

The doubly stochastic advection-diffusion-decay model: numerical scheme

Michael Tsyrlunikov and Alexander Rakitko

(michael.tsyrlunikov@gmail.com)

April 14, 2019

1 Model

$$\frac{\partial \xi(t, x)}{\partial t} + U(t, x) \frac{\partial \xi(t, x)}{\partial x} + \rho(t, x) \xi(t, x) - \nu(t, x) \frac{\partial^2 \xi(t, x)}{\partial x^2} = \sigma(t, x) \alpha(t, x). \quad (1)$$

2 Discretization

In an implicit scheme, in which the time derivative $\partial \xi / \partial t$ is approximated as $(\xi_k - \xi_{k-1}) / \Delta t$, all other occurrences of ξ (that is, $\partial \xi / \partial x \xi$, and $\partial^2 \xi / \partial x^2$) are evaluated at the current time step k . We also choose to evaluate the coefficients at the current time step k . This leads to the finite-difference scheme

$$\frac{\boldsymbol{\xi}^k - \boldsymbol{\xi}^{k-1}}{\Delta t} + \text{diag}(\mathbf{U}^k) \mathbf{D}^{upwind} \boldsymbol{\xi}^k + \text{diag}(\boldsymbol{\rho}^k) \boldsymbol{\xi}^k - \text{diag}(\boldsymbol{\nu}^k) \mathbf{D}_2 \boldsymbol{\xi}^k = \text{diag}(\boldsymbol{\sigma}^k) \boldsymbol{\alpha}^k, \quad (2)$$

where the superscript k labels the time step, \mathbf{D}^{upwind} is the upwind non-centered (directed) finite difference operator (matrix) (upwind means that the second point in the stencil besides the central point is the point towards the direction from which the wind blows, see subsection 2.1 below), \mathbf{D}_2 is the usual three-point approximation of the Laplacian in 1D (a matrix) such that the action of \mathbf{D}_2 on the vector \mathbf{y} returns

$$(\mathbf{D}_2 \mathbf{y})_j = \frac{y_{j-1} - 2y_j + y_{j+1}}{(\Delta x)^2} \quad (3)$$

(where the subscript j labels the spatial grid points), $\boldsymbol{\alpha}_j^k = \frac{1}{\sqrt{\Delta t \Delta x}} N(0, 1)$ is the time and space discretized white noise, and $\text{diag}(\mathbf{z})$ is the diagonal matrix with the vector \mathbf{z} on its main diagonal.

2.1 Upwind finite difference operator

If $U_j > 0$, then the j th row of the matrix \mathbf{D}^{upwind} has the following two non-zero entries:

$$\mathbf{D}_{j,j-1}^{upwind} = -\frac{1}{\Delta x} \quad (4)$$

$$\mathbf{D}_{j,j}^{upwind} = +\frac{1}{\Delta x} \quad (5)$$

If $U_j < 0$, then the j th row of the matrix \mathbf{D}^{upwind} has the following two non-zero entries:

$$\begin{aligned} \mathbf{D}_{j,j}^{upwind} &= -\frac{1}{\Delta x} \\ \mathbf{D}_{j,j+1}^{upwind} &= +\frac{1}{\Delta x} \end{aligned}$$

3 Time stepping

We rewrite Eq.(2) as

$$\mathbf{G}_k \cdot \boldsymbol{\xi}^k = \boldsymbol{\xi}^{k-1} + \Delta t \operatorname{diag}(\boldsymbol{\sigma}^k) \cdot \boldsymbol{\alpha}^k, \quad (6)$$

where

$$\boxed{\mathbf{G}_k = \mathbf{I} + \Delta t \operatorname{diag}(\mathbf{U}^k) \cdot \mathbf{D}^{upwind} + \Delta t \operatorname{diag}(\boldsymbol{\rho}^k) - \Delta t \operatorname{diag}(\boldsymbol{\nu}^k) \cdot \mathbf{D}_2.} \quad (7)$$

So,

$$\boxed{\boldsymbol{\xi}^k = \mathbf{F}_k [\boldsymbol{\xi}^{k-1} + \Delta t \operatorname{diag}(\boldsymbol{\sigma}^k) \cdot \boldsymbol{\alpha}^k],} \quad (8)$$

where

$$\boxed{\mathbf{F}_k = \mathbf{G}_k^{-1}.} \quad (9)$$

4 Technical details of the solution of the system of linear algebraic equations with the cyclic tridiagonal matrix

4.1 The system

Let

$$\mathbf{M} := \frac{1}{\Delta t} \mathbf{G} \quad (10)$$

By $h = \Delta s$ denote the mesh size.

Conventions.

- (i) By $i+1$ we always mean $i+1$ if $i < n$ and 1 if $i = n$.
- (ii) Likewise, by $i-1$ we mean $i-1$ if $i > 1$ and n if $i = 1$.
- (iii) We drop the time index in this section.
- (iv) By $\mathbf{a}, \mathbf{b}, \mathbf{c}$ we denote the lower, main, and upper (sub-)diagonals of \mathbf{M} , respectively.
- (v) We also often drop the row index if this does not cause an ambiguity.

Then, in any row, we have (NB: the indexes are dropped!):

- If $U > 0$:

$$a = -\frac{U}{h} - \frac{\nu}{h^2} \quad (11)$$

$$b = \frac{1}{\Delta t} + \frac{U}{h} + \rho + 2\frac{\nu}{h^2} \quad (12)$$

$$c = -\frac{\nu}{h^2} \quad (13)$$

- If $U \leq 0$:

$$a = -\frac{\nu}{h^2} \quad (14)$$

$$b = \frac{1}{\Delta t} - \frac{U}{h} + \rho + 2\frac{\nu}{h^2} \quad (15)$$

$$c = \frac{U}{h} - \frac{\nu}{h^2} \quad (16)$$

The diagonal dominance of \mathbf{M} is easily seen.

Then, the system to be solved is

$$\boxed{\mathbf{M} \cdot \mathbf{x} = \mathbf{y}} \quad (17)$$

where

$$\mathbf{x} := \boldsymbol{\xi}^k \quad (18)$$

and

$$\mathbf{y} := \frac{\boldsymbol{\xi}^{k-1}}{\Delta t} + \boldsymbol{\sigma}^k \cdot \boldsymbol{\alpha}^k \quad (19)$$

and

$$\boldsymbol{\alpha}^k = \frac{1}{\sqrt{\Delta t \Delta s}} \boldsymbol{\sigma}^k \cdot \mathbf{N}(\mathbf{0}, \mathbf{I}). \quad (20)$$

4.2 Solution

4.2.1 Single forecast

To solve the system of linear algebraic equations with the *cyclic tridiagonal matrix* we reduce it to the system with the ordinary *tridiagonal matrix* (i.e. without the “corners”). The latter (banded) system is then solved using the function `Solve.tridiag` from the R package `limSolve` (Soetaert et al., 2009).

The reduction of the *cyclic tridiagonal matrix* to the *non-cyclic tridiagonal matrix* is performed following Ahlberg et al. (2016, p.15). Specifically,

1. From the matrix \mathbf{M} , we remove the first row and the first column. The remainder is the $(n-1) \times (n-1)$ truly banded tridiagonal matrix, let us denote it $\widetilde{\mathbf{M}}$.

2. Then, we remove the first element of the r.h.s. \mathbf{y} and denote the rest by $\tilde{\mathbf{y}}$, so that

$$\mathbf{y}^\top = (y_1, \tilde{\mathbf{y}}^\top). \quad (21)$$

3. Correspondingly, we remove the first element of the solution \mathbf{x} , denote the rest by $\tilde{\mathbf{x}}$, and denote x_1 by χ , so that

$$\mathbf{x}^\top = (\chi, \tilde{\mathbf{x}}^\top). \quad (22)$$

As a result, we obtain

$$\begin{pmatrix} b_1 | & c_1 & 0 & \dots & 0 & a_1 \\ a_2 | & b_2 & c_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ c_n | & 0 & 0 & \dots & a_n & b_n \end{pmatrix} \cdot \begin{pmatrix} \chi \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad (23)$$

which is regarded as the system of the two equations. The first equation is

$$\boxed{b_1\chi + c_1x_2 + a_1x_n = y_1.} \quad (24)$$

The second equation is

$$\widetilde{\mathbf{M}}\tilde{\mathbf{x}} + \boldsymbol{\varphi}\chi = \tilde{\mathbf{y}}, \quad (25)$$

where

$$\boldsymbol{\varphi} := (a_2, 0, \dots, 0, c_n)^\top. \quad (26)$$

Let us rewrite the latter equation as

$$\boxed{\widetilde{\mathbf{M}}\tilde{\mathbf{x}} = \tilde{\mathbf{y}} - \boldsymbol{\varphi}\chi} \quad (27)$$

and note that the solution to *linear* Eq.(27) can be found by applying the superposition principle. We write down two systems (with the same non-cyclic tridiagonal matrix $\widetilde{\mathbf{M}}$) in which the right-hand sides come from the two components of the r.h.s. in Eq.(27) as follows.

System 1 is

$$\widetilde{\mathbf{M}}\mathbf{u} = \tilde{\mathbf{y}}. \quad (28)$$

System 2 is

$$\widetilde{\mathbf{M}}\mathbf{v} = -\boldsymbol{\varphi}. \quad (29)$$

Note that \mathbf{u} and \mathbf{v} can be explicitly found (e.g. using function `Solve.tridiag`). Then, it is easily verified that the solution to the system Eq.(27) is

$$\boxed{\tilde{\mathbf{x}} = \mathbf{u} + \mathbf{v} \cdot \chi} \quad (30)$$

Thus, due to Eq.(30), the vector $\tilde{\mathbf{x}}$ is known up to the one single unknown, χ . To find χ , we use Eq.(24). Specifically, from Eq.(30), we express

$$x_2 = u_1 + \chi \cdot v_1 \quad (31)$$

and

$$x_n = u_{n-1} + \chi \cdot v_{n-1} \quad (32)$$

(where, we recall, all elements of the vectors \mathbf{u} and \mathbf{v} are known at this point), so that Eq.(24) becomes

$$b_1\chi + c_1(u_1 + \chi \cdot v_1) + a_1(u_{n-1} + \chi \cdot v_{n-1}) = y_1, \quad (33)$$

whence

$$(b_1 + c_1v_1 + a_1v_{n-1}) \cdot \chi + c_1u_1 + a_1u_{n-1} = y_1 \quad (34)$$

and so

$$\chi = \frac{y_1 - c_1u_1 - a_1u_{n-1}}{b_1 + c_1v_1 + a_1v_{n-1}} \quad (35)$$

Having found χ , we substitute it into Eq.(30), getting $\tilde{\mathbf{x}}$, and form the final solution \mathbf{x} by concatenating χ and $\tilde{\mathbf{x}}$ using Eq.(22).

4.2.2 Ensemble of forecasts

If we wish to generate an ensemble of N members that share the secondary (coefficient) fields, then, in the above notation, \mathbf{x} and \mathbf{y} become $n \times N$ matrices, \mathbf{X} and \mathbf{Y} , respectively.

Correspondingly, $\tilde{\mathbf{x}}$, $\tilde{\mathbf{y}}$, \mathbf{u} become $(n-1) \times N$ matrices, $\tilde{\mathbf{X}}$, $\tilde{\mathbf{Y}}$, \mathbf{U} respectively. The vector \mathbf{v} remains a vector. The vector $\boldsymbol{\varphi}$ also remains unchanged.

The scalar χ becomes an N -vector, $\boldsymbol{\chi}$.

Computationally:

(1) System Eq.(28) is straightforwardly solved using `Solve.tridiag` with N right-hand sides,

$$\tilde{\mathbf{M}}\mathbf{U} = \tilde{\mathbf{Y}}. \quad (36)$$

(2) Eq.(29) is still solved only once.

(3) Eq.(27) becomes

$$\tilde{\mathbf{M}}\tilde{\mathbf{X}} = \tilde{\mathbf{Y}} - \boldsymbol{\varphi} \cdot \boldsymbol{\chi}^\top \quad (37)$$

(4) Eq.(30) becomes

$$\tilde{\mathbf{X}} = \mathbf{U} + \mathbf{v} \cdot \boldsymbol{\chi}^\top \quad (38)$$

(5) Eq.(24) becomes an element-wise equation, in which χ, x_2, x_n, y_1 are N -vectors. Correspondingly, Eq.(35) is now understood as a element-wise expression for the N -vectors $\mathbf{y}_1, \mathbf{u}_1, \mathbf{u}_{n-1}$:

$$\boldsymbol{\chi} = \frac{\mathbf{y}_1 - c_1\mathbf{u}_1 - a_1\mathbf{u}_{n-1}}{b_1 + c_1v_1 + a_1v_{n-1}} \quad (39)$$

5 Stability

Statement. *If the coefficient fields ρ and ν are positive, then the numerical scheme Eqs.(6)–(9) is **unconditionally stable** (i.e. stable for any time step Δt).*

To prove this statement, we note that if $\mathbf{F}_k = \mathbf{F} = \text{const}$, then the stability is equivalent to the *boundedness of norms of powers of the transition operator \mathbf{F}_k* (Godunov and Ryabenki, 1964). If $\mathbf{F}_k \neq \text{const}$, then it is sufficient to prove that $\|\mathbf{F}_1 \cdot \mathbf{F}_2 \cdots \mathbf{F}_N\| < C$, where N is the total number of time steps and $C > 0$ is a constant. The latter condition is, in turn, satisfied if $\forall k$,

$$\boxed{\|\mathbf{F}_k\| < 1}. \quad (40)$$

Now, we show that this latter condition is indeed met in our case.

To this end, let us introduce the vector norm for the spatial vector $\boldsymbol{\xi}$. Let us select the *maximal* (or l_∞) norm:

$$\|\boldsymbol{\xi}\|_\infty := \max_i |\xi_i|. \quad (41)$$

Then, we show that $\forall \boldsymbol{\xi}$,

$$\|\mathbf{G}_k \boldsymbol{\xi}\|_\infty > \|\boldsymbol{\xi}\|_\infty. \quad (42)$$

This condition would entail that the condition Eq.(40) because if we denote $\boldsymbol{\eta} := \mathbf{G}_k \boldsymbol{\xi}$ and note that $\mathbf{F}_k = \mathbf{G}_k^{-1}$. Indeed. Eq.(42) implies that $\forall \boldsymbol{\eta}$,

$$\|\mathbf{F}_k \boldsymbol{\eta}\|_\infty < \|\boldsymbol{\eta}\|_\infty, \quad (43)$$

which, by definition of the operator norm, would mean that

$$\|\mathbf{F}_k\|_\infty \leq 1. \quad (44)$$

So, to prove that the scheme is unconditionally stable, it remains to prove that the condition Eq.(42) is satisfied with our numerical scheme Eqs.(6)–(9).

Take a vector $\boldsymbol{\xi}$. Find its maximal (in modulus) element, let it be ξ_I . We have $\forall i$

$$\sum_j G_{ij} \xi_j = G_{i,i-1} \xi_{i-1} + G_{i,i} \xi_i + G_{i,i+1} \xi_{i+1} \quad (45)$$

(the time index is dropped till the end of the proof).

In Eq.(45), we have in case $U_i > 0$:

$$G_{i,i-1} = -\frac{\Delta t}{\Delta x} U_i - \frac{\Delta t}{(\Delta x)^2} \nu_i \quad (46)$$

$$G_{i,i} = 1 + \frac{\Delta t}{\Delta x} U_i + \Delta t \rho + 2 \frac{\Delta t}{(\Delta x)^2} \nu_i \quad (47)$$

$$G_{i,i+1} = -\frac{\Delta t}{(\Delta x)^2} \nu_i \quad (48)$$

For $U_i \leq 0$, we have similarly

$$G_{i,i-1} = -\frac{\Delta t}{(\Delta x)^2} \nu_i \quad (49)$$

$$G_{i,i} = 1 - \frac{\Delta t}{\Delta x} U_i + \Delta t \rho + 2 \frac{\Delta t}{(\Delta x)^2} \nu_i \quad (50)$$

$$G_{i,i+1} = \frac{\Delta t}{\Delta x} U_i - \frac{\Delta t}{(\Delta x)^2} \nu_i \quad (51)$$

We have

$$\|\mathbf{G}\boldsymbol{\xi}\|_\infty = \max_i \left| \sum_j G_{ij} \xi_j \right| \geq \left| \sum_j G_{Ij} \xi_j \right|, \quad (52)$$

where I (we recall) corresponds to the maximal $|\xi_i|$. Here

$$\left| \sum_j G_{Ij} \xi_j \right| = |G_{i,i-1} \xi_{I-1} + G_{i,i} \xi_I + G_{i,i+1} \xi_{I+1}| \quad (53)$$

For concreteness and without loss of generality, let $U_I > 0$ and $\xi_I > 0$. Then, from Eqs.(48)–(48), $G_{I,I} > 0$, $G_{I,I-1} < 0$, $G_{I,I+1} < 0$, and $|\xi_{I+1}| \leq |\xi_I|$ and $|\xi_{I-1}| \leq |\xi_I|$. This implies that the minimal value of $|\sum_j G_{Ij} \xi_j|$ is attained if $\xi_{I-1} > 0$ and $\xi_{I+1} > 0$. Therefore, finally,

$$\|\mathbf{G}\boldsymbol{\xi}\|_\infty \geq \left| \sum_j G_{Ij} \xi_j \right| \geq (G_{i,i-1} + G_{i,i} + G_{i,i+1}) \cdot |\xi_I| = (1 + \Delta t \rho) \cdot \|\boldsymbol{\xi}\|_\infty \geq \|\boldsymbol{\xi}\|_\infty, \quad (54)$$

QED.

As a result, we have proven that **the scheme is unconditionally stable if $\rho > 0$ and $\nu > 0$.**

If the condition of positivity of ρ or ν is violated, then the scheme can become unstable, and this is intentional (the intermittent instability is an important feature of the DSADM).

6 Model error

By model error in Eq.(1) we mean its r.h.s. $\sigma(t, x)\alpha(t, x)$. Rewrite Eq.(8) as

$$\boldsymbol{\xi}^k = \mathbf{F}\boldsymbol{\xi}^{k-1} + \Delta t \mathbf{F} \text{diag}(\boldsymbol{\sigma}^k) \cdot \boldsymbol{\alpha}^k, \quad (55)$$

so that the time discrete model error is

$$\boxed{\boldsymbol{\varepsilon}_k = \Delta t \cdot \mathbf{F} \cdot (\text{diag}(\boldsymbol{\sigma}^k) \cdot \boldsymbol{\alpha}^k)} \quad (56)$$

Note the application of \mathbf{F} .

Now, find the model error covariance matrix. Since

$$\text{Var} \boldsymbol{\alpha}_i^k = \frac{1}{\Delta t \cdot \Delta x}$$

we obtain

$$\boxed{\mathbf{Q}_k = \frac{\Delta t}{\Delta x} \cdot \mathbf{F} \cdot (\text{diag}(\boldsymbol{\sigma}^k))^2 \cdot \mathbf{F}^\top} \quad (57)$$

From this equation, it follows that when computing $\mathbf{B}^k = \mathbf{F}\mathbf{A}^{k-1}\mathbf{F}^\top + \mathbf{Q}^k$, it suffices to compute

$$\mathbf{A}_+^{k-1} = \mathbf{A}^{k-1} + \frac{\Delta t}{\Delta x} \cdot (\text{diag}(\boldsymbol{\sigma}^k))^2 \quad (58)$$

and then compute

$$\mathbf{B}^k = \mathbf{F}\mathbf{A}_+^{k-1}\mathbf{F}^\top. \quad (59)$$

In other words, we should add \mathbf{Q}^k to \mathbf{A}^{k-1} — rather than to $\mathbf{F}\mathbf{A}^{k-1}\mathbf{F}^\top$ (i.e. *before* the forecast rather than *after* the forecast).

References

- J. H. Ahlberg, E. N. Nilson, and J. L. Walsh. *The Theory of Splines and Their Applications*, volume 38. Elsevier, 2016.
- S. K. Godunov and V. S. Ryabenki. *Difference schemes*. Amsterdam: North Holland, 1964.
- K. Soetaert, K. Van den Meersche, and D. van Oevelen. limsolve: Solving linear inverse models. *R package version*, 1(1), 2009. URL <https://CRAN.R-project.org/package=limSolve>.