

# SOLUTION OF SPECIALIZED SYLVESTER EQUATION

ONDRA KAMENIK

Given the following matrix equation

$$AX + BX \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} C \right) = D,$$

where  $A$  is regular  $n \times n$  matrix,  $X$  is  $n \times m^i$  matrix of unknowns,  $B$  is singular  $n \times n$  matrix,  $C$  is  $m \times m$  regular matrix with  $|\beta(C)| < 1$  (i.e. modulus of largest eigenvalue is less than one),  $i$  is an order of Kronecker product, and finally  $D$  is  $n \times m^i$  matrix.

First we multiply the equation from the left by  $A^{-1}$  to obtain:

$$X + A^{-1}BX \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} C \right) = A^{-1}D$$

Then we find real Schur decomposition  $K = UA^{-1}BU^T$ , and  $F = VCV^T$ . The equation can be written as

$$UX \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} V^T \right) + KUX \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} V^T \right) \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} F \right) = UA^{-1}D \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} V^T \right)$$

This can be rewritten as

$$Y + KY \left( \begin{smallmatrix} i \\ \otimes \end{smallmatrix} F \right) = \widehat{D},$$

and vectorized

$$\left( I + \begin{smallmatrix} i \\ \otimes \end{smallmatrix} F^T \otimes K \right) \text{vec}(Y) = \text{vec}(\widehat{D})$$

Let  ${}^iF$  denote  $\begin{smallmatrix} i \\ \otimes \end{smallmatrix} F^T$  for the rest of the text.

**Lemma 1.** *For any  $n \times n$  matrix  $A$  and  $\beta_1\beta_2 > 0$ , if there is exactly one solution of*

$$\left( I_2 \otimes I_n + \begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix} \otimes A \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} d_1 \\ d_2 \end{pmatrix},$$

*then it can be obtained as solution of*

$$\begin{aligned} (I_n + 2\alpha A + (\alpha^2 + \beta^2)A^2) x_1 &= \widehat{d}_1 \\ (I_n + 2\alpha A + (\alpha^2 + \beta^2)A^2) x_2 &= \widehat{d}_2 \end{aligned}$$

Typeset by  $\mathcal{A}\mathcal{M}\mathcal{S}$ - $\text{\texttt{TeX}}$

where  $\beta = \sqrt{\beta_1\beta_2}$ , and

$$\begin{pmatrix} \hat{d}_1 \\ \hat{d}_2 \end{pmatrix} = \left( I_2 \otimes I_n + \begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix} \otimes A \right) \begin{pmatrix} d_1 \\ d_2 \end{pmatrix}$$

*Proof.* Since

$$\begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix} \begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix} = \begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix} \begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix} = \begin{pmatrix} \alpha^2 + \beta^2 & 0 \\ 0 & \alpha^2 + \beta^2 \end{pmatrix},$$

it is easy to see that if the equation is multiplied by

$$I_2 \otimes I_n + \begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix} \otimes A$$

we obtain the result. We only need to prove that the matrix is regular. But this is clear because matrix

$$\begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix}$$

collapses an eigenvalue of  $A$  to  $-1$  iff the matrix

$$\begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix}$$

does.  $\square$

**Lemma 2.** For any  $n \times n$  matrix  $A$  and  $\delta_1\delta_2 > 0$ , if there is exactly one solution of

$$\left( I_2 \otimes I_n + 2\alpha \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes A + (\alpha^2 + \beta^2) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix}^2 \otimes A^2 \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} d_1 \\ d_2 \end{pmatrix}$$

it can be obtained as

$$\begin{aligned} (I_n + 2a_1A + (a_1^2 + b_1^2)A^2) (I_n + 2a_2A + (a_2^2 + b_2^2)A^2) x_1 &= \hat{d}_1 \\ (I_n + 2a_1A + (a_1^2 + b_1^2)A^2) (I_n + 2a_2A + (a_2^2 + b_2^2)A^2) x_2 &= \hat{d}_2 \end{aligned}$$

where

$$\begin{pmatrix} \hat{d}_1 \\ \hat{d}_2 \end{pmatrix} = \left( I_2 \otimes I_n + 2\alpha \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes A + (\alpha^2 + \beta^2) \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix}^2 \otimes A^2 \right) \begin{pmatrix} d_1 \\ d_2 \end{pmatrix}$$

and

$$\begin{aligned} a_1 &= \alpha\gamma - \beta\delta \\ b_1 &= \alpha\delta + \gamma\beta \\ a_2 &= \alpha\gamma + \beta\delta \\ b_2 &= \alpha\delta - \gamma\beta \\ \delta &= \sqrt{\delta_1\delta_2} \end{aligned}$$

*Proof.* The matrix can be written as

$$\left( I_2 \otimes I_n + (\alpha + i\beta) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes A \right) \left( I_2 \otimes I_n + (\alpha - i\beta) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes A \right).$$

Note that the both matrices are regular since their product is regular. For the same reason as in the previous proof, the following matrix is also regular

$$\left( I_2 \otimes I_n + (\alpha + i\beta) \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes A \right) \left( I_2 \otimes I_n + (\alpha - i\beta) \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes A \right),$$

and we may multiply the equation by this matrix obtaining  $\widehat{d}_1$  and  $\widehat{d}_2$ . Note that the four matrices commute, that is why we can write the whole product as

$$\begin{aligned} & \left( I_2 \otimes I_n + (\alpha + i\beta) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes A \right) \cdot \left( I_2 \otimes I_n + (\alpha + i\beta) \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes A \right) \cdot \\ & \left( I_2 \otimes I_n + (\alpha - i\beta) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes A \right) \cdot \left( I_2 \otimes I_n + (\alpha - i\beta) \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes A \right) = \\ & \left( I_2 \otimes I_n + 2(\alpha + i\beta) \begin{pmatrix} \gamma & 0 \\ 0 & \gamma \end{pmatrix} \otimes A + (\alpha + i\beta)^2 \begin{pmatrix} \gamma^2 + \delta^2 & 0 \\ 0 & \gamma^2 + \delta^2 \end{pmatrix} \otimes A^2 \right) \cdot \\ & \left( I_2 \otimes I_n + 2(\alpha - i\beta) \begin{pmatrix} \gamma & 0 \\ 0 & \gamma \end{pmatrix} \otimes A + (\alpha - i\beta)^2 \begin{pmatrix} \gamma^2 + \delta^2 & 0 \\ 0 & \gamma^2 + \delta^2 \end{pmatrix} \otimes A^2 \right) \end{aligned}$$

The product is a diagonal consisting of two  $n \times n$  blocks, which are the same. The block can be rewritten as product:

$$\begin{aligned} & (I_n + (\alpha + i\beta)(\gamma + i\delta)A) \cdot (I_n + (\alpha + i\beta)(\gamma - i\delta)A) \cdot \\ & (I_n + (\alpha - i\beta)(\gamma + i\delta)A) \cdot (I_n + (\alpha - i\beta)(\gamma - i\delta)A) \end{aligned}$$

and after reordering

$$\begin{aligned} & (I_n + (\alpha + i\beta)(\gamma + i\delta)A) \cdot (I_n + (\alpha - i\beta)(\gamma - i\delta)A) \cdot \\ & (I_n + (\alpha + i\beta)(\gamma - i\delta)A) \cdot (I_n + (\alpha - i\beta)(\gamma + i\delta)A) = \\ & (I_n + 2(\alpha\gamma - \beta\delta)A + (\alpha^2 + \beta^2)(\gamma^2 + \delta^2)A^2) \cdot \\ & (I_n + 2(\alpha\gamma + \beta\delta)A + (\alpha^2 + \beta^2)(\gamma^2 + \delta^2)A^2) \end{aligned}$$

Now it suffices to compare  $a_1 = \alpha\gamma - \beta\delta$  and verify that

$$\begin{aligned} b_1^2 &= (\alpha^2 + \beta^2)(\gamma^2 + \delta^2) - a_1^2 = \\ &= \alpha^2\gamma^2 + \beta^2\gamma^2 + \alpha^2\beta^2 + \beta^2\delta^2 - (\alpha\gamma)^2 + 2\alpha\beta\gamma\delta - (\beta\delta)^2 = \\ &= (\beta\gamma)^2 + (\alpha\beta)^2 + 2\alpha\beta\gamma\delta = \\ &= (\beta\gamma + \alpha\beta)^2 \end{aligned}$$

For  $b_2$  it is done in the same way.  $\square$

## THE ALGORITHM

Below we define three functions of which  $\text{vec}(Y) = \mathbf{solv1}(1, \text{vec}(\widehat{D}), i)$  provides the solution  $Y$ .  $X$  is then obtained as  $X = U^T Y \begin{pmatrix} i \\ \otimes V \end{pmatrix}$ .

**Synopsis.**

$F^T$  is  $m \times m$  lower quasi-triangular matrix. Let  $m_r$  be a number of real eigenvalues,  $m_c$  number of complex pairs. Thus  $m = m_r + 2m_c$ . Let  $F_j$  denote  $j$ -th diagonal block of  $F^T$  ( $1 \times 1$  or  $2 \times 2$  matrix) for  $j = 1, \dots, m_r + m_c$ . For a fixed  $j$ , let  $\bar{j}$  denote index of the first column of  $F_j$  in  $F^T$ . Whenever we write something like  $(I_{m^i} \otimes I_n + r {}^i F \otimes K)x = d$ ,  $x$  and  $d$  denote column vectors of appropriate dimensions, and  $x_{\bar{j}}$  is  $\bar{j}$ -th partition of  $x$ , and  $x_j = (x_{\bar{j}} \ x_{\bar{j}+1})^T$  if  $j$ -th eigenvalue is complex, and  $x_j = x_{\bar{j}}$  if  $j$ -th eigenvalue is real.

**Function solv1.**

The function  $x = \mathbf{solv1}(r, d, i)$  solves equation

$$(I_{m^i} \otimes I_n + r {}^i F \otimes K) x = d.$$

The function proceeds as follows:

If  $i = 0$ , the equation is solved directly,  $K$  is upper quasi-triangular matrix, so this is easy.

If  $i > 0$ , then we go through diagonal blocks  $F_j$  for  $j = 1, \dots, m_r + m_c$  and perform:

- (1) if  $F_j = (f_{\bar{j}\bar{j}}) = (f)$ , then we return  $x_j = \mathbf{solv1}(rf, d_{\bar{j}}, i - 1)$ . Then precalculate  $y = d_j - x_j$ , and eliminate guys below  $F_j$ . This is, for each  $k = \bar{j} + 1, \dots, m$ , we put

$$d_k = d_k - rf_{\bar{j}k} ({}^{i-1}F \otimes K) x_{\bar{j}} = d_k - \frac{f_{\bar{j}k}}{f} y$$

- (2) if  $F_j = \begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix}$ , we return  $x_j = \mathbf{solv2}(r\alpha, r\beta_1, r\beta_2, d_j, i - 1)$ . Then we precalculate

$$y = \left( \begin{pmatrix} \alpha & -\beta_1 \\ \beta_2 & \alpha \end{pmatrix} \otimes I_{m^{i-1}} \otimes I_n \right) \begin{pmatrix} d_{\bar{j}} - x_{\bar{j}} \\ d_{\bar{j}+1} - x_{\bar{j}+1} \end{pmatrix}$$

and eliminate guys below  $F_j$ . This is, for each  $k = \bar{j} + 2, \dots, n$  we put

$$\begin{aligned} d_k &= d_k - r(f_{\bar{j}k} \ f_{\bar{j}+1k}) \otimes ({}^{i-1}F \otimes K) x_j \\ &= d_k - \frac{1}{\alpha^2 + \beta_1\beta_2} ((f_{\bar{j}k} \ f_{\bar{j}+1k}) \otimes I_{m^{i-1}} \otimes I_n) y \end{aligned}$$

**Function solv2.**

The function  $x = \mathbf{solv2}(\alpha, \beta_1, \beta_2, d, i)$  solves equation

$$\left( I_2 \otimes I_{m^i} \otimes I_n + \begin{pmatrix} \alpha & \beta_1 \\ -\beta_2 & \alpha \end{pmatrix} \otimes {}^i F \otimes K \right) x = d$$

According to **Lemma 1** the function returns

$$x = \begin{pmatrix} \mathbf{solv2p}(\alpha, \beta_1\beta_2, \widehat{d}_1, i) \\ \mathbf{solv2p}(\alpha, \beta_1\beta_2, \widehat{d}_2, i) \end{pmatrix},$$

where  $\widehat{d}_1$ , and  $\widehat{d}_2$  are partitions of  $\widehat{d}$  from the lemma.

**Function solv2p.**

The function  $x = \text{solv2p}(\alpha, \beta^2, d, i)$  solves equation

$$(I_{m^i} \otimes I_n + 2\alpha {}^iF \otimes K + (\alpha^2 + \beta^2) {}^iF^2 \otimes K^2) x = d$$

The function proceeds as follows:

If  $i = 0$ , the matrix  $I_n + 2\alpha K + (\alpha^2 + \beta^2)K^2$  is calculated and the solution is obtained directly.

Now note that diagonal blocks of  $F^{2T}$  are of the form  $F_j^2$ , since if the  $F^T$  is block partitioned according to diagonal blocks, then it is lower triangular.

If  $i > 0$ , then we go through diagonal blocks  $F_j$  for  $j = 1, \dots, m_r + m_c$  and perform:

- (1) if  $F_j = (f_{\bar{j}\bar{j}}) = (f)$  then  $j$ -th diagonal block of

$$I_{m^i} \otimes I_n + 2\alpha {}^iF \otimes K + (\alpha^2 + \beta^2) {}^iF^2 \otimes K^2$$

takes the form

$$I_{m^{i-1}} \otimes I_n + 2\alpha f {}^{i-1}F \otimes K + (\alpha^2 + \beta^2) f^2 {}^{i-1}F^2 \otimes K^2$$

and we can put  $x_j = \text{solv2p}(f\alpha, f^2\beta^2, d_j, i-1)$ .

Then we need to eliminate guys below  $F_j$ . Note that  $|f^2| < |f|$ , therefore we precalculate  $y_2 = (\alpha^2 + \beta^2) f^2 ({}^{i-1}F^2 \otimes K^2) x_j$ , and then precalculate

$$y_1 = 2\alpha f ({}^{i-1}F \otimes K) x_j = d_j - x_j - y_2.$$

Let  $g_{pq}$  denote element of  $F^{2T}$  at position  $(q, p)$ . The elimination is done by going through  $k = \bar{j} + 1, \dots, m$  and putting

$$\begin{aligned} d_k &= d_k - (2\alpha f_{\bar{j}k} {}^{i-1}F \otimes K + (\alpha^2 + \beta^2) g_{\bar{j}k} {}^{i-1}F^2 \otimes K^2) x_j \\ &= d_k - \frac{f_{\bar{j}k}}{f} y_1 - \frac{g_{\bar{j}k}}{f^2} y_2 \end{aligned}$$

- (2) if  $F_j = \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix}$ , then  $j$ -th diagonal block of

$$I_{m^i} \otimes I_n + 2\alpha {}^iF \otimes K + (\alpha^2 + \beta^2) {}^iF^2 \otimes K^2$$

takes the form

$$I_{m^{i-1}} \otimes I_n + 2\alpha \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} {}^{i-1}F \otimes K + (\alpha^2 + \beta^2) \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix}^2 {}^{i-1}F^2 \otimes K^2$$

According to **Lemma 2**, we need to calculate  $\widehat{d}_{\bar{j}}$ , and  $\widehat{d}_{\bar{j}+1}$ , and  $a_1, b_1, a_2, b_2$ . Then we obtain

$$\begin{aligned} x_{\bar{j}} &= \text{solv2p}(a_1, b_1^2, \text{solv2p}(a_2, b_2^2, \widehat{d}_{\bar{j}}, i-1), i-1) \\ x_{\bar{j}+1} &= \text{solv2p}(a_1, b_1^2, \text{solv2p}(a_2, b_2^2, \widehat{d}_{\bar{j}+1}, i-1), i-1) \end{aligned}$$

Now we need to eliminate guys below  $F_j$ . Since  $\|F_j^2\| < \|F_j\|$ , we precalculate

$$\begin{aligned} y_2 &= (\alpha^2 + \beta^2)(\gamma^2 + \delta^2) (I_2 \otimes {}^{i-1}F^2 \otimes K^2) x_j \\ y_1 &= 2\alpha(\gamma^2 + \delta^2) (I_2 \otimes {}^{i-1}F \otimes K) x_j \\ &= (\gamma^2 + \delta^2) (F_j^{-1} \otimes I_{m^{i-1}n}) \left( d_j - x_j - \frac{1}{(\gamma^2 + \delta^2)} (F_j^2 \otimes I_{m^{i-1}n}) y_2 \right) \\ &= \left( \begin{pmatrix} \gamma & -\delta_1 \\ \delta_2 & \gamma \end{pmatrix} \otimes I_{m^{i-1}n} \right) (d_j - x_j) - \left( \begin{pmatrix} \gamma & \delta_1 \\ -\delta_2 & \gamma \end{pmatrix} \otimes I_{m^{i-1}n} \right) y_2 \end{aligned}$$

Then we go through all  $k = \bar{j} + 2, \dots, m$ . For clearer formulas, let  $\mathbf{f}_k$  denote pair of  $F^T$  elements in  $k$ -th line below  $F_j$ , this is  $\mathbf{f}_k = (f_{\bar{j}k} \ f_{\bar{j}+1k})$ . And let  $\mathbf{g}_k$  denote the same for  $F^{2T}$ , this is  $\mathbf{g}_k = (g_{\bar{j}k} \ g_{\bar{j}+1k})$ . For each  $k$  we put

$$\begin{aligned} d_k &= d_k - (2\alpha\mathbf{f}_k \otimes {}^{i-1}F \otimes K + (\alpha^2 + \beta^2)\mathbf{g}_k \otimes {}^{i-1}F^2 \otimes K^2) x_j \\ &= d_k - \frac{1}{\gamma^2 + \delta^2} (\mathbf{f}_k \otimes I_{m^{i-1}n}) y_1 - \frac{1}{\gamma^2 + \delta^2} (\mathbf{g}_k \otimes I_{m^{i-1}n}) y_2 \end{aligned}$$

#### FINAL NOTES

##### Numerical Issues of $A^{-1}B$ .

We began the solution of the Sylvester equation with multiplication by  $A^{-1}$ . This can introduce numerical errors, and we need more numerically stable supplement. Its aim is to make  $A$  and  $B$  commutative, this is we need to find a regular matrix  $P$ , such that  $(PA)(PB) = (PB)(PA)$ . Recall that this is necessary in solution of

$$(I_2 \otimes I_{m^i} \otimes (PA) + (D + C) \otimes {}^iF \otimes (PB))x = d,$$

since this equation is multiplied by  $I_2 \otimes I_{m^i} \otimes (PA) + (D - C) \otimes {}^iF \otimes PB$ , and the diagonal result

$$I_2 \otimes I_{m^i} \otimes (PAPA) + 2D \otimes {}^iF \otimes (PAPB) + (D^2 - C^2) \otimes {}^iF^2 \otimes (PBPB)$$

is obtained only if  $(PA)(PB) = (PB)(PA)$ .

Finding regular solution of  $(PA)(PB) = (PB)(PA)$  is equivalent to finding regular solution of  $APB - BPA = 0$ . Numerical error of the former equation is  $P$ -times greater than the numerical error of the latter equation. And the numerical error of the latter equation also grows with the size of  $P$ . On the other hand, truncation error in  $P$  multiplication decreases with growing the size of  $P$ . By intuition, stability analysis will show that the best choice is some orthonormal  $P$ .

Obviously, since  $A$  is regular, the equation  $APB - BPA = 0$  has solution of the form  $P = \alpha A^{-1}$ , where  $\alpha \neq 0$ . There is a vector space of all solutions  $P$  (including singular ones). In precise arithmetics, its dimension is  $\sum n_i^2$ , where  $n_i$  is number of repetitions of  $i$ -th eigenvalue of  $A^{-1}B$  which is similar to  $BA^{-1}$ . In floating point arithmetics, without any further knowledge about  $A$ , and  $B$ , we are only sure about dimension  $n$  which is implied by similarity of  $A^{-1}B$  and  $BA^{-1}$ . Now we try to find the base of the vector space of solutions.

Let  $L$  denote the following linear operator:

$$L(X) = (AXB - BXA)^T.$$

Let  $\text{vec}(X)$  denote a vector made by stacking all the columns of  $X$ . Let  $T_n$  denote  $n^2 \times n^2$  matrix representing operator  $\text{vec}(X) \mapsto \text{vec}(X^T)$ . And finally let  $M$  denote  $n^2 \times n^2$  matrix representing the operator  $L$ . It is not difficult to verify that:

$$M = T_n(B^T \otimes A - A^T \otimes B)$$

Now we show that  $M$  is skew symmetric. Recall that  $T_n(X \otimes Y) = (Y \otimes X)T_n$ , we have:

$$M^T = (B^T \otimes A - A^T \otimes B)^T T_n = (B \otimes A^T - A \otimes B^T) T_n = T_n(A^T \otimes B - B^T \otimes A) = -M$$

We try to solve  $M \text{vec}(X) = T_n(0) = 0$ . Since  $M$  is skew symmetric, there is real orthonormal matrix  $Q$ , such that  $M = Q\widehat{M}Q^T$ , and  $\widehat{M}$  is block diagonal matrix consisting of  $2 \times 2$  blocks of the form

$$\begin{pmatrix} 0 & \alpha_i \\ -\alpha_i & 0 \end{pmatrix},$$

and of additional zero, if  $n^2$  is odd.

Now we solve equation  $\widehat{M}y = 0$ , where  $y = Q^T \text{vec}(X)$ . Now there are  $n$  zero rows in  $\widehat{M}$  coming from similarity of  $A^{-1}B$  and  $BA^{-1}$  (structural zeros). Note that the additional zero for odd  $n^2$  is already included in that number, since for odd  $n^2$  is  $n^2 - n$  even. Besides those, there are also zeros (esp. in floating point arithmetics), coming from repetitive (or close) eigenvalues of  $A^{-1}B$ . If we are able to select the rows with the structural zeros, a solution is obtained by picking arbitrary numbers for the same positions in  $y$ , and put  $\text{vec}(X) = Qy$ .

The following questions need to be answered:

- (1) How to recognize the structural rows?
- (2) Is  $A^{-1}$  generated by a  $y$ , which has non-zero elements only on structural rows? Note that  $A$  can have repetitive eigenvalues. The positive answer to the question implies that in each  $n$ -partition of  $y$  there is exactly one structural row.
- (3) And very difficult one: How to pick  $y$  so that  $X$  would be regular, or even close to orthonormal (pure orthonormality overdeterminates  $y$ )?

### Making Zeros in $F$ .

It is clear that the numerical complexity of the proposed algorithm strongly depends on a number of non-zero elements in the Schur factor  $F$ . If we were able to find  $P$ , such that  $PF P^{-1}$  has substantially less zeros than  $F$ , then the computation would be substantially faster. However, it seems that we have to pay price for any additional zero in terms of less numerical stability of  $PF P^{-1}$  multiplication. Consider  $P$ , and  $F$  in form

$$P = \begin{pmatrix} I & X \\ 0 & I \end{pmatrix}, \quad F = \begin{pmatrix} A & C \\ 0 & B \end{pmatrix}$$

we obtain

$$PFP^{-1} = \begin{pmatrix} A & C + XB - AX \\ 0 & B \end{pmatrix}$$

Thus, we need to solve  $C = AX - XB$ . It's clear that numerical stability of operator  $Y \mapsto PYP^{-1}$  and its inverse  $Y \mapsto P^{-1}YP$  is worse with growing norm  $\|X\|$ . The norm can be as large as  $\|F\|/\delta$ , where  $\delta$  is a distance of eigenspectra of  $A$  and  $B$ . Also, a numerical error of the solution is proportional to  $\|C\|/\delta$ .

Although, these difficulties cannot be overcome completely, we may introduce an algorithm, which works on  $F$  with ordered eigenvalues on diagonal, and seeks such partitioning to maximize  $\delta$  and minimize  $C$ . If the partitioning is found, the algorithm finds  $P$  and then is run for  $A$  and  $B$  blocks. It stops when further partitioning is not possible without breaking some user given limit for numerical errors. We have to keep in mind that the numerical errors are accumulated in product of all  $P$ 's of every step.

### Exploiting constant rows in $F$ .

If some of  $F$ 's rows consists of the same numbers, or a number of distinct values within a row is small, then this structure can be easily exploited in the algorithm. Recall, that in both functions **solv1**, and **solv2p**, we eliminate guys below diagonal element (or block) (of  $F^T$ ), by multiplying solution of the diagonal and cancelling it from right side. If the elements below the diagonal block are the same, we save one vector multiplication. Note that in **solv2p** we still need to multiply by elements below diagonal of the matrix  $F^{T^2}$ , which obviously has not the property. However, the heaviest elimination is done at the very top level, in the first call to **solv1**.

Another way of exploitation the property is to proceed all calculations in complex numbers. In that case, only **solv1** is run.

How the structure can be introduced into the matrix? Following the same notation as in previous section, we solve  $C = AX - XB$  in order to obtain zeros at place of  $C$ . If it is not possible, we may relax the equation by solving  $C - R = AX - XB$ , where  $R$  is suitable matrix with constant rows. The matrix  $R$  minimizes  $\|C - R\|$  in order to minimize  $\|X\|$  if  $A$ , and  $B$  are given. Now, in the next step we need to introduce zeros (or constant rows) to matrix  $A$ , so we seek for regular matrix  $P$ , doing the job. If found, the product looks like:

$$\begin{pmatrix} P & 0 \\ 0 & I \end{pmatrix} \begin{pmatrix} A & R \\ 0 & B \end{pmatrix} \begin{pmatrix} P^{-1} & 0 \\ 0 & I \end{pmatrix} = \begin{pmatrix} PAP^{-1} & PR \\ 0 & B \end{pmatrix}$$

Now note, that matrix  $PR$  has also constant rows. Thus, preconditioning of the matrix in upper left corner doesn't affect the property. However, a preconditioning of the matrix in lower right corner breaks the property, since we would obtain  $RP^{-1}$ .