

# report sylvester equation

sicheng4

February 2024

## Main

This section contains main theorems and proofs, as well as some errors in original paper.

### Proposition 2.1

Let  $\mathbf{S}_U \mathbf{U}_{d+1} = \mathbf{Q}_U \mathbf{U}_{d+1} \mathbf{T}_U$  be a reduced QR decomposition with

$$\mathbf{Q}_U \mathbf{U}_{d+1} = [\mathbf{Q}_U, \mathbf{Q}_{U,d+1}] \text{ and } \mathbf{T}_U \mathbf{U}_{d+1} = \begin{bmatrix} \mathbf{T}_{U,d} & \mathbf{T}_{H,d+1} \\ \mathbf{0}^\top & \boldsymbol{\tau}_{d+1} \end{bmatrix}$$

Then, for the sketched method, the following Arnoldi-like formula holds:

$$\mathbf{S}_U \mathbf{A} \mathbf{U}_d = \mathbf{S}_U \mathbf{U}_d (\mathbf{H}_d + R_H E_d^\top) + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top$$

with  $R_H = \mathbf{T}_{U,d}^{-1} \mathbf{T}_H \mathbf{h}_{d+1,d}$  and  $\mathbf{Q}_U \mathbf{U}_{d+1} \perp \mathbf{S}_U \mathbf{U}_d$ . Similarly, if

$$\mathbf{S}_V \mathbf{V}_{d+1} = \mathbf{Q}_V \mathbf{V}_{d+1} \mathbf{T}_V, \quad \mathbf{T}_V \mathbf{V}_{d+1} = \begin{bmatrix} \mathbf{T}_{V,d} & \mathbf{T}_{G,d+1} \\ \mathbf{0}^\top & \boldsymbol{\theta}_{d+1} \end{bmatrix}$$

then

$$\mathbf{S}_V \mathbf{B}^\top \mathbf{V}_d = \mathbf{S}_V \mathbf{V}_d (\mathbf{G}_d + R_G E_d^\top) + \mathbf{Q}_V \mathbf{V}_{d+1} \boldsymbol{\theta}_{d+1} \mathbf{g}_{d+1,d} E_d^\top,$$

where  $R_G := \mathbf{T}_{V,d}^{-1} \mathbf{T}_G \mathbf{g}_{d+1,d}$ .

### Proof

$$\begin{aligned} \mathbf{S}_U \mathbf{A} \mathbf{U}_d &= \mathbf{S}_U \mathbf{U}_{d+1} \underline{\mathbf{H}}_d = \mathbf{S}_U \mathbf{U}_{d+1} \mathbf{H}_d = \mathbf{S}_U \mathbf{U}_{d+1} [\mathbf{H}_d; \mathbf{H}^\top] \\ &= \mathbf{Q}_U \mathbf{U}_{d+1} \mathbf{T}_U \mathbf{U}_{d+1} [\mathbf{H}_d; \mathbf{H}^\top] \\ &= [\mathbf{Q}_U, \mathbf{Q}_{U,d+1}] \begin{bmatrix} \mathbf{T}_{U,d} & \mathbf{T}_{H,d+1} \\ \mathbf{0}^\top & \boldsymbol{\tau}_{d+1} \end{bmatrix} [\mathbf{H}_d; \mathbf{H}^\top] \\ &= [\mathbf{Q}_U, \mathbf{Q}_{U,d} \mathbf{T}_U, \mathbf{Q}_{U,d} \mathbf{T}_{H,d+1} + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1}] [\mathbf{H}_d; \mathbf{H}^\top] \\ &= \mathbf{Q}_U \mathbf{U}_{d+1} \mathbf{T}_U \mathbf{H}_d + \mathbf{Q}_U \mathbf{U}_{d+1} \mathbf{T}_{H,d+1} \mathbf{H}^\top + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1} \mathbf{H}^\top \\ &= \mathbf{Q}_U \mathbf{U}_{d+1} \mathbf{T}_U (\mathbf{H}_d + \mathbf{T}_{U,d}^{-1} \mathbf{T}_{H,d+1} \mathbf{H}^\top) + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \\ &= \mathbf{S}_U \mathbf{U}_d (\mathbf{H}_d + \mathbf{T}_{U,d}^{-1} \mathbf{T}_H \mathbf{h}_{d+1,d} E_d^\top) + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \\ &= \mathbf{S}_U \mathbf{U}_d (\mathbf{H}_d + R_H E_d^\top) + \mathbf{Q}_U \mathbf{U}_{d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \end{aligned}$$

where letting  $R_H = \mathbf{T}_{U,d}^{-1} T_H \mathbf{h}_{d+1,d}$  and  $Q_{U,d+1} \perp \mathbf{S}_U \mathbf{U}_d$  The same process holds for  $\mathbf{S}_V$

The second part is the verification of whitened-sketched Arnoldi relations. The bases are changed to

$$\widehat{\mathbf{U}}_d := \mathbf{U}_d \mathbf{T}_{U,d}^{-1}, \quad \widehat{\mathbf{V}}_d := \mathbf{V}_d \mathbf{T}_{V,d}^{-1}$$

. Again, we will only show the  $\widehat{\mathbf{U}}_d$ , since the other one could be derived by same process. By last proposition, we have

$$\mathbf{S}_U \mathbf{A} \mathbf{U}_d = \mathbf{S}_U \mathbf{U}_d \left( \mathbf{H}_d + \mathbf{T}_{U,d}^{-1} T_H \mathbf{h}_{d+1,d} E_d^\top \right) + Q_{U,d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top$$

Times  $\mathbf{T}_{U,d}^{-1}$  on both side, and just do simple calculations. Note that,

$$\mathbf{T}_{U,d}^{-1} = \begin{bmatrix} \mathbf{T}_{U,d-1}^{-1} & -\mathbf{T}_{U,d-1}^{-1} T_{H,d-1} \boldsymbol{\tau}_d^{-1} \\ & -\boldsymbol{\tau}_d^{-1} \end{bmatrix}$$

This gives us:

$$\begin{aligned} \mathbf{S}_U \mathbf{A} \mathbf{U}_d \mathbf{T}_{U,d}^{-1} &= \mathbf{S}_U \mathbf{U}_d \left( \mathbf{H}_d + \mathbf{T}_{U,d}^{-1} T_H \mathbf{h}_{d+1,d} E_d^\top \right) \mathbf{T}_{U,d}^{-1} + Q_{U,d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \mathbf{T}_{U,d}^{-1} \\ &= \mathbf{S}_U \mathbf{U}_d \mathbf{T}_{U,d}^{-1} \mathbf{T}_{U,d} \left( \mathbf{H}_d + \mathbf{T}_{U,d}^{-1} T_H \mathbf{h}_{d+1,d} E_d^\top \right) \mathbf{T}_{U,d}^{-1} + Q_{U,d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \mathbf{T}_{U,d}^{-1} \\ &= \mathbf{S}_U \widehat{\mathbf{U}}_d \left( \widehat{\mathbf{H}}_d + \widehat{\mathbf{H}} E_d^\top \right) + Q_{U,d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} E_d^\top \mathbf{T}_{U,d}^{-1} \\ &= \mathbf{S}_U \widehat{\mathbf{U}}_d \left( \widehat{\mathbf{H}}_d + \widehat{\mathbf{H}} E_d^\top \right) + \textcolor{blue}{Q_{U,d+1} \boldsymbol{\tau}_{d+1} \mathbf{h}_{d+1,d} \boldsymbol{\tau}_d^{-1} E_d^\top} \end{aligned}$$

where

$$\widehat{\mathbf{H}}_d = \mathbf{T}_{U,d} \mathbf{H}_d \mathbf{T}_{U,d}^{-1}, \widehat{\mathbf{H}} = T_{H,d+1} \mathbf{h}_{d+1,d} \boldsymbol{\tau}_d^{-1}$$

**Question1** ??? Note here, the last term is different with the item in paper. I could not figure out why the last term should be  $Q_{U,d+1} \mathbf{h}_{d+1,d} E_d^\top$

## Typos on AL1

1. Updae  $\widehat{\mathbf{H}}_{d+1}$  At d-th iteration, we could get  $\widehat{\mathbf{H}}_{d+1}$  which would be used in next iteration as  $\widehat{\mathbf{H}}_d$ . Note,  $\mathbf{T}_{\mathbf{U},d+1}^{-1} = \begin{bmatrix} \mathbf{T}_{\mathbf{U},d}^{-1} & -\mathbf{T}_{\mathbf{U},d}^{-1}T_{H,d+1}\boldsymbol{\tau}_{d+1}^{-1} \\ \boldsymbol{\tau}_{d+1}^{-1} \end{bmatrix}$

By  $\widehat{\mathbf{H}}_{d+1} = \mathbf{T}_{\mathbf{U},d+1}\mathbf{H}_{d+1}\mathbf{T}_{\mathbf{U},d+1}^{-1}$ , we could imply

$$\begin{aligned} \widehat{\mathbf{H}}_{d+1} &= \begin{bmatrix} \mathbf{T}_{\mathbf{U},d} & T_{H,d+1} \\ & \boldsymbol{\tau}_{d+1} \end{bmatrix} \begin{bmatrix} \mathbf{H}_{d+1} & H \\ 0 \dots 0, \mathbf{h}_{d+1,d} & \mathbf{h}_{d+1,d+1} \end{bmatrix} \begin{bmatrix} \mathbf{T}_{\mathbf{U},d}^{-1} & -\mathbf{T}_{\mathbf{U},d}^{-1}T_{H,d+1}\boldsymbol{\tau}_{d+1}^{-1} \\ \boldsymbol{\tau}_{d+1}^{-1} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{T}_{\mathbf{U},d}\mathbf{H}_{d+1} + T_{H,d+1}\mathbf{h}_{d+1,d}\mathbf{E}_d^\top & \mathbf{T}_{\mathbf{U},d}H + T_{H,d+1}\mathbf{h}_{d+1,d+1} \\ \boldsymbol{\tau}_{d+1}\mathbf{h}_{d+1,d}\mathbf{E}_d^\top & \boldsymbol{\tau}_{d+1}\mathbf{h}_{d+1,d+1} \end{bmatrix} \begin{bmatrix} \mathbf{T}_{\mathbf{U},d}^{-1} & -\mathbf{T}_{\mathbf{U},d}^{-1}T_{H,d+1}\boldsymbol{\tau}_{d+1}^{-1} \\ \boldsymbol{\tau}_{d+1}^{-1} \end{bmatrix} \\ &= \begin{bmatrix} \widehat{\mathbf{H}}_d + T_H\mathbf{h}_{d+1,d}\boldsymbol{\tau}_d^{-1}\mathbf{E}_d^\top & \widehat{\mathbf{H}}_{\text{new}} \\ \boldsymbol{\tau}_{d+1}\mathbf{h}_{d+1,d}\boldsymbol{\tau}_d^{-1}\mathbf{E}_d^\top & \boldsymbol{\tau}_{d+1}(-\mathbf{h}_{d+1,d}\boldsymbol{\tau}_d^{-1}\mathbf{E}_d^\top T_{H,d+1} + \mathbf{h}_{d+1,d+1})\boldsymbol{\tau}_{d+1}^{-1} \end{bmatrix} \end{aligned}$$

where  $\widehat{\mathbf{H}}_{\text{new}} = -\left(\widehat{\mathbf{H}}_d + T_H\mathbf{h}_{d+1,d}\boldsymbol{\tau}_d^{-1}\mathbf{E}_d^\top\right)T_{H,d+1}\boldsymbol{\tau}_{d+1}^{-1} + \mathbf{T}_{\mathbf{U},d}H\boldsymbol{\tau}_{d+1}^{-1} + T_{H,d+1}\mathbf{h}_{d+1,d+1}\boldsymbol{\tau}_{d+1}^{-1}$  and  $H = \left[\mathbf{h}_{1,d+1}^\top, \dots, \mathbf{h}_{d,d+1}^\top\right]^\top \in R^{dr \times r}$ . A corresponding update for  $\widehat{\mathbf{G}}_{d+1}$  is performed.

2.  $\mathbf{E}_d = [\text{zeros}((d-1)*r, r); \text{eye}(r)]$  and  $\mathbf{E}_1 = [\text{eye}(r); \text{zeros}((d-1)*r, r)]$

I could run the code, but the solution is too "bad", the residue satisfies tolerance while the solution NOT

## Report on 12/2

I have tried AL1 and AL2 , as well as the direct computation to get  $\mathbf{H}$ .

### Parameter

There are some parameters play important roles in AL.

1.  $\mathbf{k}$   $k \neq 1$ , otherwise the first for loop would be fail.
2.  $\mathbf{s}$  The dimension of sketched matrix can not be too small.  $dr \leq s$
3.  $\mathbf{p}$  A larger  $p$  will lead to shorter running time

If we do direct computation instead of updating  $\mathbf{H}$ , and let the tol be  $10^{-6}$ , the accuracy would be also lie in it. However, if we use a tighter tol, the accuracy tends to be worse. More specifically, it seems not converges well. The following is the my experiment results with direction computation of  $\mathbf{H}$  and  $s = 400$  and  $p = 10, n1 = n2 = 1200$ :

tol	residue
1e-06	3.1566e-08
1e-07	3.214e-09
1e-08	4.2937
1e-09	4.2937
1e-10	4.2937

I think the possible reason could be the limitation of  $s$ . Intutively, when you require a higher tol, the iteration should be longer. However, once  $s$  is determined, the max number of iterations could not be higher than  $s/r$ . So, I let  $s$  to be larger . I let  $n1 = n2 = 120$ , and tried to let  $s = 100$ , and 120, the accuracy is bad, while  $s = 60, 80$  are good, and  $s = 20, 40$  are getting worse again. This gives us hints that  $s$  should be determined very carefully  
The problem may also arise because of the accuracy of `lyap()`.