# Discretizing Linear and Affine Operators Overview

## 1  Overview of Notation

Some general notation, independent of operators and discretization

- For a given variable $q$, define the notation $q^- \equiv \min\{q, 0\}$ and $q^+ \equiv \max\{q, 0\}$, which will be useful for defining finite-differences with an upwind scheme. This can apply to vectors as well. For example, $q_m^- = q_m$ if $q_m < 0$ and 0 if $q_m > 0$, and $q^- \equiv \{q_m^-\}_{m=1}^M$.

- Let $W_t$ be the Wiener process with the integral defined by the Ito interpretation

- Derivatives are denoted by the operator $\partial$ and univariate derivatives such as $\partial_x \tilde{u}(x) \equiv u'(x)$.

- We will denote continuous functions, prior to discretization, like $\tilde{u}(x)$. The discretization of $\tilde{u}(x)$ on $x \in \mathbb{R}^M$ is denoted $u \in \mathbb{R}^M$ and the discretization of $\tilde{u}(x)$ on $\bar{x}$ is $\bar{u} \in \mathbb{R}^{\bar{M}}$.

- Some special matrices to help in the composition notation:

  - $\mathbf{I}_N$ is the $N \times N$ identity matrix. Always drop the subscript when the dimensions are unambiguous, as it would be the same in the code
  - $\mathbf{0}_N$ is the column vector of $N$ 0s, and $\mathbf{0}_N^\top$ a row vector
  - $\mathbf{0}_{N \times M}$ is the $N \times M$ matrix of 0s

- An *affine* operator $A$ can be decomposed into a *linear operator*, denoted $A_L$, and a *bias* denoted $A_b$, such that for all $x$ in the domain,

$$Ax = A_L x + A_b$$

- In the case of an affine operator on a discrete space, $A : \mathbb{R}^M \to \mathbb{R}^N$,

  - We can decompose this into the linear operator $A_L \in \mathbb{R}^{N \times M}$ and bias vector $A_b \in \mathbb{R}^N$ such that for all $x \in \mathbb{R}^M$, $Ax = A_L x + A_b$

The purpose of these notes is to discretize affine or linear operators with finite differences.[1] To set some notation and definitions on operators,

- Denote a typical affine operator on the space of continuous functions as $\tilde{A}$, and $\tilde{L}$ for a linear operator. When these are discretized on a particular grid, denote them as $A$ and $L$ accordingly.

- The baseline domain of the operator is on $[x^{\min}, x^{\max}]$.

---

[1]These are often the infinitesimal generator of a stochastic process. See https://en.wikipedia.org/wiki/Infinitesimal_generator_(stochastic_processes) for some formulas and interpretation for diffusions, and https://en.wikipedia.org/wiki/Transition_rate_matrix

- Form a grid on the domain with $M$ points, $\{x_m\}_{m=1}^M$ with $x_1 = x^{\min}$ and $x_M = x^{\max}$ when. After discretizing, we can sometimes denote the grid with the variable name, i.e. $x \equiv \{x_m\}_{m=1}^M$. In the simple case of a uniform grid, $\Delta \equiv x_{m+1} - x_m$ for all $m < M$.

- A core part of the discretization process will be to expand the variable onto the *extension* (i.e. including any boundary points required for the boundary conditions). The set of boundary points will be referred to as $S_E$. If there are $M$ points in the grid, and $M_E$ points required for the boundary conditions (i.e. $|S_E| = M_E$), then define $\bar{M} = M + M_E$ as the total set of points on the extended domain. We will denote the extended domain as $\bar{x} \in \mathbb{R}^{\bar{M}}$.

- For any arbitrary continuous function $\tilde{y}(x)$ defined in the whole space of $x$, we define $\bar{y}$ as its discretization on the whole domain of $x$ and $y$ as the discretization only for interior points of the domain. So while $\bar{y}$ has length $\bar{M}$, $y$ has length $M$.

## 2 General Overview of Discretization and Boundary Values

### 2.1 Simple Differential Operators

Take a simple linear or affine differential operator $\tilde{A}$, (possibly affine) boundary conditions $\tilde{B}$, boundary value function $\tilde{b}(x)$ and the function of interest $\tilde{u}(x)$. The general problem to solve is to find the $\tilde{u}(x)$ such that.

$$\tilde{A}\tilde{u}(x) = 0 \tag{1}$$
$$\tilde{B}\tilde{u}(x) = \tilde{b}(x) \tag{2}$$

For linear $\tilde{A}$, we will denote it as $\tilde{L}$ to emphasize the fact that it's not affine. The discretization process generates the following objects:

- $B \in \mathbb{R}^{M_E \times \bar{M}}$ is the (possibly affine) *boundary condition operator*, which satisfies the equation

$$B\bar{u} = \mathbf{0}_{M_E} \tag{3}$$

   for any $\bar{u}$ in the space of functions that satisfy the discretized boundary conditions. [2] For affine $B$, we can write out (3) as

$$B_L\bar{u} = -B_b \tag{4}$$

- $R \in \mathbb{R}^{M \times \bar{M}}$ is the linear *restriction operator* which is defined by the domain. It removes columns which are not in the interior. It fulfills

$$R\bar{u} = u \tag{5}$$

- $Q^B : \mathbb{R}^M \to \mathbb{R}^{\bar{M}}$ is the (potentially affine) *boundary extrapolation operator* associated with $B$. The operator $Q^B$ is defined as fulfilling the following relationships (keeping in mind that $Q^B$ is affine and $R$ is linear)

$$Q^B R\bar{u} = \bar{u} \tag{6}$$
$$BQ^B u = \mathbf{0}_{M_E} \tag{7}$$

---

[2]Notice that $B$ is not necessarily unique. The choice of $B$ is exactly the choice of boundary value discretization. For instance, choosing to do first or second order Neumann border conditions is simply the choice of the operator $B$.

To give intuition, for any $\bar{u}$ that satisfies the border conditions:[3] (6) says that finding the restriction of the function and then extrapolating to extension yields the same function, and (7) says that the boundary extrapolation of the interior of the function, $u$, fulfills the boundary value.

- $A : \mathbb{R}^{\bar{M}} \to \mathbb{R}^{M}$ is the (possibly affine) *stencil operator*. It maps the extended domain to the interior by applying a stencil, which is determined by the derivative operator and the numerical differentiation scheme. As with the continous case, we will denote the stencil operator as $L$ if it is linear.

- The *discretized derivative operator* is $A^{B} : \mathbb{R}^{M} \to \mathbb{R}^{M}$. We use the $B$ superscript to emphasize the operator's dependence on the boundary condition.

  The operator is composed as $A^{B} = AQ^{B}$. The intuition is that first $Q^{B}$ is applied to the interior points to add the "ghost nodes" corresponding to the boundary condition, and then the stencil operator $A$ is applied to the whole domain, including the ghost nodes. $A^{B}$ is in general affine if $A$ and/or $Q^{B}$ are affine, and linear if both of them are linear, in which case we will denote it as $L^{B}$ to emphasize the linearity.

## 2.2   Composite Differential Operators

The discretization of a composite differential operator follows the same framework as 2.1. However, instead of deriving the stencil operator $A$ directly it is customary to think of it as being composed of several component operators. For this document we will only consider the simplest of compositions: linear combinations. For more complex examples, especially high-dimensional ones, we may encounter more advanced form of compositions such as tensor products/sums.

Let $\tilde{A}$ be the linear combination of several differential operators: $\tilde{A} = \tilde{A}_1 + \tilde{A}_2 + \cdots + \tilde{A}_n$. To get the discretized entities of 2.1, we will proceed as follows:

- First, discretize each $\tilde{A}_k$ separately and get its $B_k$, $R_k$, $Q^{B_k}$ and $A_k$ along with the extended domain $\bar{u}_k$.[4]

- Let $\bar{u}$ be the union of all the $\bar{u}_k$, i.e. the largest common extended domain. Usually this is just the $\bar{u}_k$ corresponding to the "biggest" component operator. We need to work out the $B$, $R$ and $Q_B$ for $\bar{u}$, which is trivial if it coincides with one of the $\bar{u}_k$.

- The composed stencil operator defined on $\bar{u}$ is $A = A_1 E_1 + A_2 E_2 + \cdots + A_n E_n$, with $E_k \in \mathbb{R}^{M_{E,k} \times M_E}$ the *extension operator* for $A_k$. In practice we don't need $E_k$ explicitly as $A_k E_k$ can be constructed easily by padding $A_k$ with zeros at columns corresponding to boundary points that are not used. For example, if $L_1 \in \mathbb{R}^{2 \times 3}$ but the common extended domain include an extra boundary point at the end, then

$$L_k E_k = \begin{bmatrix} L_k & \mathbf{0}_2 \end{bmatrix} \in \mathbb{R}^{2 \times 4}$$

  which extends $L_k$ to the common extended domain.[5] The extended stencil operators can now be linearly combined as they are of the same shape.

---

[3]Notice that $Q = R^{-1}$ if R is square, and this is only true as maps on functions which satisfy the boundary. Furthermore, in order for (6) to hold on trivial $u$, we need that the interior of $Q$ is identity, so it is defined by its first and last rows.

[4]Even if the physical boundary condition is the same for all components, the required boundary points may still be different because of the order of differentiation and/or numerical approximation. For example, a second-order approximation to $\partial_{xx}$ requires one boundary point at each end, whereas a fourth-order approximation requires two.

[5]For affine $A_k$, the rule of affine algebra applies: $(A_k E_k)_L = A_{k,L} E_k$, $(A_k E_k)_b = A_{k,b}$. In other words, the bias term remains unchanged and the linear part gets padded with zeros.

Alternatively, $E_k$ can be viewed as the *restriction operator* mapping $\bar{u}$ to $\bar{u}_k$ by stripping away unused boundary points. In other words, instead of viewing $A_k E_k \bar{u}$ as $(A_k E_k)\bar{u}$ we can view it as $A_k(E_k \bar{u}) = A_k \bar{u}_k$.[6]

- The discretized derivative operator is still $A^B = AQ^B$.

# 3 Time-Invariant Stochastic Process Examples

## 3.1 Definitions and Notation for Examples

Let $x_t$ be a stochastic process for a univariate function defined on a continuous domain $x \in (x^{\min}, x^{\max})$ where $-\infty < x^{\min} < x^{\max} < \infty$. We will assume throughout that the domain is time-invariant.

For a given $\tilde{L}^s$ as the infinitesimal generate for a stochastic process. Then, if the payoff in state $x$ is $\tilde{p}(x)$, and payoffs are discounted at rate $r > 0$, then the simple HJBE for $\tilde{u}(x)$ is,

$$r\tilde{u}(x) = \tilde{p}(x) + \tilde{L}^s \tilde{u}(x) \tag{8}$$

Rearranging and defining an intermediate,

$$\tilde{L}\tilde{u}(x) = \tilde{p}(x) \tag{9}$$

where the differential operator is

$$\tilde{L} \equiv r - \tilde{L}^s \tag{10}$$

subject to $\partial_x \tilde{u}(x^{\min}) = 0$ and $\partial_x \tilde{u}(x^{\max}) = 0$ for reflecting barriers. If it is a lower absorbing barrier, then denote $\partial_x \tilde{u}(x^{\min}) = \underline{u}$ which may be non-zero.

For a simple example of a payoff, choose $\tilde{p}(x) = x$.

Since many of the examples will use the same boundary values and discretizations of the operators: define the discretization of the second-order derivative with central differences as

$$L_2 \equiv \frac{1}{\Delta^2} \begin{bmatrix} 1 & -2 & 1 & \dots & 0 & 0 & 0 \\ 0 & 1 & -2 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & -2 & 1 & 0 \\ 0 & 0 & 0 & \cdots & 1 & -2 & 1 \end{bmatrix} \tag{11}$$

Next, define the discretization of the first-derivative (forward and backward) as

$$L_1^+ = \frac{1}{\Delta} \begin{bmatrix} 0 & -1 & 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & -1 & 1 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & \cdots & 0 & -1 & 1 \end{bmatrix} \tag{12}$$

$$L_1^- = \frac{1}{\Delta} \begin{bmatrix} -1 & 1 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & -1 & 1 & 0 \end{bmatrix} \tag{13}$$

---

[6]This can make a big difference in application depending on the details of implemnetation: whether $A_k$ is sparese or not, how the grid is implemented, etc.

Next, notice that many of the $R$ matrix follow a simple pattern. When $M_E = 2$ with a boundary point at both sides,

$$R_2 \equiv \begin{bmatrix} \mathbf{0}_M & \mathbf{I}_M & \mathbf{0}_M \end{bmatrix} \tag{14}$$

A common $Q$ setup for this is,

$$Q_{A,2} \equiv \begin{bmatrix} \mathbf{0}_{1 \times M} \\ \mathbf{I}_M \\ \mathbf{0}_{1 \times M} \end{bmatrix} \tag{15}$$

Next, define the $B$ associated with a reflection at both sides

$$B_{RR} \equiv \begin{bmatrix} -1 & 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 1 \end{bmatrix}_{2 \times (M+2)} \tag{16}$$

Another for the $B$ associated with an absorbing barrier at both sides

$$B_{AA} \equiv \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & 1 \end{bmatrix}_{2 \times (M+2)} \tag{17}$$

And another for an absorbing at the bottom and reflecting at the top,

$$B_{AR} \equiv \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 1 \end{bmatrix}_{2 \times (M+2)} \tag{18}$$

## 3.2 Stationary HJBE with Reflecting Barriers

Take the stochastic process

$$dx_t = dW_t$$

with reflecting barriers at $x^{\min}$ and $x^{\max}$. The partial differential operator (infinitesimal generator) associated with the stochastic process is

$$\tilde{L}^s \equiv \frac{1}{2} \partial_{xx} \tag{19}$$

For this process, we derive below all of the matrices of Section 2 and the system of equations to solve for $\tilde{u}(x)$ in (8). We still have **to do**:

- Check that the code `operator_examples\simple_stationary_HJBE_reflecting.jl` is correct

Consider

$$r\tilde{u}(x) = \tilde{L}^s \tilde{u}(x) - x \tag{20}$$

Define the operator $\tilde{L}$ and rearrange,

$$\tilde{L}\tilde{u}(x) = x \tag{21}$$

$$\tilde{L} \equiv r - \frac{1}{2} \partial_{xx} \tag{22}$$

5

We first consider a one-dimension case where $x \in [x^{\min}, x^{\max}]$. Let $M_E = 2$, $S_E = \{1, \bar{M}\}$, thereby $\Delta = \frac{x^{max} - x^{min}}{\bar{M}}$, and $\bar{M} = M + 2$. From (22) and given (11) from the previous section, the matrix form of operator $\tilde{L}$ can be defined, given as $A = \frac{L_2}{2}$.

By reflecting barriers, we can define $B$ just like as $B_{RR}$ in (16) and then we have

$$B\bar{u} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \tag{23}$$

Therefore, $\bar{u}(x_0) = \bar{u}(x_1)$ and $\bar{u}(x_{M+1}) = \bar{u}(x_M)$. It is important to notice that our choice for $B$ defines the linear relationship between the interior points and the boundary conditions.

Moreover, $R$ is again defined as in (14) and $Q$ is defined by $QR\bar{u} \equiv Q_L R\bar{u} + Q_b = \bar{u}$, where

$$Q_L = \begin{bmatrix} 1 & 0 & \ldots & 0 & 0 \\ & & \mathbf{I}_M & & \\ 0 & 0 & \ldots & 0 & 1 \end{bmatrix} \quad , \quad Q_b = \mathbf{0}_{\bar{M}} \tag{24}$$

Then, it is easy to verify that (83) and (84) hold in this case. Additionally, it is worth to note that, since $Q_b = \mathbf{0}_{\bar{M}}$, $Q$ is a linear operator.

To solve $\bar{u}(x)$, we first solve interiors according to (20) and the definition of operator $Q$, which provides us with two conditions:

$$L\bar{u} = x \tag{25}$$
$$QR\bar{u} = Qu = Q_L u + Q_b = \bar{u} \tag{26}$$

where the discretization of the linear operator $\tilde{L}$, defined on the extension, is

$$L \equiv ([\mathbf{0}_M \ \mathbf{I}_M \ \mathbf{0}_M]r - L_2) \tag{27}$$

Note that in doing the composition of the operator in (27), we need to combine the $L_2$ (defined on $M + 2$ for the extension) with $rI$, defined on $M$ points. In order to compose these, we need to extend the $rI$ operator with 0s. Substitute (26) into (25), we get

$$LQu = L(Q_L u + Q_b) = x \tag{28}$$

Since $Q_b = 0$ in this case, interiors are solved as

$$u = (LQ_L)^{-1} x \tag{29}$$

## 3.3 Stationary HJBE with a Lower Absorbing Barrier

Take the stochastic process
$$dx_t = dW_t$$

with an absorbing barrier at $x^{\min}$, and a reflecting barrier at $x^{\max}$. Again, the partial differential operator (infinitesimal generator) associated with the stochastic process is (19)

For the absorbing barrier, when solving for the HJBE assume that $u(x^{\min}) = b^{\min}$ and $u'(x^{\max}) = 0$.

For this process, we derive below all of the matrices of Section 2 and the system of equations to solve for $\tilde{u}(x)$ in (8).

Again, we first consider a one-dimension case where $x \in [x^{\min}, x^{\max}]$. Let $M_E = 2$, $S_E = \{1, \bar{M}\}$, thereby $\Delta = \frac{x^{max} - x^{min}}{\bar{M}}$, and $\bar{M} = M + 2$. From (22) and given (11) from the Section 3, the matrix form of operator $\tilde{L}$ is again defined (27)

According to lower absorbing barrier, we can define $B$ just like as $B_{AR}$ in (18) and then we have

$$B\bar{u} = \begin{bmatrix} b^{\min} \\ 0 \end{bmatrix} \equiv b \tag{30}$$

Therefore, $\bar{u}(x^{\min}) = x^{\min}$ and $\bar{u}(x_{M+1}) = \bar{u}(x_M)$.

Moreover, $R$ is again defined as in (14) and $Q$ is defined by defined by $QR\bar{u} \equiv Q_L R\bar{u} + Q_b = \bar{u}$, where

$$Q_L = \begin{bmatrix} & \mathbf{0}_{1 \times M} & \\ & \mathbf{I}_M & \\ 0 \ 0 & \dots & 0 \ 1 \end{bmatrix} \quad , \quad Q_b = \begin{bmatrix} b^{\min} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \tag{31}$$

Then, it is easy to verify that (83) and (84) hold in this case. Additionally, it is worth to note that, if the absorbing boundary was of Dirichlet$_0$ type, then $b^{\min} = 0$ and $Q_b = \mathbf{0}_{\bar{M}}$ and $Q$ would be a linear operator.

In this case, $L$ is given by (27). According to conditions (25) and (26), again we get

$$L(Q_L R\bar{u} + Q_b) = x \tag{32}$$

With $Q_b \neq 0$, we can solve interiors as

$$R\bar{u} = (LQ_L)^{-1}(x - LQ_b) \tag{33}$$

Then, the extended state vector $\bar{u}$ again can be similarly solved by **??**.

## 3.4 Stationary HJBE with Only Drift

Now, do the same after adding in constant drift (and manually choose the correct upwind direction!)

$$dx_t = \mu dt$$

With a generator

$$\tilde{L}^s \equiv \mu \partial_x \tag{34}$$

For this process, we derive below all of the matrices of Section 2, paying special attention to the upwind direction, and the system of equations to solve for $\tilde{u}(x)$ in (8). We still have **to do**:

- Write julia code to solve for $\tilde{u}(x)$ with the grid

- Check these for $\mu < 0$ and $\mu > 0$.

Since the choice of the first difference depends on the sign of drift $\mu$, we define $\mu^- = \min\{\mu, 0\}$ and $\mu^- = \max\{\mu, 0\}$. Consider

$$\tilde{L}\tilde{u}(x) = x \tag{35}$$
$$\text{where } \tilde{L} \equiv r - \mu \partial_x \tag{36}$$

We first consider a one-dimension case where $x \in [x^{\min}, x^{\max}]$. Let $M_E = 2$ and thereby $\Delta = \frac{x^{\max} - x^{\min}}{\bar{M}}$ and $\bar{M} = M + 2$.

Considering $\mu > 0$, we must choose the forward first difference, thus the matrix form of operator $\tilde{L}$ can be defined as $L = \mu L_1^+$ as in (12). Analogously, for $\mu < 0$, we must choose the backward first difference, which implies that the matrix form of operator $\tilde{L}$ can be defined as $L = \mu L_1^-$ as in (13).

Considering the absorbing barriers, we can define $B$ just like as $B_{AA}$ in (17) and then we have

$$B\bar{u} = \begin{bmatrix} b^{\min} \\ b^{\max} \end{bmatrix} \tag{37}$$

.

Therefore, $\bar{u}(x_0) = x^{\min}$ and $\bar{u}(x_{M+1}) = x^{\max}$.

Moreover, $R$ is again defined as in (14) and $Q$ is defined by $QR\bar{u} \equiv Q_L R\bar{u} + Q_b = \bar{u}$, where

$$Q_L = \begin{bmatrix} \mathbf{0}_{1 \times M} \\ \mathbf{I}_M \\ \mathbf{0}_{1 \times M} \end{bmatrix} \quad , \; Q_b = \begin{bmatrix} b^{\min} \\ 0 \\ \vdots \\ b^{\max} \end{bmatrix} \tag{38}$$

Then, it is easy to verify that (83) and (84) hold in this case.

Similarly, as $L$ is defined by (27) and the interiors are solved by (25) and (26):

$$u = (LQ_L)^{-1}(x - LQ_b) \tag{39}$$

Thus, the extended state vector $\bar{u}$ again can be similarly solved by **??**.

## 3.5 Stationary HJBE with Reflecting Barriers and Drift

Now, do the same after adding in constant drift (and manually choose the correct upwind direction!)

$$dx_t = \mu dt + \sigma dW_t$$

With a generator

$$\tilde{L}^s \equiv \mu \partial_x + \frac{\sigma^2}{2} \partial_{xx}$$

For this process, we derive below all of the matrices of Section 2, paying special attention to the upwind direction, and the system of equations to solve for $\tilde{u}(x)$ in (8). We still have **to do**:

- Write julia code to solve for $\tilde{u}(x)$ with the grid

- Check these for $\mu < 0$ and $\mu > 0$.

We first consider a one-dimension case where $x \in [x^{\min}, x^{\max}]$. Let $M_E = 2$ and thereby $\Delta = \frac{x^{\max} - x^{\min}}{\bar{M}}$ and $\bar{M} = M + 2$.

By combining operators from previous sections, in this case $L^s$ is defined as

$$L^s = \mu L_1^- + \frac{\sigma^2}{2} L_2 \quad \text{if } \mu < 0 \tag{40}$$

$$L^s = \mu L_1^+ + \frac{\sigma^2}{2} L_2 \quad \text{if } \mu > 0 \tag{41}$$

And the composed operator,

$$L \equiv ([\mathbf{0}_M \; \mathbf{I}_M \; \mathbf{0}_M]r - L^s \tag{42}$$

Since barriers are reflecting, we can have the same boundary conditions as what we had in the case with reflecting barriers but no drifts. Hence, operators $R$, $B$ and $Q$ are defined by (14), (23) and (24), respectively. Also, with some simple algebras, we can easily verify that (83) and (84) hold in this case.

Given $L$ defined by (27), the rest steps for solving interiors, $u$, and the extended state vector $\bar{u}$, are similar with what we did for previous examples.

## 3.6 Stationary Bellman Equation with Reflecting Barriers State Varying Drift/Variance

Now, do the same after adding in constant drift (and manually choose the correct upwind direction!)

$$dx_t = \tilde{\mu}(x_t)dt + \tilde{\sigma}(x_t)dW_t$$

With a generator

$$\tilde{L}^s \equiv \tilde{\mu}(x)\partial_x + \frac{\tilde{\sigma}(x)^2}{2}\partial_{xx}$$

For this process, we derive below all of the matrices of Section 2, paying special attention to the upwind direction, and the system of equations to solve for $\tilde{u}(x)$ in (8). We still have **to do**:

- Write julia code to solve for $\tilde{u}(x)$ with the grid.

    - Choose a $\tilde{u}(x)$ and $\tilde{\sigma}(x)$ functions, consider using geometric brownian motion as a test. That is:

$$\tilde{L}^s \equiv \bar{\mu}x\partial_x + \frac{\bar{\sigma}^2}{2}x^2\partial_{xx} \tag{43}$$

We first consider a one-dimension case where $x \in [x^{\min}, x^{\max}]$. Let $M_E = 2$ and thereby $\Delta = \frac{x^{\max}-x^{\min}}{M}$ and $\bar{M} = M + 2$.

Again, consider $\bar{\mu}^- = \min\{\bar{\mu}, 0\}$ and $\bar{\mu}^+ = \max\{\bar{\mu}, 0\}$.

This case is similar with the previous one but with variable drift and variance. By combining operators $L$ from previous sections, in this case $L^s, L$ is defined as

$$L^s = \mu(x)^- L_1^- + \mu(x)^+ L_1^+ + \sigma(x)L_2 \tag{44}$$

where

$$\mu(x)^- = \begin{bmatrix} \bar{\mu}^- x_1 & 0 & \cdots & 0 \\ 0 & \bar{\mu}^- x_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \bar{\mu}^- x_M \end{bmatrix}$$

$$\mu(x)^+ = \begin{bmatrix} \bar{\mu}^+ x_1 & 0 & \cdots & 0 \\ 0 & \bar{\mu}^+ x_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \bar{\mu}^+ x_M \end{bmatrix}$$

$$\sigma(x) = \begin{bmatrix} \frac{(\bar{\sigma}x_1)^2}{2} & 0 & \cdots & 0 \\ 0 & \frac{(\bar{\sigma}x_2)^2}{2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{(\bar{\sigma}x_M)^2}{2} \end{bmatrix}$$

And the composed operator,

$$L \equiv ([\mathbf{0}_M \ \mathbf{I}_M \ \mathbf{0}_M]r - L^s$$

Again, since barriers are reflecting, we can have the same boundary conditions as what we had in the case with reflecting barriers but no drifts. Hence, operators $R$, $B$ and $Q$ are defined by (14), (23) and (24), respectively. Also, with some simple algebras, we can easily verify that (83) and (84) hold in this case.

Again, given $L$ defined by (27), the remaining steps for solving interiors, $u$, and the extended state vector, $\bar{u}$, are similar with what we did for previous examples.

# A   Extended Equation and Gaussian Elimination

## A.1   Overview

The document focuses on solving the discretized equation on the interior $u$, and the information for the boundary condition is encoded in the boundary extrapolation operator $Q^B$. Alternatively, we can consider the discretized equation on the extended domain $\bar{u}$. We wish to demonstrate in this section that the two equations are indeed equivalent.

The starting point is again the continuous equation (1) and (2). We discretize the domain and get the stencil operator $L$ and boundary operator $B$. For simplicity, we shall assume $L$ to be linear in this section, but the boundary conditions need not be homogenous. The discretized equations on the extended domain $\bar{u}$ are then:

- From (1): $L\bar{u} = p$.

- From (2): $B\bar{u} = b$ (the same as (3)).

The two linear equations have the same number of unknowns, so they can stacked up:

$$\begin{bmatrix} L \\ B \end{bmatrix} \bar{u} = \begin{bmatrix} p \\ b \end{bmatrix} \tag{45}$$

In the examples we shall see that the extended equation (45) can be transformed to an equation on the interior by way of Gaussian elimination. We wish to show that the $Q^B$ matrix can be naturally generated using the elementary matrices associated with the elimination process, however the algebra is not in place yet. Nevertheless, for a known $Q^B$ the equivalence between the equations can be proved easily. [7]

## A.2   Example 1: Diffusion with Reflecting Boundaries

### A.2.1   Continuous Equation

We consider a slightly modified example from Section 3.2. The stochastic process in question is $dx_t = \sqrt{2}dW_t$ and a reflecting boundary is present at $x^{\min} = 0$ and $x^{\max} = 2$. The infinitesimal generator is then simply $\tilde{L}^s = \partial_{xx}$.

For discount rate $r > 0$ and payoff $\tilde{p}(x)$, the stationiary HJBE is then

$$\tilde{L}u(x) = \tilde{p}(x) \tag{46}$$

$$\tilde{L} \equiv r\tilde{I} - \partial_{xx} \tag{47}$$

$$\partial_x \tilde{u}(0) = \partial_x \tilde{u}(2) = 0 \tag{48}$$

Here $\tilde{I}$ is the (continuous) identity operator.

### A.2.2   Discretized Equation

We'll be using a uniform grid with $\Delta x = 1$, in other words $M = 3$ interior nodes. A second-order approximation to $\partial_{xx}$ is used which require $M_E = 2$ addition nodes, and a total of $\bar{M} = 5$ nodes on the whole domain. The discretized grid entities are:

$$p = \begin{bmatrix} \tilde{p}(0) & \tilde{p}(1) & \tilde{p}(2) \end{bmatrix}^\top \tag{49}$$

$$u = \begin{bmatrix} \tilde{u}(0) & \tilde{u}(1) & \tilde{u}(2) \end{bmatrix}^\top \tag{50}$$

$$\bar{u} = \begin{bmatrix} \tilde{u}(-1) & \tilde{u}(0) & \tilde{u}(1) & \tilde{u}(2) & \tilde{u}(3) \end{bmatrix}^\top \tag{51}$$

---

[7]Note that the equation (45) does not necessarily have a unique solution. Nevertheless the reduced equation from Gaussian elimination should be the same we get from $LQ^B u = p$.

The boundary operator for simple reflecting boundaries is given in (16). In particular, for this case we have

$$B = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \in \mathbb{R}^{M_E \times \bar{M}}, \quad b = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \tag{52}$$

For the discretization of $\tilde{L} = r\tilde{I} - \partial_{xx}$, since it is composed of two parts we shall first discretize the components:

- The scaling operator $r\tilde{I}$ is discretized simply as $rI_M$, defined on the interior

- The discrete stencil operator for $\partial_{xx}$ is defined in (11), specifically for this case it is

$$L_2 = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 \\ 0 & 1 & -2 & 1 & 0 \\ 0 & 0 & 1 & -2 & 1 \end{bmatrix} \in \mathbb{R}^{M \times \bar{M}} \tag{53}$$

However, we cannot simply add the two operators since they're defined on different domains. In cases like this, we need to first extend the "smaller" component operators to the largest common grid and then combine them. In this case, $L_2$'s domain is the largest so we need only extend $I_M$.

While we could add in arbitrary points in the extension, the algebra will be easier if we add 0s. Define the extension operator adding one point to the left and one to the right of the grid as

$$E_{11} \equiv \begin{bmatrix} \mathbf{0}_M^\top \\ I_M \\ \mathbf{0}_M^\top \end{bmatrix} \in \mathbb{R}^{\bar{M} \times M} \tag{54}$$

Using this, we can extend the identity operator to

$$I^E \equiv E_{11}I_M = \begin{bmatrix} \mathbf{0}_M & I_M & \mathbf{0}_M \end{bmatrix} \tag{55}$$

$$= \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} \tag{56}$$

and get the composed stencil as

$$L = rI^E - L_2 \tag{57}$$

$$= \begin{bmatrix} -1 & 2+r & -1 & 0 & 0 \\ 0 & -1 & 2+r & -1 & 0 \\ 0 & 0 & -1 & 2+r & -1 \end{bmatrix} \tag{58}$$

### A.2.3 Solving the Stacked Equation

Substituting $L$, $p$, $B$ and $b$ into the stacked extended equation (45), we get

$$\begin{bmatrix} -1 & 2+r & -1 & 0 & 0 \\ 0 & -1 & 2+r & -1 & 0 \\ 0 & 0 & -1 & 2+r & -1 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \bar{u} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ 0 \\ 0 \end{bmatrix} \tag{59}$$

This is a well defined linear system. As an example, when $r = 0.25$, the solution is $u \approx \begin{bmatrix} 5.2 & 6.5 & 8.4 \end{bmatrix}$.

We will now show that the rows corresponding to $B$ can be used in Gaussian elimination to reduce the system to one defined in the interior. First, add the 4th row to the first row to get

$$\begin{bmatrix} 1-1 & -1+2+r & -1 & 0 & 0 \\ 0 & -1 & 2+r & -1 & 0 \\ 0 & 0 & -1 & 2+r & -1 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \bar{u} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ 0 \\ 0 \end{bmatrix} \tag{60}$$

Next, add the 5th row to the 3rd row and simplify,

$$\begin{bmatrix} 0 & 1+r & -1 & 0 & 0 \\ 0 & -1 & 2+r & -1 & 0 \\ 0 & 0 & -1 & 1+r & 0 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} \bar{u}(-1) \\ \bar{u}(0) \\ \bar{u}(1) \\ \bar{u}(2) \\ \bar{u}(3) \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ 0 \\ 0 \end{bmatrix} \tag{61}$$

Notice that we have eliminated the extension nodes from all of the equations involving the $L$. Consequently, can just take out the sub-matrix between columns 2-4 and rows 1-3 to get

$$\begin{bmatrix} 1+r & -1 & 0 \\ -1 & 2+r & -1 \\ 0 & -1 & 1+r \end{bmatrix} u = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} \tag{62}$$

On the other hand, for the reflecting boundaries, we know the boundary extrapolation operator is

$$Q^B \equiv \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \tag{63}$$

and the associated discretized equation on the interior is $LQ^B u = p$. It is easy to check that plugging $L$ from (58) and $Q^B$ from (63) also gives (62), which proves the equivalence between the two equations.

## A.3 Example 2: Diffusion and Drift with Reflecting Boundaries

### A.3.1 Continuous Equation

Let's add a drift term to A.2 with constant rate $\mu < 0$ and keep everything else the same. The stochastic process is now $dx_t = \mu dt + \sqrt{2}dW_t$ and the corresponding stationary HJBE becomes

$$\tilde{L}u(x) = \tilde{p}(x) \tag{64}$$

$$\tilde{L} \equiv r\tilde{I} - \partial_{xx} - \mu\partial_x \tag{65}$$

$$\partial_x \tilde{u}(0) = \partial_x \tilde{u}(2) = 0 \tag{66}$$

### A.3.2 Discretized Equation

We have the same $B$ and $b$ as (52). For the stencils, $rI_M$ along with its extension $rI^E$ and $L_2$ are also the same. Because $\mu < 0$, backward difference should be used to discretize $\partial_x$ and that gives

the stencil operator (13), which in this case is

$$
L_1^- = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix} \tag{67}
$$

and again we need to extend $L_1^-$ to the whole domain, which include one extra node at the top end. This gives the extended operator

$$
L_1^{-E} = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 \end{bmatrix} \tag{68}
$$

Finally, the composed operator is

$$
L = rI^E - L_2 - \mu L_1^{-E} \tag{69}
$$

$$
= \begin{bmatrix} -1+\mu & 2-\mu+r & -1 & 0 & 0 \\ 0 & -1+\mu & 2-\mu+r & -1 & 0 \\ 0 & 0 & -1+\mu & 2-\mu+r & -1 \end{bmatrix} \tag{70}
$$

and the stacked equation (45) is

$$
\begin{bmatrix} -1+\mu & 2-\mu+r & -1 & 0 & 0 \\ 0 & -1+\mu & 2-\mu+r & -1 & 0 \\ 0 & 0 & -1+\mu & 2-\mu+r & -1 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \bar{u} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ 0 \\ 0 \end{bmatrix} \tag{71}
$$

### A.3.3 Solving the Stacked Equation

Again we solve the stacked equation (45) using Gaussian elimination on the $B$ rows. First add $(1 - \mu)$ times row 4 to row 1, and then add row 5 to row 3. This gives

$$
\begin{bmatrix} 0 & 1+r & -1 & 0 & 0 \\ 0 & -1+\mu & 2-\mu+r & -1 & 0 \\ 0 & 0 & -1+\mu & 1-\mu+r & 0 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \bar{u} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ 0 \\ 0 \end{bmatrix} \tag{72}
$$

Extract the interior of the matrix to get

$$
\begin{bmatrix} 1+r & -1 & 0 \\ -1+\mu & 2-\mu+r & -1 \\ 0 & -1+\mu & 1-\mu+r \end{bmatrix} u = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} \tag{73}
$$

The $Q^B$ in this example is the same as (63). It is easy to check that the interior equation $LQ^B u = p$ also gives (73), again confirming that the extended equation gives the same resuls.

## A.4 Example 3: Diffusion with Inhomogeneous Boundaries

### A.4.1 Continuous Equation

Let's consider A.2 again but change the boudnary conditions to be inhomogeneous. The stationary HJBE:

$$\tilde{L}u(x) = \tilde{p}(x) \tag{74}$$

$$\tilde{L} \equiv r\tilde{I} - \partial_{xx} \tag{75}$$

$$\partial_x \tilde{u}(0) = b^{\min} \tag{76}$$

$$\partial_x \tilde{u}(2) = b^{\max} \tag{77}$$

### A.4.2 Discretized Equation

The discretized $L$, $B$ and $p$ are the same as A.2. Since the boundary conditions are now inhomogeneous $b$ is no longer the zero vector. Recalling (16), for this example we have

$$b = \begin{bmatrix} -b^{\min} \\ b^{\max} \end{bmatrix} \tag{78}$$

and the stacked equation (45) is now

$$\begin{bmatrix} -1 & 2+r & -1 & 0 & 0 \\ 0 & -1 & 2+r & -1 & 0 \\ 0 & 0 & -1 & 2+r & -1 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix} \bar{u} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ -b^{\min} \\ b^{\max} \end{bmatrix} \tag{79}$$

### A.4.3 Solving the Stacked Equation

Since the left hand side coefficient matrix is the same, we can use the same Gaussian elimination procedure as A.2. This gives the reduced equation

$$\begin{bmatrix} 1+r & -1 & 0 \\ -1 & 2+r & -1 \\ 0 & -1 & 1+r \end{bmatrix} u = \begin{bmatrix} p_1 - b^{\min} \\ p_2 \\ p_3 + b^{\max} \end{bmatrix} \tag{80}$$

For the $LQ^B u = p$ route, $Q^B$ is now affine and from (15) we have

$$Q_L^B = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \qquad Q_b^B = \begin{bmatrix} -b^{\min} \\ 0 \\ 0 \\ 0 \\ b^{\max} \end{bmatrix} \tag{81}$$

and the equation is

$$LQ_L^B u = p - LQ_b^B \tag{82}$$

substituting $L$, $p$, $Q_L^B$ and $Q_b^B$, we get back (80), again proving the equivalence.

# B Affine Relations and Intuition

The following provides intuition on the relationships above:[8]

- Define $S_E$ as the set of non interior grid point indexes, e.g, if we have a univariate problem with just three non-interior point, let's say, the first one and the last two grid points, then $M_E = 3$ and $S_E = \{1, \bar{M} - 1, \bar{M}\}$[9].

- Now let's focus on solving an expression like $\bar{u} = A\bar{u}$ where $A$ is affine.[10] Consider, with some abuse of notation, that such expression means both the discretized and non-discretized forms. Define, for any arbitrary matrix J with at least $M_E$ columns, that $J[:, S_E]$ is a submatrix whose columns are the concatenated vectors $J[:, s], \forall s \in S_E$. Then we have two relations[11]:

$$A[:, S] \left( B[:, S_E]^{-1} b \right) = A Q_b \tag{83}$$

$$(A - A[:, S_E](B[:, S_E]^{-1} B)) R^\top = A Q_L \tag{84}$$

Our intuition is that the main idea here is using interiors to recover a relation that boundary conditions should satisfy. Since $Q_b$ is a length $\bar{M}$ vector containing zeros excepts two ends, the two non-zero elements in $Q_b$ capture partial information of boundary nodes (the part that is "independent" of interiors).[12]

Recall (3), so $\left( B[:, S_E]^{-1} b \right)$ recovers the "independent" part of boundary nodes. Then it is reasonable to expect that (83) holds.

Multiply both sides of (84) by $\bar{u}$, we can roughly rewrite the relation as

$$(A\bar{u} - A Q_b) R^\top = A Q_L u \tag{85}$$

so $A\bar{u} - A Q_b$ will be a discretized $\bar{u}$ which contains the entire information of interiors and the rest part of boundary information that is not covered by $A Q_b$.

However, we are not sure if $R$ should exist on the left of (84) since R by definition is a restriction operator and $(A\bar{u} - A Q_b) R^\top$ only contains information from interiors. For now, while we work on better understanding those expressions, we will take them as given.

Given those relationships, in order to solve the differential equation, we now only have to solve for the interior. Otherwise, including the boundary values would imply having more points than there are degrees of freedom in the problem - thus making the numerical solution unstable. Moreover, the boundaries are given directly by the interior $\bar{u} = Qu$.

Therefore, we actually want to solve $\bar{u} = AQu$. Notice that the discretized A maps from the full domain to the interior[13]. Notably, that means it's not square. Additionally, consider that, as described above, since $Q$ is in general affine, thus:

$$\bar{u} = A Q_L u + A Q_b \tag{86}$$

---

[8]**TODO: Fernando/Steven** I think you will need to rewrite this with the modified notation and go through it carefully. I don't quite get it, and the notation was slightly different than the rest of the text... Also, I think that abusing notation for the discretized and non-discretized is part of the problem. We might want to rewrite this a little after the expanding operator setup is dine.

[9]Notice that, by construction, the number of elements in $S_E$ is always $M_E$

[10]**Fernando/Steven**: Is this a particular operator you have in mind from our setup, or a general affine operator you are going through? Point it out from above, and differentiate the $\tilde{A}$ from the discretized $A$

[11]We could not find a way to clearly show two relations above are correct, but some intuitions are provided in the text

[12]**Typo From Chris** (83) is actually defined from (84) which is defined from the next one. Looking at it like that, it's clear to see the error since 89 is just saying $(A - A Q b) * R^T = A Q L$ substitute in 88 for AQb) you see the substitution was done incorrectly has a big B instead of a little b.

[13]Notice that the PDE is only defined on the interior

Those are the linear equations which define the ODE or whose solution is the solution to the PDE.