

## Rational approximations

What is the idea? We are going to build an approximation that is written as the quotient of polynomials; i.e.  $f(x) \approx \frac{p(x)}{q(x)}$  where  $p$  and  $q$  are polynomials.

## Padé approximations

The general idea is that we know from calculus how to create a Taylor expansion but we also know computing with a Taylor expansion is not a good plan. Instead we will write an rational approximation and match terms.

In other words, we will match the following

$$\frac{a_0 + a_1x + \cdots + a_mx^m}{1 + b_1x + \cdots + b_nx^n} = \mathbb{T}_{m+n}$$

where  $\mathbb{T}_{m+n}$  is the Taylor polynomial off degree  $m + n$ .

Note: in the left hand side there are  $n + m + 1$  unknown coefficients and the right hand side has  $m + n + 1$  known coefficients.

While there are the same number of terms in the two expansions, the rational approximation is more stable and thus more powerful.

For example, a polynomial will always diverge for large arguments and it cannot be nearly constant in any subinterval. Rational approximations can do have either of these problems.

**Example:** Consider the function  $\frac{1}{1+x^2}$ . It is nearly constant outside of an interval around the origin so a Taylor polynomial is going to struggle to approximate this function away from  $x = 0$ .

There are several types of rational approximations.

- Padé = easy to create and a powerful generalization of Taylor polynomials.
- Continued fractions
- Optimal rational

The latter two are very accurate but not very easy to construct.

In this class, we will focus on Padé approximations.

## How do we create a Padé approximation?

1. Decide on the orders for the polynomials in the rational approximation.
2. Write the approximation.

$$P_M^N(x) = \frac{\sum_{n=0}^N a_n x^n}{\sum_{m=0}^M b_m x^m}$$

where we set  $b_0 = 1$ .

3. Write the Taylor polynomial of order  $M + N$ .

$$T_{M+N}(x) = \sum_{n=0}^{M+N} c_n x^n$$

This Taylor polynomial is centered at 0. We know how to write the Taylor polynomial from calculus.

4. Set the  $P_M^N(x) = T_{M+N}$  and solve for the coefficients  $a_l$  and  $b_k$ . i.e.

$$c_0 + c_1x + \cdots + c_{N+M}x^{M+N} = \frac{a_0 + a_1x + \cdots + a_Nx^N}{1 + b_1x + \cdots + b_Mx^M}$$

(The  $c_j$  for  $j = 0, \dots, M + N$  are known.) Multiply both sides by  $1 + b_1x + \cdots + b_Mx^M$  and match terms.

You will end up having to write a linear system to solve for the constants  $b_k$ .

**Example:** Create  $P_3^2(x)$  for approximating  $f(x) = e^{-x}$ .

**Soln:** We know that  $P_3^2(x)$  will have the form

$$P_3^2(x) = \frac{a_0 + a_1x + a_2x^2 + a_3x^3}{1 + b_1x + b_2x^2}.$$

This is roughly an  $O(x^{3/2})$  approximation.

We also know that we should match this with

$$T_5(x) = 1 - x + \frac{x^2}{2} - \frac{x^3}{6} + \frac{x^4}{24} - \frac{x^5}{120}.$$

Setting these equal we find

$$\frac{a_0 + a_1x + a_2x^2 + a_3x^3}{1 + b_1x + b_2x^2} = 1 - x + \frac{x^2}{2} - \frac{x^3}{6} + \frac{x^4}{24} - \frac{x^5}{120}$$

Multiplying the polynomial with  $b$  coefficients on both sides results in the following equation.

$$a_0 + a_1x + a_2x^2 + a_3x^3 = (1 - x + \frac{x^2}{2} - \frac{x^3}{6} + \frac{x^4}{24} - \frac{x^5}{120})(1 + b_1x + b_2x^2)$$

Now we will build a table to match coefficients.

Term	equation for coefficient
constant	$a_0 = 1$
$x$	$a_1 = b_1 - 1$
$x^2$	$a_2 = 1/2 - b_1 + b_2$
$x^3$	$a_3 = -\frac{1}{6} + \frac{b_1}{2} - b_2$
$x^4$	$0 = \frac{1}{24} - \frac{1}{6}b_1 + \frac{1}{2}b_2$
$x^5$	$0 = -\frac{1}{120} + \frac{1}{24}b_1 - \frac{1}{6}b_2$

The last two equations are independent of  $a_j$  for  $j = 0, 3$  so we can solve them to find  $b_1$  and  $b_2$ .

$$\begin{bmatrix} -\frac{1}{6} & \frac{1}{2} \\ \frac{1}{24} & -\frac{1}{6} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} -\frac{1}{24} \\ \frac{1}{120} \end{bmatrix}$$

Solving this we find  $b_1 = \frac{2}{5}$  and  $b_2 = \frac{1}{20}$ . Plugging these into the equations for  $a_1$ ,  $a_2$  and  $a_3$ , we find  $a_1 = -\frac{3}{5}$  and  $a_2 = \frac{3}{20}$  and  $a_3 = -\frac{1}{60}$ .

Thus the rational approximation is

$$P_2^3(x) = \frac{1 - \frac{3}{5}x + \frac{3}{20}x^2 - \frac{1}{60}x^3}{1 + \frac{2}{5}x + \frac{1}{20}x^2}.$$

To investigate the performance of this technique, let's plot the error in the two approximation techniques. (see code) Also compare with interpolation.

## Fourier series

The idea is to approximate functions as sum of sine and cosines on  $[-\pi, \pi]$ .

In other words, we seek to find  $\{a_k\}_{k=0}^{\infty}$  and  $\{b_k\}_{k=1}^{\infty}$  such that

$$f(x) = a_0 + \sum_{k=1}^{\infty} a_k \cos(kx) + \sum_{k=1}^{\infty} b_k \sin(kx).$$

We know  $\{1, \sin(x), \dots, \sin(nx), \dots, \cos(x), \dots, \cos(nx), \dots\}$  are linearly independent on  $[-\pi, \pi]$  and that they are orthogonal. i.e.

$$\int_{-\pi}^{\pi} \cos(kx) \cos(jx) dx = 0$$

for  $j \neq k$ .

This is great. We need only take inner products to find the constants  $a_k$  and  $b_k$ .

Let  $l \neq 0$ , then

$$\begin{aligned} \int_{-\pi}^{\pi} \cos(lx) f(x) dx &= \int_{-\pi}^{\pi} \left( a_0 + \sum_{k=1}^{\infty} a_k \cos(kx) + \sum_{k=1}^{\infty} b_k \sin(kx) \right) dx \\ &= \int_{-\pi}^{\pi} a_0 \cos(lx) dx + \sum_{k=1}^{\infty} a_k \int_{-\pi}^{\pi} \cos(kx) \cos(lx) dx + \sum_{k=1}^{\infty} b_k \int_{-\pi}^{\pi} \sin(kx) \cos(lx) dx \\ &= a_l \int_{-\pi}^{\pi} \cos^2(lx) dx \\ &= a_l \int_{-\pi}^{\pi} \frac{1}{2} (1 + \cos(2lx)) dx \\ &= a_l \pi \end{aligned}$$

Thus

$$a_l = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(lx) dx$$

Likewise you can find formulas for the other coefficients;

$$a_0 = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) dx$$

and

$$b_l = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(lx) dx.$$

Note: This matches doing least squares approximation with trig. polynomials.  
There are some possible problems.

- adding infinitely many things is not possible in practice
- there are limitations on the function  $f(x)$  for which this will work.
- We need to use quadrature (most of the time) to evaluate the coefficients. (We will talk about quadrature next week.)

In practice we will use a partial sum.

$$S_N(x) = a_0 + \sum_{k=1}^N a_k \cos(kx) + \sum_{k=1}^N b_k \sin(kx).$$

The reasonable question is: does  $S_N(x) \rightarrow f(x)$  as  $N \rightarrow \infty$ ?  
This depends on the function  $f(x)$ .

Up until 1876 people thought it was enough for  $f$  to be continuous.

This was not enough. We need  $f$  to be continuous and piecewise differentiable on  $[-\pi, \pi]$ .

Why? What happens if the derivative blows up at a point in  $[-\pi, \pi]$ . Then we are approximating something that is continuous but not infinitely differentiable by something that is smooth and infinitely differentiable.

**Example:** Evaluate the Fourier series of  $f(x) = |x|$  on  $[-\pi, \pi]$ .

**Soln:** The coefficients are

$$a_0 = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) dx = \frac{\pi}{2}$$
$$a_k = \frac{1}{\pi} \int_{-\pi}^{\pi} |x| \cos(kx) dx = \frac{2}{\pi} \int_0^{\pi} x \cos(kx) dx = \frac{2}{\pi k^2} \left( (-1)^k - 1 \right)$$

and

$$b_k = \frac{1}{\pi} \int_{-\pi}^{\pi} |x| \sin(kx) dx = 0$$

since  $|x| \sin(kx)$  is odd on  $[-\pi, \pi]$ .

Thus

$$f(x) = \frac{\pi}{2} + \frac{2}{\pi} \sum_{k=1}^{\infty} \frac{1}{k^2} \left( (-1)^k - 1 \right) \cos(kx)$$

Look at the numerical example to explore convergence.

**Example:** (Do on your own) Evaluate the Fourier series of

$$g(x) = \begin{cases} -1/2 & \text{for } x \in (-\pi, 0) \\ 1/2 & \text{for } x \in [0, \pi) \end{cases}$$

**Soln:** Coefficients are in the code. Since  $g$  is odd,  $a_k = 0$  for all  $k$ . Thus the Fourier series is a sine series. What is happening?

## Gibbs' phenomena

### Brief history

- It was first noted in 1848 by Wilbraham.
- Michelson and Stratten found traces of the overshoots in 1898.
- Michelson wrote a note to Nature asking if anyone understood the convergence problems of Fourier series.
- Gibbs (Chemist) provided the answer. His first answer was flawed but then he corrected it.

How big is the error due to Gibbs?

In the max norm...

Near the jump (discontinuity), the error is  $O(1)$ .

Away from the jump, the error is  $O(1/N)$