

7.5

Mathematical Methods

Padé Approximants

1 Introduction

Programming Task: Writing Programs A and B

The programs written for this task can be found on pages ?? and ??. From the numpy package in Python, I use the `lstsq` function as an equivalent of `mldivide` from Matlab. I first tested Program A for basic functions such as $f(x) = 0$ with different values of L and M . Then I carried out testing with more complex functions such as $f(x) = \sin(x)$ and confirming the $O(x^{L+M+1})$ accuracy via polynomial division of the results.

Question 1

NEED TO TAKE $M = L + 1$ SOMEWHERE IN THIS PROJECT. IN QUESTION 4? - WHAT IS THE DETERMINANT IN PROGRAM A HERE? - COULD IT AFFECT ERROR IN QUESTION 3?

Using the binomial expansion, we obtain,

$$f_1(x) = 1 + \frac{x}{2} - \frac{x^2}{8} + \frac{x^3}{16}$$

with the following formula for the coefficients,

$$c_0 = 1, \quad c_1 = \frac{1}{2}, \quad c_k = \frac{(-1)^{k-1}(2k-3)!}{2^{2k-2}k!(k-2)!} \text{ for } k \geq 1.$$

The radius of convergence can be found via the ratio test.

$$\begin{aligned} \text{Radius} &= \lim_{k \rightarrow \infty} \left| \frac{c_k}{c_{k+1}} \right| \\ &= \lim_{k \rightarrow \infty} \left| \frac{(2k-1)(2k-2)}{4(k+1)(k-1)} \right| \\ &= 1 \end{aligned}$$

This means that the power series will only provide a useful estimate in the disk of radius 1 centred at 0 in the complex plane. Also, the further from 0 you go, the more terms of the power series are required for an equally precise result. Since a Padé approximant approximates the power series we would expect approximants with a higher value of $L + M + 1$ to be much more useful further from 0 (and potentially beyond 1).

Taking $x = 1$ in the power series, we obtain $\sum_{k=0}^{\infty} c_k$. Since this converges, the sequence of partial sums $\sum_{k=0}^N c_k$ converges. This convergence is illustrated in Figure 1.

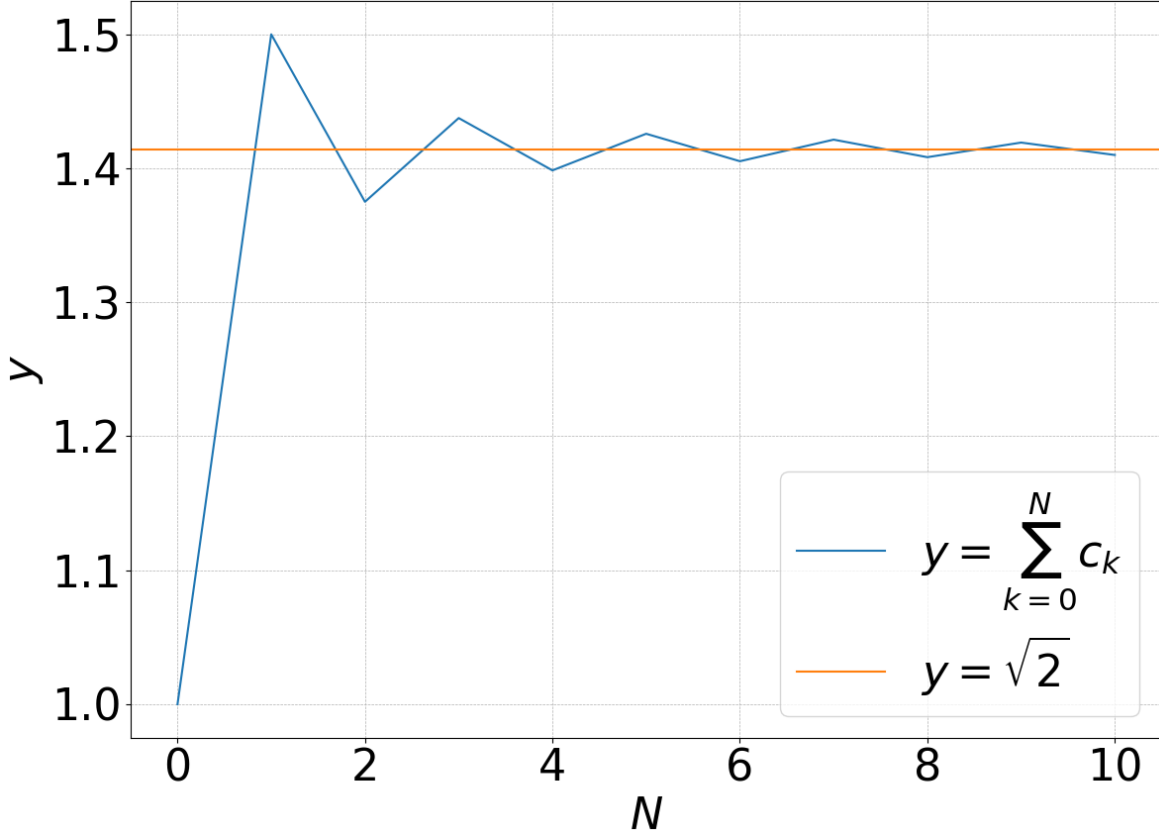


Figure 1: Line graph of partial sums

We notice that the partial sums oscillate above and below $\sqrt{2}$, with the error getting smaller as N increases.

We use the Lagrange form of the remainder for a Taylor expansion to estimate the error of the partial sum as an estimate of $\sqrt{2}$. The remainder at $x = 1$ will be,

$$R_N(1) = \frac{f_1^{(N+1)}(\xi)}{(N+1)!}$$

for some real number ξ between 0 and 1 and following simplifications, we have,

$$|R_N(1)| = \frac{(2N-1)!}{2^{2N}(N-1)!(N+1)!} (1+\xi)^{-N+\frac{1}{2}}$$

Using the function `find_xi` in the program on page ??, we find that ξ is small and decreases as N increases. For example, for $N = 3$, $\xi = 0.230$; for $N = 10$, $\xi = 0.0681$ and for $N = 50$, $\xi = 0.0138$. Then the `print_xi_factor` function also in the program on page ?? can be used to show that for $N < 80$, $0.5 \leq (1+\xi)^{-N+\frac{1}{2}} \leq 0.7$ and $(1+\xi)^{-N+\frac{1}{2}}$ is decreasing. Hence, we have,

$$\begin{aligned} R_N(1) &\leq 0.7 \cdot \frac{(2N-1)!}{2^{2N}(N-1)!(N+1)!} \\ &= 0.7 \cdot \frac{1 \cdot 3 \cdot 5 \cdot \dots \cdot (2N-1)}{2 \cdot 4 \cdot 6 \cdot \dots \cdot 2N} \cdot \frac{1}{2N+2} \end{aligned} \quad (1)$$

We prove the following lemma:

$$\frac{1 \cdot 3 \cdot 5 \cdot \dots \cdot (2N-1)}{2 \cdot 4 \cdot 6 \cdot \dots \cdot 2N} \leq \frac{1}{\sqrt{2N}}$$

$$\begin{aligned} \text{Proof: } \left(\frac{1 \cdot 3 \cdot 5 \cdot \dots \cdot (2N-1)}{2 \cdot 4 \cdot 6 \cdot \dots \cdot 2N} \right)^2 &= \frac{1 \cdot 1}{2 \cdot 2} \cdot \frac{3 \cdot 3}{4 \cdot 4} \cdots \frac{(2N-1) \cdot (2N-1)}{2N \cdot 2N} \\ &= \frac{1 \cdot 3}{2 \cdot 2} \cdot \frac{3 \cdot 5}{4 \cdot 4} \cdots \frac{(2N-1)}{(2N)^2} \\ &\leq \frac{2N-1}{(2N)^2} \quad \text{since } (n-1)(n+1) = n^2 - 1 \leq n^2 \end{aligned}$$

$$\leq 1/2N \quad \square$$

Applying the lemma to (1), we obtain the following overestimate for the error as N increases:

$$|R_N(1)| \approx \frac{0.7}{(2N+2)\sqrt{2N}} \quad (2)$$

And for large $N > 80$ the following overestimate is more accurate:

$$|R_N(1)| \approx \frac{0.505}{(2N+2)\sqrt{2N}}$$

Figure 2 illustrates (2) as an error bound. We can see that this gives an accurate approximation of the error.

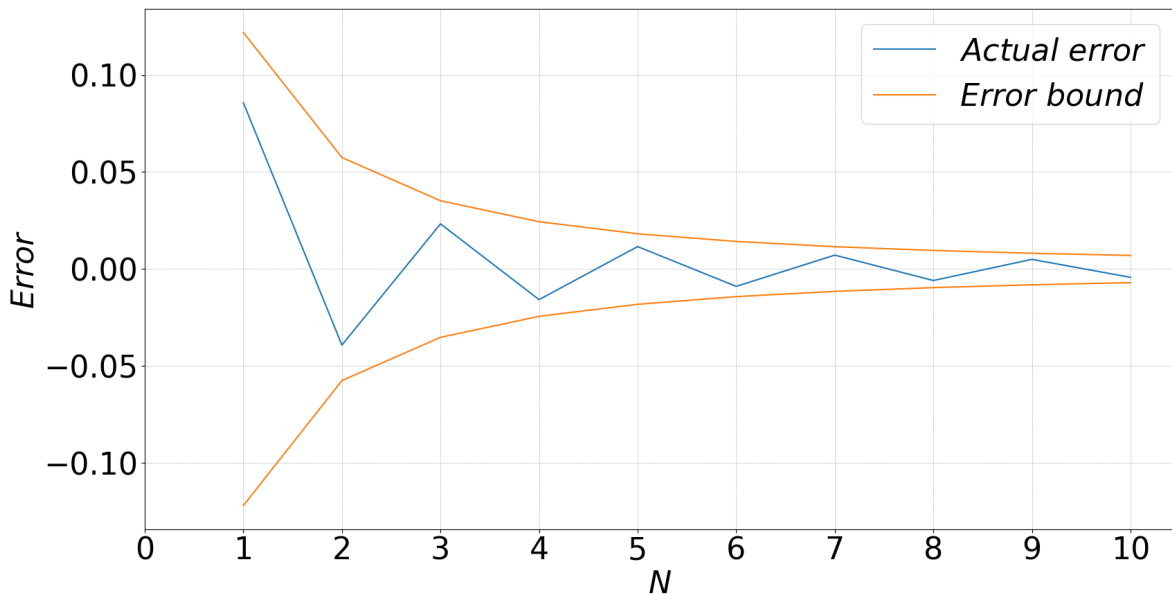


Figure 2: Actual error and estimated error of the partial sum as an estimate of $\sqrt{2}$

We know $(1 + \xi)^{-N+\frac{1}{2}}$ is relatively unchanged when N is increased by 1. Thus taking the ratio of $R_N(1)$ and $R_{N+1}(1)$ shows that incrementing N by 1 decreases the error by almost exactly $\frac{2N+1}{2N+4}$. Hence the larger N is, the less the percentage decrease of error is for increasing N .

Question 2

It is now much more difficult to obtain a theoretical result for the error, so we investigate $R_{L,L}(1)$ as an estimate of $\sqrt{2}$ numerically. The results in the table below show the error of the approximant for different values of L .

L	Error
0	0.41421356237309515
1	0.014213562373095234
2	0.00042045892481934466
3	1.2378941142587863e-05
4	3.644035522221145e-07
5	1.072704058913132e-08
6	3.1577518377901015e-10
7	9.29567534058151e-12
8	2.737809978725636e-13
9	7.993605777301127e-15
10	4.440892098500626e-16
11	2.220446049250313e-16
12	2.220446049250313e-16
13	6.661338147750939e-16
14	6.661338147750939e-16
15	2.220446049250313e-16

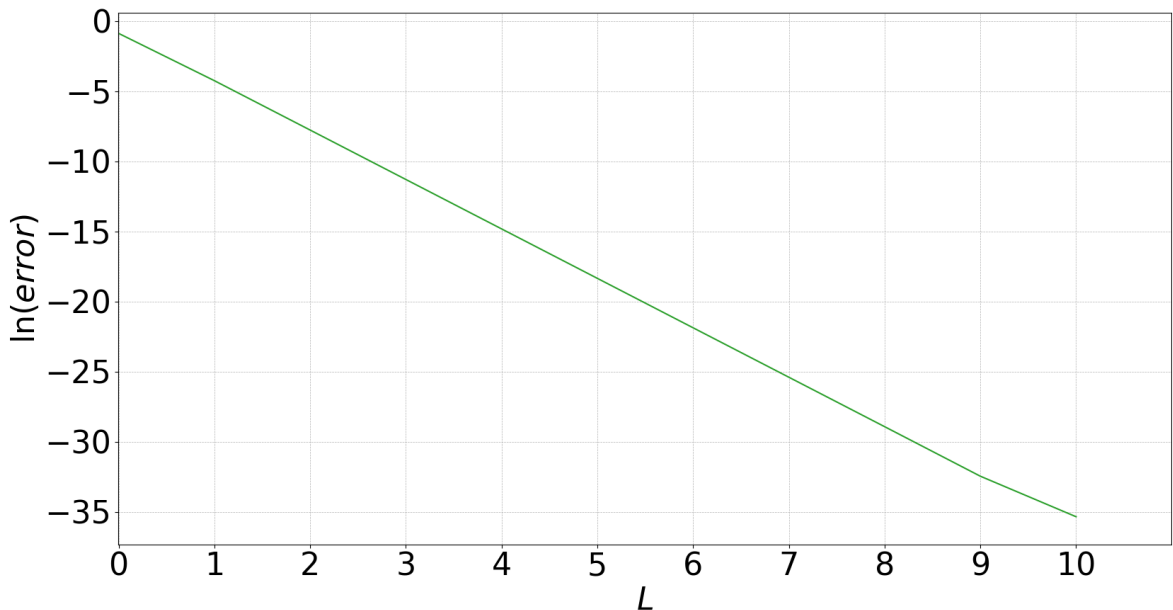


Figure 3: Plot of L against $\ln(\text{Error})$

From Figure 3, we see that for $L \leq 10$, the error decreases exponentially. We observe from the table that the minimum value of the error for $L \leq 15$ is $2.220446049250313\text{e-}16$. This must be true in general since this is exactly the machine precision, i.e. it is the difference between 1.0 and the smallest 64-bit double precision floating-point value larger than 1.0. Hence, the machine precision determines the smallest error.

In cases where the matrix used to solve equation (4) is non-singular, iterative improvement will make no difference since the exact solution is found so no more improvements can be made. The determinant of the matrix in this case approaches 0 as L increases. The `lstsq` function which I have used as the Python equivalent of `mldivide` finds the least-squares solution of the equation $A\mathbf{x} = \mathbf{b}$. Suppose for some L the determinant of the matrix used to solve equation (4) were 0 and let the least squares solution for the q_k be \mathbf{y} . Then $\mathbf{b} - A\mathbf{y}$ will be orthogonal to $A\delta\mathbf{y}$ for any $\delta\mathbf{y}$. This is because for least square problems, we have $A^T(\mathbf{b} - A\mathbf{y}) = 0$ which can be derived by differentiation of $\|\mathbf{b} - A\mathbf{y}\|^2$. Hence, no more improvements can be made in this case either.

In addition, the limit on the error is caused by the machine precision, not the solution to equation (4). Thus iterative improvement would have no effect on the minimum error.

In the power series of $R_{L,L}(x)$, the first $2L + 1$ terms match that of $f_1(x)$. The error

of $R_{L,L}(1)$ is much less than the power series estimate of $\sqrt{2}$ for the same number of matching terms. For instance, with $L = 5$, the error of the Padé approximant is 1.07×10^{-8} while for $N = 10$ the error from the partial sum is 4.28×10^{-2} . This is surprising because using the same amount of information, a much more accurate estimate is obtained. This can be explained by the fact that as we approach the radius of convergence the error power series expansion of $f_1(x)$ diverges at a faster rate than the error term from the power series expansion of $R_{L,L}(x)$.

It is then clear that the Padé approximant should be used as an estimate of $\sqrt{2}$ to specified accuracy in all cases. The error estimation (2) shows that to have an error of 2.22×10^{-16} , which only requires $L = 11$ for the Padé approximant, you would need N to be close to 100,000. Even if you wanted more accuracy than this, that wouldn't be possible with 64-bit floats since 2.22×10^{16} is the machine precision.

Question 3

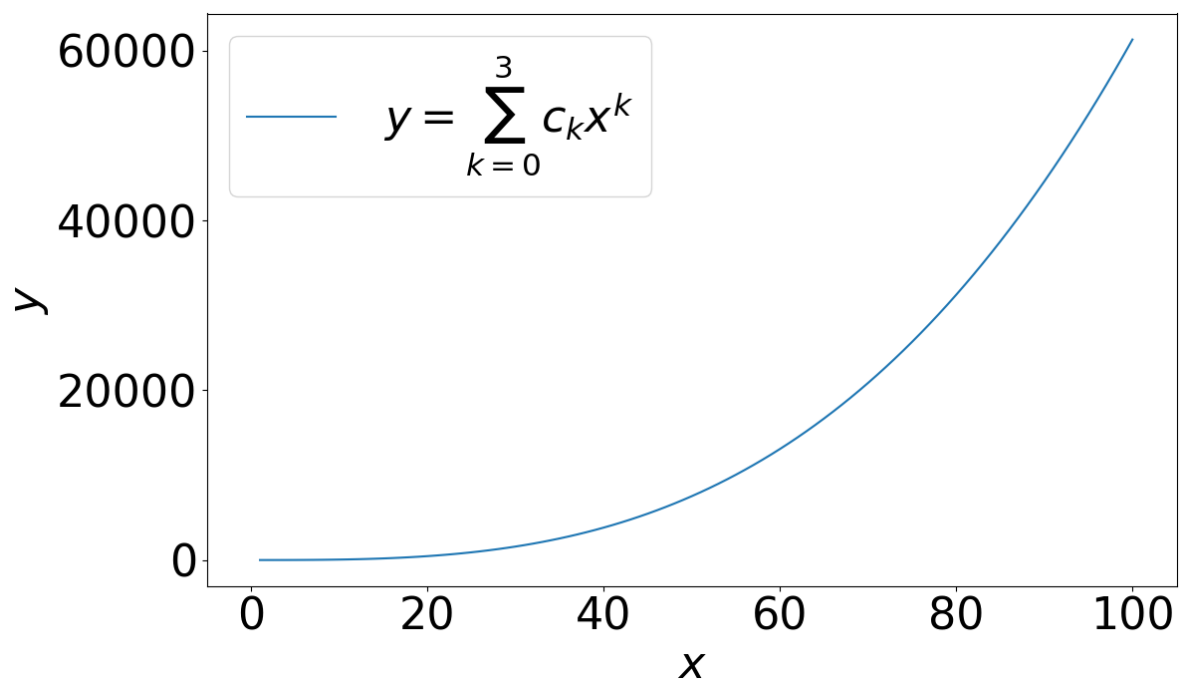
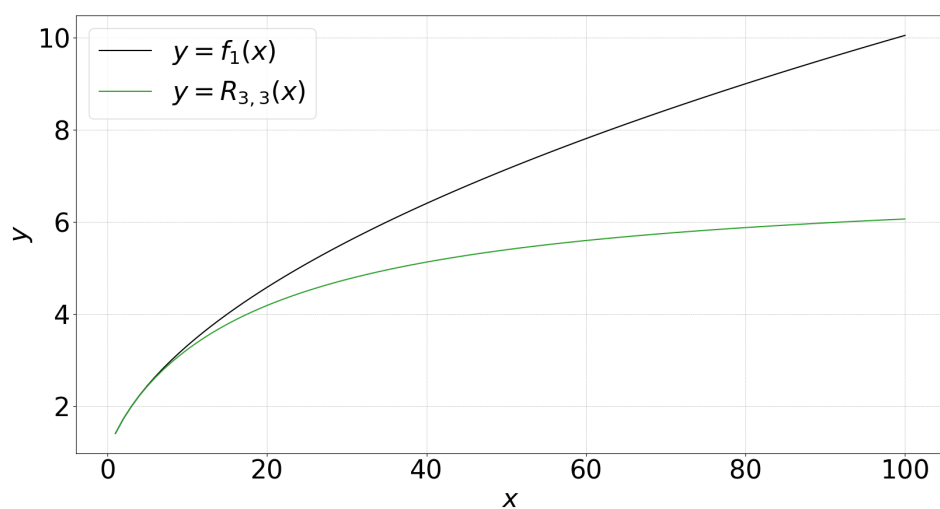


Figure 4: Plot of power series estimate of $f_1(x)$ for $N = 3$



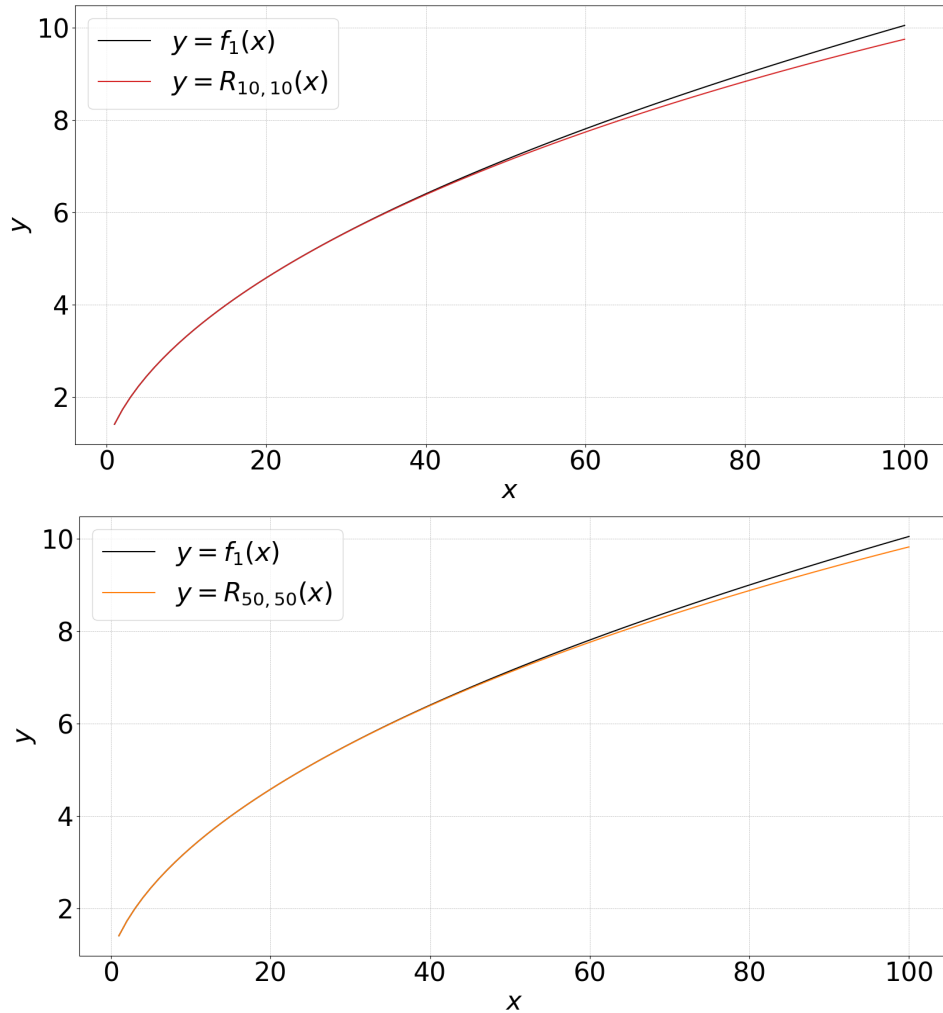


Figure 5: Plots of Padé approximant estimates of $f_1(x)$ for $L = 3, 10$ and 50

We can see from Figure 4 that for $N = 3$, the estimate increases at a rate of x^N . This means it diverges from $f_1(x)$ and for larger N the estimate diverges even more quickly. This is because of the x^N term in the power series which becomes very large for $x > 1$.

In comparison, the diagonal Padé approximant with $L = 3$ stays much closer to $f_1(x)$ than the power series estimate. This is due to the fact that the Padé approximant is a fraction so its limiting behaviour as $x \rightarrow \infty$ is much more similar to $f_1(x)$ than the power series' behaviour is. However, $L = 3$ does not give a good estimate in the range $1 \leq x \leq 100$. We see that for $L = 10$ and $L = 50$, we obtain much closer estimates while the power series would only diverge further.

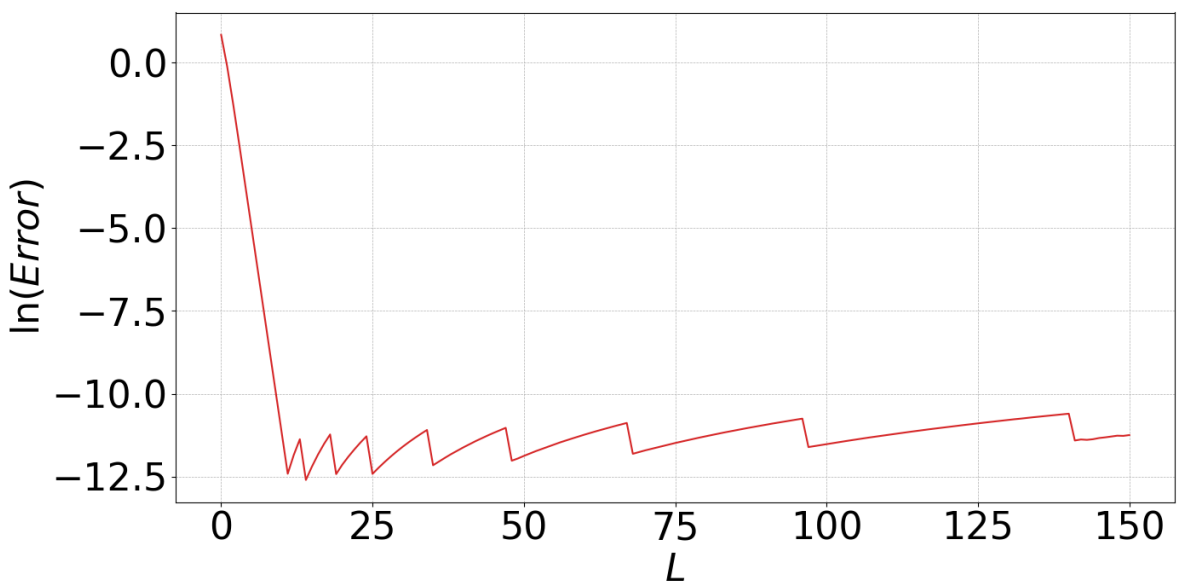


Figure 6: Plot of L against $\ln(\text{Error})$ for $x = 10$

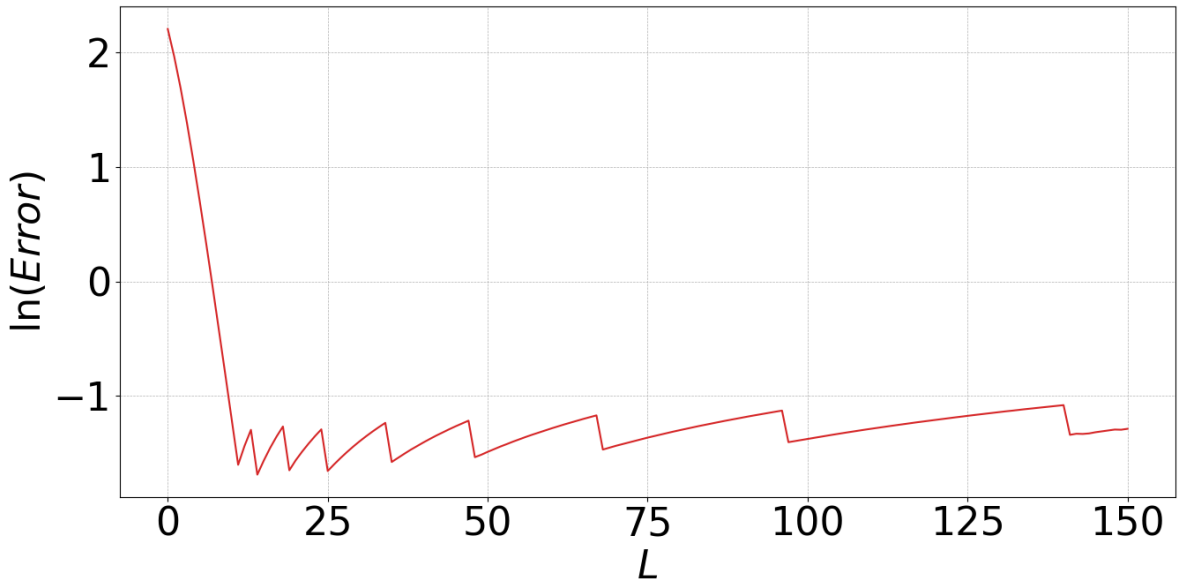


Figure 7: Plot of L against $\ln(\text{Error})$ for $x = 100$

From Figures 6 and 7, we see that the error decreases exponentially for $L \leq 11$ since the graphs are linear for this range of L . Beyond this range, the error remains in the same area as it increases and decreases in a zig-zag pattern. The reason the error stops decreasing exponentially is that the Padé approximant calculated by Program A is not accurate for large L (i.e. the q_k and p_k calculated are slightly off). I know that the approximant is inaccurate since it should give the same result as the $(2L)^{\text{th}}$ continued fraction from:

$$\begin{aligned} \sqrt{1+x} &= 1 + \frac{x}{1 + \sqrt{1+x}} \\ \Rightarrow \sqrt{1+x} &= 1 + \frac{x}{2 + \frac{x}{2 + \dots}} \end{aligned}$$

Truncating this fraction after the $(2L)^{\text{th}}$ two and simplifying, we obtain the diagonal Padé approximant. The program `continued_fraction.py` on page ?? demonstrates that the results of Program A and what the approximant should be (using the continued fraction) are different; using $L = 15$ as an example, the error from the Padé approximant is 5.03×10^{-6} while the error should be 2.77×10^{-8} . This must be due to the limitations of the 64-bit float arithmetic used for Program A and hence explains why the error stops decreasing as L is increased.

Overall, we can see that the Padé approximant almost converges to $f_1(x)$ for large x despite the fact that this is outside the radius of convergence of the power series. However, there is a limitation on the accuracy of the estimate you can get where increasing the value of L will not improve the result.

Question 4

I first present some graphs which give the error for different L using a diagonal approximant and the error for different orders using the power series. This is so that the optimal such L and order can be found for calculation across the whole range $0 \leq x \leq 20$.

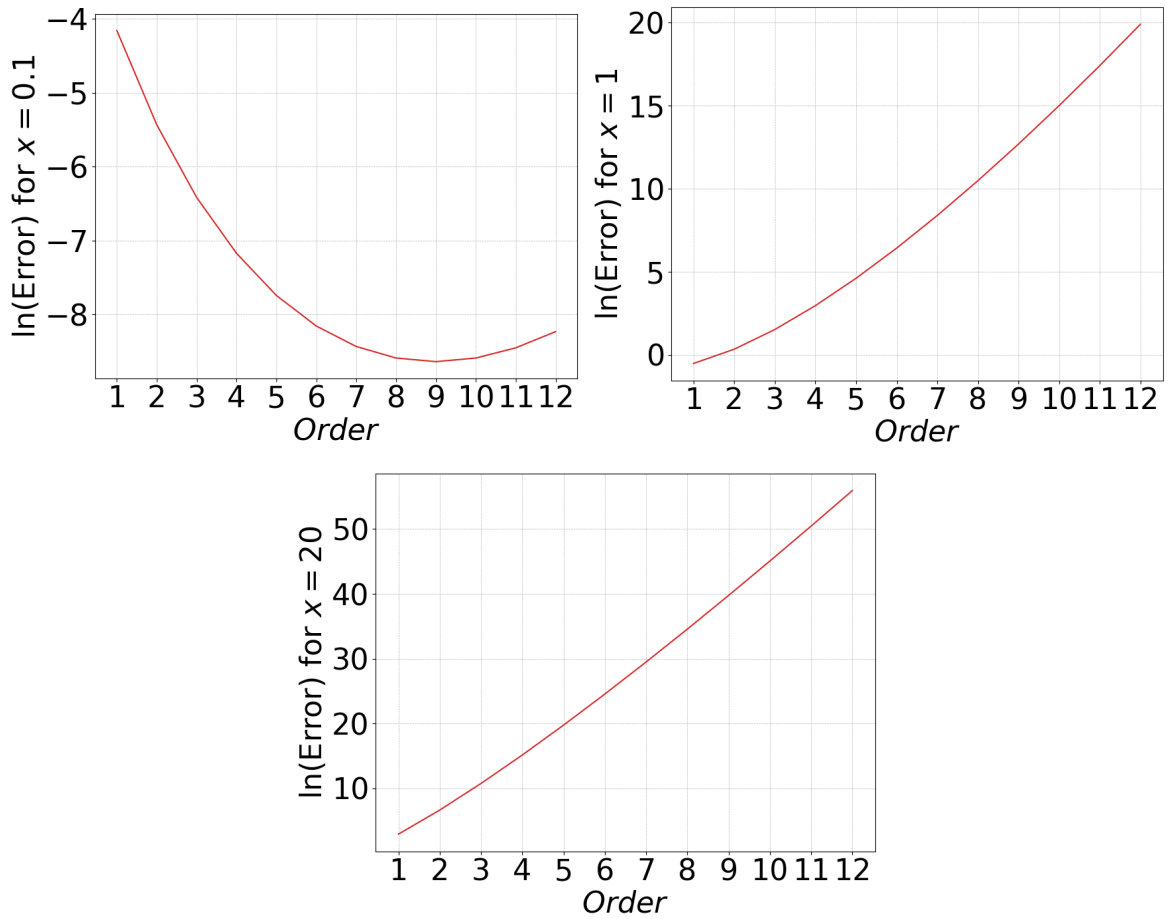


Figure 8: Plots of power series order against $\ln(\text{Error})$ for $x = 0.1, 1$ and 20

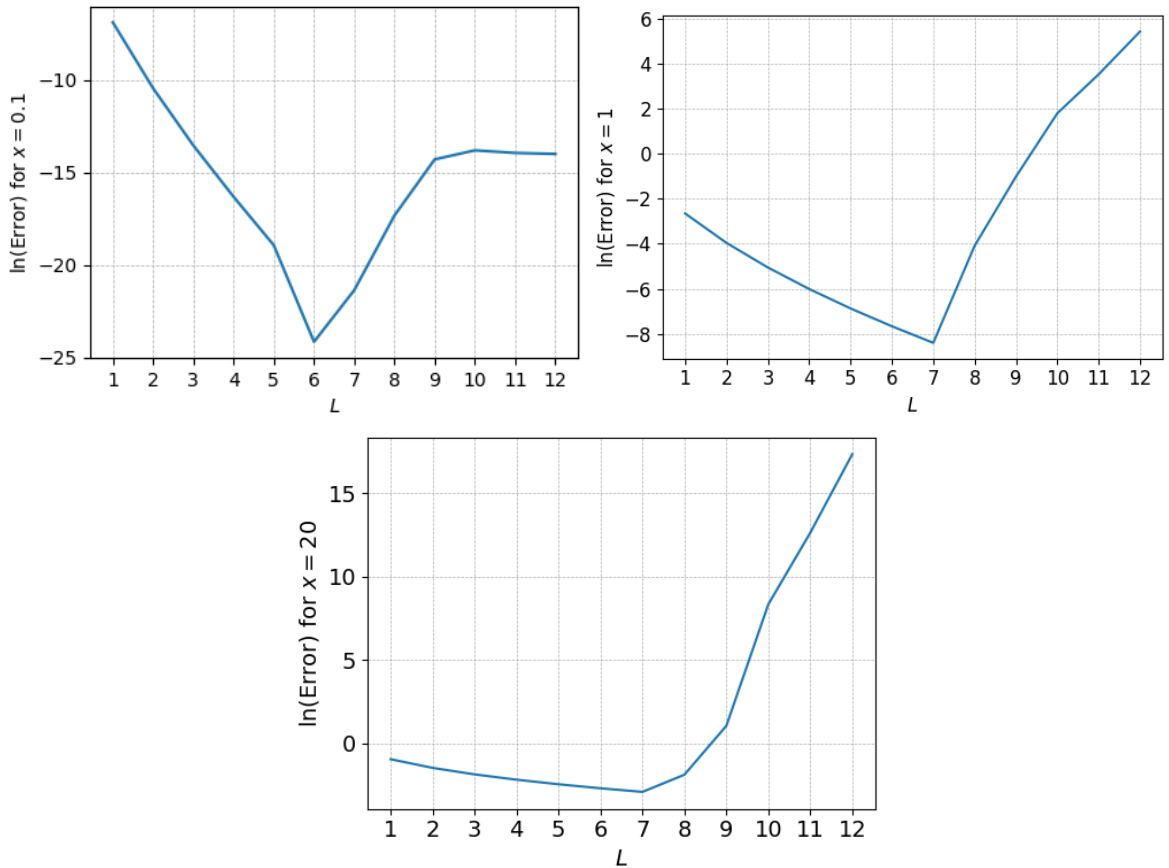


Figure 9: Plots of L against $\ln(\text{Error})$ for $x = 0.1, 1$ and 20

From Figure 8 we can see that for large $x \geq 1$ the error of the power series estimates diverges exponentially so the power series is not useful for any order here. On the other hand, the graph with $x = 0.1$ demonstrates that the series gives a good estimate for much smaller x as the x^n terms do not diverge. In particular, the order of 9 gives the minimum error.

Figure 9 shows that the Padé approximant can give a relatively small error for different

values of x in the range $[0, 20]$. The minimum error is given by either $L = 6$ or $L = 7$ and taking the diagonal approximant.

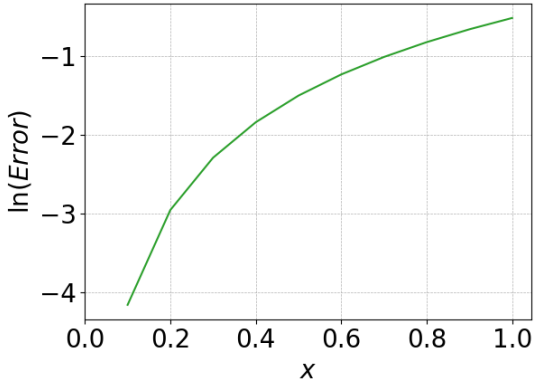


Figure 10: Error from the order 1 power series

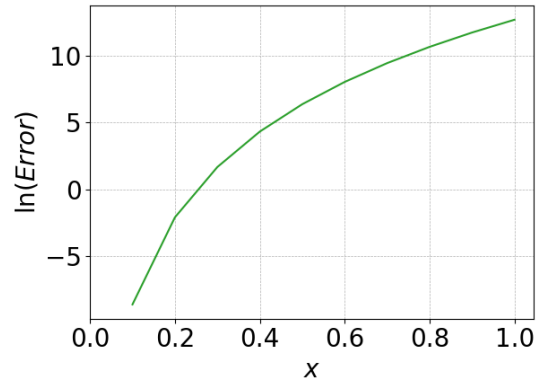


Figure 11: Error from the order 9 power series

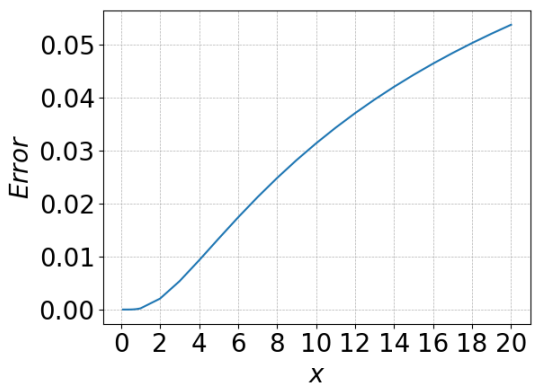


Figure 12: Error of Padé approximant with $L = 7$ in range the $[0, 20]$

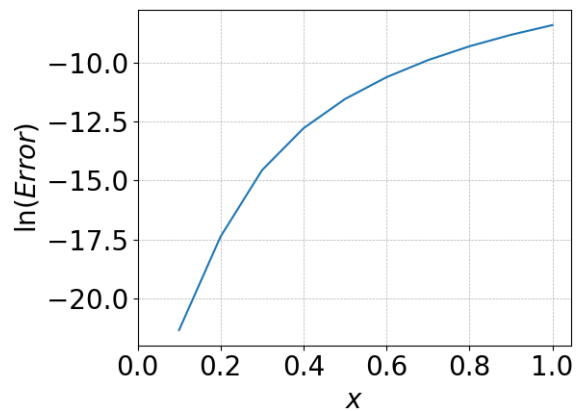


Figure 13: Error of Padé approximant with $L = 7$ zoomed in on the range $[0, 1]$

Since the power series diverges for larger x , only the Padé approximant is a good basis for calculating $f_2(x)$ in the range $0 \leq x \leq 20$. Figure 12 shows that the error remains small up to $x = 20$. In the range $0 \leq x \leq 1$, the Padé approximant also gives far more accurate calculations than the power series comparing Figures 10, 11 and 13. In addition, while the order 9 power series gives a small error for $x = 0.1$, this quickly diverges.

The reason that the Padé approximant is more accurate than the truncated power series is similar to the reasoning for $f_1(x)$; the limiting behaviour is better because the approximant is a fraction and for small x the power series expansion of the Padé approximant is better than the truncated power series with the matching first $L + M + 1$ terms.

Question 5

Poles and zeros of $f_1(x)$

I will display the values of x at the poles and zeros of the Padé approximant of $f_1(x)$ for selected L . Here, I will only consider $L \leq 11$ as we explore larger L for $f_1(x)$ in depth when discussing ‘anomalous’ poles.

Poles,

$$L = 1 : \quad x \approx -4$$

$$L = 3 : \quad x \approx -20.196, -2.572, -1.232$$

$$L = 10 : \quad x \approx -179.079, -20.197, \dots, -1.095, -1.023$$

Zeros,

$$L = 1 : \quad x \approx -1.333$$

$L = 3 : x \approx -5.312, -1.636, -1.052$
 $L = 10 : x \approx -45.021, -11.511, \dots, -1.052, -1.006$

All poles and zeros are real. Additionally, there are exactly L poles and zeros for $R_{L,L}(x)$ meaning that there is one more pole and zero when L is increased by one. All poles and zeros are negative, with the least negative values approaching 1 as L is increased. Meanwhile, the largest negative increases as L is increased and the remaining poles/zeros lie between this value and 1.

Poles and zeros of $f_3(x)$

The Padé approximant of $f_3(x)$ is the inverse of the Padé approximant of $f_1(x)$. Therefore, the poles of this approximant are the zeros of the approximant of $f_1(x)$ and vice versa. Hence the results above can be used to describe their behaviour.

Poles and zeros of $f_4(x)$

We have the following results for $f_4(x)$.

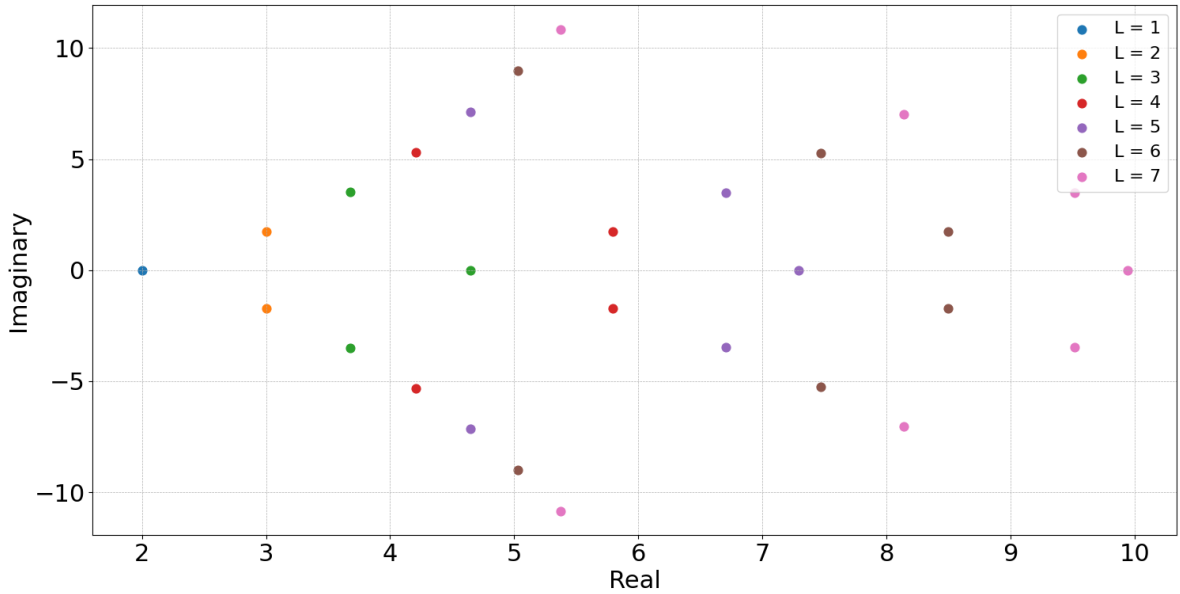


Figure 14: Graph showing poles of $R_{L,L}(x)$ for $f_4(x)$

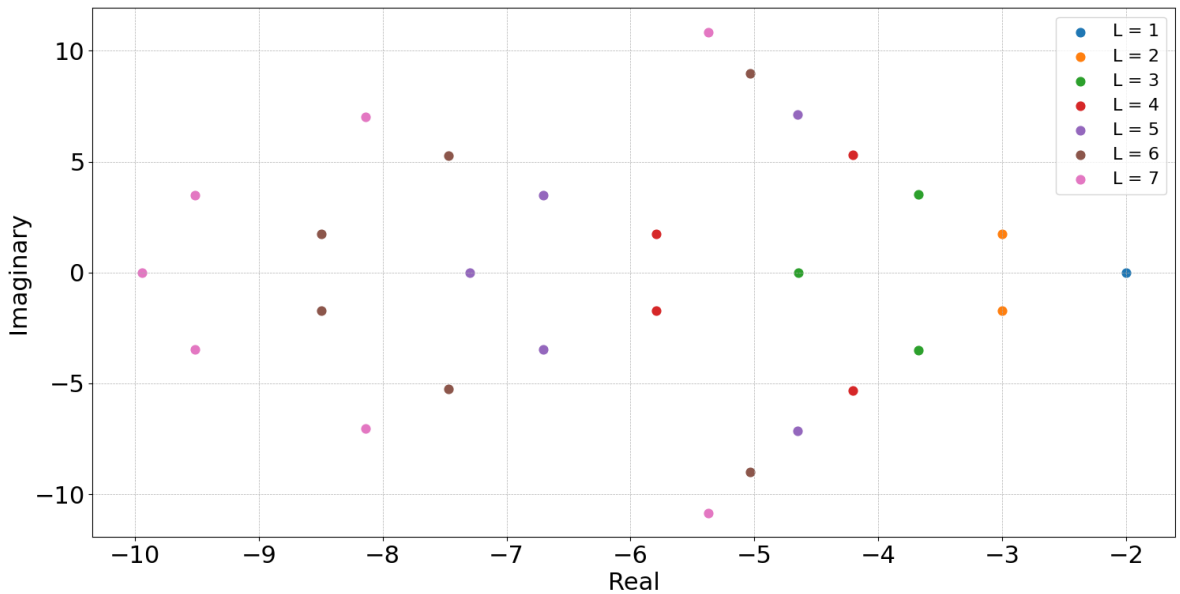


Figure 15: Graph showing zeros of $R_{L,L}(x)$ for $f_4(x)$

From Figures 14 and 15, we notice that there are mostly complex roots in conjugate pairs with the magnitude increasing as L increases. There are no real poles or zeros

for $L = 2, 4$ or 6 however there is exactly one for all odd L excluding 1, and exactly two for the remaining even $L \geq 8$.

Poles and zeros of $f_5(x)$

We can obtain similar graphs for $f_5(x)$ as for $f_4(x)$. The number of real poles and zeros remains the same as above apart from the fact that there are two real poles for $L = 4$ and $L = 6$. When there are real poles, there is always one very close to -1 and for even values of L there is an additional positive real pole. For instance, for $L = 14$ there are poles -1.00 and 1.48. As L increases above 10, additional poles and zeros are within a distance of 2 to the origin.

Poles and zeros of $f_6(x)$

For $f_6(x)$, there is a real pole for $L \geq 2$ with increasingly large magnitude. For example, for $L = 9$ there is a pole at $x = 7390.11$ and for $L = 10$ there is one at $x = -22133$. This large real pole alternates between positive and negative values for odd and even values of L respectively and seems to increase in magnitude by a factor of approximately 3 each time L is increased by 1. There is a second real pole for even L which is just less than -2. Then for $L > 16$, there can be even more real poles than the ones described. The remaining poles tend towards $-\frac{1}{2} \pm \frac{\sqrt{3}}{2}$ as outline in Figure 16.

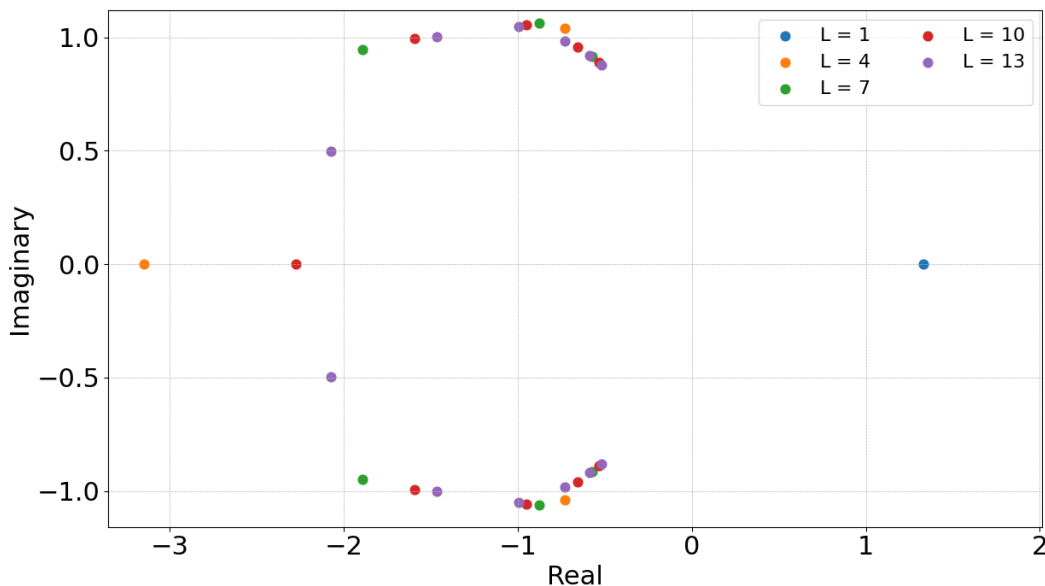


Figure 16: Graph showing poles of approximants of $f_6(x)$ close to the origin

The zeros of $f_6(x)$ are positioned in an identical fashion to Figure 16 with an increasing number of roots approaching $-\frac{1}{2} \pm \frac{\sqrt{3}}{2}$.

Correspondence to poles, zeros, branch points and branch cuts

The smallest in magnitude negative poles and zeros of the approximants of $f_1(x)$ and $f_3(x)$ tending to -1 corresponds to a branch point of both functions at $x = -1$. The remaining poles and zeros can be considered to be lying on a branch cut along the negative real axis.

$f_4(x)$ has no zeros, poles or branch points so there is no correspondence with the approximants. $f_5(x)$ has a pole at $z = -1$ which corresponds to the real poles we find very close to -1 in the approximants. Table 1 demonstrates this in more detail as it displays the value of the pole close to -1 for $L \leq 10$.

L	Pole
2	-1.0216272797495005
3	-0.9998189279229968
4	-1.0000008547034955
5	-0.999999997472825
6	1.0000000000050808
7	-0.9999999999999929
8	-0.9999999999999992
9	-0.9999999999999993
10	-0.9999999999999997

Table 1

We find that $-\frac{1}{2} \pm \frac{\sqrt{3}}{2}$ are zeros and branch points of $f_6(x)$ which explains why the zeros approach these values. ∞ is the third and final branch point so any branch cut requires two segments. Thus the poles lying along the real axis can be considered to be lying along a segment of a branch cut which goes from one of the finite branch points and then to ∞ along the real axis.

‘Anomalous’ poles and zeros

All of the poles and zeros of $f_4(x)$ shown in Figures 14 and 15 are anomalous. The same is true of $f_5(x)$, excluding the real pole close to -1.

For $f_1(x)$ we have anomalous poles and zeros when these values do not lie along the branch cut of the negative real axis, i.e. they are positive real or complex. The first time we get such an anomalous result is for $L = 12$ where there is a pole at $x = 1.126575649706596$ and a zero at $x = 1.1265756497065988$. Next, I will present figures for results of the poles but the exact same kind of results are obtained examining the zeros. Figure 17 gives an overview of where the poles/zeros are located with the different colours corresponding to different values of L .

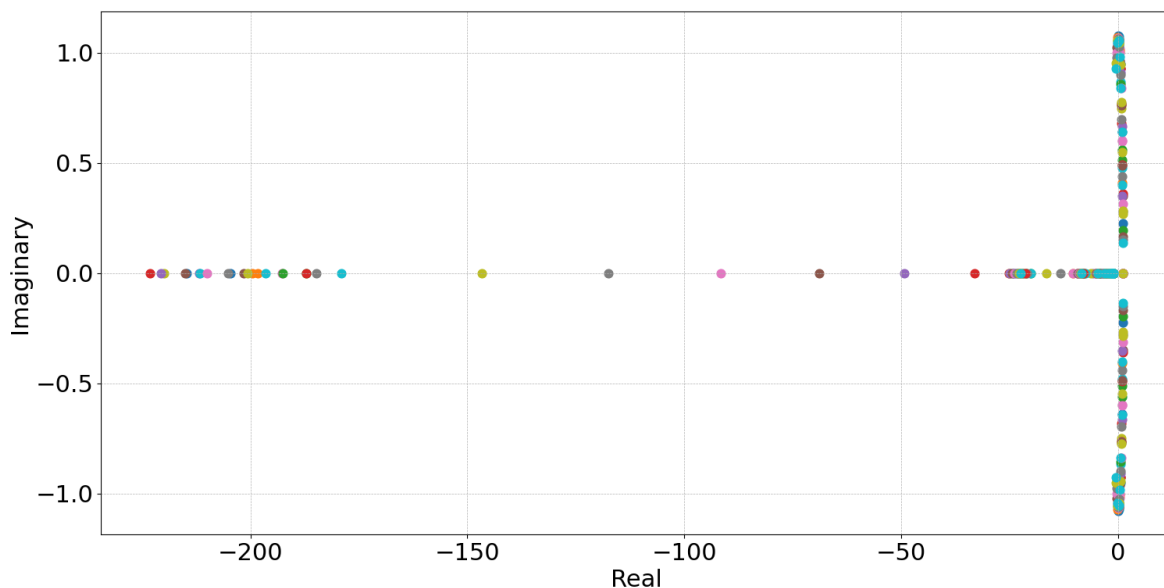


Figure 17: Graph showing poles of approximants of $f_1(x)$ for $L \leq 30$

We then focus on the anomalous poles which have a real part near 0 as in Figure 17.

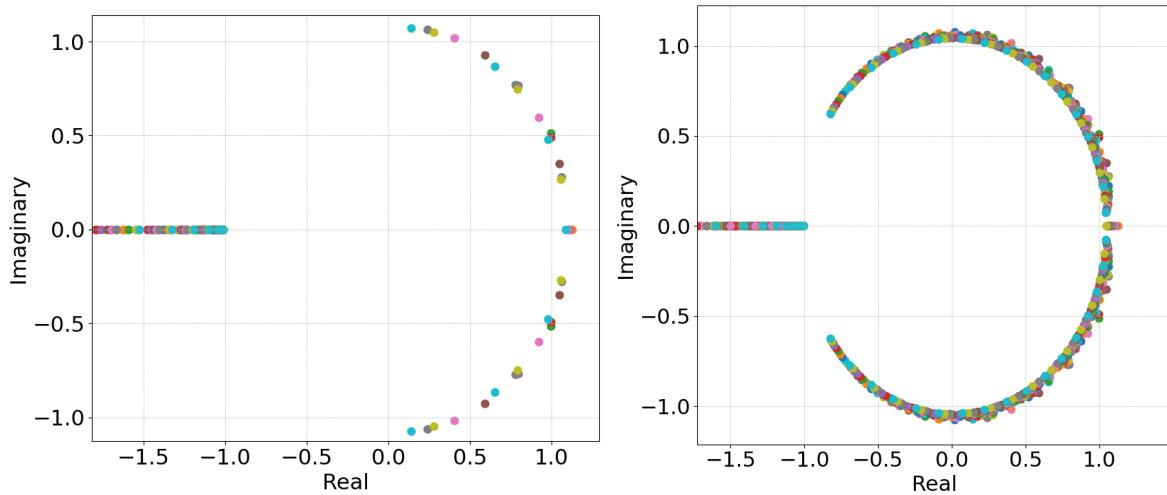


Figure 18: Graphs showing poles of approximants of $f_1(x)$ for $L \leq 20$ (left) and $L \leq 50$ (right)

We can see from Figure 18 that for large L we get close to a unit circle of poles/zeros for the diagonal Padé approximant of $f_1(x)$. As L becomes larger, this circle gets closer to -1 and there are always points close to 1 .

Problems with approximants to f_6

A major issue that may be encountered using Padé approximants to estimate $f_6(x)$ along the real x -axis is that the approximant is incorrect for most negative reals. For instance, in Figure 19 we see that for two different diagonal approximants the approximants diverge from the actual value for $x < -1.5$.

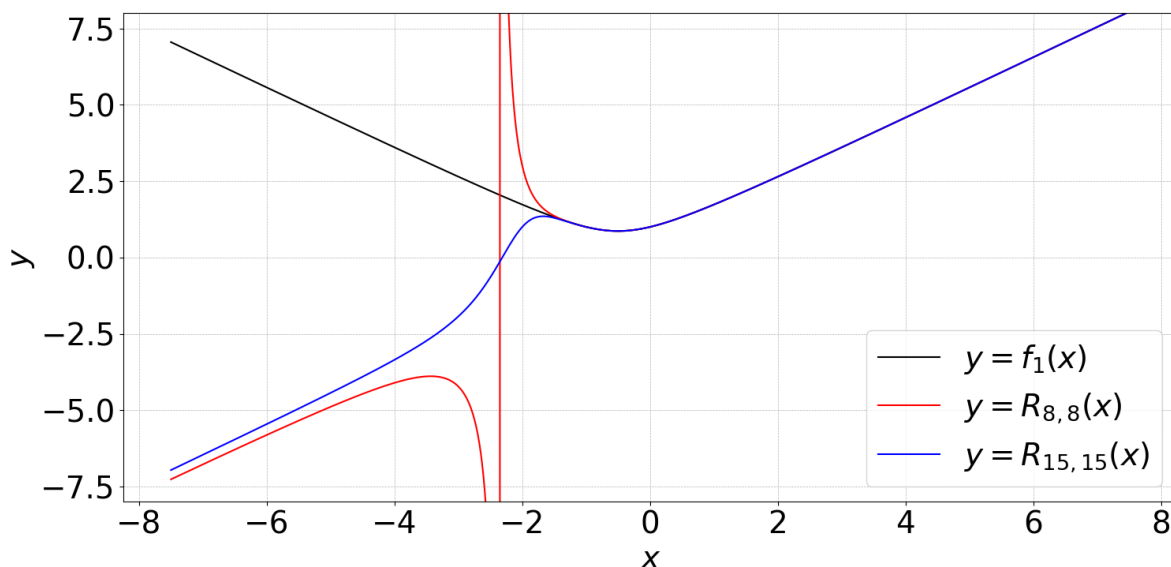


Figure 19: Graph with $f_6(x)$ and two different approximants

We see that for negative reals beyond about $x = -3$, the approximant is the negative of what it should be. For odd L this is because the leading coefficients have opposite signs e.g. for $L = 15$ we have ...

For the even case we also have negative result because ... Explain with leading coefficient for odd (pole for even?)