\begin{pmatrix} \cos\gamma & -\sin\gamma \\ \sin\gamma & \cos\gamma \end{pmatrix} = \frac{1}{\sqrt{1 + \tan\gamma^2}} \begin{pmatrix} 1 & -\tan\gamma \\ \tan\gamma & 1 \end{pmatrix}$$

By taking taking smaller and smaller values of $\gamma$ we can create an iterative process to find $\mathbf{x}_n$ which converges, for a given $\beta \in (-\frac{\pi}{2}, \frac{\pi}{2})$, to

$$\begin{pmatrix} \cos\beta \\ \sin\beta \end{pmatrix}$$

# 5 Root Functions

In this section of the document we will consider several methods for approximating root functions. For our purposes we are only going to consider roots of $N \in \mathbb{R}_0^+$, this is because if $N \in \mathbb{R}^-$ then it follows that $\sqrt{N} = i\sqrt{|N|}$.

## 5.1 Digit by Digit Method

The first method we will examine is an old method, that has been observed in Babylonian Mathematics over 2000 years ago, which is used to accurately generate the square root of numbers one digit at a time. This method differs from others discussed as it generates each digit of the root with perfect accuracy, one at a time, thus in a theoretical sense this algorithm is the most accurate of the methods we will view; we will see however that this method is slow.

Now suppose we are looking for $\sqrt{N}$, then we know that $\sqrt{N} = a_0 10^n + a_1 10^{n-1} + a_2 10^{n-2} + \ldots$ for some $n \in \mathbb{Z}$; it then follows that $N = (a_0 10^n + a_1 10^{n-1} + a_2 10^{n-1} + \ldots)^2$. By expanding the quadratic value we get that

$$N = a_0^2 10^{2n} + (20a_0 + a_1)a_1 10^{2n-2} + (20(a_0 10 + a_1) + a_2)a_2 10^{2n-4} + \cdots + (20\sum_{i=0}^{k-1} a_i 10^{k-i-1} + a_k)a_k 10^{2n-2k}$$

An observation should be made regarding the value of $n$ that we use for the theorem. We could of course try different values of $n$, in some structured procedure, that will find the largest $n$ such that $10^n \leq N$. However we can note that $log_{10}(\sqrt{N}) = \frac{1}{2}log_{10}(N)$, thus $10^{\frac{1}{2}log_{10}(N)} = \sqrt{N}$. Using this information, and the fact that $n \in \mathbb{Z}$, we can have $n := \lfloor \frac{1}{2}log_{10}(N) \rfloor$.

This allows us to get successive apporximations of $N$ where $N_0 = a_0^2 10^{2n}$, $N_1 = N_0 + (20a_0 + a_1)a_1 10^{2n-2}$, $N_2 = N_1 + (20(a_0 10 + a_1) + a_2)a_2 10^{2n-4}$. This will alllow us to create an algorithm that will give successive approximations of $sqrtN = a_0 10^n + a_1 10^{n-1} + \ldots$, more importantly each approximation will give us the exact next digit in the decimal representation of $\sqrt{N}$.

Thus we can have an iterative method to solve the problem, where at each stage we are trying to find the largest digit which satisfies the inequality $(20\sum_{i=0}^{k-1} a_i 10^{k-i-1} + a_k)a_k 10^{2n-2k} \leq N - N_{k-1}$. Thus we get the following pseudocode, which outputs two sequences, one indicating the digits before the decimal point and one afterwards. I will use set notation to indicate the sequences, but in this case order is important and repetition is allowed.

Algorithm 5.1.1: Exact Digit by Digits Square Root

```
1   exactRootDigits(N ∈ ℝ₀⁺, d ∈ ℕ):
2       Digitsₐ := ∅
3       Digit_b := ∅
4       k := 0
5       n := ⌊½log₁₀(N)⌋
6       while  k < d:
7           aₖ := max{t ∈ [0,9] ∩ ℤ : (20∑ᵢ₌₀ᵏ⁻¹ aᵢ10ᵏ⁻ⁱ⁻¹ + t) t10²ⁿ⁻²ᵏ ≤ N}
8           N ↦ N − (20∑ᵢ₌₀ᵏ⁻¹ aᵢ10ᵏ⁻ⁱ⁻¹ + aₖ) aₖ10²ⁿ⁻²ᵏ
9           if  n − k < 0:
```

```
10                          Digits_b ↦ Digits_b ∪ {a_k}
11                  else :
12                          Digits_a ↦ Digits_a ∪ {a_k}
13              k ↦ k + 1
14          if  Digits_a = ∅ :
15              Digits_a := {0}
16          if  Digits_b = ∅ :
17              Digits_b := {0}
18          return  (Digits_a, Digits_b)
```

This method has a computational complexity of $\mathcal{O}(d^2)$, as each loop requires the operations of summing $k$ elements, and the loop is repeated for $k = 0 \to d$. We will see that by considering some changes to the algorithm we can change the complexity class to be $\mathcal{O}(d)$.

First we will note that line 5 is not an issue, as if we only care about the first significant digit of $\frac{1}{2} log_{10}(N)$, then this is $\mathcal{O}(|log(N)|)$. This can be seen as if we start from $n = 0$ we can either count up or down until a we find $10^{2n}$ at most or at least N, respectively. This obviously takes at most $|log_{10}(N)|$ steps, giving us our stated complexity. We will also assume that $\mathcal{O}(|log(N)|) \le \mathcal{O}(d)$, as we have already seen that we can manipulate our input N to be within a reasonable range.

Second we note that on line 7 we calculate $\sum_{i=0}^{k-1} a_i 10^{k-i-1}$ for each value of $t$; we can reduce the complexity of this line by pre-calculating this value. However we can do even better if we consider that at step $k + 1$ we are calculating $\sum_{i=0}^{k} a_i 10^{k-i} = a_k + 10 \sum_{i=0}^{k-1} a_i 10^{k-i-1}$. Thus if we introduce $P_0 := 0$, and fore each k we calculate $P_{k+1} := 10P_k + a_k$, then we can reduce the complexity from $\mathcal{O}(k)$ to $\mathcal{O}(1)$.

This calculation of $P_k$, then carries over to reduce the complexity of line 8 to be $\mathcal{O}(1)$ instead of $\mathcal{O}(k)$. Combining this we can create the modified algorithm below:

Algorithm 5.1.2: Exact Digit by Digits Square Root version 2

```
1    exactRootDigits_v2 (N ∈ ℝ₀⁺, d ∈ ℕ):
2        Digits_a := ∅
3        Digits_b := ∅
4        k := 0
5        n := ⌊½log_{10}(N)⌋
6        P_0 := 0
7        while  k < d:
8            a_k := max {t ∈ [0,9] ∩ ℤ : (20P_k + t) t10^{2n-2k} ≤ N}
9            N ↦ N − (20P_k + a_k) a_k 10^{2n-2k}
10           P_{k+1} := 10P_k + a_k
11           if  n − k < 0:
12               Digits_b ↦ Digits_b ∪ {a_k}
13           else :
14               Digits_a ↦ Digits_a ∪ {a_k}
15           k ↦ k + 1
16       if  Digits_a = ∅:
17           Digits_a := {0}
18       if  Digits_b = ∅:
19           Digits_b := {0}
```

```
20        return (Digits_a, Digits_b)
```

This method is usefull, but can be difficult to implement as it requires high precision for the representation of the real value of $N$. In my implementation using C, I utilised the MPFR library to utilise high precision integers, but still encountered issues regarding loss of precision.

As an example the table below shows the number of digits of accuracy I was able to calculate for $\sqrt{2}$ using the above algorithm, compared to the number of bits of precision used in the calculations.

| Bits of Precision | Maximum Accuracy |
|---|---|
| 8 | 2 |
| 16 | 5 |
| 32 | 9 |
| 64 | 18 |
| 128 | 39 |
| 256 | 77 |
| 512 | 154 |
| 1024 | 308 |
| 2048 | 615 |
| 4096 | 1234 |
| 8192 | 2466 |

This data is highly structured and so we can hope to create a simple function that would allow us to calculate how much precision would be needed for a given number of digits of accuracy, at least for single digit inputs for $N$. We can see that the average ratio of Precision to Accuracy is 3.41259..., which ranges from 3.31928... to 4.0. From this we can draw a general trend that Digits of Accuracy $\approx 3.4 \times$ Bits of Precision; thus if we take the more generous assumption that Digits of Accuracy $4 \times$ Bits of Precision, we can use this to pre-determine the accuracy needed.

It should be noted that to ensure accuracy we should over-estimate the required precision, however if we overestimate the precision, then our calculations will be performed using unnecessarily large data structures and thus computation time will increase.

One particular use of this technique is to find an approximation of a squar root to it's integer part, calculated in base 2. This algorithm is of note as we will see that it has a computation time of $(1)$.

The algorithm uses the same basis as the base 10 version, for it's calculations, but due to the nature of being in binary several changes can be made for computational efficiency. To do this we will view the problem as follows: if we know some $r \in \mathbb{Z}_0^+$ which is our current approximation of our root, we are looking for some $e \in \mathbb{Z}_0^+$ such that $(r + e)^2 \le N$. Expanding this out we get $r^2 + 2re + e^2 \le N$, and if we keep track of $M = N - r^2$, we can test if $2re + e^2 \le M$.

Now we can consider our choice of $e$, the most practical method is to test successeive $e_m := 2^m$, where $m$ is descending starting with $m = \max m \in \mathbb{Z}_0^+ : 4^m \le N$. We can use an iterative

formula to build up the integer square root, where we start with $r = 0, M = N$ and have $rr + e_m$ whenever $2re_m + e_m^2 \leq M$, stopping when $m < 0$. This is then implemented as follows:

Algorithm 5.1.3: Integer Square Root Algorithm

```
1   integerSquareRoot(N ∈ ℤ₀⁺):
2       M := N
3       m := max m ∈ ℤ₀⁺ : 4ᵐ ≤ M
4       r := 0
5       while m ≥ 0:
6           if 2r(2ᵐ) + 4ᵐ ≤ M:
7               M ↦ M − 2r(2ᵐ) + 4ᵐ
8               r ↦ r + 2ᵐ
9           m ↦ m − 1
10      return r
```

If we now conisder an implementation of the above algorithm using an unsigned integer system with $K$ bits, where $2|K$. We will use `res` to represent $2re_m$, which means at the start of the algorihtm we will have `res = 0`; similarly we can use `bit` to represent $e_m^2$. As we know that $K$ bits are used and $2|K$, it then follows that the largest power of 4 less than the maximum representable value ($2^K - 1$ is $2^{K-2}$, which can be calculated as `bit = 1 << (K - 2)` using bitshift operations. Finally we will use `num` to represent $M$.

Now that we have discussed the setup we can consider how to implement some of the steps above. First to implement line 3 we can simply keep dividing `bit` by 4 while `bit > num`, which can be efficiently implemented as `bit >> 2` by using bitshifts in place of division by powers of 2. The same technique can be used in place of line 9, which leads us to re-evaluating our usage of line 5. As we are using bitshifting and a bitshift that would take a number past 0 instead results in 0, we also know that $2|K$ and so eventually we will reach `bit == 1`, which represents $m = 0$; therefore we can use `bit > 0` as our stopping criteria on line 5.

Line 6 is easy to convert, given our definitions of `res`, `bit` and {num, as is line 7. All that remains is to consider how to update `res`, which has two different ways of being updated depending on whether `res + bit <= num`. If it is false that `res + bit <= num`, then we wish for `res` to represent $2re_{m-1}$; this is easily acheived if we consider that $2re_{m-1} = \frac{1}{2}(2re_m)$, which prompts the update `res = res >> 1`. For the second case, when `res + bit <= num` is true, we want `res` to represent $2(r + e_m)e_{m-1}$; to implement this we consider the following derivation:

$$
\begin{aligned}
2(r + e_m)e_{m-1} &= \frac{1}{2} \cdot 2(r + e_m)e_m \\
&= \frac{1}{2} \cdot 2(re_m + e_m^2) \\
&= \frac{1}{2}(2re_m) + e_m^2
\end{aligned}
$$

Using this above derivation we see that we can calculate this as `res = (res >> 1) + bit`. Below is a simple implementaion of this in C using the unsigned 32 bit integer type `uint32_t`. A more commented and slightly modified version can be found in Appendix **??**, File **??**.

```
uint32_t int_sqrt(uint32_t num)
```

```
{
        uint32_t res = 0, bit = (1 << 30);

        while (bit > num)
                bit = bit >> 2;

        while (bit > 0)
        {
                if (res + bit <= num)
                {
                        num = num - (res + bit);
                        res = (res >> 1) + bit;
                }
                else
                        res = res >> 1;

                bit = bit >> 2;
        }

        return res;
}
```

We should consider the final step of the loop, when `bit == 1`. In this case when `res` is updated we have `res` represent either $2(r+e_0)e_{-1} = r+e_0$, or $2re_{-1} = r$; thus the algorithm exits with the correct value.

Now that the algorithm is correctly constructed using simple unsigned integer addition, subtraction and bitshifting (which we can assume all have computational time of $\mathcal{O}1$), we can look at the worst case complexity of the algorithm:

- The complexity of the set up of variables is constant time.

- The worst case complexity would be to to have `bit <= num` at the start.

- The loop would execute 16 times for our 32 bit integers, and contains a single operation which is $\mathcal{O}(1)$ complexity.

    - The worst case within the loop is to have `res + bit <= num` for each iteration.
    - Within the first `if` branch there are a constant 4 operations.
    - Each loop has an additional operation operation to update `bit`.
    - This makes 5 operations per loop, giving $\mathcal{O}(1)$ complexity within the loops.

Therefore we see that the algorithm has $\mathcal{O}(1)$ time complexity, and even has the same in storage complexity. In particular our 32 bit example requires 163 opertaions, including assignments, comparrisons and caluicuations. This means that the integer square root of any number up to 4294967295 can be calculated extremely quickly.

## 5.2  Bisection Method

The Bisection Method is a general method for approximating the zero, $\alpha$, of a function, $f$, on a bounded interval, $I := [a, b]$, where $f$ has the property $f(x)f(y) < 0 \forall (x, y) \in [a, \alpha) \times (\alpha, b]$; we may assume, without loss of generality, that $f(x) < 0 \forall x \in [a, \alpha]$.

The bisection method starts with initial bounds $a_0 = a, b_0 = b$, where the initial approximation for the root is $x_0 = \frac{1}{2}(a + b)$. We will consider pseudocode of the iteration process, that uses $b_n - a_a < \tau$ or $f(x_n) = 0$ as exit criteria. Here $\tau$ is a tolerance threshold, and if the exit criteria is met it means that $|x_n - \alpha| \leq \frac{\tau}{2}$, while the other exit criteria means we have reached an exact solution.

Algorithm 5.2.1: General Bisection Method

```
bisectionMethod (a ∈ ℝ, b ∈ (a, ∞), f ∈ C[a, b], τ ∈ ℝ⁺)
    a₀ := a
    b₀ := b
    x₀ := ½(a + b)
        n := 0
        while  f(xₙ) ≠ 0  AND  bₙ - aₙ > τ:
        if  f(xₙ) < 0:
            aₙ₊₁ := xₙ
            bₙ₊₁ := bₙ
        else :
            aₙ₊₁ := aₙ
            bₙ₊₁ := xₙ
        n ↦ n + 1
        xₙ := ½(aₙ + bₙ)
    return  xₙ
```

For our purposes we are trying to find the zero of $f(x) = x^2 - N$, which is a strictly increasing function on $\mathbb{R}_0^+$. If $N >= 1$, then $\sqrt{N} \in [0, N]$, while $N < 1 \implies \sqrt{N} \in [0, 1]$. It is obvious that our function has the required property, and thus we get the following method for finding the square root of $N$:

Algorithm 5.2.2: Bisection Method for Square Roots

```
bisectionSquareRoot (N ∈ ℝ₀⁺, τ ∈ ℝ⁺)
    a₀ := 0
    b₀ := max 1, N
    x₀ := ½(a₀ + b₀)
    n := 0
    while  xₙ² - N ≠ 0  AND  bₙ - aₙ > τ:
        if  xₙ² - N < 0:
            aₙ₊₁ := xₙ
            bₙ₊₁ := bₙ
        else :
            aₙ₊₁ := aₙ
            bₙ₊₁ := xₙ
        n ↦ n + 1
        xₙ := ½(aₙ + bₙ)
    return  xₙ
```

The implementation of this method is efficiently acheived in C using only addition, subtraction and multiplication by a constant. Before this method is implemented, however, we must first consider if and or when it converges to the correct answer. From an intuitive standpoint we

would assume that if there is only one root in the interval, it would follow that we would converge to the root.

**Proposition 5.2.1.** $\lim_{n \to \infty} x_n = \sqrt{N}$ *for Algorithm 5.2.2*

*Proof.* To prove this statement it suffices to prove that $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$ and $\lim_{n \to \infty} |x_n - \sqrt{N}| = 0$.

*Claim 1:* $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$

*Proof.* $a_0 := 0 \implies a_0 \leq \sqrt{N}$
$b_0 := \max\{1, N\} \implies b_0 \geq \sqrt{N}$
Therefore it is obvious that $\sqrt{N} \in [a_0, b_0]$
Now suppose $\sqrt{N} \in [a_n, b_n]$ for some $n \in \mathbb{N}$
It should be noted that $a_n, b_n, x_n \in \mathbb{R}_0^+ \forall n \in \mathbb{N}$ as $a_0, b_0 \in \mathbb{R}_0^+$ and all the subsequent values are derived from these using only addition and multiplication by positive factors.
We then see that $x_n := \frac{1}{2}(a_n + b_n)$, and we consider the two cases that $x_n^2 - N \leq 0$ or $x_n^2 - N \geq 0$.

**Case** $x_n^2 - N \leq 0$: $a_{n+1} := x_n, b_{n+1} := b_n$
It is therefore obvious that $\sqrt{N} \leq b_{n+1}$.
Now we see that $x_n^2 - N \leq 0 \implies x_n^2 \leq N \implies x_n \leq N$ as all the values are non-negative.
Thus $\sqrt{N} \in [a_{n+1}, b_{n+1}]$.

**Case** $x_n^2 - N \geq 0$: $a_{n+1} := a_n, b_{n+1} := x_n$
It is therefore obvious that $\sqrt{N} \geq a_{n+1}$.
Now we see that $x_n^2 - N \geq 0 \implies x_n^2 \geq N \implies x_n \geq N$ as all the values are non-negative.
Thus $\sqrt{N} \in [a_{n+1}, b_{n+1}]$.

Hence $\sqrt{N} \in [a_n, b_n] \implies \sqrt{N} \in [a_{n+1}, b_{n+1}] \forall n \in \mathbb{N}$
As $sqrtN \in [a_0, b_0]$ then we see that $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$ ∎

*Claim 2:* $\lim_{n \to \infty} |x_n - \sqrt{N}| = 0$

*Proof.* Let $n \in \mathbb{N}$ be arbitrary.
As $x_n := \frac{1}{2}(a_n + b_n)$ then we see that $|a_n - x_n| = |b_n - x_n| = \frac{1}{2}(b_n - a_n)$.
Now as $\sqrt{N} \in [a_n, b_n]$ it follows that $|\sqrt{N} - x_n| \leq \frac{1}{2}(b_n - a_n)$.
As the modulas function is a mapping from $\mathbb{R}$ to $\mathbb{R}_0^+$, it is clear that $|\sqrt{N} - x_n|$ is bounded below by 0.
Now as for each $n \in \mathbb{N}$, either $a_{n+1} = x_n$ or $b_{n+1} = x_n$, we see that $b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$.
Further we can see that $b_n - a_n \geq 0 \forall n \in \mathbb{N}$ because $b_n \geq a_n$.
Therefore the sequence of $\frac{1}{2}(b_n - a_n)$ is a strictly decreasing sequence that is bounded below, by 0. Thus $\lim_{n \to \infty} \frac{1}{2}(b_n - a_n) = 0$
Therefore $\lim_{n \to \infty} |x_n - \sqrt{N}| = \lim_{n \to \infty} \frac{1}{2}(b_n - a_n) = 0$ ∎

By using our two claims above we see that $\lim_{n \to \infty} x_n = \sqrt{N}$. □

The algorithm can be generalised to search for $\sqrt[k]{N}$, where $k \in [2, \infty) \cap \mathbb{Z}$. We can do this by using the integer power function discussed previously in section **??**. This gives the following algorithm:

Algorithm 5.2.3: Bisection Method for General Roots

```
kRootBisectionMethod (N ∈ ℝ₀⁺, k ∈ [2, ∞) ∩ ℤ, τ ∈ ℝ⁺)
      a₀ := 0
      b₀ := max 1, N
      x₀ := ½(a₀ + b₀)
      n := 0
      while  intPow(xₙ, k) − N ≠ 0  AND  bₙ − aₙ > τ :
            if  intPow(xₙ, k) − N < 0 :
                  aₙ₊₁ := xₙ
                  bₙ₊₁ := bₙ
            else :
                  aₙ₊₁ := aₙ
                  bₙ₊₁ := xₙ
            n ↦ n + 1
            xₙ := ½(aₙ + bₙ)
      return  xₙ
```

The proof that this converges to the correct value is very similar to the proof for square roots.

We can now consider the accuracy that can be acheived by our algorithm, for our purposes we will be considering $\sqrt{N}$, though the same applies for $\sqrt[k]{N}$. We know that $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$, and in particular we know that either $\sqrt{N} \in [a_n, x_n]$ or $\sqrt{N} \in [x_n, b_n] \forall n \in \mathbb{N}$; therefore we know that $\epsilon_n := \left| x_n - \sqrt{N} \right| \leq \frac{1}{2}(b_n - a_n) \forall n \in \mathbb{N}$. Then as we know that $b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$, we know that $\epsilon_n \leq \frac{1}{2^n}(b_0 - a_0)$.

We can consider that $\forall N \in \mathbb{R}_0^+ \exists (r, k) \in [\frac{1}{4}, 1) \times \mathbb{Z} : N = r \cdot 2^{2k}$; using this we know that $\sqrt{N} = \sqrt{r} \cdot 2^k$. As we have the fixed initial bounds of $a_0 = 0$ and $b_0 = 1$, then if we are finding $\sqrt{r}$ we know that $\epsilon_n \leq \frac{1}{2^n} \forall n \in \mathbb{N}$. Hence we can calculate the precision of our current estimate beforehand for any $n \in \mathbb{N}$, and thus we can guarantee $d$ significant digits of accuracy for $r \in [\frac{1}{4}, 1)$.

To get this accuracy must find $n \in \mathbb{N}$ such that $\epsilon_n \leq \frac{1}{10^d}$, to acheive this we must find $n \in \mathbb{N}$ such that $2^n \geq 10^d$. For example the following table indicates the required $n$, required for certain significant digits of accuracy.

| $d$ | $n : 2^n \geq 10^n$ |
| --- | --- |
| 1 | 0 |
| 5 | 15 |
| 10 | 30 |
| 20 | 64 |
| 50 | 163 |
| 100 | 329 |

Now ususally finding $r$ and $k$ as above would be as hard as calculating the logorithm of $N$; however due to the way that C stores real numbers as either `double` or in the MPFR

library, finding these values are actually fairly trivial. Both provide a functionality to find $(a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z} : N = a \cdot 2^b$, and from this we merely require a simple comparison and division by 2 if $b$ is not even. This leads to the following algorithm, which has the above maximum number of iterations for a required accuracy:

Algorithm 5.2.4: Bisection Method for Square Roots with fixed bounds

$$\begin{aligned}
&\mathrm{bisectionSquareRoot}\,(N \in \mathbb{R}_0^+, \tau \in \mathbb{R}^+) \\
&\quad \text{Let} \quad (r, e) \in [\tfrac{1}{2}, 1) : N = r \cdot 2^e \\
&\quad \text{if} \ \ 2 \nmid e: \\
&\quad\quad r \mapsto \tfrac{r}{2} \\
&\quad\quad e \mapsto e - 1 \\
&\quad a_0 := 0 \\
&\quad b_0 := 1 \\
&\quad x_0 := \tfrac{1}{2}(a_0 + b_0) \\
&\quad n := 0 \\
&\quad \text{while} \ \ x_n^2 - N \neq 0 \ \ \text{AND} \ \ b_n - a_n > \tau: \\
&\quad\quad \text{if} \ \ x_n^2 - N < 0: \\
&\quad\quad\quad a_{n+1} := x_n \\
&\quad\quad\quad b_{n+1} := b_n \\
&\quad\quad \text{else}: \\
&\quad\quad\quad a_{n+1} := a_n \\
&\quad\quad\quad b_{n+1} := x_n \\
&\quad\quad n \mapsto n + 1 \\
&\quad\quad x_n := \tfrac{1}{2}(a_n + b_n) \\
&\quad \text{return} \ \ x_n \cdot 2^{\frac{e}{2}}
\end{aligned}$$

## 5.3 Newton's Method for Square Roots

If we consider $f(x) = x^2 - N$ then if $x^*$ is a solution to $f(x) = 0$ we see that $x^* = \sqrt{N}$. As $f'(x) = 2x$, then the Newton's Method, will give $x_{n+1} = x_n - \frac{x^2 - N}{2x}$, where $x_0$ is a given initial guess.

We can see that, in C, each iteration will calculate `x = x - (x*x - N) / (2*x)`, which requires 5 operations; however if we re-arrange our equation, we instead get $x_{n+1} = \frac{1}{2}x_n + \frac{N}{x}$. Implementing our new iterative formula we get `x = 0.5 * (x + N/x)`, which now uses only 3 operations.

We can then use the following pseudocode as the basis of our implementaions of the Newton-Raphson Method for Square Roots:

Algorithm 5.3.1: Basic Newton Method for Square Root

$$\begin{aligned}
&\mathrm{NewtonSquareRoot}\,(N \in \mathbb{R}, x_0 \in \mathbb{R}, \tau \in (0, 1)): \\
&\quad n := 0 \\
&\quad \text{loop}: \\
&\quad\quad x_{n+1} := \tfrac{1}{2}(x_n + \tfrac{N}{x_n}) \\
&\quad\quad \delta_n := |x_{n+1} - x_n| \\
&\quad\quad \text{if} \ \ \delta_n \leq \tau: \\
&\quad\quad\quad \text{return} \ \ x_{n+1}
\end{aligned}$$

$$n \mapsto n + 1$$

Next we want to consider our initial estimate $x_0$; it is prudent to first consider when our initial estimate will converge to the correct root. By looking at a graph of the function, and in particular the tangents to the curve, it would seem reasonable to wonder if $\lim_{n \to \infty} x_n = \sqrt{N}$.

**Proposition 5.3.1.** *If $x_0 \in \sqrt{N}, \infty$ and $\{x_n : n \in \mathbb{N}\}$ is a sequence of approximations of $\sqrt{N}$ found via the Newton-Raphson Method, as detailed above, then:*

$$\lim_{n \to \infty} x_n = \sqrt{N}$$

*Proof.* Suppose $x_n > \sqrt{N}$, then

$$
\begin{aligned}
x_{n+1} &= \frac{1}{2}\left(x_n + \frac{N}{x_n}\right) \\
&< \frac{1}{2}\left(x_n + \frac{N}{\sqrt{N}}\right) \qquad\qquad \text{as} \sqrt{N} < x_n \implies \frac{1}{x_n} < \frac{1}{\sqrt{N}} \\
&= \frac{1}{2}\left(x_n + \sqrt{N}\right) \\
&< \frac{1}{2}(2x_n) \\
&= x_n
\end{aligned}
$$

Therefore we see that $\{x_k : k \in [n, \infty) \cap \mathbb{Z}\}$ is a strictly decreasing sequence.
Now suppose that $x_n \geq \sqrt{N}$ and then, for a contradiction, assume that $x_{n+1} < \sqrt{N}$. We then see that:

$$
\begin{aligned}
\frac{1}{2}\left(x_n + \frac{N}{x_n}\right) &< \sqrt{N} \\
\implies x_n + \frac{N}{x_n} &< 2\sqrt{N} \\
\implies x_n^2 + N &< 2\sqrt{N}x_n \\
\implies x_n^2 - 2\sqrt{N}x_n + N &< 0 \\
\implies \left(x_n - \sqrt{N}\right)^2 &< 0
\end{aligned}
$$

This is a contradiction as $x_n, \sqrt{n} \in \mathbb{R} \implies \left(x_n - \sqrt{N}\right)^2 \geq 0$.
Therefore $x_n \geq \sqrt{N} \implies x_{n+1} \geq \sqrt{N}$.
Hence if $x_0 > \sqrt{N}$, then it follows that $\{x_n : n \in \mathbb{N}\}$ is a strictly decreasing sequence that is bounded below. Therefore by an elementary result from limit theory, we see that $\lim_{n \to} x_n = \inf\{x_n : n \in \mathbb{N}\}$. $\qquad\square$

The most obvious choice for $x_0$ would be $N$, but we see that $N \in (0, 1)$, then $N < \sqrt{N}$. In this case, we could choose $x_0 = 1$ for the case that $N \in (0, 1)$. Therefore we can choose

$$x_0 := \begin{cases} N &:\ N \in (1, \infty) \\ 1 &:\ N \in (0, 1) \end{cases}$$

In our choice of $x_0$, we have so far left out the cases where $N \in \{0, 1\}$. In both of these case we already know the correct answer, namely $\sqrt{N} = N$ provided $N \in 0, 1$. Therefore we can exclude them from our calculations, as we can pre-asses the value of $N$, simply returning the correct answer if one of these cases is encountered.

This then leads to an updated version of the above pseudocode:

Algorithm 5.3.2: Basic Newton Method for Square Root

```
NewtonSquareRoot ( N ∈ ℝ₀⁺, τ ∈ (0,1) ):
    if  N ∈ {0,1}:
        return  N
    if  N > 1:
        x₀ := N
    else:
        x₀ := 1
    n := 0
    loop:
        xₙ₊₁ := ½(xₙ + N/xₙ)
        δₙ := |xₙ₊₁ − xₙ|
        if  δₙ ≤ τ:
            return  xₙ₊₁
        n ↦ n + 1
```

TODO: Write up examination and implementation of this pseudocode

An alternative would be to use the integer square root method discussed in Section 5.1 to improve our initial choice of $x_0$. We will start by showing, that for intervals $I \subset \mathbb{R}^+$, the first two criteria for quadratic convergence of the Newton Raphson method are met.

**Proposition 5.3.2.** *If $I \subset \mathbb{R}^+$ then $NR_1$ and $NR_2$ are satisfied for $f(x) = x^2 - N$*

*Proof.* $f(x) = x^2 - N \implies f'(x) = 2x \implies f''(x) = 2$
Now as $x \in \mathbb{R}^+ \forall x \in I$, then it is obvious that $f'(x) > 0$
Therefore $f'(x) \neq 0 \forall x \in I$, and so $NR_1$ is satisfied.
As $f''(x)$ is a constant function, then it is continuous on all of $\mathbb{R}$.
Hence $f''(x)$ is continuous $\forall x \in I$ and so $NR_2$ is satisfied. $\qquad \square$

Now the integer square root function will always produce a root that is at most a distance of 1 from $\sqrt{N}$; therefore we can consider $I = [\sqrt{N} - 1, \sqrt{N} + 1]$. Now if $N \leq 1$, then $I\mathbb{R}^+$ and so we cannot guarantee the satisfaction of $NR_1$. Therefore we can proceed with our analysis of the case that $N > 1$.

If $N > 1$ we need to find when we can satisfy $NR_3$. First, we remember that $M := \sup \left| \frac{f''(x)}{f'(x)} \right| : x \in I$ and $\epsilon_0 := \left| x_0 - \sqrt{(N)} \right|$. Then to satisfy $NR_3$, we must have that $M\epsilon_0 < 1$.

We can guarantee that $\epsilon_0 \leq 1$ because $x_0 \in I$ from the integer square root algortihm; therefore it suffices to find the situation where $M < 1$. As both $f'$ and $f''$ are continuous and non-zero

on $I$ it follows that $M = \sup x^{-1} : x \in I = (\sqrt{N} - 1)^{-1}$. We then see that:

$$M < 1 \iff \sqrt{N} - 1 > 1$$
$$\iff \sqrt{N} > 2$$
$$\iff N > 4$$

Therefore we can get the following new choice for $x_0$, and thus new pseudocode:

$$x_0 := \begin{cases} 1 & : \ N \in (0, 1) \\ N & : \ N \in (1, 4] \\ intSqrt(N) & : \ N \in (4, \infty) \end{cases}$$

Algorithm 5.3.3: Basic Newton Method for Square Root

```
NewtonSquareRoot(N ∈ ℝ₀⁺, τ ∈ (0,1)):
    if N ∈ {0,1}:
        return N
    if N < 1:
        x₀ := 1
    else:
        if N ≤ 4:
            x₀ := N
        else:
            x₀ := IntSqrt(N)
    n := 0
    loop:
        xₙ₊₁ := ½(xₙ + N/xₙ)
        δₙ := |xₙ₊₁ - xₙ|
        if δₙ ≤ τ:
            return xₙ₊₁
        n ↦ n + 1
```

TODO: Write up examination of different versions tried, such as using $x_0 = N$, etc...

If we consider any $N \in \mathbb{R}_0^+$, then $\exists a \in \left[\frac{1}{2}, 1\right), b \in \mathbb{Z} : N = a \times 2^b$. Finding this value would be a hard as finding the logarithm of $N$ base 2, but due to the representation of numbers within C, both standard C and MPFR have functions that allow us to extract these two values with minimal computational expenditure.

This helps as we can then narrow our problem, to only finding $\sqrt{a} : a \in \left[\frac{1}{2}, 1\right)$, and then calculating

$$\sqrt{N} = \sqrt{a} \times 2^{\lfloor \frac{b}{2} \rfloor} \times \alpha \text{ where } \alpha = \begin{cases} 1 & : \ b \in 2\mathbb{Z} \\ \sqrt{2} & : \ b \in \mathbb{Z}^+ \setminus 2\mathbb{Z} \\ \frac{1}{sqrt2} & : \ b \in \mathbb{Z}^- \setminus \mathbb{Z} \end{cases}$$

We then get the following algorithm, which implements this:

Algorithm 5.3.4: Newton Method for Square Root v3

```
NewtonSquareRoot(N ∈ ℝ₀⁺, τ ∈ (0,1)):
    Let (a,b) :∈ [½,1) × ℤ s.t. N = a * 2ᵇ
```

```
    x_0 := 1
    if  b ≡ 0 mod 2 :
        α := 1
    else :
        if  b > 0 :
            α := √2
        else :
            α := 1/√2
    n := 0
    loop :
        x_{n+1} := 1/2 (x_n + a/x_n)
        δ_n := |x_{n+1} − x_n|
        if  δ_n ≤ τ :
            return  α · x_{n+1} · 2^{⌊b/2⌋}
        n ↦ n + 1
```

We must first consider the fact that the algorithm requires the pre-calculation of both $\sqrt{2}$ and $\frac{1}{\sqrt{2}}$, to be able to calculate all values. However, it turns out we can use the algorithm itself to generate these values as $2 = \frac{1}{2} \cdot 2^2$, and as the exponent of 2 is even then the algorithm does not require $\sqrt{2}$ for this computation. Similarly $\frac{1}{2} = \frac{1}{2} \cdot 2^0$, which again is an even exponent. We can thus run our algorithm to find an arbitrarily accurate values for $\sqrt{2}$ and $\frac{1}{\sqrt{2}}$ to allow us to run the algorithm for other values.

With this observation can then consider $N \in \left[\frac{1}{2}, 1\right)$. As this is a small range and, as per our previous algorithm, we use an initial guess of $x_0 = 1$, then we can prove that our algorithm will converge quadratically to $\sqrt{N}$.

**Proposition 5.3.3.** *Algorithm 5.3.4, satisfies the criteria of Theorem 2.3.2, and thus has quadratic convergence to $\sqrt{N}$.*

*Proof.* To fulfill the criteria of Theorem 2.3.2, we must find and interval $I := [\sqrt{N}-r, \sqrt{N}+r]$ for some $r \geq \epsilon_0$.

Consider $\epsilon_0 = |\sqrt{N} - x_0| = 1 - \sqrt{N}$. We see that as $N \geq \frac{1}{2}$ then $\sqrt{N} \geq \sqrt{2}^{-1}$, and thus $\epsilon_0 \leq 1 - \sqrt{2}^{-1}$. Let us have $r := 1 - \frac{1}{\sqrt{2}}$, and $I$ as defined above.

If we look at the lower bound of $I$, then we see that:

$$\sqrt{N} - r \geq \frac{1}{\sqrt{2}} - \left(1 - \frac{1}{\sqrt{2}}\right)$$
$$= \frac{2}{\sqrt{2}} - 1$$
$$= \sqrt{2} - 1$$
$$> 0$$

Therefore we see that $I \subset \mathbb{R}^+$, and so by Proposition 5.3.2 we get that $\mathrm{NR}_1$ and $\mathrm{NR}_2$ ar satisfied. It then remains to show that $\mathrm{NR}_3$ is satisfied on $I$.

Now by the definition in Theorem 2.3.2, we have that $M = \sup\left\{\frac{1}{2}\left|\frac{f''(x)}{f'(y)}\right| : x, y \in I\right\}$. We know that $I$ is bounded, $f''(x) = 2$ and $f'(x) = 2x$ meaning that $\frac{1}{2}\left|\frac{f''(x)}{f'(y)}\right| = \frac{1}{f'(x)}$ as $x \in \mathbb{R}^+$.

Therefore our problem is reduced to finding $\max\left\{\frac{1}{2x} : x \in I\right\}$, which is equivalent to finding $\min\{x : x \in I\} = \sqrt{N} - r$. Therefore by passing this information back up the chain we get that

$$M = \frac{1}{2(\sqrt{N} - r)}$$

Then we see that:

$$
\begin{aligned}
M\epsilon_0 &= \frac{1 - \sqrt{N}}{2(\sqrt{N} - r)} \\
&\leq \frac{1 - \frac{1}{\sqrt{2}}}{2(\sqrt{N} - r)} \qquad \text{as } \sqrt{N} \geq \frac{1}{\sqrt{2}} \\
&\leq \frac{1 - \frac{1}{\sqrt{2}}}{2(\frac{1}{\sqrt{2}} - r)} \qquad \text{as } \sqrt{N} \geq \frac{1}{\sqrt{2}} \\
&= \frac{1 - \frac{1}{\sqrt{2}}}{2(\frac{2}{\sqrt{2}} - 1)} \\
&= \frac{1 - \frac{1}{\sqrt{2}}}{2\sqrt{2}(1 - \frac{1}{\sqrt{2}})} \\
&= \frac{1}{2\sqrt{2}} \\
&< 1 \qquad\qquad \text{as } 2\sqrt{2} > 1
\end{aligned}
$$

As we have confirmed that $M\epsilon_0 < 1$, then we have confirmed that $\mathrm{NR}_3$ is satisfied on $I$, and so the algorithm converges quadratically to the desired root. □

Using the previous proposition we can, similar to our previous methods, consider how many iterations would be needed to reach a required tolerance. To start we consider that, as mentioned in the proof or Theorem 2.3.2, that $\epsilon_n \leq (M\epsilon_0)^{2^n - 1}\epsilon_0$.

We know that $M\epsilon_0 \leq \frac{1}{2\sqrt{2}}$ and that $\epsilon_0 \leq 1 - \frac{1}{\sqrt{2}}$, giving:

$$\epsilon_n \leq \left(\frac{1}{2\sqrt{2}}\right)^{2^n - 1}\left(1 - \frac{1}{\sqrt{2}}\right)$$

Thus if we want to acheive a tolerance of $\epsilon_n \leq \tau$, then it suffices to find $n \in \mathbb{N}_0$ such that:

$$\left(\frac{1}{2\sqrt{2}}\right)^{2^n - 1} \leq \tau$$

Then,

$$(2^n - 1)\log\left(\frac{1}{2\sqrt{2}}\right) \leq \log\left(\frac{\tau}{1 - \frac{1}{\sqrt{2}}}\right)$$

By noting that $\log(\frac{1}{a}) = -\log(a)$, then we get

$$(1 - 2^n)\log(2\sqrt{2}) \leq \log\left(\frac{\tau}{1 - \frac{1}{\sqrt{2}}}\right)$$

Once this is rearranged we get the following inequality:

$$2^n \geq \frac{\log\left(\frac{2(\sqrt{2}-1)}{\tau}\right)}{\log(2\sqrt{2})}$$

By taking logarithms again and re-arranging we get that

$$n \geq \frac{\log\left(\frac{\log\left(\frac{2(\sqrt{2}-1)}{\tau}\right)}{\log(2\sqrt{2})}\right)}{\log(2)} = \log_2\left(\log_{2\sqrt{2}}\left(2\frac{\sqrt{2}-1}{\tau}\right)\right)$$

Now for an example, suppose we want to know how many iterations we need to perform to find $\sqrt{N}$ to within 10 decimal places, i.e. $\tau = 10^{-}10 = 0.0000000001$. We remember that $\sqrt{N} \in [\frac{1}{2}, 1)$, and then we will apply transformations to this value afterwards, therefore this is equivalent to finding 10 significant digits of accuracy for our square root (ignoring any loss of accuracy that may arrise from multiplications afterwards).

Now in this case we want to find $n \in \mathbb{N}$ such that $n \geq log_2(log_{2\sqrt{2}}(2 \cdot 10^{10}(\sqrt{2} - 1)))$. Using Wolfram Alpha to calculate this value we get that we need $n \geq 4.457144...$ and so we can take $n = 5$. This means that we could modify our algorithm and implementation to do 5 fixed iterations of Newton's Method to guarantee at least 10 decimal places of accuracy.

In terms of efficiency versus accuracy tradeoff modifying the problem thus would improve it's efficiency by removing, now unneccesary, calculation and comparrison of $\delta_n$ at each stage. However this does need a fixed guaranteed accuracy, and therefore such a program would no longer be suitable if we needed to calculate a square root accurate to 15 decimal places.

Below is a table that lists the minimum $n \in \mathbb{N}$ such that $n$ satisfies our inequality, where our tolerance is $10^k$ for some $k \in \mathbb{N}$. This will give us the maximum number of iterations that must be performed for the required accuracy.

| $k : \tau = 10^k$ | $n$ |
|---|---|
| 5 | 4 |
| 10 | 5 |
| 100 | 8 |
| 1,000 | 12 |
| 1,000,000 | 22 |

## 5.4  Newton's Inverse Square Root Method

As discussed in Section **??**, computers are more efficient at multiplication over division. We would therefore prefer to find a way of utilising Neton's Method without having to perform any costly division operations.

If we consider $f(x) = N - \frac{1}{x^2}$ then if $x^*$ is a solution to $f(x) = 0$ we see that $x^* = \frac{1}{\sqrt{N}}$. As $f'(x) = \frac{2}{x^3}$, then the Newton's Method, will give

$$x_{n+1} = x_n - \frac{N - \frac{1}{x_n^2}}{\frac{2}{x_n^3}} = x_n \left( \frac{3}{2} - \frac{N}{2} x_n^2 \right)$$

where $x_0$ is a given initial guess. As can be seen this algorithm requires no division if we multiply by real constants rather than the division implied above.

We can then consider that, similar to Algorithm 5.3.4, any $N$ can be represented as $a \cdot 2^b$ where $a \in \left[ \frac{1}{2}, 1 \right)$. This will, again allow us to narrow our problem to a known range of values, by using the following transormations.

$$N = a \cdot 2^b \implies \frac{1}{N} = \frac{1}{a} \cdot 2^{-b}$$

$$\implies \frac{1}{\sqrt{N}} = \frac{1}{a} \cdot 2^{\lfloor \frac{-b}{2} \rceil} \cdot \alpha \qquad \alpha := \begin{cases} 1 & : & b \equiv 0 \mod 2 \\ \sqrt{2} & : & b \equiv 1 \mod 2, b \in \mathbb{Z}^- \\ \frac{1}{\sqrt{2}} & : & b \equiv 1 \mod 2, b \in \mathbb{Z}^+ \end{cases}$$

$$\implies \sqrt{N} = N \cdot \frac{1}{\sqrt{a}} \cdot 2^{\lfloor \frac{-b}{2} \rceil} \cdot \alpha$$

Therefore we only need to calculate inverse square roots for values of $N$ in the range $\left[ \frac{1}{2}, 1 \right)$. Thus giving us the following algorithm:

Algorithm 5.4.1: Newton Inverse Square Root Method

$$\begin{array}{l}
\mathrm{NewtonSquareRoot}\,(\,N \in \mathbb{R}_0^+, \tau \in (0,1)\,)\colon \\
\quad \mathrm{Let}\ (a,b) :\in \left[\frac{1}{2}, 1\right) \times \mathbb{Z}\ \mathrm{s.t.}\ N = a * 2^b \\
\quad x_0 := 1 \\
\quad \mathrm{if}\ b \equiv 0 \bmod 2\colon \\
\qquad \alpha := 1 \\
\quad \mathrm{else}\colon \\
\qquad \mathrm{if}\ b > 0\colon \\
\qquad\quad \alpha := \frac{1}{\sqrt{2}} \\
\qquad \mathrm{else}\colon \\
\qquad\quad \alpha := \sqrt{2} \\
\quad n := 0 \\
\quad \mathrm{loop}\colon \\
\qquad x_{n+1} := x_n\left(\frac{3}{2} + \frac{a}{2} x_n^2\right) \\
\qquad \delta_n := |x_{n+1} - x_n| \\
\qquad \mathrm{if}\ \delta_n \leq \tau\colon \\
\qquad\quad \mathrm{return}\ N \cdot \alpha \cdot x_{n+1} \cdot 2^{\lfloor \frac{-b}{2} \rceil} \\
\qquad n \mapsto n + 1
\end{array}$$

With this method we can once again consider it's convergence properties, in particular does it satisfy the criteria for quadratic convergence in Theorem 2.3.2.

**Proposition 5.4.1.** *Algorithm 5.4.1 satisfies the criteria of Theorem 2.3.2, and thus has quadratic convergence to $\sqrt{N}$.*

*Proof.* We know that we only need to consider $N \in [\frac{1}{2}, 1)$, and therefore $\sqrt{N}^{-1} \in (1, \sqrt{2}]$. Also $x_0 = 1$ and so we see that

$$\epsilon_0 = |x_0 - \sqrt{N}^{-1}| = \sqrt{N}^{-1} - x_0 \leq \sqrt{2} - 1$$

Now let $r := \epsilon_0 = \sqrt{N} - 1$ and $I := [\sqrt{N}^{-1} - r, \sqrt{N}^{-1}]$. If we consider the lower bound of I we see that $\sqrt{N}^{-1} - (\sqrt{N}^{-1} - 1) = 1$, and in particular $0 \notin I$.

Next we know that $f(x) = N - x^{-2}$, and therefore we get $f'(x) = 2x^{-3}$, $f''(x) = -6x^{-4}$. It is obvious that $\nexists x \in \mathbb{R} : f'(x) = 0$, which means that $f'(x) \neq 0 \forall x \in I$ and so $\mathrm{NR}_1$ is satisfied. Also as $f''$ is only discontinuous at $x = 0$ and $0 \notin I$, then $f''(x)$ is continuous $\forall x \in I$, meaning this satisfies $\mathrm{NR}_2$.

Now $M = \sup \left\{ \frac{1}{2} \left| \frac{2x^3}{6y^4} \right| : x, y \in I \right\}$, we can simplify the function we are trying to minimise to get $\frac{1}{6} \frac{x^3}{y^4}$. It is obvious that in order to maximise this function we should find the largest possibe $x$ and smallest possible $y$, as both are positive. Hence by taking $x = \sqrt{N}^{-1} + r$ and $y = 1$, then $M = \frac{1}{6}(2\sqrt{N}^{-1} - 1)^3 \leq \frac{1}{6}(2\sqrt{2} - 1)^3$.

Now we consider $M\epsilon_0$:

$$\begin{aligned} M\epsilon_0 &= \frac{1}{6}(2\sqrt{N}^{-1} - 1)^3(\sqrt{N} - 1) \\ &\leq \frac{1}{6}(2\sqrt{2} - 1)^3(\sqrt{2} - 1) \\ &\approx 0.42199376\ldots \\ &< 1 \end{aligned}$$

Therefore as $M\epsilon_0 < 1$ we have satisfied $\mathrm{NR}_3$, and as such we have quadratic convergence of our method to $\sqrt{N}^{-1}$. $\qquad\qquad\square$

## 5.5   Comparrison of Methods

We have observed several methods that can be used to calculate Square Roots, and so now we will see how the methods compare to each other in theory and practice. The exact root method that we first discussed is the hardest to compare to the other methods as it works in a very different manner to the others. For now we will merely observe that it is an inefficient method that can be will be shown to tae longer than the others.

Now for the other methods we will be comparing Algorithms 5.2.4, 5.3.4 and 5.4.1; we will first consider their computational complexity. As all of these methods have a variable number of iterations that they can perform then we need to consider the complexity fore each iteration step.

First for the bisection method if we calculate $x_n^2 - N$ (which takes 3 operations) whenever we generate $n$, and as we also calculate $x_n := \frac{1}{2}(a_n + b_n)$ each iteration. Then the total number of operations for each iteration is 6.

Next the regular newton method calculates $\frac{1}{2}(x_n + \frac{a}{x_n})$ and $|x_n + 1 - x_n|$ each iteration whic take 3 and 3 operations respectively, giving a total of 6 operations per iteration; however as previously noted one of those operations is a costly division operation. Finally the inverse newton method calculates $x_n(\frac{3}{2} + \frac{a}{2}x_n^2$ and $|x_n + 1 - x_n|$ each iteration which take 4 and 3 opertaions respectively, giving a total of 7 operations per iteration.

As we can see the number of operations per iteration are similar for each of these methods and so the number of required iterations will have a large impact on the speed of these methdos. Now we have seen the required number of iterations for a tolerance $\tau = 10^{-k} : k \in \mathbb{N}$, for both the bisection and basic newton square root methods, and similar to the basic newton method, we can show that for the inverse newton method we are looking for $n \in \mathbb{N}$ that satisfies the following inequality:

$$n > \log_2 \left( \log_{\frac{1}{6}(\sqrt{2}-1)(2\sqrt{2}-1)^3} \left( \frac{\tau}{\sqrt{2}-1} \right) \right) - 1$$

This gives the following table:

| $k : \tau = 10^k$ | Bisection Method | Newton Method | Inverse Newton |
|---|---|---|---|
| 5 | 16 | 4 | 4 |
| 10 | 33 | 5 | 5 |
| 100 | 332 | 8 | 9 |
| 1,000 | 3321 | 12 | 12 |
| 1,000,000 | 3219280 | 22 | 22 |

# 6 Logarithms and Exponentials

## 6.1 Taylor Series Method

## 6.2 Hyperbolic Series Method

## 6.3 CORDIC

TODO: Fill this out with stuff

# 7 Preliminary References

http://math.exeter.edu/rparris/peanut/cordic.pdf
Inside your Calculator by Gerald R Rising
Wolfram Alpha