

Numerical Analysis and Optimization Homework Project 1

Marco Riggirello, Francesco Vaselli

This assignment is done using the Julia programming language, a high-level, high-performance language for technical computing. To ease the task, we load some Julia packages:

- **LinearAlgebra.jl**: the standard library for LinearAlgebra routines and utilities, leveraging BLAS subroutines;
- **Random.jl**: for generating random numbers;
- **Statistics.jl**: for computing statistics of samples;
- **Plots.jl**: for plotting our results;
- **BenchmarkTools.jl**: to provide benchmarking of our implementations.

```
1 using LinearAlgebra, Plots, Statistics, Random, BenchmarkTools
```

Problem 1

Task 1

Here we define the function `lufact()` that takes as input a square matrix and computes the non-pivoted LU factorization and its growth factor.

The \mathbf{L} and \mathbf{U} matrices are computed via Gaussian elimination, as presented in the class, while the growth factor γ is defined as the ratio of the largest entry in the matrix $\mathbf{G} \equiv |\mathbf{L}||\mathbf{U}|$ and the largest entry in the matrix $|\mathbf{A}|$ (absolute values are element-wise).

In our implementation, Gaussian elimination is carried on by storing both \mathbf{L} and \mathbf{U} in a single working matrix \mathbf{A}_k . In the last iteration n , the upper triangular part of \mathbf{A}_n represents the nonzero elements of \mathbf{U} while the lower part represent the \mathbf{L} matrix, diagonal excluded. The \mathbf{A}_k update is done using the slicing syntax of Julia's arrays.

We made our implementation more robust in various ways:

- in the first place, the Julia idiomatic type declaration (function `lufact(A::AbstractMatrix{T})` where `T <: Union{Real, Complex}`) allowed us to know how to correctly initialize \mathbf{A}_k and γ based on the matrix type. Specifically, for every subtype of `Real`, we initialize the copy to `Float64`, so that integers matrix factorization will not fail due to the creation of floating-point values in an integer matrix;
- furthermore, a preliminary check on the matrix squareness is performed;
- At every \mathbf{A}_k update, we check that the pivot is greater than the machine epsilon to ensure a sensible result;
- At the end, we also check for the presence of NaNs in \mathbf{A}_k , sign of numerical instability in this non-pivoted variant of the LU factorization.

lufact (generic function with 1 method)

```
1 function lufact(A::AbstractMatrix{T}) where T <: Union{Real, Complex}
2     m, n = size(A)
3     # Ensures the matrix is square.
4     if m != n
5         throw(ArgumentError("Matrix is not square."))
6     end
7
8     # Creates a copy of A
9     U = T <: Real ? Float64 : ComplexF64
10
11     A_k = map(U, A)
12     γ = zero(U)
13
14     # Performs the LU factorization.
15     for k in 1:n-1
16         pivot = A_k[k, k]
17         # Sanity check
18         if abs(pivot) <= eps(Float64)
19             throw(DomainError("Matrix is ill-conditioned."))
20         end
21         # Compute i-th column of L
22         A_k[k+1:n,k] = A_k[k+1:n,k] / pivot
23         # Update A_k
24         A_k[k+1:n,k+1:n] -= A_k[k+1:n,k] * (A_k[k,k+1:n])'
25
26     end
27     # check if we ended up with NaNs
28     if any(isnan, A_k)
29         throw(DomainError("Factorization is numerically unstable for this
30             matrix."))
31     end
32     # Constructs L and U from the updated matrix A_k.
33     L = UnitLowerTriangular(A_k)
34     U = UpperTriangular(A_k)
35
36     # Computes the growth factor γ.
37     G = abs.(L) * abs.(U)
38     γ = maximum(G) / maximum(abs.(A))
39
40     return L, U, γ
end
```

A note on growth factor

Ideally, the growth factor should be close to 1. This indicates that the LU decomposition did not significantly increase the magnitude of the matrix's elements, suggesting a stable decomposition. It's important to note that a low growth factor does not guarantee a good decomposition. It's just one of several indicators of the quality and stability of the decomposition.

Task 2 and 3

Now we want to test the performances of this algorithm over various matrix types. Before we proceed, we define a bunch of utility functions:

- introdurre in elenco?

RICORDIAMOCI DI COMMENTARE I RISULTATI SULLE VARIE MATRICI

```
1 md"""
2 Now we want to test the performances of this algorithm over various matrix
  types. Before we proceed, we define a bunch of utility functions:
3 - introdurre in elenco?
4
5 RICORDIAMOCI DI COMMENTARE I RISULTATI SULLE VARIE MATRICI
6 """
```

relative_backward_error (generic function with 1 method)

```
1 function relative_backward_error(A, L, U)
2     return opnorm(A - L * U, Inf)/opnorm(A, Inf)
3 end
```

print_summary (generic function with 1 method)

```
1 function print_summary(g, b, f, t)
2     println("Failure rate: $(100*f/t) %")
3     # add check if g, b are empty
4     if length(g) == 0 || length(b) == 0
5         println("No data to show.")
6         return
7     end
8     println("GROWTH FACTOR")
9     println("Min:      $(minimum(g))")
10    println("Max:      $(maximum(g))")
11    println("Mean:     $(mean(g))")
12    println("Median:  $(median(g))")
13    println("StdDev:   $(std(g))")
14    println("RELATIVE BACKWARD ERROR")
15    println("Min:      $(minimum(b))")
16    println("Max:      $(maximum(b))")
17    println("Mean:     $(mean(b))")
18    println("Median:  $(median(b))")
19    println("StdDev:   $(std(b))")
20 end
```

plot_summary (generic function with 2 methods)

```
1 function plot_summary(g, b, pq=1.)
2   # Adjust the binning for histograms if necessary to make the x-axis less
3   crowded
4   # Adjust the size of the plot or the labels to prevent squeezing
5   hg = histogram(g[g .<= quantile(g, pq)], xlabel="Growth factor",
6   ylabel="Frequency",
7   legend=false, xrotation=45, bins=30) # Rotate x-axis labels
8   and set bins
9   hb = histogram(b[b .<= quantile(g, pq)], xlabel="Relative backward error",
10  ylabel="Frequency",
11  legend=false, xrotation=45, bins=30) # Rotate x-axis labels
12  and set bins
13
14  # Combine the two histograms into one plot without a legend
15  plot(hg, hb, layout=(1,2), size=(800, 400)) # Adjust the size as needed
16 end
```

test_dataset (generic function with 2 methods)

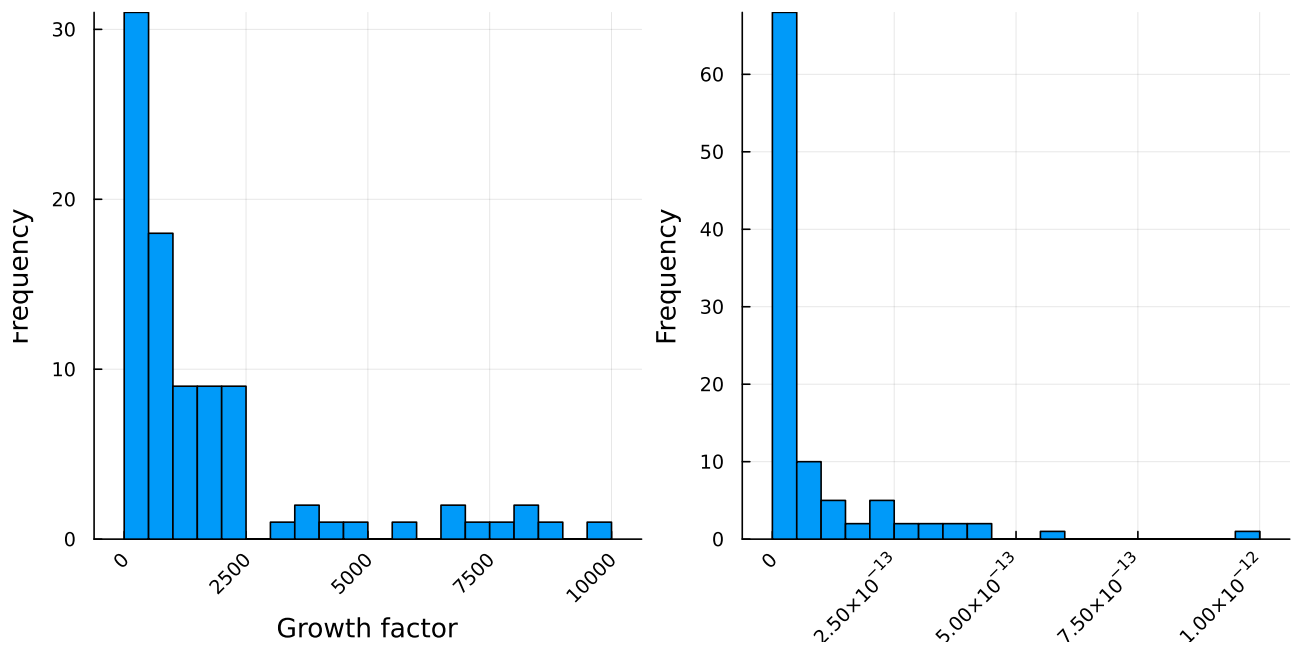
```
1 function test_dataset(dataset_generator, t, plot_quantile=1.)
2   Random.seed!(42)
3   f = 0
4   g = Float64[]
5   b = Float64[]
6   for i in 1:t
7     A = dataset_generator(i)
8     try
9       L, U,  $\gamma$  = lufact(A)
10       $\beta$  = relative_backward_error(A, L, U)
11      push!(g,  $\gamma$ )
12      push!(b,  $\beta$ )
13    catch e
14      if isa(e, DomainError)
15        f += 1
16      else
17        throw(e)
18      end
19    end
20  end
21  print_summary(g, b, f, t)
22  plot_summary(g, b, plot_quantile)
23  #return g, b
24 end
```

Random matrices

We generate real matrices of uniformly distributed sizes in the range 2:100 and normal distributed elements.

generate_random (generic function with 1 method)

```
1 function generate_random(i)
2     N = rand(2:100)
3     return randn(Float64, N, N)
4 end
```



```
1 test_dataset(generate_random, 100, 0.9)
```



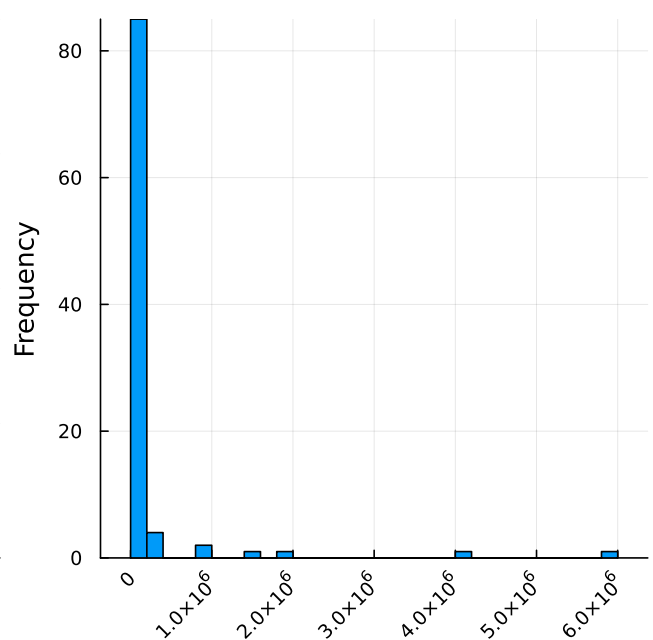
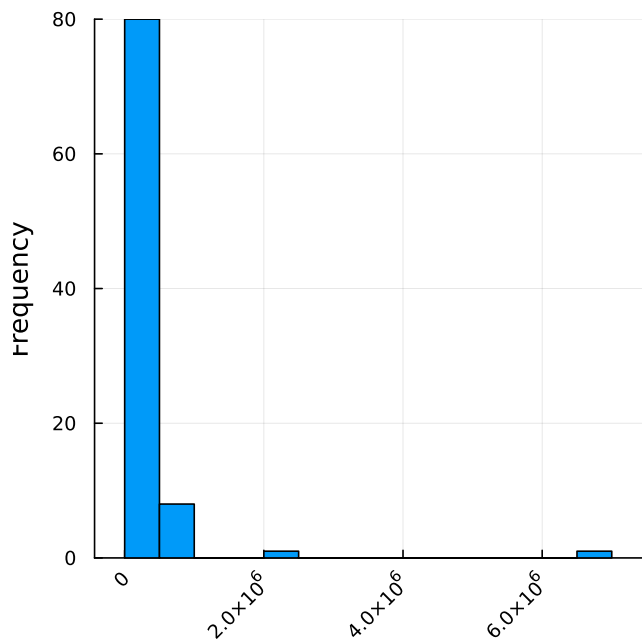
```
Failure rate: 0.0 %
GROWTH FACTOR
Min: 1.0
Max: 47101.04950458209
Mean: 3405.46298922573
Median: 1131.8943097590077
StdDev: 6362.426046577013
RELATIVE BACKWARD ERROR
Min: 0.0
Max: 9.664073393522062e-13
Mean: 8.028283558195911e-14
Median: 2.040374369416543e-14
StdDev: 1.4448853727177071e-13
```

Furthermore, we generate complex matrix with the same size range and elements distribution.

```
1 md"""
2 Furthermore, we generate complex matrix with the same size range and elements
  distribution.
3 """
```

generate_random_complex (generic function with 1 method)

```
1 function generate_random_complex(i)
2     N = rand(2:100)
3     return randn(ComplexF64, N, N)
4 end
```



```
1 test_dataset(generate_random_complex, 100, 0.9)
```

```
> Failure rate: 0.0 %
GROWTH FACTOR
Min: 2.671102767863563
Max: 1.4204629724095856e10
Mean: 1.9728159765561047e8
Median: 19854.886046128588
StdDev: 1.4595643568192155e9
RELATIVE BACKWARD ERROR
Min: 1.2498785104196468
Max: 2.2868260696960516e9
Mean: 2.82932720756088e7
Median: 4442.241962951809
StdDev: 2.3071620968102214e8
```

FOR RANDOM IT SEEMS OK. SAY IT BETTER

```
1 md"""
2 FOR RANDOM IT SEEMS OK. SAY IT BETTER
3 """
```

Hilbert matrices

Hilber matrices are defined per-element as

$$H_{ij} = \frac{1}{i + j - 1}$$

so for example the 4×4 Hilbert matrix is

$$\begin{bmatrix} 1 & 1/2 & 1/3 & 1/4 \\ 1/2 & 1/3 & 1/4 & 1/5 \\ 1/3 & 1/4 & 1/5 & 1/6 \\ 1/4 & 1/5 & 1/6 & 1/7 \end{bmatrix}$$

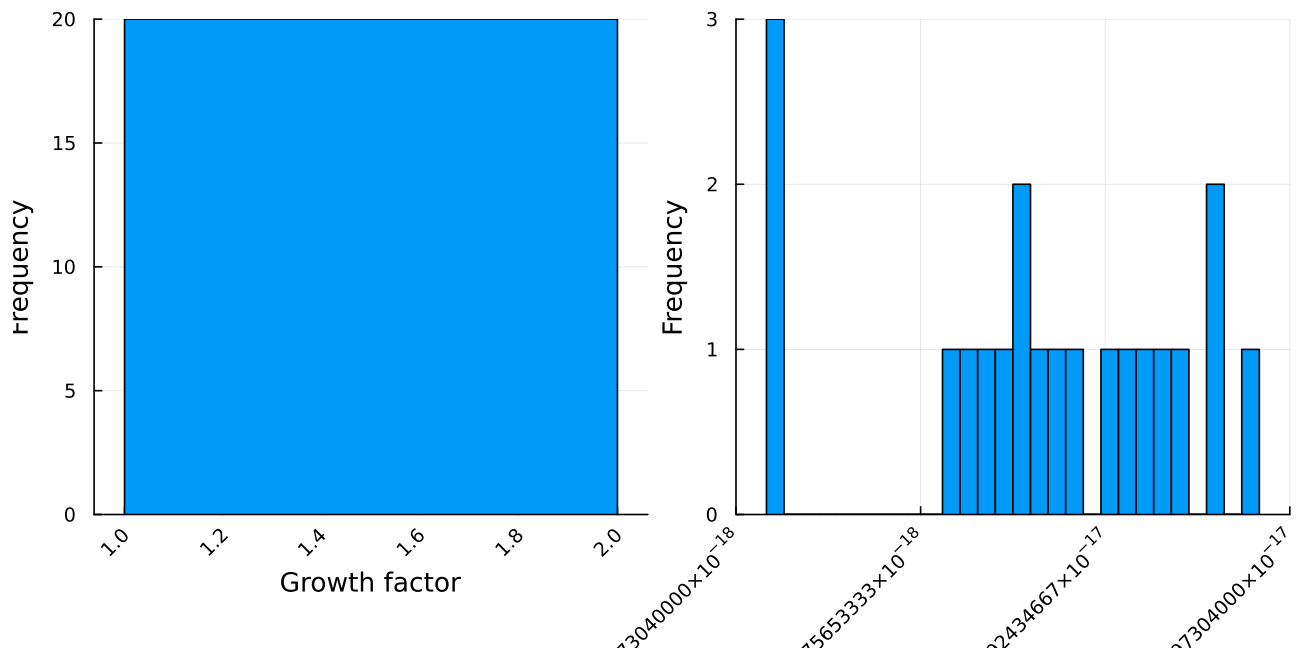
MENTION ILL-CONDITIONED NATURE OF THESE MATRICES

Hence we define the proper generator function.

```
1 md"""
2 Hence we define the proper generator function.
3 """
```

generate_hilbert (generic function with 1 method)

```
1 function generate_hilbert(N)
2     return [1 / (i + j - 1) for i in 1:N, j in 1:N]
3 end
```

```
1 test_dataset(generate_hilbert, 300)
```

```
> Failure rate: 93.33333333333333 %
GROWTH FACTOR
Min: 1.0
Max: 1.0
Mean: 1.0
Median: 1.0
StdDev: 0.0
RELATIVE BACKWARD ERROR
Min: 0.0
Max: 2.7001540943022594e-17
Mean: 1.5533358019968334e-17
Median: 1.5701316941667965e-17
StdDev: 8.271771338970318e-18
```

FANNO CAAA qed

Diagonally dominant matrices

A square matrix \mathbf{A} is said to be diagonally dominant (by rows) if for every row of the matrix, the magnitude of the diagonal entry in a row is greater than or equal to the sum of the magnitudes of all the other (non-diagonal) entries in that row. Mathematically, this can be written as:

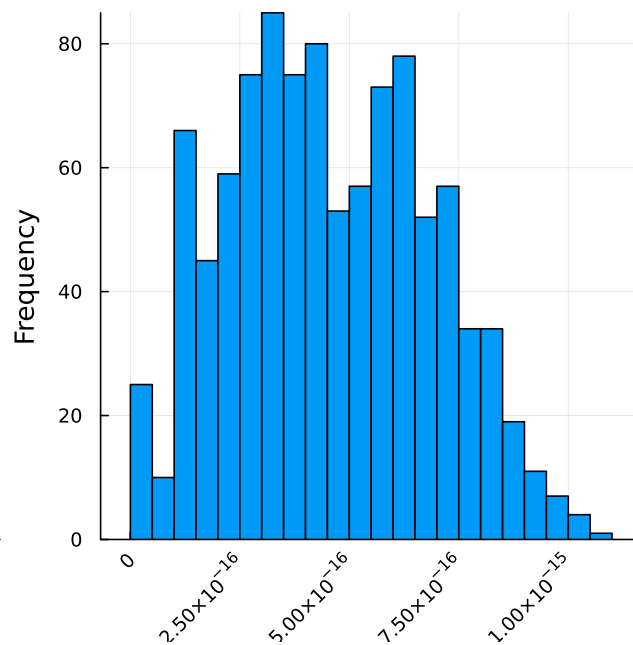
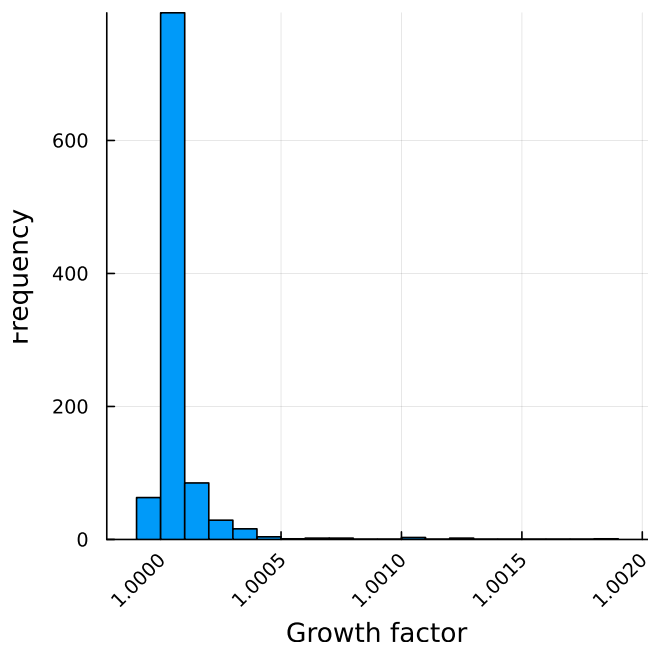
$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}| \quad \text{for all } i.$$

This condition must hold for each i from 1 to n , where n is the size of the $n \times n$ matrix \mathbf{A} . If the inequality is strict for all rows, then \mathbf{A} is said to be strictly diagonally dominant. The LU factorization without pivoting is known to be backward stable for this irreducible DD matrices.

We built a diagonally dominant generator function by first generating a $N \times N$ matrix with element uniformly distributed in the range 0:1 and then we summed $N + 1$ to the diagonal element.

generate_dd (generic function with 1 method)

```
1 function generate_dd(i)
2   N = rand(2:100)
3   A = rand(N, N)
4   A += diagm([N+1 for _ in 1:N])
5   return A
6 end
```



```
1 test_dataset(generate_dd, 1000)
```



```
Failure rate: 0.0 %
GROWTH FACTOR
Min: 0.9999999999999996
Max: 1.0018434432259373
Mean: 1.0000555146830499
Median: 1.000013938106836
StdDev: 0.00012768933341942098
RELATIVE BACKWARD ERROR
Min: 0.0
Max: 1.0565702407697936e-15
Mean: 4.563304322455912e-16
Median: 4.327370074077119e-16
StdDev: 2.2881824703137433e-16
```



stable as qed

UndefVarError: `plot_correlation` not defined

1. top-level scope @ (Local: 1)

```
1 plot_correlation(g3, b3)
```

SPD matrices

A symmetric matrix \mathbf{A} of size $n \times n$ is said to be positive definite if the following two conditions are met:

1. The matrix is symmetric;
2. For any non-zero vector \mathbf{x} in \mathbb{R}^n , the quadratic form $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is positive, i.e., $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$.

Mathematically, this can be expressed as:

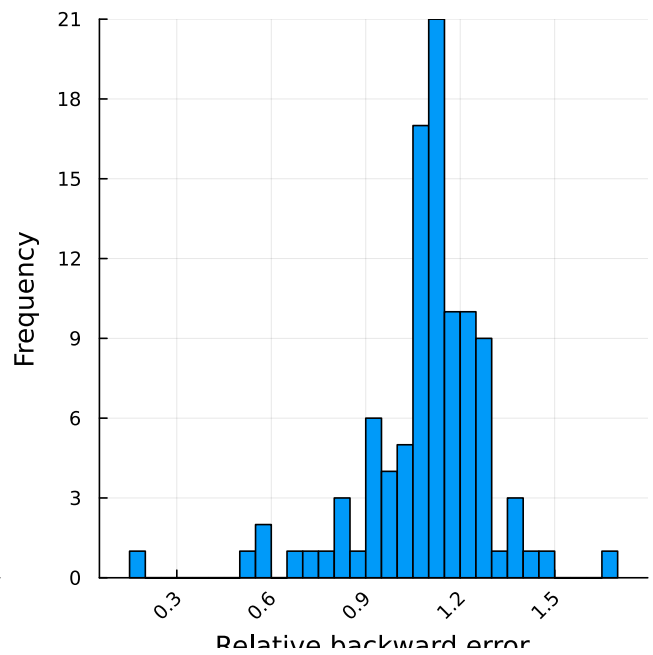
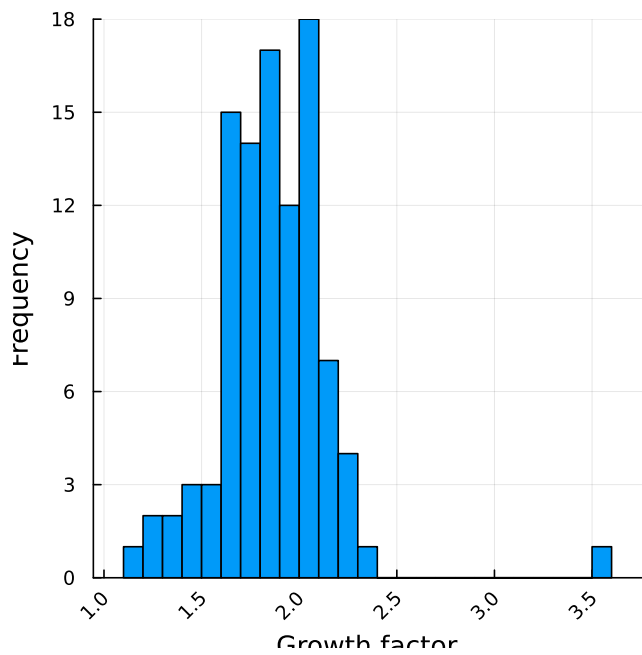
$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}.$$

Additionally, a matrix is positive definite if and only if all its eigenvalues are positive.

```
1 md"""
2 A symmetric matrix  $\mathbf{A}$  of size  $n \times n$  is said to be positive definite if
  the following two conditions are met:
3
4 1. The matrix is symmetric;
5 2. For any non-zero vector  $\mathbf{x}$  in  $\mathbb{R}^n$ , the quadratic form
   $\mathbf{x}^T \mathbf{A} \mathbf{x}$  is positive, i.e.,  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ .
6
7 Mathematically, this can be expressed as:
8
9  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}.$ 
10
11 Additionally, a matrix is positive definite if and only if all its eigenvalues
  are positive.
12
13 """
```

generate_spd (generic function with 1 method)

```
1 function generate_spd(i)
2     N = rand(2:100)
3     A = randn(Complex{Float64}, N, N)
4     return A*A'
5 end
```



```
1 test_dataset(generate_spd, 100)
```

```
> Failure rate: 0.0 %
GROWTH FACTOR
Min: 1.1683629837591885
Max: 3.584486788468825
Mean: 1.8637815560595936
Median: 1.866012813410923
StdDev: 0.2892846992223377
RELATIVE BACKWARD ERROR
Min: 0.19547137206541038
Max: 1.6932937729927284
Mean: 1.0965237628900593
Median: 1.1176707045489345
StdDev: 0.2010187457276718
```

commento risultati su stabilita'

an example of performance profiling

```
1 md"""
2 an example of performance profiling
3 """
```

UndefVarError: `A` not defined

```
1. var"##core#409"() @ execution.jl:547
2. var"##sample#410"(::Tuple{}, ::BenchmarkTools.Parameters) @ execution.jl:556
3. var"#_run#48"(::Bool, ::String, ::Base.Pairs{Symbol, Integer, NTuple{4,
  Symbol}, NamedTuple{(:samples, :evals, :gctrail, :gcsample), Tuple{Int64,
  Int64, Bool, Bool}}}, ::typeof(BenchmarkTools._run),
  ::BenchmarkTools.Benchmark, ::BenchmarkTools.Parameters) @ execution.jl:109
4. #invokelatest#2 @ essentials.jl:821 [inlined]
5. invokelatest @ essentials.jl:816 [inlined]
6. #run_result#45 @ execution.jl:41 [inlined]
7. run_result @ execution.jl:40 [inlined]
8. var"#run#49"(::Nothing, ::Float64, ::Float64, ::Base.Pairs{Symbol, Integer,
  NTuple{5, Symbol}, NamedTuple{(:verbose, :samples, :evals, :gctrail,
  :gcsample), Tuple{Bool, Int64, Int64, Bool, Bool}}}, ::typeof(run),
  ::BenchmarkTools.Benchmark, ::BenchmarkTools.Parameters) @ execution.jl:134
9. run @ execution.jl:126 [inlined]
10. #warmup#54 @ execution.jl:189 [inlined]
11. warmup(::BenchmarkTools.Benchmark) @ execution.jl:188
12. macro expansion @ { Local: 432 [inlined]
13. top-level scope @ { Local: 12
```

```
1 begin
2
3 # Define the range of matrix sizes
4 n_values = 400:400:4000
5 # Array to store the benchmark results
6 timings = []
7
8 # Benchmark the lufact() function for different sizes of N
9 for N in n_values
10     A = randn(N, N)
11     # Benchmark and get the median time
12     bench = @benchmark lufact(A)
13     median_time = median(bench).time / 1e9 # Convert to seconds
14     push!(timings, median_time)
15 end
16
17 # Plot the timings on a log-log graph
18 scatter(n_values, timings, label="data", legend=:topleft, xscale=:log10,
19         yscale=:log10, xlabel="n", ylabel="elapsed time (s)")
20 plot!(n_values, timings[end]*n_values./(n_values[end]).^3, label="O(n^3)",
21       linestyle=:dash, xscale=:log10, yscale=:log10)
22 end
```

Problem 2

The Wilkinson matrix is

$$\begin{bmatrix} 1 & 0 & 0 & \cdots & 1 \\ -1 & 1 & 0 & \cdots & 1 \\ -1 & -1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \cdots & -1 & 1 \end{bmatrix}$$

e qui ci mettiamo due o tre conticini. Sono le otto di sabato, non lo farò ora.

Task 1

LU decomposition of a Wilkinson Matrix W_5

$$W_5 = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 \\ -1 & 1 & 0 & 0 & 1 \\ -1 & -1 & 1 & 0 & 1 \\ -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & 1 \end{bmatrix} = LU = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ -1 & -1 & 1 & 0 & 0 \\ -1 & -1 & -1 & 1 & 0 \\ -1 & -1 & -1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 2 \\ 0 & 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 1 & 8 \\ 0 & 0 & 0 & 0 & 16 \end{bmatrix}$$

Task 2

Guess of the LU factorization of W_n for a generic n

$$L = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ -1 & 1 & 0 & \cdots & 0 \\ -1 & -1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & -1 & \cdots & 1 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 1 & 0 & \cdots & 0 & 2 \\ 0 & 0 & 1 & \cdots & 0 & 2^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 2^{n-1} \end{bmatrix}$$

or

$$L = I_n - \sum_{i=2}^n \sum_{j=1}^{i-1} e_i e_j^T$$

$$U = I_n + \sum_{i=1}^{n-1} 2^{i-1} e_i e_n^T$$

Task 3

introduce function. unite in one?

wilkinson_element (generic function with 1 method)

```
1 function wilkinson_element(i, j, N)
2     if i == j || j == N
3         return 1.
4     elseif i < j
5         return 0.
6     else
7         return -1.
8     end
9 end
```

wilkin (generic function with 1 method)

```
1 function wilkin(N::Integer)
2     return [wilkinson_element(i,j,N) for i in 1:N, j in 1:N]
3 end
```

Task 4

The vector \mathbf{b} formed by $\mathbf{b} = \mathbf{A}\mathbf{e}$ for the Wilkinson matrix \mathbf{W}_n where \mathbf{e} is a column vector of ones is given by:

$$\mathbf{b} = \begin{bmatrix} 2 \\ 1 \\ 0 \\ -1 \\ -2 \\ \vdots \\ -(n-3) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ -1 \end{bmatrix}$$

With each entry b_i defined as:

$$b_i = \begin{cases} -(i-3) & \text{for } 1 \leq i < n \\ -(i-2) & \text{for } i = n \end{cases}$$

We thus define a function to store in memory the exact solution to the problem:

```

1 md"""
2 ### Task 4
3 The vector  $\mathbf{b}$  formed by  $\mathbf{b} = \mathbf{A}\mathbf{e}$  for the Wilkinson
  matrix  $\mathbf{W}_n$  where  $\mathbf{e}$  is a column vector of ones is given by:
4
5  $\mathbf{b} = \begin{bmatrix}$ 
6 2 \\
7 1 \\
8 0 \\
9 -1 \\
10 -2 \\
11  $\vdots$  \\
12  $-(n-3)$ 
13  $\end{bmatrix} + \begin{bmatrix}$ 
14 0 \\
15 0 \\
16 0 \\
17 0 \\
18 0 \\
19  $\vdots$  \\
20  $-1$ 
21  $\end{bmatrix}$ 
22
23 With each entry  $b_i$  defined as:
24
25  $b_i =$ 
26  $\begin{cases}$ 
27  $-(i-3)$  &  $\text{for } 1 \leq i < n$  \\
28  $-(i-2)$  &  $\text{for } i = n$ 
29  $\end{cases}$ 
30
31 We thus define a function to store in memory the exact solution to the problem:
32 """

```


generate_b_vector (generic function with 1 method)

```
1 function generate_b_vector(n)
2     b = [-(i-3) for i in 1:n]
3     b[1] = 2
4     if n > 1
5         b[end] = - (n - 2)
6     end
7     return b
8 end
```

```
1 begin
2 for n in 2:100
3     @assert wilkin(n)*[1 for i in 1:n] == generate_b_vector(n)
4     #@show wilkin(n)*[1 for i in 1:n]
5 end
6
7 # see which one is faster
8
9 #@btime b100 = generate_b_vector(100)
10 #@btime b100v = generate_b_vector_third(100)
11 end
```

1 Enter cell code...

Not final , to revise

In task 4, we are asked to perform a numerical experiment with the Wilkinson matrix (W_n) and the vector (e), which leads to the following steps:

1. Generate ($A = W_n$), the Wilkinson matrix of size ($n \times n$).
2. Let (e) be the column vector with all entries equal to 1, then form ($b = Ae$), resulting in a specific pattern as described above.
3. Solve the linear system ($Ax = b$) using the backslash operator, which in Julia (or MATLAB) leverages optimized algorithms for solving linear equations.
4. The computed solution (x) is then compared to the exact solution, which should be the vector (e).

Given the specific structure of (W_n), the backslash operator should ideally find the exact solution without difficulty for smaller values of (n). However, as (n) becomes large, numerical instability can occur due to the ill-conditioning of the Wilkinson matrix. This can lead to a computed solution (x) that deviates from the exact solution (e), especially in the presence of round-off errors in floating-point arithmetic.

The success of this numerical experiment heavily depends on the numerical stability of the algorithms used by the backslash operator and the conditioning of the matrix (W_n). For (W_{60}), we would need to assess the accuracy of the computed solution by comparing it to (e) and examining the residual ($r = b - Ax$), which should be close to the zero vector for an accurate solution.

AssertionError: x == e

```
1. top-level scope @ Local: 10 [inlined]
```

```

1 begin
2   A = wilkin(60)
3   local b = generate_b_vector(60)
4   # a 60 dim vector of ones
5   e = ones(60)
6   x = A \ b
7   @show x
8   diff = x - e
9   @show norm(diff, Inf)
10  @assert x == e
11 end

```

```
x = [1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.  
0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.  
0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.  
0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.  
0, 1.0]  
norm(diff, Inf) = 1.0
```

Task 5 and 6

Repeat the experiment for smaller values of n . What is the largest value of n for which $W_n x = b$ can be solved accurately by GEPP when $b = W_n e$? Provide an explanation of the observed behavior.

DomainError with relative error introduced at step n = 55:

1. top-level scope @ (Local: 14 [inlined])

```
1 # repeat the experiment for smaller values of n
2 # NOTE: for 1 it is false
3 begin
4   for n in 2:60
5     W_n = wilkin(n)
6     b = generate_b_vector(n)
7     e = ones(n)
8     x = W_n \ b
9
10    Δ = x - e
11    ε_rel = norm(Δ, Inf) / norm(e, Inf)
12    if ε_rel != 0
13      @show ε_rel
14      throw(DomainError("relative error introduced at step n = $n"))
15    end
16    @assert x == e
17    println("n = $n")
18  end
19 end
```



```
n = 2
n = 3
n = 4
n = 5
n = 6
n = 7
n = 8
n = 9
n = 10
n = 11
n = 12
n = 13
n = 14
n = 15
n = 16
n = 17
n = 18
n = 19
n = 20
n = 21
n = 22
n = 23
n = 24
n = 25
n = 26
n = 27
n = 28
n = 29
n = 30
n = 31
n = 32
n = 33
```



The Wilkinson matrix is specifically designed to be a challenging test case for numerical algorithms due to its structure, which induces a large growth factor (2^{n-1}) in the elements of the LU decomposition without pivoting. As the size of the matrix grows, the condition number of the Wilkinson matrix increases exponentially, making it more susceptible to round-off errors.

At $n = 55$ we have a transition: we get a first non-zero entry for the relative error, which becomes 1 under the Infinity norm. Why is that the case? We consider the analytical solution to the problem using the known LU factorization for a Wilkinson matrix. Let $Lz = b$:

We have the following set of equations:

$$\begin{aligned} z_1 &= b_1, \\ z_2 &= b_2 + z_1 = b_2 + b_1, \\ z_3 &= b_3 + z_2 + z_1 = b_3 + b_2 + 2b_1, \\ z_4 &= b_4 + z_3 + z_2 + z_1 = b_4 + b_3 + 2b_2 + 4b_1, \\ &\vdots \end{aligned}$$

This pattern continues, and we can generalize it to the equation for any z_i :

$$z_i = b_i + b_{i-1} + 2b_{i-2} + \dots + 2^{i-2}b_1.$$

This can be written more compactly as:

$$z_i = b_i + \sum_{k=1}^{i-1} 2^{k-1} b_{i-k}.$$

We now use the analytical expression for b to rewrite z_i as:

$$z_i = -(i-3) - \sum_{k=1}^{i-1} 2^{k-1} (i-k-3).$$

except for the last z_n where $b_n = -(i-2)$.

We continue from the previous step:

$$-(i-3) - (i-3) \left(\sum_{k=1}^{i-1} 2^{k-1} \right) + \sum_{k=1}^{i-1} k 2^{k-1} = \dots$$

This simplifies to:

$$-(i-3)(2^i - 1) + \sum_{k=1}^{i-1} k 2^{k-1} = \dots$$

And we can conclude with:

$$z_i = 2^{i-1} + 1$$

and

$$z_n = 2^{n-1}$$

DA RIMUOVERE? Finally, we can relate this to the geometric progression sum formula:

$$\sum_{k=0}^N 2^k = 2^{N+1} - 1.$$

For example, adding powers of 2 like $1 + 2 + 4 + 8 + 16 + 32$ can be calculated using the sum formula for a geometric series:

$$\sum_{k=0}^N 2^k = 2^{N+1} - 1.$$

Now we continue with the backward substitution step, $U\mathbf{x} = \mathbf{z}$, and we know that the following holds for the specific for of U in the Wilkinson matrix:

$$\begin{aligned} x_n &= \frac{z_n}{2^{n-1}} = \frac{2^{n-1}}{2^{n-1}} = 1, \\ &\vdots \\ x_i + 2^{i-1}x_n &= z_i, \rightarrow x_i = z_i - 2^{i-1} = 2^{i-1} + 1 - 2^{i-1} \end{aligned}$$

CONTINUARE QUI

The previous is an exact explanation of the loss of accuracy we observe. However, in order to investigate further this problem, we can recall Theorem 4.29 of the lectures, which states that the relative error for our case is bounded by:

$$\frac{\|\mathbf{x}^* - \tilde{\mathbf{x}}\|}{\|\mathbf{x}^*\|} \leq O(\gamma u) \kappa(W_n)$$

where $\gamma = 2^{n-1}$ for the Wilkinson matrix. We can thus compute the upper bound of ϵ_{rel} in the vicinity of $n = 55$ to convince ourselves that the problem is unstable in this regime:

```

1 begin
2   for n in 45:60
3     W_n = wilkin(n)
4     κ = cond(W_n, Inf)
5     γ = 2^(n-1)
6     ε_bound = γ * eps()/2 * κ
7     println("At n = $n ε_bound ≈ $ε_bound")
8   end
9 end

```

```

>
At n = 45 ε_bound ≈ 0.087890625
At n = 46 ε_bound ≈ 0.1796875
At n = 47 ε_bound ≈ 0.3671875
At n = 48 ε_bound ≈ 0.75
At n = 49 ε_bound ≈ 1.53125
At n = 50 ε_bound ≈ 3.125
At n = 51 ε_bound ≈ 6.375
At n = 52 ε_bound ≈ 13.0
At n = 53 ε_bound ≈ 26.5
At n = 54 ε_bound ≈ 54.0
At n = 55 ε_bound ≈ 110.0
At n = 56 ε_bound ≈ 224.0
At n = 57 ε_bound ≈ 456.0
At n = 58 ε_bound ≈ 928.0
At n = 59 ε_bound ≈ 1888.0
At n = 60 ε_bound ≈ 3840.0

```

We see that the relative error bound, ϵ_{bound} , is of order 1 at about $n = 48$. After that, numerical optimizations make so that the solution can still be computed accurately, but we do not have any guarantee from theory.

Problem 3

Task 1

Point a)

We want to prove that the matrix $\tilde{A} = A + \mathbf{u}\mathbf{v}^T$ is nonsingular if and only if $\mathbf{v}^T A^{-1} \mathbf{u} \neq -1$.

To start, let's remember that a matrix is nonsingular (or invertible) if it has a nonzero determinant.

First we introduce the matrix determinant lemma which states that for any invertible $n \times n$ matrix A and column vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, the determinant of $A + \mathbf{u}\mathbf{v}^T$ is given by:

$$\det(A + \mathbf{u}\mathbf{v}^T) = \det(A) \cdot (1 + \mathbf{v}^T A^{-1} \mathbf{u})$$

Proof of the Matrix Determinant Lemma

We start from the special case $A = I$. Let's define the bordered matrix

$$\begin{pmatrix} I + \mathbf{u}\mathbf{v}^T & \mathbf{u} \\ 0 & 1 \end{pmatrix}$$

the following identity holds:

$$\begin{pmatrix} I & 0 \\ \mathbf{v}^T & 1 \end{pmatrix} \begin{pmatrix} I + \mathbf{u}\mathbf{v}^T & \mathbf{u} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} I & 0 \\ -\mathbf{v}^T & 1 \end{pmatrix} = \begin{pmatrix} I & \mathbf{u} \\ 0 & 1 + \mathbf{v}^T \mathbf{u} \end{pmatrix}$$

The first and third matrix on the left hand side of the equation are unit lower triangular matrices, thus their determinants are 1. Since the determinant of a block triangular matrix is the product of the determinants of the diagonal blocks, we have:

$$\det(I + \mathbf{u}\mathbf{v}^T) = (1 + \mathbf{v}^T \mathbf{u}).$$

It is simple the to derive the general case:

$$\begin{aligned} \det(A + \mathbf{u}\mathbf{v}^T) &= \det(A) \det(I + (A^{-1} \mathbf{u}) \mathbf{v}^T) \\ &= \det(A) (1 + \mathbf{v}^T (A^{-1} \mathbf{u})). \end{aligned}$$

And this ends the proof of the Matrix Determinant Lemma.

For \tilde{A} to be nonsingular, $\det(\tilde{A})$ must be nonzero. Based on the matrix determinant lemma, this will be the case if and only if:

$$1 + \mathbf{v}^T A^{-1} \mathbf{u} \neq 0$$

Then, for $\tilde{\mathbf{A}}$ to be invertible,

$$\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq -1$$

And this concludes the proof.

Point b)

We want to derive the Sherman Morrison formula:

$$\tilde{\mathbf{A}}^{-1} = \mathbf{A}^{-1} - \alpha \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}, \quad \alpha = \frac{1}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}}$$

To do so, let's assume that $\tilde{\mathbf{A}}^{-1}$ is of the form $\mathbf{A}^{-1} + \mathbf{B}$. Then

$$\begin{aligned} \tilde{\mathbf{A}} \tilde{\mathbf{A}}^{-1} &= \mathbf{I} = (\mathbf{A} + \mathbf{u} \mathbf{v}^T)(\mathbf{A}^{-1} + \mathbf{B}) \\ &= \mathbf{I} + \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1} + (\mathbf{A} + \mathbf{u} \mathbf{v}^T) \mathbf{B} \end{aligned}$$

and hence

$$\begin{aligned} \mathbf{B} &= -(\mathbf{A} + \mathbf{u} \mathbf{v}^T)^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1} \\ &= -(\mathbf{A}^{-1} + \mathbf{B}) \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}. \end{aligned}$$

After some manipulation, we find that

$$(1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}) \mathbf{B} = -\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}$$

Therefore, we can isolate \mathbf{B} on one side of the equation:

$$\mathbf{B} = -\alpha \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}$$

With α defined as above: the Sherman Morrison formula is found.

Point c)

Since we have the LU factorization of A , let's define \mathbf{x} as the solution of the problem $A\mathbf{x} = \tilde{\mathbf{b}}$ which can be found by backward-forward substitution with a complexity n^2 .

Now, our goal is to find the solution $\tilde{\mathbf{x}}$ to the problem $\tilde{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$. Formally, we have $\tilde{\mathbf{x}} = \tilde{A}^{-1}\tilde{\mathbf{b}}$ but by using the Sherman Morrison formula we can write

$$\begin{aligned}\tilde{\mathbf{x}} &= (A^{-1} - \alpha A^{-1} \mathbf{u} \mathbf{v}^T A^{-1}) \tilde{\mathbf{b}} \\ &= A^{-1} \tilde{\mathbf{b}} - \alpha A^{-1} \mathbf{u} \mathbf{v}^T A^{-1} \tilde{\mathbf{b}} \\ &= \mathbf{x} - \alpha A^{-1} \mathbf{u} \mathbf{v}^T \mathbf{x} \\ A^{-1} \mathbf{u} &\equiv \mathbf{y}\end{aligned}$$

can be found using the LU decomposition with backward forward substitution as well (quadratic complexity). The remaining operations are vector dot products or vector sums, which are of linear complexity.

To conclude, the following algorithm then finds $\tilde{\mathbf{x}}$ with $\mathcal{O}(n^2)$ flops:

1. Find \mathbf{x} by forward backward substitution;
2. Find \mathbf{y} by forward backward substitution;
3. Compute $\alpha = 1/(1 + \mathbf{v}^T \mathbf{y})$;
4. Compute $\tilde{\mathbf{x}} = \mathbf{x} - \alpha \mathbf{y}(\mathbf{v}^T \mathbf{x})$.

Task 2

The bordered system

QUA CI PENSA GPTGPTGPTGPT

corresponds to the system

ANCHE QUA BRUMBRUMBRUM

Then, by substitution of \mathbf{z} , we can recounduce to the previous case, that has quadratic complexity.