

# A. Introduction to Julia

References: [The Julia Documentation](#), [The Julia–Matlab–Python Cheatsheet](#), [Think Julia](#)

These notes give an overview of Julia. In these notes we focus on the aspects of Julia and computing that are essential to numerical computing:

1. Integers: We discuss briefly how to create and manipulate integers, and how to see the underlying bit representation.
2. Strings and parsing: We discuss how to create and manipulate strings and characters, and how we can convert a string

of 0's and 1's to an integer or other type. 3. Vectors and matrices: We discuss how to build and manipulate vectors and matrices (which are both types of *arrays*). Later lectures will discuss linear algebra. 4. Types: In Julia everything has a type, which plays a similar role to classes in Python. Here we discuss how to make new types, for example, a complex number in radial format. 5. Loops and branches: We discuss `if`, `for` and `while`, which work similar to Python. 6. Functions: We discuss the construction of named and anonymous functions. Julia allows overloading functions for different types, for example, we can overload `*` for our radial complex type. 7. Modules, Packages, and Plotting: We discuss how to load external packages, in particular, for plotting.

## 1. Integers

Julia uses a math-like syntax for manipulating integers:

```
In [1]: 1 + 1 # Addition
```

```
Out[1]: 2
```

```
In [2]: 2 * 3 # Multiplication
```

```
Out[2]: 6
```

```
In [3]: 2 / 3 # Division
```

```
Out[3]: 0.6666666666666666
```

```
In [4]: x = 5; # semicolon is optional but supresses output if used in the last line  
x^2 # Powers
```

```
Out[4]: 25
```

In Julia everything has a type. This is similar in spirit to a class in Python, but much more lightweight. An integer defaults to a type `Int`, which is either 32-bit (`Int32`) or 64-

bit (`Int64`) depending on the processor of the machine. There are also 8-bit (`Int8`), 16-bit (`Int16`), and 128-bit (`Int128`) integer types, which we can construct by converting an `Int`, e.g. `Int8(3)`.

These are all "primitive types", instances of the type are stored in memory as a fixed length sequence of bits. We can find the type of a variable as follows:

```
In [5]: typeof(x)
```

```
Out[5]: Int64
```

For a primitive type we can see the bits using the function `bitstring`:

```
In [6]: bitstring(Int8(1))
```

```
Out[6]: "00000001"
```

Negative numbers may be surprising:

```
In [7]: bitstring(-Int8(1))
```

```
Out[7]: "11111111"
```

This is explained in detail in Chapter [Numbers](#).

There are other primitive integer types: `UInt8`, `UInt16`, `UInt32`, and `UInt64` are unsigned integers, e.g., we do not interpret the number as negative if the first bit is `1`. As they tend to be used to represent bit sequences they are displayed in hexadecimal, that is base-16, using digits 0-9a-c, e.g.,  $12 = (c)_{16}$ :

```
In [8]: UInt16(12)
```

```
Out[8]: 0x000c
```

A non-primitive type is `BigInt` which allows arbitrary length integers (using an arbitrary amount of memory):

```
In [9]: factorial(big(100))^10
```

## 2. Strings and parsing

We have seen that `bitstring` returns a string of bits. Strings can be created with quotation marks

```
In [10]: str = "hello world 😊"
```

Out[10]: "hello world 😊"

We can access characters of a string with brackets:

```
In [11]: str[1], str[13]
```

```
Out[11]: ('h', '😊')
```

Each character is a primitive type, in this case using 32 bits/4 bytes:

```
In [12]: typeof(str[6]), length(bitstring(str[6]))
```

Out[12]: (Char, 32)

Strings are not primitive types, but rather point to the start of a sequence of `Char`s in memory. In this case, there are  $32 * 13 = 416$  bits/52 bytes in memory.

Strings are *immutable*: once created they cannot be changed. But a new string can be created that modifies an existing string. The simplest example is `*`, which concatenates two strings:

```
In [13]: "hi" * "bye"
```

```
Out[13]: "hibye"
```

(Why `*`? Because concatenation is non-commutive.) We can combine this with indexing to, for example, create a new string with a different last character:

```
In [14]: str[1:end-1] * "😊"
```

```
Out[14]: "hello world 😊"
```

## Parsing strings

We can use the command `parse` to turn a string into an integer:

```
In [15]: parse(Int, "123")
```

```
Out[15]: 123
```

We can specify base 2 as an optional argument:

```
In [16]: parse(Int, "-101"; base=2)
```

```
Out[16]: -5
```

If we are specifying bits its safer to parse as an `UInt32`, otherwise the first bit is not recognised as a sign:

```
In [17]: bts = "1111000010011111001100110001010"
x = parse(UInt32, bts; base=2)
```

```
Out[17]: 0xf09f998a
```

The function `reinterpret` allows us to reinterpret the resulting sequence of 32 bits as a different type. For example, we can reinterpret as an `Int32` in which case the first bit is taken to be the sign bit and we get a negative number:

```
In [18]: reinterpret(Int32, x)
```

```
Out[18]: -257975926
```

We can also reinterpret as a `Char`:

```
In [19]: reinterpret(Char, x)
```

```
Out[19]: '😊': Unicode U+1F64A (category So: Symbol, other)
```

We will use `parse` and `reinterpret` as it allows one to easily manipulate bits. This is not actually how one should do it as it is slow.

## Bitwise operations (non-examinable)

In practice, one should manipulate bits using bitwise operations. These will not be required in this course and are not examinable, but are valuable to know if you have a career involving high performance computing. The `p << k` shifts the bits of `p` to the left `k` times inserting zeros, while `p >> k` shifts to the right:

```
In [20]: println(bitstring(23));  
        println(bitstring(23 << 2));  
        println(bitstring(23 >> 2));
```

The operations `&`, `|` and `^` do bitwise and, or, and xor.

## 3. Vectors, Matrices, and Arrays

We can create a vector using brackets:

```
In [21]: v = [11, 24, 32]
```

```
Out[21]: 3-element Vector{Int64}:
          11
          24
          32
```

Like a string, elements are accessed via brackets. Julia uses 1-based indexing (like Matlab and Mathematica, unlike Python and C which use 0-based indexing):

In [22]: v[1], v[3]

Out[22]: (11, 32)

Accessing outside the range gives an error:

In [23]: v[4]

BoundsError: attempt to access 3-element Vector{Int64} at index [4]

## Stacktrace:

```
[1] getindex(A::Vector{Int64}, i1::Int64)
    @ Base ./array.jl:924
[2] top-level scope
    @ In[23]:1
```

Vectors can be made with different types, for example, here is a vector of three 8-bit integers:

```
In [24]: v = [Int8(11), Int8(24), Int8(32)]
```

```
Out[24]: 3-element Vector{Int8}:
```

```
11  
24  
32
```

Just like strings, Vectors are not primitive types, but rather point to the start of sequence of bits in memory that are interpreted in the corresponding type. In this last case, there are  $3 * 8 = 24$  bits/3 bytes in memory.

The easiest way to create a vector is to use `zeros` to create a zero `Vector` and then modify its entries:

```
In [25]: v = zeros(Int, 5)  
v[2] = 3  
v
```

```
Out[25]: 5-element Vector{Int64}:
```

```
0  
3  
0  
0  
0
```

Note: we can't assign a non-integer floating point number to an integer vector:

```
In [26]: v[2] = 3.5
```

```
InexactError: Int64(3.5)
```

```
Stacktrace:
```

```
[1] Int64  
  @ ./float.jl:788 [inlined]  
[2] convert  
  @ ./number.jl:7 [inlined]  
[3] setindex!(A::Vector{Int64}, x::Float64, i1::Int64)  
  @ Base ./array.jl:966  
[4] top-level scope  
  @ In[26]:1
```

We can also create vectors with `ones` (a vector of all ones), `rand` (a vector of random numbers between `0` and `1`) and `randn` (a vector of samples of normal distributed quasi-random numbers).

When we create a vector whose entries are of different types, they are mapped to a type that can represent every entry. For example, here we input a list of one `Int32` followed by three `Int64`s, which are automatically converted to all be `Int64`:

```
In [27]: [Int32(1), 2, 3, 4]
```

```
Out[27]: 4-element Vector{Int64}:
```

```
1  
2  
3  
4
```

In the event that the types cannot automatically be converted, it defaults to an `Any` vector, which is similar to a Python list. This is bad performancewise as it does not know how many bits each element will need, so should be avoided.

```
In [28]: [1.0, 1, "1"]
```

```
Out[28]: 3-element Vector{Any}:
```

```
1.0  
1  
"1"
```

We can also specify the type of the Vector explicitly by writing the desired type before the first bracket:

```
In [29]: Int32[1, 2, 3]
```

```
Out[29]: 3-element Vector{Int32}:
```

```
1  
2  
3
```

We can also create an array using comprehensions:

```
In [30]: [k^2 for k = 1:5]
```

```
Out[30]: 5-element Vector{Int64}:
```

```
1  
4  
9  
16  
25
```

Matrices are created similar to vectors, but by specifying two dimensions instead of one. Again, the simplest way is to use `zeros` to create a matrix of all zeros:

```
In [31]: zeros(Int, 4, 5) # creates a 4 × 5 matrix of Int zeros
```

```
Out[31]: 4×5 Matrix{Int64}:
```

```
0 0 0 0 0  
0 0 0 0 0  
0 0 0 0 0  
0 0 0 0 0
```

We can also create matrices by hand. Here, spaces delimit the columns and semicolons delimit the rows:

```
In [32]: A = [1 2; 3 4; 5 6]
```

```
Out[32]: 3x2 Matrix{Int64}:
```

```
1 2  
3 4  
5 6
```

We can also create matrices using brackets, a formula, and a `for` command:

```
In [33]: [k^2+j for k=1:4, j=1:5]
```

```
Out[33]: 4x5 Matrix{Int64}:
```

```
2 3 4 5 6  
5 6 7 8 9  
10 11 12 13 14  
17 18 19 20 21
```

Matrices are really vectors in disguise. They are still stored in memory in a consecutive sequence of bits. We can see the underlying vector using the `vec` command:

```
In [34]: vec(A)
```

```
Out[34]: 6-element Vector{Int64}:
```

```
1  
3  
5  
2  
4  
6
```

The only difference between matrices and vectors from the computers perspective is that they have a `size` which changes the interpretation of what's stored in memory:

```
In [35]: size(A)
```

```
Out[35]: (3, 2)
```

Matrices can be manipulated easily on a computer. We can multiply a matrix times vector:

```
In [36]: x = [8; 9]  
A * x
```

```
Out[36]: 3-element Vector{Int64}:
```

```
26  
60  
94
```

or a matrix times matrix:

```
In [37]: A * [4 5; 6 7]
```

```
Out[37]: 3x2 Matrix{Int64}:
```

```
16 19  
36 43  
56 67
```

If you use `.*`, it does entrywise multiplication:

```
In [38]: [1 2; 3 4] .* [4 5; 6 7]
```

```
Out[38]: 2x2 Matrix{Int64}:
```

```
4 10  
18 28
```

We can take the transpose of a real vector as follows:

```
In [39]: a = [1, 2, 3]  
a'
```

```
Out[39]: 1x3 adjoint(::Vector{Int64}) with eltype Int64:
```

```
1 2 3
```

Note for complex-valued vectors this is the conjugate-transpose, and so one may need to use `transpose(a)`. Both `a'` and `transpose(a)` should be thought of as "dual-vectors", and so multiplication with a transposed vector with a normal vector gives a constant:

```
In [40]: b = [4, 5, 6]  
a' * b
```

```
Out[40]: 32
```

One important note: a vector is not the same as an `n × 1` matrix, and a transposed vector is not the same as a `1 × n` matrix.

## Accessing and altering subsections of arrays

We will use the following notation to get at the columns and rows of matrices:

```
A[a:b,k]      # returns the a-th through b-th rows of the k-th  
               column of A as a Vector of length (b-a+1)  
A[k,a:b]      # returns the a-th through b-th columns of the k-  
               th row of A as a Vector of length (b-a+1)  
A[:,k]         # returns all rows of the k-th column of A as a  
               Vector of length size(A,1)  
A[k,:]        # returns all columns of the k-th row of A as a  
               Vector of length size(A,2)  
A[a:b,c:d]    # returns the a-th through b-th rows and c-th  
               through d-th columns of A  
               # as a (b-a+1) x (d-c+1) Matrix
```

The ranges `a:b` and `c:d` can be replaced by any `AbstractVector{Int}`. For example:

```
In [41]: A = [1 2 3; 4 5 6; 7 8 9; 10 11 12]
A[[1,3,4],2] # returns the 1st, 3rd and 4th rows of the 2nd column of A
```

Out[41]: 3-element Vector{Int64}:

```
2
8
11
```

**Exercise** Can you guess what `A[2, [1,3,4]]` returns, using the definition of `A` as above?

What about `A[1:2, [1,3]]`? And `A[1,B[1:2,1]]`? And `vec(A[1,B[1:2,1]])`?

We can also use this notation to modify entries of the matrix. For example, we can set the `1:2 x 2:3` subblock of `A` to `[1 2; 3 4]` as follows:

```
In [42]: A[1:2,2:3] = [1 2; 3 4]
A
```

Out[42]: 4x3 Matrix{Int64}:

```
1 1 2
4 3 4
7 8 9
10 11 12
```

## Broadcasting

It often is necessary to apply a function to every entry of a vector. By adding `.` to the end of a function we "broadcast" the function over a vector:

```
In [43]: x = [1,2,3]
cos.(x) # equivalent to [cos(1), cos(2), cos(3)]
```

Out[43]: 3-element Vector{Float64}:

```
0.5403023058681398
-0.4161468365471424
-0.9899924966004454
```

Broadcasting has some interesting behaviour for matrices. If one dimension of a matrix (or vector) is 1, it automatically repeats the matrix (or vector) to match the size of another example.

### Example

```
In [44]: [1,2,3] .* [4,5]'
```

```
Out[44]: 3x2 Matrix{Int64}:
```

```
4 5  
8 10  
12 15
```

Since `size([1,2,3],2) == 1` it repeats the same vector to match the size `size([4,5]',2) == 2`. Similarly, `[4,5]'` is repeated 3 times. So the above is equivalent to:

```
In [45]: [1 1; 2 2; 3 3] .* [4 5; 4 5; 4 5]
```

```
Out[45]: 3x2 Matrix{Int64}:
```

```
4 5  
8 10  
12 15
```

Note we can also use broadcasting with our own functions (construction discussed later):

```
In [46]: f = (x,y) -> cos(x + 2y)  
f.([1,2,3], [4,5]')
```

```
Out[46]: 3x2 Matrix{Float64}:
```

```
-0.91113 0.0044257  
-0.839072 0.843854  
0.0044257 0.907447
```

## Ranges

We have already seen that we can represent a range of integers via `a:b`. Note we can convert it to a `Vector` as follows:

```
In [47]: Vector(2:6)
```

```
Out[47]: 5-element Vector{Int64}:
```

```
2  
3  
4  
5  
6
```

We can also specify a step:

```
In [48]: Vector(2:2:6), Vector(6:-1:2)
```

```
Out[48]: ([2, 4, 6], [6, 5, 4, 3, 2])
```

Finally, the `range` function gives more functionality, for example, we can create 4 evenly spaced points between `-1` and `1`:

```
In [49]: Vector(range(-1, 1; length=4))
```

```
Out[49]: 4-element Vector{Float64}:
```

```
-1.0  
-0.333333333333333  
0.333333333333333  
1.0
```

Note that `Vector` is mutable but a range is not:

```
In [50]: r = 2:6  
r[2] = 3 # Not allowed
```

```
CanonicalIndexError: setindex! not defined for UnitRange{Int64}
```

Stacktrace:

```
[1] error_if_canonical_setindex(#unused#::IndexLinear, A::UnitRange{Int64}, #unused#::Int64)  
    @ Base ./abstractarray.jl:1352  
[2] setindex!(A::UnitRange{Int64}, v::Int64, I::Int64)  
    @ Base ./abstractarray.jl:1343  
[3] top-level scope  
    @ In[50]:2
```

## 4. Types

Julia has two different kinds of types: primitive types (like `Int64`, `Int32`, `UInt32` and `Char`) and composite types. Here is an example of an in-built composite type representing complex numbers, for example,  $z = 1 + 2i$ :

```
In [51]: z = 1 + 2im  
typeof(z)
```

```
Out[51]: Complex{Int64}
```

A complex number consists of two fields: a real part (denoted `re`) and an imaginary part (denoted `im`). Fields of a type can be accessed using the `.` notation:

```
In [52]: z.re, z.im
```

```
Out[52]: (1, 2)
```

We can make our own types. Let's make a type to represent complex numbers in the format

$$z = r \exp(i\theta)$$

That is, we want to create a type with two fields: `r` and `θ`. This is done using the `struct` syntax, followed by a list of names for the fields, and finally the keyword `end`.

```
In [53]: struct RadialComplex  
    r  
    θ
```

```
end  
z = RadialComplex(1,0.1)
```

Out[53]: RadialComplex(1, 0.1)

We can access the fields using `.`:

```
In [54]: z.r, z.θ
```

Out[54]: (1, 0.1)

Note that the fields are immutable: we can create a new `RadialComplex` but we cannot modify an existing one. To make a mutable type we use the command `mutable struct`:

```
In [55]: mutable struct MutableRadialComplex  
    r  
    θ  
end  
  
z = MutableRadialComplex(1,2)  
z.r = 2  
z.θ = 3  
z
```

Out[55]: MutableRadialComplex(2, 3)

## Abstract types

Every type is a sub-type of an *abstract type*, which can never be instantiated on its own. For example, every integer and floating point number is a real number.

Therefore, there is an abstract type `Real`, which encapsulates many other types, including `Float64`, `Float32`, `Int64` and `Int32`.

We can test if type `T` is part of an abstract type `V` using the syntax `T <: V`:

```
In [56]: Float64 <: Real, Float32 <: Real, Int64 <: Real
```

Out[56]: (true, true, true)

Every type has one and only one super type, which is *always* an abstract type.

The function `supertype` applied to a type returns its super type:

```
In [57]: supertype(Int32) # returns Signed, which represents all signed integers.
```

Out[57]: Signed

```
In [58]: supertype(Float32) # returns `AbstractFloat`, which is a subtype of `Real`
```

Out[58]: AbstractFloat

An abstract type also has a super type:

```
In [59]: supertype(Real)
```

```
Out[59]: Number
```

## Type annotation and templating

The types `RadialComplex` and `MutableRadialComplex` won't be efficient as we have not told the compiler the type of `r` and `θ`. For the purposes of this module, this is fine as we are not focussing on high performance computing. However, it may be of interest how to rectify this.

We can impose a type on the field name with `::`:

```
In [60]: struct FastRadialComplex
    r::Float64
    θ::Float64
end
z = FastRadialComplex(1,0.1)
z.r, z.θ
```

```
Out[60]: (1.0, 0.1)
```

In this case `z` is stored using precisely 128-bits.

Sometimes we want to support multiple types. For example, we may wish to support 32-bit floats. This can be done as follows:

```
In [61]: struct TemplatizedRadialComplex{T}
    r::T
    θ::T
end
z = TemplatizedRadialComplex(1f0,0.1f0) # f0 creates a `Float32`
```

```
Out[61]: TemplatizedRadialComplex{Float32}(1.0f0, 0.1f0)
```

This is stored in precisely 64-bits.

## Relationship with C structs, heap and stack (advanced)

For those familiar with C, a `struct` in Julia whose fields are primitive types or composite types built from primitive types, is exactly equivalent to a `struct` C, and can in fact be passed to C functions without any performance cost. Behind the scenes Julia uses the LLVM compiler and so C and Julia can be freely mixed.

Another thing to note is that from a programmers perspective there are three types of memory: `registers`, the `stack` and the `heap`. Registers only live on the CPU and are extremely fast. The stack lives in memory and has fixed memory length and is much

faster than the heap as it avoids dynamic allocation and deallocation of memory. So an instance of a type with a known fixed length (like `FastRadialComplex`) will typically live either in registers or in the stack (the compiler sorts out the details) and be much faster than an instance of a type with unknown or variable length (like `RadialComplex` or `Vector`), which must be on the heap.

## 5. Loops and branches

Loops such as `for` work essentially the same as in Python. The one caveat is to remember we are using 1-based indexing, e.g., `1:5` is a range consisting of `[1,2,3,4,5]`:

```
In [62]: for k = 1:5
          println(k^2)
end
```

```
1
4
9
16
25
```

There are also `while` loops:

```
In [63]: x = 1
while x < 5
    println("x is $x which is less than 5, incrementing!")
    x += 1
end
x
```

```
x is 1 which is less than 5, incrementing!
x is 2 which is less than 5, incrementing!
x is 3 which is less than 5, incrementing!
x is 4 which is less than 5, incrementing!
```

```
Out[63]: 5
```

If-elseif-else statements look like:

```
In [64]: x = 5
if isodd(x)
    println("it's odd")
elseif x == 2
    println("it's 2")
else
    println("it's even")
end
```

```
it's odd
```

## 6. Functions

Functions are created in a number of ways. The most standard way is using the keyword `function`, followed by a name for the function, and in parentheses a list of arguments. Let's make a function that takes in a single number  $x$  and returns  $x^2$ .

```
In [65]: function sq(x)
    x^2
end
sq(2), sq(3)
```

```
Out[65]: (4, 9)
```

There is also a convenient syntax for defining functions on one line, e.g., we can also write

```
In [66]: sq(x) = x^2
```

```
Out[66]: sq (generic function with 1 method)
```

Multiple arguments to the function can be included with `,`.

Here's a function that takes in 3 arguments and returns the average.

(We write it on 3 lines only to show that functions can take multiple lines.)

```
In [67]: function av(x, y, z)
    ret = x + y
    ret = ret + z
    ret/3
end
av(1, 2, 3)
```

```
Out[67]: 2.0
```

Variables live in different scopes. In the previous example, `x`, `y`, `z` and `ret` are *local variables*: they only exist inside of `av`.

So this means `x` and `z` are *not* the same as our complex number `x` and `z` defined above.

**Warning:** if you reference variables not defined inside the function, they will use the outer scope definition.

The following example shows that if we mistype the first argument as `xx`, then it takes on the outer scope definition `x`, which is a complex number:

```
In [68]: function av2(xx, y, z)
    (x + y + z)/3
end
```

```
Out[68]: av2 (generic function with 1 method)
```

You should almost never use this feature!!

We should ideally be able to predict the output of a function from knowing just the

inputs.

**Example** Let's create a function that calculates the average of the entries of a vector.

```
In [69]: function vecaverage(v)
    ret=0
    for k = 1:length(v)
        ret = ret + v[k]
    end
    ret/length(v)
end
vecaverage([1,5,2,3,8,2])
```

```
Out[69]: 3.5
```

Julia has an inbuilt `sum` command that we can use to check our code:

```
In [70]: sum([1,5,2,3,8,2])/6
```

```
Out[70]: 3.5
```

## Functions with type signatures

functions can be defined only for specific types using `::` after the variable name. The same function name can be used with different type signatures.

The following defines a function `mydot` that calculates the dot product, with a definition changing depending on whether it is an `Integer` or a `Vector`.

Note that `Integer` is an abstract type that includes all integer types: `mydot` is defined for pairs of `Int64`'s, `Int32`'s, etc.

```
In [71]: function mydot(a::Integer, b::Integer)
    a*b
end

function mydot(a::AbstractVector, b::AbstractVector)
    # we assume length(a) == length(b)
    ret = 0
    for k = 1:length(a)
        ret = ret + a[k]*b[k]
    end
    ret
end

mydot(5, 6) # calls the first definition
```

```
Out[71]: 30
```

```
In [72]: mydot(Int8(5), Int8(6)) # also calls the first definition
```

```
Out[72]: 30
```

```
In [73]: mydot(1:3, [4,5,6])      # calls the second definition
```

```
Out[73]: 32
```

We should actually check that the lengths of `a` and `b` match.

Let's rewrite `mydot` using an `if`, `else` statement. The following code only does the for loop if the length of `a` is equal to the length of `b`, otherwise, it throws an error.

If we name something with the exact same signature (name, and argument types), previous definitions get overridden. Here we correct the implementation of `mydot` to throw an error if the lengths of the inputs do not match:

```
In [74]: function mydot(a::AbstractVector, b::AbstractVector)
    ret=0
    if length(a) == length(b)
        for k = 1:length(a)
            ret = ret + a[k]*b[k]
        end
    else
        error("arguments have different lengths")
    end
    ret
end
mydot([1,2,3], [5,6,7,8])
```

```
arguments have different lengths
```

Stacktrace:

```
[1] error(s::String)
    @ Base ./error.jl:35
[2] mydot(a::Vector{Int64}, b::Vector{Int64})
    @ Main ./In[74]:8
[3] top-level scope
    @ In[74]:12
```

## Anonymous (lambda) functions

Just like Python it is possible to make anonymous functions, with two variants on syntax:

```
In [75]: f = x -> x^2
g = function(x)
    x^2
end
```

```
Out[75]: #9 (generic function with 1 method)
```

There is not much difference between named and anonymous functions, both are compiled in the same manner. The only difference is named functions are in a sense "permanent". One can essentially think of named functions as "global constant anonymous functions".

## Tuples

`Tuple`s are similar to vectors but written with the notation `(x,y,z)` instead of `[x,y,z]`. They allow the storage of *different types*. For example:

```
In [76]: t = (1,2.0,"hi")
```

```
Out[76]: (1, 2.0, "hi")
```

On the surface, this is very similar to a `Vector{Any}`:

```
In [77]: v=[1,2.0,"hi"]
```

```
Out[77]: 3-element Vector{Any}:
```

```
1  
2.0  
"hi"
```

The main difference is that a `Tuple` knows the type of its arguments:

```
In [78]: typeof(t)
```

```
Out[78]: Tuple{Int64, Float64, String}
```

The main benefit of tuples for us is that they provide a convenient way to return multiple arguments from a function. For example, the following returns both `cos(x)` and `x^2` from a single function:

```
In [79]: function mytuplereturn(x)  
    (cos(x), x^2)  
end  
mytuplereturn(5)
```

```
Out[79]: (0.28366218546322625, 25)
```

We can also employ the convenient syntax to create two variables at once:

```
In [80]: x,y = mytuplereturn(5)
```

```
Out[80]: (0.28366218546322625, 25)
```

## Modules, Packages, and Plotting

Julia, like Python, has modules and packages. For example to load support for linear algebra functionality like `norm` and `det`, we need to load the `LinearAlgebra` module:

```
In [81]: using LinearAlgebra  
norm([1,2,3]), det([1 2; 3 4])
```

```
Out[81]: (3.7416573867739413, -2.0)
```

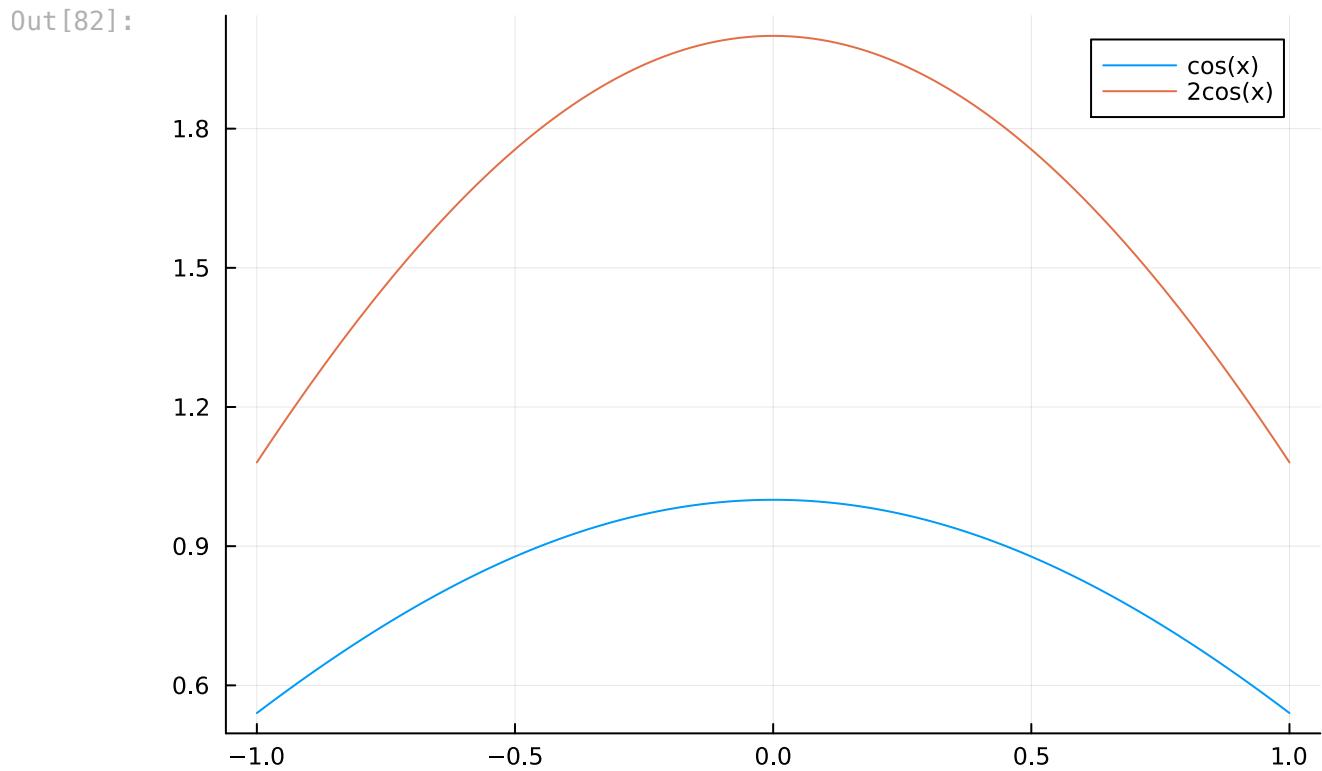
It is fairly straightforward to create one's own modules and packages, however, we will not need modules in this....module.

## Plotting

Some important functionality such as plotting requires non-built-in packages. There are many packages such as [PyPlot.jl](#), which wraps Python's [matplotlib](#) and [Makie.jl](#), which is a state-of-the-art GPU based 3D plotting package. We will use [Plots.jl](#), which is an umbrella package that supports different backends.

For example, we can plot a simple function as follows:

```
In [82]: using Plots  
x = range(-1, 1; length=1000) # Create a range of a 1000 evenly spaced numbers  
y = cos.(x) # Create a new vector with `cos` applied to each entry of `x`  
plot(x, y; label="cos(x)")  
plot!(x, 2y; label="2cos(x)")
```



Note the `!` is just a convention: any function that modifies its input or global state should have `!` at the end of its name.

## Installing packages (advanced)

If you choose to use Julia on your own machine, you may need to install packages. This can be done by typing the following, either in Jupyter or in the REPL: `] add Plots`.



## B. Asymptotics and Computational Cost

We introduce Big-O, little-o and asymptotic notation and see how they can be used to describe computational cost.

1. Asymptotics as  $n \rightarrow \infty$
2. Asymptotics as  $x \rightarrow x_0$
3. Computational cost

### 1. Asymptotics as $n \rightarrow \infty$

Big-O, little-o, and "asymptotic to" are used to describe behaviour of functions at infinity.

#### Definition 1 (Big-O)

$$f(n) = O(\phi(n)) \quad (\text{as } n \rightarrow \infty)$$

means  $\left| \frac{f(n)}{\phi(n)} \right|$  is bounded for sufficiently large  $n$ . That is, there exist constants  $C$  and  $N_0$  such that, for all  $n \geq N_0$ ,  $\left| \frac{f(n)}{\phi(n)} \right| \leq C$ .

#### Definition 2 (little-O)

$$f(n) = o(\phi(n)) \quad (\text{as } n \rightarrow \infty)$$

means  $\lim_{n \rightarrow \infty} \frac{f(n)}{\phi(n)} = 0$ .

#### Definition 3 (asymptotic to)

$$f(n) \sim \phi(n) \quad (\text{as } n \rightarrow \infty)$$

means  $\lim_{n \rightarrow \infty} \frac{f(n)}{\phi(n)} = 1$ .

#### Examples

1. \$

$$\{\cos n \mid \text{over } n^2 - 1\} = O(n^{-2}) \text{ as } \left| \frac{\cos n}{n^2 - 1} \right| \leq \left| \frac{n^2}{n^2 - 1} \right| \leq 2 \text{ for } n \geq N_0 = 2.$$

2. \$

$$\log n = o(n) \text{ as } \lim_{n \rightarrow \infty} \{\log n \mid \text{over } n\} = 0. \$$$

3. \$

$$n^2 + 1 \sim n^2 \text{ as } \{n^2 + 1\} / n^2 \rightarrow 1.$$

Note we sometimes write  $f(O(\phi(n)))$  for a function of the form  $f(g(n))$  such that  $g(n) = O(\phi(n))$ .

## Rules

We have some simple algebraic rules:

### Proposition 1 (Big-O rules)

$$\begin{aligned} O(\phi(n))O(\psi(n)) &= O(\phi(n)\psi(n)) && (\text{as } n \rightarrow \infty) \\ O(\phi(n)) + O(\psi(n)) &= O(|\phi(n)| + |\psi(n)|) && (\text{as } n \rightarrow \infty). \end{aligned}$$

## 2. Asymptotics as $x \rightarrow x_0$

We also have Big-O, little-o and "asymptotic to" at a point:

### Definition 4 (Big-O)

$$f(x) = O(\phi(x)) \quad (\text{as } x \rightarrow x_0)$$

means  $\left| \frac{f(x)}{\phi(x)} \right|$  is bounded in a neighbourhood of  $x_0$ . That is, there exist constants  $C$  and  $r$  such that, for all  $0 \leq |x - x_0| \leq r$ ,  $\left| \frac{f(x)}{\phi(x)} \right| \leq C$ .

### Definition 5 (little-O)

$$f(x) = o(\phi(x)) \quad (\text{as } x \rightarrow x_0)$$

means  $\lim_{x \rightarrow x_0} \frac{f(x)}{\phi(x)} = 0$ .

### Definition 6 (asymptotic to)

$$f(x) \sim \phi(x) \quad (\text{as } x \rightarrow x_0)$$

means  $\lim_{x \rightarrow x_0} \frac{f(x)}{\phi(x)} = 1$ .

### Example

$$\exp x = 1 + x + O(x^2) \quad \text{as } x \rightarrow 0$$

since  $\exp x = 1 + x + \frac{\exp t}{2}x^2$  for some  $t \in [0, x]$  and

$$\left| \frac{\frac{\exp t}{2}x^2}{x^2} \right| \leq \frac{3}{2}$$

provided  $x \leq 1$ .

### 3. Computational cost

We will use Big-O notation to describe the computational cost of algorithms. Consider the following simple sum

$$\sum_{k=1}^n x_k^2$$

which we might implement as:

```
In [1]: function sumsq(x)
    n = length(x)
    ret = 0.0
    for k = 1:n
        ret = ret + x[k]^2
    end
    ret
end

n = 100
x = randn(n)
sumsq(x)
```

```
Out[1]: 90.95882254617737
```

Each step of this algorithm consists of one memory look-up ( $z = x[k]$ ), one multiplication ( $w = z*z$ ) and one addition ( $ret = ret + w$ ). We will ignore the memory look-up in the following discussion. The number of CPU operations per step is therefore 2 (the addition and multiplication). Thus the total number of CPU operations is  $2n$ . But the constant 2 here is misleading: we didn't count the memory look-up, thus it is more sensible to just talk about the asymptotic complexity, that is, the *computational cost* is  $O(n)$ .

Now consider a double sum like:

$$\sum_{k=1}^n \sum_{j=1}^k x_j^2$$

which we might implement as:

```
In [2]: function sumsq2(x)
    n = length(x)
    ret = 0.0
    for k = 1:n
        for j = 1:k
            ret = ret + x[j]^2
        end
    end
    ret
end
```

```
n = 100
x = randn(n)
sumsq2(x)
```

Out[2]: 5377.916656174812

Now the inner loop is  $O(1)$  operations (we don't try to count the precise number), which we do  $k$  times for  $O(k)$  operations as  $k \rightarrow \infty$ . The outer loop therefore takes

$$\sum_{k=1}^n O(k) = O\left(\sum_{k=1}^n k\right) = O\left(\frac{n(n+1)}{2}\right) = O(n^2)$$

operations.

# C. Adjoints and Normal Matrices

Here we review

1. Complex inner-products
2. Adjoints
3. Normal Matrices and the spectral theorem

## 1. Complex-inner products

We will use bars to denote the complex-conjugate: if  $z = x + iy$  then  $\bar{z} = x - iy$ .

**Definition 1 (inner-product)** An *inner product*  $\langle \cdot, \cdot \rangle$  over  $\mathbb{C}$  satisfies, for  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{C}^n$  and  $a, b \in \mathbb{C}$ ,

1. Conjugate symmetry:  $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$
2. Linearity:  $\langle \mathbf{x}, a\mathbf{y} + b\mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{y} \rangle + b\langle \mathbf{x}, \mathbf{z} \rangle$ . (Some authors use linearity in the first argument.)
3. Positive-definiteness: for  $\mathbf{x} \neq 0$  we have  $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ .

We will usually use the standard inner product defined on  $\mathbb{C}^n$ :

$$\langle \mathbf{x}, \mathbf{y} \rangle := \bar{\mathbf{x}}^\top \mathbf{y} = \sum_{k=1}^n \bar{x}_k y_k$$

Note that  $\overline{zw} = \bar{z}\bar{w}$  and  $\overline{z+w} = \bar{z} + \bar{w}$  together imply that:

$$\overline{A\mathbf{x}} = A\bar{\mathbf{x}}.$$

## 2. Adjoints

**Definition 2 (adjoint)** Given an inner product, an adjoint of a matrix  $A \in \mathbb{C}^{m \times n}$  is the unique matrix  $A^*$  satisfying

$$\langle \mathbf{x}, A\mathbf{y} \rangle = \langle A\mathbf{x}, \mathbf{y} \rangle$$

for all  $\mathbf{x} \in \mathbb{C}^m$ ,  $\mathbf{y} \in \mathbb{C}^n$ . (Note this definition can be extended to other inner products.)

**Proposition 1 (adjoints are conjugate-transposes)** For the standard inner product,  $A^* = \bar{A}^\top$ . If  $A \in \mathbb{R}^{m \times n}$  then it reduces to the transpose  $A^* = A^\top$ .

### Proof

$$A_{k,j} = \langle \mathbf{e}_k, A\mathbf{e}_j \rangle = \langle A^*\mathbf{e}_k, \mathbf{e}_j \rangle = \mathbf{e}_k^\top \overline{(A^*)^\top} \mathbf{e}_j = \overline{A_{j,k}}$$

In this module we will assume the standard inner product (unless otherwise stated), that is, we will only use the standard adjoint  $A^* = \bar{A}^\top$ . Note if  $A = A^*$  a matrix is called *Hermitian*. If it is real it is also called *Symmetric*.

### 3. Normal matrices

**Definition 3 (normal)** A (square) *normal matrix* commutes with its adjoint:

$$AA^* = A^*A$$

Examples of normal matrices include:  
 2. Symmetric and Hermitian: ( $A^*A = A^2 = AA^*$ )  
 3. Orthogonal and Unitary: ( $Q^*Q = I = QQ^*$ )

An important property of normal matrices is that they are diagonalisable using unitary eigenvectors:

**Theorem 1 (spectral theorem for normal matrices)** If  $A \in \mathbb{C}^{n \times n}$  is normal then it is diagonalisable with unitary eigenvectors:

$$A = Q\Lambda Q^*$$

where  $Q \in U(n)$  and  $\Lambda$  is diagonal.

We will prove this later in the module. There is an important corollary for symmetric matrices that you may have seen before:

**Corollary 1 (spectral theorem for symmetric matrices)** If  $A \in \mathbb{R}^{n \times n}$  is symmetric then it is diagonalisable with orthogonal eigenvectors:

$$A = Q\Lambda Q^\top$$

where  $Q \in O(n)$  and  $\Lambda$  is real and diagonal.

#### Proof

$A = Q\Lambda Q^*$  since its normal hence we find that:

$$\Lambda = \Lambda^* = (Q^*AQ)^* = Q^*A^\top Q = Q^*AQ = \Lambda$$

which shows that  $\Lambda$  is real. The fact that  $Q$  is real follows since the column  $\mathbf{q}_k = Q\mathbf{e}_k$  is in the null space of the real matrix  $A - \lambda_k I$ .

# C. Spectral theorem for symmetric and normal matrices

Here we review the proof of the spectral theorem for symmetric and normal matrices, as well as adjoints (complex-conjugation). Here we use the standard inner product defined on  $\mathbb{C}^n$ :

$$\langle \mathbf{x}, \mathbf{y} \rangle := \bar{\mathbf{x}}^\top \mathbf{y} = \sum_{k=1}^n \bar{x}_k y_k$$

where the bars indicate complex conjugate: if  $z = x + iy$  then  $\bar{z} = x - iy$ . Note that  $\overline{zw} = \bar{z}\bar{w}$  and  $\overline{z+w} = \bar{z} + \bar{w}$  together imply that:

$$\overline{A\mathbf{x}} = A\bar{\mathbf{x}}.$$

## 1. Adjoints

**Definition 1 (adjoint)** An adjoint of a matrix  $A \in \mathbb{C}^{m \times n}$  is its conjugate transpose:  $A^* := A^\top$ . If  $A \in \mathbb{R}^{m \times n}$  then it reduces to the transpose  $A^* = A^\top$ .

Note adjoints have the important product that for the standard inner product they satisfy:

$$\langle \mathbf{x}, A\mathbf{y} \rangle = \bar{\mathbf{x}}^\top (A\mathbf{y}) = (A^\top \bar{\mathbf{x}})^\top \mathbf{y} = (\overline{A^\top \mathbf{x}})^\top \mathbf{y} = \langle A\mathbf{x}, \mathbf{y} \rangle$$

## 2. Spectral theorem for symmetric matrices

**Theorem (spectral theorem for symmetric matrices)** If  $A$  is symmetric ( $A = A^\top$ ) then

$$A = Q\Lambda Q^\top$$

where  $Q$  is orthogonal ( $Q^\top Q = I$ ) and  $\Lambda$  is real.

### Proof

Recall every eigenvalue  $\lambda$  has at least one eigenvector  $\mathbf{q}$  which we can normalize, but this can be complex. Thus we have

$$\lambda = \lambda \mathbf{q}^* \mathbf{q} = \mathbf{q}^* A \mathbf{q} =$$

■

# I.1 Integers

In this chapter we discuss the following:

1. Binary representation: Any real number can be represented in binary, that is, by an infinite sequence of 0s and 1s (bits). We review binary representation. 2. Unsigned integers: We discuss how computers represent non-negative integers using only  $p$ -bits, via [modular arithmetic](#). 3. Signed integers: we discuss how negative integers are handled using the [Two's-complement](#) format. 4. As an advanced (non-examinable) topic we discuss `BigInt`, which uses variable bit length storage.

Before we begin it's important to have a basic model of how a computer works. Our simplified model of a computer will consist of a [Central Processing Unit \(CPU\)](#)—the brains of the computer—and [Memory](#)—where data is stored. Inside the CPU there are [registers](#), where data is temporarily stored after being loaded from memory, manipulated by the CPU, then stored back to memory.

Memory is a sequence of bits: `1`s and `0`s, essentially "on/off" switches. These are grouped into bytes, which consist of 8 bits. Each byte has a memory address: a unique number specifying its location in memory. The number of possible addresses is limited by the processor: if a computer has a  $p$ -bit CPU then each address is represented by  $p$  bits, for a total of  $2^p$  addresses (on a modern 64-bit CPU this is  $2^{64} \approx 1.8 \times 10^{19}$  bytes). Further, each register consists of exactly  $p$ -bits.

A CPU has the following possible operations:

1. load data from memory addresses (up to  $p$ -bits) to a register
2. store data from a register to memory addresses (up to  $p$ -bits)
3. Apply some basic functions ("+", "-", etc.) to the bits in one or two registers

and write the result to a register.

Mathematically, the important point is CPUs only act on  $2^p$  possible sequences of bits at a time. That is, essentially all functions  $f$  implemented on a CPU are either of the form  $f : \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}$  or  $f : \mathbb{Z}_{2^p} \times \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}$ , where we use the following notation:

**Definition 1 ( $\mathbb{Z}_m$ , signed integers)** Denote the

$$\mathbb{Z}_m := \{0, 1, \dots, m - 1\}$$

The limitations this imposes on representing integers is substantial. If we have an implementation of `+`, which we shall denote  $\oplus_m$ , how can we possibly represent  $m + 1$  in this implementation when the result is above the largest possible integer?

The solution that is used is straightforward: the CPU uses modular arithmetic. E.g., we have

$$(m - 1) \oplus_m 1 = m \pmod{m} = 0.$$

In this chapter we discuss the implications of this approach and how it works with negative numbers.

We will use Julia in these notes to explore what is happening as a computer does integer arithmetic. We load an external package which implements functions `printbits` (and `printlnbits`) to print the bits (and with a newline) of numbers in colour:

```
In [1]: using ColorBitstring
```

## 1. Binary representation

Any integer can be presented in binary format, that is, a sequence of `0`s and `1`s.

**Definition 2 (binary format)** For  $B_0, \dots, B_p \in \{0, 1\}$  denote an integer in *binary format* by:

$$\pm(B_p \dots B_1 B_0)_2 := \pm \sum_{k=0}^p B_k 2^k$$

**Example 1 (integers in binary)** A simple integer example is  $5 = 2^2 + 2^0 = (101)_2$ . On the other hand, we write  $-5 = -(101)_2$ . Another example is  $258 = 2^8 + 2 = (1000000010)_2$ .

## 2. Unsigned Integers

Computers represent integers by a finite number of  $p$  bits, with  $2^p$  possible combinations of 0s and 1s. For *unsigned integers* (non-negative integers) these bits dictate the first  $p$  binary digits:  $(B_{p-1} \dots B_1 B_0)_2$ .

Integers on a computer follow [modular arithmetic](#): Integers represented with  $p$ -bits on a computer actually represent elements of  $\mathbb{Z}_{2^p}$  and integer arithmetic on a computer is equivalent to arithmetic modulo  $2^p$ . We denote modular arithmetic with  $m = 2^p$  as follows:

$$\begin{aligned} x \oplus_m y &:= (x + y) \pmod{m} \\ x \ominus_m y &:= (x - y) \pmod{m} \\ x \otimes_m y &:= (x * y) \pmod{m} \end{aligned}$$

When  $m$  is implied by context we just write  $\oplus, \ominus, \otimes$ .

**Example 2 (arithmetic with 8-bit unsigned integers)** If arithmetic lies between 0 and  $m = 2^8 = 256$  works as expected. For example,

$$17 \oplus_{256} 3 = 20 \pmod{256} = 20$$

$$17 \ominus_{256} 3 = 14 \pmod{256} = 14$$

This can be seen in Julia:

```
In [2]: x = UInt8(17) # An 8-bit representation of the number 255, i.e. with bits 0
y = UInt8(3) # An 8-bit representation of the number 1, i.e. with bits 0
printbits(x); println(" + "); printbits(y); println(" = ")
printlnbits(x + y) # + is automatically modular arithmetic
printbits(x); println(" - "); printbits(y); println(" = ")
printbits(x - y) # - is automatically modular arithmetic

00010001 +
00000011 =
00010100
00010001 -
00000011 =
00001110
```

**Example 3 (overflow with 8-bit unsigned integers)** If we go beyond the range the result "wraps around". For example, with integers we have

$$255 + 1 = (11111111)_2 + (00000001)_2 = (100000000)_2 = 256$$

However, the result is impossible to store in just 8-bits! So as mentioned instead it treats the integers as elements of  $\mathbb{Z}_{256}$ :

$$255 \oplus_{256} 1 = 255 + 1 \pmod{256} = (00000000)_2 \pmod{256} = 0 \pmod{256}$$

We can see this in code:

```
In [3]: x = UInt8(255) # An 8-bit representation of the number 255, i.e. with bits 1
y = UInt8(1) # An 8-bit representation of the number 1, i.e. with bits 0
printbits(x); println(" + "); printbits(y); println(" = ")
printbits(x + y) # + is automatically modular arithmetic

11111111 +
00000001 =
00000000
```

On the other hand, if we go below 0 we wrap around from above:

$$3 \ominus_{256} 5 = -2 \pmod{256} = 254 = (11111110)_2$$

```
In [4]: x = UInt8(3) # An 8-bit representation of the number 3, i.e. with bits 000
y = UInt8(5) # An 8-bit representation of the number 5, i.e. with bits 000
printbits(x); println(" - "); printbits(y); println(" = ")
printbits(x - y) # + is automatically modular arithmetic
```

```
00000011 -  
00000101 =  
11111110
```

**Example 4 (multiplication of 8-bit unsigned integers)** Multiplication works similarly: for example,

$$254 \otimes_{256} 2 = 254 * 2 \pmod{256} = 252 \pmod{256} = (11111100)_2 \pmod{256}$$

We can see this behaviour in code by printing the bits:

```
In [5]: x = UInt8(254) # An 8-bit representation of the number 254, i.e. with bits 1  
y = UInt8(2) # An 8-bit representation of the number 2, i.e. with bits 0  
printbits(x); println(" * "); printbits(y); println(" = ")  
printbits(x * y)  
  
11111110 *  
00000010 =  
11111100
```

## Hexadecimal and binary format

In Julia unsigned integers are displayed in hexadecimal form: that is, in base-16. Since there are only 10 standard digits ( 0–9 ) it uses 6 letters ( a–f ) to represent 11–16. For example,

```
In [6]: UInt8(250)
```

```
Out[6]: 0xfa
```

because f corresponds to 15 and a corresponds to 10, and we have

$$15 * 16 + 10 = 250.$$

The reason for this is that each hex-digit encodes 4 bits (since 4 bits have  $2^4 = 16$  possible values) and hence two hex-digits are encode 1 byte, and thus the digits correspond exactly with how memory is divided into addresses. We can create unsigned integers either by specifying their hex format:

```
In [7]: 0xfa
```

```
Out[7]: 0xfa
```

Alternatively, we can specify their digits. For example, we know  $(f)_{16} = 15 = (1111)_2$  and  $(a)_{16} = 10 = (1010)_2$  and hence  $250 = (fa)_{16} = (11111010)_2$  can be written as

```
In [8]: 0b11111010
```

```
Out[8]: 0xfa
```

## 3. Signed integer

Signed integers use the [Two's complement](#) convention. The convention is if the first bit is 1 then the number is negative: the number  $2^p - y$  is interpreted as  $-y$ . Thus for  $p = 8$  we are interpreting  $2^7$  through  $2^8 - 1$  as negative numbers. More precisely:

**Definition 3 ( $\mathbb{Z}_{2^p}^s$ , unsigned integers)**

$$\mathbb{Z}_{2^p}^s := \{-2^{p-1}, \dots, -1, 0, 1, \dots, 2^{p-1} - 1\}$$

**Definition 4 (Shifted mod)** Define for  $y = x \pmod{2^p}$

$$x \pmod{s}{2^p} := \begin{cases} y & 0 \leq y \leq 2^{p-1} - 1 \\ y - 2^p & 2^{p-1} \leq y \leq 2^p - 1 \end{cases}$$

Note that if  $R_p(x) = x \pmod{s}{2^p}$  then it can be viewed as a map  $R_p : \mathbb{Z} \rightarrow \mathbb{Z}_{2^p}^s$  or a one-to-one map  $R_p : \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}^s$  whose inverse is  $R_p^{-1}(x) = x \pmod{2^p}$ .

**Example 5 (converting bits to signed integers)** What 8-bit integer has the bits

`01001001`? Because the first bit is 0 we know the result is positive. Adding the corresponding decimal places we get:

In [9]: `2^0 + 2^3 + 2^6`

Out[9]: 73

What 8-bit (signed) integer has the bits `11001001`? Because the first bit is 1 we know it's a negative number, hence we need to sum the bits but then subtract  $2^p$ :

In [10]: `2^0 + 2^3 + 2^6 + 2^7 - 2^8`

Out[10]: -55

We can check the results using `printbits`:

In [11]: `printlnbits(Int8(73)) # Int8 is an 8-bit representation of the signed integer`  
`printbits(-Int8(55))`

`01001001`  
`11001001`

Arithmetic works precisely the same for signed and unsigned integers, e.g. we have

$$x \oplus_{2^p}^s y := x + y \pmod{s}{2^p}$$

**Example 6 (addition of 8-bit integers)** Consider `(-1) + 1` in 8-bit arithmetic. The number  $-1$  has the same bits as  $2^8 - 1 = 255$ . Thus this is equivalent to the previous question and we get the correct result of `0`. In other words:

$$-1 \oplus_{256} 1 = -1 + 1 \pmod{2^8} = 2^8 - 1 + 1 \pmod{2^8} = 2^8 \pmod{2^8} = 0 \pmod{2^8}$$

**Example 7 (multiplication of 8-bit integers)** Consider `(-2) * 2`.  $-2$  has the same bits as  $2^{256} - 2 = 254$  and  $-4$  has the same bits as  $2^{256} - 4 = 252$ , and hence from the previous example we get the correct result of  $-4$ . In other words:

$$(-2) \otimes_{2^p}^s 2 = (-2) * 2 \pmod{2^p} = (2^p - 2) * 2 \pmod{2^p} = 2^{p+1} - 4 \pmod{2^p} = -4$$

**Example 8 (overflow)** We can find the largest and smallest instances of a type using `typemax` and `typemin`:

```
In [12]: printlnbits(typemax(Int8)) # 2^7-1 = 127
printlnbits(typemin(Int8)) # -2^7 = -128
01111111
10000000
```

As explained, due to modular arithmetic, when we add `1` to the largest 8-bit integer we get the smallest:

```
In [13]: typemax(Int8) + Int8(1) # returns typemin(Int8)
Out[13]: -128
```

This behaviour is often not desired and is known as *overflow*, and one must be wary of using integers close to their largest value.

## Division

In addition to `+`, `-`, and `*` we have integer division `÷`, which rounds towards zero:

```
In [14]: 5 ÷ 2 # equivalent to div(5,2)
Out[14]: 2
```

Standard division `/` (or `\` for division on the right) creates a floating-point number, which will be discussed in the next chapter:

```
In [15]: 5 / 2 # alternatively 2 | 5
Out[15]: 2.5
```

We can also create rational numbers using `//`:

```
In [16]: (1//2) + (3//4)
Out[16]: 5//4
```

Rational arithmetic often leads to overflow so it is often best to combine `big` with rationals:

```
In [17]: big(102324)//132413023 + 23434545//4243061 + 23434545//42430534435
```

```
Out[17]: 26339037835007648477541540//4767804878707544364596461
```

## 4. Variable bit representation (non-examinable)

An alternative representation for integers uses a variable number of bits, with the advantage of avoiding overflow but with the disadvantage of a substantial speed penalty.

In Julia these are `BigInt`s, which we can create by calling `big` on an integer:

```
In [18]: x = typemax(Int64) + big(1) # Too big to be an `Int64`
```

```
Out[18]: 9223372036854775808
```

Note in this case addition automatically promotes an `Int64` to a `BigInt`. We can create very large numbers using `BigInt`:

```
In [19]: x^100
```

```
Out[19]: 308299402527763474570010682154566572137179853330569745885534227792109373198  
447640470596653941241089824056172991237203850122889314192108015240464239377  
659907729443406151990542412460139422694360143091643438371471672472022733159  
695061370166103454894838872109766727543876375812850840329719945826027770730  
120246098009381841416708056334276148239586243518509394244354072236315177002  
222178324395959253133606299849420991475240801906072080512453438264605109361  
38148486460620386624234875043260443612037084304893058642343380140154714002  
337629571838339036072866290023067143715171661582628684226791756074958601816  
573949210192042971926128564012559683306389156286526215702602395591987379284  
682309585448452092050934594471287167569179082769090777848505882924858894568  
168528817978796393118106206809246398429622597308249405630795808918972670167  
873557636539414623207691708807594905363669045958112877309721274696727649649  
601081087800063823914375007554316324004987448998664232743644123445804025448  
082503822047990459461530060239055638579924527680558002493780472302931956594  
201351581704871454345525023520878974570116527956902624814539521898506299183  
170783021797439315846606778519958103771496882062824105186711983296636153004  
791033906572655026074103671610093220596965508325771424407112022165467934046  
108400156032167602544380124835543930597492387362414798072811058145280610901  
173900506006060422808766749928885121870507880736423792545581389057525756998  
145009099711769746929923409439498484057402540146394209901941336109623390905  
61174276634397649549164015925656511157141476925718770456826870124308204483  
840020135761385100647110424482884227023263774739896271187541348841577264708  
85711252729324907172174682636046833259334695562978550702077536636800275361  
270990152624845632820964329212289967743661388636076587788674818529924999492  
184318357313040349631189661494939940979601130119128006720905325934191881396  
7552543176532349157376
```

Note the number of bits is not fixed, the larger the number, the more bits required to represent it, so while overflow is impossible, it is possible to run out of memory if a number is astronomically large: go ahead and try `x^x` (at your own risk).

## I.2 Reals

Reference: [Overton \(<https://cs.nyu.edu/~overton/book/>\)](https://cs.nyu.edu/~overton/book/)

In this chapter, we introduce the [IEEE Standard for Floating-Point Arithmetic \(\[https://en.wikipedia.org/wiki/IEEE\\\_754\]\(https://en.wikipedia.org/wiki/IEEE\_754\)\)](https://en.wikipedia.org/wiki/IEEE_754). There are multiple ways of representing real numbers on a computer, as well as the precise behaviour of operations such as addition, multiplication, etc.: one can use

1. [Fixed-point arithmetic \(\[https://en.wikipedia.org/wiki/Fixed-point\\\_arithmetic\]\(https://en.wikipedia.org/wiki/Fixed-point\_arithmetic\)\):](https://en.wikipedia.org/wiki/Fixed-point_arithmetic) essentially representing a real number as integer where a decimal point is inserted at a fixed point. This turns out to be impractical in most applications, e.g., due to loss of relative accuracy for small numbers.
2. [Floating-point arithmetic \(\[https://en.wikipedia.org/wiki/Floating-point\\\_arithmetic\]\(https://en.wikipedia.org/wiki/Floating-point\_arithmetic\)\):](https://en.wikipedia.org/wiki/Floating-point_arithmetic) essentially scientific notation where an exponent is stored alongside a fixed number of digits. This is what is used in practice.
3. [Level-index arithmetic \(\[https://en.wikipedia.org/wiki/Symmetric\\\_level-index\\\_arithmetic\]\(https://en.wikipedia.org/wiki/Symmetric\_level-index\_arithmetic\)\):](https://en.wikipedia.org/wiki/Symmetric_level-index_arithmetic) stores numbers as iterated exponents. This is the most beautiful mathematically but unfortunately is not as useful for most applications and is not implemented in hardware.

Before the 1980s each processor had potentially a different representation for floating-point numbers, as well as different behaviour for operations. IEEE introduced in 1985 was a means to standardise this across processors so that algorithms would produce consistent and reliable results.

This chapter may seem very low level for a mathematics course but there are two important reasons to understand the behaviour of floating-point numbers in details:

1. Floating-point arithmetic is very precisely defined, and can even be used in rigorous computations as we shall see in the labs. But it is not exact and its important to understand how errors in computations can accumulate.
2. Failure to understand floating-point arithmetic can cause catastrophic issues in practice, with the extreme example being the [explosion of the Ariane 5 rocket \(<https://youtu.be/N6PWATvLQCY?t=86>\)](https://youtu.be/N6PWATvLQCY?t=86).

In this chapter we discuss the following:

1. Real numbers in binary: we discuss how binary digits can be used to represent real numbers.
2. Floating-point numbers: Real numbers are stored on a computer with a finite number of bits. There are three types of floating-point numbers: *normal numbers*, *subnormal numbers*, and *special numbers*.
3. Arithmetic: Arithmetic operations in floating-point are exact up to rounding, and how the rounding mode can be set. This allows us to bound errors computations.

4. High-precision floating-point numbers: As an advanced (non-examinable) topic, we discuss how the precision of floating-point arithmetic can be increased arbitrary using `BigFloat`.

Before we begin, we load two external packages. `SetRounding.jl` allows us to set the rounding mode of floating-point arithmetic. `ColorBitstring.jl` implements functions `printbits` (and `printlnbits`) which print the bits (and with a newline) of floating-point numbers in colour.

```
In [1]: using SetRounding, ColorBitstring
```

## 1. Real numbers in binary

Reals can also be presented in binary format, that is, a sequence of 0 s and 1 s alongside a decimal point:

**Definition 1 (real binary format)** For  $b_1, b_2, \dots \in \{0, 1\}$ , Denote a non-negative real number in *binary format* by:

$$(B_p \dots B_0. b_1 b_2 b_3 \dots)_2 := (B_p \dots B_0)_2 + \sum_{k=1}^{\infty} \frac{b_k}{2^k}.$$

**Example 1 (rational in binary)** Consider the number  $1/3$ . In decimal recall that:

$$1/3 = 0.3333 \dots = \sum_{k=1}^{\infty} \frac{3}{10^k}$$

We will see that in binary

$$1/3 = (0.010101 \dots)_2 = \sum_{k=1}^{\infty} \frac{1}{2^{2k}}$$

Both results can be proven using the geometric series:

$$\sum_{k=0}^{\infty} z^k = \frac{1}{1-z}$$

provided  $|z| < 1$ . That is, with  $z = \frac{1}{4}$  we verify the binary expansion:

$$\sum_{k=1}^{\infty} \frac{1}{4^k} = \frac{1}{1-1/4} - 1 = \frac{1}{3}$$

A similar argument with  $z = 1/10$  shows the decimal case.

## 2. Floating-point numbers

Floating-point numbers are a subset of real numbers that are representable using a fixed number of bits.

**Definition 2 (floating-point numbers)** Given integers  $\sigma$  (the "exponential shift")  $Q$  (the number of exponent bits) and  $S$  (the precision), define the set of *Floating-point numbers* by dividing into *normal*, *sub-normal*, and *special number* subsets:

$$F_{\sigma,Q,S} := F_{\sigma,Q,S}^{\text{normal}} \cup F_{\sigma,Q,S}^{\text{sub}} \cup F^{\text{special}}.$$

The *normal numbers*  $F_{\sigma,Q,S}^{\text{normal}} \subset \mathbb{R}$  are

$$F_{\sigma,Q,S}^{\text{normal}} := \{\pm 2^{q-\sigma} \times (1.b_1 b_2 b_3 \dots b_S)_2 : 1 \leq q < 2^Q - 1\}.$$

The *sub-normal numbers*  $F_{\sigma,Q,S}^{\text{sub}} \subset \mathbb{R}$  are

$$F_{\sigma,Q,S}^{\text{sub}} := \{\pm 2^{1-\sigma} \times (0.b_1 b_2 b_3 \dots b_S)_2\}.$$

The *special numbers*  $F^{\text{special}}$   $\not\subset \mathbb{R}$  are

$$F^{\text{special}} := \{\infty, -\infty, \text{NaN}\}$$

where NaN is a special symbol representing "not a number", essentially an error flag.

Note this set of real numbers has no nice *algebraic structure*: it is not closed under addition, subtraction, etc. On the other hand, we can control errors effectively hence it is extremely useful for analysis.

Floating-point numbers are stored in  $1 + Q + S$  total number of bits, in the format

$$sq_{Q-1} \dots q_0 b_1 \dots b_S$$

The first bit ( $s$ ) is the **sign bit**: 0 means positive and 1 means negative. The bits  $q_{Q-1} \dots q_0$  are the **exponent bits**: they are the binary digits of the unsigned integer  $q$ :

$$q = (q_{Q-1} \dots q_0)_2.$$

Finally, the bits  $b_1 \dots b_S$  are the **significand bits**. If  $1 \leq q < 2^Q - 1$  then the bits represent the normal number

$$x = \pm 2^{q-\sigma} \times (1.b_1 b_2 b_3 \dots b_S)_2.$$

If  $q = 0$  (i.e. all bits are 0) then the bits represent the sub-normal number

$$x = \pm 2^{1-\sigma} \times (0.b_1 b_2 b_3 \dots b_S)_2.$$

If  $q = 2^Q - 1$  (i.e. all bits are 1) then the bits represent a special number, discussed later.

## IEEE floating-point numbers

**Definition 3 (IEEE floating-point numbers)** IEEE has 3 standard floating-point formats: 16-bit (half precision), 32-bit (single precision) and 64-bit (double precision) defined by:

$$F_{16} := F_{15,5,10}$$

$$F_{32} := F_{127,8,23}$$

$$F_{64} := F_{1023,11,52}$$

In Julia these correspond to 3 different floating-point types:

1. `Float64` is a type representing double precision ( $F_{64}$ ). We can create a `Float64` by including a decimal point when writing the number: `1.0` is a `Float64`. Alternatively, one can use scientific notation: `1e0`. `Float64` is the default format for scientific computing (on the *Floating-Point Unit*, FPU).
2. `Float32` is a type representing single precision ( $F_{32}$ ). We can create a `Float32` by including a `f0` when writing the number: `1f0` is a `Float32` (this is in fact scientific notation so `1f1`  $\equiv$  `10f0`). `Float32` is generally the default format for graphics (on the *Graphics Processing Unit*, GPU), as the difference between 32 bits and 64 bits is indistinguishable to the eye in visualisation, and more data can be fit into a GPU's limited memory.

3. **Float16** is a type representing half-precision ( $F_{16}$ ). It is important in machine learning where one wants to maximise the amount of data and high accuracy is not necessarily helpful.

**Example 2 (rational in `Float32` )** How is the number  $1/3$  stored in `Float32`? Recall that

$$\frac{1}{3} = (0.010101 \dots)_2 = 2^{-2}(1.0101 \dots)_2 = 2^{125-127}(1.0101 \dots)_2$$

and since  $125 = (1111101)_2$  the exponent bits are **01111101**. . For the significand we round the last bit to the nearest element of  $F_{32}$  , (this is explained in detail in the section on rounding), so we have

$1.010101010101010101010101 \dots \approx 1.01010101010101010101011 \in F_{32}$   
 and the significand bits are 010101010101010101011. Thus the Float32 bits for  $1/3$  are:

```
In [2]: printbits(1f0/3)
```

For sub-normal numbers, the simplest example is zero, which has  $q = 0$  and all significand bits zero:

```
In [3]: printbits(0.0)
```

Unlike integers, we also have a negative zero:

In [4]: printbits(-0.0)

This is treated as identical to  $0.0$  (except for degenerate operations as explained in special numbers).

## Special normal numbers

When dealing with normal numbers there are some important constants that we will use to bound errors.

**Definition 4 (machine epsilon/smallest positive normal number/largest normal number)** *Machine epsilon* is denoted

$$\epsilon_{m,S} := 2^{-S}.$$

When  $S$  is implied by context we use the notation  $\epsilon_m$ . The *smallest positive normal number* is  $q = 1$  and  $b_k$  all zero:

$$\min |F_{\sigma,Q,S}^{\text{normal}}| = 2^{1-\sigma}$$

where  $|A| := \{|x| : x \in A\}$ . The *largest (positive) normal number* is

$$\max F_{\sigma,Q,S}^{\text{normal}} = 2^{2^Q-2-\sigma}(1.11\dots)_2 = 2^{2^Q-2-\sigma}(2 - \epsilon_m)$$

We confirm the simple bit representations:

```
In [5]: σ,Q,S = 127,8,23 # Float32
ε_m = 2.0^(-S)
printlnbits(Float32(2.0^(1-σ))) # smallest positive normal Float32
printlnbits(Float32(2.0^(2^Q-2-σ) * (2-ε_m))) # largest normal Float32

00000000100000000000000000000000000000000
0111111101111111111111111111111111111111
```

For a given floating-point type, we can find these constants using the following functions:

```
In [6]: eps(Float32), floatmin(Float32), floatmax(Float32)
```

```
Out[6]: (1.1920929f-7, 1.1754944f-38, 3.4028235f38)
```

**Example 3 (creating a sub-normal number)** If we divide the smallest normal number by two, we get a subnormal number:

```
In [7]: mn = floatmin(Float32) # smallest normal Float32
printlnbits(mn)
printbits(mn/2)

0000000010000000000000000000000000000000
0000000010000000000000000000000000000000
```

Can you explain the bits?

## Special numbers

The special numbers extend the real line by adding  $\pm\infty$  but also a notion of "not-a-number" NaN. Whenever the bits of  $q$  of a floating-point number are all 1 then they represent an element of  $F^{\text{special}}$ . If all  $b_k = 0$ , then the number represents either  $\pm\infty$ , called Inf and -Inf for 64-bit floating-point numbers (or Inf16 , Inf32 for 16-bit and 32-bit, respectively):

In [8]: `printlnbits(Inf16)  
printbits(-Inf16)`

```
011111000000000000  
111111000000000000
```

All other special floating-point numbers represent NaN. One particular representation of NaN is denoted by NaN for 64-bit floating-point numbers (or NaN16 , NaN32 for 16-bit and 32-bit, respectively):

In [9]: `printbits(NaN16)`

```
011111000000000000
```

These are needed for undefined algebraic operations such as:

In [10]: `0/0`

Out [10]: NaN

Essentially it is a CPU's way of indicating an error has occurred.

**Example 4 (many NaN s)** What happens if we change some other  $b_k$  to be nonzero? We can create bits as a string and see:

In [11]: `i = 0b011110000010001 # an UInt16  
reinterpret(Float16, i)`

Out [11]: NaN16

Thus, there are more than one `NaN`s on a computer.

### 3. Arithmetic

Arithmetic operations on floating-point numbers are *exact up to rounding*. There are three basic rounding strategies: round up/down/nearest. Mathematically we introduce a function to capture the notion of rounding:

**Definition 6 (rounding)**  $\text{fl}_{\sigma,Q,S}^{\text{up}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  denotes the function that rounds a real number up to the nearest floating-point number that is greater or equal.

$\text{fl}_{\sigma,Q,S}^{\text{down}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  denotes the function that rounds a real number down to the nearest floating-point number that is greater or equal.  $\text{fl}_{\sigma,Q,S}^{\text{nearest}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  denotes the function that rounds a real number to the nearest floating-point number. In case of a tie, it returns the floating-point number whose least significant bit is equal to zero. We use the notation  $\text{fl}$  when  $\sigma, Q, S$  and the rounding mode are implied by context, with  $\text{fl}^{\text{nearest}}$  being the default rounding mode.

In Julia, the rounding mode is specified by tags `RoundUp`, `RoundDown`, and `RoundNearest`. (There are also more exotic rounding strategies `RoundToZero`, `RoundNearestTiesAway` and `RoundNearestTiesUp` that we won't use.)

Let's try rounding a `Float64` to a `Float32`.

The default rounding mode can be changed:

Or alternatively we can change the rounding mode for a chunk of code using `setrounding`. The following computes upper and lower bounds for  $\sqrt{2}$ :

```
In [14]: x = 1f0
setrounding(Float32, RoundDown) do
    x/3
end,
setrounding(Float32, RoundUp) do
    x/3
end
```

Out[14]: (0.3333333f0, 0.33333334f0)

**WARNING (compiled constants, non-examinable):** Why did we first create a variable `x` instead of typing `1f0/3`? This is due to a very subtle issue where the compiler is *too clever for its own good*: it recognises `1f0/3` can be computed at compile time, but failed to recognise the rounding mode was changed.

In IEEE arithmetic, the arithmetic operations `+`, `-`, `*`, `/` are defined by the property that they are exact up to rounding. Mathematically we denote these operations as  $\oplus, \ominus, \otimes, \oslash : F \otimes F \rightarrow F$  as follows:

$$\begin{aligned}x \oplus y &:= \text{fl}(x + y) \\x \ominus y &:= \text{fl}(x - y) \\x \otimes y &:= \text{fl}(x * y) \\x \oslash y &:= \text{fl}(x/y)\end{aligned}$$

Note also that `^` and `sqrt` are similarly exact up to rounding. Also, note that when we convert a Julia command with constants specified by decimal expansions we first round the constants to floats, e.g., `1.1 + 0.1` is actually reduced to

$$\text{fl}(1.1) \oplus \text{fl}(0.1)$$

This includes the case where the constants are integers (which are normally exactly floats but may be rounded if extremely large).

**Example 5 (decimal is not exact)** The Julia command `1.1+0.1` gives a different result than `1.2`:

```
In [15]: x = 1.1
y = 0.1
x + y - 1.2 # Not Zero???
```

Out[15]: 2.220446049250313e-16

This is because  $\text{fl}(1.1) \neq 1 + 1/10$  and  $\text{fl}(1.1) \neq 1/10$  since their expansion in *binary* is not finite, but rather:

Thus when we add them we get

On the other hand,

which differs by 1 bit.

**WARNING (non-associative)** These operations are not associative! E.g.  $(x \oplus y) \oplus z$  is not necessarily equal to  $x \oplus (y \oplus z)$ . Commutativity is preserved, at least. Here is a surprising example of non-associativity:

```
In [16]: (1.1 + 1.2) + 1.3, 1.1 + (1.2 + 1.3)
```

**Out[16]:** (3.5999999999999996, 3.6)

Can you explain this in terms of bits?

# Bounding errors in floating point arithmetic

Before we discuss bounds on errors, we need to talk about the two notions of errors:

**Definition 7 (absolute/relative error)** If  $\tilde{x} = x + \delta_a = x(1 + \delta_r)$  then  $|\delta_a|$  is called the *absolute error* and  $|\delta_r|$  is called the *relative error* in approximating  $x$  by  $\tilde{x}$ .

We can bound the error of basic arithmetic operations in terms of machine epsilon, provided a real number is close to a normal number:

**Definition 8 (normalised range)** The *normalised range*  $\mathcal{N}_{\sigma,Q,S} \subset \mathbb{R}$  is the subset of real numbers that lies between the smallest and largest normal floating-point number:

$$\mathcal{N}_{\sigma,Q,S} := \{x : \min |F_{\sigma,Q,S}^{\text{normal}}| \leq |x| \leq \max |F_{\sigma,Q,S}^{\text{normal}}|\}$$

When  $\sigma, O, S$  are implied by context we use the notation  $\mathcal{N}$ .

We can use machine epsilon to determine bounds on rounding:

**Proposition 1 (round bound)** If  $x \in \mathcal{N}$  then

$$f_l^{\text{mode}}(x) \equiv x(1 + \delta_x^{\text{mode}})$$

where the *relative error* is

$$|\delta_x^{\text{nearest}}| \leq \frac{\epsilon_m}{2}$$

$$|\delta_x^{\text{up/down}}| < \epsilon_m.$$

This immediately implies relative error bounds on all IEEE arithmetic operations, e.g., if  $x + y \in \mathcal{N}$  then we have

$$x \oplus y = (x + y)(1 + \delta_1)$$

where (assuming the default nearest rounding)  $|\delta_1| \leq \frac{\epsilon_m}{2}$ .

**Example 6 (bounding a simple computation)** We show how to bound the error in computing

$$(1.1 + 1.2) + 1.3$$

using floating-point arithmetic. First note that 1.1 on a computer is in fact  $\text{fl}(1.1)$ . Thus this computation becomes

$$(\text{fl}(1.1) \oplus \text{fl}(1.2)) \oplus \text{fl}(1.3)$$

First we find

$$(\text{fl}(1.1) \oplus \text{fl}(1.2)) = (1.1(1 + \delta_1) + 1.2(1 + \delta_2))(1 + \delta_3) = 2.3 + \underbrace{1.1\delta_1 + 1.2\delta_2 + 2}_{\delta_4}$$

In this module we will never ask for precise bounds: that is, we will always want bounds of the form  $C\epsilon_m$  for a specified constant  $C$  but the choice of  $C$  need not be sharp. Thus we will tend to round up to integers. Further, while  $\delta_1\delta_3$  and  $\delta_2\delta_3$  are absolutely tiny we will tend to bound them rather naively by  $|\epsilon_m/2|$ . Using these rules we have the bound

$$|\delta_4| \leq (1 + 1 + 2 + 1 + 1)\epsilon_m = 6\epsilon_m$$

Thus the computation becomes

$$((2.3 + \delta_4) + 1.3(1 + \delta_5))(1 + \delta_6) = 3.6 + \underbrace{\delta_4 + 1.3\delta_5 + 3.6\delta_6 + \delta_4\delta_6 + 1.3\delta_5\delta_6}_{\delta_7}$$

where the *absolute error* is

$$|\delta_7| \leq (6 + 1 + 2 + 1 + 1)\epsilon_m = 11\epsilon_m$$

Indeed, this bound is bigger than the observed error:

In [17]: `abs(3.6 - (1.1+1.2+1.3)), 11eps()`

Out[17]: (4.440892098500626e-16, 2.4424906541753444e-15)

## Arithmetic and special numbers

Arithmetic works differently on Inf and NaN and for undefined operations. In particular we have:

```
In [18]: 1/0.0      # Inf
          1/(-0.0)    # -Inf
          0.0/0.0      # NaN

          Inf*0       # NaN
          Inf+5       # Inf
          (-1)*Inf    # -Inf
          1/Inf        # 0.0
          1/(-Inf)     # -0.0
          Inf - Inf    # NaN
          Inf == Inf   # true
          Inf == -Inf  # false

          NaN*0       # NaN
          NaN+5       # NaN
          1/NaN        # NaN
          NaN == NaN   # false
          NaN != NaN   # true
```

Out [18]: true

## Special functions (non-examinable)

Other special functions like `cos`, `sin`, `exp`, etc. are *not* part of the IEEE standard. Instead, they are implemented by composing the basic arithmetic operations, which accumulate errors. Fortunately many are designed to have *relative accuracy*, that is,  $s = \sin(x)$  (that is, the Julia implementation of  $\sin x$ ) satisfies

$$s = (\sin x)(1 + \delta)$$

where  $|\delta| < c\epsilon_m$  for a reasonably small  $c > 0$ , provided that  $x \in F^{\text{normal}}$ . Note these special functions are written in (advanced) Julia code, for example, `sin` (<https://github.com/JuliaLang/julia/blob/d08b05df6f01cf4ec6e4c28ad94cedda76cc62e8/b/>

**WARNING (`sin(fl(x))` is not always close to `sin(x)`)** This is possibly a misleading statement when one thinks of  $x$  as a real number. Consider  $x = \pi$  so that  $\sin x = 0$ . However, as  $\text{fl}(\pi) \neq \pi$ . Thus we only have relative accuracy compared to the floating point approximation:

```
In [19]: π₆₄ = Float64(π)
          πᵢ = big(π₆₄) # Convert 64-bit approximation of π to higher precision
          abs(sin(π₆₄)), abs(sin(π₆₄) - sin(πᵢ)) # only has relative accuracy
```

Out [19]: (1.2246467991473532e-16, 2.994769809718339860754263822337778811430  
799841054596882794158676581342467643355e-33)

Another issue is when  $x$  is very large:

```
In [20]: ε = eps() # machine epsilon, 2^(-52)
          x = 2*10.0^100
          abs(sin(x) - sin(big(x))) ≤ abs(sin(big(x))) * ε
```

Out[20]: true

But if we instead compute  $10^{100}$  using `BigFloat` we get a completely different answer that even has the wrong sign!

```
In [21]: x = 2*big(10.0)^100  
sin(x), sin(x)
```

**Out[21]:** (-0.703969872087777, 0.6911910845037462219623751594978914260403966  
392716944990360937340001300242965408)

This is because we commit an error on the order of roughly

$$2 * 10^{100} * \epsilon_m \approx 4.44 * 10^{84}$$

when we round  $2 * 10^{100}$  to the nearest float.

**Example 7 (polynomial near root)** For general functions we do not generally have relative accuracy. For example, consider a simple polynomial  $1 + 4x + x^2$  which has a root at  $\sqrt{3} - 2$ . But

```
In [22]: f = x -> 1 + 4x + x^2
        x = sqrt(3) - 2
        aberr = abs(f(big(x)) - f(x))
        relerr = aberr/abs(f(x))
        aberr, relerr # very large relative error
```

We can see this in the error bound (note that  $4x$  is exact for floating point numbers and adding 1 is exact for this particular  $x$ ):

$$(x \otimes x) \oplus 4x + 1 = (x^2(1 + \delta_1) + 4x)(1 + \delta_2) + 1 = x^2 + 4x + 1 + \delta_1 x^2 + 4x\delta_2 +$$

Using a simple bound  $|x| < 1$  we get a (pessimistic) bound on the absolute error of  $3\epsilon_m$ . Here  $f(x)$  itself is less than  $2\epsilon_m$  so this does not imply relative accuracy. (Of course, a bad upper bound is not the same as a proof of inaccuracy, but here we observe the inaccuracy in practice.)

## **4. High-precision floating-point numbers (non-examinable)**

It is possible to set the precision of a floating-point number using the `BigFloat` type, which results from the usage of `big` when the result is not an integer. For example, here is an approximation of  $1/3$  accurate to 77 decimal digits:

```
In [23]: big(1)/3
```

Note we can set the rounding mode as in `Float64`, e.g., this gives (rigorous) bounds on  $1/3$ :

```
In [24]: setrounding(BigFloat, RoundDown) do
    big(1)/3
end, setrounding(BigFloat, RoundUp) do
    big(1)/3
end
```

We can also increase the precision, e.g., this finds bounds on  $1/3$  accurate to more than 1000 decimal places:

```
In [25]: setprecision(4_000) do # 4000 bit precision
    setrounding(BigFloat, RoundDown) do
        big(1)/3
    end, setrounding(BigFloat, RoundUp) do
        big(1)/3
    end
end
```

In the labs we shall see how this can be used to rigorously bound  $e$ , accurate to 1000 digits.



## I.3 Divided Differences

We now get to our first computational problem: given a function, how can we approximate its derivative at a point? We consider an intuitive approach to this problem using *Divided Differences*:

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

for  $h$  small. From the definition of a derivative we know that as  $h \rightarrow 0$  the approximation (the right-hand side) converges to  $f'(x)$ , hence a natural algorithm is to use this formula for very small  $h$ . Unfortunately, the round-off errors of floating point arithmetic introduces typically limit its accuracy.

In this chapter we view  $f$  as a *Black-box function* which we can only evaluate *pointwise*: Consider a floating-point valued function  $f^{\text{FP}} : D \rightarrow F$  where  $D \subset F \equiv F_{\sigma,Q,S}$ . For example, if we have `f(x::Float64) = sqrt(x)` then our function is defined only for positive `Float64` (otherwise it throws an error) hence  $D = [0, \infty) \cap F_{64}$ . This is the situation if we have a function that relies on a compiled C library, which hides the underlying operations (which are usually composition of floating point operations). Since the set of floating point numbers  $F$  is discrete,  $f^{\text{FP}}$  cannot be differentiable in an obvious way, therefore we need to assume that  $f^{\text{FP}}$  approximates a differentiable function  $f$  with controllable error bounds in order to state anything precise.

In the next chapter (I.4) we will introduce a more accurate approach based on *dual numbers*, which requires a more general notion of a function where a formula (i.e., code) is available.

Note there are other techniques for differentiation that we don't discuss:

1. Symbolic differentiation: A tree is built representing a formula which is differentiated using

the product and chain rule. 2. Adoints and back-propagation (reverse-mode automatic differentiation): This is similar to symbolic differentiation but automated, where the adjoints of Jacobians of each operation are constructed in a way that they can be composed (essentially the chain rule). It's outside the scope of this module but is computationally preferred for computing gradients of large dimensional functions which is critical in machine learning. 4. Interpolation and differentiation: We can also differentiate functions *globally*, that is, in an interval instead of only a single point, which will be discussed in Part III of the module.

In [1]: `using ColorBitstring`

# 1. Divided differences

The definition

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

tells us that

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

provided that  $h$  is sufficiently small.

It's important to note that approximation uses only the *black-box* notion of a function but to obtain bounds we need more.

If we know a bound on  $f''(x)$  then Taylor's theorem tells us a precise bound:

**Proposition 1** The error in approximating the derivative using divided differences is

$$\left| f'(x) - \frac{f(x+h) - f(x)}{h} \right| \leq \frac{M}{2} h$$

where  $M = \sup_{x \leq t \leq x+h} |f''(t)|$ .

**Proof** Follows immediately from Taylor's theorem:

$$f(x+h) = f(x) + f'(x)h + \frac{f''(t)}{2}h^2$$

for some  $x \leq t \leq x+h$ .



There are also alternative versions of divided differences. Leftside divided differences:

$$f'(x) \approx \frac{f(x) - f(x-h)}{h}$$

and central differences:

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$

Composing these approximations is useful for higher-order derivatives as we discuss in the problem sheet.

Note this is assuming *real arithmetic*, the answer is drastically different with *floating point arithmetic*.

## 2. Does divided differences work with floating point arithmetic?

Let's try differentiating two simple polynomials  $f(x) = 1 + x + x^2$  and  $g(x) = 1 + x/3 + x^2$  by applying the divided difference approximation to their floating point implementations  $f^{\text{FP}}$  and  $g^{\text{FP}}$ :

```
In [2]: f = x -> 1 + x + x^2      # we treat f and g as black-boxes
g = x -> 1 + x/3 + x^2
h = 0.000001
(f(h)-f(0))/h, (g(h)-g(0))/h
```

```
Out[2]: (1.000001000006634, 0.33333433346882657)
```

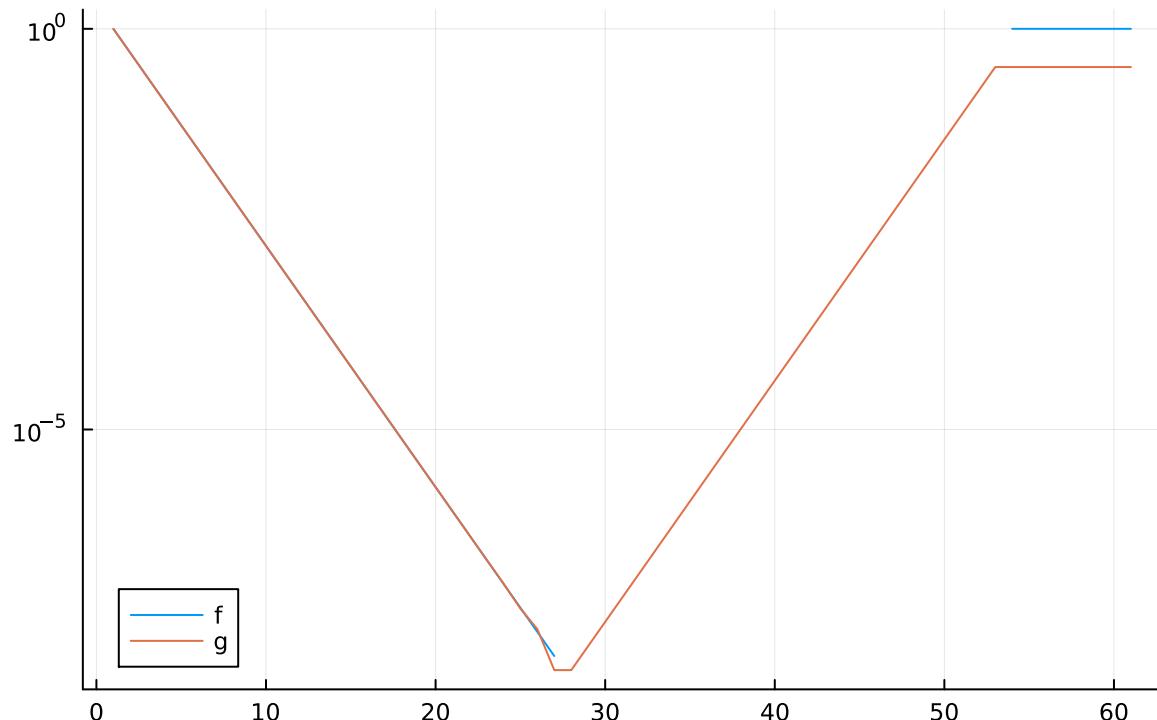
Both seem to roughly approximate the true derivatives (1 and  $1/3$ ). We can do a plot to see how fast the error goes down as we let  $h$  become small.

```
In [3]: using Plots
h = 2.0 .^ (0:-1:-60) # [1, 1/2, 1/4, ...]
nanabs = x -> iszero(x) ? NaN : abs(x) # avoid 0's in log scale plot
plot(nanabs.((f.(h) .- f(0)) ./ h .- 1),yscale=:log10, title="convergence of f")
plot!(abs.((g.(h) .- g(0)) ./ h .- 1/3),yscale=:log10, label = "g")
```

```
[ Info: Precompiling Plots [91a5bcdd-55d7-5caf-9e0b-520d859cae80]
[ Info: GR
```

```
Out[3]:
```

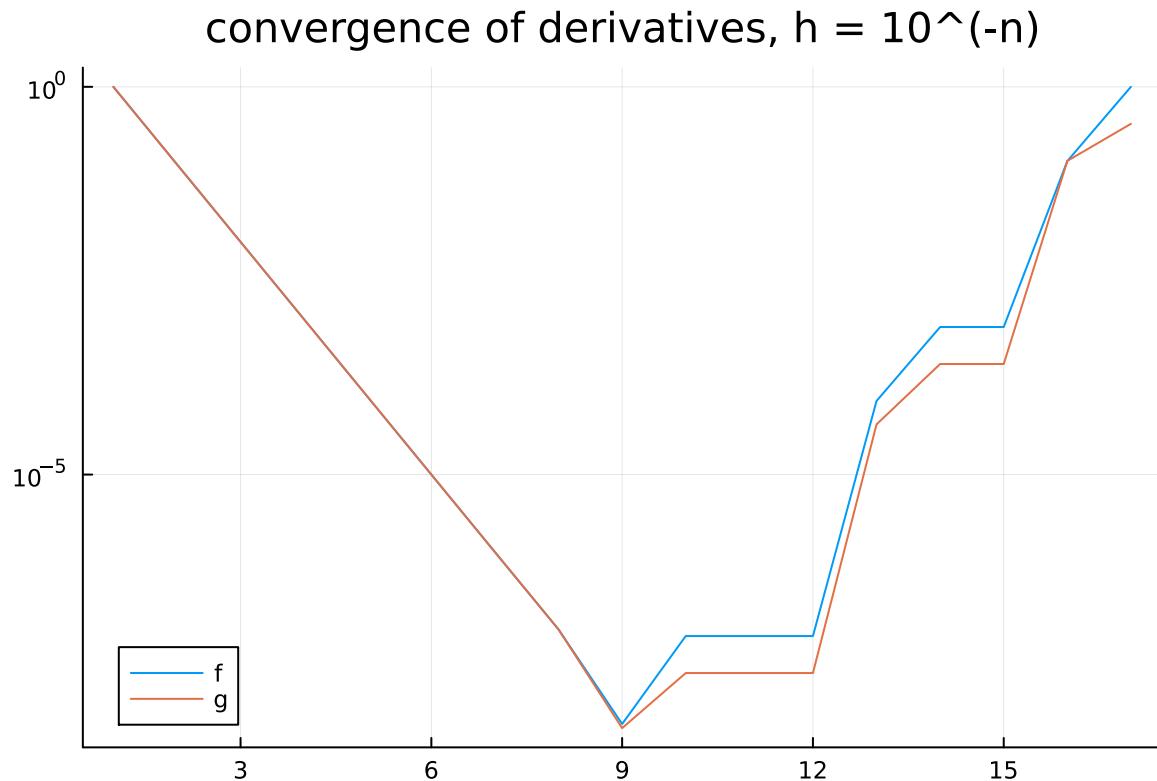
convergence of derivatives,  $h = 2^{-n}$



It is clear that  $f$  is extremely special. Most functions will behave like  $g$ , and had we not taken  $h$  to be a power of two we also see divergence for differentiating  $f$ :

```
In [4]: h = 10.0 .^ (0:-1:-16) # [1, 1/10, 1/100,...]
plot(abs.((f.(h) .- f(0)) ./ h .- 1);yscale=:log10, title="convergence of d
plot!(abs.((g.(h) .- g(0)) ./ h .- 1/3);yscale=:log10, label = "g")
```

Out [4]:



For these two simple examples, we can understand why we see very different behaviour.

**Example 1 (convergence?) of divided differences** Consider differentiating  $f(x) = 1 + x + x^2$  at 0 with  $h = 2^{-n}$ . We consider 4 different cases with different behaviour, where  $S$  is the number of significand bits:

1.  $0 \leq n \leq S/2$
2.  $S/2 < n \leq S$
3.  $S \leq n < S + \sigma$
4.  $S + \sigma \leq n$

Note that  $f^{\text{FP}}(0) = f(0) = 1$ . Thus we wish to understand the error in approximating  $f'(0) = 1$  by

$$(f^{\text{FP}}(h) \ominus 1) \oslash h \quad \text{where} \quad f^{\text{FP}}(x) = 1 \oplus x \oplus x \otimes x.$$

**Case 1** ( $0 \leq n \leq S/2$ ): note that  $f^{\text{FP}}(h) = f(h) = 1 + 2^{-n} + 2^{-2n}$  as each computation is precisely a floating point number (hence no rounding). We can see this in half-precision, with  $n = 3$  we have a 1 in the 3rd and 6th decimal place:

```
In [5]: S = 10 # 10 significant bits
n = 3 # 3 ≤ S/2 = 5
h = Float16(2)^(-n)
printbits(f(h))
```

`0011110010010000`

Subtracting 1 and dividing by  $h$  will also be exact, hence we get

$$(f^{\text{FP}}(h) \ominus 1) \oslash h = 1 + 2^{-n}$$

which shows exponential convergence.

*Case 2 ( $S/2 < n \leq S$ ):* Now we have (using round-to-nearest)

$$f^{\text{FP}}(h) = (1 + 2^{-n}) \oplus 2^{-2n} = 1 + 2^{-n}$$

Then

$$(f^{\text{FP}}(h) \ominus 1) \oslash h = 1 = f'(0)$$

We have actually performed better than true real arithmetic and converged without a limit!

*Case 3 ( $S < n < \sigma + S$ ):* If we take  $n$  too large, then  $1 \oplus h = 1$  and we have  $f^{\text{FP}}(h) = 1$ , that is and

$$(f^{\text{FP}}(h) \ominus 1) \oslash h = 0 \neq f'(0)$$

*Case 4 ( $\sigma + S \leq n$ ):* In this case  $h = 2^{-n}$  is rounded to zero and hence we get *NaN*.

**Example 2 (divergence of divided differences)** Consider differentiating  $g(x) = 1 + x/3 + x^2$  at 0 with  $h = 2^{-n}$  and assume  $n$  is even for simplicity and consider half-precision with  $S = 10$ . Note that  $g^{\text{FP}}(0) = g(0) = 1$ . Recall

$$h \oslash 3 = 2^{-n-2} * (1.0101010101)_2$$

Note we lose two bits each time in the computation of  $1 \oplus (h \oslash 3)$ :

```
In [6]: n = 0; h = Float16(2)^(-n); printlnbits(1 + h/3)
n = 2; h = Float16(2)^(-n); printlnbits(1 + h/3)
n = 4; h = Float16(2)^(-n); printlnbits(1 + h/3)
n = 8; h = Float16(2)^(-n); printlnbits(1 + h/3)
```

`0011110101010101  
0011110001010101  
0011110000010101  
0011110000000001`

It follows if  $S/2 < n < S$  that

$$1 \oplus (h \oslash 3) = 1 + h/3 - 2^{-10}/3$$

Therefore

$$(g^{\text{FP}}(h) \ominus 1) \oslash h = 1/3 - 2^{n-10}/3$$

Thus the error grows exponentially with  $n$ .

If  $S \leq n < \sigma + S$  then  $1 \oplus (h \oslash 3) = 1$  and we have

$$(g^{\text{FP}}(h) \ominus 1) \oslash h = 0$$

### 3. Bounding the error

We can bound the error using the bounds on floating point arithmetic.

**Theorem 1 (divided difference error bound)** Let  $f$  be twice-differentiable in a neighbourhood of  $x$  and assume that

$$f^{\text{FP}}(x) = f(x) + \delta_x^f$$

has uniform absolute accuracy in that neighbourhood, that is:

$$|\delta_x^f| \leq c\epsilon_m$$

for a fixed constant  $c \geq 0$ . Assume for simplicity  $h = 2^{-n}$  where  $n \leq S$  and  $|x| \leq 1$ . Assuming that all calculations result in normal floating point numbers, the divided difference approximation satisfies

$$(f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)) \oslash h = f'(x) + \delta_{x,h}^{\text{FD}}$$

where

$$|\delta_{x,h}^{\text{FD}}| \leq \frac{|f'(x)|}{2}\epsilon_m + Mh + \frac{4c\epsilon_m}{h}$$

for  $M = \sup_{x \leq t \leq x+h} |f''(t)|$ .

#### Proof

We have (noting by our assumptions  $x \oplus h = x + h$  and that dividing by  $h$  will only change the exponent so is exact)

$$\begin{aligned} (f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)) \oslash h &= \frac{f(x+h) + \delta_{x+h}^f - f(x) - \delta_x^f}{h} (1 + \delta_1) \\ &= \frac{f(x+h) - f(x)}{h} (1 + \delta_1) + \frac{\delta_{x+h}^f - \delta_x^f}{h} (1 + \delta_1) \end{aligned}$$

where  $|\delta_1| \leq \epsilon_m/2$ . Applying Taylor's theorem we get

$$(f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)) \oslash h = f'(x) + f'(x)\delta_1 + \underbrace{\frac{f''(t)}{2}h(1+\delta_1)}_{\delta_{x,h}^{\text{FD}}} + \frac{\delta_{x+h}^f - \delta_x^f}{h}(1+\delta_1)$$

The bound then follows, using the very pessimistic bound  $|1 + \delta_1| \leq 2$ .

■

The three-terms of this bound tell us a story: the first term is a fixed (small) error, the second term tends to zero as  $h \rightarrow 0$ , while the last term grows like  $\epsilon_m/h$  as  $h \rightarrow 0$ . Thus we observe convergence while the second term dominates, until the last term takes over. Of course, a bad upper bound is not the same as a proof that something grows, but it is a good indication of what happens *in general* and suffices to motivate the following heuristic to balance the two sources of errors:

**Heuristic (divided difference with floating-point step)** Choose  $h$  proportional to  $\sqrt{\epsilon_m}$  in divided differences.

In the case of double precision  $\sqrt{\epsilon_m} \approx 1.5 \times 10^{-8}$ , which is close to when the observed error begins to increase in our examples.

**Remark** While divided differences is of debatable utility for computing derivatives, it is extremely effective in building methods for solving differential equations, as we shall see later. It is also very useful as a "sanity check" if one wants something to compare with for other numerical methods for differentiation.

## I.4 Dual Numbers

In this chapter we introduce a mathematically beautiful alternative to divided differences for computing derivatives: *dual numbers*. As we shall see, these are a commutative ring that *exactly* compute derivatives, which when implemented in floating point give very high-accuracy approximations to derivatives. They underpin forward-mode [automatic differentiation](#).

Before we begin, we must be clear what a "function" is. Consider three possible scenarios:

1. *Black-box function*: this was considered last chapter where we can only input floating-point numbers

and get out a floating-point number. 2. *Generic function*: Consider a function that is a formula (or, equivalently, a *piece of code*) that we can evaluate on arbitrary types, including custom types that we create. An example is a polynomial:

$$p(x) = p_0 + p_1x + \cdots + p_nx^n$$

which can be evaluated for  $x$  in the reals, complexes, or any other ring. More generally, if we have a function defined in Julia that does not call any C libraries it can be evaluated on different types. 3. *Graph function*: The function is built by composing different basic "kernels" with known differentiability properties. We won't consider this situation in this module, though it is the model used by Python machine learning toolbox's like [PyTorch](#) and [TensorFlow](#).

**Definition 1 (Dual numbers)** Dual numbers  $\mathbb{D}$  are a commutative ring (over  $\mathbb{R}$ ) generated by 1 and  $\epsilon$  such that  $\epsilon^2 = 0$ . Dual numbers are typically written as  $a + b\epsilon$  where  $a$  and  $b$  are real.

This is very much analogous to complex numbers, which are a field generated by 1 and  $i$  such that  $i^2 = -1$ . Compare multiplication of each number type:

$$\begin{aligned}(a + bi)(c + di) &= ac + (bc + ad)i + bdi^2 = ac - bd + (bc + ad)i \\(a + b\epsilon)(c + d\epsilon) &= ac + (bc + ad)\epsilon + bde^2 = ac + (bc + ad)\epsilon\end{aligned}$$

And just as we view  $\mathbb{R} \subset \mathbb{C}$  by equating  $a \in \mathbb{R}$  with  $a + 0i \in \mathbb{C}$ , we can view  $\mathbb{R} \subset \mathbb{D}$  by equating  $a \in \mathbb{R}$  with  $a + 0\epsilon \in \mathbb{D}$ .

### 1. Differentiating polynomials

Applying a polynomial to a dual number  $a + b\epsilon$  tells us the derivative at  $a$ :

**Theorem 1 (polynomials on dual numbers)** Suppose  $p$  is a polynomial. Then

$$p(a + b\epsilon) = p(a) + bp'(a)\epsilon$$

### Proof

It suffices to consider  $p(x) = x^n$  for  $n \geq 1$  as other polynomials follow from linearity.

We proceed by induction: The case  $n = 1$  is trivial. For  $n > 1$  we have

$$(a + b\epsilon)^n = (a + b\epsilon)(a + b\epsilon)^{n-1} = (a + b\epsilon)(a^{n-1} + (n-1)ba^{n-2}\epsilon) = a^n + bna^{n-1}\epsilon.$$

■

**Example 1 (differentiating polynomial)** Consider computing  $p'(2)$  where

$$p(x) = (x - 1)(x - 2) + x^2.$$

We can use Dual numbers to differentiating, avoiding expanding in monomials or rules of differentiating:

$$p(2 + \epsilon) = (1 + \epsilon)\epsilon + (2 + \epsilon)^2 = \epsilon + 4 + 4\epsilon = 4 + \underbrace{5}_{p'(2)}\epsilon$$

## 2. Differentiating other functions

We can extend real-valued differentiable functions to dual numbers in a similar manner. First, consider a standard function with a Taylor series (e.g.  $\cos$ ,  $\sin$ ,  $\exp$ , etc.)

$$f(x) = \sum_{k=0}^{\infty} f_k x^k$$

so that  $a$  is inside the radius of convergence. This leads naturally to a definition on dual numbers:

$$\begin{aligned} f(a + b\epsilon) &= \sum_{k=0}^{\infty} f_k (a + b\epsilon)^k = f_0 + \sum_{k=1}^{\infty} f_k (a^k + ka^{k-1}b\epsilon) = \sum_{k=0}^{\infty} f_k a^k + \sum_{k=1}^{\infty} f_k a^{k-1} b \\ &= f(a) + bf'(a)\epsilon \end{aligned}$$

More generally, given a differentiable function we can extend it to dual numbers:

**Definition 2 (dual extension)** Suppose a real-valued function  $f$  is differentiable at  $a$ . If

$$f(a + b\epsilon) = f(a) + bf'(a)\epsilon$$

then we say that it is a *dual extension at  $a$* .

Thus, for basic functions we have natural extensions:

$$\begin{aligned}
\exp(a + b\epsilon) &:= \exp(a) + b \exp(a)\epsilon \\
\sin(a + b\epsilon) &:= \sin(a) + b \cos(a)\epsilon \\
\cos(a + b\epsilon) &:= \cos(a) - b \sin(a)\epsilon \\
\log(a + b\epsilon) &:= \log(a) + \frac{b}{a}\epsilon \\
\sqrt{a + b\epsilon} &:= \sqrt{a} + \frac{b}{2\sqrt{a}}\epsilon \\
|a + b\epsilon| &:= |a| + b \operatorname{sign} a \epsilon
\end{aligned}$$

provided the function is differentiable at  $a$ . Note the last example does not have a convergent Taylor series (at 0) but we can still extend it where it is differentiable.

Going further, we can add, multiply, and compose such functions:

**Lemma 1 (product and chain rule)** If  $f$  is a dual extension at  $g(a)$  and  $g$  is a dual extension at  $a$ , then  $q(x) := f(g(x))$  is a dual extension at  $a$ . If  $f$  and  $g$  are dual extensions at  $a$  then  $r(x) := f(x)g(x)$  is also dual extensions at  $a$ . In other words:

$$\begin{aligned}
q(a + b\epsilon) &= q(a) + bq'(a)\epsilon \\
r(a + b\epsilon) &= r(a) + br'(a)\epsilon
\end{aligned}$$

**Proof** For  $q$  it follows immediately:

$$\begin{aligned}
q(a + b\epsilon) &= f(g(a + b\epsilon)) = f(g(a) + bg'(a)\epsilon) \\
&= f(g(a)) + bg'(a)f'(g(a))\epsilon = q(a) + bq'(a)\epsilon.
\end{aligned}$$

For  $r$  we have

$$\begin{aligned}
r(a + b\epsilon) &= f(a + b\epsilon)g(a + b\epsilon) = (f(a) + bf'(a)\epsilon)(g(a) + bg'(a)\epsilon) \\
&= f(a)g(a) + b(f'(a)g(a) + f(a)g'(a))\epsilon = r(a) + br'(a)\epsilon.
\end{aligned}$$

■

A simple corollary is that any function defined in terms of addition, multiplication, composition, etc. of functions that are dual with differentiation will be differentiable via dual numbers.

### Example 2 (differentiating non-polynomial)

Consider  $f(x) = \exp(x^2 + e^x)$  by evaluating on the duals:

$$f(1 + \epsilon) = \exp(1 + 2\epsilon + e + e\epsilon) = \exp(1 + e) + \exp(1 + e)(2 + e)\epsilon$$

and therefore we deduce that

$$f'(1) = \exp(1 + e)(2 + e).$$

## 3. Implementation as a special type

We now consider a simple implementation of dual numbers that works on general polynomials:

```
In [1]: # Dual(a,b) represents a + b*ε
struct Dual{T}
    a::T
    b::T
end

# Dual(a) represents a + 0*ε
Dual(a::Real) = Dual(a, zero(a)) # for real numbers we use a + 0ε

# Allow for a + b*ε syntax
const ε = Dual(0, 1)

import Base: +, *, -, /, ^, zero, exp

# support polynomials like 1 + x, x - 1, 2x or x*2 by reducing to Dual
+(x::Real, y::Dual) = Dual(x) + y
+(x::Dual, y::Real) = x + Dual(y)
-(x::Real, y::Dual) = Dual(x) - y
-(x::Dual, y::Real) = x - Dual(y)
*(x::Real, y::Dual) = Dual(x) * y
*(x::Dual, y::Real) = x * Dual(y)

# support x/2 (but not yet division of duals)
/(x::Dual, k::Real) = Dual(x.a/k, x.b/k)

# a simple recursive function to support x^2, x^3, etc.
function ^(x::Dual, k::Integer)
    if k < 0
        error("Not implemented")
    elseif k == 1
        x
    else
        x^(k-1) * x
    end
end

# Algebraic operations for duals
-(x::Dual) = Dual(-x.a, -x.b)
+(x::Dual, y::Dual) = Dual(x.a + y.a, x.b + y.b)
-(x::Dual, y::Dual) = Dual(x.a - y.a, x.b - y.b)
*(x::Dual, y::Dual) = Dual(x.a*y.a, x.a*y.b + x.b*y.a)

exp(x::Dual) = Dual(exp(x.a), exp(x.a) * x.b)
```

```
Out[1]: exp (generic function with 15 methods)
```

We can also try it on the two polynomials as above:

```
In [2]: f = x -> 1 + x + x^2
g = x -> 1 + x/3 + x^2
f(ε).b, g(ε).b
```

```
Out[2]: (1, 0.333333333333333)
```

The first example exactly computes the derivative, and the second example is exact up to the last bit rounding! It also works for higher order polynomials:

```
In [3]: f = x -> 1 + 1.3x + 2.1x^2 + 3.1x^3  
f(0.5 + ε).b - 5.725
```

```
Out[3]: 8.881784197001252e-16
```

It is indeed "accurate to (roughly) 16-digits", the best we can hope for using floating point.

We can use this in "algorithms" as well as simple polynomials. Consider the polynomial  $1 + \dots + x^n$ :

```
In [4]: function s(n, x)  
    ret = 1 + x # first two terms  
    for k = 2:n  
        ret += x^k  
    end  
    ret  
end  
s(10, 0.1 + ε).b
```

```
Out[4]: 1.2345678999999998
```

This matches exactly the "true" (up to rounding) derivative:

```
In [5]: sum((1:10) .* 0.1 .^(0:9))
```

```
Out[5]: 1.2345678999999998
```

Finally, we can try the more complicated example:

```
In [6]: f = x -> exp(x^2 + exp(x))  
f(1 + ε)
```

```
Out[6]: Dual{Float64}(41.193555674716116, 194.362805189629)
```

What makes dual numbers so effective is that, unlike divided differences, they are not prone to disastrous growth due to round-off errors.

## II.1 Structured Matrices

We have seen how algebraic operations (`+`, `-`, `*`, `/`) are defined exactly in terms of rounding ( $\oplus$ ,  $\ominus$ ,  $\otimes$ ,  $\oslash$ ) for floating point numbers. Now we see how this allows us to do (approximate) linear algebra operations on matrices.

A matrix can be stored in different formats. Here we consider the following structures:

1. *Dense*: This can be considered unstructured, where we need to store all entries in a vector or matrix. Matrix multiplication reduces directly to standard algebraic operations. Solving linear systems with dense matrices will be discussed later. 2. *Triangular*: If a matrix is upper or lower triangular, we can immediately invert using back-substitution. In practice we store a dense matrix and ignore the upper/lower entries. 3. *Banded*: If a matrix is zero apart from entries a fixed distance from the diagonal it is called banded and this allows for more efficient algorithms. We discuss diagonal, tridiagonal and bidiagonal matrices.

In the next chapter we consider more complicated orthogonal matrices.

```
In [1]: # LinearAlgebra contains routines for doing linear algebra
# BenchmarkTools is a package for reliable timing
using LinearAlgebra, Plots, BenchmarkTools, Test
```

### 1. Dense vectors and matrices

A `Vector` of a primitive type (like `Int` or `Float64`) is stored consecutively in memory: that is, a vector consists of a memory address (a *pointer*) to the first entry and a length. E.g. if we have a `Vector{Int8}` of length `n` then it is stored as  $8n$  bits (`n` bytes) in a row. That is, if the memory address of the first entry is `k` and the type is `T`, the memory address of the second entry is `k + sizeof(T)`.

**Remark (advanced)** We can actually experiment with this (NEVER DO THIS IN PRACTICE!!), beginning with an 8-bit type:

```
In [2]: a = Int8[2, 4, 5]
p = pointer(a) # pointer(a) returns memory address of the first entry, which
# We can think of a pointer as simply a UInt64 alongside a Type to interpret
```

```
Out[2]: Ptr{Int8} @0x000000014539fd78
```

We can see what's stored at a pointer as follows:

```
In [3]: Base.unsafe_load(p) # loads data at `p`. Knows its an `Int8` because of type
```

```
Out[3]: 2
```

Adding an integer to a pointer gives a new pointer with the address incremented:

```
In [4]: p + 1 # memory address of next entry, which is 1 more than first
```

```
Out[4]: Ptr{Int8} @0x000000014539fd79
```

We see that this gives us the next entry:

```
In [5]: Base.unsafe_load(p) # loads data at `p+1`, which is second entry of the vect
```

```
Out[5]: 2
```

For other types we need to increment the address by the size of the type:

```
In [6]: a = [2.0, 1.3, 1.4]
p = pointer(a)
Base.unsafe_load(p + 8) # sizeof(Float64) == 8
```

```
Out[6]: 1.3
```

Why not do this in practice? It's unsafe because there's nothing stopping us from going past the end of an array:

```
In [7]: Base.unsafe_load(p + 3 * 8) # whatever bits happened to be next in memory, u
```

```
Out[7]: 2.405888678e-314
```

This may even crash Julia! (I got lucky that it didn't when producing the notes.)

---

A `Matrix` is stored consecutively in memory, going down column-by- column (*column-major*). That is,

```
In [8]: A = [1 2;
            3 4;
            5 6]
```

```
Out[8]: 3×2 Matrix{Int64}:
         1 2
         3 4
         5 6
```

Is actually stored equivalently to a length 6 vector:

```
In [9]: vec(A)
```

```
Out[9]: 6-element Vector{Int64}:
```

```
1  
3  
5  
2  
4  
6
```

which in this case would be stored using in  $8 * 6 = 48$  consecutive memory addresses. That is, a matrix is a pointer to the first entry alongside two integers dictating the row and column sizes.

---

**Remark (advanced)** Note that transposing `A` is done lazily and so `transpose(A)` (which is equivalent to the adjoint/conjugate-transpose `A'` when the entries are real), is just a special type with a single field: `transpose(A).parent == A`. This is equivalent to *row-major* format, where the next address in memory of `transpose(A)` corresponds to moving along the row.

---

Matrix-vector multiplication works as expected:

```
In [10]: x = [7, 8]  
A * x
```

```
Out[10]: 3-element Vector{Int64}:
```

```
23  
53  
83
```

Note there are two ways this can be implemented:

**Algorithm 1 (matrix-vector multiplication by rows)** For a ring  $R$  (typically  $\mathbb{R}$  or  $\mathbb{C}$ ),  $A \in R^{m \times n}$  and  $\mathbf{x} \in R^n$  we have

$$A\mathbf{x} = \begin{bmatrix} \sum_{j=1}^n a_{1,j}x_j \\ \vdots \\ \sum_{j=1}^n a_{m,j}x_j \end{bmatrix}.$$

In code this can be implemented for any types that support `*` and `+` as follows:

```
In [11]: function mul_rows(A, x)  
    m,n = size(A)  
    # promote_type type finds a type that is compatible with both types, elt  
    T = promote_type(eltype(x), eltype(A))  
    c = zeros(T, m) # the returned vector, begins of all zeros  
    for k = 1:m, j = 1:n  
        c[k] += A[k, j] * x[j] # equivalent to c[k] = c[k] + A[k, j] * x[j]  
    end  
    c  
end
```

```
Out[11]: mul_rows (generic function with 1 method)
```

**Algorithm 2 (matrix-vector multiplication by columns)** For a ring  $R$  (typically  $\mathbb{R}$  or  $\mathbb{C}$ ),  $A \in R^{m \times n}$  and  $\mathbf{x} \in R^n$  we have

$$A\mathbf{x} = x_1\mathbf{a}_1 + \cdots + x_n\mathbf{a}_n$$

where  $\mathbf{a}_j := A\mathbf{e}_j \in R^m$  (that is, the  $j$ -th column of  $A$ ). In code this can be implemented for any types that support `*` and `+` as follows:

```
In [12]: function mul_cols(A, x)
    m,n = size(A)
    # promote_type type finds a type that is compatible with both types, elt
    T = promote_type(eltype(x),eltype(A))
    c = zeros(T, m) # the returned vector, begins of all zeros
    for j = 1:n, k = 1:m
        c[k] += A[k, j] * x[j] # equivalent to c[k] = c[k] + A[k, j] * x[j]
    end
    c
end
```

```
Out[12]: mul_cols (generic function with 1 method)
```

Both implementations match exactly for integer inputs:

```
In [13]: mul_rows(A, x), mul_cols(A, x) # also matches `A*x`
```

```
Out[13]: ([23, 53, 83], [23, 53, 83])
```

Either implementation will be  $O(mn)$  operations. However, the implementation `mul_cols` accesses the entries of `A` going down the column, which happens to be *significantly faster* than `mul_rows`, due to accessing memory of `A` in order. We can see this by measuring the time it takes using `@btime`:

```
In [14]: n = 1000
A = randn(n,n) # create n x n matrix with random normal entries
x = randn(n) # create length n vector with random normal entries

@btime mul_rows(A,x)
@btime mul_cols(A,x)
@btime A*x; # built-in, high performance implementation. USE THIS in practice
```

1.667 ms (1 allocation: 7.94 KiB)  
755.646 μs (1 allocation: 7.94 KiB)  
220.887 μs (1 allocation: 7.94 KiB)

Here `ms` means milliseconds ( $0.001 = 10^{-3}$  seconds) and `μs` means microseconds ( $0.000001 = 10^{-6}$  seconds). So we observe that `mul` is roughly 3x faster than `mul_rows`, while the optimised `*` is roughly 5x faster than `mul`.

**Remark (advanced)** For floating point types, `A*x` is implemented in BLAS which is generally multi-threaded and is not identical to `mul_cols(A,x)`, that is, some inputs will differ in how the computations are rounded.

---

Note that the rules of floating point arithmetic apply here: matrix multiplication with floats will incur round-off error (the precise details of which are subject to the implementation):

```
In [15]: A = [1.4 0.4;
            2.0 1/2]
A * [1, -1] # First entry has round-off error, but 2nd entry is exact
```

```
Out[15]: 2-element Vector{Float64}:
          0.9999999999999999
          1.5
```

And integer arithmetic will be subject to overflow:

```
In [16]: A = fill(Int8(2^6), 2, 2) # make a matrix whose entries are all equal to 2^6
A * Int8[1,1] # we have overflowed and get a negative number -2^7
```

```
Out[16]: 2-element Vector{Int8}:
          -128
          -128
```

Solving a linear system is done using `\`:

```
In [17]: A = [1 2 3;
            1 2 4;
            3 7 8]
b = [10; 11; 12]
A \ b
```

```
Out[17]: 3-element Vector{Float64}:
          41.00000000000036
          -17.00000000000014
          1.0
```

Despite the answer being integer-valued, here we see that it resorted to using floating point arithmetic, incurring rounding error. But it is "accurate to (roughly) 16-digits". As we shall see, the way solving a linear system works is we first write `A` as a product of matrices that are easy to invert, e.g., a product of triangular matrices or a product of an orthogonal and triangular matrix.

## 2. Triangular matrices

Triangular matrices are represented by dense square matrices where the entries below the diagonal are ignored:

```
In [18]: A = [1 2 3;
            4 5 6;
            7 8 9]
U = UpperTriangular(A)
```

```
Out[18]: 3x3 UpperTriangular{Int64, Matrix{Int64}}:
 1 2 3
 · 5 6
 · · 9
```

We can see that `U` is storing all the entries of `A` in a field called `data`:

```
In [19]: U.data
```

```
Out[19]: 3x3 Matrix{Int64}:
 1 2 3
 4 5 6
 7 8 9
```

Similarly we can create a lower triangular matrix by ignoring the entries above the diagonal:

```
In [20]: L = LowerTriangular(A)
```

```
Out[20]: 3x3 LowerTriangular{Int64, Matrix{Int64}}:
 1 · ·
 4 5 ·
 7 8 9
```

If we know a matrix is triangular we can do matrix-vector multiplication in roughly half the number of operations by skipping over the entries we know are zero:

### Algorithm 3 (upper-triangular matrix-vector multiplication by columns)

```
In [21]: function mul_cols(U::UpperTriangular, x)
    n = size(U,1)
    # promote_type type finds a type that is compatible with both types, elt
    T = promote_type(eltype(x),eltype(U))
    b = zeros(T, n) # the returned vector, begins of all zeros
    for j = 1:n, k = 1:j # k = 1:j instead of 1:m since we know U[k,j] = 0 i
        b[k] += U[k, j] * x[j]
    end
    b
end

x = [10, 11, 12]
# matches built-in *
@test mul_cols(U, x) == U*x
```

```
Out[21]: Test Passed
```

Moreover, we can easily invert matrices. Consider a simple  $3 \times 3$  example, which can be solved with `\`:

```
In [22]: b = [5, 6, 7]
x = U \ b # Excercise: why does this return a float vector?
```

```
Out[22]: 3-element Vector{Float64}:
2.133333333333333
0.2666666666666666
0.7777777777777778
```

Behind the scenes, `\` is doing back-substitution: considering the last row, we have all zeros apart from the last column so we know that `x[3]` must be equal to:

```
In [23]: b[3] / U[3,3]
```

```
Out[23]: 0.7777777777777778
```

Once we know `x[3]`, the second row states `U[2,2]*x[2] + U[2,3]*x[3] == b[2]`, rearranging we get that `x[2]` must be:

```
In [24]: (b[2] - U[2,3]*x[3])/U[2,2]
```

```
Out[24]: 0.2666666666666666
```

Finally, the first row states `U[1,1]*x[1] + U[1,2]*x[2] + U[1,3]*x[3] == b[1]` i.e. `x[1]` is equal to

```
In [25]: (b[1] - U[1,2]*x[2] - U[1,3]*x[3])/U[1,1]
```

```
Out[25]: 2.133333333333333
```

More generally, we can solve the upper-triangular system using *back-substitution*:

**Algorithm 4 (back-substitution)** Let  $\mathbb{F}$  be a field (typically  $\mathbb{R}$  or  $\mathbb{C}$ ). Suppose  $U \in \mathbb{F}^{n \times n}$  is upper-triangular and invertible. Then for  $\mathbf{b} \in \mathbb{F}^n$  the solution  $\mathbf{x} \in \mathbb{F}^n$  to  $U\mathbf{x} = \mathbf{b}$ , that is,

$$\begin{bmatrix} u_{11} & \cdots & u_{1n} \\ \ddots & \ddots & \vdots \\ & u_{nn} & \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

is given by computing  $x_n, x_{n-1}, \dots, x_1$  via:

$$x_k = \frac{b_k - \sum_{j=k+1}^n u_{kj}x_j}{u_{kk}}$$

In code this can be implemented for any types that support `*`, `+` and `/` as follows:

```
In [26]: # ldiv(U, b) is our implementation of U\b
function ldiv(U::UpperTriangular, b)
    n = size(U,1)
```

```

if length(b) != n
    error("The system is not compatible")
end

x = zeros(n) # the solution vector

for k = n:-1:1 # start with k=n, then k=n-1, ...
    r = b[k] # dummy variable
    for j = k+1:n
        r -= U[k,j]*x[j] # equivalent to r = r - U[k,j]*x[j]
    end
    # after this for loop, r = b[k] - \sum_{j=k+1}^n U[k,j]*x[j]
    x[k] = r/U[k,k]
end
x
end

@test ldiv(U, x) ≈ U\x

```

Out[26]: **Test Passed**

The problem sheet will explore implementing multiplication and forward substitution for lower triangular matrices. The cost of multiplying and solving linear systems with a triangular matrix is  $O(n^2)$ .

---

### 3. Banded matrices

A *banded matrix* is zero off a prescribed number of diagonals. We call the number of (potentially) non-zero diagonals the *bandwidths*:

**Definition 1 (bandwidths)** A matrix  $A$  has *lower-bandwidth*  $l$  if  $A[k, j] = 0$  for all  $k - j > l$  and *upper-bandwidth*  $u$  if  $A[k, j] = 0$  for all  $j - k > u$ . We say that it has *strictly lower-bandwidth*  $l$  if it has lower-bandwidth  $l$  and there exists a  $j$  such that  $A[j + l, j] \neq 0$ . We say that it has *strictly upper-bandwidth*  $u$  if it has upper-bandwidth  $u$  and there exists a  $k$  such that  $A[k, k + u] \neq 0$ .

#### Diagonal

**Definition 2 (Diagonal)** *Diagonal matrices* are square matrices with bandwidths  $l = u = 0$ .

Diagonal matrices in Julia are stored as a vector containing the diagonal entries:

In [27]:

```

x = [1,2,3]
D = Diagonal(x) # the type Diagonal has a single field: D.diag

```

```
Out[27]: 3x3 Diagonal{Int64, Vector{Int64}}:
```

```
1 . .
. 2 .
. . 3
```

It is clear that we can perform diagonal-vector multiplications and solve linear systems involving diagonal matrices efficiently (in  $O(n)$  operations).

## Bidiagonal

**Definition 3 (Bidiagonal)** If a square matrix has bandwidths  $(l, u) = (1, 0)$  it is *lower-bidiagonal* and if it has bandwidths  $(l, u) = (0, 1)$  it is *upper-bidiagonal*.

We can create Bidiagonal matrices in Julia by specifying the diagonal and off-diagonal:

```
In [28]: L = Bidiagonal([1,2,3], [4,5], :L) # the type Bidiagonal has three fields: L
```

```
Out[28]: 3x3 Bidiagonal{Int64, Vector{Int64}}:
```

```
1 . .
4 2 .
. 5 3
```

```
In [29]: Bidiagonal([1,2,3], [4,5], :U)
```

```
Out[29]: 3x3 Bidiagonal{Int64, Vector{Int64}}:
```

```
1 4 .
. 2 5
. . 3
```

Multiplication and solving linear systems with Bidiagonal systems is also  $O(n)$  operations, using the standard multiplications/back-substitution algorithms but being careful in the loops to only access the non-zero entries.

## Tridiagonal

**Definition 4 (Tridiagonal)** If a square matrix has bandwidths  $l = u = 1$  it is *tridiagonal*.

Julia has a type `Tridiagonal` for representing a tridiagonal matrix from its sub-diagonal, diagonal, and super-diagonal:

```
In [30]: T = Tridiagonal([1,2], [3,4,5], [6,7]) # The type Tridiagonal has three fields
```

```
Out[30]: 3x3 Tridiagonal{Int64, Vector{Int64}}:
```

```
3 6 .
1 4 7
. 2 5
```

Tridiagonal matrices will come up in solving second-order differential equations and orthogonal polynomials. We will later see how linear systems involving tridiagonal matrices can be solved in  $O(n)$  operations.

## II.2 Orthogonal and Unitary Matrices

A very important class of matrices are *orthogonal* and *unitary* matrices:

**Definition 1 (orthogonal/unitary matrix)** A square real matrix is *orthogonal* if its inverse is its transpose:

$$O(n) = \{Q \in \mathbb{R}^{n \times n} : Q^\top Q = I\}$$

A square complex matrix is *unitary* if its inverse is its adjoint:

$$U(n) = \{Q \in \mathbb{C}^{n \times n} : Q^* Q = I\}.$$

Here the adjoint is the same as the conjugate-transpose:  $Q^* := Q^\top$ .

Note that  $O(n) \subset U(n)$  as for real matrices  $Q^* = Q^\top$ . Because in either case  $Q^{-1} = Q^*$  we also have  $QQ^* = I$  (which for real matrices is  $QQ^\top = I$ ). These matrices are particularly important for numerical linear algebra for a number of reasons (we'll explore these properties in the problem sheets):

1. They are norm-preserving: for any vector  $\mathbf{x} \in \mathbb{C}^n$  we have

$\|Q\mathbf{x}\| = \|\mathbf{x}\|$  where  $\|\mathbf{x}\|^2 := \sum_{k=1}^n x_k^2$  (i.e. the 2-norm). 2. All eigenvalues have absolute value equal to 1. 3. For  $Q \in O(n)$ ,  $\det Q = \pm 1$ . 2. They are trivially invertible (just take the transpose). 3. They are generally "stable": errors are controlled. 4. They are *normal matrices*: they commute with their adjoint ( $QQ^* = QQQ^* = Q$ ). See Chapter C for why this is important.

On a computer there are multiple ways of representing orthogonal/unitary matrices, and it is almost never to store a dense matrix storing the entries. We shall therefore investigate three classes:

1. *Permutation*: A permutation matrix permutes the rows of a vector and is a representation of the symmetric group.
2. *Rotations*: The simple rotations are also known as  $2 \times 2$  special orthogonal matrices ( $SO(2)$ ) and correspond to rotations in 2D.
3. *Reflections*: Reflections are  $n \times n$  orthogonal matrices that have simple definitions in terms of a single vector.

We remark a very similar concept are rectangular matrices with orthogonal columns, e.g.

$$U = [\mathbf{u}_1 | \cdots | \mathbf{u}_n] \in \mathbb{R}^{m \times n}$$

where  $m \geq n$  such that  $U^\top U = I_n$  (the  $n \times n$  identity matrix). In this case we must have  $UU^\top \neq I_m$  as the rank of  $U$  is  $n$ . These will play an important role in the Singular Value Decomposition.

# 1. Permutation Matrices

Permutation matrices are matrices that represent the action of permuting the entries of a vector, that is, matrix representations of the symmetric group  $S_n$ , acting on  $\mathbb{R}^n$ . Recall every  $\sigma \in S_n$  is a bijection between  $\{1, 2, \dots, n\}$  and itself. We can write a permutation  $\sigma$  in *Cauchy notation*:

$$\begin{pmatrix} 1 & 2 & 3 & \cdots & n \\ \sigma_1 & \sigma_2 & \sigma_3 & \cdots & \sigma_n \end{pmatrix}$$

where  $\{\sigma_1, \dots, \sigma_n\} = \{1, 2, \dots, n\}$  (that is, each integer appears precisely once). We denote the *inverse permutation* by  $\sigma^{-1}$ , which can be constructed by swapping the rows of the Cauchy notation and reordering.

We can encode a permutation in vector  $\sigma = [\sigma_1, \dots, \sigma_n]$ . This induces an action on a vector (using indexing notation)

$$\mathbf{v}[\sigma] = \begin{bmatrix} v_{\sigma_1} \\ \vdots \\ v_{\sigma_n} \end{bmatrix}$$

**Example 1 (permutation of a vector)** Consider the permutation  $\sigma$  given by

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 1 & 4 & 2 & 5 & 3 \end{pmatrix}$$

We can apply it to a vector:

```
In [1]: σ = [1, 4, 2, 5, 3]
v = [6, 7, 8, 9, 10]
v[σ] # we permute entries of v
```

```
Out[1]: 5-element Vector{Int64}:
```

```
6
9
7
10
8
```

Its inverse permutation  $\sigma^{-1}$  has Cauchy notation coming from swapping the rows of the Cauchy notation of  $\sigma$  and sorting:

$$\begin{pmatrix} 1 & 4 & 2 & 5 & 3 \\ 1 & 2 & 3 & 4 & 5 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 2 & 4 & 3 & 5 \\ 1 & 3 & 2 & 5 & 4 \end{pmatrix}$$

Julia has the function `invperm` for computing the vector that encodes the inverse permutation: And indeed:

```
In [2]: σ⁻¹ = invperm(σ) # note that ⁻¹ are just unicode characters in the variable
```

```
Out[2]: 5-element Vector{Int64}:
```

```
1  
3  
5  
2  
4
```

And indeed permuting the entries by  $\sigma$  and then by  $\sigma^{-1}$  returns us to our original vector:

```
In [3]: v[σ][σ⁻¹] # permuting by σ and then σi gets us back
```

```
Out[3]: 5-element Vector{Int64}:
```

```
6  
7  
8  
9  
10
```

Note that the operator

$$P_\sigma(\mathbf{v}) = \mathbf{v}[\sigma]$$

is linear in  $\mathbf{v}$ , therefore, we can identify it with a matrix whose action is:

$$P_\sigma \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} v_{\sigma_1} \\ \vdots \\ v_{\sigma_n} \end{bmatrix}.$$

The entries of this matrix are

$$P_\sigma[k, j] = \mathbf{e}_k^\top P_\sigma \mathbf{e}_j = \mathbf{e}_k^\top \mathbf{e}_{\sigma_j^{-1}} = \delta_{k, \sigma_j^{-1}} = \delta_{\sigma_k, j}$$

where  $\delta_{k,j}$  is the *Kronecker delta*:

$$\delta_{k,j} := \begin{cases} 1 & k = j \\ 0 & \text{otherwise} \end{cases}.$$

This construction motivates the following definition:

**Definition 2 (permutation matrix)**  $P \in \mathbb{R}^{n \times n}$  is a permutation matrix if it is equal to the identity matrix with its rows permuted.

**Example 2 (5×5 permutation matrix)** We can construct the permutation representation for  $\sigma$  as above as follows:

```
In [4]: P = I(5)[σ, :]
```

```
Out[4]: 5×5 Matrix{Bool}:
```

```
1 0 0 0 0  
0 0 0 1 0  
0 1 0 0 0  
0 0 0 0 1  
0 0 1 0 0
```

And indeed, we see its action is as expected:

```
In [5]: P * v
```

```
Out[5]: 5-element Vector{Int64}:
```

```
6  
9  
7  
10  
8
```

**Remark (advanced)** Note that `P` is a special type `SparseMatrixCSC`. This is used to represent a matrix by storing only the non-zero entries as well as their location. This is an important data type in high-performance scientific computing, but we will not be using general sparse matrices in this module.

**Proposition 1 (permutation matrix inverse)** Let  $P_\sigma$  be a permutation matrix corresponding to the permutation  $\sigma$ . Then

$$P_\sigma^\top = P_{\sigma^{-1}} = P_\sigma^{-1}$$

That is,  $P_\sigma$  is *orthogonal*:

$$P_\sigma^\top P_\sigma = P_\sigma P_\sigma^\top = I.$$

### Proof

We prove orthogonality via:

$$\mathbf{e}_k^\top P_\sigma^\top P_\sigma \mathbf{e}_j = (P_\sigma \mathbf{e}_k)^\top P_\sigma \mathbf{e}_j = \mathbf{e}_{\sigma_k^{-1}}^\top \mathbf{e}_{\sigma_j^{-1}} = \delta_{k,j}$$

This shows  $P_\sigma^\top P_\sigma = I$  and hence  $P_\sigma^{-1} = P_\sigma^\top$ .

■

## 2. Rotations

We begin with a general definition:

**Definition 3 (Special Orthogonal and Rotations)** *Special Orthogonal Matrices* are

$$SO(n) := \{Q \in O(n) \mid \det Q = 1\}$$

And (simple) *rotations* are  $SO(2)$ .

In what follows we use the following for writing the angle of a vector:

**Definition 4 (two-arg arctan)** The two-argument arctan function gives the angle  $\theta$  through the point  $[a, b]^\top$ , i.e.,

$$\sqrt{a^2 + b^2} \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix}$$

It can be defined in terms of the standard arctan as follows:

$$\text{atan}(b, a) := \begin{cases} \text{atan} \frac{b}{a} & a > 0 \\ \text{atan} \frac{b}{a} + \pi & a < 0 \text{ and } b > 0 \\ \text{atan} \frac{b}{a} - \pi & a < 0 \text{ and } b < 0 \\ \pi/2 & a = 0 \text{ and } b > 0 \\ -\pi/2 & a = 0 \text{ and } b < 0 \end{cases}$$

This is available in Julia via the function `atan(y, x)`.

We show  $SO(2)$  are exactly equivalent to standard rotations:

**Proposition 2 (simple rotation)** A  $2 \times 2$  rotation matrix through angle  $\theta$  is

$$Q_\theta := \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

We have  $Q \in SO(2)$  iff  $Q = Q_\theta$  for some  $\theta \in \mathbb{R}$ .

### Proof

We will write  $c = \cos \theta$  and  $s = \sin \theta$ . Then we have

$$Q_\theta^\top Q_\theta = \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} c & -s \\ s & c \end{pmatrix} = \begin{pmatrix} c^2 + s^2 & 0 \\ 0 & c^2 + s^2 \end{pmatrix} = I$$

and  $\det Q_\theta = c^2 + s^2 = 1$  hence  $Q_\theta \in SO(2)$ .

Now suppose  $Q = [\mathbf{q}_1, \mathbf{q}_2] \in SO(2)$  where we know its columns have norm 1  $\|\mathbf{q}_k\| = 1$  and are orthogonal. Write  $\mathbf{q}_1 = [c, s]$  where we know  $c = \cos \theta$  and  $s = \sin \theta$  for  $\theta = \text{atan}(s, c)$ . Since  $\mathbf{q}_1 \cdot \mathbf{q}_2 = 0$  we can deduce  $\mathbf{q}_2 = \pm[-s, c]$ . The sign is positive as  $\det Q = \pm(c^2 + s^2) = \pm 1$ .

■

We can rotate an arbitrary vector in  $\mathbb{R}^2$  to the unit axis using rotations, which are useful in linear algebra decompositions. Interestingly it only requires basic algebraic functions (no trigonometric functions):

**Proposition 3 (rotation of a vector)** The matrix

$$Q = \frac{1}{\sqrt{a^2 + b^2}} \begin{bmatrix} a & b \\ -b & a \end{bmatrix}$$

is a rotation matrix ( $Q \in SO(2)$ ) satisfying

$$Q \begin{bmatrix} a \\ b \end{bmatrix} = \sqrt{a^2 + b^2} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

### Proof

The last equation is trivial so the only question is that it is a rotation matrix. This follows immediately:

$$Q^\top Q = \frac{1}{a^2 + b^2} \begin{bmatrix} a^2 + b^2 & 0 \\ 0 & a^2 + b^2 \end{bmatrix} = I$$

and  $\det Q = 1$ .

■

## 3. Reflections

In addition to rotations, another type of orthogonal/unitary matrix are reflections:

**Definition 5 (reflection matrix)** Given a unit vector  $\mathbf{v} \in \mathbb{C}^n$  (satisfying  $\|\mathbf{v}\| = 1$ ), the *reflection matrix*

$$Q_{\mathbf{v}} := I - 2\mathbf{v}\mathbf{v}^*$$

These are reflections in the direction of  $\mathbf{v}$ . We can show this as follows:

**Proposition 4 (Householder properties)**  $Q_{\mathbf{v}}$  satisfies:

1.  $Q_{\mathbf{v}} = Q_{\mathbf{v}}^*$  (Symmetry)
2.  $Q_{\mathbf{v}}^* Q_{\mathbf{v}} = I$  (Orthogonality  $Q_{\mathbf{v}} \in U(n)$ )
3.  $\mathbf{v}$  is an eigenvector of  $Q_{\mathbf{v}}$  with eigenvalue  $-1$
4.  $Q_{\mathbf{v}}$  is a rank-1 perturbation of  $I$
5.  $\det Q_{\mathbf{v}} = -1$  ( $Q_{\mathbf{v}} \notin SO(n)$ )

### Proof

Property 1 follows immediately. Property 2 follows from

$$Q_{\mathbf{v}}^* Q_{\mathbf{v}} = Q_{\mathbf{v}}^2 = I - 4\mathbf{v}\mathbf{v}^* + 4\mathbf{v}\mathbf{v}^*\mathbf{v}\mathbf{v}^* = I$$

Property 3 follows since

$$Q_{\mathbf{v}}\mathbf{v} = -\mathbf{v}$$

Property 4 follows since  $\mathbf{v}\mathbf{v}^\top$  is a rank-1 matrix as all rows are linear combinations of each other. To see property 5, note there is a dimension  $n - 1$  space  $W$  orthogonal to  $\mathbf{v}$ , that is, for all  $\mathbf{w} \in W$  we have  $\mathbf{w}^\star\mathbf{v} = 0$ , which implies that

$$Q_{\mathbf{v}}\mathbf{w} = \mathbf{w}$$

In other words, 1 is an eigenvalue with multiplicity  $n - 1$  and  $-1$  is an eigenvalue with multiplicity 1, and thus the product of the eigenvalues is  $-1$ .

■

**Example 3 (reflection through 2-vector)** Consider reflection through  $\mathbf{x} = [1, 2]^\top$ . We first need to normalise  $\mathbf{x}$ :

$$\mathbf{v} = \frac{\mathbf{x}}{\|\mathbf{x}\|} = \begin{bmatrix} \frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{bmatrix}$$

Note this indeed has unit norm:

$$\|\mathbf{v}\|^2 = \frac{1}{5} + \frac{4}{5} = 1.$$

Thus the reflection matrix is:

$$Q_{\mathbf{v}} = I - 2\mathbf{v}\mathbf{v}^\top = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \frac{2}{5} \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ -4 & -3 \end{bmatrix}$$

Indeed it is symmetric, and orthogonal. It sends  $\mathbf{x}$  to  $-\mathbf{x}$ :

$$Q_{\mathbf{v}}\mathbf{x} = \frac{1}{5} \begin{bmatrix} 3 - 8 \\ -4 - 6 \end{bmatrix} = -\mathbf{x}$$

Any vector orthogonal to  $\mathbf{x}$ , like  $\mathbf{y} = [-2, 1]^\top$ , is left fixed:

$$Q_{\mathbf{v}}\mathbf{y} = \frac{1}{5} \begin{bmatrix} -6 - 4 \\ 8 - 3 \end{bmatrix} = \mathbf{y}$$

Note that *building* the matrix  $Q_{\mathbf{v}}$  will be expensive ( $O(n^2)$  operations), but we can *apply*  $Q_{\mathbf{v}}$  to a vector in  $O(n)$  operations using the expression:

$$Q_{\mathbf{v}}\mathbf{x} = \mathbf{x} - 2\mathbf{v}(\mathbf{v}^\star\mathbf{x}) = \mathbf{x} - 2\mathbf{v}(\mathbf{v} \cdot \mathbf{x}).$$

Just as rotations can be used to rotate vectors to be aligned with coordinate axis, so can reflections, but in this case it works for vectors in  $\mathbb{C}^n$ , not just  $\mathbb{R}^2$ :

**Definition 6 (Householder reflection, real case)** For a given vector  $\mathbf{x} \in \mathbb{R}^n$ , define the Householder reflection

$$Q_{\mathbf{x}}^{\pm, H} := Q_{\mathbf{w}}$$

for  $\mathbf{y} = \mp \|\mathbf{x}\| \mathbf{e}_1 + \mathbf{x}$  and  $\mathbf{w} = \frac{\mathbf{y}}{\|\mathbf{y}\|}$ . The default choice in sign is:

$$Q_{\mathbf{x}}^H := Q_{\mathbf{x}}^{-\text{sign}(x_1), H}.$$

**Lemma 1 (Householder reflection maps to axis)** For  $\mathbf{x} \in \mathbb{R}^n$ ,

$$Q_{\mathbf{x}}^{\pm, H} \mathbf{x} = \pm \|\mathbf{x}\| \mathbf{e}_1$$

**Proof** Note that

$$\begin{aligned}\|\mathbf{y}\|^2 &= 2\|\mathbf{x}\|^2 \mp 2\|\mathbf{x}\|x_1, \\ \mathbf{y}^\top \mathbf{x} &= \|\mathbf{x}\|^2 \mp \|\mathbf{x}\|x_1\end{aligned}$$

where  $x_1 = \mathbf{e}_1^\top \mathbf{x}$ . Therefore:

$$Q_{\mathbf{x}}^{\pm, H} \mathbf{x} = (I - 2\mathbf{w}\mathbf{w}^\top) \mathbf{x} = \mathbf{x} - 2 \frac{\mathbf{y}\|\mathbf{x}\|}{\|\mathbf{y}\|^2} (\|\mathbf{x}\| \mp x_1) = \mathbf{x} - \mathbf{y} = \pm \|\mathbf{x}\| \mathbf{e}_1.$$

■

Why do we choose the opposite sign of  $x_1$  for the default reflection? For stability. We demonstrate the reason for this by numerical example. Consider  $\mathbf{x} = [1, h]$ , i.e., a small perturbation from  $\mathbf{e}_1$ . If we reflect to  $\text{norm}(\mathbf{x})\mathbf{e}_1$  we see a numerical problem:

```
In [6]: h = 10.0^(-10)
x = [1, h]
y = -norm(x)*[1, 0] + x
w = y/norm(y)
Q = I - 2w*w'
Q*x
```

```
Out[6]: 2-element Vector{Float64}:
 1.0
 -1.0e-10
```

It didn't work! Even worse is if  $h = 0$ :

```
In [7]: h = 0
x = [1, h]
y = -norm(x)*[1, 0] + x
w = y/norm(y)
Q = I - 2w*w'
Q*x
```

```
Out[7]: 2-element Vector{Float64}:
 NaN
 NaN
```

This is because  $y$  has large relative error due to cancellation from floating point errors in computing the first entry  $x[1] - \text{norm}(x)$ . (Or has norm zero if  $h=0$ .) We avoid this cancellation by using the default choice:

```
In [8]: h = 10.0^(-10)
x = [1,h]
y = sign(x[1])*norm(x)*[1,0] + x
w = y/norm(y)
Q = I - 2w*w'
Q*x
```

Out[8]: 2-element Vector{Float64}:

```
-1.0
0.0
```

We can extend this definition for complexes:

**Definition 7 (Householder reflection, complex case)** For a given vector  $\mathbf{x} \in \mathbb{C}^n$ , define the Householder reflection as

$$Q_{\mathbf{x}}^H := Q_{\mathbf{w}}$$

for  $\mathbf{y} = \text{csign}(x_1)\|\mathbf{x}\|\mathbf{e}_1 + \mathbf{x}$  and  $\mathbf{w} = \frac{\mathbf{y}}{\|\mathbf{y}\|}$ , for  $\text{csign}(z) = e^{i \arg z}$ .

**Lemma 2 (Householder reflection maps to axis, complex case)** For  $\mathbf{x} \in \mathbb{C}^n$ ,

$$Q_{\mathbf{x}}^H \mathbf{x} = -\text{csign}(x_1)\|\mathbf{x}\|\mathbf{e}_1$$

**Proof** Denote  $\alpha := \text{csign}(x_1)$ . Note that  $\bar{\alpha}x_1 = e^{-i \arg x_1}x_1 = |x_1|$ . Now we have

$$\begin{aligned} \|\mathbf{y}\|^2 &= (\alpha\|\mathbf{x}\|\mathbf{e}_1 + \mathbf{x})^*(\alpha\|\mathbf{x}\|\mathbf{e}_1 + \mathbf{x}) = |\alpha|\|\mathbf{x}\|^2 + \|\mathbf{x}\|\alpha\bar{x}_1 + \bar{\alpha}x_1\|\mathbf{x}\| + \|\mathbf{x}\|^2 \\ &= 2\|\mathbf{x}\|^2 + 2|x_1|\|\mathbf{x}\| \\ \mathbf{y}^*\mathbf{x} &= \bar{\alpha}x_1\|\mathbf{x}\| + \|\mathbf{x}\|^2 = \|\mathbf{x}\|^2 + |x_1|\|\mathbf{x}\| \end{aligned}$$

Therefore:

$$Q_{\mathbf{x}}^H \mathbf{x} = (I - 2\mathbf{w}\mathbf{w}^*)\mathbf{x} = \mathbf{x} - 2\frac{\mathbf{y}}{\|\mathbf{y}\|^2}(\|\mathbf{x}\|^2 + |x_1|\|\mathbf{x}\|) = \mathbf{x} - \mathbf{y} = -\alpha\|\mathbf{x}\|\mathbf{e}_1.$$

■

## II.3 QR factorisation

Let  $A \in \mathbb{C}^{m \times n}$  be a rectangular or square matrix such that  $m \geq n$  (i.e. more rows than columns). In this chapter we consider two closely related factorisations:

## 1. The *QR* factorisation

$$A = QR = \underbrace{[ \mathbf{q}_1 | \cdots | \mathbf{q}_m ]}_{Q \in U(m)} \left[ \begin{array}{ccc} \times & \cdots & \times \\ & \ddots & \vdots \\ & & \times \\ & & 0 \\ & & \vdots \\ & & 0 \end{array} \right] R \in \mathbb{C}^{m \times n}$$

where  $Q$  is unitary (i.e.,  $Q \in U(m)$ , satisfying  $Q^*Q = I$ , with columns  $\mathbf{q}_j \in \mathbb{C}^m$ ) and  $R$  is *right triangular*, which means it is only nonzero on or to the right of the diagonal ( $r_{kj} = 0$  if  $k > j$ ).

## 2. The reduced QR factorisation

$$A = \hat{Q}\hat{R} = \underbrace{[\mathbf{q}_1 | \cdots | \mathbf{q}_n]}_{\hat{Q} \in \mathbb{C}^{m \times n}} \underbrace{\begin{bmatrix} \times & \cdots & \times \\ & \ddots & \vdots \\ & & \times \end{bmatrix}}_{\hat{R} \in \mathbb{C}^{n \times n}}$$

where  $Q$  has orthonormal columns ( $Q^*Q = I$ ,  $\mathbf{q}_j \in \mathbb{C}^m$ ) and  $\hat{R}$  is upper triangular.

Note for a square matrix the reduced QR factorisation is equivalent to the QR factorisation, in which case  $R$  is *upper triangular*. The importance of these decompositions for square matrices is that their component pieces are easy to invert:

$$A = QR \quad \Rightarrow \quad A^{-1}\mathbf{b} = R^{-1}Q^\top\mathbf{b}$$

and we saw in the last two chapters that triangular and orthogonal matrices are easy to invert when applied to a vector  $\mathbf{b}$ , e.g., using forward/back-substitution.

For rectangular matrices we will see that they lead to efficient solutions to the *least squares problem*: find  $\mathbf{x}$  that minimizes the 2-norm

$$\|A\mathbf{x} - \mathbf{b}\|.$$

Note in the rectangular case the QR decomposition contains within it the reduced QR decomposition:

$$A = QR = [\hat{Q}|\mathbf{q}_{n+1}| \cdots |\mathbf{q}_m] \begin{bmatrix} \hat{R} \\ \mathbf{0}_{m-n \times n} \end{bmatrix} = \hat{Q}\hat{R}.$$

In this lecture we discuss the following:

1. QR and least squares: We discuss the QR decomposition and its usage in solving least squares problems.
2. Reduced QR and Gram–Schmidt: We discuss computation of the Reduced QR decomposition using Gram–Schmidt.
3. Householder reflections and QR: We discuss computing the QR decomposition using Householder reflections.

```
In [1]: using LinearAlgebra, Plots, BenchmarkTools
```

## 1. QR and least squares

Here we consider rectangular matrices with more rows than columns. Given  $A \in \mathbb{C}^{m \times n}$  and  $\mathbf{b} \in \mathbb{C}^m$ , least squares consists of finding a vector  $\mathbf{x} \in \mathbb{C}^n$  that minimises the 2-norm:  $\|A\mathbf{x} - \mathbf{b}\|$ .

**Theorem 1 (least squares via QR)** Suppose  $A \in \mathbb{C}^{m \times n}$  has full rank. Given a QR decomposition  $A = QR$  then

$$\mathbf{x} = \hat{R}^{-1} \hat{Q}^* \mathbf{b}$$

minimises  $\|A\mathbf{x} - \mathbf{b}\|$ .

### Proof

The norm-preserving property (see PS4 Q3.1) of unitary matrices tells us

$$\|A\mathbf{x} - \mathbf{b}\| = \|QR\mathbf{x} - \mathbf{b}\| = \|Q(R\mathbf{x} - Q^*\mathbf{b})\| = \|R\mathbf{x} - Q^*\mathbf{b}\| = \left\| \begin{bmatrix} \hat{R} \\ \mathbf{0}_{m-n \times n} \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{q}_{n+1} \\ \vdots \\ \mathbf{q}_m \\ \vdots \end{bmatrix} \right\|$$

Now note that the rows  $k > n$  are independent of  $\mathbf{x}$  and are a fixed contribution. Thus to minimise this norm it suffices to drop them and minimise:

$$\|\hat{R}\mathbf{x} - \hat{Q}^*\mathbf{b}\|$$

This norm is minimised if it is zero. Provided the column rank of  $A$  is full,  $\hat{R}$  will be invertible (Exercise: why is this?).

■

**Example 1 (quadratic fit)** Suppose we want to fit noisy data by a quadratic

$$p(x) = p_0 + p_1x + p_2x^2$$

That is, we want to choose  $p_0, p_1, p_2$  at data samples  $x_1, \dots, x_m$  so that the following is true:

$$p_0 + p_1x_k + p_2x_k^2 \approx f_k$$

where  $f_k$  are given by data. We can reinterpret this as a least squares problem: minimise the norm

$$\left\| \begin{bmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_m & x_m^2 \end{bmatrix} \begin{bmatrix} p_0 \\ p_1 \\ p_2 \end{bmatrix} - \begin{bmatrix} f_1 \\ \vdots \\ f_m \end{bmatrix} \right\|$$

We can solve this using the QR decomposition:

```
In [2]: m, n = 100, 3

x = range(0,1; length=m) # 100 points
f = 2 .+ x .+ 2x.^2 .+ 0.1 .* randn() # Noisy quadratic

A = x .^ (0:2)' # 100 x 3 matrix, equivalent to [ones(m) x x.^2]
Q, R = qr(A)
Q̂ = Q[:,1:n] # Q represents full orthogonal matrix so we take first 3 columns

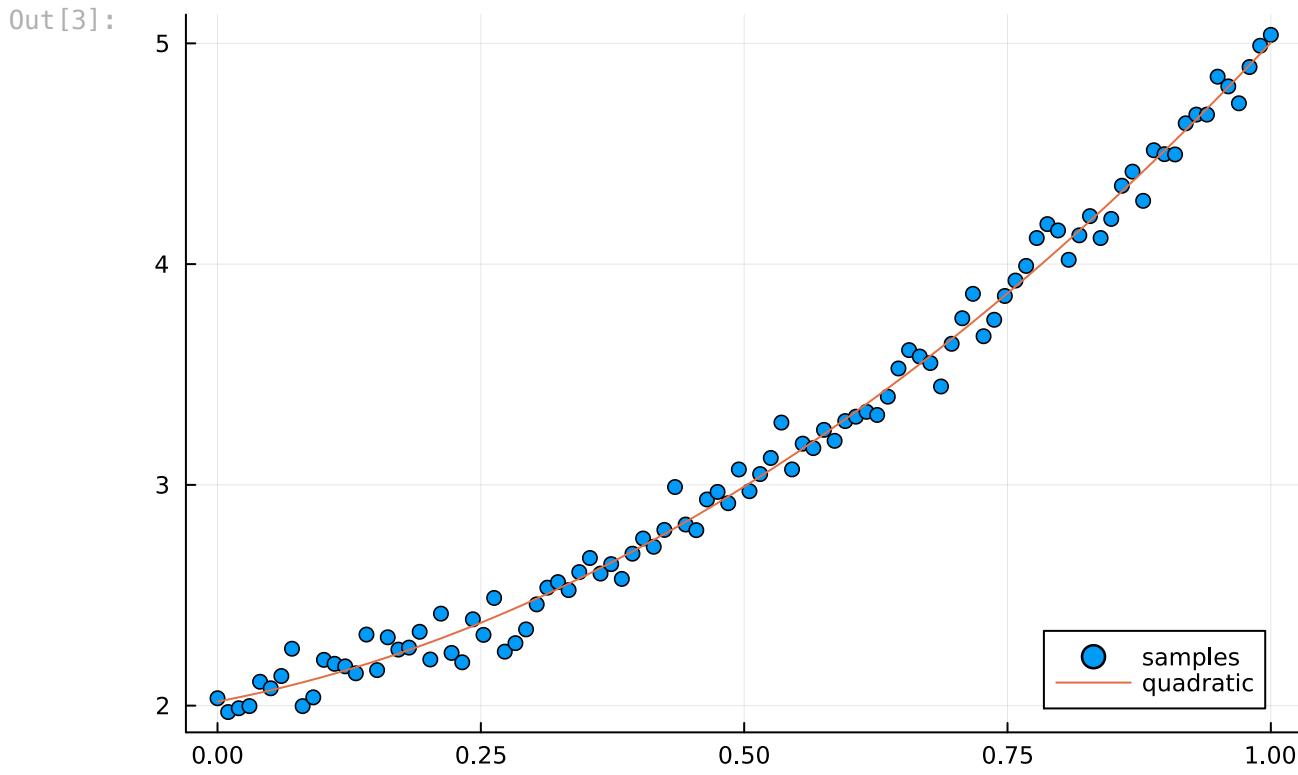
p₀, p₁, p₂ = R \ Q̂'f
```

```
Out[2]: 3-element Vector{Float64}:
2.0184037164216138
0.9091960695680688
2.075701832630235
```

We can visualise the fit:

```
In [3]: p = x -> p₀ + p₁*x + p₂*x.^2

scatter(x, f; label="samples", legend=:bottomright)
plot!(x, p.(x); label="quadratic")
```



Note that `\` with a rectangular system does least squares by default:

In [4]: `A \ f`

Out [4]: 3-element Vector{Float64}:

2.018403716421614
0.9091960695680675
2.075701832630236

## 2. Reduced QR and Gram–Schmidt

How do we compute the QR decomposition? We begin with a method you may have seen before in another guise. Write

$$A = [\mathbf{a}_1 | \cdots | \mathbf{a}_n]$$

where  $\mathbf{a}_k \in \mathbb{C}^m$  and assume they are linearly independent ( $A$  has full column rank).

**Proposition 1 (Column spaces match)** Suppose  $A = \hat{Q}\hat{R}$  where  $\hat{Q} = [\mathbf{q}_1 | \dots | \mathbf{q}_n]$  has orthonormal columns and  $\hat{R}$  is upper-triangular, and  $A$  has full rank. Then the first  $j$  columns of  $\hat{Q}$  span the same space as the first  $j$  columns of  $A$ :

$$\text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j) = \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_j).$$

**Proof**

Because  $A$  has full rank we know  $\hat{R}$  is invertible, i.e. its diagonal entries do not vanish:  $r_{jj} \neq 0$ . If  $\mathbf{v} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j)$  we have for  $\mathbf{c} \in \mathbb{C}^j$

$$\mathbf{v} = [\mathbf{a}_1 | \cdots | \mathbf{a}_j] \mathbf{c} = [\mathbf{q}_1 | \cdots | \mathbf{q}_j] \hat{R}[1:j, 1:j] \mathbf{c} \in \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_j)$$

while if  $\mathbf{w} \in \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_j)$  we have for  $\mathbf{d} \in \mathbb{R}^j$

$$\mathbf{w} = [\mathbf{q}_1 | \cdots | \mathbf{q}_j] \mathbf{d} = [\mathbf{a}_1 | \cdots | \mathbf{a}_j] \hat{R}[1:j, 1:j]^{-1} \mathbf{d} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_j).$$

■

It is possible to find  $\hat{Q}$  and  $\hat{R}$  the using the *Gram–Schmidt algorithm*. We construct it column-by-column:

**Algorithm 1 (Gram–Schmidt)** For  $j = 1, 2, \dots, n$  define

$$\begin{aligned}\mathbf{v}_j &:= \mathbf{a}_j - \sum_{k=1}^{j-1} \underbrace{\mathbf{q}_k^* \mathbf{a}_j}_{r_{kj}} \mathbf{q}_k \\ r_{jj} &:= \|\mathbf{v}_j\| \\ \mathbf{q}_j &:= \frac{\mathbf{v}_j}{r_{jj}}\end{aligned}$$

**Theorem 2 (Gram–Schmidt and reduced QR)** Define  $\mathbf{q}_j$  and  $r_{kj}$  as in Algorithm 1 (with  $r_{kj} = 0$  if  $k > j$ ). Then a reduced QR decomposition is given by:

$$A = \underbrace{[\mathbf{q}_1 | \cdots | \mathbf{q}_n]}_{\hat{Q} \in \mathbb{C}^{m \times n}} \underbrace{\begin{bmatrix} r_{11} & \cdots & r_{1n} \\ & \ddots & \vdots \\ & & r_{nn} \end{bmatrix}}_{\hat{R} \in \mathbb{C}^{n \times n}}$$

### Proof

We first show that  $\hat{Q}$  has orthonormal columns. Assume that  $\mathbf{q}_\ell^* \mathbf{q}_k = \delta_{\ell k}$  for  $k, \ell < j$ . For  $\ell < j$  we then have

$$\mathbf{q}_\ell^* \mathbf{v}_j = \mathbf{q}_\ell^* \mathbf{a}_j - \sum_{k=1}^{j-1} \mathbf{q}_\ell^* \mathbf{q}_k \mathbf{q}_k^* \mathbf{a}_j = 0$$

hence  $\mathbf{q}_\ell^* \mathbf{q}_j = 0$  and indeed  $\hat{Q}$  has orthonormal columns. Further: from the definition of  $\mathbf{v}_j$  we find

$$\mathbf{a}_j = \mathbf{v}_j + \sum_{k=1}^{j-1} r_{kj} \mathbf{q}_k = \sum_{k=1}^j r_{kj} \mathbf{q}_k = \hat{Q} \hat{R} \mathbf{e}_j$$

■

### Gram–Schmidt in action

We are going to compute the reduced QR of a random matrix

```
In [5]: m, n = 5,4
A = randn(m,n)
Q, R = qr(A)
Q_hat = Q[:,1:n]
```

```
Out[5]: 5x4 Matrix{Float64}:
-0.389812 -0.584855 -0.538419  0.0057763
 0.745394 -0.285764 -0.278455 -0.531648
-0.148359  0.566612 -0.753843 -0.0908584
 0.022578  0.492735  0.117206 -0.368433
-0.519532 -0.111563  0.224836 -0.757178
```

The first column of  $\hat{Q}$  is indeed a normalised first column of  $A$ :

```
In [6]: R = zeros(n,n)
Q = zeros(m,n)
R[1,1] = norm(A[:,1])
Q[:,1] = A[:,1]/R[1,1]
```

```
Out[6]: 5-element Vector{Float64}:
-0.3898123894153746
 0.745394206206071
-0.14835938605298077
 0.02257801294843124
-0.519532005091102
```

We now determine the next entries as

```
In [7]: R[1,2] = Q[:,1]'A[:,2]
v = A[:,2] - Q[:,1]*R[1,2]
R[2,2] = norm(v)
Q[:,2] = v/R[2,2]
```

```
Out[7]: 5-element Vector{Float64}:
 0.5848549497934579
 0.2857643934456054
-0.5666124529207136
-0.4927345617481724
 0.11156334072967596
```

And the third column is then:

```
In [8]: R[1,3] = Q[:,1]'A[:,3]
R[2,3] = Q[:,2]'A[:,3]
v = A[:,3] - Q[:,1:2]*R[1:2,3]
R[3,3] = norm(v)
Q[:,3] = v/R[3,3]
```

```
Out[8]: 5-element Vector{Float64}:
-0.5384185231775488
-0.27845509681334496
-0.7538434527793718
0.11720557043152843
0.2248358421592807
```

(Note the signs may not necessarily match.)

We can clean this up as a simple algorithm:

```
In [9]: function gramschmidt(A)
    m,n = size(A)
    m ≥ n || error("Not supported")
    R = zeros(n,n)
    Q = zeros(m,n)
    for j = 1:n
        for k = 1:j-1
            R[k,j] = Q[:,k]'*A[:,j]
        end
        v = A[:,j] - Q[:,1:j-1]*R[1:j-1,j]
        R[j,j] = norm(v)
        Q[:,j] = v/R[j,j]
    end
    Q,R
end

Q,R = gramschmidt(A)
norm(A - Q*R)
```

```
Out[9]: 5.581498056676073e-16
```

## Complexity and stability

We see within the `for j = 1:n` loop that we have  $O(mj)$  operations. Thus the total complexity is  $O(mn^2)$  operations.

Unfortunately, the Gram–Schmidt algorithm is *unstable*: the rounding errors when implemented in floating point accumulate in a way that we lose orthogonality:

```
In [10]: A = randn(300,300)
Q,R = gramschmidt(A)
norm(Q'Q-I)
```

```
Out[10]: 3.1871092146215285e-12
```

## 3. Householder reflections and QR

As an alternative, we will consider using Householder reflections to introduce zeros below the diagonal. Thus, if Gram–Schmidt is a process of *triangular orthogonalisation*

(using triangular matrices to orthogonalise), Householder reflections is a process of *orthogonal triangularisation* (using orthogonal matrices to triangularise).

Consider multiplication by the Householder reflection corresponding to the first column, that is, for

$$Q_1 := Q_{\mathbf{a}_1}^H,$$

consider

$$Q_1 A = \begin{bmatrix} \times & \times & \cdots & \times \\ & \times & \cdots & \times \\ & & \ddots & \\ \vdots & & & \vdots \\ \times & \cdots & \times \end{bmatrix} = \begin{bmatrix} \alpha & \mathbf{w}^\top \\ & A_2 \end{bmatrix}$$

where

$$\alpha := -\text{csign}(a_{11})\|\mathbf{a}_1\|, \mathbf{w} = (Q_1 A)[1, 2 : n] \quad \text{and} \quad A_2 = (Q_1 A)[2 : m, 2 : n],$$

$\text{csign}(z) := e^{i \arg z}$ . That is, we have made the first column triangular. In terms of an algorithm, we then introduce zeros into the first column of  $A_2$ , leaving an  $A_3$ , and so-on. But we can wrap this iterative algorithm into a simple proof by induction:

**Theorem 3 (QR)** Every matrix  $A \in \mathbb{C}^{m \times n}$  has a QR factorisation:

$$A = QR$$

where  $Q \in U(m)$  and  $R \in \mathbb{C}^{m \times n}$  is right triangular.

### Proof

Assume  $m \geq n$ . If  $A = [\mathbf{a}_1] \in \mathbb{C}^{m \times 1}$  then we have for the Householder reflection  $Q_1 = Q_{\mathbf{a}_1}^H$

$$Q_1 A = [\alpha \mathbf{e}_1]$$

which is right triangular, where  $\alpha = -\text{sign}(a_{11})\|\mathbf{a}_1\|$ . In other words

$$A = \underbrace{Q_1}_{Q} \underbrace{[\alpha \mathbf{e}_1]}_{R}.$$

For  $n > 1$ , assume every matrix with less columns than  $n$  has a QR factorisation. For  $A = [\mathbf{a}_1 | \dots | \mathbf{a}_n] \in \mathbb{C}^{m \times n}$ , let  $Q_1 = Q_{\mathbf{a}_1}^H$  so that

$$Q_1 A = \begin{bmatrix} \alpha & \mathbf{w}^\top \\ & A_2 \end{bmatrix}$$

where  $A_2 = (Q_1 A)[2 : m, 2 : n]$  and  $\mathbf{w} = (Q_1 A)[1, 2 : n]$ . By assumption  $A_2 = \tilde{Q} \tilde{R}$ . Thus we have

$$A = Q_1 \begin{bmatrix} \alpha & \mathbf{w}^\top \\ \tilde{Q}\tilde{R} \end{bmatrix}$$

$$= \underbrace{Q_1 \begin{bmatrix} 1 & \\ & \tilde{Q} \end{bmatrix}}_Q \underbrace{\begin{bmatrix} \alpha & \mathbf{w}^\top \\ & \tilde{R} \end{bmatrix}}_R.$$

This proof by induction leads naturally to an iterative algorithm. Note that  $\tilde{Q}$  is a product of all Householder reflections that come afterwards, that is, we can think of  $Q$  as:

$$Q = Q_1 \tilde{Q}_2 \tilde{Q}_3 \cdots \tilde{Q}_n \quad \text{for} \quad \tilde{Q}_j = \begin{bmatrix} I_{j-1} & \\ & Q_j \end{bmatrix}$$

where  $Q_j$  is a single Householder reflection corresponding to the first column of  $A_j$ .

This is stated cleanly in Julia code:

**Algorithm 2 (QR via Householder)** For  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$ , the QR factorisation can be implemented as follows:

```
In [11]: function householderreflection(x)
    y = copy(x)
    if x[1] == 0
        y[1] += norm(x)
    else # note sign(z) = exp(im*angle(z)) where `angle` is the argument of
        y[1] += sign(x[1])*norm(x)
    end
    w = y/norm(y)
    I = 2*w*w'
end
function householderqr(A)
    T = eltype(A)
    m,n = size(A)
    if n > m
        error("More columns than rows is not supported")
    end

    R = zeros(T, m, n)
    Q = Matrix(one(T)*I, m, m)
    A_j = copy(A)

    for j = 1:n
        a1 = A_j[:,1] # first columns of A_j
        Q1 = householderreflection(a1)
        Q1A_j = Q1*A_j
        α, w = Q1A_j[1,1], Q1A_j[1,2:end]
        A_j+1 = Q1A_j[2:end,2:end]

        # populate returned data
        R[j,j] = α
        R[j,j+1:end] = w
    end
end
```

```

# following is equivalent to Q = Q*[I 0 ; 0 Q_1]
Q[:,j:end] = Q[:,j:end]*Q1

A_j = A_{j+1} # this is the "induction"
end
Q,R
end

m,n = 100,50
A = randn(m,n)
Q,R = householderqr(A)
@test Q'Q ≈ I
@test Q*R ≈ A

```

Out[11]: **Test Passed**

Note because we are forming a full matrix representation of each Householder reflection this is a slow algorithm, taking  $O(n^4)$  operations. The problem sheet will consider a better implementation that takes  $O(n^3)$  operations.

**Example 2 (non-examinable)** We will now do an example by hand. Consider the  $4 \times 3$  matrix

$$A = \begin{bmatrix} 2 & 3 & 0 \\ 0 & 0 & 1 \\ -2 & -3 & 0 \\ -1 & -3 & -3 \end{bmatrix}$$

For the first column we have

$$\mathbf{y}_1 := [-1, 0, -2, -1]$$

where  $\|\mathbf{y}_1\|^2 = 6$ . Hence

$$Q_1 := I - \frac{1}{3} \begin{bmatrix} -1 \\ 0 \\ -2 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & 0 & -2 & -1 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 2 & 0 & -2 & -1 \\ 0 & 3 & 0 & 0 \\ -2 & 0 & -1 & -2 \\ -1 & 0 & -2 & 2 \end{bmatrix}$$

so that

$$Q_1 A = \begin{bmatrix} 3 & 5 & 1 \\ 0 & 1 \\ 1 & 2 \\ -1 & -2 \end{bmatrix}$$

For the second column we have

$$\mathbf{y}_2 := [-\sqrt{2}, 1, -1]$$

where  $\|\mathbf{y}_2\|^2 = 4$ . Thus we have

$$Q_2 := I - \frac{1}{2} \begin{bmatrix} -\sqrt{2} \\ 1 \\ -1 \end{bmatrix} \begin{bmatrix} -\sqrt{2} & 1 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/2 & 1/2 \\ -1/\sqrt{2} & 1/2 & 1/2 \end{bmatrix}$$

so that

$$\tilde{Q}_2 Q_1 A = \begin{bmatrix} 3 & 5 & 1 \\ \sqrt{2} & 2\sqrt{2} & 0 \\ 0 & 1/\sqrt{2} & 0 \\ 0 & -1/\sqrt{2} & 0 \end{bmatrix}$$

The final vector is

$$\mathbf{y}_3 := [1/\sqrt{2} - 1, -1/\sqrt{2}]$$

where  $\|\mathbf{y}_3\|^2 = 2 - 2/\sqrt{2}$ . Hence

$$Q_3 := I - \frac{\sqrt{2}}{\sqrt{2}-1} \begin{bmatrix} 1/\sqrt{2}-1 \\ -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2}-1 & -1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} \sqrt{2} & -\sqrt{2} \\ -\sqrt{2} & -\sqrt{2} \end{bmatrix}$$

so that

$$\tilde{Q}_3 \tilde{Q}_2 Q_1 A = \begin{bmatrix} 3 & 5 & 1 \\ \sqrt{2} & 2\sqrt{2} & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} =: R$$

and

$$Q := Q_1 \tilde{Q}_2 \tilde{Q}_3 = \begin{bmatrix} 2/3 & -1/(3\sqrt{2}) & 0 & 1/\sqrt{2} \\ 0 & 0 & 1 & 0 \\ -2/3 & 1/(3\sqrt{2}) & 0 & 1/\sqrt{2} \\ -1/3 & -4/(3\sqrt{2}) & 0 & 0 \end{bmatrix}.$$

## II.4 PLU and Cholesky factorisations

In this chapter we consider the following factorisations for square invertible matrices  $A$ :

1. The *LU factorisation*:

$A = LU$  where  $L$  is lower triangular and  $U$  is upper triangular. This is equivalent to Gaussian elimination without pivoting, so may not exist (e.g. if  $A[1, 1] = 0$ ).

1. The *PLU factorisation*:

$A = P^T LU$  where  $P$  is a permutation matrix,  $L$  is lower triangular and  $U$  is upper triangular. This is equivalent to Gaussian elimination with pivoting. It always exists but may in extremely rare cases be unstable. 2. For a real square *symmetric positive definite* ( $A \in \mathbb{R}^{n \times n}$  such that  $A^T = A$  and  $\mathbf{x}^T A \mathbf{x} > 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \neq 0$ ) matrix the LU decompostion has a special form which is called the *Cholesky factorisation*:  $A = LL^T$ . This provides an algorithmic way to *prove* that a matrix is symmetric positive definite. 3. We also discuss timing and stability of the different factorisations.

```
In [1]: using LinearAlgebra, Plots, BenchmarkTools
```

```
[ Info: Precompiling Plots [91a5bcdd-55d7-5caf-9e0b-520d859cae80]
[ Info: GR
```

### 1. LU Factorisation

Just as Gram–Schmidt can be reinterpreted as a reduced QR factorisation, Gaussian elimination can be interpreted as an LU factorisation. Write a matrix  $A \in \mathbb{C}^{n \times n}$  as follows:

$$A = \begin{bmatrix} \alpha_1 & \mathbf{w}_1^\top \\ \mathbf{v}_1 & A_2 \end{bmatrix}$$

where  $\alpha_1 = a_{11}$ ,  $\mathbf{v}_1 = A[2 : n, 1]$  and  $\mathbf{w}_1 = A[1, 2 : n]$ . Gaussian elimination consists of taking the first row, dividing by  $\alpha$  and subtracting from all other rows. That is equivalent to multiplying by a lower triangular matrix:

$$\begin{bmatrix} 1 & \\ -\mathbf{v}_1/\alpha_1 & I \end{bmatrix} A = \begin{bmatrix} \alpha_1 & \mathbf{w}_1^\top \\ & K - \mathbf{v}_1 \mathbf{w}_1^\top / \alpha_1 \end{bmatrix}$$

where  $A_2 := K - \mathbf{v}_1 \mathbf{w}_1^\top / \alpha_1$  happens to be a rank-1 perturbation of  $K$ . We can write this another way:

$$A = \underbrace{\begin{bmatrix} 1 \\ \mathbf{v}_1/\alpha_1 & I \end{bmatrix}}_{L_1} \begin{bmatrix} \alpha_1 & \mathbf{w}_1^\top \\ & A_2 \end{bmatrix}$$

Now assume we continue this process and manage to deduce  $A_2 = L_2 U_2$ . Then

$$A = L_1 \begin{bmatrix} \alpha_1 & \mathbf{w}_1^\top \\ & L_2 U_2 \end{bmatrix} = \underbrace{L_1 \begin{bmatrix} 1 & \\ & L_2 \end{bmatrix}}_L \underbrace{\begin{bmatrix} \alpha_1 & \mathbf{w}_1^\top \\ & U_2 \end{bmatrix}}_U$$

Note we can multiply through to find

$$L = \begin{bmatrix} 1 \\ \mathbf{v}_1/\alpha_1 & L_2 \end{bmatrix}$$

This procedure implies an algorithm:

### Algorithm 1 (LU)

```
In [2]: function mylu(A)
    n,m = size(A)
    if n ≠ m
        error("Matrix must be square")
    end
    T = eltype(A)
    L = LowerTriangular(zeros(T,n,n))
    U = UpperTriangular(zeros(T,n,n))

    σ = Vector(1:n)

    A_j = copy(A)

    for j = 1:n-1
        α,v,w = A_j[1,1],A_j[2:end,1],A_j[1,2:end]
        K = A_j[2:end,2:end]

        # populate data
        L[j,j] = 1
        L[j+1:end,j] = v/α
        U[j,j] = α
        U[j,j+1:end] = w

        # this is the "recursion": A_j is now the next block
        # We use transpose(w) instead of w' incase w is complex

        A_j = K - v*transpose(w)/α
    end
    # j = n case
    L[n,n] = 1
    U[n,n] = A_j[1,1]

    L,U
end
```

```
A = randn(5,5) + 100I # need + 100I so that the matrix is (probably) diagonal
L,U = mylu(A)
@test A ≈ L*U
```

Out[2]: **Test Passed**

### Example 1 (by-hand)

Consider the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 4 & 8 \\ 1 & 4 & 9 \end{bmatrix}$$

We write

$$\begin{aligned} A &= \underbrace{\begin{bmatrix} 1 & & \\ 2 & 1 & \\ 1 & & 1 \end{bmatrix}}_{L_1} \underbrace{\begin{bmatrix} 1 & 1 & 1 \\ 0 & 2 & 6 \\ 0 & 3 & 8 \end{bmatrix}}_{\text{ }} = L_1 \underbrace{\begin{bmatrix} 1 & & \\ & 1 & \\ & 3/2 & 1 \end{bmatrix}}_{L_2} \underbrace{\begin{bmatrix} 1 & 1 & 1 \\ 0 & 2 & 6 \\ 0 & 0 & -1 \end{bmatrix}}_{\text{ }} \\ &= \underbrace{\begin{bmatrix} 1 & & \\ 2 & 1 & \\ 1 & 3/2 & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} 1 & 1 & 1 \\ 0 & 2 & 6 \\ 0 & 0 & -1 \end{bmatrix}}_U \end{aligned}$$

## 2. PLU Factorisation

We learned in first year linear algebra that if a diagonal entry is zero when doing Gaussian elimination one has to *row pivot*. For stability, in implementation one *always* pivots: swap the largest in magnitude entry for the entry on the diagonal.

We will see this is equivalent to a PLU decomposition:

**Theorem 1 (PLU)** A matrix  $A \in \mathbb{C}^{n \times n}$  is invertible if and only if it has a PLU decomposition:

$$A = P^\top LU$$

where the diagonal of  $L$  are all equal to 1 and the diagonal of  $U$  are all non-zero.

### Proof

If we have a PLU decomposition of this form then  $L$  and  $U$  are invertible and hence the inverse is simply  $A^{-1} = U^{-1}L^{-1}P$ .

If  $A \in \mathbb{C}^{1 \times 1}$  we trivially have an LU decomposition  $A = [1] * [a_{11}]$  as all  $1 \times 1$  matrices are triangular. We now proceed by induction: assume all invertible matrices of lower

dimension have a PLU factorisation. As  $A$  is invertible not all entries in the first column are zero. Therefore there exists a permutation  $P_1$  so that  $\alpha := (P_1 A)[1, 1] \neq 0$ . Hence we write

$$P_1 A = \begin{bmatrix} \alpha & \mathbf{w}^\top \\ \mathbf{v} & K \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & \\ \mathbf{v}/\alpha & I \end{bmatrix}}_{L_1} \begin{bmatrix} \alpha & \mathbf{w}^\top \\ & K - \mathbf{v}\mathbf{w}^\top/\alpha \end{bmatrix}$$

We deduce that  $A_2 := K - \mathbf{v}\mathbf{w}^\top/\alpha$  is invertible because  $A$  and  $L_1$  are invertible (Exercise).

By assumption we can write  $A_2 = P^\top LU$ . Thus we have:

$$\begin{aligned} \underbrace{\begin{bmatrix} 1 & \\ P & \end{bmatrix}}_P P_1 A &= \begin{bmatrix} 1 & \\ & P \end{bmatrix} \begin{bmatrix} \alpha & \mathbf{w}^\top \\ \mathbf{v} & A_2 \end{bmatrix} = \begin{bmatrix} 1 & \\ & P \end{bmatrix} L_1 \begin{bmatrix} \alpha & \mathbf{w}^\top \\ & P^\top LU \end{bmatrix} \\ &= \begin{bmatrix} 1 & \\ P\mathbf{v}/\alpha & P \end{bmatrix} \begin{bmatrix} 1 & \\ & P^\top L \end{bmatrix} \begin{bmatrix} \alpha & \mathbf{w}^\top \\ & U \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} 1 & \\ P\mathbf{v}/\alpha & L \end{bmatrix}}_L \underbrace{\begin{bmatrix} \alpha & \mathbf{w}^\top \\ & U \end{bmatrix}}_U. \end{aligned}$$

■

In the above we neglected to state which permutation is used as for the proof of existence it is immaterial. For *stability* however, we choose one that puts the largest entry: let  $\sigma_{\max} : \mathbb{R}^n \rightarrow S_n$  be the permutation that swaps the first row with the row of  $\mathbf{a}$  whose absolute value is maximised. In cycle notation we then have:

$$\sigma_{\max}(\mathbf{a}) = (1, \text{indmax}|\mathbf{a}|)$$

where `indmax` gives the index of the entry of a vector which is its maximum.

This inductive proof encodes an algorithm. Note that in the above, just like in Householder QR,  $P$  is a product of all permutations that come afterwards, that is, we can think of  $P$  as:

$$P = P_{n-1} \cdots P_3 P_2 \quad \text{for} \quad P_j = \begin{bmatrix} I_{j-2} & \\ & P_j \end{bmatrix}$$

where  $P_j$  is a single permutation corresponding to the first column of  $A_j$ . That is, we have

$$P\mathbf{v} = P_{n-1} \cdots P_3 P_2 \mathbf{v}.$$

**Algorithm 2 (PLU)** This can be implemented in Julia as follows:

```
In [3]: function σ_max(a)
    n = length(a)
    mx, ind = findmax(abs.(a)) # finds the index of the maximum entry
    if ind == 1
        1:n
    else
        [ind; 2:ind-1; 1; ind+1:n]
    end
end

function plu(A)
    n,m = size(A)
    if n ≠ m
        error("Matrix must be square")
    end
    T = eltype(A)
    L = LowerTriangular(zeros(T,n,n))
    U = UpperTriangular(zeros(T,n,n))

    σ = Vector(1:n)

    A_j = copy(A)

    for j = 1:n-1
        σ₁ = σ_max(A_j[:,1])
        P₁A_j = A_j[σ₁,:] # permute rows of A_j
        α,v,w = P₁A_j[1,1],P₁A_j[2:end,1],P₁A_j[1,2:end]
        K = P₁A_j[2:end,2:end]

        # populate data
        L[j,j] = 1
        L[j+1:end,j] = v/α
        U[j,j] = α
        U[j,j+1:end] = w

        # apply permutation to previous L
        # and compose the permutations
        L[j:n,1:j-1] = L[(j:n)[σ₁],1:j-1]
        σ[j:n] = σ[j:n][σ₁]

        # this is the "recursion": A_j is now the next block
        # We use transpose(w) instead of w' incase w is complex
        A_j = K - v*transpose(w)/α
    end
    # j = n case
    L[n,n] = 1
    U[n,n] = A_j[1,1]

    L,U,σ
end

A = randn(5,5)
L,U,σ = plu(A)
@test L*U ≈ A[σ,:]
```

Out[3]: Test Passed

## Example 2

Again we consider the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 4 & 8 \\ 1 & 4 & 9 \end{bmatrix}$$

Even though  $a_{11} = 1 \neq 0$ , we still pivot: placing the maximum entry on the diagonal to mitigate numerical errors. That is, we first pivot and upper triangularise the first column:

$$\underbrace{\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ & 1 \end{bmatrix}}_{P_1} A = \begin{bmatrix} 2 & 4 & 8 \\ 1 & 1 & 1 \\ 1 & 4 & 9 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & & \\ 1/2 & 1 & \\ 1/2 & & 1 \end{bmatrix}}_{L_1} \begin{bmatrix} 2 & 4 & 8 \\ -1 & -3 & \\ 2 & 5 & \end{bmatrix}$$

That is we have  $\alpha_1 = 2$ ,  $\mathbf{v}_1 = [1, 1]$ , and  $\mathbf{w}_1 = [4, 8]$ . We now pivot for  $A_2$ :

$$\underbrace{\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}}_{P_2} \underbrace{\begin{bmatrix} -1 & -3 \\ 2 & 5 \end{bmatrix}}_{A_2} = \begin{bmatrix} 2 & 5 \\ -1 & -3 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & \\ -1/2 & 1 \end{bmatrix}}_{L_2} \begin{bmatrix} 2 & 5 \\ -\frac{1}{2} & \end{bmatrix}$$

Note that  $P_2 \mathbf{v}_1 = \mathbf{v}_1$  and

$$P = P_2 P_1 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

Hence we have

$$PA = \begin{bmatrix} 1 & & \\ 1/2 & 1 & \\ 1/2 & -1/2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & 8 \\ 2 & 5 \\ -1/2 \end{bmatrix}$$

We see how this example is done on a computer:

```
In [4]: A = [1 1 1;
           2 4 8;
           1 4 9]
L,U,sigma = lu(A) # sigma is a vector encoding the permutation

@test L == [1      0      0 ;
            1/2    1      0 ;
            1/2   -1/2   1 ]
@test U == [2 4 8 ;
            0 2 5 ;
            0 0 -1/2]
@test I(3)[sigma,:] == [0 1 0 ;
```

```

0 0 1 ;
1 0 0

```

Out[4]: **Test Passed**

To invert a system we can do:

```
In [5]: b = randn(3)
@test U\(\mathbf{L}\backslash b[\sigma]) == A\b
```

Out[5]: **Test Passed**

Note the entries match exactly because this is precisely what `\` is using.

### 3. Cholesky Factorisation

Cholesky Factorisation is a form of Gaussian elimination (without pivoting) that exploits symmetry in the problem, resulting in a substantial speedup. It is only relevant for *symmetric positive definite* (SPD) matrices.

**Definition 1 (positive definite)** A square matrix  $A \in \mathbb{R}^{n \times n}$  is *positive definite* if for all  $\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq 0$  we have

$$\mathbf{x}^\top A \mathbf{x} > 0$$

First we establish some basic properties of positive definite matrices:

**Proposition 3 (conj. pos. def.)** If  $A \in \mathbb{R}^{n \times n}$  is positive definite and  $V \in \mathbb{R}^{n \times n}$  is non-singular then

$$V^\top A V$$

is positive definite.

**Proposition 4 (diag positivity)** If  $A \in \mathbb{R}^{n \times n}$  is positive definite then its diagonal entries are positive:  $a_{kk} > 0$ .

**Theorem 1 (subslice pos. def.)** If  $A \in \mathbb{R}^{n \times n}$  is positive definite and  $\mathbf{k} \in \{1, \dots, n\}^m$  is a vector of  $m$  integers where any integer appears only once, then  $A[\mathbf{k}, \mathbf{k}] \in \mathbb{R}^{m \times m}$  is also positive definite.

We leave the proofs to the problem sheets. Here is the key result:

**Theorem 2 (Cholesky and SPD)** A matrix  $A$  is symmetric positive definite if and only if it has a Cholesky factorisation

$$A = LL^\top$$

where the diagonals of  $L$  are positive.

**Proof** If  $A$  has a Cholesky factorisation it is symmetric ( $A^\top = (LL^\top)^\top = A$ ) and for  $\mathbf{x} \neq 0$  we have

$$\mathbf{x}^\top A \mathbf{x} = (\mathbf{L}\mathbf{x})^\top \mathbf{L}\mathbf{x} = \|\mathbf{L}\mathbf{x}\|^2 > 0$$

where we use the fact that  $L$  is non-singular.

For the other direction we will prove it by induction, with the  $1 \times 1$  case being trivial. Assume all lower dimensional symmetric positive definite matrices have Cholesky decompositions. Write

$$A = \begin{bmatrix} \alpha & \mathbf{v}^\top \\ \mathbf{v} & K \end{bmatrix} = \underbrace{\begin{bmatrix} \sqrt{\alpha} & \\ \frac{\mathbf{v}}{\sqrt{\alpha}} & I \end{bmatrix}}_{L_1} \underbrace{\begin{bmatrix} 1 & \\ & K - \frac{\mathbf{v}\mathbf{v}^\top}{\alpha} \end{bmatrix}}_{A_1} \underbrace{\begin{bmatrix} \sqrt{\alpha} & \frac{\mathbf{v}^\top}{\sqrt{\alpha}} \\ & I \end{bmatrix}}_{L_1^\top}.$$

Note that  $A_2 := K - \frac{\mathbf{v}\mathbf{v}^\top}{\alpha}$  is a subslice of  $L_1^{-1}AL_1^{-\top}$ , hence by the previous propositions is itself symmetric positive definite. Thus we can write

$$A_2 = K - \frac{\mathbf{v}\mathbf{v}^\top}{\alpha} = LL^\top$$

and hence  $A = LL^\top$  for

$$L = L_1 \begin{bmatrix} 1 & \\ & L \end{bmatrix} = \begin{bmatrix} \sqrt{\alpha} & \\ \frac{\mathbf{v}}{\sqrt{\alpha}} & L \end{bmatrix}$$

satisfies  $A = LL^\top$ . ■

Note hidden in this proof is a simple algorithm for computing the Cholesky factorisation.

### Algorithm 3 (Cholesky)

```
In [6]: function mycholesky(A)
    T = eltype(A)
    n,m = size(A)
    if n ≠ m
        error("Matrix must be square")
    end
    if A ≠ A'
        error("Matrix must be symmetric")
    end
    T = eltype(A)
    L = LowerTriangular(zeros(T,n,n))
    A_j = copy(A)
    for j = 1:n
        α,v = A_j[1,1], A_j[2:end,1]
        if α ≤ 0
            error("Matrix is not SPD")
        end
```

```

L[j,j] = sqrt(alpha)
L[j+1:end,j] = v/sqrt(alpha)

# induction part
A_j = A_j[2:end,2:end] - v*v'/alpha
end
L
end

A = Symmetric(rand(100,100) + 100I)
L = mycholesky(A)
@test A ≈ L*L'

```

Out[6]: **Test Passed**

This algorithm succeeds if and only if  $A$  is symmetric positive definite.

**Example 3 (Cholesky by hand)** Consider the matrix

$$A = \begin{bmatrix} 2 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 \\ 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 2 \end{bmatrix}$$

Then  $\alpha_1 = 2$ ,  $\mathbf{v}_1 = [1, 1, 1]$ , and

$$A_2 = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^T = \frac{1}{2} \begin{bmatrix} 3 & 1 & 1 \\ 1 & 3 & 1 \\ 1 & 1 & 3 \end{bmatrix}.$$

Continuing, we have  $\alpha_2 = 3/2$ ,  $\mathbf{v}_2 = [1/2, 1/2]$ , and

$$A_3 = \frac{1}{2} \left( \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix}^T \right) = \frac{1}{3} \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix}$$

Next,  $\alpha_3 = 4/3$ ,  $\mathbf{v}_3 = [1]$ , and

$$A_4 = [4/3 - 3/4 * (1/3)^2] = [5/4]$$

i.e.  $\alpha_4 = 5/4$ .

Thus we get

$$L = \begin{bmatrix} \sqrt{\alpha_1} & & & \\ \frac{\mathbf{v}_1[1]}{\sqrt{\alpha_1}} & \sqrt{\alpha_2} & & \\ \frac{\mathbf{v}_1[2]}{\sqrt{\alpha_1}} & \frac{\mathbf{v}_2[1]}{\sqrt{\alpha_2}} & \sqrt{\alpha_3} & \\ \frac{\mathbf{v}_1[3]}{\sqrt{\alpha_1}} & \frac{\mathbf{v}_2[2]}{\sqrt{\alpha_2}} & \frac{\mathbf{v}_3[1]}{\sqrt{\alpha_3}} & \sqrt{\alpha_4} \end{bmatrix} = \begin{bmatrix} \sqrt{2} & & & \\ \frac{1}{\sqrt{2}} & \sqrt{\frac{3}{2}} & & \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{3}} & \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{12}} & \frac{\sqrt{5}}{2} \end{bmatrix}$$

We can check if our answer is correct:

```
In [7]: A = ones(4,4) + I
# The inbuilt cholesky returns a special type whose field L is the factor
@test cholesky(A).L ≈ [ sqrt(2) 0 0 0;
                         1/sqrt(2) sqrt(3/2) 0 0;
                         1/sqrt(2) 1/sqrt(6) 2/sqrt(3) 0;
                         1/sqrt(2) 1/sqrt(6) 1/sqrt(12) sqrt(5)/2]
```

Out[7]: **Test Passed**

## 4. Timings and Stability

The different factorisations have trade-offs between speed and stability. First we compare the speed of the different factorisations on a symmetric positive definite matrix, from fastest to slowest:

```
In [8]: n = 100
A = Symmetric(rand(n,n)) + 100I # shift by 10 ensures positivity
@btime cholesky(A);
@btime lu(A);
@btime qr(A);
```

```
292.085 μs (3 allocations: 78.20 KiB)
79.052 μs (4 allocations: 79.08 KiB)
256.735 μs (7 allocations: 134.55 KiB)
```

On my machine, `cholesky` is ~1.5x faster than `lu`, which is ~2x faster than QR.

In terms of stability, QR computed with Householder reflections (and Cholesky for positive definite matrices) are stable, whereas LU is usually unstable (unless the matrix is diagonally dominant). PLU is a very complicated story: in theory it is unstable, but the set of matrices for which it is unstable is extremely small, so small one does not normally run into them.

Here is an example matrix that is in this set.

```
In [9]: function badmatrix(n)
    A = Matrix(1I, n, n)
    A[:,end] .= 1
    for j = 1:n-1
        A[j+1:end,j] .= -1
    end
    A
end
A = badmatrix(5)
```

```
Out[9]: 5×5 Matrix{Int64}:
```

```
1 0 0 0 1
-1 1 0 0 1
-1 -1 1 0 1
-1 -1 -1 1 1
-1 -1 -1 -1 1
```

Note that pivoting will not occur (we do not pivot as the entries below the diagonal are the same magnitude as the diagonal), thus the PLU Factorisation is equivalent to an LU factorisation:

```
In [10]: L, U = lu(A)
```

```
Out[10]: LU{Float64, Matrix{Float64}, Vector{Int64}}
```

L factor:

5×5 Matrix{Float64}:

```
1.0 0.0 0.0 0.0 0.0
-1.0 1.0 0.0 0.0 0.0
-1.0 -1.0 1.0 0.0 0.0
-1.0 -1.0 -1.0 1.0 0.0
-1.0 -1.0 -1.0 -1.0 1.0
```

U factor:

5×5 Matrix{Float64}:

```
1.0 0.0 0.0 0.0 1.0
0.0 1.0 0.0 0.0 2.0
0.0 0.0 1.0 0.0 4.0
0.0 0.0 0.0 1.0 8.0
0.0 0.0 0.0 0.0 16.0
```

But here we see an issue: the last column of `U` is growing exponentially fast! Thus when `n` is large we get very large errors:

```
In [11]: n = 100
b = randn(n)
A = badmatrix(n)
norm(A\b - qr(A)\b) # A \ b still uses lu
```

```
Out[11]: 4.474199133706833
```

Note `qr` is completely fine:

```
In [12]: norm(qr(A)\b - qr(big.(A))\b) # roughly machine precision
```

```
Out[12]: 7.8671462675752303843621536723335856774082232211278287852647793570089031949
00461e-15
```

Amazingly, PLU is fine if applied to a small perturbation of `A`:

```
In [13]: ε = 0.000001
Aε = A .+ ε .* randn_()
norm(Aε \ b - qr(Aε)\b) # Now it matches!
```

```
Out[13]: 1.1117203027706248e-14
```

The big *open problem* in numerical linear algebra is to prove that the set of matrices for which PLU fails has extremely small measure.

## II.5 Norms

In this lecture we discuss matrix and vector norms.

1. Vector norms: we discuss the standard  $p$ -norm for vectors in  $\mathbb{R}^n$ .
2. Matrix norms: we discuss how two vector norms can be used to induce a norm on matrices. These

satisfy an additional multiplicative inequality.

### 1. Vector norms

Recall the definition of a (vector-)norm:

**Definition 1 (vector-norm)** A norm  $\|\cdot\|$  on a vector space  $V$  (e.g.  $\mathbb{R}^n$  or  $\mathbb{C}^n$ ) over a field  $\mathbb{F}$  (e.g.  $\mathbb{R}$  or  $\mathbb{C}$ )

is a function that satisfies the following, for  $\mathbf{x}, \mathbf{y} \in V$  and  $c \in \mathbb{F}$ :

1. Triangle inequality:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
2. Homogeneity:  $\|c\mathbf{x}\| = |c|\|\mathbf{x}\|$
3. Positive-definiteness:  $\|\mathbf{x}\| = 0$  implies that  $\mathbf{x} = 0$ .

Consider the following example:

**Definition 2 (p-norm)** For  $1 \leq p < \infty$  and  $\mathbf{x} \in \mathbb{C}^n$ , define the  $p$ -norm:

$$\|\mathbf{x}\|_p := \left( \sum_{k=1}^n |x_k|^p \right)^{1/p}$$

where  $x_k$  is the  $k$ -th entry of  $\mathbf{x}$ . For  $p = \infty$  we define

$$\|\mathbf{x}\|_\infty := \max_k |x_k|$$

**Theorem 1 (p-norm)**  $\|\cdot\|_p$  is a norm for  $1 \leq p \leq \infty$ .

#### Proof

We will only prove the case  $p = 1, 2, \infty$  as general  $p$  is more involved.

Homogeneity and positive-definiteness are straightforward: e.g.,

$$\|c\mathbf{x}\|_p = \left( \sum_{k=1}^n |cx_k|^p \right)^{1/p} = \left( |c|^p \sum_{k=1}^n |x_k|^p \right)^{1/p} = |c|\|\mathbf{x}\|$$

and if  $\|\mathbf{x}\|_p = 0$  then all  $|x_k|^p$  are have to be zero.

For  $p = 1, \infty$  the triangle inequality is also straightforward:

$$\|\mathbf{x} + \mathbf{y}\|_\infty = \max_k(|x_k + y_k|) \leq \max_k(|x_k| + |y_k|) \leq \|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty$$

and

$$\|\mathbf{x} + \mathbf{y}\|_1 = \sum_{k=1}^n |x_k + y_k| \leq \sum_{k=1}^n (|x_k| + |y_k|) = \|\mathbf{x}\|_1 + \|\mathbf{y}\|_1$$

For  $p = 2$  it can be proved using the Cauchy–Schwartz inequality:

$$|\mathbf{x}^\star \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$$

That is, we have

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\mathbf{x}^\top \mathbf{y} + \|\mathbf{y}\|^2 \leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 = (\|\mathbf{x}\| + \|\mathbf{y}\|)^2$$

■

In Julia, one can use the inbuilt `norm` function to calculate norms:

```
norm([1,-2,3]) == norm([1,-2,3], 2) == sqrt(1^2 + 2^2 + 3^2) ==
sqrt(14);
norm([1,-2,3], 1) == 1 + 2 + 3 == 6;
norm([1,-2,3], Inf) == 3;
```

## 2. Matrix norms

Just like vectors, matrices have norms that measure their "length". The simplest example is the Fröbenius norm:

**Definition 3 (Fröbenius norm)** For  $A \in \mathbb{C}^{m \times n}$  define

$$\|A\|_F := \sqrt{\sum_{k=1}^m \sum_{j=1}^n |a_{kj}|^2}$$

This is available as `norm` in Julia:

```
In [1]: A = randn(5,3)
norm(A) == norm(vec(A))
```

```
Out[1]: true
```

While this is the simplest norm, it is not the most useful. Instead, we will build a matrix norm from a vector norm:

**Definition 4 (matrix-norm)** Suppose  $A \in \mathbb{C}^{m \times n}$  and consider two norms  $\|\cdot\|_X$  on  $\mathbb{C}^m$  and  $\|\cdot\|_Y$  on  $\mathbb{C}^n$ . Define the (*induced*) *matrix norm* as:

$$\|A\|_{X \rightarrow Y} := \sup_{\mathbf{v}: \|\mathbf{v}\|_X=1} \|A\mathbf{v}\|_Y$$

Also define

$$\|A\|_X := \|A\|_{X \rightarrow X}$$

For the induced  $p$ -norm we use the notation  $\|A\|_p$ .

Note an equivalent definition of the induced norm:

$$\|A\|_{X \rightarrow Y} = \sup_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|_Y}{\|\mathbf{x}\|_X}$$

This follows since we can scale  $\mathbf{x}$  by its norm so that it has unit norm, that is,  $\frac{\mathbf{x}}{\|\mathbf{x}\|_X}$  has unit norm.

**Lemma 1 (matrix norms are norms)** Induced matrix norms are norms, that is for  $\|\cdot\| = \|\cdot\|_{X \rightarrow Y}$  we have:

1. Triangle inequality:  $\|A + B\| \leq \|A\| + \|B\|$
2. Homogeneity:  $\|cA\| = |c|\|A\|$
3. Positive-definiteness:  $\|A\| = 0 \Rightarrow A = 0$

In addition, they satisfy the following additional properties:

1.  $\|A\mathbf{x}\|_Y \leq \|A\|_{X \rightarrow Y} \|\mathbf{x}\|_X$
2. Multiplicative inequality:  $\|AB\|_{X \rightarrow Z} \leq \|A\|_{Y \rightarrow Z} \|B\|_{X \rightarrow Y}$

### Proof

First we show the *triangle inequality*:

$$\|A + B\| \leq \sup_{\mathbf{v}: \|\mathbf{v}\|_X=1} (\|A\mathbf{v}\|_Y + \|B\mathbf{v}\|_Y) \leq \|A\| + \|B\|.$$

Homogeneity is also immediate. Positive-definiteness follows from the fact that if  $\|A\| = 0$  then  $A\mathbf{x} = 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ . The property  $\|A\mathbf{x}\|_Y \leq \|A\|_{X \rightarrow Y} \|\mathbf{x}\|_X$  follows from the definition. Finally, the multiplicative inequality follows from

$$\|AB\| = \sup_{\mathbf{v}: \|\mathbf{v}\|_X=1} \|AB\mathbf{v}\|_Z \leq \sup_{\mathbf{v}: \|\mathbf{v}\|_X=1} \|A\|_{Y \rightarrow Z} \|B\mathbf{v}\|_Y = \|A\|_{Y \rightarrow Z} \|B\|_{X \rightarrow Y}$$

■

We have some simple examples of induced norms:

**Example 1 (1-norm)** We claim

$$\|A\|_1 = \max_j \|\mathbf{a}_j\|_1$$

that is, the maximum 1-norm of the columns. To see this use the triangle inequality to find for  $\|\mathbf{x}\|_1 = 1$

$$\|A\mathbf{x}\|_1 \leq \sum_{j=1}^n |x_j| \|\mathbf{a}_j\|_1 \leq \max_j \|\mathbf{a}_j\| \sum_{j=1}^n |x_j| = \max_j \|\mathbf{a}_j\|_1.$$

But the bound is also attained since if  $j$  is the column that maximises the norms then

$$\|A\mathbf{e}_j\|_1 = \|\mathbf{a}_j\|_1 = \max_j \|\mathbf{a}_j\|_1.$$

In the problem sheet we see that

$$\|A\|_\infty = \max_k \|A[k, :]\|_1$$

that is, the maximum 1-norm of the rows.

Matrix norms are available via `opnorm`:

```
In [2]: m,n = 5,3
A = randn(m,n)
opnorm(A,1) == maximum(norm(A[:,j],1) for j = 1:n)
opnorm(A,Inf) == maximum(norm(A[k,:],1) for k = 1:m)
opnorm(A) # the 2-norm
```

Out[2]: 2.6041739084301563

An example that does not have a simple formula is  $\|A\|_2$ , but we do have two simple cases:

**Proposition 1 (diagonal/orthogonal 2-norms)** If  $\Lambda$  is diagonal with entries  $\lambda_k$  then  $\|\Lambda\|_2 = \max_k |\lambda_k|$ . If  $Q$  is orthogonal then  $\|Q\| = 1$ .

In the next chapter we see how the 2-norm for a matrix can be defined in terms of the *Singular Value Decomposition*.

## II.6 Singular Value Decomposition

In this chapter we discuss the *Singular Value Decomposition (SVD)*: a matrix factorisation that encodes how much a matrix "stretches" a random vector. This includes *singular values*, the largest of which dictates the 2-norm of the matrix.

**Definition 1 (singular value decomposition)** For  $A \in \mathbb{C}^{m \times n}$  with rank  $r > 0$ , the (reduced) singular value decomposition (SVD) is

$$A = U\Sigma V^*$$

where  $U \in \mathbb{C}^{m \times r}$  and  $V \in \mathbb{C}^{n \times r}$  have orthonormal columns and  $\Sigma \in \mathbb{R}^{r \times r}$  is diagonal whose diagonal entries, which we call *singular values*, are all positive and non-increasing:  $\sigma_1 \geq \dots \geq \sigma_r > 0$ . The full singular value decomposition (SVD) is

$$A = U\Sigma V^*$$

where  $U \in U(m)$  and  $V \in U(n)$  are unitary matrices and  $\Sigma \in \mathbb{R}^{m \times n}$  has only diagonal non-zero entries, i.e., if  $m > n$ ,

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & 0 & \\ & & \vdots & \\ & & 0 & \end{bmatrix}$$

and if  $m < n$ ,

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_m & 0 & \cdots & 0 \end{bmatrix}$$

where  $\sigma_k = 0$  if  $k > r$ .

In particular, we discuss:

1. Existence of the SVD: we show that an SVD exists by relating it to the eigenvalue Decomposition of  $A^*A$  and  $AA^*$ .
2. 2-norm and SVD: the 2-norm of a matrix is defined in terms of the largest singular value.
3. Best rank- $k$  approximation and compression: the best approximation of a matrix by a smaller rank matrix can be constructed

using the SVD, which gives an effective way to compress matrices.

```
In [1]: using LinearAlgebra, Plots
```

## 1. Existence

To show the SVD exists we first establish some properties of a *Gram matrix* ( $A^*A$ ):

**Proposition 1 (Gram matrix kernel)** The kernel of  $A$  is also the kernel of  $A^*A$ .

**Proof** If  $A^*Ax = 0$  then we have

$$0 = \mathbf{x}^* A^* A \mathbf{x} = \|A\mathbf{x}\|^2$$

which means  $A\mathbf{x} = 0$  and  $\mathbf{x} \in \ker(A)$ . ■

**Proposition 2 (Gram matrix diagonalisation)** The Gram-matrix satisfies

$$A^*A = Q\Lambda Q^* \in \mathbb{C}^{n \times n}$$

is a Hermitian matrix where  $Q \in U(n)$  and the eigenvalues  $\lambda_k$  are real and non-negative. If  $A \in \mathbb{R}^{m \times n}$  then  $Q \in O(n)$ .

**Proof**  $A^*A$  is Hermitian so we appeal to the spectral theorem for the existence of the decomposition, and the fact that the eigenvalues are real. For the corresponding (orthonormal) eigenvector  $\mathbf{q}_k$ ,

$$\lambda_k = \lambda_k \mathbf{q}_k^* \mathbf{q}_k = \mathbf{q}_k^* A^* A \mathbf{q}_k = \|A\mathbf{q}_k\|^2 \geq 0.$$

■

This connection allows us to prove existence:

**Theorem 1 (SVD existence)** Every  $A \in \mathbb{C}^{m \times n}$  has an SVD.

**Proof** Consider

$$A^*A = Q\Lambda Q^*.$$

Assume (as usual) that the eigenvalues are sorted in decreasing modulus, and so  $\lambda_1, \dots, \lambda_r$  are an enumeration of the non-zero eigenvalues and

$$V := [\mathbf{q}_1 | \cdots | \mathbf{q}_r]$$

the corresponding (orthonormal) eigenvectors, with

$$K = [\mathbf{q}_{r+1} | \cdots | \mathbf{q}_n]$$

the corresponding kernel. Define

$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_r} \end{bmatrix}$$

Now define

$$U := AV\Sigma^{-1}$$

which is orthogonal since  $A^*AV = V\Sigma^2$ :

$$U^*U = \Sigma^{-1}V^*A^*AV\Sigma^{-1} = I.$$

Thus we have

$$U\Sigma V^* = AVV^* = A\begin{bmatrix} V | K \end{bmatrix} \begin{bmatrix} V^* \\ K^* \end{bmatrix}$$

where we use the fact that  $AK = 0$  so that concatenating  $K$  does not change the value.

■

## 2. 2-norm and SVD

Singular values tell us the 2-norm:

**Corollary 1 (singular values and norm)**

$$\|A\|_2 = \sigma_1$$

and if  $A \in \mathbb{C}^{n \times n}$  is invertible, then

$$\|A^{-1}\|_2 = \sigma_n^{-1}$$

**Proof**

First we establish the upper-bound:

$$\|A\|_2 \leq \|U\|_2 \|\Sigma\|_2 \|V^*\|_2 = \|\Sigma\|_2 = \sigma_1$$

This is attained using the first right singular vector:

$$\|A\mathbf{v}_1\|_2 = \|\Sigma V^* \mathbf{v}_1\|_2 = \|\Sigma \mathbf{e}_1\|_2 = \sigma_1$$

The inverse result follows since the inverse has SVD

$$A^{-1} = V\Sigma^{-1}U^* = (VW)(W\Sigma^{-1}W)(WU)^*$$

is the SVD of  $A^{-1}$ , i.e.  $VW \in U(n)$  are the left singular vectors and  $WU$  are the right singular vectors, where

$$W := P_\sigma = \begin{bmatrix} & & & 1 \\ & \ddots & & \\ 1 & & & \end{bmatrix}$$

is the permutation that reverses the entries, that is,  $\sigma$  has Cauchy notation

$$\begin{pmatrix} 1 & 2 & \cdots & n \\ n & n-1 & \cdots & 1 \end{pmatrix}.$$

■

We will not discuss in this module computation of singular value decompositions or eigenvalues: they involve iterative algorithms (actually built on a sequence of QR decompositions).

### 3. Best rank- $k$ approximation and compression

One of the main usages for SVDs is low-rank approximation:

**Theorem 2 (best low rank approximation)** The matrix

$$A_k := [\mathbf{u}_1 | \cdots | \mathbf{u}_k] \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \end{bmatrix} [\mathbf{v}_1 | \cdots | \mathbf{v}_k]^*$$

is the best 2-norm approximation of  $A$  by a rank  $k$  matrix, that is, for all rank- $k$  matrices  $B$ , we have

$$\|A - A_k\|_2 \leq \|A - B\|_2.$$

**Proof** We have

$$A - A_k = U \begin{bmatrix} 0 & & & \\ & \ddots & & \\ & & 0 & \\ & & & \sigma_{k+1} \\ & & & & \ddots \\ & & & & & \sigma_r \end{bmatrix} V^*.$$

Suppose a rank- $k$  matrix  $B$  has

$$\|A - B\|_2 < \|A - A_k\|_2 = \sigma_{k+1}.$$

For all  $\mathbf{w} \in \ker(B)$  we have

$$\|A\mathbf{w}\|_2 = \|(A - B)\mathbf{w}\|_2 \leq \|A - B\|\|\mathbf{w}\|_2 < \sigma_{k+1}\|\mathbf{w}\|_2$$

But for all  $\mathbf{u} \in \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_{k+1})$ , that is,  $\mathbf{u} = V[:, 1:k+1]\mathbf{c}$  for some  $\mathbf{c} \in \mathbb{C}^{k+1}$  we have

$$\|A\mathbf{u}\|_2^2 = \|U\Sigma_k\mathbf{c}\|_2^2 = \|\Sigma_k\mathbf{c}\|_2^2 = \sum_{j=1}^{k+1} (\sigma_j c_j)^2 \geq \sigma_{k+1}^2 \|\mathbf{c}\|^2,$$

i.e.,  $\|A\mathbf{u}\|_2 \geq \sigma_{k+1}\|\mathbf{u}\|$ , where we use the fact that

$$\|\mathbf{u}\|^2 = \|V[:, 1:k+1]\mathbf{c}\|^2 = \mathbf{c}^* V[:, 1:k+1]^* V[:, 1:k+1] \mathbf{c} = \mathbf{c}^* \mathbf{c} = \|\mathbf{c}\|^2.$$

Thus  $\mathbf{w}$  cannot be in this span.

The dimension of the span of  $\ker(B)$  is at least  $n - k$ , but the dimension of  $\text{span}(\mathbf{v}_1, \dots, \mathbf{v}_{k+1})$  is at least  $k + 1$ . Since these two spaces cannot intersect (apart from at 0) we have a contradiction, since  $(n - r) + (r + 1) = n + 1 > n$ . ■

**Example 1 (Hilbert matrix)** Here we show an example of a simple low-rank approximation using the SVD. Consider the Hilbert matrix:

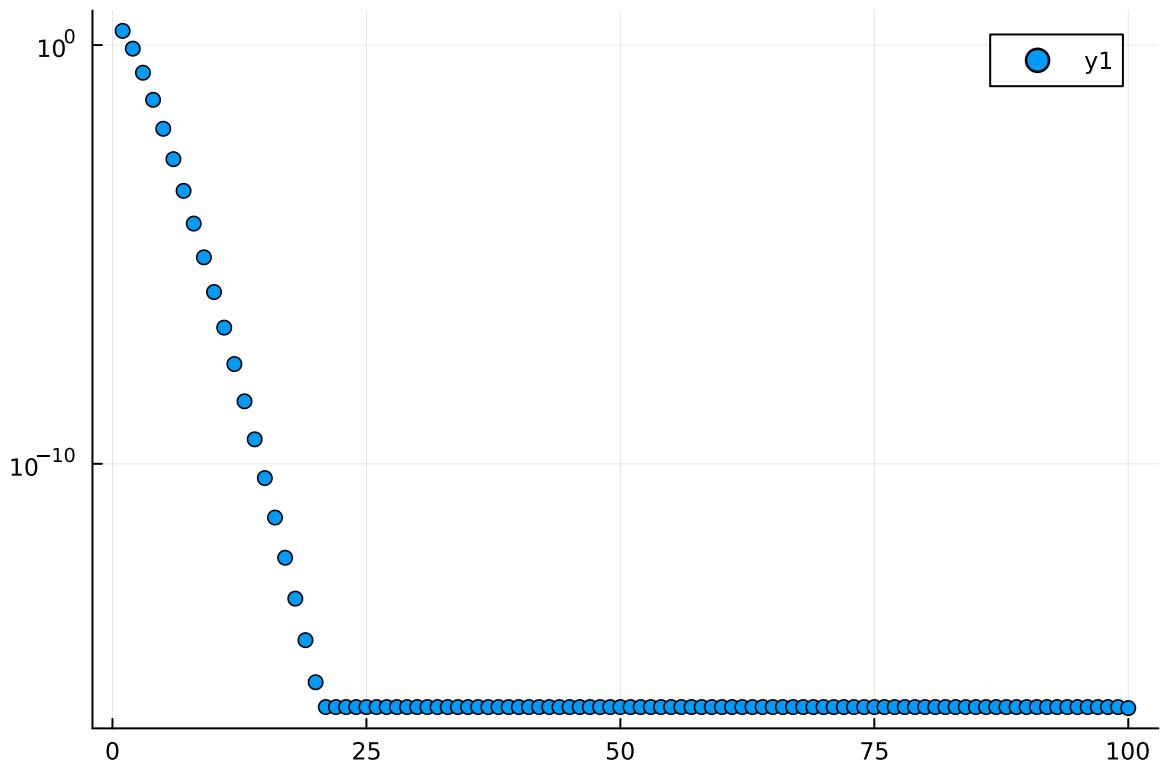
```
In [2]: hilbertmatrix(n) = [1/(k+j-1) for j = 1:n, k=1:n]
hilbertmatrix(5)
```

```
Out[2]: 5×5 Matrix{Float64}:
 1.0      0.5      0.333333  0.25      0.2
 0.5      0.333333  0.25      0.2      0.166667
 0.333333  0.25      0.2      0.166667  0.142857
 0.25      0.2      0.166667  0.142857  0.125
 0.2      0.166667  0.142857  0.125      0.111111
```

That is, the  $H[k, j] = 1/(k + j - 1)$ . This is a famous example of matrix with rapidly decreasing singular values:

```
In [3]: H = hilbertmatrix(100)
U,σ,V = svd(H)
scatter(σ;yscale=:log10)
```

Out[3]:



Note numerically we typically do not get exactly zero singular values so the rank is always treated as  $\min(m, n)$ . Because the singular values decay rapidly we can approximate the matrix very well with a rank 20 matrix:

In [4]:

```
k = 20 # rank
Σ_k = Diagonal(σ[1:k])
U_k = U[:,1:k]
V_k = V[:,1:k]
opnorm(U_k * Σ_k * V_k' - H)
```

Out[4]: 8.20222266307798e-16

Note that this can be viewed as a *compression* algorithm: we have replaced a matrix with  $100^2 = 10,000$  entries by two matrices and a vector with 4,000 entries without losing any information. In the problem sheet we explore the usage of low rank approximation to smooth functions and to compress images.

## II.7 Condition numbers

We have seen that floating point arithmetic induces errors in computations, and that we can typically bound the absolute errors to be proportional to  $C\epsilon_m$ . We want a way to bound the effect of more complicated calculations like computing  $A\mathbf{x}$  or  $A^{-1}\mathbf{y}$  without having to deal with the exact nature of floating point arithmetic, as it will depend on the *data A and x*. That is, we want to reduce floating point stability to a more fundamental property: *mathematical stability*: how does a mathematical operation like  $A\mathbf{x}$  change in the presence of small perturbations (random noise or structured floating point errors)?

1. Backward error analysis: We introduce the concept of *backward error analysis*, which is a more practical

way of understanding and bounding floating point errors. 2. Condition numbers: We introduce a *condition numbers*, which can capture the effect of perturbations in  $A$  for linear algebra operations. More precisely: matrix operations are mathematically *stable* when the condition number is small. 3. Bounding floating point errors for linear algebra: we see how simple operations like  $A\mathbf{x}$  can be put into a backward error analysis framework, leading to bounds on the errors in terms of the condition number.

### 1. Backward error analysis

So far we have done forward error analysis, e.g., to understand  $f(x) \approx f^{\text{FP}}(x)$  we consider either the absolute

$$f^{\text{FP}}(x) = f(x) + \delta_a$$

or relative

$$f^{\text{FP}}(x) = f(x)(1 + \delta_r)$$

errors of the *output*. More generally, for two vector spaces  $V$  and  $W$  (e.g.  $V = \mathbb{R}^n$  and  $W = \mathbb{R}^m$ ) consider functions  $\mathbf{f} = V \rightarrow W$ . We write

$$\mathbf{f}^{\text{FP}}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) + \delta_f$$

where we bound a norm of  $\delta_f \in W$  either *absolutely*:

$$\|\delta_f\|_W \leq C\varepsilon$$

or *relative* to the true result:

$$\|\delta_f\|_W \leq C\|\mathbf{f}(\mathbf{x})\|_W\varepsilon$$

(which is similar to PS4, Q1.3).

On the other hand, *backward error analysis* considers computations as errors in the *input*. That is, one writes the approximate function as the true function with a perturbed input: e.g. find  $\tilde{\mathbf{x}} \in V$  such that

$$\mathbf{f}^{\text{FP}}(\mathbf{x}) = \mathbf{f}(\tilde{\mathbf{x}}).$$

We study the *backward error*, the error in the input

$$\tilde{\mathbf{x}} = \mathbf{x} + \delta_b$$

where  $\delta_b \in \mathbb{R}^n$  by bounding either the absolute error

$$\|\delta_b\|_V \leq C\varepsilon$$

or relative error:

$$\delta_b\|_V \leq C\|\mathbf{x}\|_V\varepsilon$$

We shall see that some algorithms (like `mul_rows`) lead naturally to backward error results.

## 2. Condition numbers

So now we get to a mathematical question independent of floating point: can we bound the *relative error* in approximating

$$A\mathbf{x} \approx (A + \delta A)\mathbf{x}$$

if we know a bound on the relative backward error  $\|\delta A\|$ ? It turns out we can in terms of the *condition number* of the matrix:

**Definition 2 (condition number)** For a square matrix  $A$ , the *condition number* (in  $p$ -norm) is

$$\kappa_p(A) := \|A\|_p \|A^{-1}\|_p$$

with the default being the 2-norm condition number, writable in terms of the singular values as:

$$\kappa(A) := \kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_1}{\sigma_n}.$$

**Theorem 1 (relative forward error for matrix-vector)** Assume we have the relative backward error bound  $\|\delta A\| \leq \|A\|\varepsilon$ . Then for

$$(A + \delta A)\mathbf{x} = A\mathbf{x} + \delta_f$$

the forward error has the relative error bound

$$\|\delta_f\| \leq \|A\mathbf{x}\| \kappa(A) \varepsilon$$

**Proof** We can assume  $A$  is invertible (as otherwise  $\kappa(A) = \infty$ ). Denote  $\mathbf{y} = A\mathbf{x}$  and we have

$$\frac{\|\mathbf{x}\|}{\|A\mathbf{x}\|} = \frac{\|A^{-1}\mathbf{y}\|}{\|\mathbf{y}\|} \leq \|A^{-1}\|$$

Thus we have:

$$\frac{\|\delta_f\|}{\|A\mathbf{x}\|} \leq \frac{\|\delta A\| \|\mathbf{x}\|}{\|A\mathbf{x}\|} \leq \underbrace{\|A\| \|A^{-1}\|}_{\kappa(A)} \varepsilon.$$

■

### 3. Bounding floating point errors for linear algebra

We now observe that errors in implementing matrix-vector multiplication using floating points can be captured by considering the multiplication to be exact on the wrong matrix: that is, `A*x` (implemented with floating point as `mul_rows`) is precisely  $A + \delta A$  where  $\delta A$  has small norm, relative to  $A$ . That is, we have a bound on the *backward relative error*.

To discuss floating point errors we need to be precise which order the operations happened. We will use the definition `mul_rows(A, x)` (which is equivalent to `mul_cols(A, x)`). Note that each entry of the result is in fact a dot-product of the corresponding rows so we first consider the error in the dot product `dot(x, y)` as implemented in floating-point:

$$\text{dot}(\mathbf{x}, \mathbf{y}) = \bigoplus_{k=1}^n (x_k \otimes y_k).$$

We first need a helper proposition:

**Proposition 1 [PS2 Q2.1]** If  $|\epsilon_i| \leq \epsilon$  and  $n\epsilon < 1$ , then

$$\prod_{k=1}^n (1 + \epsilon_i) = 1 + \theta_n$$

for some constant  $\theta_n$  satisfying

$$|\theta_n| \leq \underbrace{\frac{n\epsilon}{1 - n\epsilon}}_{E_{n,\epsilon}}.$$

**Lemma 1 (dot product backward error)** For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$\text{dot}(\mathbf{x}, \mathbf{y}) = (\mathbf{x} + \delta\mathbf{x})^\top \mathbf{y}$$

where, assuming  $n\epsilon_m < 2$ , the entries satisfy

$$|\delta x_k| \leq E_{n,\epsilon_m/2} |x_k|.$$

### Proof

This is related to PS2 Q2.3 but asks for the *backward error* instead of the *forward error*. Note

$$\text{dot}(\mathbf{x}, \mathbf{y}) = \bigoplus_{j=1}^n (x_j \otimes y_j) = \bigoplus_{j=1}^n (x_j y_j)(1 + \delta_j) = x_1 y_1 (1 + \theta_n) + \sum_{j=2}^n x_j y_j (1 + \theta_{n-j+2})$$

where  $|\theta_n|, |\theta_k| \leq E_{n,\epsilon_m/2}$  (the subscript denotes the number of terms bounded by  $\epsilon_m/2$ ). Thus we can define

$$\delta \mathbf{x} := \begin{bmatrix} x_1 \theta_n \\ x_2 \theta_n \\ \vdots \\ x_n \theta_2 \end{bmatrix}$$

where

$$|\delta x_k| \leq E_{n,\epsilon_m/2} |x_k|.$$

■

We can use this to get a relative backward error bound on `mul_rows` :

**Theorem 2 (matrix-vector backward error)** For  $A \in \mathbb{R}^{m \times n}$  and  $\mathbf{x} \in \mathbb{R}^n$  (both with normal float entries) we have

$$\text{mul\_rows}(A, \mathbf{x}) = (A + \delta A)\mathbf{x}$$

where, assuming  $n\epsilon_m < 2$  and all operations are in the normalised range, the entries (denoting  $\delta a_{kj} = \delta A[k, j] = \mathbf{e}_k^\top \delta A \mathbf{e}_j$ ) satisfy

$$|\delta a_{kj}| \leq E_{n,\epsilon_m/2} |a_{kj}|.$$

**Proof** The bound on the entries of  $\delta A$  is implied by the previous lemma since each row is equivalent to a dot product. ■

### Corollary 1 (Norms)

$$\begin{aligned} \|\delta A\|_1 &\leq E_{n,\epsilon_m/2} \|A\|_1 \\ \|\delta A\|_2 &\leq \sqrt{\min(m, n)} E_{n,\epsilon_m/2} \|A\|_2 \\ \|\delta A\|_\infty &\leq E_{n,\epsilon_m/2} \|A\|_\infty \end{aligned}$$

In particular,

$$\text{mul\_rows}(A, \mathbf{x}) = A\mathbf{x} + \delta_f$$

where

$$\|\delta_f\| \leq \|A\mathbf{x}\| \kappa(A) E_{n, \epsilon_m/2}$$

### Proof

The 1-norm follows since

$$\|\delta A\|_1 = \max_j \sum_{k=1}^m |\delta a_{kj}| \leq E_{n, \epsilon_m/2} \max_j \sum_{k=1}^m |a_{kj}| = E_{n, \epsilon_m/2} \|A\|_1$$

and the proof for the  $\infty$ -norm is similar.

This leaves the 2-norm, which is a bit more challenging. We will prove the result by going through the Fröbenius norm and using

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{r} \|A\|_2$$

where  $r$  is rank of  $A$  (see PS6 Q5.2). So we deduce

$$\begin{aligned} \|\delta A\|_2^2 &\leq \|\delta A\|_F^2 = \sum_{k=1}^m \sum_{j=1}^n |\delta a_{kj}|^2 \leq E_{n, \epsilon_m/2}^2 \sum_{k=1}^m \sum_{j=1}^n |a_{kj}|^2 \\ &= E_{n, \epsilon_m/2}^2 \|A\|_F^2 \leq E_{n, \epsilon_m/2}^2 r \|A\|_2^2. \end{aligned}$$

and the rank of  $A$  is bounded by  $\min(m, n)$ . The bound on the forward error then follows from Theorem 1.

■

We can also bound the error of back-substitution in terms of the condition number (see PS7). If one uses QR to solve  $A\mathbf{x} = \mathbf{y}$  the condition number also gives a meaningful bound on the error. As we have already noted, there are some matrices where PLU decompositions introduce large errors, so in that case well-conditioning is not a guarantee of accuracy (but it still usually works).

## III.1 Fourier expansions

In Part III, Computing with Functions, we work with approximating functions by expansions in bases: that is, instead of approximating at a grid (as in the Differential Equations chapter), we approximate functions by other, simpler, functions. The ultimate goal is to use such expansions to numerically approximate solutions to ordinary and partial differential equations.

1. Review of Fourier series
2. Trapezium rule and discrete orthogonality
3. Convergence of approximate Fourier expansions

### 1. Review of Fourier series

The most fundamental basis is (complex) Fourier: we have  $e^{ik\theta}$  are orthogonal with respect to the inner product

$$\langle f, g \rangle := \frac{1}{2\pi} \int_0^{2\pi} f(\theta)g(\theta)d\theta,$$

where we conjugate the first argument to be consistent with the vector inner product  $\mathbf{x}^*\mathbf{y}$ . We can (typically) expand functions in this basis:

**Definition 1 (Fourier)** A function  $f$  has a Fourier expansion if

$$f(\theta) = \sum_{k=-\infty}^{\infty} f_k e^{ik\theta}$$

where

$$f_k := \langle e^{ik\theta}, f \rangle = \frac{1}{2\pi} \int_0^{2\pi} e^{-ik\theta} f(\theta) d\theta$$

A basic observation is if a Fourier expansion has no negative terms it is equivalent to a Taylor series if we write  $z = e^{i\theta}$ :

**Definition 2 (Fourier–Taylor)** A function  $f$  has a Fourier–Taylor expansion if

$$f(\theta) = \sum_{k=0}^{\infty} f_k e^{ik\theta}$$

where  $f_k := \langle e^{ik\theta}, f \rangle$ .

In numerical analysis we try to build on the analogy with linear algebra as much as possible. Therefore we can write this this as:

$$f(\theta) = \underbrace{[1|e^{i\theta}|e^{2i\theta}|\dots]}_{T(\theta)} \underbrace{\begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_k \end{bmatrix}}_{\mathbf{f}}.$$

Essentially, expansions in bases are viewed as a way of turning *functions* into (infinite) vectors. And (differential) operators into matrices.

In analysis one typically works with continuous functions and relates results to continuity. In numerical analysis we inherently have to work with *vectors*, so it is more natural to focus on the case where the *Fourier coefficients*  $f_k$  are *absolutely convergent*:

**Definition 3 (absolute convergent)** We write  $\mathbf{f} \in \ell^1$  if it is absolutely convergent, or in otherwords, the 1-norm of  $\mathbf{f}$  is bounded:

$$\|\mathbf{f}\|_1 := \sum_{k=-\infty}^{\infty} |f_k| < \infty$$

We first state a couple basic results (whose proof is beyond the scope of this module):

**Theorem 1 (2-norm convergence)** A function  $f$  has bounded 2-norm

$$\|f\|_2 := \sqrt{\langle f, f \rangle} = \sqrt{\int_0^{2\pi} |f(\theta)|^2 d\theta} < \infty$$

if and only if

$$\|\mathbf{f}\|_2 = \sqrt{\mathbf{f}^* \mathbf{f}} = \sqrt{\sum_{k=-\infty}^{\infty} |f_k|^2} < \infty.$$

In this case

$$\|f - \lim_{m \rightarrow \infty} \sum_{k=-m}^m f_k e^{ik\theta}\| \rightarrow 0$$

**Theorem 2 (Absolute convergence)** If  $\mathbf{f} \in \ell^1$  then

$$f(\theta) = \sum_{k=-\infty}^{\infty} f_k e^{ik\theta},$$

which converges uniformly.

Continuity gives us sufficient (though not necessary) conditions for absolute convergence:

**Lemma 1 (differentiability and absolutely convergence)** If  $f : \mathbb{R} \rightarrow \mathbb{C}$  and  $f'$  are periodic and  $f''$  is uniformly bounded, then  $\mathbf{f} \in \ell^1$ .

**Proof** Integrate by parts twice using the fact that  $f(0) = f(2\pi)$ ,  $f'(0) = f'(2\pi)$ :

$$\begin{aligned}
2\pi f_k &= \int_0^{2\pi} f(\theta) e^{-ik\theta} d\theta = [f(\theta) e^{-ik\theta}]_0^{2\pi} + \frac{1}{ik} \int_0^{2\pi} f'(\theta) e^{-ik\theta} d\theta \\
&= \frac{1}{ik} [f'(\theta) e^{-ik\theta}]_0^{2\pi} - \frac{1}{k^2} \int_0^{2\pi} f''(\theta) e^{-ik\theta} d\theta \\
&= -\frac{1}{k^2} \int_0^{2\pi} f''(\theta) e^{-ik\theta} d\theta
\end{aligned}$$

thus uniform boundedness of  $f''$  guarantees  $|f_k| \leq M|k|^{-2}$  for some  $M$ , and we have

$$\sum_{k=-\infty}^{\infty} |f_k| \leq |f_0| + 2M \sum_{k=1}^{\infty} |k|^{-2} < \infty.$$

using the dominant convergence test.

■

This condition can be weakened to Lipschitz continuity but the proof is beyond the scope of this module. Of more practical importance is the other direction: the more times differentiable a function the faster the coefficients decay, and thence the faster Fourier expansions converge. In fact, if a function is smooth and  $2\pi$ -periodic its Fourier coefficients decay faster than algebraically: they decay like  $O(k^{-\lambda})$  for any  $\lambda$ . This will be explored in the problem sheet.

**Remark (advanced)** Going further, if we let  $z = e^{i\theta}$  then if  $f(z)$  is *analytic* in a neighbourhood of the unit circle the Fourier coefficients decay *exponentially fast*. And if  $f(z)$  is entire they decay even faster than exponentially.

## 2. Trapezium rule and discrete Fourier coefficients

**Definition 4 (Trapezium Rule)** Let  $\theta_j = 2\pi j/n$  for  $j = 0, 1, \dots, n$  denote  $n+1$  evenly spaced points over  $[0, 2\pi]$ . The *Trapezium rule* over  $[0, 2\pi]$  is the approximation:

$$\int_0^{2\pi} f(\theta) d\theta \approx \frac{2\pi}{n} \left[ \frac{f(0)}{2} + \sum_{j=1}^{n-1} f(\theta_j) + \frac{f(2\pi)}{2} \right]$$

But if  $f$  is periodic we have  $f(0) = f(2\pi)$  we get the *periodic Trapezium rule*:

$$\int_0^{2\pi} f(\theta) d\theta \approx 2\pi \underbrace{\frac{1}{n} \sum_{j=0}^{n-1} f(\theta_j)}_{\Sigma_n[f]}$$

We know that  $e^{ik\theta}$  are orthogonal with respect to the continuous inner product. The following says that this property is maintained (up to "aliasing") when we replace the continuous integral with a trapezium rule approximation:

**Lemma 2 (Discrete orthogonality)** We have:

$$\sum_{j=0}^{n-1} e^{ik\theta_j} = \begin{cases} n & k = \dots, -2n, -n, 0, n, 2n, \dots \\ 0 & \text{otherwise} \end{cases}$$

In other words,

$$\Sigma_n[e^{i(k-\ell)\theta}] = \begin{cases} 1 & k - \ell = \dots, -2n, -n, 0, n, 2n, \dots \\ 0 & \text{otherwise} \end{cases}.$$

### Proof

Consider  $\omega := e^{i\theta_1} = e^{\frac{2\pi i}{n}}$ . This is an  $n$  th root of unity:  $\omega^n = 1$ . Note that  $e^{i\theta_j} = e^{\frac{2\pi ij}{n}} = \omega^j$ .

(Case 1:  $k = pn$  for an integer  $p$ ) We have

$$\sum_{j=0}^{n-1} e^{ik\theta_j} = \sum_{j=0}^{n-1} \omega^{kj} = \sum_{j=0}^{n-1} (\omega^{pn})^j = \sum_{j=0}^{n-1} 1 = n$$

(Case 2  $k \neq pn$  for an integer  $p$ ) Recall that

$$\sum_{j=0}^{n-1} z^j = \frac{z^n - 1}{z - 1}.$$

Then we have

$$\sum_{j=0}^{n-1} e^{ik\theta_j} = \sum_{j=0}^{n-1} (\omega^k)^j = \frac{\omega^{kn} - 1}{\omega^k - 1} = 0.$$

where we use the fact that  $k$  is not a multiple of  $n$  to guarantee that  $\omega^k \neq 1$ .

■

### 3. Convergence of Approximate Fourier expansions

We will now use the Trapezium rule to approximate Fourier coefficients and expansions:

**Definition 5 (Discrete Fourier coefficients)** Define the Trapezium rule approximation to the Fourier coefficients by:

$$f_k^n := \Sigma_n[e^{-ik\theta} f(\theta)] = \frac{1}{n} \sum_{j=0}^{n-1} e^{-ik\theta_j} f(\theta_j)$$

A remarkable fact is that the discrete Fourier coefficients can be expressed as a sum of the true Fourier coefficients:

**Theorem 3 (discrete Fourier coefficients)** If  $f \in \ell^1$  (absolutely convergent Fourier coefficients) then

$$f_k^n = \dots + f_{k-2n} + f_{k-n} + f_k + f_{k+n} + f_{k+2n} + \dots$$

### Proof

$$\begin{aligned}
f_k^n &= \sum_n [f(\theta) e^{-ik\theta}] = \sum_{\ell=-\infty}^{\infty} f_\ell \sum_n [e^{i(\ell-k)\theta}] \\
&= \sum_{\ell=-\infty}^{\infty} f_\ell \begin{cases} 1 & \ell - k = \dots, -2n, -n, 0, n, 2n, \dots \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

Note that there is redundancy:

**Corollary 1 (aliasing)** For all  $p \in \mathbb{Z}$ ,  $f_k^n = f_{k+pn}^n$ .

In other words if we know  $f_0^n, \dots, f_{n-1}^n$ , we know  $f_k^n$  for all  $k$  via a permutation, for example if  $n = 2m + 1$  we have

$$\begin{bmatrix} f_{-m}^n \\ \vdots \\ f_m^n \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 1 & & & & & \\ & & & \ddots & & & & \\ & & & & 1 & & & \\ & & & & & \ddots & & \\ & & & & & & 1 & \\ & & & & & & & \end{bmatrix}}_{P_\sigma} \begin{bmatrix} f_0^n \\ \vdots \\ f_{n-1}^n \end{bmatrix}$$

where  $\sigma$  has Cauchy notation (*Careful*: we are using 1-based indexing here):

$$\begin{pmatrix} 1 & 2 & \cdots & m & m+1 & m+2 & \cdots & n \\ m+2 & m+3 & \cdots & n & 1 & 2 & \cdots & m+1 \end{pmatrix}.$$

We first discuss the case when all negative coefficients are zero. That is,  $f_0^n, \dots, f_{n-1}^n$  are approximations of the Fourier–Taylor coefficients by evaluating on the boundary.

We can prove convergence whenever of this approximation whenever  $f$  has absolutely summable coefficients. We will prove the result here in the special case where the negative coefficients are zero.

**Theorem 4 (Approximate Fourier–Taylor expansions converge)** If  $0 = f_{-1} = f_{-2} = \dots$  and  $\mathbf{f}$  is absolutely convergent then

$$f_n(\theta) = \sum_{k=0}^{n-1} f_k^n e^{ik\theta}$$

converges uniformly to  $f(\theta)$ .

### Proof

$$\begin{aligned}
|f(\theta) - f_n(\theta)| &= \left| \sum_{k=0}^{n-1} (f_k - f_k^n) e^{ik\theta} + \sum_{k=n}^{\infty} f_k e^{ik\theta} \right| = \left| \sum_{k=n}^{\infty} f_k (e^{ik\theta} - e^{i\text{mod}(k,n)}) \right| \\
&\leq 2 \sum_{k=n}^{\infty} |f_k|
\end{aligned}$$

which goes to zero as  $n \rightarrow \infty$ . ■

For the general case we need to choose a range of coefficients that includes roughly an equal number of negative and positive coefficients (preferring negative over positive in a tie as a convention):

$$f_n(\theta) = \sum_{k=-\lceil n/2 \rceil}^{\lfloor n/2 \rfloor} f_k e^{ik\theta}$$

In the problem sheet we will prove this converges provided the coefficients are absolutely convergent.

## III.2 Discrete Fourier Transform

In the last lecture we explored using the trapezium rule for approximating Fourier coefficients. This is a linear map from function values to coefficients and thus can be reinterpreted as a matrix-vector product, called the the Discrete Fourier Transform. It turns out the matrix is unitary which leads to important properties including interpolation. Finally, we discuss how a clever way of decomposing the DFT leads to a fast way of applying and inverting it, which is one of the most influencial algorithms of the 20th century: the Fast Fourier Transform.

1. The Discrete Fourier Transform (DFT): We discuss the map from values to approximate Fourier coefficients, and back.
2. Interpolation: We show that the approximate Fourier expansion *exactly* interpolates the values at the sample grid.
3. The Fast Fourier Transform (FFT): We discuss how the DFT can be applied in  $O(n \log n)$  operations.

### 1. The Discrete Fourier transform

**Definition 1 (DFT)** The *Discrete Fourier Transform (DFT)* is defined as:

$$Q_n := \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & e^{-i\theta_1} & e^{-i\theta_2} & \cdots & e^{-i\theta_{n-1}} \\ 1 & e^{-i2\theta_1} & e^{-i2\theta_2} & \cdots & e^{-i2\theta_{n-1}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & e^{-i(n-1)\theta_1} & e^{-i(n-1)\theta_2} & \cdots & e^{-i(n-1)\theta_{n-1}} \end{bmatrix}$$

$$= \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega^{-1} & \omega^{-2} & \cdots & \omega^{-(n-1)} \\ 1 & \omega^{-2} & \omega^{-4} & \cdots & \omega^{-2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{-(n-1)} & \omega^{-2(n-1)} & \cdots & \omega^{-(n-1)^2} \end{bmatrix}$$

for the  $n$ -th root of unity  $\omega = e^{i\pi/n}$ . Note that

$$Q_n^* = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & e^{i\theta_1} & e^{i2\theta_1} & \cdots & e^{i(n-1)\theta_1} \\ 1 & e^{i\theta_2} & e^{i2\theta_2} & \cdots & e^{i(n-1)\theta_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & e^{i\theta_{n-1}} & e^{i2\theta_{n-1}} & \cdots & e^{i(n-1)\theta_{n-1}} \end{bmatrix}$$

$$= \frac{1}{\sqrt{n}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega^1 & \omega^2 & \cdots & \omega^{(n-1)} \\ 1 & \omega^2 & \omega^4 & \cdots & \omega^{2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{(n-1)} & \omega^{2(n-1)} & \cdots & \omega^{(n-1)^2} \end{bmatrix}$$

Note that

$$\underbrace{\begin{bmatrix} f_0^n \\ \vdots \\ f_{n-1}^n \end{bmatrix}}_{\mathbf{f}^n} = \frac{1}{\sqrt{n}} Q_n \underbrace{\begin{bmatrix} f(\theta_0) \\ \vdots \\ f(\theta_n) \end{bmatrix}}_{\mathbf{f}}$$

The choice of normalisation constant is motivated by the following:

**Proposition 1 (DFT is Unitary)**  $Q_n \in U(n)$ , that is,  $Q_n^* Q_n = Q_n Q_n^* = I$ .

**Proof**

$$Q_n Q_n^* = \begin{bmatrix} \Sigma_n[1] & \Sigma_n[e^{i\theta}] & \cdots & \Sigma_n[e^{i(n-1)\theta}] \\ \Sigma_n[e^{-i\theta}] & \Sigma_n[1] & \cdots & \Sigma_n[e^{i(n-2)\theta}] \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_n[e^{-i(n-1)\theta}] & \Sigma_n[e^{-i(n-2)\theta}] & \cdots & \Sigma_n[1] \end{bmatrix} = I$$

■

In other words,  $Q_n$  is easily inverted and we also have a map from discrete Fourier coefficients back to values:

$$\sqrt{n} Q_n^* \mathbf{f}^n = \mathbf{f}$$

## 2. Interpolation

We investigated briefly interpolation and least squares using polynomials at evenly spaced points, observing that there were issues with stability. We now show that the DFT actually gives coefficients that interpolate using Fourier expansions. As the DFT is a

unitary matrix its (2-norm) condition number is 1, hence this is a stable process. Thus we arrive at the main result:

### Corollary 1 (Interpolation)

$$f_n(\theta) := \sum_{k=0}^{n-1} f_k^n e^{ik\theta}$$

interpolates  $f$  at  $\theta_j$ :

$$f_n(\theta_j) = f(\theta_j)$$

**Proof** We have

$$f_n(\theta_j) = \sum_{k=0}^{n-1} f_k^n e^{ik\theta_j} = \sqrt{n} \mathbf{e}_j^\top Q_n^* \mathbf{f}^n = \mathbf{e}_j^\top Q_n^* Q_n \mathbf{f}^n = f(\theta_j).$$

■

We will leave extending this result to the general non-Taylor case to the problem sheet. Note that regardless of choice of coefficients we interpolate provided we have  $n$  consecutive coefficients, though some interpolations are better than others:

In [1]:

```
using Plots, LinearAlgebra

# evaluates f_n at a point
function finitefourier(f_n, θ)
    m = n ÷ 2 # use coefficients between -m:m
    ret = 0.0 + 0.0im # coefficients are complex so we need complex arithmetic
    for k = 0:m
        ret += f_n[k+1] * exp(im*k*θ)
    end
    for k = -m:-1
        ret += f_n[end+k+1] * exp(im*k*θ)
    end
    ret
end

function finitetaylor(f_n, θ)
    ret = 0.0 + 0.0im # coefficients are complex so we need complex arithmetic
    for k = 0:n-1
        ret += f_n[k+1] * exp(im*k*θ)
    end
    ret
end

f = θ -> exp(cos(θ))
n = 7
θ = range(0, 2π; length=n+1)[1:end-1] # θ_0, ..., θ_{n-1}, dropping θ_n == 2π
Q_n = 1/sqrt(n) * [exp(-im*(k-1)*θ[j]) for k = 1:n, j=1:n]
f_n = 1/sqrt(n) * Q_n * f.(θ)
```

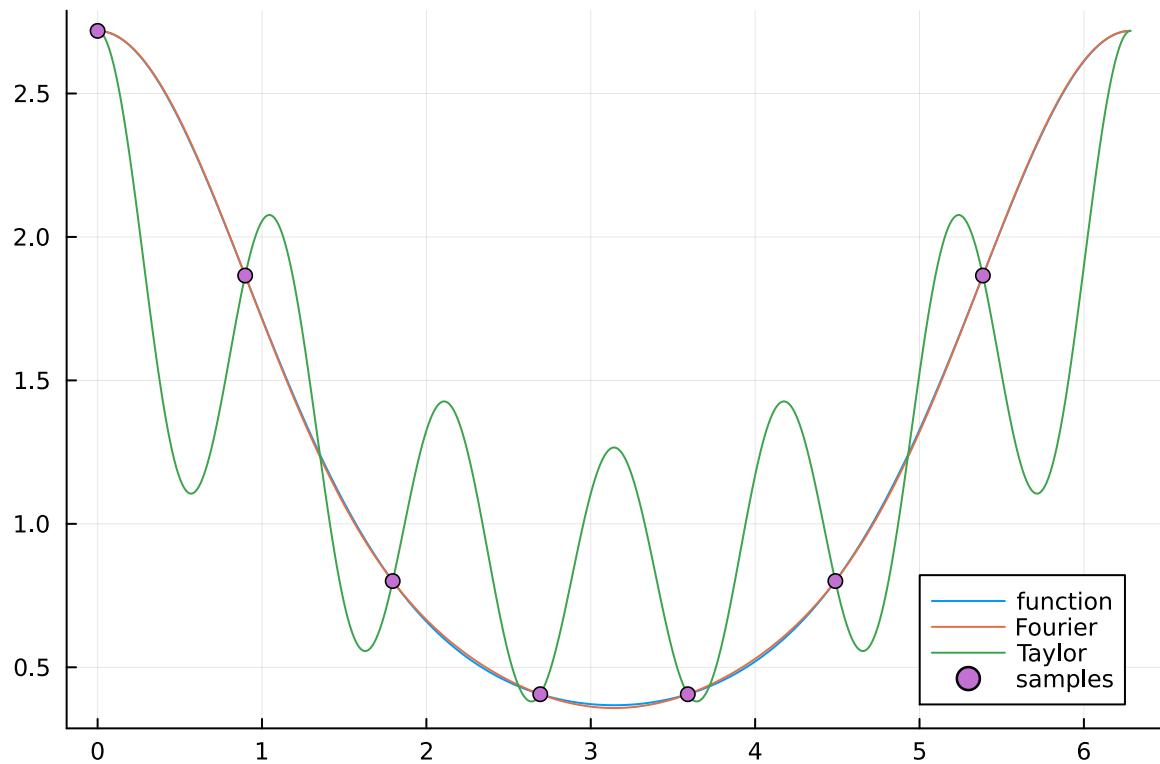
```

f_n = θ -> finitefourier(f_n, θ)
t_n = θ -> finitetaylor(f_n, θ)

g = range(0, 2π; length=1000) # plotting grid
plot(g, f.(g); label="function", legend=:bottomright)
plot!(g, real.(f_n.(g)); label="Fourier")
plot!(g, real.(t_n.(g)); label="Taylor")
scatter!(θ, f.(θ); label="samples")

```

Out[1]:



We now demonstrate the relationship of Taylor and Fourier coefficients and their discrete approximations for some examples:

**Example 1** Consider the function

$$f(\theta) = \frac{2}{2 - e^{i\theta}}$$

Under the change of variables  $z = e^{i\theta}$  we know for  $z$  on the unit circle this becomes (using the geometric series with  $z/2$ )

$$\frac{2}{2 - z} = \sum_{k=0}^{\infty} \frac{z^k}{2^k}$$

i.e.,  $f_k = 1/2^k$  which is absolutely summable:

$$\sum_{k=0}^{\infty} |f_k| = f(0) = 2.$$

If we use an  $n$  point discretisation we get (using the geoemtric series with  $2^{-n}$ )

$$f_k^n = f_k + f_{k+n} + f_{k+2n} + \dots = \sum_{p=0}^{\infty} \frac{1}{2^{k+pn}} = \frac{2^{n-k}}{2^n - 1}$$

We can verify this numerically:

```
In [2]: f = θ -> 2/(2 - exp(im*θ))
n = 7
θ = range(0, 2π; length=n+1)[1:end-1] # θ_0, ..., θ_{n-1}, dropping θ_n == 2π
Q_n = 1/sqrt(n) * [exp(-im*(k-1)*θ[j]) for k = 1:n, j=1:n]

Q_n/sqrt(n)*f.(θ) ≈ 2 .^ (n .- (0:n-1)) / (2^(n-1))
```

Out[2]: true

**Example 2** Define the following infinite sum (which has no name apparently, according to Mathematica):

$$S_n(k) := \sum_{p=0}^{\infty} \frac{1}{(k+pn)!}$$

We can use the DFT to compute  $S_n(0), \dots, S_n(n-1)$ . Consider

$$f(\theta) = \exp(e^{i\theta}) = \sum_{k=0}^{\infty} \frac{e^{ik\theta}}{k!}$$

where we know the Fourier coefficients from the Taylor series of  $e^z$ . The discrete Fourier coefficients satisfy for  $0 \leq k \leq n-1$ :

$$f_k^n = f_k + f_{k+n} + f_{k+2n} + \dots = S_n(k)$$

Thus we have

$$\begin{bmatrix} S_n(0) \\ \vdots \\ S_n(n-1) \end{bmatrix} = \frac{1}{\sqrt{n}} Q_n \begin{bmatrix} 1 \\ \exp(e^{2i\pi/n}) \\ \vdots \\ \exp(e^{2i(n-1)\pi/n}) \end{bmatrix}$$

### 3. Fast Fourier Transform (non-examinable)

Applying  $Q_n$  or its adjoint  $Q_n^*$  to a vector naively takes  $O(n^2)$  operations. Both can be reduced to  $O(n \log n)$  using the celebrated *Fast Fourier Transform* (FFT), which is one of the [Top 10 Algorithms of the 20th Century](#) (You won't believe number 7!).

The key observation is that hidden in  $Q_{2n}$  are 2 copies of  $Q_n$ . We will work with multiple  $n$  we denote the  $n$ -th root as  $\omega_n = \exp(2\pi/n)$ . Note that we can relate a vector of powers of  $\omega_{2n}$  to two copies of vectors of powers of  $\omega_n$ :

$$\underbrace{\begin{bmatrix} 1 \\ \omega_{2n} \\ \vdots \\ \omega_{2n}^{2n-1} \end{bmatrix}}_{\vec{\omega}_{2n}} = P_\sigma^\top \begin{bmatrix} I_n \\ \omega_{2n} I_n \end{bmatrix} \underbrace{\begin{bmatrix} 1 \\ \omega_n \\ \vdots \\ \omega_n^{n-1} \end{bmatrix}}_{\vec{\omega}_n}$$

where  $\sigma$  has the Cauchy notation

$$\begin{pmatrix} 1 & 2 & 3 & \cdots & n & n+1 & \cdots & 2n \\ 1 & 3 & 5 & \cdots & 2n-1 & 2 & \cdots & 2n \end{pmatrix}$$

That is,  $P_\sigma$  is the following matrix which takes the even entries and places them in the first  $n$  entries and the odd entries in the last  $n$  entries:

```
In [3]: n = 4
σ = [1:2:2n-1; 2:2:2n]
P_σ = I(2n)[σ,:]
```

```
Out[3]: 8×8 Matrix{Bool}:
1 0 0 0 0 0 0 0
0 0 1 0 0 0 0 0
0 0 0 0 1 0 0 0
0 0 0 0 0 0 1 0
0 1 0 0 0 0 0 0
0 0 0 1 0 0 0 0
0 0 0 0 0 1 0 0
0 0 0 0 0 0 0 1
```

and so  $P_\sigma^\top$  reverses the process. Thus we have

$$\begin{aligned} Q_{2n}^* &= \frac{1}{\sqrt{2n}} [\mathbf{1}_{2n} | \vec{\omega}_{2n} | \vec{\omega}_{2n}^2 | \cdots | \vec{\omega}_{2n}^{2n-1}] \\ &= \frac{1}{\sqrt{2n}} P_\sigma^\top \begin{bmatrix} \mathbf{1}_n & \vec{\omega}_n & \vec{\omega}_n^2 & \cdots & \vec{\omega}_n^{n-1} & \vec{\omega}_n^n & \cdots & \vec{\omega}_n^{2n-1} \\ \mathbf{1}_n & \omega_{2n} \vec{\omega}_n & \omega_{2n}^2 \vec{\omega}_n^2 & \cdots & \omega_{2n}^{n-1} \vec{\omega}_n^{n-1} & \omega_{2n}^n \vec{\omega}_n^n & \cdots & \omega_{2n}^{2n-1} \vec{\omega}_n^{2n-1} \end{bmatrix} \\ &= \frac{1}{\sqrt{2}} P_\sigma^\top \begin{bmatrix} Q_n^* & Q_n^* \\ Q_n^* D_n & -Q_n^* D_n \end{bmatrix} = \frac{1}{\sqrt{2}} P_\sigma^\top \begin{bmatrix} Q_n^* & \\ & Q_n^* \end{bmatrix} \begin{bmatrix} I_n & I_n \\ D_n & -D_n \end{bmatrix} \end{aligned}$$

In other words, we reduced the DFT to two DFTs applied to vectors of half the dimension.

We can see this formula in code:

```
In [4]: function fftmatrix(n)
    θ = range(0, 2π; length=n+1)[1:end-1] # θ_0, ..., θ_{n-1}, dropping θ_n == 2
    [exp(-im*(k-1)*θ[j]) for k = 1:n, j=1:n]/sqrt(n)
end

Q_{2n} = fftmatrix(2n)
Q_n = fftmatrix(n)
```

```
D_n = Diagonal([exp(im*k*pi/n) for k=0:n-1])
(P_sigma'*[Q_n' Q_n'; Q_n'*D_n -Q_n'*D_n])[1:n,1:n] ≈ sqrt(2)Q_2n'[1:n,1:n]
```

Out [4]: true

Now assume  $n = 2^q$  so that  $\log_2 n = q$ . To see that we get  $O(n \log n) = O(nq)$  operations we need to count the operations. Assume that applying  $F_n$  takes  $\leq 3nq$  additions and multiplications. The first  $n$  rows takes  $n$  additions. The last  $n$  has  $n$  multiplications and  $n$  additions. Thus we have

$6nq + 3n \leq 6n(q+1) = 3(2n) \log_2(2n)$  additions/multiplications, showing by induction that we have  $O(n \log n)$  operations.

**Remark** The FFTW.jl package wraps the FFTW (Fastest Fourier Transform in the West) library, which is a highly optimised implementation of the FFT that also works well even when  $n$  is not a power of 2. (As an aside, the creator of FFTW [Steven Johnson](#) is now a Julia contributor and user.) Here we approximate  $\exp(\cos(\theta - 0.1))$  using 31 nodes:

In [5]:

```
using FFTW
f = θ -> exp(cos(θ-0.1))
n = 31
m = n÷2
# evenly spaced points from 0:2π, dropping last node
θ = range(0, 2π; length=n+1)[1:end-1]

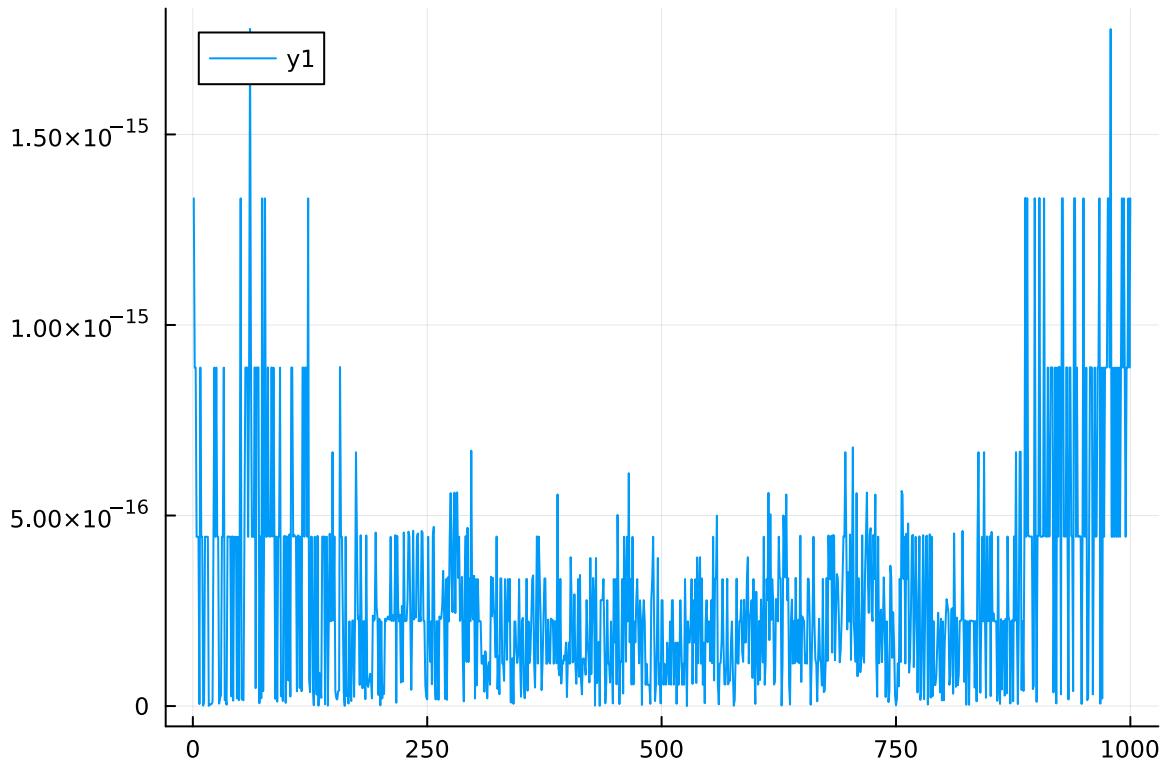
# fft returns discrete Fourier coefficients n*[f̂_0, ..., f̂_(n-1)]
fc = fft(f.(θ))/n

# We reorder using [f̂_(-m), ..., f̂_(-1)] == [f̂_(n-m), ..., f̂_(n-1)]
# == [f̂_(m+1), ..., f̂_(n-1)]
f̂ = [fc[m+2:end]; fc[1:m+1]]

# equivalent to f̂_(-m)*exp(-im*m*θ) + ... + f̂_(m)*exp(im*m*θ)
f_n = θ -> transpose([exp(im*k*θ) for k=-m:m]) * f̂

# plotting grid
g = range(0, 2π; length=1000)
plot(abs.(f_n.(g) - f.(g)))
```

Out[5]:



Thus we have successfully approximate the function to roughly machine precision. The magic of the FFT is because it's  $O(n \log n)$  we can scale it to very high orders. Here we plot the Fourier coefficients for a function that requires around 100k coefficients to resolve:

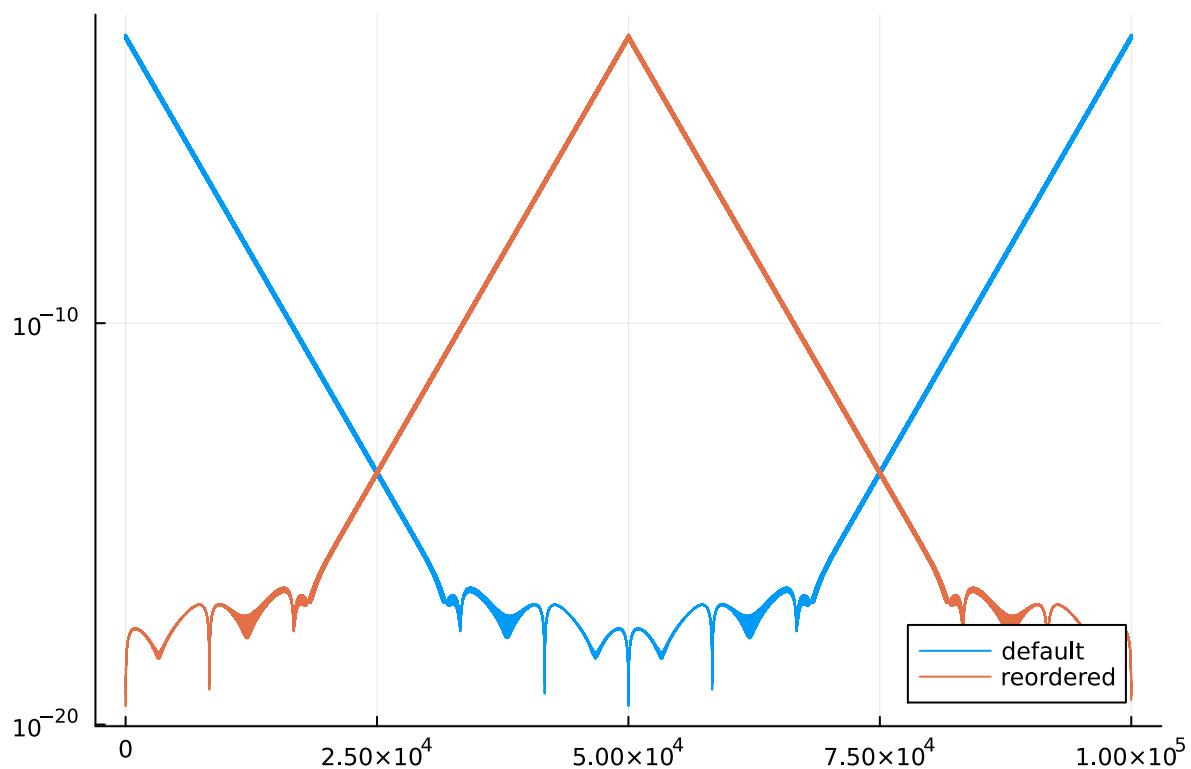
```
In [6]: f = θ -> exp(sin(θ))/(1+1e6cos(θ)^2)
n = 100_001
m = n÷2
# evenly spaced points from 0:2π, dropping last node
θ = range(0, 2π; length=n+1)[1:end-1]

# fft returns discrete Fourier coefficients n*[f^n_0, ..., f^n_(n-1)]
fc = fft(f.(θ))/n

# We reorder using [f^n_(-m), ..., f^n_(-1)] == [f^n_(n-m), ..., f^n_(n-1)]
# == [f^n_(m+1), ..., f^n_(n-1)]
f̂ = [fc[m+2:end]; fc[1:m+1]]

plot(abs.(fc);yscale=:log10, legend=:bottomright, label="default")
plot!(abs.(f̂);yscale=:log10, label="reordered")
```

Out[6]:



## III.3 Orthogonal Polynomials

Fourier series proved very powerful for approximating periodic functions. If periodicity is lost, however, uniform convergence is lost. In this chapter we introduce alternative bases, *Orthogonal Polynomials (OPs)* built on polynomials that are applicable in the non-periodic setting. That is we consider expansions of the form

$$f(x) = \sum_{k=0}^{\infty} c_k p_k(x)$$

where  $p_k(x)$  are special families of polynomials,  $c_k$  are expansion coefficients. The computation of the coefficients  $c_k \approx c_k^n$  will wait until the next section.

Why not use monomials as in Taylor series? Hidden in the previous lecture was that we could effectively compute Taylor coefficients by evaluating on the unit circle in the complex plane, *only* if the radius of convergence was 1. Many functions are smooth on say  $[-1, 1]$  but have non-convergent Taylor series, e.g.:

$$\frac{1}{25x^2 + 1}$$

While orthogonal polynomials span the same space as monomials, and therefore we can in theory write an approximation in monomials, orthogonal polynomials are *much* more stable.

In addition to numerics, OPs play a very important role in many mathematical areas including functional analysis, integrable systems, singular integral equations, complex analysis, and random matrix theory.

1. General properties: we define orthogonal polynomials
2. Three-term recurrences
3. Jacobi matrices

**Remark (advanced)** Going beyond what we discuss here are many other beautiful properties of orthogonal polynomials that are useful in computation, analysis, representation theory, quantum mechanics, etc. These include sparse recurrence relationships for derivatives, that they are eigenfunctions of differential equations, their asymptotics, generating functions, etc.

### 1. General properties

**Definition 1 (graded polynomial basis)** A set of polynomials  $\{p_0(x), p_1(x), \dots\}$  is *graded* if  $p_n$  is precisely degree  $n$ : i.e.,

$$p_n(x) = k_n x^n + k_n^{(1)} x^{n-1} + \cdots + k_n^{(n-1)} x + k_n^{(n)}$$

for  $k_n \neq 0$ .

Note that if  $p_n$  are graded then  $\{p_0(x), \dots, p_n(x)\}$  are a basis of all polynomials of degree  $n$ .

**Definition 2 (orthogonal polynomials)** Given an (integrable) weight  $w(x) > 0$  for  $x \in (a, b)$ , which defines a continuous inner product

$$\langle f, g \rangle = \int_a^b f(x)g(x)w(x)dx$$

a graded polynomial basis  $\{p_0(x), p_1(x), \dots\}$  are *orthogonal polynomials (OPs)* if

$$\langle p_n, p_m \rangle = 0$$

whenever  $n \neq m$ . We assume through that integrals of polynomials are finite:

$$\int_a^b x^k w(x)dx < \infty.$$

Note in the above

$$h_n := \langle p_n, p_n \rangle = \|p_n\|^2 = \int_a^b p_n(x)^2 w(x)dx > 0.$$

Multiplying any orthogonal polynomial by a nonzero constant necessarily is also an orthogonal polynomial. We have two standard normalisations:

**Definition 3 (orthonormal polynomials)** A set of orthogonal polynomials  $\{q_0(x), q_1(x), \dots\}$  are *orthonormal* if  $\|q_n\| = 1$ .

**Definition 4 (monic orthogonal polynomials)** A set of orthogonal polynomials  $\{\pi_0(x), \pi_1(x), \dots\}$  are *monic* if  $k_n = 1$ .

**Proposition 1 (existence)** Given a weight  $w(x)$ , monic orthogonal polynomials exist.

### Proof

Existence follows immediately from the Gram–Schmidt procedure. That is, define  $\pi_0(x) := 1$  and

$$\pi_n(x) := x^n - \sum_{k=0}^{n-1} \frac{\langle x^n, \pi_k \rangle}{\|\pi_k\|^2} \pi_k(x)$$

■

We are primarily concerned with the usage of orthogonal polynomials in approximating functions. First we observe the following:

**Proposition 2 (expansion)** If  $r(x)$  is a degree  $n$  polynomial and  $\{p_n\}$  are orthogonal then

$$r(x) = \sum_{k=0}^n \frac{\langle p_k, r \rangle}{\|p_k\|^2} p_k(x)$$

Note for  $\{q_n\}$  orthonormal we have

$$r(x) = \sum_{k=0}^n \langle q_k, r \rangle q_k(x).$$

**Proof** Because  $\{p_0, \dots, p_n\}$  are a basis of polynomials we can write

$$r(x) = \sum_{k=0}^n r_k p_k(x)$$

for constants  $r_k \in \mathbb{R}$ . By linearity we have

$$\langle p_m, r \rangle = \sum_{k=0}^n r_k \langle p_m, p_k \rangle = r_m \langle p_m, p_m \rangle$$

■

**Corollary 1 (zero inner product)** If a degree  $n$  polynomial  $r$  satisfies

$$0 = \langle p_0, r \rangle = \dots = \langle p_n, r \rangle$$

then  $r = 0$ .

**Corollary 2 (uniqueness)** Monic orthogonal polynomials are unique.

**Proof** If  $p_n(x)$  and  $\pi_n(x)$  are both monic orthogonal polynomials then  $r(x) = p_n(x) - \pi_n(x)$  is degree  $n - 1$  but satisfies

$$\langle r, \pi_k \rangle = \langle p_n, \pi_k \rangle - \langle \pi_n, \pi_k \rangle = 0$$

for  $k = 0, \dots, n - 1$ . Note  $\langle p_n, \pi_k \rangle = 0$  can be seen by expanding

$$\pi_k(x) = \sum_{j=0}^k c_j p_j(x).$$

■

OPs are uniquely defined (up to a constant) by the property that they are orthogonal to all lower degree polynomials.

**Theorem 1 (orthogonal to lower degree)** Given a weight  $w(x)$ , a polynomial

$$p(x) = k_n x^n + O(x^{n-1})$$

with  $k_n \neq 0$  satisfies

$$\langle p, f_m \rangle = 0$$

for all polynomials  $f_m$  of degree  $m < n$  if and only if  $p(x) = k_n \pi_n(x)$  where  $\pi_n(x)$  are the monic orthogonal polynomials. Therefore an orthogonal polynomial is uniquely defined by the weight and leading order coefficient  $k_n$ .

We leave this proof to the problem sheets. A consequence of this is that orthonormal polynomials are always a constant multiple of orthogonal polynomials.

## 2. 3-term recurrence

The most *fundamental* property of orthogonal polynomials is their three-term recurrence.

**Theorem 2 (3-term recurrence, 2nd form)** If  $\{p_n\}$  are OPs then there exist real constants  $a_n, b_n \neq 0, c_{n-1} \neq 0$  such that

$$\begin{aligned} xp_0(x) &= a_0 p_0(x) + b_0 p_1(x) \\ xp_n(x) &= c_{n-1} p_{n-1}(x) + a_n p_n(x) + b_n p_{n+1}(x) \end{aligned}$$

**Proof** The  $n = 0$  case is immediate since  $\{p_0, p_1\}$  are a basis of degree 1 polynomials. The  $n > 0$  case follows from

$$\langle xp_n, p_k \rangle = \langle p_n, xp_k \rangle = 0$$

for  $k < n - 1$  as  $xp_k$  is of degree  $k + 1 < n$ .

Note that

$$b_n = \frac{\langle p_{n+1}, xp_n \rangle}{\|p_{n+1}\|^2} \neq 0$$

since  $xp_n = k_n x^{n+1} + O(x^n)$  is precisely degree  $n$ . Further,

$$c_{n-1} = \frac{\langle p_{n-1}, xp_n \rangle}{\|p_{n-1}\|^2} = \frac{\langle p_n, xp_{n-1} \rangle}{\|p_{n-1}\|^2} = b_{n-1} \frac{\|p_n\|^2}{\|p_{n-1}\|^2} \neq 0.$$

■

Clearly if  $\pi_n$  is monic then so is  $x\pi_n$  which leads to the following:

**Corollary 3 (monic 3-term recurrence)** If  $\{\pi_n\}$  are monic then  $b_n = 1$ .

**Corollary 4 (Gram–Schmidt revisited)** We can define  $\pi_{n+1}(x)$  in terms of  $\pi_{n-1}$  and  $\pi_n$ :

$$\pi_{n+1}(x) = x\pi_n(x) - a_n \pi_n(x) - c_{n-1} \pi_{n-1}(x)$$

where

$$a_n = \frac{\langle x\pi_n, \pi_n \rangle}{\|\pi_n\|^2} \quad \text{and} \quad c_{n-1} = \frac{\langle x\pi_n, \pi_{n-1} \rangle}{\|\pi_{n-1}\|^2}$$

**Example 1** What are the monic OPs  $\pi_0(x), \dots, \pi_3(x)$  with respect to  $w(x) = 1$  on  $[0, 1]$ ? We can construct these using Gram–Schmidt, but exploiting the 3-term recurrence to reduce the computational cost. We have  $\pi_0(x) = 1$ , which we see is orthogonal:

$$\|\pi_0\|^2 = \langle \pi_0, \pi_0 \rangle = \int_0^1 dx = 1.$$

We know from the 3-term recurrence that

$$x\pi_0(x) = a_0\pi_0(x) + \pi_1(x)$$

where

$$a_0 = \frac{\langle \pi_0, x\pi_0 \rangle}{\|\pi_0\|^2} = \int_0^1 x dx = 1/2.$$

Thus

$$\begin{aligned} \pi_1(x) &= x\pi_0(x) - a_0\pi_0(x) = x - 1/2 \\ \|\pi_1\|^2 &= \int_0^1 (x^2 - x + 1/4) dx = 1/12 \end{aligned}$$

From

$$x\pi_1(x) = c_0\pi_0(x) + a_1\pi_1(x) + \pi_2(x)$$

we have

$$\begin{aligned} c_0 &= \frac{\langle \pi_0, x\pi_1 \rangle}{\|\pi_0\|^2} = \int_0^1 (x^2 - x/2) dx = 1/12 \\ a_1 &= \frac{\langle \pi_1, x\pi_1 \rangle}{\|\pi_1\|^2} = 12 \int_0^1 (x^3 - x^2 + x/4) dx = 1/2 \\ \pi_2(x) &= x\pi_1(x) - c_0 - a_1\pi_1(x) = x^2 - x + 1/6 \\ \|\pi_2\|^2 &= \int_0^1 (x^4 - 2x^3 + 4x^2/3 - x/3 + 1/36) dx = \frac{1}{180} \end{aligned}$$

Finally, from

$$x\pi_2(x) = c_1\pi_1(x) + a_2\pi_2(x) + \pi_3(x)$$

we have

$$\begin{aligned}
c_1 &= \frac{\langle \pi_1, x\pi_2 \rangle}{\|\pi_1\|^2} = 12 \int_0^1 (x^4 - 3x^3/2 + 2x^2/3 - x/12) dx = 1/15 \\
a_2 &= \frac{\langle \pi_2, x\pi_2 \rangle}{\|\pi_2\|^2} = 180 \int_0^1 (x^5 - 2x^4 + 4x^3/3 - x^2/3 + x/36) dx = 1/2 \\
\pi_3(x) &= x\pi_2(x) - c_1\pi_1(x) - a_2\pi_2(x) = x^3 - x^2 + x/6 - x/15 + 1/30 - x^2/2 + x/2 \\
&= x^3 - 3x^2/2 + 3x/5 - 1/20
\end{aligned}$$

### 3. Jacobi matrices

The three-term recurrence can also be interpreted as a matrix known as the Jacobi matrix:

**Corollary 5 (Jacobi matrix)** For

$$P(x) := [p_0(x) | p_1(x) | \dots]$$

then we have

$$xP(x) = P(x) \underbrace{\begin{bmatrix} a_0 & c_0 & & \\ b_0 & a_1 & c_1 & \\ & b_1 & a_2 & \ddots \\ & \ddots & \ddots & \ddots \end{bmatrix}}_X$$

More generally, for any polynomial  $a(x)$  we have

$$a(x)P(x) = P(x)a(X).$$

The Jacobi matrix is  $J := X^\top$ .

For the special cases of orthonormal and monic polynomials we have extra structure:

**Corollary 6 (orthonormal 3-term recurrence)** The Jacobi matrix of a family of orthogonal polynomials  $p_n(x)$  is symmetric:

$$J = X = \begin{bmatrix} a_0 & b_0 & & \\ b_0 & a_1 & b_1 & \\ & b_1 & a_2 & \ddots \\ & \ddots & \ddots & \ddots \end{bmatrix}$$

if and only if  $p_n(x)$  is up-to-sign a fixed constant scaling of orthonormal: for  $q_n(x) := \pi_n(x)/\|\pi_n\|^2$  we have for a fixed  $\alpha \in \mathbb{R}$  and  $s_n \in \{-1, 1\}$

$$p_n(x) = \alpha s_n q_n(x).$$

**Proof** Noting that  $\|q_n\|^2 = 1$  and thence  $\|p_n\|^2 = \alpha^2$ , if  $p_n(x) = \alpha s_n q_n(x)$  we have

$$b_n = \frac{\langle xp_n, p_{n+1} \rangle}{\|p_{n+1}\|^2} = s_n s_{n+1} \langle xq_n, q_{n+1} \rangle = s_n s_{n+1} \langle q_n, xq_{n+1} \rangle = \frac{\langle p_n, xp_{n+1} \rangle}{\|p_n\|^2} = c_{n-1}.$$

Conversely, suppose  $X = X^\top$ , i.e.,  $b_n = c_{n-1}$  and write the corresponding orthogonal polynomials as  $p_n(x) = \alpha_n q_n(x)$ . We have

$$\begin{aligned} b_n &= \frac{\langle xp_n, p_{n+1} \rangle}{\|p_{n+1}\|^2} = \frac{\alpha_n}{\alpha_{n+1}} \langle xq_n, q_{n+1} \rangle = \frac{\alpha_n}{\alpha_{n+1}} \langle q_n, xq_{n+1} \rangle = \frac{\alpha_n^2}{\alpha_{n+1}^2} \frac{\langle p_n, xp_{n+1} \rangle}{\|p_n\|^2} = \frac{\alpha_n^2}{\alpha_{n+1}^2} \\ &= \frac{\alpha_n^2}{\alpha_{n+1}^2} b_n. \end{aligned}$$

Hence  $\alpha_n^2 = \alpha_{n+1}^2$  which implies that  $\alpha_{n+1} = \pm \alpha_n$ . By induction the result follows, where  $\alpha := \alpha_0$ . ■

**Remark (advanced)** If you are worried about multiplication of infinite matrices/vectors note it is well-defined by the standard definition because it is banded. It can also be defined in terms of functional analysis where one considers these as linear operators (functions of functions) between vector spaces.

**Remark (advanced)** Every integrable weight generates a family of orthonormal polynomials, which in turn generates a symmetric Jacobi matrix. There is a "Spectral Theorem for Jacobi matrices" that says one can go the other way: every tridiagonal symmetric matrix with bounded entries is a Jacobi matrix for some integrable weight with compact support. This is an example of what [Barry Simon](#) calls a "Gem of spectral theory", that is.

**Example 2 (uniform weight Jacobi matrix)** Consider the monic orthogonal polynomials  $\pi_0(x), \pi_1(x), \dots, \pi_3(x)$  for  $w(x) = 1$  on  $[0, 1]$  constructed above. We can write the 3-term recurrence coefficients we have computed above as the Jacobi matrix:

$$x[\pi_0(x)|\pi_1(x)|\dots] = [\pi_0(x)|\pi_1(x)|\dots] \underbrace{\begin{bmatrix} 1/2 & 1/12 & & \\ 1 & 1/2 & 1/15 & \\ & 1 & 1/2 & \ddots \\ & & \ddots & \ddots \end{bmatrix}}_X$$

We can compute the orthonormal polynomials, using

$$\|\pi_3\|^2 = \int_0^1 (x^6 - 3x^5 + 69x^4/20 - 19x^3/10 + 51x^2/100 - 3x/50 + 1/400) dx = \frac{1}{280}$$

as:

$$\begin{aligned}
q_0(x) &= \pi_0(x) \\
q_1(x) &= \sqrt{12}\pi_1(x) = \sqrt{3}(2x - 1) \\
q_2(x) &= \sqrt{180}\pi_2(x) = \sqrt{5}(6x^2 - 6x + 1) \\
q_3(x) &= \sqrt{2800}\pi_3(x) = \sqrt{7}(20x^3 - 30x^2 + 12x - 1)
\end{aligned}$$

which have the Jacobi matrix

$$\begin{aligned}
x[q_0(x)|q_1(x)|\dots] &= x[\pi_0(x)|\pi_1(x)|\dots] \underbrace{\begin{bmatrix} 1 & & & \\ & 2\sqrt{3} & & \\ & & 6\sqrt{5} & \\ & & & 20\sqrt{7} \\ & & & \ddots \end{bmatrix}}_D \\
&= [q_0(x)|q_1(x)|\dots] D^{-1} X D = \begin{bmatrix} 1/2 & 1/\sqrt{12} & & \\ 1/\sqrt{12} & 1/2 & 1/\sqrt{15} & \\ & 1/\sqrt{15} & 1/2 & \ddots \\ & \ddots & \ddots & \ddots \end{bmatrix}
\end{aligned}$$

which is indeed symmetric. The problem sheet explores a more elegant way of doing this.

**Example 3 (expansion)** Consider expanding a low degree polynomial like  $f(x) = x^2$  in  $\pi_n(x)$ . We have

$$\begin{aligned}
\langle \pi_0, f \rangle &= \int_0^1 x^2 dx = 1/3 \langle \pi_1, f \rangle = \int_0^1 x^2(x - 1/2) dx = 1/12 \langle \pi_2, f \rangle = \\
&\quad \int_0^1 x^2(x^2 - x + 1/6) dx = 1/180
\end{aligned}$$

Thus we have:

$$f(x) = \frac{\pi_0(x)}{3} + \pi_1(x) + \pi_2(x) = [\pi_0(x)|\pi_1(x)|\pi_2(x)|\dots] \begin{bmatrix} 1/3 \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

We multiply (using that  $b_2 = 1$  for monic OPs) to deduce:

$$\begin{aligned}
xf(x) &= x[\pi_0(x)|\pi_1(x)|\pi_2(x)|\dots] \begin{bmatrix} 1/3 \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \\
&= [\pi_0(x)|\pi_1(x)|\pi_2(x)|\dots] X \begin{bmatrix} 1/3 \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} = [\pi_0(x)|\pi_1(x)|\pi_2(x)|\dots] \begin{bmatrix} 1/4 \\ 9/10 \\ 3/2 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \\
&= \frac{\pi_0(x)}{4} + \frac{9\pi_1(x)}{10} + \frac{3\pi_2(x)}{2} + \pi_3(x)
\end{aligned}$$

## III.4 Classical Orthogonal Polynomials

Classical orthogonal polynomials are special families of orthogonal polynomials with a number of beautiful properties, for example

1. Their derivatives are also OPs
2. They are eigenfunctions of simple differential operators

As stated above orthogonal polynomials are uniquely defined by the weight  $w(x)$  and the constant  $k_n$ . The classical orthogonal polynomials are:

1. Chebyshev polynomials (1st kind)  $T_n(x)$ :  $w(x) = 1/\sqrt{1-x^2}$  on  $[-1, 1]$ .
2. Chebyshev polynomials (2nd kind)  $U_n(x)$ :  $\sqrt{1-x^2}$  on  $[-1, 1]$ .
3. Legendre polynomials  $P_n(x)$ :  $w(x) = 1$  on  $[-1, 1]$ .
4. Ultrappherical polynomials (my fav!):  $C_n^{(\lambda)}(x)$ :  $w(x) = (1-x^2)^{\lambda-1/2}$  on  $[-1, 1]$ ,  $\lambda \neq 0, \lambda > -1/2$ .
5. Jacobi polynomials:  $P_n^{(a,b)}(x)$ :  $w(x) = (1-x)^a(1+x)^b$  on  $[-1, 1]$ ,  $a, b > -1$ .
6. Laguerre polynomials:  $L_n(x)$ :  $w(x) = \exp(-x)$  on  $[0, \infty)$ .
7. Hermite polynomials  $H_n(x)$ :  $w(x) = \exp(-x^2)$  on  $(-\infty, \infty)$ .

In the notes we will discuss:

1. Chebyshev polynomials: These are closely linked to Fourier series and are one of the most powerful tools in numerics.
2. Legendre polynomials: These have no simple closed-form expression but can be defined in terms of a Rodriguez formula, a feature that

applies to all other classical families.

### 1. Chebyshev

**Definition 1 (Chebyshev polynomials, 1st kind)**  $T_n(x)$  are orthogonal with respect to  $1/\sqrt{1-x^2}$  and satisfy:

$$\begin{aligned} T_0(x) &= 1, \\ T_n(x) &= 2^{n-1}x^n + O(x^{n-1}) \end{aligned}$$

**Definition 2 (Chebyshev polynomials, 2nd kind)**  $U_n(x)$  are orthogonal with respect to  $\sqrt{1-x^2}$ .

$$U_n(x) = 2^n x^n + O(x^{n-1})$$

**Theorem 1 (Chebyshev T are cos)** For  $-1 \leq x \leq 1$

$$T_n(x) = \cos n \arccos x.$$

In other words

$$T_n(\cos \theta) = \cos n\theta.$$

### Proof

We need to show that  $p_n(x) := \cos n\theta \cos x$  are

1. graded polynomials
2. orthogonal w.r.t.  $1/\sqrt{1-x^2}$  on  $[-1, 1]$ , and
3. have the right normalisation constant  $k_n = 2^{n-1}$  for  $n = 2, \dots$

Property (2) follows under a change of variables:

$$\int_{-1}^1 \frac{p_n(x)p_m(x)}{\sqrt{1-x^2}} dx = \int_{-\pi}^{\pi} \frac{\cos(n\theta)\cos(m\theta)}{\sqrt{1-\cos^2\theta}} \sin\theta d\theta = \int_{-\pi}^{\pi} \cos(n\theta)\cos(m\theta) dx = 0$$

if  $n \neq m$ .

To see that they are graded we use the fact that

$$xp_n(x) = \cos\theta \cos n\theta = \frac{\cos(n-1)\theta + \cos(n+1)\theta}{2} = \frac{p_{n-1}(x) + p_{n+1}(x)}{2}$$

In other words  $p_{n+1}(x) = 2xp_n(x) - p_{n-1}(x)$ . Since each time we multiply by  $2x$  and  $p_0(x) = 1$  we have

$$p_n(x) = (2x)^n + O(x^{n-1})$$

which completes the proof.

■

Buried in the proof is the 3-term recurrence:

### Corollary 1 (Chebyshev 3-term recurrence)

$$xT_0(x) = T_1(x)$$

$$xT_n(x) = \frac{T_{n-1}(x) + T_{n+1}(x)}{2}$$

Chebyshev polynomials are particularly powerful as their expansions are cosine series in disguise: for

$$f(x) = \sum_{k=0}^{\infty} f_k T_k(x)$$

we have

$$f(\cos\theta) = \sum_{k=0}^{\infty} f_k \cos k\theta.$$

Thus the coefficients can be recovered fast using FFT-based techniques as discussed in the problem sheet.

In the problem sheet we will also show the following:

**Theorem 2 (Chebyshev U are sin)** For  $x = \cos \theta$ ,

$$U_n(x) = \frac{\sin((n+1)\theta)}{\sin \theta}$$

which satisfy:

$$\begin{aligned} xU_0(x) &= U_1(x)/2 \\ xU_n(x) &= \frac{U_{n-1}(x)}{2} + \frac{U_{n+1}(x)}{2}. \end{aligned}$$

## 2. Legendre

**Definition 3 (Legendre)** Legendre polynomials  $P_n(x)$  are orthogonal polynomials with respect to  $w(x) = 1$  on  $[-1, 1]$ , with

$$k_n = \frac{1}{2^n} \binom{2n}{n} = \frac{(2n)!}{2^n (n!)^2}$$

The reason for this complicated normalisation constant is both historical and that it leads to simpler formulae for recurrence relationships.

Classical orthogonal polynomials have *Rodriguez formulae*, defining orthogonal polynomials as high order derivatives of simple functions. In this case we have:

**Lemma 1 (Legendre Rodriguez formula)**

$$P_n(x) = \frac{1}{(-2)^n n!} \frac{d^n}{dx^n} (1 - x^2)^n$$

**Proof** We need to verify:

1. graded polynomials
2. orthogonal to all lower degree polynomials on  $[-1, 1]$ , and
3. have the right normalisation constant  $k_n = \frac{1}{2^n} \binom{2n}{n}$ .

(1) follows since its a degree  $n$  polynomial (the  $n$ -th derivative of a degree  $2n$  polynomial). (2) follows by integration by parts. Note that  $(1 - x^2)^n$  and its first  $n - 1$  derivatives vanish at  $\pm 1$ . If  $r_m$  is a degree  $m < n$  polynomial we have:

$$\int_{-1}^1 \frac{d^n}{dx^n} (1-x^2)^n r_m(x) dx = - \int_{-1}^1 \frac{d^{n-1}}{dx^{n-1}} (1-x^2)^n r'_m(x) dx = \dots = (-)^n \int_{-1}^1 (1-x^2)^n r_m^{(n)}(x) dx = 0.$$

(3) follows since:

$$\begin{aligned} \frac{d^n}{dx^n} [(-)^n x^{2n} + O(x^{2n-1})] &= (-)^n 2n \frac{d^{n-1}}{dx^{n-1}} x^{2n-1} + O(x^{2n-1}) \\ &= (-)^n 2n(2n-1) \frac{d^{n-2}}{dx^{n-2}} x^{2n-2} + O(x^{2n-2}) = \dots \\ &= (-)^n 2n(2n-1) \dots (n+1)x^n + O(x^{n-1}) = (-)^n \frac{(2n)!}{n!} \end{aligned}$$

■  
This allows us to determine the coefficients  $k_n^{(\lambda)}$  which are useful in proofs. In particular we will use  $k_n^{(2)}$ :

### Lemma 2 (Legendre monomial coefficients)

$$\begin{aligned} P_0(x) &= 1 \\ P_1(x) &= x \\ P_n(x) &= \underbrace{\frac{(2n)!}{2^n(n!)^2} x^n}_{k_n} - \underbrace{\frac{(2n-2)!}{2^n(n-2)!(n-1)!} x^{n-2}}_{k_n^{(2)}} + O(x^{n-4}) \end{aligned}$$

(Here the  $O(x^{n-4})$  is as  $x \rightarrow \infty$ , which implies that the term is a polynomial of degree  $\leq n-4$ . For  $n=2, 3$  the  $O(x^{n-4})$  term is therefore precisely zero.)

### Proof

The  $n=0$  and  $1$  case are immediate. For the other case we expand  $(1-x^2)^n$  to get:

$$\begin{aligned} (-)^n \frac{d^n}{dx^n} (1-x^2)^n &= \frac{d^n}{dx^n} [x^{2n} - nx^{2n-2} + O(x^{2n-4})] \\ &= (2n) \dots (2n-n+1)x^n - n(2n-2) \dots (2n-2-n+1)x^{n-2} \\ &= \frac{(2n)!}{n!} x^n - \frac{n(2n-2)!}{(n-2)!} x^{n-2} + O(x^{n-4}) \end{aligned}$$

Multiplying through by  $\frac{1}{2^n(n!)}$  completes the derivation.

■  
**Theorem 3 (Legendre 3-term recurrence)**

$$xP_0(x) = P_1(x)$$

$$(2n+1)xP_n(x) = nP_{n-1}(x) + (n+1)P_{n+1}(x)$$

**Proof** The  $n = 0$  case is immediate (since  $w(x) = w(-x)$   $a_n = 0$ , from PS8). For the other cases we match terms:

$$(2n+1)xP_n(x) - nP_{n-1}(x) - (n+1)P_{n+1}(x) = [(2n+1)k_n - (n+1)k_{n+1}]x^{n+1}$$

$$+ [(2n+1)k_n^{(2)} - nk_{n-1} - (n+1)k_{n+1}^{(2)}]x^{n-1} + O(x^{n-3})$$

Using the expressions for  $k_n$  and  $k_n^{(2)}$  above we have (leaving the manipulations as an exercise):

$$(2n+1)k_n - (n+1)k_{n+1} = \frac{(2n+1)!}{2^n(n!)^2} - (n+1)\frac{(2n+2)!}{2^{n+1}((n+1)!)^2}$$

$$(2n+1)k_n^{(2)} - nk_{n-1} - (n+1)k_{n+1}^{(2)} = -(2n+1)\frac{(2n-2)!}{2^n(n-2)!(n-1)!} - n\frac{(2n-2)}{2^{n-1}((n-1)!)^2}$$

$$+ (n+1)\frac{(2n)!}{2^{n+1}(n-1)!n!} = 0$$

Thus

$$(2n+1)xP_n(x) - nP_{n-1}(x) - (n+1)P_{n+1}(x) = O(x^{n-3})$$

But as it is orthogonal to  $P_k(x)$  for  $0 \leq k \leq n-3$  it must be zero. ■

## III.5 Interpolation and quadrature

*Polynomial interpolation* is the process of finding a polynomial that equals data at a precise set of points. *Quadrature* is the act of approximating an integral by a weighted sum:

$$\int_a^b f(x)w(x)dx \approx \sum_{j=1}^n w_j f(x_j).$$

In these notes we see that the two concepts are intrinsically linked: interpolation leads naturally to quadrature rules.

1. Polynomial Interpolation: we describe how to interpolate a function by a polynomial and a set of points.
2. Interpolatory quadrature rule: polynomial interpolation leads naturally to ways to integrate

functions, but only realisable in the simplest cases.

### 1. Polynomial Interpolation

We already saw a special case of polynomial interpolation, where we saw that the polynomial

$$f(z) \approx \sum_{k=0}^{n-1} f_k^n z^k$$

equaled  $f$  at evenly spaced points on the unit circle:  $e^{i2\pi j/n}$ . But here we consider the following:

**Definition 1 (interpolatory polynomial)** Given  $n$  distinct points  $x_1, \dots, x_n \in \mathbb{R}$  and  $n$  samples  $f_1, \dots, f_n \in \mathbb{R}$ , a degree  $n - 1$  *interpolatory polynomial*  $p(x)$  satisfies

$$p(x_j) = f_j$$

The easiest way to solve this problem is to invert the Vandermonde system:

**Definition 2 (Vandermonde)** The *Vandermonde matrix* associated with  $n$  distinct points  $x_1, \dots, x_n \in \mathbb{R}$  is the matrix

$$V := \begin{bmatrix} 1 & x_1 & \cdots & x_1^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^{n-1} \end{bmatrix}$$

**Proposition 1 (interpolatory polynomial uniqueness)** The interpolatory polynomial is unique, and the Vandermonde matrix is invertible.

**Proof** Suppose  $p$  and  $\tilde{p}$  are both interpolatory polynomials. Then  $p(x) - \tilde{p}(x)$  vanishes at  $n$  distinct points  $x_j$ . By the fundamental theorem of algebra it must be zero, i.e.,  $p = \tilde{p}$ .

For the second part, if  $V\mathbf{c} = 0$  for  $\mathbf{c} \in \mathbb{R}$  then for  $q(x) = c_1 + \cdots + c_n x^{n-1}$  we have

$$q(x_j) = \mathbf{e}_j^\top V\mathbf{c} = 0$$

hence  $q$  vanishes at  $n$  distinct points and is therefore 0, i.e.,  $\mathbf{c} = 0$ .

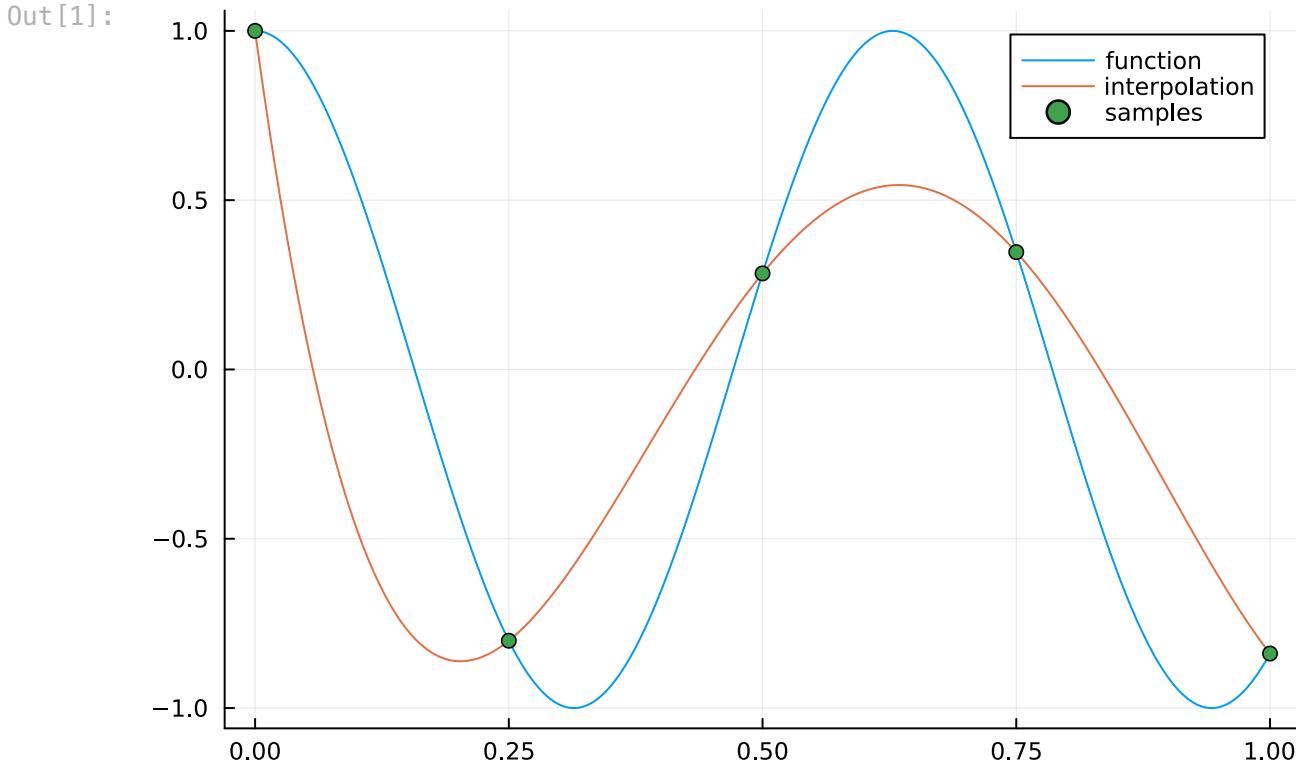
■

Thus a quick-and-dirty way to do interpolation is to invert the Vandermonde matrix (which we saw in the least squares setting with more samples than coefficients):

```
In [1]: using Plots, LinearAlgebra
f = x -> cos(10x)
n = 5

x = range(0, 1; length=n)# evenly spaced points (BAD for interpolation)
V = x .^ (0:n-1)' # Vandermonde matrix
c = V \ f.(x) # coefficients of interpolatory polynomial
p = x -> dot(c, x .^ (0:n-1))

g = range(0,1; length=1000) # plotting grid
plot(g, f.(g); label="function")
plot!(g, p.(g); label="interpolation")
scatter!(x, f.(x); label="samples")
```



But it turns out we can also construct the interpolatory polynomial directly. We will use the following which equal 1 at one grid point and zero at the others:

**Definition 3 (Lagrange basis polynomial)** The *Lagrange basis polynomial* is defined as

$$\ell_k(x) := \prod_{j \neq k} \frac{x - x_j}{x_k - x_j} = \frac{(x - x_1) \cdots (x - x_{k-1})(x - x_{k+1}) \cdots (x - x_n)}{(x_k - x_1) \cdots (x_k - x_{k-1})(x_k - x_{k+1}) \cdots (x_k - x_n)}$$

Plugging in the grid points verifies the following:

**Proposition 2 (delta interpolation)**

$$\ell_k(x_j) = \delta_{kj}$$

We can use these to construct the interpolatory polynomial:

**Theorem 1 (Lagrange interpolation)** The unique polynomial of degree at most  $n - 1$  that interpolates  $f$  at  $n$  distinct points  $x_j$  is:

$$p(x) = f(x_1)\ell_1(x) + \cdots + f(x_n)\ell_n(x)$$

**Proof** Note that

$$p(x_j) = \sum_{j=1}^n f(x_j)\ell_k(x_j) = f(x_j)$$

so we just need to show it is unique. Suppose  $\tilde{p}(x)$  is a polynomial of degree at most  $n - 1$  that also interpolates  $f$ . Then  $\tilde{p} - p$  vanishes at  $n$  distinct points. Thus by the fundamental theorem of algebra it must be zero.

■

**Example 1** We can interpolate  $\exp(x)$  at the points 0, 1, 2:

$$\begin{aligned} p(x) &= \ell_1(x) + e\ell_2(x) + e^2\ell_3(x) = \frac{(x-1)(x-2)}{(-1)(-2)} + e\frac{x(x-2)}{(-1)} + e^2\frac{x(x-1)}{2} \\ &= (1/2 - e + e^2/2)x^2 + (-3/2 + 2e - e^2/2)x + 1 \end{aligned}$$

**Remark** Interpolating at evenly spaced points is a really **bad** idea: interpolation is inherently ill-conditioned. The labs have explored this issue experimentally.

## 2. Interpolatory quadrature rules

By integrating an interpolant exactly we get a simple approach to approximating integrals. Using the Lagrange basis we can rewrite this procedure as a simple weighted sum:

**Definition 4 (interpolatory quadrature rule)** Given a set of points  $\mathbf{x} = [x_1, \dots, x_n]$  the interpolatory quadrature rule is:

$$\Sigma_n^{w,\mathbf{x}}[f] := \sum_{j=1}^n w_j f(x_j)$$

where

$$w_j := \int_a^b \ell_j(x) w(x) dx$$

**Proposition 3 (interpolatory quadrature is exact for polynomials)** Interpolatory quadrature is exact for all degree  $n - 1$  polynomials  $p$ :

$$\int_a^b p(x) w(x) dx = \Sigma_n^{w,\mathbf{x}}[f]$$

**Proof** The result follows since, by uniqueness of interpolatory polynomial:

$$p(x) = \sum_{j=1}^n p(x_j) \ell_j(x)$$

■

**Example 2 (arbitrary points)** Find the interpolatory quadrature rule for  $w(x) = 1$  on  $[0, 1]$  with points  $[x_1, x_2, x_3] = [0, 1/4, 1]$ ? We have:

$$\begin{aligned}
w_1 &= \int_0^1 w(x)\ell_1(x)dx = \int_0^1 \frac{(x-1/4)(x-1)}{(-1/4)(-1)}dx = -1/6 \\
w_2 &= \int_0^1 w(x)\ell_2(x)dx = \int_0^1 \frac{x(x-1)}{(1/4)(-3/4)}dx = 8/9 \\
w_3 &= \int_0^1 w(x)\ell_3(x)dx = \int_0^1 \frac{x(x-1/4)}{3/4}dx = 5/18
\end{aligned}$$

That is we have

$$\Sigma_n^{w,\mathbf{x}}[f] = -\frac{f(0)}{6} + \frac{8f(1/4)}{9} + \frac{5f(1)}{18}$$

This is indeed exact for polynomials up to degree 2 (and no more):

$$\Sigma_n^{w,\mathbf{x}}[1] = 1, \Sigma_n^{w,\mathbf{x}}[x] = 1/2, \Sigma_n^{w,\mathbf{x}}[x^2] = 1/3, \Sigma_n^{w,\mathbf{x}}[x^3] = 7/24 \neq 1/4.$$

**Example 3 (Chebyshev roots)** Find the interpolatory quadrature rule for  $w(x) = 1/\sqrt{1-x^2}$  on  $[-1, 1]$  with points equal to the roots of  $T_3(x)$ . This is a special case of Gaussian quadrature which we will approach in another way below. We use:

$$\int_{-1}^1 w(x)dx = \pi, \int_{-1}^1 xw(x)dx = 0, \int_{-1}^1 x^2w(x)dx = \pi/2$$

Recall from before that  $x_1, x_2, x_3 = \sqrt{3}/2, 0, -\sqrt{3}/2$ . Thus we have:

$$\begin{aligned}
w_1 &= \int_{-1}^1 w(x)\ell_1(x)dx = \int_{-1}^1 \frac{x(x+\sqrt{3}/2)}{(\sqrt{3}/2)\sqrt{3}\sqrt{1-x^2}}dx = \frac{\pi}{3} \\
w_2 &= \int_{-1}^1 w(x)\ell_2(x)dx = \int_{-1}^1 \frac{(x-\sqrt{3}/2)(x+\sqrt{3}/2)}{(-3/4)\sqrt{1-x^2}}dx = \frac{\pi}{3} \\
w_3 &= \int_{-1}^1 w(x)\ell_3(x)dx = \int_{-1}^1 \frac{(x-\sqrt{3}/2)x}{(-\sqrt{3})(-\sqrt{3}/2)\sqrt{1-x^2}}dx = \frac{\pi}{3}
\end{aligned}$$

(It's not a coincidence that they are all the same but this will differ for roots of other OPs.) That is we have

$$\Sigma_n^{w,\mathbf{x}}[f] = \frac{\pi}{3}(f(\sqrt{3}/2) + f(0) + f(-\sqrt{3}/2))$$

This is indeed exact for polynomials up to degree  $n-1=2$ , but it goes all the way up to  $2n-1=5$ :

$$\begin{aligned}
\Sigma_n^{w,\mathbf{x}}[1] &= \pi, \Sigma_n^{w,\mathbf{x}}[x] = 0, \Sigma_n^{w,\mathbf{x}}[x^2] = \frac{\pi}{2}, \\
\Sigma_n^{w,\mathbf{x}}[x^3] &= 0, \Sigma_n^{w,\mathbf{x}}[x^4] = \frac{3\pi}{8}, \Sigma_n^{w,\mathbf{x}}[x^5] = 0 \\
\Sigma_n^{w,\mathbf{x}}[x^6] &= \frac{9\pi}{32} \neq \frac{5\pi}{16}
\end{aligned}$$

We shall explain this miracle in the next chapter.

## III.6 Gaussian quadrature

In this chapter we see that a special quadrature rule can be constructed by using the roots of orthogonal polynomials, leading to a method that is exact for polynomials of twice the expected degree. Importantly, we can use quadrature to compute expansions in orthogonal polynomials that interpolate, mirroring the link between the Trapezium rule, Fourier series, and interpolation but now for orthogonal polynomials.

1. Truncated Jacobi matrices: we see that truncated Jacobi matrices are diagonalisable

in terms of orthogonal polynomials and their zeros. 2. Gaussian quadrature: Using roots of orthogonal polynomials and truncated Jacobi matrices leads naturally to an efficiently computable interpolatory quadrature rule. The *miracle* is its exact for twice as many polynomials as expected.

### 1. Roots of orthogonal polynomials and truncated Jacobi matrices

We now consider roots (zeros) of orthogonal polynomials  $p_n(x)$ . This is important as we shall see they are useful for interpolation and quadrature. For interpolation to be well-defined we first need to guarantee that the roots are distinct.

**Lemma 1** An orthogonal polynomial  $p_n(x)$  has exactly  $n$  distinct roots.

#### Proof

Suppose  $x_1, \dots, x_j$  are the roots where  $q_n(x)$  changes sign, that is,

$$p_n(x) = c_k(x - x_k)^{2p+1} + O((x - x_k)^{2p+2})$$

for  $c_k \neq 0$  and  $k = 1, \dots, j$  and  $p \in \mathbb{Z}$ , as  $x \rightarrow x_k$ . Then

$$p_n(x)(x - x_1) \cdots (x - x_j)$$

does not change signs: it behaves like  $c_k(x - x_k)^{2p+2} + O(x - x_k)^{2p+3}$  as  $x \rightarrow x_k$ . In other words:

$$\langle p_n, (x - x_1) \cdots (x - x_j) \rangle = \int_a^b p_n(x)(x - x_1) \cdots (x - x_j) w(x) dx \neq 0.$$

where  $w(x)$  is the weight of orthogonality. This is only possible if  $j = n$  as  $p_n(x)$  is orthogonal w.r.t. all lower degree polynomials.

■

**Definition 1 (truncated Jacobi matrix)** Given a symmetric Jacobi matrix  $X$  (associated with a family of orthonormal polynomials), the *truncated Jacobi matrix* is

$$J_n := \begin{bmatrix} a_0 & b_0 & & \\ b_0 & \ddots & \ddots & \\ & \ddots & a_{n-2} & b_{n-2} \\ & & b_{n-2} & a_{n-1} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

**Lemma 2 (zeros)** The zeros  $x_1, \dots, x_n$  of an orthonormal polynomial  $q_n(x)$  are the eigenvalues of the truncated Jacobi matrix  $J_n$ . More precisely,

$$J_n Q_n = Q_n \begin{bmatrix} x_1 & & \\ & \ddots & \\ & & x_n \end{bmatrix}$$

for the orthogonal matrix

$$Q_n = \underbrace{\begin{bmatrix} q_0(x_1) & \cdots & q_0(x_n) \\ \vdots & \cdots & \vdots \\ q_{n-1}(x_1) & \cdots & q_{n-1}(x_n) \end{bmatrix}}_{V_n^\top} \begin{bmatrix} \alpha_1^{-1} & & \\ & \ddots & \\ & & \alpha_n^{-1} \end{bmatrix}$$

where  $\alpha_j = \sqrt{q_0(x_j)^2 + \cdots + q_{n-1}(x_j)^2}$ .

### Proof

We construct the eigenvector (noting  $b_{n-1}q_n(x_j) = 0$ ):

$$J_n \begin{bmatrix} q_0(x_j) \\ \vdots \\ q_{n-1}(x_j) \end{bmatrix} = \begin{bmatrix} a_0 q_0(x_j) + b_0 q_1(x_j) \\ b_0 q_0(x_j) + a_1 q_1(x_j) + b_1 q_2(x_j) \\ \vdots \\ b_{n-3} q_{n-3}(x_j) + a_{n-2} q_{n-2}(x_j) + b_{n-2} q_{n-1}(x_j) \\ b_{n-2} q_{n-2}(x_j) + a_{n-1} q_{n-1}(x_j) + b_{n-1} q_n(x_j) \end{bmatrix} = x_j \begin{bmatrix} q_0(x_j) \\ q_1(x_j) \\ \vdots \\ q_{n-1}(x_j) \end{bmatrix}$$

The result follows from normalising the eigenvectors. Since  $J_n$  is symmetric the eigenvector matrix is orthogonal.

■

**Example 1 (Chebyshev roots)** Consider  $T_n(x) = \cos n \arccos x$ . The roots are  $x_j = \cos \theta_j$  where  $\theta_j = (j - 1/2)\pi/n$  for  $j = 1, \dots, n$  are the roots of  $\cos n\theta$  that are inside  $[0, \pi]$ .

Consider the  $n = 3$  case where we have

$$x_1, x_2, x_3 = \cos(\pi/6), \cos(\pi/2), \cos(5\pi/6) = \sqrt{3}/2, 0, -\sqrt{3}/2$$

We also have from the 3-term recurrence:

$$\begin{aligned} T_0(x) &= 1 \\ T_1(x) &= x \\ T_2(x) &= 2xT_1(x) - T_0(x) = 2x^2 - 1 \\ T_3(x) &= 2xT_2(x) - T_1(x) = 4x^3 - 3x \end{aligned}$$

We orthonormalise by rescaling

$$\begin{aligned} q_0(x) &= 1/\sqrt{\pi} \\ q_k(x) &= T_k(x)\sqrt{2}/\sqrt{\pi} \end{aligned}$$

so that the Jacobi matrix is symmetric:

$$x[q_0(x)|q_1(x)|\dots] = [q_0(x)|q_1(x)|\dots] \underbrace{\begin{bmatrix} 0 & 1/\sqrt{2} & & \\ 1/\sqrt{2} & 0 & 1/2 & \\ & 1/2 & 0 & 1/2 \\ & & 1/2 & 0 & \ddots \\ & & & \ddots & \ddots \end{bmatrix}}_X$$

We can then confirm that we have constructed an eigenvector/eigenvalue of the  $3 \times 3$  truncation of the Jacobi matrix, e.g. at  $x_2 = 0$ :

$$\begin{bmatrix} 0 & 1/\sqrt{2} & \\ 1/\sqrt{2} & 0 & 1/2 \\ & 1/2 & 0 \end{bmatrix} \begin{bmatrix} q_0(0) \\ q_1(0) \\ q_2(0) \end{bmatrix} = \frac{1}{\sqrt{\pi}} \begin{bmatrix} 0 & 1/\sqrt{2} & \\ 1/\sqrt{2} & 0 & 1/2 \\ & 1/2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -\sqrt{2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

## 2. Gaussian quadrature

Gaussian quadrature is the interpolatory quadrature rule corresponding to the grid  $x_j$  defined as the roots of the orthonormal polynomial  $q_n(x)$ . We shall see that it is exact for polynomials up to degree  $2n - 1$ , i.e., double the degree of other interpolatory quadrature rules from other grids.

**Definition 2 (Gauss quadrature)** Given a weight  $w(x)$ , the Gauss quadrature rule is:

$$\int_a^b f(x)w(x)dx \approx \underbrace{\sum_{j=1}^n w_j f(x_j)}_{\Sigma_n^w[f]}$$

where  $x_1, \dots, x_n$  are the roots of the orthonormal polynomials  $q_n(x)$  and

$$w_j := \frac{1}{\alpha_j^2} = \frac{1}{q_0(x_j)^2 + \cdots + q_{n-1}(x_j)^2}.$$

Equivalently,  $x_1, \dots, x_n$  are the eigenvalues of  $J_n$  and

$$w_j = \int_a^b w(x) dx Q_n[1, j]^2.$$

(Note we have  $\int_a^b w(x) dx q_0(x)^2 = 1$ .)

In analogy to how Fourier series are orthogonal with respect to the Trapezium rule, Orthogonal polynomials are orthogonal with respect to Gaussian quadrature:

**Lemma 3 (Discrete orthogonality)** For  $0 \leq \ell, m \leq n - 1$ , the orthonormal polynomials  $q_n(x)$  satisfy

$$\Sigma_n^w[q_\ell q_m] = \delta_{\ell m}$$

### Proof

$$\Sigma_n^w[q_\ell q_m] = \sum_{j=1}^n \frac{q_\ell(x_j) q_m(x_j)}{\alpha_j^2} = [q_\ell(x_1)/\alpha_1 | \cdots | q_\ell(x_n)/\alpha_n] \begin{bmatrix} q_m(x_1)/\alpha_1 \\ \vdots \\ q_m(x_n)/\alpha_n \end{bmatrix} = \mathbf{e}_\ell Q_n \delta_{\ell m}$$

■

Just as approximating Fourier coefficients using Trapezium rule gives a way of interpolating at the grid, so does Gaussian quadrature:

**Theorem 1 (interpolation via quadrature)** For the orthonormal polynomials  $q_n(x)$ ,

$$f_n(x) := \sum_{k=0}^{n-1} c_k^n q_k(x) \text{ for } c_k^n := \Sigma_n^w[f q_k]$$

interpolates  $f(x)$  at the Gaussian quadrature points  $x_1, \dots, x_n$ .

### Proof

Consider the Vandermonde-like matrix from above:

$$V_n := \begin{bmatrix} q_0(x_1) & \cdots & q_{n-1}(x_1) \\ \vdots & \ddots & \vdots \\ q_0(x_n) & \cdots & q_{n-1}(x_n) \end{bmatrix}$$

and define

$$Q_n^w := V_n^\top \begin{bmatrix} w_1 & & \\ & \ddots & \\ & & w_n \end{bmatrix} = \begin{bmatrix} q_0(x_1)w_1 & \cdots & q_0(x_n)w_n \\ \vdots & \ddots & \vdots \\ q_{n-1}(x_1)w_1 & \cdots & q_{n-1}(x_n)w_n \end{bmatrix}$$

so that

$$\begin{bmatrix} c_0^n \\ \vdots \\ c_{n-1}^n \end{bmatrix} = Q_n^w \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}.$$

Note that if  $p(x) = [q_0(x)|\cdots|q_{n-1}(x)]\mathbf{c}$  then

$$\begin{bmatrix} p(x_1) \\ \vdots \\ p(x_n) \end{bmatrix} = V_n \mathbf{c}$$

But we see that (similar to the Fourier case)

$$Q_n^w V_n = \begin{bmatrix} \Sigma_n^w[q_0q_0] & \cdots & \Sigma_n^w[q_0q_{n-1}] \\ \vdots & \ddots & \vdots \\ \Sigma_n^w[q_{n-1}q_0] & \cdots & \Sigma_n^w[q_{n-1}q_{n-1}] \end{bmatrix} = I_n$$

■

**Example 2 (Chebyshev expansions)** Consider the construction of Gaussian quadrature associated with the Chebyshev weight for  $n = 3$ . To determine the weights we need:

$$w_j^{-1} = \alpha_j^2 = q_0(x_j)^2 + q_1(x_j)^2 + q_2(x_j)^2 = \frac{1}{\pi} + \frac{2}{\pi}x_j^2 + \frac{2}{\pi}(2x_j^2 - 1)^2$$

We can check each case and deduce that  $w_j = \pi/3$ . Thus we recover the interpolatory quadrature rule. Further, we can construct the transform

$$\begin{aligned} Q_3^w &= \begin{bmatrix} w_1q_0(x_1) & w_2q_0(x_2) & w_3q_0(x_3) \\ w_1q_1(x_1) & w_2q_1(x_2) & w_3q_1(x_3) \\ w_1q_3(x_1) & w_2q_3(x_2) & w_3q_3(x_3) \end{bmatrix} \\ &= \frac{\pi}{3} \begin{bmatrix} 1/\sqrt{\pi} & 1/\sqrt{\pi} & 1/\sqrt{\pi} \\ x_1\sqrt{2/\pi} & x_2\sqrt{2/\pi} & x_3\sqrt{2/\pi} \\ (2x_1^2 - 1)\sqrt{2/\pi} & (2x_2^2 - 1)\sqrt{2/\pi} & (2x_3^2 - 1)\sqrt{2/\pi} \end{bmatrix} \\ &= \frac{\sqrt{\pi}}{3} \begin{bmatrix} 1 & 1 & 1 \\ \sqrt{6}/2 & 0 & -\sqrt{6}/2 \\ 1/\sqrt{2} & -\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \end{aligned}$$

We can use this to expand a polynomial, e.g.  $x^2$ :

$$Q_3^w \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_3^2 \end{bmatrix} = \frac{\sqrt{\pi}}{3} \begin{bmatrix} 1 & 1 & 1 \\ \sqrt{6}/2 & 0 & -\sqrt{6}/2 \\ 1/\sqrt{2} & -\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3/4 \\ 0 \\ 3/4 \end{bmatrix} = \begin{bmatrix} \sqrt{\pi}/2 \\ 0 \\ \sqrt{\pi}/(2\sqrt{2}) \end{bmatrix}$$

In other words:

$$x^2 = \frac{\sqrt{\pi}}{2} q_0(x) + \frac{\sqrt{\pi}}{2\sqrt{2}} q_2(x) = \frac{1}{2} T_0(x) + \frac{1}{2} T_2(x)$$

which can be easily confirmed.

**Corollary 1** Gaussian quadrature is an interpolatory quadrature rule with the interpolation points equal to the roots of  $q_n$ :

$$\Sigma_n^w[f] = \int_a^b f_n(x)w(x)dx$$

**Proof** We want to show that its the same as integrating the interpolatory polynomial:

$$\int_a^b f_n(x)w(x)dx = \frac{1}{q_0(x)} \sum_{k=0}^{n-1} c_k^n \int_a^b q_k(x)q_0(x)w(x)dx = \frac{c_0^n}{q_0} = \Sigma_n^w[f].$$

■

**Example 3 (Chebyshev quadrature via Lagrange polynomials)** The connection with interpolatory quadrature means theres another way to compute the quadrature: integrate the Lagrange polynomials associated with the zeros. For example, with  $n = 3$  recall that the roots of  $T_3(x)$  are  $\pm\sqrt{3}/2$  and 0. Thus we have

$$\begin{aligned} \ell_1(x) &= x(x + \sqrt{3}/2)/(3/2) = \frac{(2x^2 + \sqrt{3}x)}{3} \\ \ell_2(x) &= (x - \sqrt{3}/2)(x + \sqrt{3}/2)/(-3/4) = -\frac{4}{3}x^2 + 1 \\ \ell_3(x) &= (x - \sqrt{3}/2)x/(3/2) = \frac{(2x^2 - \sqrt{3}x)}{3} \end{aligned}$$

A quick check confirms that

$$w_j = \int_{-1}^1 \frac{\ell_j(x)}{\sqrt{1-x^2}} dx = \frac{\pi}{3}.$$

A consequence of being an interpolatory quadrature rule is that it is exact for all polynomials of degree  $n - 1$ . The *miracle* of Gaussian quadrature is it is exact for twice as many!

**Theorem 2 (Exactness of Gauss quadrature)** If  $p(x)$  is a degree  $2n - 1$  polynomial then Gauss quadrature is exact:

$$\int_a^b p(x)w(x)dx = \Sigma_n^w[p].$$

**Proof** Using polynomial division algorithm (e.g. by matching terms) we can write

$$p(x) = q_n(x)s(x) + r(x)$$

where  $s$  and  $r$  are degree  $n - 1$  and  $q_n(x)$  is the degree  $n$  orthonormal polynomial.

Then we have:

$$\begin{aligned} \Sigma_n^w[p] &= \underbrace{\Sigma_n^w[q_n s]}_{0 \text{ since evaluating } q_n \text{ at zeros}} + \Sigma_n^w[r] = \int_a^b r(x)w(x)dx \\ &= \underbrace{\int_a^b q_n(x)s(x)w(x)dx}_{0 \text{ since } s \text{ is degree } < n} + \int_a^b r(x)w(x)dx \\ &= \int_a^b p(x)w(x)dx. \end{aligned}$$

■

**Example 4 (Double exactness)** We are exact for all polynomials of degree  $2n - 1$ , so for our  $n = 3$  rule consider integrating  $x^5$ . We correctly get:

$$\Sigma_n^w[x^5] = \frac{\pi}{3} \left( \frac{9\sqrt{3}}{32} - \frac{9\sqrt{3}}{2} \right) = 0.$$

We are also correct for  $x^4$ :

$$\Sigma_n^w[x^4] = \frac{\pi}{3} \left( \frac{9}{16} + \frac{9}{16} \right) = \frac{3\pi}{8}.$$

However,

$$\Sigma_n^w[x^6] = \frac{9\pi}{64} \neq \frac{5\pi}{16}$$

hence it is incorrect for larger degree polynomials.