# Contents

# Appendix A
## Linear algebra basics

## A.1 Basic operations

We frequently envision or frame problems as linear systems of equations.[1] It is useful to write this compactly in matrix notation, say

$$Ay = x$$

where $A$ is an $m \times n$ (rows $\times$ columns) *matrix* (a rectangular array of elements), $y$ is an $n$-element vector, and $x$ is an $m$-element vector. This statement compares the result on the left with that on the right, element-by-element. The operation on the left is *matrix multiplication* or each element is recovered by a vector inner product of the corresponding row from $A$ with the vector $y$. That is, the first element of the product vector $Ay$ is the vector inner product of the first row $A$ with $y$, the second element of the product vector is the inner product of the second row $A$ with $y$, and so on. A vector *inner product* multiplies the same position element of the leading row and trailing column and sums over the products. Of course, this means that the operation is only well-defined if the number of columns in the leading matrix, $A$, equals the number of rows of the trailing, $y$. Further, the product matrix has the same number of rows as the leading matrix and

---

[1] G. Strang, *Linear Algebra and its Applications*, Harcourt Brace Jovanovich College Publishers, or *Introduction to Linear Algebra*, Wellesley-Cambridge Press offers a mesmerizing discourse on linear algebra.

columns of the trailing. For example, let

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix},$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix},$$

and

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

then

$$Ay = \begin{bmatrix} a_{11}y_1 + a_{12}y_2 + a_{13}y_3 + a_{14}y_4 \\ a_{21}y_1 + a_{22}y_2 + a_{23}y_3 + a_{24}y_4 \\ a_{31}y_1 + a_{32}y_2 + a_{33}y_3 + a_{34}y_4 \end{bmatrix}$$

The system of equations also covers *matrix addition* and *scalar multiplication* by a matrix in the sense that we can rewrite the equations as

$$Ay - x = 0$$

First, multiplication by a scalar or constant simply multiplies each element of the matrix by the scalar. In this instance, we multiple the elements of $x$ by $-1$.

$$\begin{bmatrix} a_{11}y_1 + a_{12}y_2 + a_{13}y_3 + a_{14}y_4 \\ a_{21}y_1 + a_{22}y_2 + a_{23}y_3 + a_{24}y_4 \\ a_{31}y_1 + a_{32}y_2 + a_{33}y_3 + a_{34}y_4 \end{bmatrix} - \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} a_{11}y_1 + a_{12}y_2 + a_{13}y_3 + a_{14}y_4 \\ a_{21}y_1 + a_{22}y_2 + a_{23}y_3 + a_{24}y_4 \\ a_{31}y_1 + a_{32}y_2 + a_{33}y_3 + a_{34}y_4 \end{bmatrix} + \begin{bmatrix} -x_1 \\ -x_2 \\ -x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Then, we add the $m$-element vector $x$ to the $m$-element vector $Ay$ where same position elements are summed.

$$\begin{bmatrix} a_{11}y_1 + a_{12}y_2 + a_{13}y_3 + a_{14}y_4 - x_1 \\ a_{21}y_1 + a_{22}y_2 + a_{23}y_3 + a_{24}y_4 - x_2 \\ a_{31}y_1 + a_{32}y_2 + a_{33}y_3 + a_{34}y_4 - x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Again, the operation is only well-defined for equal size matrices and, unlike matrix multiplication, matrix addition always commutes. Of course, the $m$-element vector (the additive identity) on the right has all zero elements.

By convention, vectors are represented in columns. So how do we represent an inner product of a vector with itself? We create a row vector by

transposing the original. *Transposition* simply puts columns of the original into same position rows of the transposed. For example, $y^T y$ represents the vector inner product (the product is a scalar) of $y$ with itself where the superscript $T$ represents transposition.

$$y^T y \;=\; \begin{bmatrix} y_1 & y_2 & y_3 & y_4 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$=\; y_1^2 + y_2^2 + y_3^2 + y_4^2$$

Similarly, we might be interested in $A^T A$.

$$A^T A = \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \\ a_{14} & a_{24} & a_{34} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix}$$

$$= \begin{bmatrix}
\left( \begin{smallmatrix} a_{11}^2 + a_{21}^2 \\ + a_{31}^2 \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{11}a_{12} + a_{21}a_{22} \\ + a_{31}a_{32} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{11}a_{13} + a_{21}a_{23} \\ + a_{31}a_{33} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{11}a_{14} + a_{21}a_{24} \\ + a_{31}a_{34} \end{smallmatrix} \right) \\[2mm]
\left( \begin{smallmatrix} a_{11}a_{12} + a_{21}a_{22} \\ + a_{31}a_{32} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{12}^2 + a_{22}^2 \\ + a_{32}^2 \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{12}a_{13} + a_{22}a_{23} \\ + a_{32}a_{33} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{12}a_{14} + a_{22}a_{24} \\ + a_{32}a_{34} \end{smallmatrix} \right) \\[2mm]
\left( \begin{smallmatrix} a_{11}a_{13} + a_{21}a_{23} \\ + a_{31}a_{33} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{12}a_{13} + a_{22}a_{23} \\ + a_{32}a_{33} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{13}^2 + a_{23}^2 \\ + a_{33}^2 \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{13}a_{14} + a_{23}a_{24} \\ + a_{33}a_{34} \end{smallmatrix} \right) \\[2mm]
\left( \begin{smallmatrix} a_{11}a_{14} + a_{21}a_{24} \\ + a_{31}a_{34} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{12}a_{14} + a_{22}a_{24} \\ + a_{32}a_{34} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{13}a_{14} + a_{23}a_{24} \\ + a_{33}a_{34} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{14}^2 + a_{24}^2 \\ + a_{34}^2 \end{smallmatrix} \right)
\end{bmatrix}$$

This yields an $n \times n$ symmetric product matrix. A matrix is *symmetric* if the matrix equals its transpose, $A = A^T$. Or, $AA^T$ which yields an $m \times m$ product matrix.

$$AA^T = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix} \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \\ a_{14} & a_{24} & a_{34} \end{bmatrix}$$

$$= \begin{bmatrix}
\left( \begin{smallmatrix} a_{11}^2 + a_{12}^2 \\ + a_{13}^2 + a_{14}^2 \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{11}a_{21} + a_{12}a_{22} \\ + a_{13}a_{23} + a_{14}a_{24} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{11}a_{31} + a_{12}a_{22} \\ + a_{13}a_{33} + a_{14}a_{34} \end{smallmatrix} \right) \\[2mm]
\left( \begin{smallmatrix} a_{11}a_{21} + a_{12}a_{22} \\ + a_{13}a_{23} + a_{14}a_{24} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{21}^2 + a_{22}^2 \\ + a_{23}^2 + a_{24}^2 \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{21}a_{31} + a_{22}a_{32} \\ + a_{23}a_{33} + a_{24}a_{34} \end{smallmatrix} \right) \\[2mm]
\left( \begin{smallmatrix} a_{11}a_{31} + a_{12}a_{22} \\ + a_{13}a_{33} + a_{14}a_{34} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{21}a_{31} + a_{22}a_{32} \\ + a_{23}a_{33} + a_{24}a_{34} \end{smallmatrix} \right) &
\left( \begin{smallmatrix} a_{31}^2 + a_{32}^2 \\ + a_{33}^2 + a_{34}^2 \end{smallmatrix} \right)
\end{bmatrix}$$

## A.2    Fundamental theorem of linear algebra

With these basic operations in hand, return to

$$Ay = x$$

When is there a unique solution, $y$? The answer lies in the *fundamental theorem of linear algebra*. The theorem has two parts.

### A.2.1    Part one

First, the theorem says that for every matrix the number of linearly independent rows equals the number of linearly independent columns. *Linearly independent* vectors are the set of vectors such that no one of them can be duplicated by a linear combination of the other vectors in the set. A *linear combination* of vectors is the sum of scalar-vector products where each vector may have a different scalar multiplier. For example, $Ay$ is a linear combination of the columns of $A$ with the scalars in $y$. Therefore, if there exists some $(n-1)$-element vector, $w$, when multiplied by an $(m \times (n-1))$ submatrix of $A$, call it $B$, such that $Bw$ produces the dropped column from $A$ then the dropped column is not linearly independent of the other columns. To reiterate, if the matrix $A$ has $r$ linearly independent columns it also has $r$ linearly independent rows. $r$ is referred to as the *rank* of the matrix and *dimension* of the *rowspace* and *columnspace* (the spaces *spanned* by all possible linear combination of the rows and columns, respectively). Further, $r$ linearly independent rows of $A$ form a basis for its rowspace and $r$ linearly independent columns of $A$ form a basis for its columnspace.

*Accounting example*

Consider an incidence matrix describing the journal entry properties of accounting in its columns (each column has one $+1$ and $-1$ in it with the remaining elements equal to zero) and T accounts in its rows. The rows capture changes in account balances when multiplied by a transaction amounts vector $y$. By convention, we assign $+1$ for a debit entry and $-1$ for a credit entry. Suppose

$$A = \begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix}$$

where the rows represent cash, noncash assets, liabilities, and owners' equity, for instance. Notice, $-1$ times the sum of any three rows produces the remaining row. Since we cannot produce another row from the remaining two rows, the number of linearly independent rows is 3. By the fundamental theorem, the number of linearly independent columns must also be 3.

Let's check. Suppose the first three columns is a basis for the columnspace. Column 4 is the negative of column 3, column 5 is the negative of the sum of columns 1 and 3, and column 6 is the negative of the sum of columns 2 and 3. Can any of columns 1, 2 and 3 be produced as a linearly combination of the remaining two columns? No, the zeroes in rows 2 through 4 rule it out. For this matrix, we've confirmed the number of linearly independent rows and columns is the same.

### A.2.2   Part two

The second part of the fundamental theorem describes the orthogonal complements to the rowspace and columnspace. Two vectors are *orthogonal* if they are perpendicular to one another. As their vector inner product is proportional to the cosine of the angle between them, if their vector inner product is zero they are orthogonal.[2] $n$-space is spanned by the rowspace (with dimension $r$) plus the $n - r$ dimension orthogonal complement, the *nullspace* where

$$AN^T = 0$$

$N$ is an $(n - r) \times n$ matrix whose rows are orthogonal to the rows of $A$ and $0$ is an $m \times (n - r)$ matrix of zeroes.

*Accounting example continued*

For the A matrix above, a basis for the nullspace is

$$N = \begin{bmatrix} 1 & -1 & 0 & 0 & 1 & -1 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

and

$$AN^T = \begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ -1 & 1 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

---

[2] The inner product of a vector with itself is the squared length (or squared norm) of the vector.

Similarly, $m$-space is spanned by the columnspace (with dimension $r$) plus the $m - r$ dimension orthogonal complement, the *left nullspace* where

$$(LN)^T A = 0$$

$LN$ is an $m \times (m - r)$ matrix whose rows are orthogonal to the columns of $A$ and 0 is an $(m - r) \times n$ matrix of zeroes. The origin is the only point in common to the four subspaces: rowspace, columnspace, nullspace, and left nullspace.

$$LN = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

and

$$(LN)^T A = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

## A.3   Nature of the solution

### A.3.1   Exactly-identified

If $r = m = n$, then there is a unique solution, $y$, to $By = x$, and the problem is said to be exactly-identified.[3] Since $B$ is square and has a full set of linearly independent rows and columns, the rows and columns of $B$ span $r$ space (including $x$) and the two nullspaces have dimension zero. Consequently, there exists a matrix $B^{-1}$, the *inverse of B*, when multiplied by $B$ produces the *identity matrix, I*. The identity matrix is a matrix when multiplied (on the left or on the right) by any other vector or matrix leaves that vector or matrix unchanged. The identity matrix is a square matrix with ones along the principal diagonal and zeroes on the off-diagonals.

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

Hence,

$$\begin{aligned} B^{-1}By &= B^{-1}x \\ Iy &= B^{-1}x \\ y &= B^{-1}x \end{aligned}$$

Suppose

$$B = \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix},$$

and

$$x = \begin{bmatrix} 6 \\ 5 \end{bmatrix}$$

then

$$y = B^{-1}x$$

$$\begin{bmatrix} \frac{4}{10} & -\frac{1}{10} \\ -\frac{2}{10} & \frac{3}{10} \end{bmatrix}\begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix}\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \frac{4}{10} & -\frac{1}{10} \\ -\frac{2}{10} & \frac{3}{10} \end{bmatrix}\begin{bmatrix} 6 \\ 5 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \frac{24-5}{10} \\ \frac{-12+15}{10} \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \frac{19}{10} \\ \frac{3}{10} \end{bmatrix}$$

---

[3] Both $y$ and $x$ are $r$-element vectors.

## A.3.2   Under-identified

However, it is more common for $r \leq m, n$ with one inequality strict. In this case, spanning $m$-space draws upon both the columnspace and left nullspace and spanning $n$-space draws from both the rowspace and the nullspace. If the dimension of the nullspace is greater than zero, then it is likely that there are many solutions, $y$, that satisfy $Ay = x$. On the other hand, if the dimension of the left nullspace is positive and the dimension of the nullspace is zero, then typically there is no exact solution, $y$, for $Ay = x$. When $r < n$, the problem is said to be *under-identified* (there are more unknown parameters than equations) and a complete set of solutions can be described by the solution that lies entirely in the rows of $A$ (this is often called the *row component* as it is a linear combination of the rows) plus arbitrary weights on the nullspace of $A$. The row component, $y^{RS(A)}$, can be found by projecting any consistent solution, $y^p$, onto a *basis* for the rows (any linearly independent set of $r$ rows) of $A$. Let $A^r$ be a submatrix derived from $A$ with $r$ linearly independent rows. Then,

$$
\begin{aligned}
y^{RS(A)} &= (A^r)^T \left( A^r (A^r)^T \right)^{-1} A^r y^p \\
&= (P_{A^r}) y^p
\end{aligned}
$$

and

$$
y^p = y^{RS(A)} + N^T k
$$

where $P_{A^r}$ is the projection matrix, $(A^r)^T \left( A^r (A^r)^T \right)^{-1} A^r$, onto the rows of $A$ and $k$ is any $n$-element vector of arbitrary weights.

Utilizing $y^p = y^{RS(A)} + y^{NS(A)} = (A^r)^T b + N^T k$, we have two immediate ways to derive projection matrices. First, $y^{RS(A)} = (A^r)^T b$ says the row component of $y^p$ is a linear combination of the rows of $A^r$ with weights $b$ and $y^{NS(A)} = N^T k$ says the null component of $y^p$ is a linear combination of the rows of $N$ with weights $k$. Projecting into the rows of $A^r$ follows from orthogonality of the row and null components.

$$
y^p = (A^r)^T b + y^{NS(A)}
$$

where

$$
A^r y^{NS(A)} = 0
$$

Since

$$
y^{NS(A)} = y^p - (A^r)^T b
$$

we have by substitution

$$
A^r \left( y^p - (A^r)^T b \right) = 0
$$

or

$$
A^r y^p = A^r (A^r)^T b
$$

As $A^r$ has linearly independent rows, the inverse of $A^r \left(A^r\right)^T$ exists and we can solve for the weights

$$\left(A^r \left(A^r\right)^T\right)^{-1} A^r y^p = \left(A^r \left(A^r\right)^T\right)^{-1} A^r \left(A^r\right)^T b = Ib = b$$

Now that we have $b$, we can immediately identify the row component of $y^p$

$$\begin{aligned} y^{RS(A)} &= \left(A^r\right)^T b \\ &= \left(A^r\right)^T \left(A^r \left(A^r\right)^T\right)^{-1} A^r y^p \\ &= \left(P_{A^r}\right) y^p \end{aligned}$$

The projection is matrix is symmetric $\left(\left(P_{A^r}\right)^T = P_{A^r}\right)$ and idempotent $\left(\left(P_{A^r}\right)\left(P_{A^r}\right) = P_{A^r}\right)$. Idempotency is appealing since if $y^p = y^{RS(A)}$ and we project $y^p$ into the rows of $A^r$ then it doesn't change rather it remains $y^{RS(A)}$ (the row component is unique).

Notice from above we have

$$\begin{aligned} y^{NS(A)} &= y^p - \left(A^r\right)^T b \\ &= y^p - \left(A^r\right)^T \left(A^r \left(A^r\right)^T\right)^{-1} A^r y^p \\ &= \left(I - P_{A^r}\right) y^p \end{aligned}$$

which implies the projection matrix into the rows of the nullspace of $A^r$ can be described by $P_{A^n} = \left(I - P_{A^r}\right)$. Alternatively (and equivalently), $P_{A^n} = N^T \left(NN^T\right) N$. This representation of the projection matrix is derived in analogous fashion to $P_{A^r}$ above.

$$y^p = y^{RS(A)} + N^T k$$

where

$$N y^{RS(A)} = 0$$

Since

$$y^{RS(A)} = y^p - N^T k$$

we have by substitution

$$N \left(y^p - N^T k\right) = 0$$

or

$$N y^p = N N^T k$$

As $N$ has linearly independent rows, the inverse of $NN^T$ exists and we can solve for the weights

$$\left(NN^T\right)^{-1} N y^p = \left(NN^T\right)^{-1} NN^T k = Ik = k$$

Now that we have $k$, we can immediately identify the null component of $y^p$

$$
\begin{aligned}
y^{NS(A)} &= N^T k \\
&= N^T \left(N N^T\right)^{-1} N y^p \\
&= (P_{A^n})\, y^p
\end{aligned}
$$

$P_{A^n}$ is also symmetric and idempotent. Further, from the above analysis it's clear $P_{A^r} + P_{A^n} = I$ (the entire $n$-dimensional space is spanned by linear combinations of the rowspace and nullspace).

Return to our accounting example above. Suppose the changes in account balances are

$$
x = \begin{bmatrix} 2 \\ 1 \\ -1 \\ -2 \end{bmatrix}
$$

Then, a particular solution can be found by setting, for example, the last three elements of $y$ equal to zero and solving for the remaining elements.

$$
y^p = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}
$$

so that

$$
A y^p = x
$$

$$
\begin{bmatrix}
-1 & -1 & 1 & -1 & 0 & 0 \\
1 & 0 & 0 & 0 & -1 & 0 \\
0 & 1 & 0 & 0 & 0 & -1 \\
0 & 0 & -1 & 1 & 1 & 1
\end{bmatrix}
\begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 0 \end{bmatrix}
=
\begin{bmatrix} 2 \\ 1 \\ -1 \\ -2 \end{bmatrix}
$$

Let $A^r$ be the first three rows.

$$
A^r = \begin{bmatrix}
-1 & -1 & 1 & -1 & 0 & 0 \\
1 & 0 & 0 & 0 & -1 & 0 \\
0 & 1 & 0 & 0 & 0 & -1
\end{bmatrix}
$$

$$
P_{A^r} = \frac{1}{12}
\begin{bmatrix}
7 & 1 & -2 & 2 & -5 & 1 \\
1 & 7 & -2 & 2 & 1 & -5 \\
-2 & -2 & 4 & -4 & -2 & -2 \\
2 & 2 & -4 & 4 & 2 & 2 \\
-5 & 1 & -2 & 2 & 7 & 1 \\
1 & -5 & -2 & 2 & 1 & 7
\end{bmatrix}
$$

and

$$y^{RS(A)} = (P_{A^r})\, y^p$$

$$= \frac{1}{12} \begin{bmatrix} 7 & 1 & -2 & 2 & -5 & 1 \\ 1 & 7 & -2 & 2 & 1 & -5 \\ -2 & -2 & 4 & -4 & -2 & -2 \\ 2 & 2 & -4 & 4 & 2 & 2 \\ -5 & 1 & -2 & 2 & 7 & 1 \\ 1 & -5 & -2 & 2 & 1 & 7 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \frac{1}{6} \begin{bmatrix} 1 \\ -5 \\ 4 \\ -4 \\ -5 \\ 1 \end{bmatrix}$$

The complete solution, with arbitrary weights $k$, is

$$y = y^{RS(A)} + N^T k$$

$$y = \frac{1}{6} \begin{bmatrix} 1 \\ -5 \\ 4 \\ -4 \\ -5 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 1 \\ -5 \\ 4 \\ -4 \\ -5 \\ 1 \end{bmatrix} + \begin{bmatrix} k_1 \\ -k_1 + k_2 \\ k_2 + k_3 \\ k_3 \\ k_1 \\ -k_1 + k_2 \end{bmatrix}$$

### A.3.3  Over-identified

In the case where there is no exact solution, $m > r = n$, the vector that lies entirely in the columns of $A$ which is nearest $x$ is frequently identified as the best approximation. This case is said to be *over-identified* (there are more equations than unknown parameters) and this best approximation is the column component, $y^{CS(A)}$, and is found via projecting $x$ onto the columns of $A$. A common variation on this theme is described by

$$Y = X\beta$$

where $Y$ is an $n$-element vector and $X$ is an $n \times p$ matrix. Typically, no exact solution for $\beta$ (a $p$-element vector) exists, $p = r$ ($X$ is composed of

linearly independent columns), and

$$b = \beta^{CS(A)} = \left(X^T X\right)^{-1} X^T Y$$

is known as the *ordinary least squares* (*OLS*) estimator of $\beta$ and the estimated conditional expectation function is the projection of $Y$ into the columns of $X$

$$Xb = X \left(X^T X\right)^{-1} X^T Y = P_X Y$$

For example, let $X = (A^r)^T$ and $Y = y^P$. then

$$b = \frac{1}{6} \begin{bmatrix} 4 \\ 5 \\ -1 \end{bmatrix}$$

and $Xb = P_X Y = y^{RS(A)}$.

$$Xb = P_X Y = \frac{1}{6} \begin{bmatrix} 1 \\ -5 \\ 4 \\ -4 \\ -5 \\ 1 \end{bmatrix}$$

## A.4   Matrix decomposition and inverse operations

Inverse operations are inherently related to the fundamental theorem and matrix decomposition. There are a number of important decompositions, we'll focus on four: *LU* factorization, Cholesky decomposition, singular value decomposition, and spectral decomposition.

### A.4.1   LU factorization

Gaussian elimination is the key to solving systems of linear equations and gives us *LU decomposition.*

*Nonsingular case*

Any square, *nonsingular* matrix $A$ (has linearly independent rows and columns) can be written as the product of a lower triangular matrix, $L$, times an upper triangular matrix, $U$.[4]

$$A = LU$$

where $L$ is *lower triangular* meaning that it has all zero elements above the main diagonal and $U$ is *upper triangular* meaning that it has all zero elements below the main diagonal. *Gaussian elimination* says we can write any system of linear equations in triangular form so that by backward recursion we solve a series of one equation, one variable problems. This is accomplished by row operations: row eliminations and row exchanges. *Row eliminations* involve a series of operations where a scalar multiple of one row is added to a target row so that a revised target row is produced until a triangular matrix, $L$ or $U$, is generated. As the same operation is applied to both sides (the same row(s) of $A$ and $x$) equality is maintained. *Row exchanges* simply revise the order of both sides (rows of $A$ and elements of $x$) to preserve the equality. Of course, the order in which equations are written is flexible.

In principle then, Gaussian elimination on

$$Ay = x$$

involves, for instance, multiplication of both sides by the inverse of $L$, provided the inverse exists $(m = r)$,

$$
\begin{aligned}
L^{-1}Ay &= L^{-1}x \\
L^{-1}LUy &= L^{-1}x \\
Uy &= L^{-1}x
\end{aligned}
$$

---

[4] The general case, $A$ is a $m \times n$ matrix, is discussed below.

As Gaussian elimination is straightforward, we have a simple approach for finding whether the inverse of the lower triangular matrix exists and, if so, its elements. Of course, we can identify $L$ in similar fashion

$$
\begin{aligned}
L &= AU^{-1} \\
&= LUU^{-1}
\end{aligned}
$$

Let $A = A^r (A^r)^T$ the $3 \times 3$ full rank matrix from the accounting example above.

$$
A = \begin{bmatrix} 4 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix}
$$

We find $U$ by Gaussian elimination. Multiply row 1 by $1/4$ and add to rows 2 and 3 to revise rows 2 and 3 as follows.

$$
\begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & -1/4 & 7/4 \end{bmatrix}
$$

Now, multiply row 2 by $1/7$ and add this result to row 3 to identify $U$.

$$
U = \begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & 0 & 12/7 \end{bmatrix}
$$

Notice we have constructed $L^{-1}$ in the process.

$$
\begin{aligned}
L^{-1} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1/7 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 1/4 & 1 & 0 \\ 1/4 & 0 & 1 \end{bmatrix} \\
&= \begin{bmatrix} 1 & 0 & 0 \\ 1/4 & 1 & 0 \\ 2/7 & 1/7 & 1 \end{bmatrix}
\end{aligned}
$$

so that

$$
L^{-1}A = U
$$

$$
\begin{bmatrix} 1 & 0 & 0 \\ 1/4 & 1 & 0 \\ 2/7 & 1/7 & 1 \end{bmatrix} \begin{bmatrix} 4 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix} = \begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & 0 & 12/7 \end{bmatrix}
$$

Also, we have $L$ in hand.

$$
L = \left(L^{-1}\right)^{-1}
$$

and

$$
L^{-1}L = I
$$

$$
\begin{bmatrix} 1 & 0 & 0 \\ 1/4 & 1 & 0 \\ 2/7 & 1/7 & 1 \end{bmatrix} \begin{bmatrix} \ell_{11} & 0 & 0 \\ \ell_{21} & \ell_{22} & 0 \\ \ell_{31} & \ell_{32} & \ell_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}
$$

From the first row-first column, $\ell_{11} = 1$. From the second row-first column, $1/4\ell_{11} + 1\ell_{21} = 0$, or $\ell_{21} = -1/4$. From the third row-first column, $2/7\ell_{11} + 1/7\ell_{21} + 1\ell_{31} = 0$, or $\ell_{31} = -2/7 + 1/28 = -1/4$. From the second row-second column, $\ell_{22} = 1$. From the third row-second column, $1/7\ell_{22} + 1\ell_{32} = 0$, or $\ell_{32} = -1/7$. And, from the third row-third column, $\ell_{33} = 1$. Hence,

$$
L = \begin{bmatrix} 1 & 0 & 0 \\ -1/4 & 1 & 0 \\ -1/4 & -1/7 & 1 \end{bmatrix}
$$

and

$$
LU = A
$$

$$
\begin{bmatrix} 1 & 0 & 0 \\ -1/4 & 1 & 0 \\ -1/4 & -1/7 & 1 \end{bmatrix} \begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & 0 & 12/7 \end{bmatrix} = \begin{bmatrix} 4 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix}
$$

For

$$
x = \begin{bmatrix} -1 \\ 3 \\ 5 \end{bmatrix}
$$

the solution to $Ay = x$ is

$$
\begin{aligned}
Ay &= x \\
LUy &= x \\
Uy &= L^{-1}x
\end{aligned}
$$

$$
\begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & 0 & 12/7 \end{bmatrix} \begin{bmatrix} y1 \\ y2 \\ y3 \end{bmatrix} = \begin{bmatrix} -1 \\ 3 - 1/4 \\ 5 - 1/4 + 11/28 \end{bmatrix} = \begin{bmatrix} -1 \\ 11/4 \\ 36/7 \end{bmatrix}
$$

Backward recursive substitution solve for $y$. From row three, $12/7y_3 = 36/7$ or $y_3 = 7/12 \times 36/7 = 3$. From row two, $7/4y_2 - 1/4y_3 = 11/4$, or $y_2 = 4/7\,(11/4 + 3/4) = 2$. And, from row one, $4y_1 - y_2 - y_3 = -1$, or $y_1 = 1/4\,(-1 + 2 + 3) = 1$. Hence,

$$
y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}
$$

*General case*

If the inverse of $A$ doesn't exist (the matrix is singular), we find some equations after elimination are $0 = 0$, and possibly, some elements of $y$ are not uniquely determinable as discussed above for the under-identified case.

For an $m \times n$ matrix $A$, the general form of $LU$ factorization may involve row exchanges via a permutation matrix, $P$.

$$
PA = LU
$$

where $L$ is lower triangular with ones on the diagonal and $U$ is an $m \times n$ upper echelon matrix with the pivots along the main diagonal.

   $LU$ decomposition can also be written as $LDU$ *factorization* where, as before, $L$ and $U$ are lower and upper triangular matrices but now have ones along their diagonals and $D$ is a diagonal matrix with the pivots of $A$ along its diagonal.

   Returning to the accounting example, we utilize $LU$ factorization to solve for $y$, a set of transactions amounts that are consistent with the financial statements.

$$Ay = x$$

$$\begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \\ -1 \\ -2 \end{bmatrix}$$

For this $A$ matrix, $P = I_4$ (no row exchanges are called for), and row one is added to row two, the revised row two is added to row three, and the revised row three is added to row four, which gives

$$\begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 0 & -1 & 1 & -1 & -1 & 0 \\ 0 & 0 & 1 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \\ 2 \\ 0 \end{bmatrix}$$

The last row conveys no information and the third row indicates we have three free variables. Recall, for our solution, $y^p$ above, we set $y_4 = y_5 = y_6 = 0$ and solved. From row three, $y_3 = 2$. From row two, $y_2 = -(3-2) = -1$. And, from row one, $y_1 = -(2-1-2) = 1$. Hence,

$$y^p = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

## A.4.2  Cholesky decomposition

If the matrix $A$ is symmetric, *positive definite*[5] as well as nonsingular, then we have $A = LDL^T$ as $U = L^T$. In this symmetric case, we identify another useful factorization, Cholesky decomposition. *Cholesky decomposition* writes

$$A = \Gamma\Gamma^T$$

where $\Gamma = LD^{\frac{1}{2}}$ and $D^{\frac{1}{2}}$ has the square root of the pivots on the diagonal. Since $A$ is positive definite, all of its pivots are positive and their square root is real so, in turn, $\Gamma$ is real. Of course, we now have

$$
\begin{aligned}
\Gamma^{-1}A &= \Gamma^{-1}\Gamma\Gamma^T \\
&= \Gamma^T
\end{aligned}
$$

or

$$
\begin{aligned}
A\left(\Gamma^T\right)^{-1} &= \Gamma\Gamma^T\left(\Gamma^T\right)^{-1} \\
&= \Gamma
\end{aligned}
$$

For the example above, $A = A^r\left(A^r\right)^T$, $A$ is symmetric, positive definite and we found

$$
\begin{aligned}
A &= LU \\
\begin{bmatrix} 4 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix} &= \begin{bmatrix} 1 & 0 & 0 \\ -1/4 & 1 & 0 \\ -1/4 & -1/7 & 1 \end{bmatrix} \begin{bmatrix} 4 & -1 & -1 \\ 0 & 7/4 & -1/4 \\ 0 & 0 & 12/7 \end{bmatrix}
\end{aligned}
$$

Factoring the pivots from $U$ gives $D$ and $LDL^T$.

$$
\begin{aligned}
A &= LDL^T \\
&= \begin{bmatrix} 1 & 0 & 0 \\ -1/4 & 1 & 0 \\ -1/4 & -1/7 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & 7/4 & 0 \\ 0 & 0 & 12/7 \end{bmatrix} \begin{bmatrix} 1 & -1/4 & -1/4 \\ 0 & 1 & -1/7 \\ 0 & 0 & 1 \end{bmatrix}
\end{aligned}
$$

And, the Cholesky decomposition is

$$
\begin{aligned}
\Gamma &= LD^{\frac{1}{2}} \\
&= \begin{bmatrix} 1 & 0 & 0 \\ -1/4 & 1 & 0 \\ -1/4 & -1/7 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & \sqrt{7/4} & 0 \\ 0 & 0 & \sqrt{12/7} \end{bmatrix} \\
&= \begin{bmatrix} 2 & 0 & 0 \\ -1/2 & \frac{\sqrt{7}}{2} & 0 \\ -1/2 & -\frac{1}{2\sqrt{7}} & 2\sqrt{\frac{3}{7}} \end{bmatrix}
\end{aligned}
$$

---

[5] A matrix, $A$, is *positive definite* if its *quadratic form* is strictly positive

$$x^T A x > 0$$

for all nonzero $x$.

so that

$$A \;=\; \Gamma\Gamma^T$$

$$\begin{bmatrix} 4 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ -1/2 & \frac{\sqrt{7}}{2} & 0 \\ -1/2 & -\frac{1}{2\sqrt{7}} & 2\sqrt{\frac{3}{7}} \end{bmatrix} \begin{bmatrix} 2 & -1/2 & -1/2 \\ 0 & \frac{\sqrt{7}}{2} & -\frac{1}{2\sqrt{7}} \\ 0 & 0 & 2\sqrt{\frac{3}{7}} \end{bmatrix}$$

## A.4.3   Eigenvalues and eigenvectors

A square $n \times n$ matrix $A$ times a characteristic vector $x$ can be written as a characteristic scalar $\lambda$ times the same vector.

$$Ax = \lambda x$$

The characteristic scalar is called an eigenvalue and the characteristic vector is called an eigenvector. There are $n$ (not necessarily unique) eigenvalues and associated eigenvectors.[6] Rewriting the above as

$$(A - \lambda I)\,x = 0$$

reveals the key subspace feature. That is, we choose $\lambda$ such that $A - \lambda I$ has a nullspace. Then, $x$ is a vector in the nullspace of $A - \lambda I$.

Example

Now, we explore construction of eigenvalues and eigenvectors via our accounting example.

$$A = \begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix}$$

In particular, focus attention on

$$AA^T = \begin{bmatrix} 4 & -1 & -1 & -2 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ -2 & -1 & -1 & 4 \end{bmatrix}$$

First, we know due to the balancing property of accounting this matrix has a nullspace. Hence, at least one of its eigenvalues equals zero. We'll

---

[6] For instance, an $n \times n$ identity matrix has $n$ eigenvalues equal to one and any orthogonal (or unitary) matrix is a basis for the eigenvectors.

verify this by $AA^T = LU = LDL^T$ and then utilize this result to find the eigenvalues.

First, utilize row operations to put $AA^T$ in row echelon form and find its pivots. Row operations on the first column are

$$L_1^{-1}AA^T = \begin{bmatrix} 4 & -1 & -1 & -2 \\ 0 & \frac{7}{4} & -\frac{1}{4} & -\frac{3}{2} \\ 0 & -\frac{1}{4} & \frac{7}{4} & -\frac{3}{2} \\ 0 & -\frac{3}{2} & -\frac{3}{2} & 3 \end{bmatrix}$$

where

$$L_1^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ \frac{1}{4} & 1 & 0 & 0 \\ \frac{1}{4} & 0 & 1 & 0 \\ \frac{1}{2} & 0 & 0 & 1 \end{bmatrix}$$

Combine this with row operations on the second column.

$$L_2^{-1}L_1^{-1}AA^T = \begin{bmatrix} 4 & -1 & -1 & -2 \\ 0 & \frac{7}{4} & -\frac{1}{4} & -\frac{3}{2} \\ 0 & 0 & \frac{12}{7} & -\frac{12}{7} \\ 0 & 0 & -\frac{12}{7} & \frac{12}{7} \end{bmatrix}$$

where

$$L_2^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & \frac{1}{7} & 1 & 0 \\ 0 & \frac{6}{7} & 0 & 1 \end{bmatrix}$$

Combining this with row operations on the third column yields the upper triangular result we're after.

$$L_3^{-1}L_2^{-1}L_1^{-1}AA^T = U = \begin{bmatrix} 4 & -1 & -1 & -2 \\ 0 & \frac{7}{4} & -\frac{1}{4} & -\frac{3}{2} \\ 0 & 0 & \frac{12}{7} & -\frac{12}{7} \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

where

$$L_3^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

Hence,

$$L^{-1} = L_3^{-1}L_2^{-1}L_1^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ \frac{1}{4} & 1 & 0 & 0 \\ \frac{2}{7} & \frac{1}{7} & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

and

$$L^{-1}AA^T = U$$
$$= \begin{bmatrix} 4 & -1 & -1 & -2 \\ 0 & \frac{7}{4} & -\frac{1}{4} & -\frac{3}{2} \\ 0 & 0 & \frac{12}{7} & -\frac{12}{7} \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Clearly, the rank of $U$ is three and column four is free as its pivot (main diagonal element in row echelon form) is zero. This means, as suggested before, one eigenvalue equals zero. To find its associated eigenvector replace row four with the row vector $\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}$, call this $U^a$ and solve

$$U^a x = b$$
$$\begin{bmatrix} 4 & -1 & -1 & -2 \\ 0 & \frac{7}{4} & -\frac{1}{4} & -\frac{3}{2} \\ 0 & 0 & \frac{12}{7} & -\frac{12}{7} \\ 0 & 0 & 0 & 1 \end{bmatrix} x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

for $x$. A solution is

$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

Since eigenvectors are scale-free, $AA^T x = \lambda x$ accommodates any rescaling of $x$, it is often convenient to make this vector unit length. Accordingly, define the unit length eigenvector associated with the zero eigenvalue ($\lambda_1 = 0$) as

$$x_1 = \frac{x}{\sqrt{x^T x}} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

What are the remaining three eigenvalues? Clearly from $U$ (the first row is unchanged by our row operations), therefore at least one of the remaining eigenvalues is $\lambda = 4$ (we could have repeated eigenvalues). A more general approach is find $\lambda$ such that the matrix $AA^T - \lambda I$ is singular or equivalently, its determinant is zero. Determinants are messy but we'll utilize two facts: the determinant of a triangular matrix is the product of its pivots (main diagonal elements) and the product of determinants equals the determinant of the products, $\det(L)\det(U) = \det(LU)$. Since $L$ has ones along the main diagonal it's determinant is one, the determinant of $U$ is the determinant of $AA^T - \lambda I$. Finding eigenvalues of $AA^T$ boils down to

finding the roots of the product of the main diagonal elements of $U$ where $AA^T - \lambda I = LU$.

Following similar steps to those above, we find

$$
U = \begin{bmatrix}
4 - \lambda & -1 & -1 & -2 \\
0 & \frac{7-6\lambda+\lambda^2}{4-\lambda} & -\frac{1}{4-\lambda} & \frac{\lambda-6}{4-\lambda} \\
0 & 0 & \frac{12-18\lambda+8\lambda^2-\lambda^3}{7-6\lambda+\lambda^2} & \frac{-12+8\lambda-\lambda^2}{7-6\lambda+\lambda^2} \\
0 & 0 & 0 & \frac{-24\lambda+10\lambda^2-\lambda^3}{6-6\lambda+\lambda^2}
\end{bmatrix}
$$

$$
\begin{aligned}
\det(U) &= (4-\lambda)\left(\frac{7-6\lambda+\lambda^2}{4-\lambda}\right)\left(\frac{12-18\lambda+8\lambda^2-\lambda^3}{7-6\lambda+\lambda^2}\right) \times \\
&\quad \left(\frac{-24\lambda+10\lambda^2-\lambda^3}{6-6\lambda+\lambda^2}\right) \\
&= -48\lambda + 44\lambda^2 - 12\lambda^3 + \lambda^4
\end{aligned}
$$

The roots are $\lambda = 0$, 2, 4, and 6.

The next step is to find eigenvectors for $\lambda = 2$, 4, and 6. For $\lambda = 2$, $U = \begin{bmatrix} 2 & -1 & -1 & -2 \\ 0 & -\frac{1}{2} & -\frac{1}{2} & -2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 8 \end{bmatrix}$. Since the third pivot equals zero its a free

variable and we replace row three with $\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$ and solve

$$
U^a x = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}
$$

$$
\begin{bmatrix}
2 & -1 & -1 & -2 \\
0 & -\frac{1}{2} & -\frac{1}{2} & -2 \\
0 & 0 & 1 & 0 \\
0 & 0 & 0 & 8
\end{bmatrix} x = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}
$$

This yields

$$
x = \begin{bmatrix} 0 \\ -1 \\ 1 \\ 0 \end{bmatrix}
$$

which can be unitized as follows

$$
x_2 = \frac{x}{\sqrt{x^T x}} = \begin{bmatrix} 0 \\ -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}
$$

Notice, $x_2$ is orthogonal to $x_1$.

$$x_1^T x_2 = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 \\ -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix} = 0$$

It works largely the same for $\lambda = 6$. For $\lambda = 6$,

$$U = \begin{bmatrix} -2 & -1 & -1 & -2 \\ 0 & -\frac{7}{2} & \frac{1}{2} & 0 \\ 0 & 0 & -\frac{24}{7} & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Since the fourth pivot equals zero its a free variable and we replace row four with $\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}$ and solve

$$U^a x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} -2 & -1 & -1 & -2 \\ 0 & -\frac{7}{2} & \frac{1}{2} & 0 \\ 0 & 0 & -\frac{24}{7} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

This yields

$$x = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

which can be unitized as follows

$$x_4 = \frac{x}{\sqrt{x^T x}} = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ 0 \\ 0 \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

Notice, this eigenvector is orthogonal to both $x_1$ and $x_2$.

Unfortunately, we can't just plug $\lambda = 4$ into our expression for $U$ as it produces infinities. Rather, we return to $AA^T - 4I$ and factor into its own

$LU$. First, we apply a permutation (row exchanges[7]) to $AA^T - 4I$

$$P\left(AA^T - 4I\right) = LU$$

$$\begin{bmatrix} -1 & -2 & 0 & -1 \\ 0 & -1 & -1 & -2 \\ -1 & 0 & -2 & -1 \\ -2 & -1 & -1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & -2 & 1 & 0 \\ 2 & -3 & 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & -2 & 0 & -1 \\ 0 & -1 & -1 & -2 \\ 0 & 0 & -4 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

where

$$P = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

swaps rows one and two. Then, we follow similar row operations as described above to produce the $LU$ factors where

$$U = \begin{bmatrix} -1 & -2 & 0 & -1 \\ 0 & -1 & -1 & -2 \\ 0 & 0 & -4 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

As the fourth pivot is zero it's a free variable and we replace row four with $\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}$ to solve

$$U^a x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & -2 & 0 & -1 \\ 0 & -1 & -1 & -2 \\ 0 & 0 & -4 & -4 \\ 0 & 0 & 0 & 1 \end{bmatrix} x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

This yields $x = \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix}$ and is unitized as

$$x_3 = \frac{x}{\sqrt{x^T x}} = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

---

[7] Row exchanges can change the sign of the determinant but that is of consequence here because we've chosen the eigenvalue to make the determinant zero.

Now, as $AA^T$ is symmetric all four eigenvectors are orthonormal. Hence, when we construct a matrix $Q$ of eigenvectors in its columns

$$Q = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

and multiply by its transpose

$$QQ^T = Q^T Q = I$$

Further,

$$Q\Sigma Q^T = AA^T$$

$$\begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{\sqrt{2}} & -\frac{1}{2} & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 6 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ 0 & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$= \begin{bmatrix} 4 & -1 & -1 & -2 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ -2 & -1 & -1 & 4 \end{bmatrix}$$

where $\Sigma$ is a diagonal matrix of eigenvalues and the order of the eigenvectors matches the order of the eigenvalues.

### A.4.4   Singular value decomposition

Now, we introduce a matrix factorization that exists for every matrix. *Singular value decomposition* says every $m \times n$ matrix, $A$, can be written as the product of a $m \times m$ orthogonal matrix, $U$, multiplied by a diagonal $m \times n$ matrix, $\Sigma$, and finally multiplied by the transpose of a $n \times n$ orthogonal matrix, $V$.[8] $U$ is composed of the eigenvectors of $AA^T$, $V$ is composed of the eigenvectors of $A^T A$, and $\Sigma$ contains the singular values (the square root of the eigenvalues of $AA^T$ or $A^T A$) along the diagonal.

$$A = U\Sigma V^T$$

Further, singular value decomposition allows us to define a general inverse or *pseudo-inverse*, $A^+$.

$$A^+ = V\Sigma^+ U^T$$

---

[8] An *orthogonal* (or *unitary*) matrix is comprised of orthonormal vectors. That is, mutually orthogonal, unit length vectors so that $U^{-1} = U^T$ and $UU^T = U^T U = I$.

where $\Sigma^+$ is an $n \times m$ diagonal matrix with nonzero elements equal to the reciprocal of those for $\Sigma$. This implies

$$AA^+A = A$$

$$A^+AA^+ = A^+$$

$$\left(A^+A\right)^T = A^+A$$

and

$$\left(AA^+\right)^T = AA^+$$

Also, for the system of equations

$$Ay = x$$

the least squares solution is

$$y^{CS(A)} = A^+x$$

and $AA^+$ is always the projection onto the columns of $A$. Hence,

$$AA^+ = P_A = A\left(A^TA\right)^{-1}A^T$$

if $A$ has linearly independent columns. Or,

$$
\begin{aligned}
A^T\left(A^T\right)^+ &= \left(A^+A\right)^T \\
&= \left(V\Sigma^+U^TU\Sigma V^T\right)^T \\
&= V\Sigma^TU^TU\left(\Sigma^+\right)^TV^T \\
&= A^+A \\
&= P_{A^T} = A^T\left(AA^T\right)^{-1}A
\end{aligned}
$$

if $A$ has linearly independent rows (if $A^T$ has linearly independent columns).

For the accounting example, recall the row component is the consistent solution to $Ay = x$ that is only a linearly combination of the rows of $A$; that is, it is orthogonal to the nullspace. Utilizing the pseudo-inverse we have

$$
\begin{aligned}
y^{RS(A)} &= A^T\left(A^T\right)^+y^p \\
&= P_{A^T}y^p \\
&= \left(A^r\right)^T\left(A^r\left(A^r\right)^T\right)^{-1}A^ry^p \\
&= A^+Ay^p
\end{aligned}
$$

or simply, since $Ay^p = x$

$$
\begin{aligned}
y^{RS(A)} &= A^+x \\
&= \frac{1}{24}
\begin{bmatrix}
-5 & 9 & -3 & -1 \\
-5 & -3 & 9 & -1 \\
4 & 0 & 0 & -4 \\
-4 & 0 & 0 & 4 \\
1 & -9 & 3 & 5 \\
1 & 3 & -9 & 5
\end{bmatrix}
\begin{bmatrix}
2 \\
1 \\
-1 \\
-2
\end{bmatrix} \\
&= \frac{1}{6}
\begin{bmatrix}
1 \\
-5 \\
4 \\
-4 \\
-5 \\
1
\end{bmatrix}
\end{aligned}
$$

The beauty of singular value decomposition is that any $m \times n$ matrix, $A$, can be factored as

$$AV = U\Sigma$$

since

$$AVV^T = A = U\Sigma V^T$$

where $U$ and $V$ are $m \times m$ and $n \times n$ orthogonal matrices (of eigenvectors), respectively, and $\Sigma$ is a $m \times n$ matrix with singular values along its main diagonal.

Eigenvalues are characteristic values or singular values of a square matrix and eigenvectors are characteristic vectors or singular vectors of the matrix such that

$$AA^T u = \lambda u$$

or we can work with

$$A^T Av = \lambda v$$

where $u$ is an $m$-element unitary ($u^T u = 1$) eigenvector (component of $Q_1$), $v$ is an $n$-element unitary ($v^T v = 1$) eigenvector (component of $Q_2$), and $\lambda$ is an eigenvalue of $AA^T$ or $A^T A$. We can write $AA^T u = \lambda u$ as

$$
\begin{aligned}
AA^T u &= \lambda I u \\
\left(AA^T - \lambda I\right) u &= 0
\end{aligned}
$$

or write $A^T Av = \lambda v$ as

$$
\begin{aligned}
A^T Av &= \lambda I v \\
\left(A^T A - \lambda I\right) v &= 0
\end{aligned}
$$

then solve for unitary vectors $u$, $v$, and and roots $\lambda$. For instance, once we have $\lambda_i$ and $u_i$ in hand. We find $v_i$ by

$$u_i^T A = \lambda_i v_i$$

such that $v_i$ is unit length, $v_i^T v_i = 1$.

The sum of the eigenvalues equals the *trace* of the matrix (sum of the main diagonal elements) and the product of the eigenvalues equals the *determinant* of the matrix. A *singular matrix* has some zero eigenvalues and pivots (the $\det(A) = \pm[\text{product of the pivots}]$), hence the determinant of a singular matrix, $\det(A)$, is zero.[9] The eigenvalues can be found by solving $\det\left(AA^T - \lambda I\right) = 0$. Since this is an $m$ order polynomial, there are $m$ eigenvalues associated with an $m \times m$ matrix.

*Accounting example*

Return to the accounting example for an illustration. The singular value decomposition ($SVD$) of $A$ proceeds as follows. We'll work with the square, symmetric matrix $AA^T$. Notice, by $SVD$,

$$\begin{aligned}
AA^T &= U\Sigma V^T \left(U\Sigma V^T\right)^T \\
&= U\Sigma V^T V\Sigma^T U^T \\
&= U\Sigma\Sigma^T U^T
\end{aligned}$$

so that the eigenvalues of $AA^T$ are the squared singular values of $A$, $\Sigma\Sigma^T$. The eigenvalues are found by solving for the roots of[10]

$$\det\left(AA^T - \lambda I_m\right) = 0$$

$$\det\begin{bmatrix} 4-\lambda & -1 & -1 & -2 \\ -1 & 2-\lambda & 0 & -1 \\ -1 & 0 & 2-\lambda & -1 \\ -2 & -1 & -1 & 4-\lambda \end{bmatrix} = 0$$

$$-48\lambda + 44\lambda^2 - 12\lambda^3 + \lambda^4 = 0$$

---

[9] The *determinant* is a value associated with a square matrix with many (some useful) properties. For instance, the determinant provides a test of invertibility (linear independence). If $\det(A) = 0$, then the matrix is singular and the inverse doesn't exist; otherwise $\det(A) \neq 0$, the matrix is nonsingular and the inverse exists. The determinant is the volume of a parallelpiped in n-dimensions where the edges come from the rows of $A$. The determinant of a triangular matrix is the product of the main diagonal elements. Determinants are unchanged by row eliminations and their sign is changed by row exchanges. The determinant of the transpose of a matrix equals the determinant of the matrix, $\det(A) = \det\left(A^T\right)$. The determinant of the product of matrices is the product of their determinants, $\det(AB) = \det(A)\det(B)$. Some useful determinant identities are reported in section five of the appendix.

[10] Below we show how to find the determinant of a square matrix and illustrate with this example.

Immediately, we see that one of the roots is zero,[11] and

$$
\begin{aligned}
-48\lambda + 44\lambda^2 - 12\lambda^3 + \lambda^4 &= 0 \\
\lambda\,(\lambda - 2)\,(\lambda - 4)\,(\lambda - 6) &= 0
\end{aligned}
$$

or

$$
\lambda = \{6, 4, 2, 0\}
$$

for $AA^T$.[12] The eigenvectors for $AA^T$ are found by solving (employ Gaussian elimination and back substitution)

$$
\left(AA^T - \lambda_i I_4\right) u_i = 0
$$

Since there is freedom in the solution, we can make the vectors orthonormal (see Gram-Schmidt discussion below). For instance, $\left(AA^T - 6I_4\right) u_1 = 0$ leads to $u_1^T = \begin{bmatrix} -a & 0 & 0 & a \end{bmatrix}$, so we make $a = \frac{1}{\sqrt{2}}$ and $u_1^T = \begin{bmatrix} -\frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$. Now, the complementary right hand side eigenvector, $v_1$, is found by

$$
u_1^T A = \sqrt{\lambda_1} v_1
$$

$$
v_1 = \frac{1}{\sqrt{6}} u_1^T A = \begin{bmatrix} \frac{1}{2\sqrt{3}} \\ \frac{1}{2\sqrt{3}} \\ -\frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} \\ \frac{1}{2\sqrt{3}} \\ \frac{1}{2\sqrt{3}} \end{bmatrix}
$$

Repeating these steps for the remaining eigenvalues (in descending order; remember its important to match eigenvectors with eigenvalues) leads to

$$
U = \begin{bmatrix}
-\frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \\
0 & -\frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \\
0 & -\frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\
\frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2}
\end{bmatrix}
$$

---

[11] For $\det\left(A^T A - \lambda I_6\right) = 0$, we have $-48\lambda^3 + 44\lambda^4 - 12\lambda^5 + \lambda^6 = 0$. Hence, there are at least three zero roots. Otherwise, the roots are the same as for $AA^T$.

[12] Clearly, $\lambda = \{6, 4, 2, 0, 0, 0\}$ for $A^T A$.

and

$$V = \begin{bmatrix} \frac{1}{2\sqrt{3}} & -\frac{1}{2} & -\frac{1}{2} & 0 & \frac{\sqrt{3}}{2\sqrt{2}} & -\frac{1}{2\sqrt{6}} \\ \frac{1}{2\sqrt{3}} & -\frac{1}{2} & \frac{1}{2} & \frac{1}{\sqrt{3}} & -\frac{1}{2\sqrt{6}} & -\frac{1}{2\sqrt{6}} \\ -\frac{1}{\sqrt{3}} & 0 & 0 & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{6}} & \frac{1}{6} \\ \frac{1}{\sqrt{3}} & 0 & 0 & 0 & 0 & \sqrt{\frac{2}{3}} \\ \frac{1}{2\sqrt{3}} & \frac{1}{2} & \frac{1}{2} & 0 & \frac{\sqrt{3}}{2\sqrt{2}} & -\frac{1}{2\sqrt{6}} \\ \frac{1}{2\sqrt{3}} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{\sqrt{3}} & -\frac{1}{2\sqrt{6}} & -\frac{1}{2\sqrt{6}} \end{bmatrix}$$

where $UU^T = U^TU = I_4$ and $VV^T = V^TV = I_6$.[13] Remarkably,

$$
\begin{aligned}
A &= U\Lambda V^T \\
&= \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & -\frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \\ 0 & -\frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} \sqrt{6} & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{2} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\
&\quad \times \begin{bmatrix} \frac{1}{2\sqrt{3}} & -\frac{1}{2} & -\frac{1}{2} & 0 & \frac{\sqrt{3}}{2\sqrt{2}} & -\frac{1}{2\sqrt{6}} \\ \frac{1}{2\sqrt{3}} & -\frac{1}{2} & \frac{1}{2} & \frac{1}{\sqrt{3}} & -\frac{1}{2\sqrt{6}} & -\frac{1}{2\sqrt{6}} \\ -\frac{1}{\sqrt{3}} & 0 & 0 & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{6}} & \frac{1}{6} \\ \frac{1}{\sqrt{3}} & 0 & 0 & 0 & 0 & \sqrt{\frac{2}{3}} \\ \frac{1}{2\sqrt{3}} & \frac{1}{2} & \frac{1}{2} & 0 & \frac{\sqrt{3}}{2\sqrt{2}} & -\frac{1}{2\sqrt{6}} \\ \frac{1}{2\sqrt{3}} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{\sqrt{3}} & -\frac{1}{2\sqrt{6}} & -\frac{1}{2\sqrt{6}} \end{bmatrix}^T \\
&= \begin{bmatrix} -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 1 & 1 & 1 \end{bmatrix}
\end{aligned}
$$

where $\Lambda$ is $m \times n$ $(4 \times 6)$ with the square root of the eigenvalues (in descending order) on the main diagonal.

### A.4.5 Spectral decomposition

When $A$ is a square, symmetric matrix, singular value decomposition can be expressed as *spectral decomposition*.

$$A = U\Sigma U^T$$

---

[13] There are many choices for the eigenvectors associated with zero eigenvalues. We select them so that they orthonormal. As with the other eigenvectors, this is not unique.

where $U$ is an *orthogonal matrix*. Notice, the matrix on the right is the transpose of the matrix on the left. This follows as $AA^T = A^T A$ when $A = A^T$. We've illustrated this above if when we decomposed $AA^T$, a square symmetric matrix.

$$
\begin{aligned}
AA^T &= U\Sigma U^T \\[6pt]
&= \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & -\frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \\ 0 & -\frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}
\begin{bmatrix} 6 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}
\begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & -\frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \\ 0 & -\frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}^T \\[6pt]
&= \begin{bmatrix} 4 & -1 & -1 & -2 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ -2 & -1 & -1 & 4 \end{bmatrix}
\end{aligned}
$$

### A.4.6   *quadratic forms, eigenvalues, and positive definiteness*

A symmetric matrix $A$ is positive definite if the quadratic form $x^T A x$ is positive for every nonzero $x$. Positive semi-definiteness follows if the quadratic form is non-negative, $x^T A x \geq 0$ for every nonzero $x$. Negative definite and negative semi-definite symmetric matrices follow in analogous fashion where the quadratic form is negative or non-positive, respectively. A positive (semi-) definite matrix has positive (non-negative) eigenvalues. This result follows immediately from spectral decomposition. Let $y = Qx$ ($y$ is arbitrary since $x$ is) and write the spectral decomposition of $A$ as $Q^T \Lambda Q$ where $Q$ is an orthogonal matrix and $\Lambda$ is a diagonal matrix composed of the eigenvalues of $A$. Then the quadratic form $x^T A x > 0$ can be written as $x^T Q^T \Lambda Q x > 0$ or $y^T \Lambda y > 0$. Clearly, this is only true if $\Lambda$, the eigenvalues, are all positive.

### A.4.7   *similar matrices, Jordan form, and generalized eigenvectors*

Now, we provide some support for properties associated with eigenvalues. Namely, for any square matrix the sum of the eigenvalues equals the trace of the matrix and the product of the eigenvalues equals the determinant of the matrix. To aid with this discussion we first develop the idea of similar matrices and the Jordan form of a matrix.

Two matrices, $A$ and $B$, are similar if there exists $M$ and $M^{-1}$ such that $B = M^{-1}AM$. Similar matrices have the same eigenvalues as seen from $Ax = \lambda x$ where $x$ is an eigenvector of $A$ associated with $\lambda$.

$$
\begin{aligned}
Ax &= \lambda x \\
AMM^{-1}x &= \lambda x
\end{aligned}
$$

Since $MB = AM$, we have

$$
\begin{aligned}
MBM^{-1}x &= \lambda x \\
M^{-1}MBM^{-1}x &= \lambda M^{-1}x \\
B\left(M^{-1}x\right) &= \lambda\left(M^{-1}x\right)
\end{aligned}
$$

Hence, $A$ and $B$ have the same eigenvalues where $x$ is the eigenvector of $A$ and $M^{-1}x$ is the eigenvector of $B$.

From here we can see $A$ and $B$ have the same trace and determinant. First, we'll demonstrate, via example, the trace of a matrix equals the sum of its eigenvalues, $\sum \lambda_i = tr\left(A\right)$ for any square matrix $A$. Consider $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$ where $tr\left(A\right) = a_{11} + a_{22}$. The eigenvalues of A are determined from solving $\det\left(A - \lambda I\right) = 0$.

$$
\begin{aligned}
\left(a_{11} - \lambda\right)\left(a_{22} - \lambda\right) - a_{12}a_{21} &= 0 \\
\lambda^2 - \left(a_{11} + a_{22}\right)\lambda + a_{11}a_{22} - a_{12}a_{21} &= 0
\end{aligned}
$$

The two roots or eigenvalues are

$$
\lambda = \frac{a_{11} + a_{22} \pm \sqrt{\left(a_{11} + a_{22}\right)^2 - 4\left(a_{11}a_{22} - a_{12}a_{21}\right)}}{2}
$$

and their sum is $\lambda_1 + \lambda_2 = a_{11} + a_{22} = tr\left(A\right)$. The idea extends to any square matrix $A$ such that $\sum \lambda_i = tr\left(A\right)$. This follows as $\det\left(A - \lambda I\right)$ for any $n \times n$ matrix $A$ has coefficient on the $\lambda^{n-1}$ term equal to minus the coefficient on $\lambda^n$ times $\sum \lambda_i$, as in the $2 \times 2$ example above.[14]

We'll demonstrate the determinant result in two parts: one for diagonalizable matrices and one for non-diagonalizable matrices using their Jordan form. Any diagonalizable matrix can be written as $A = S\Lambda S^{-1}$. The determinant of $A$ is then $|A| = \left|S\Lambda S^{-1}\right| = |S||\Lambda|\left|S^{-1}\right| = |\Lambda|$ since $\left|S^{-1}\right| = \frac{1}{|S|}$ which follows from $\left|SS^{-1}\right| = \left|S^{-1}\right||S| = |I| = 1$. Now, we have $|\Lambda| = \prod \lambda_i$.

The second part follows from similar matrices and the Jordan form. When a matrix is not diagonalizable because it doesn't have a complete set of linearly independent eigenvectors, we say it is nearly diagonalizable when it's in Jordan form. Jordan form means the matrix is nearly diagonal except for perhaps ones immediately above the diagonal.

For example, the identity matrix, $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, is in Jordan form as well as being diagonalizable while $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ is in Jordan form but not diagonaliz-

---

[14] For $n$ even the coefficient on $\lambda^n$ is 1 and for $n$ odd the coefficient on $\lambda^n$ is $-1$ with the coefficient on $\lambda^{n-1}$ of opposite sign.

able. Even though both matrices have the same eigenvalues they are not similar matrices as there exists no $M$ such that $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ equals $M^{-1}IM$.

Nevertheless, the Jordan form is the characteristic form for a family of similar matrices as there exists $P$ such that $P^{-1}AP = J$ where $J$ is the Jordan form for the family. For instance, $A = \frac{1}{3}\begin{bmatrix} 1 & 4 \\ -1 & 5 \end{bmatrix}$ has Jordan form $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ with $P = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$. Consider another example, $A = \frac{1}{5}\begin{bmatrix} 13 & 6 \\ 1 & 12 \end{bmatrix}$ has Jordan form $\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$ with $P = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}$. Since they are similar matrices, $A$ and $J$ have the same eigenvalues. Plus, as in the above examples, the eigenvalues lie on the diagonal of $J$ in general. The determinant of $A = |PJP^{-1}| = |P||J||P^{-1}| = |J| = \prod \lambda_i$. This completes the argument.

To summarize, for any $n \times n$ matrix $A$:

$$(1)\ |A| = \prod \lambda_i,$$

and

$$(2)\ tr\,(A) = \sum \lambda_i.$$

Generalized eigenvectors

The idea of eigenvectors is generalized for non-diagonalizable matrices like $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ as it doesn't have a full set of regular eigenvectors. For such matrices, eigenvectors are the (nullspace or nonzero) solutions, $q$, to $(A - \lambda I)^k q = 0$ for $k \geq 1$ ($k = 1$ for diagonalizable matrices). For the above matrix $k = 2$ as there are two occurrences of $\lambda = 1$.

$$(A - \lambda I)^1 = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

therefore $q = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ is an eigenvector of $A$ but there is no other nonzero, linearly independent vector that resides in the nullspace of $A - \lambda I$. On the other hand,

$$(A - \lambda I)^2 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

and $q = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ and $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ are a basis for the nullspace of $(A - \lambda I)^2$ or generalized eigenvectors of $A$.

## A.5   Gram-Schmidt construction of an orthogonal matrix

Before we put this section to bed, we'll undertake one more task. Construction of an *orthogonal* matrix (that is, a matrix with orthogonal, unit length vectors so that $QQ^T = Q^TQ = I$). Suppose we have a square, symmetric matrix

$$A = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 3 & 2 \\ 1 & 2 & 3 \end{bmatrix}$$

with eigenvalues $\left\{ \frac{1}{2} \left( 7 + \sqrt{17} \right), \frac{1}{2} \left( 7 - \sqrt{17} \right), 1 \right\}$ and eigenvectors (in the columns)

$$\begin{bmatrix} v_1 & v_2 & v_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \left( -3 + \sqrt{17} \right) & \frac{1}{2} \left( -3 - \sqrt{17} \right) & 0 \\ 1 & 1 & -1 \\ 1 & 1 & 1 \end{bmatrix}$$

The first two columns are not orthogonal to one another and none of the columns are unit length.

First, the Gram-Schmidt procedure normalizes the length of the first vector

$$\begin{aligned} q_1 &= \frac{v_1}{\sqrt{v_1^T v_1}} \\ &= \begin{bmatrix} \frac{-3 + \sqrt{17}}{\sqrt{34 - 6\sqrt{17}}} \\ \sqrt{\frac{2}{17 - 3\sqrt{17}}} \\ \sqrt{\frac{2}{17 - 3\sqrt{17}}} \end{bmatrix} \\ &\approx \begin{bmatrix} 0.369 \\ 0.657 \\ 0.657 \end{bmatrix} \end{aligned}$$

Then, finds the residuals (null component) of the second vector projected onto $q_1$.[15]

$$\begin{aligned} r_2 &= \left( 1 - q_1 q_1^T \right) v_2 \\ &= \begin{bmatrix} \frac{1}{2} \left( -3 - \sqrt{17} \right) \\ 1 \\ 1 \end{bmatrix} \end{aligned}$$

Now, normalize $r_2$

$$q_2 = \frac{r_2}{\sqrt{r_2^T r_2}}$$

---

[15] Since the $\left( v_1^T v_1 \right)^{-1}$ term is the identity, we omit it in the development.

so that $q_1$ and $q_2$ are orthonormal vectors. Let

$$
\begin{aligned}
Q_{12} &= \begin{bmatrix} q_1 & q_2 \end{bmatrix} \\[4pt]
&= \begin{bmatrix} \dfrac{-3+\sqrt{17}}{\sqrt{34-6\sqrt{17}}} & -\dfrac{3+\sqrt{17}}{\sqrt{34+6\sqrt{17}}} \\[10pt] \sqrt{\dfrac{2}{17-3\sqrt{17}}} & \sqrt{\dfrac{2}{17+3\sqrt{17}}} \\[10pt] \sqrt{\dfrac{2}{17-3\sqrt{17}}} & \sqrt{\dfrac{2}{17+3\sqrt{17}}} \end{bmatrix} \\[10pt]
&\approx \begin{bmatrix} 0.369 & -0.929 \\ 0.657 & 0.261 \\ 0.657 & 0.261 \end{bmatrix}
\end{aligned}
$$

Finally, compute the residuals of $v_3$ projected onto $Q_{12}$

$$
\begin{aligned}
r_3 &= v_3 - Q_{12}Q_{12}^T v_3 \\[4pt]
&= \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}
\end{aligned}
$$

and normalize its length.[16]

$$
q_3 = \frac{r_3}{\sqrt{r_3^T r_3}}
$$

Then,

$$
\begin{aligned}
Q &= \begin{bmatrix} q_1 & q_2 & q_3 \end{bmatrix} \\[4pt]
&= \begin{bmatrix} \dfrac{-3+\sqrt{17}}{\sqrt{34-6\sqrt{17}}} & -\dfrac{3+\sqrt{17}}{\sqrt{34+6\sqrt{17}}} & 0 \\[10pt] \sqrt{\dfrac{2}{17-3\sqrt{17}}} & \sqrt{\dfrac{2}{17+3\sqrt{17}}} & -\dfrac{1}{\sqrt{2}} \\[10pt] \sqrt{\dfrac{2}{17-3\sqrt{17}}} & \sqrt{\dfrac{2}{17+3\sqrt{17}}} & \dfrac{1}{\sqrt{2}} \end{bmatrix} \\[10pt]
&\approx \begin{bmatrix} 0.369 & -0.929 & 0 \\ 0.657 & 0.261 & -0.707 \\ 0.657 & 0.261 & 0.707 \end{bmatrix}
\end{aligned}
$$

and $QQ^T = Q^T Q = I$. If there are more vectors then we continue along the same lines with the fourth vector made orthogonal to the first three vectors (by finding its residual from the projection onto the first three columns) and then normalized to unit length, and so on.

---

[16] Again, $\left(Q_{12}^T Q_{12}\right)^{-1} = I$ so it is omitted in the expression. In this example, $v_3$ is orthogonal to $Q_{12}$ (as well as $v_1$ and $v_2$) so it is unaltered.

### A.5.1  QR decomposition

$QR$ is another important (especially for computation) matrix decomposition. $QR$ combines Gram-Schmidt orthogonalization and Gaussian elimination to factor an $m \times n$ matrix $A$ with linearly independent columns into a matrix composed of orthonormal columns, $Q$ such that $Q^T Q = I$, multiplied by a square, invertible upper triangular matrix $R$. This provides distinct advantages when dealing with projections into the column space of $A$. Recall, this problem takes the form $Ay = b$ where the objective is to find $y$ that minimizes the distance to $b$. Since $A = QR$, we have $QRy = b$ and $R^{-1}Q^T QRy = y = R^{-1}Q^T b$. Next, we summarize the steps for two $QR$ algorithms: the Gram-Schmidt approach and the Householder approach.

### A.5.2  Gram-Schmidt QR algorithm

The Gram-Schmidt algorithm proceeds as described above to form $Q$. Let $a$ denote the first column of $A$ and construct $a_1 = \frac{a}{\sqrt{a^T a}}$ to normalize the first column. Construct the projection matrix for this column, $P_1 = a_1^T a_1$ (since $a_1$ is normalized the inverse of $a_1^T a_1$ is unity so it's dropped from the expression). Now, repeat with the second column. Let $a$ denote the second column of $A$ and make it orthogonal to $a_1$ by redefining it as $a = (I - P_1)\,a$. Then normalize via $a_2 = \frac{a}{\sqrt{a^T a}}$. Construct the projection matrix for this column, $P_2 = a_2^T a_2$. The third column is made orthonormal in similar fashion. Let $a$ denote the third column of $A$ and make it orthogonal to $a_1$ and $a_2$ by redefining it as $a = (I - P_1 - P_2)\,a$. Then normalize via $a_3 = \frac{a}{\sqrt{a^T a}}$. Construct the projection matrix for this column, $P_3 = a_3^T a_3$. Repeat this for all $n$ columns of $A$. $Q$ is constructed by combining the columns $Q = \begin{bmatrix} a_1 & a_2 & \cdots & a_n \end{bmatrix}$ such that $Q^T Q = I$. $R$ is constructed as $R = Q^T A$. To see that this is upper triangular let the columns of $A$ be denoted $A_1, A_2, \ldots, A_n$. Then,

$$Q^T A = \begin{bmatrix} a_1^T A_1 & a_1^T A_2 & \cdots & a_1^T A_n \\ a_2^T A_1 & a_2^T A_2 & \cdots & a_2^T A_n \\ \vdots & \vdots & \ddots & \vdots \\ a_n^T A_1 & a_n^T A_2 & \cdots & a_n^T A_n \end{bmatrix}$$

The terms below the main diagonal are zero since $a_j$ for $j = 2, \ldots, n$ are constructed to be orthogonal to $A_1$, $a_j$ for $j = 3, \ldots, n$ are constructed to be orthogonal to $A_2 = a_2 + P_1 A_2$, and so on.

Notice, how straightforward it is to solve $Ay = b$ for $y$.

$$\begin{aligned} Ay &= b \\ QRy &= b \\ R^{-1}Q^T QRy &= y = R^{-1}Q^T b \end{aligned}$$

### A.5.3  Accounting example

Return to the 4 accounts by 6 journal entries $A$ matrix. This matrix clearly does not have linearly independent columns (or for that matter rows) but we'll drop a redundant row (the last row) and denote the resultant matrix $A_0$. Now, we'll find the $QR$ decomposition of the $6 \times 3$ $A_0^T$, $A_0^T = QR$ by the Gram-Schmidt process.

$$A_0^T = \begin{bmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & 0 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

$$a_1 = \frac{1}{2}\begin{bmatrix} -1 \\ -1 \\ 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}, P_1 = \frac{1}{4}\begin{bmatrix} 1 & 1 & -1 & 1 & 0 & 0 \\ 1 & 1 & -1 & 1 & 0 & 0 \\ -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 1 & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

$$a_2 = \frac{1}{2}\begin{bmatrix} 0.567 \\ -0.189 \\ 0.189 \\ -0.189 \\ -0.756 \\ 0 \end{bmatrix}, P_2 = \frac{1}{4}\begin{bmatrix} 1 & 1 & -1 & 1 & 0 & 0 \\ 1 & 1 & -1 & 1 & 0 & 0 \\ -1 & -1 & 1 & -1 & 0 & 0 \\ 1 & 1 & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

and

$$a_3 = \frac{1}{2}\begin{bmatrix} -0.109 \\ 0.546 \\ 0.218 \\ -0.218 \\ -0.109 \\ -0.764 \end{bmatrix}$$

so that

$$Q = \begin{bmatrix} -0.5 & 0.567 & -0.109 \\ -0.5 & -0.189 & 0.546 \\ 0.5 & 0.189 & 0.218 \\ -0.5 & -0.189 & -0.218 \\ 0 & -0.756 & -0.109 \\ 0 & 0 & -0.764 \end{bmatrix}$$

and

$$R = Q^T A = \begin{bmatrix} 2 & -0.5 & -0.5 \\ 0 & 1.323 & -0.189 \\ 0 & 0 & 1.309 \end{bmatrix}$$

The projection solution to $Ay = x$ or $A_0 y = x_0$ where $x = \begin{bmatrix} 2 \\ 1 \\ -1 \\ -2 \end{bmatrix}$ and

$$x_0 = \begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix} \text{ is } y_{\text{row}} = Q\left(R^T\right)^{-1} x_0 = \frac{1}{6} \begin{bmatrix} 1 \\ -5 \\ 4 \\ -4 \\ -5 \\ 1 \end{bmatrix}.$$

### A.5.4   The Householder QR algorithm

The Householder algorithm is not as intuitive as the Gram-Schmidt algorithm but is computationally more stable. Let $a$ denote the first column of $A$ and $z$ be a vector of zeros except the first element is one. Define $v = a + \sqrt{a^T a} z$ and $H_1 = I - 2 * \frac{vv^T}{v^T v}$. Then, $H_1 A$ puts the first column of $A$ in upper triangular form. Now, repeat the process where $a$ is now defined to be the second column of $H_1 A$ whose first element is set to zero and $z$ is defined to be a vector of zeros except the second element is one. Utilize these components to create $v$ in the same form as before and to construct $H_2$ in the same form as $H_1$. Then, $H_2 H_1 A$ puts the first two columns of $A$ in upper triangular form. Next, we work with the third column of $H_2 H_1 A$ where the first two elements of $a$ are set to zero and repeat for all $n$ columns. When complete, $R$ is constructed from the first $n$ rows of $H_n \cdots H_2 H_1 A$. and $Q^T$ is constructed from the first $n$ rows of $H_n \cdots H_2 H_1$.

### A.5.5   Accounting example

Again, return to the 4 accounts by 6 journal entries $A$ matrix and work with $A_0$. Now, we'll find the $QR$ decomposition of the $6 \times 3$ $A_0^T$, $A_0^T = QR$ by Householder transformation.

$$A_0^T = \begin{bmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & 0 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

In the construction of $H_1$, $a = \begin{bmatrix} -1 \\ -1 \\ 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}$, $z = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $v = \begin{bmatrix} -\frac{1}{2} \\ -1 \\ 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}$, and

$H_1 = \frac{1}{2} \begin{bmatrix} 1 & 1 & -1 & 1 & 0 & 0 \\ 1 & 1 & 1 & -1 & 0 & 0 \\ -1 & 1 & 1 & 1 & 0 & 0 \\ 1 & -1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{bmatrix}$. Then, $H_1 A_0^T = \frac{1}{2} \begin{bmatrix} -4 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & -1 & 1 \\ 0 & 1 & -1 \\ 0 & -2 & 0 \\ 0 & 0 & -2 \end{bmatrix}$.

For the construction of $H_2$, $a = \frac{1}{2} \begin{bmatrix} 0 \\ 1 \\ -1 \\ 1 \\ -2 \\ 0 \end{bmatrix}$, $z = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $v = \begin{bmatrix} 0 \\ 1.823 \\ -0.5 \\ 0.5 \\ -1 \\ 0 \end{bmatrix}$,

and $H_2 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & -0.378 & 0.378 & -0.378 & 0.756 & 0 \\ 0 & 0.378 & 0.896 & 0.104 & -0.207 & 0 \\ 0 & -0.378 & 0.104 & 0.896 & 0.207 & 0 \\ 0 & 0.756 & -0.207 & 0.207 & 0.585 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$. Then, $H_2 H_1 A_0^T =$

$\begin{bmatrix} -2 & 0.5 & 0.5 \\ 0 & -1.323 & 0.189 \\ 0 & 0 & 0.585 \\ 0 & 0 & -0.585 \\ 0 & 0 & 0.171 \\ 0 & 0 & -1 \end{bmatrix}$.

For the construction of $H_3$, $a = \begin{bmatrix} 0 \\ 0 \\ 0.585 \\ -0.585 \\ 0.171 \\ -1 \end{bmatrix}$, $z = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $v = \begin{bmatrix} 0 \\ 0 \\ 1.894 \\ -0.585 \\ 0.171 \\ -1 \end{bmatrix}$,

and $H_3 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -0.447 & 0.447 & -0.130 & 0.764 \\ 0 & 0 & 0.447 & 0.862 & 0.040 & -0.236 \\ 0 & 0 & -0.130 & 0.040 & 0.988 & 0.069 \\ 0 & 0 & 0.764 & -0.236 & 0.069 & 0.597 \end{bmatrix}$. Then, $H_3 H_2 H_1 A_0^T =$

$$
\begin{bmatrix}
-2 & 0.5 & 0.5 \\
0 & -1.323 & 0.189 \\
0 & 0 & -1.309 \\
0 & 0 & 0 \\
0 & 0 & 0 \\
0 & 0 & 0
\end{bmatrix}
. \text{ This leads to } R =
\begin{bmatrix}
-2 & 0.5 & 0.5 \\
0 & -1.323 & 0.189 \\
0 & 0 & -1.309
\end{bmatrix}
$$

$$
\text{and } Q =
\begin{bmatrix}
0.5 & -0.567 & 0.109 \\
0.5 & 0.189 & -0.546 \\
-0.5 & -0.189 & -0.218 \\
0.5 & 0.189 & 0.218 \\
0 & 0.756 & 0.109 \\
0 & 0 & 0.764
\end{bmatrix}
.
$$

Finally, the projection solution to $A_0 y = x_0$ is $y_{\text{row}} = Q \left( R^T \right)^{-1} x_0 =$

$$
\begin{bmatrix}
0.5 & -0.567 & 0.109 \\
0.5 & 0.189 & -0.546 \\
-0.5 & -0.189 & -0.218 \\
0.5 & 0.189 & 0.218 \\
0 & 0.756 & 0.109 \\
0 & 0 & 0.764
\end{bmatrix}
\left( \begin{bmatrix}
-2 & 0.5 & 0.5 \\
0 & -1.323 & 0.189 \\
0 & 0 & -1.309
\end{bmatrix}^T \right)^{-1}
\begin{bmatrix}
2 \\
1 \\
-1
\end{bmatrix}
=
$$

$$
\frac{1}{6}
\begin{bmatrix}
1 \\
-5 \\
4 \\
-4 \\
-5 \\
1
\end{bmatrix}
.
$$

## A.6   Computing eigenvalues

As discussed above, eigenvalues are the characteristic values that ensure $(A - \lambda I)$ has a nullspace for square matrix $A$. That is, $(A - \lambda I) x = 0$ where $x$ is an eigenvector. If an eigenvector can be identified such that $Ax = \lambda x$ then the constant, $\lambda$, is an associated eigenvalue. For instance, if the rows of $A$ have the same sum then $x = \iota$ (a vector of ones) and $\lambda$ equals the sum of any row of $A$.

Further, since the sum of the eigenvalues equals the trace of the matrix and the product of the eigenvalues equals the determinant of the matrix, finding the eigenvalues for small matrices is relatively simple. For instance, eigenvalues of a $2 \times 2$ matrix can be found by solving

$$
\begin{aligned}
\lambda_1 + \lambda_2 &= tr\,(A) \\
\lambda_1 \lambda_2 &= \det\,(A)
\end{aligned}
$$

Alternatively, we can solve the roots or zeroes of the characteristic polynomial. That is, $\det\,(A - \lambda I) = 0$.

**Example 1** *Suppose* $A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$ *then* $tr\,(A) = 5$ *and* $\det\,(A) = 4$. *Therefore,*

$$
\begin{aligned}
\lambda_1 + \lambda_2 &= 5 \\
\lambda_1 \lambda_2 &= 4
\end{aligned}
$$

*which leads to* $\lambda_1 = 4$ *and* $\lambda_2 = 1$. *Likewise, the characteristic polynomial is* $\det\,(A - \lambda I) = (2 - \lambda)(3 - \lambda) - 2 = 0$ *leading to the same solution for* $\lambda$.

However, for larger matrices this approach proves impractical. Hence, we'll explore some alternatives.

### A.6.1   Schur's lemma

Schur's lemma says that while every square matrix may not be diagonalizable, it can be triangularized by some unitary operator $U$.

$$
\begin{aligned}
T &= U^{-1} A U \\
&= U^* A U
\end{aligned}
$$

or

$$
A = U T U^*
$$

where $A$ is the matrix of interest, $T$ is a triangular matrix, and $U$ is unitary so that $U^* U = U U^* = I$ ($U^*$ denotes the complex conjugate transpose of

$U$). Further, since $T$ and $A$ are similar matrices they have the same eigenvalues and the eigenvalues reside on the main diagonal of $T$. To see they are similar matrices recognize they have the same characteristic polynomial.

$$
\begin{aligned}
\det(A - \lambda I) &= \det\left(T - \lambda I\right) \\
&= \det\left(U^* A U - \lambda I\right) \\
&= \det\left(U^* A U - \lambda U^* I U\right) \\
&= \det\left(U^*(A - \lambda I)U\right) \\
&= \det\left(U^*\right)\det\left(A - \lambda I\right)\det\left(U\right) \\
&= 1\det\left(A - \lambda I\right)1 \\
&= \det\left(A - \lambda I\right)
\end{aligned}
$$

Before discussing construction of $T$, we introduce some eigenvalue construction algorithms.

### A.6.2   Power algorithm

The power algorithm is an iterative process for finding the largest absolute value eigenvalue.

1. Let $k_1$ be a vector of ones where the number of elements in the vector equals the number of rows or columns in $A$.

2. Let $k_{t+1} = \frac{Ak_t}{\sqrt{k_t^T A^T A k_t}}$ where $\sqrt{k_t^T A^T A k_t} = norm$.

3. iterate until $|k_{t+1} - k_t| < \varepsilon\iota$ for desired precision $\varepsilon$.

4. $norm$ is the largest eigenvalue of $A$ and $k_t = k_{t+1}$ is it's associated eigenvector.

Clearly, if $k_t = k_{t+1}$ this satisfies the property of eigenvalues and eigenvectors, $Ax = \lambda x$ or $Ak_t = \sqrt{k_t^T A^T A k_t}\, k_t$.

Alternatively, let $\mu_t \equiv \frac{k_t^T A k_t}{k_t^T k_t}$ and scale $Ak_t$ by $\mu_t$ to form $k_{t+1} = \frac{Ak_t}{\mu_t}$. Then, iterate as above. This follows as eigensystems are defined by

$$
Ak_t = \lambda k_t
$$

Now, multiply both sides by $k_t^T$ to generate a quadratic form (scalars on both sides of the equation).

$$
k_t^T A k_t = \lambda k_t^T k_t
$$

Then, isolate the eigenvalue, $\lambda$, by dividing both sides by the right-hand side scalar, $k_t^T k_t$, to produce the result. As $t \to n$,

$$
\mu_t \equiv \frac{k_t^T A k_t}{k_t^T k_t} \to \lambda
$$

**Example 2** *Continue with* $A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$. $k_2 = \frac{Ak_1}{norm_1} = \frac{1}{4\sqrt{2}} \begin{bmatrix} 4 \\ 4 \end{bmatrix} =$

$\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$ $k_3 = \frac{Ak_2}{norm_2} = \frac{1}{4} \begin{bmatrix} \frac{4}{\sqrt{2}} \\ \frac{4}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$ *Hence,* $\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$ *is an eigen-*

*vector and* $norm_2 = 4$ *is the associated (largest) eigenvalue.*

**Example 3 (complex eigenvalues)** *Suppose* $A = \begin{bmatrix} -4 & 2 \\ -2 & -4 \end{bmatrix}$. *The eigen-*

*values are* $\lambda = -4 \pm 2i$ *with norm* $= \sqrt{(-4+2i)(-4-2i)} = 4.472136$ *(not*

*a complex number). The power algorithm settles on the norm but* $Ak_n \neq$

$norm * k_n$. *Try the algorithm again except begin with* $k_1 = \begin{bmatrix} 1 \\ i \end{bmatrix}$. *The algo-*

*rithm converges to the same norm but* $k_n = \begin{bmatrix} -0.4406927 - 0.5529828i \\ 0.5529828 - 0.4406927i \end{bmatrix}$.

*Now,*

$$
\begin{aligned}
Ak_n &= \lambda k_n \\
&\phantom{=} \begin{bmatrix} -4 & 2 \\ -2 & -4 \end{bmatrix} \begin{bmatrix} -0.4406927 - 0.5529828i \\ 0.5529828 - 0.4406927i \end{bmatrix} \\
&= \lambda \begin{bmatrix} -0.4406927 - 0.5529828i \\ 0.5529828 - 0.4406927i \end{bmatrix}
\end{aligned}
$$

*solving for* $\lambda$ *yields* $-4 + 2i$. *Since complex roots always come in conjugate*
*pairs we also know the other eigenvalue,* $-4 - 2i$. *However, the second*

*power algorithm converges very quickly with initial vector* $k_1 = \begin{bmatrix} 1 \\ i \end{bmatrix}$ *to*

$\mu_2 = -4 + 2i$ *and* $k_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{i}{\sqrt{2}} \end{bmatrix}$. *This suggests the second algorithm is more*

*versatile and perhaps converges faster.*

### A.6.3   QR algorithm

The QR algorithm parallels Schur's lemma and supplies a method to com-
pute all eigenvalues.

   1. Compute the factors $Q$, an orthogonal matrix $QQ^T = Q^TQ = I$, and
$R$, a right or upper triangular matrix, such that $A = QR$.
   2. Reverse the factors and denote this $A_1$, $A_1 = RQ$.
   3. Factor $A_1$, $A_1 = Q_1R_1$ then $A_2 = R_1Q_1$.
   4. Repeat until $A_k$ is triangular.

$$
\begin{aligned}
A_{k-1} &= Q_{k-1}R_{k-1} \\
A_k &= R_{k-1}Q_{k-1}
\end{aligned}
$$

The main diagonal elements of $A_k$ are the eigenvalues of $A$.

The connection to Schur's lemma is $RQ = Q^T Q R Q = Q^T A Q = A_1$ so that $A$, $A_1$ and $A_k$ are similar matrices (they have the same eigenvalues).

**Example 4** *Continue with* $A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$. $A_1 = RQ = \begin{bmatrix} 3.4 & -1.8 \\ -0.8 & 1.6 \end{bmatrix}$ *and* $A_{11} = R_{10}Q_{10} = \begin{bmatrix} 4 & -1 \\ 0 & 1 \end{bmatrix}$.[17] *Hence, the eigenvalues of A (and also $A_{10}$) are the main diagonal elements, 4 and 1.*

**Example 5 (complex eigenvalues)** *Suppose* $A = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2 & 3 \\ 0 & -3 & 2 \end{bmatrix}$. *The QR algorithm leaves A unchanged. However, we can work in blocks to solve for the eigenvalues. The first block is simply $B_1 = 5$ (bordered by zeroes in the first row, first column) and 5 is an eigenvalue. The second block is rows 2 and 3 and columns 2 and 3 or $B_2 = \begin{bmatrix} 2 & 3 \\ -3 & 2 \end{bmatrix}$. Now solve the characteristic polynomial for this $2 \times 2$ matrix.*

$$\begin{aligned} -\lambda^2 + 4\lambda - 13 &= 0 \\ \lambda &= 2 \pm 3i \end{aligned}$$

*We can check that each of these three eigenvalues creates a nullspace for $A - \lambda I$.*

$$A - 5I = \begin{bmatrix} 0 & 0 & 0 \\ 0 & -3 & 3 \\ 0 & -3 & -3 \end{bmatrix}$$

*has rank 2 and nullspace or eigenvector* $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$.

$$A - (2 + 3i)I = \begin{bmatrix} 3 - 3i & 0 & 0 \\ 0 & -3i & 3 \\ 0 & -3 & -3i \end{bmatrix}$$

*The second row is a scalar multiple $(-i)$ of the third (and vice versa) and a nullspace or eigenvector is* $\frac{1}{\sqrt{2}} \begin{bmatrix} 0 \\ i \\ -1 \end{bmatrix}$. *Finally,*[18]

$$A - (2 - 3i)I = \begin{bmatrix} 3 - 3i & 0 & 0 \\ 0 & 3i & 3 \\ 0 & -3 & 3i \end{bmatrix}$$

---

[17] Shifting refinements are typically employed to speed convergence (see Strang).

[18] Gauss' fundamental theorem of algebra insures complex roots always come in conjugate pairs so this may be overly pedantic.

*Again, the second row is a scalar multiple ($i$) of the third (and vice versa) and a nullspace or eigenvector is* $\frac{1}{\sqrt{2}} \begin{bmatrix} 0 \\ i \\ 1 \end{bmatrix}$. *Hence, the eigenvalues are* $\lambda = 5, 2 \pm 3i$.

### A.6.4  Schur decomposition

Schur decomposition works similarly.

1. Use one of the above algorithms to find an eigenvalue of $n \times n$ matrix $A$, $\lambda_1$.

2. From this eigenvalue, construct a unit length eigenvector, $x_1$.

3. Utilize Gram-Schmidt to construct a unitary matrix $U_1$ from $n-1$ columns of $A$ where $x_1$ is the first column of $U$. This creates

$$AU_1 = U_1 \begin{bmatrix} \lambda_1 & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix}$$

or

$$U_1^* AU_1 = \begin{bmatrix} \lambda_1 & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix}$$

4. The next step works the same way except with the lower right $(n-1) \times (n-1)$ matrix. then, $U_2$ is constructed from this lower, right block with a one in the upper, left position with zeroes in its row and column.

$$U_2 = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & x_{22} & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & x_{2n} & \cdots & * \end{bmatrix}$$

$$U_2^* U_1^* AU_1 U_2 = \begin{bmatrix} \lambda_1 & * & \cdots & * \\ 0 & \lambda_2 & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & * \end{bmatrix}$$

5. Continue until $T$ is constructed.

$$T = U_{n-1}^* \cdots U_1^* AU_1 \cdots U_{n-1}$$

$$U^* AU = \begin{bmatrix} \lambda_1 & * & \cdots & * \\ 0 & \lambda_2 & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \lambda_n \end{bmatrix}$$

where $U = U_1 \cdots U_{n-1}$. When triangularization is complete, the eigenvalues reside on the main diagonal of $T$.

**Example 6 (not diagonalizable)** *Suppose* $A = \begin{bmatrix} 5 & 0 & 1 \\ 0 & 2 & -3 \\ 0 & -3 & 2 \end{bmatrix}$. *This matrix has repeated eigenvalues $(5, 5, -1)$ and lacks a full set of linearly indepedent eigenvectors therefore it cannot be expressed in diagonalizable form $A = S\Lambda S^{-1}$ (as the latter term doesn't exist). Nonetheless, the Schur decomposition can still be employed to triangularize the matrix. A unit length eigenvector associated with $\lambda = 5$ is $x_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$. Applying Gram-Schmidt to columns two and three of A yields* $U_1 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0.55470 & -0.83205 \\ 0 & -0.83205 & -0.55470 \end{bmatrix}$.

*This leads to*

$$
\begin{aligned}
T_1 &= U_1^* A U_1 \\
&= \begin{bmatrix} 5 & -0.83205 & -0.55470 \\ 0 & 4.76923 & -1.15385 \\ 0 & -1.15385 & -0.76923 \end{bmatrix}
\end{aligned}
$$

*Working with the lower, right $2 \times 2$ block gives*

$$
U_2 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & -0.98058 & -0.19612 \\ 0 & 0.19612 & -0.98058 \end{bmatrix}
$$

*Then,*

$$
\begin{aligned}
T &= U_2^* U_1^* A U_1 U_2 \\
U^* A U &= \begin{bmatrix} 5 & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ 0 & 5 & 0 \\ 0 & 0 & -1 \end{bmatrix}
\end{aligned}
$$

*where* $U = U_1 U_2 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$.

**Example 7 (complex eigenvalues)** *Suppose* $A = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2 & 3 \\ 0 & -3 & 2 \end{bmatrix}$. *We know from example 5 A has complex eigenvalues. Let's explore its Schur decomposition. Again, $\lambda = 5$ is an eigenvalue with corresponding eigenvector $x_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$. Applying Gram-Schmidt to columns two and three of A*

*yields* $U_1 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0.55470 & 0.83205 \\ 0 & -0.83205 & 0.55470 \end{bmatrix}$. *This leads to*

$$T_1 = U_1^* A U_1$$
$$= \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2 & 3 \\ 0 & -3 & 2 \end{bmatrix}$$

*Working with the lower, right $2 \times 2$ block, $\lambda = 2 + 3i$, and associated eigenvector* $x_2 = \begin{bmatrix} 0 \\ \frac{1}{\sqrt{2}}i \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$ *gives*

$$U_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{\sqrt{2}}i & \frac{1}{\sqrt{2}} \\ 0 & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}}i \end{bmatrix}$$

*where* $x_{12} = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{2}}i \\ 0 & -\frac{1}{\sqrt{2}} \end{bmatrix}$ *is applied via Gram-Schmidt to create the third (column) vector of $U_2$ from the third column of $A$, $A_{\cdot 3}$.*[19]

$$A_{\cdot 3} - x_{12}x_{12}^* A_{\cdot 3}$$
$$= \begin{bmatrix} 0 \\ 3 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{2}}i \\ 0 & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & -\frac{1}{\sqrt{2}}i & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0 \\ 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \\ -3i \end{bmatrix}$$

*before normalization and after we have* $\begin{bmatrix} 0 \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}}i \end{bmatrix}$. *Then,*

$$T = U_2^* U_1^* A U_1 U_2$$
$$U^* A U = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2+3i & 0 \\ 0 & 0 & 2-3i \end{bmatrix}$$

*where $U = U_1 U_2 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & -0.5883484 + 0.3922323i & 0.3922323 - 0.5883484i \\ 0 & -0.3922323 - 0.5883484i & -0.5883484 - 0.3922323i \end{bmatrix}$.*
*The eigenvalues lie along the main diagonal of $T$.*

---

[19] Notice, conjugate transpose is employed in the construction of the projection matrix to accommodate complex elements.

## A.7   Some determinant identities

### A.7.1   Determinant of a square matrix

We utilize the fact that

$$
\begin{aligned}
\det(A) &= \det(LU) \\
&= \det(L)\det(U)
\end{aligned}
$$

and the determinant of a triangular matrix is the product of the diagonal elements. Since $L$ has ones along its diagonal, $\det(A) = \det(U)$. Return to the example above

$$
\det\left(AA^T - \lambda I_4\right) = \det
\begin{bmatrix}
4-\lambda & -1 & -1 & -2 \\
-1 & 2-\lambda & 0 & -1 \\
-1 & 0 & 2-\lambda & -1 \\
-2 & -1 & -1 & 4-\lambda
\end{bmatrix}
$$

Factor $AA^T - \lambda I_4$ into its upper and lower triangular components via Gaussian elimination (this step can be computationally intensive).

$$
L =
\begin{bmatrix}
1 & 0 & 0 & 0 \\
\frac{1}{-4+\lambda} & 1 & 0 & 0 \\
\frac{1}{-4+\lambda} & -\frac{1}{7-6\lambda+\lambda^2} & 1 & 0 \\
\frac{2}{-4+\lambda} & \frac{-6+\lambda}{7-6\lambda+\lambda^2} & \frac{-6+\lambda}{6-6\lambda+\lambda^2} & 1
\end{bmatrix}
$$

and

$$
U =
\begin{bmatrix}
4-\lambda & -1 & -1 & -2 \\
0 & \frac{7-6\lambda+\lambda^2}{4-\lambda} & \frac{1}{-4+\lambda} & \frac{6-\lambda}{-4+\lambda} \\
0 & 0 & \frac{12-18\lambda+8\lambda^2-\lambda^3}{7-6\lambda+\lambda^2} & -\frac{12-8\lambda+\lambda^2}{7-6\lambda+\lambda^2} \\
0 & 0 & 0 & -\frac{\lambda\left(24-10\lambda+\lambda^2\right)}{6-6\lambda+\lambda^2}
\end{bmatrix}
$$

The determinant of $A$ equals the determinant of $U$ which is the product of the diagonal elements.

$$
\begin{aligned}
\det\left(AA^T - \lambda I_4\right) &= \det(U) \\
&= (4-\lambda)\left(\frac{7-6\lambda+\lambda^2}{4-\lambda}\right)\left(\frac{12-18\lambda+8\lambda^2-\lambda^3}{7-6\lambda+\lambda^2}\right) \\
&\quad \times \left(-\frac{\lambda\left(24-10\lambda+\lambda^2\right)}{6-6\lambda+\lambda^2}\right)
\end{aligned}
$$

which simplifies as

$$
\det\left(AA^T - \lambda I_4\right) = -48\lambda + 44\lambda^2 - 12\lambda^3 + \lambda^4
$$

Of course, the roots of this equation are the eigenvalues of $A$.

## A.7.2   Identities

Below the notation $|A|$ refers to the determinant of matrix $A$.

**Theorem 8** $\left|\begin{bmatrix} A_{m\times m} & B_{m\times n} \\ C_{n\times m} & D_{n\times n} \end{bmatrix}\right| = |A|\,|D - CA^{-1}B| = |D|\,|A - BD^{-1}C|$
where $A^{-1}$ and $D^{-1}$ exist.

**Proof.**

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A & 0 \\ C & I \end{bmatrix}\begin{bmatrix} I & A^{-1}B \\ 0 & D - CA^{-1}B \end{bmatrix}$$

$$= \begin{bmatrix} I & B \\ 0 & D \end{bmatrix}\begin{bmatrix} A - BD^{-1}C & 0 \\ D^{-1}C & I \end{bmatrix}$$

Since the determinant of a block triangular matrix is the product of the determinants of the diagonal blocks and the determinant of the product of matrices is the product of their determinants,

$$\left|\begin{bmatrix} A_{m\times m} & B_{m\times n} \\ C_{n\times m} & D_{n\times n} \end{bmatrix}\right| = |A|\,|I|\,|I|\,|D - CA^{-1}B| = |D|\,|I|\,|A - BD^{-1}C|\,|I|$$

$$= |A|\,|D - CA^{-1}B| = |D|\,|A - BD^{-1}C|$$

■

**Theorem 9** *For $A$ and $B$ $m \times n$ matrices,*

$$\left|I_n + A^T B\right| = \left|I_m + BA^T\right| = \left|I_n + B^T A\right| = \left|I_m + AB^T\right|$$

**Proof.** Since the determinant of the transpose of a matrix equals the determinant of the matrix,

$$\left|I_n + A^T B\right| = \left|\left(I_n + A^T B\right)^T\right| = \left|I_n + B^T A\right|$$

From theorem 8, $\left|\begin{bmatrix} I_m & -B \\ A^T & I_n \end{bmatrix}\right| = |I|\,|I + A^T I B| = |I|\,|I + BIA^T|$. Hence, $\left|I + A^T B\right| = \left|I + BA^T\right| = \left|\left(I + BA^T\right)^T\right| = \left|I + AB^T\right|$ ■

**Theorem 10** *For vectors $x$ and $y$, $\left|I + xy^T\right| = 1 + y^T x$.*

**Proof.** From theorem 9, $\left|I + xy^T\right| = \left|I + y^T x\right| = 1 + y^T x$. ■

**Theorem 11** $\left|A_{n\times n} + xy^T\right| = |A|\left(1 + y^T A^{-1} x\right)$ *where $A^{-1}$ exists.*

**Proof.**
$$\begin{bmatrix} A & -x \\ y^T & 1 \end{bmatrix} = \begin{bmatrix} A & 0 \\ y^T & 1 \end{bmatrix} \begin{bmatrix} I & -A^{-1}x \\ 0 & 1 + y^T A^{-1}x \end{bmatrix} = \begin{bmatrix} I & -x \\ 0 & 1 \end{bmatrix} \begin{bmatrix} A + x1y^T & 0 \\ 1y^T & 1 \end{bmatrix}.$$

$$\begin{aligned} \left| \begin{bmatrix} A & 0 \\ y^T & 1 \end{bmatrix} \begin{bmatrix} I & -A^{-1}x \\ 0 & 1 + y^T A^{-1}x \end{bmatrix} \right| &= |A| \left( 1 + y^T A^{-1}x \right) \\ &= \left| \begin{bmatrix} I & -x \\ 0 & 1 \end{bmatrix} \begin{bmatrix} A + x1y^T & 0 \\ 1y^T & 1 \end{bmatrix} \right| \\ &= 1 \left| A + xy^T \right| \end{aligned}$$

∎

## A.8   Matrix exponentials and logarithms

For matrices $A$ and $B$, where $e^B = A$, then $B = \ln A$. Further, the matrix exponential is

$$e^B = \sum_{k=0}^{\infty} \frac{1}{k!} B^k$$

Suppose the matrix $A$ is diagonalizable.

$$A = S\Lambda S^{-1}$$

where $\Lambda$ is a diagonal matrix with eigenvalues of $A$ on the diagonal. Then,

$$\Lambda = S^{-1} A S$$

and

$$
\begin{aligned}
\ln A &= S \ln \Lambda S^{-1} \\
e^B &= \sum_{k=0}^{\infty} \frac{1}{k!} S \left( \ln \Lambda \right)^k S^{-1}
\end{aligned}
$$

where $\ln \Lambda = \begin{bmatrix} \ln \lambda_1 & 0 & \cdots & 0 \\ 0 & \ln \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \ln \lambda_n \end{bmatrix}$. From this result we see the log-

arithm of a matrix is well-defined if and only if the matrix is full rank (has a complete set of linearly independent rows and columns or, in other words, is invertible). For example, $\ln \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = Q \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} Q^T = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ where $Q$ is any $2 \times 2$ orthogonal matrix ($QQ^T = Q^T Q = I$).

If $A$ is not diagonalizable, then we work with its Jordan form and in particular, the logarithm of Jordan blocks. A Jordan block has the form

$$B = \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ 0 & 0 & \lambda & 1 & 0 \\ 0 & \vdots & 0 & \lambda & 1 \\ 0 & \cdots & 0 & 0 & \lambda \end{bmatrix}$$

where $\lambda$ is the repeated eigenvalue. This can be written

$$B = \lambda \begin{bmatrix} 1 & \lambda^{-1} & 0 & \cdots & 0 \\ 0 & 1 & \lambda^{-1} & \cdots & 0 \\ 0 & 0 & 1 & \lambda^{-1} & 0 \\ 0 & \vdots & 0 & 1 & \lambda^{-1} \\ 0 & \cdots & 0 & 0 & 1 \end{bmatrix} = \lambda \left( I + K \right)$$

where $K = \begin{bmatrix} 0 & \lambda^{-1} & 0 & \cdots & 0 \\ 0 & 0 & \lambda^{-1} & \cdots & 0 \\ 0 & 0 & 0 & \lambda^{-1} & 0 \\ 0 & \vdots & 0 & 0 & \lambda^{-1} \\ 0 & \cdots & 0 & 0 & 0 \end{bmatrix}$. Since $\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \cdots$, we have

$$
\begin{aligned}
\ln B &= \ln \lambda (I + K) \\
&= \ln \lambda I + \ln (I + K) \\
&= \ln \lambda I + K - \frac{K^2}{2} + \frac{K^3}{3} - \frac{K^4}{4} + \cdots
\end{aligned}
$$

This may not converge for all $K$. However, in the case $B = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, $K = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ and we know from the discussion of generalized eigenvectors $K^2$ (as well as higher powers) $= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$. Hence,

$$
\begin{aligned}
\ln B &= \ln \lambda I + K \\
\ln \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} &= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \\
&= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}
\end{aligned}
$$

# Appendix B
## Iterated expectations

Along with Bayes' theorem (the glue holding consistent probability assessment together), iterated expectations is extensively employed for connecting conditional expectation (regression) results with causal effects of interest.

**Theorem 12** *Law of iterated expectations*

$$E\left[Y\right] = E_X\left[E\left[Y \mid X\right]\right]$$

**Proof.**

$$
\begin{aligned}
E_X\left[E\left[Y \mid X\right]\right] &= \int_{\underline{x}}^{\overline{x}} E\left[Y \mid X\right] f\left(x\right) dx \\
&= \int_{\underline{x}}^{\overline{x}} \left[\int_{\underline{y}}^{\overline{y}} y f\left(y \mid x\right) dy\right] f\left(x\right) dx
\end{aligned}
$$

By Fubini's theorem, we can change the order of integration

$$\int_{\underline{x}}^{\overline{x}} \left[\int_{\underline{y}}^{\overline{y}} y f\left(y \mid x\right) dy\right] f\left(x\right) dx = \int_{\underline{y}}^{\overline{y}} y \left[\int_{\underline{x}}^{\overline{x}} f\left(y \mid x\right) f\left(x\right) dx\right] dy$$

The product rule of Bayes' theorem, $f\left(y \mid x\right) f\left(x\right) = f\left(y, x\right)$, implies iterated expectations can be rewritten as

$$E_X\left[E\left[Y \mid X\right]\right] = \int_{\underline{y}}^{\overline{y}} y \left[\int_{\underline{x}}^{\overline{x}} f\left(y, x\right) dx\right] dy$$

Finally, the summation rule integrates out $X$, $\int_{\underline{x}}^{\overline{x}} f(y,x)\, dx = f(y)$, and produces the result.

$$E_X\left[E\left[Y \mid X\right]\right] = \int_{\underline{y}}^{\overline{y}} yf(y)\, dy = E\left[Y\right]$$

∎

## B.1   Decomposition of variance

**Corollary 1** *Decomposition of variance.*

$$Var\left[Y\right] = E_X\left[Var\left[Y \mid X\right]\right] + Var_X\left[E\left[Y \mid X\right]\right]$$

**Proof.**

$$
\begin{aligned}
Var\left[Y\right] &= E\left[Y^2\right] - E\left[Y\right]^2 \\
&= E_X\left[E\left[Y^2 \mid X\right]\right] - E_X\left[E\left[Y \mid X\right]\right]^2 \\
Var\left[Y\right] &= E_X\left[E\left[Y^2 \mid X\right]\right] - E\left[Y\right]^2 \\
&= E_X\left[Var\left[Y \mid X\right] + E\left[Y \mid X\right]^2\right] - E_X\left[E\left[Y \mid X\right]\right]^2 \\
&= E_X\left[Var\left[Y \mid X\right]\right] + E_X\left[E\left[Y \mid X\right]^2\right] - E_X\left[E\left[Y \mid X\right]\right]^2 \\
Var\left[Y\right] &= E_X\left[Var\left[Y \mid X\right]\right] + Var_X\left[E\left[Y \mid X\right]\right]
\end{aligned}
$$

The second line draws from iterated expectations while the fourth line is the decomposition of the second moment.  ∎

In analysis of variance language, the first term is the residual variation (or variation unexplained)and the second term is the regression variation (or variation explained).

## B.2    Jensen's inequality

For any concave function $g(x)$, $E[g(x)] \leq g(E[x])$. Likewise, for any convex function $h(x)$, $E[h(x)] \geq h(E[x])$. Hence, utility functions exhibiting concavity characterize risk aversion or positive risk premia while utility functions exhibiting convexity characterize risk-seeking preferences or negative risk premia.

Further, Jensen's inequality tells us the geometric mean, $G(x)$, is less than or equal to the arithmetic mean, $A(x)$, with equality only when all outcomes are the same.

$$G(x) \equiv \prod_{i=1}^{n} x_i^{p_i} \leq A(x) \equiv \sum_{i=1}^{n} p_i x_i$$

To see this result, let $g(\cdot)$ be the logarithm (a monotone increasing, concave function) for $x$ nonnegative (if any $x_i = 0$ then the inequality is trivially satisfied as the geometric mean is zero if any $x_i = 0$)

$$E[g(x)] = \sum_{i=1}^{n} p_i \ln x_i \leq g(E[x]) = \ln \sum_{i=1}^{n} p_i x_i$$

Let $p_i = \frac{w_i}{w}$ where $w = \sum_{i=1}^{n} w_i$, then

$$\sum_{i=1}^{n} \frac{w_i}{w} \ln x_i \leq \ln \sum_{i=1}^{n} \frac{w_i}{w} x_i$$

To recover geometric and arithmetic means exponentiate (a monotone increasing function) both sides

$$\sqrt[w]{\prod_{i=1}^{n} x_i^{w_i}} = G(x) \leq \sum_{i=1}^{n} \frac{w_i}{w} x_i = A(x)$$

# Appendix C
## Multivariate normal theory

The Gaussian or normal probability distribution is ubiquitous when the data have continuous support as seen in the Central Limit theorems and the resilience to transformation of Gaussian random variables. When confronted with a vector of variables (multivariate), it often is sensible to think of a joint normal probability assignment to describe their stochastic properties. Let $W = \begin{bmatrix} X \\ Z \end{bmatrix}$ be an $m$-element vector ($X$ has $m_1$ elements and $Z$ has $m_2$ variables such that $m_1 + m_2 = m$) with joint normal probability, then the density function is

$$f_W(w) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left[ -\frac{1}{2}(w-\mu)^T \Sigma^{-1}(w-\mu) \right]$$

where

$$\mu = \begin{bmatrix} \mu_X \\ \mu_Z \end{bmatrix}$$

is the $m$-element vector of means for $W$ and

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{XZ} \\ \Sigma_{ZX} & \Sigma_{ZZ} \end{bmatrix}$$

is the $m \times m$ variance-covariance matrix for $W$ with $m$ linearly independent rows and columns.

Of course, the density integrates to unity, $\int_X \int_Z f_W(w)\, dzdx = 1$. And, the marginal densities are found by integrating out the other variables, for example, $f_X(x) = \int_Z f_W(w)\, dz$ and $f_Z(z) = \int_X f_W(w)\, dx$.

Importantly, as it unifies linear regression, the conditional distributions are also Gaussian.

$$
\begin{aligned}
f_X\left(x \mid Z=z\right) &= \frac{f_W\left(w\right)}{f_Z\left(z\right)} \\
&\sim N\left(E\left[X \mid Z=z\right], Var\left[X \mid Z\right]\right)
\end{aligned}
$$

where

$$
E\left[X \mid Z=z\right] = \mu_X + \Sigma_{XZ}\Sigma_{ZZ}^{-1}\left(z - \mu_Z\right)
$$

and

$$
Var\left[X \mid Z\right] = \Sigma_{XX} - \Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZX}
$$

Also,

$$
f_Z\left(z \mid X=x\right) \sim N\left(E\left[Z \mid X=x\right], Var\left[Z \mid X\right]\right)
$$

where

$$
E\left[Z \mid X=x\right] = \mu_Z + \Sigma_{ZX}\Sigma_{XX}^{-1}\left(x - \mu_X\right)
$$

and

$$
Var\left[Z \mid X\right] = \Sigma_{ZZ} - \Sigma_{ZX}\Sigma_{XX}^{-1}\Sigma_{XZ}
$$

From the conditional expectation (or regression) function we see when the data are Gaussian, linearity imposes no restriction.

$$
E\left[Z \mid X=x\right] = \mu_Z + \Sigma_{ZX}\Sigma_{XX}^{-1}\left(x - \mu_X\right)
$$

is often written

$$
\begin{aligned}
E\left[Z \mid X=x\right] &= \left\{\mu_Z - \Sigma_{ZX}\Sigma_{XX}^{-1}\mu_X\right\} + \left\{\Sigma_{ZX}\Sigma_{XX}^{-1}x\right\} \\
&= \alpha + \beta^T x
\end{aligned}
$$

where $\alpha$ corresponds to an intercept (or vector of intercepts) and $\beta^T x$ corresponds to weighted regressors. Applied linear regression estimates the sample analogs to the above parameters, $\alpha$ and $\beta$.

## C.1   Conditional distribution

Next, we develop more carefully the result for $f_X \left( x \mid Z = z \right)$; the result for $f_Z \left( z \mid X = x \right)$ follows in analogous fashion.

$$
\begin{aligned}
f_X \left( x \mid Z \right) &= \frac{f_W \left( w \right)}{f_Z \left( z \right)} \\[2mm]
&= \frac{\frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left[ -\frac{1}{2} \left( w - \mu \right)^T \Sigma^{-1} \left( w - \mu \right) \right]}{\frac{1}{(2\pi)^{m_2/2} |\Sigma_{ZZ}|^{1/2}} \exp\left[ -\frac{1}{2} \left( z - \mu_Z \right)^T \Sigma_{ZZ}^{-1} \left( z - \mu_z \right) \right]} \\[2mm]
&= \frac{1}{(2\pi)^{m_1/2} |Var\left[ X \mid Z \right]|^{1/2}} \\[1mm]
&\quad \times \exp\left[ -\frac{1}{2} \left( x - E\left[ X \mid Z \right] \right)^T Var\left[ X \mid Z \right]^{-1} \left( x - E\left[ X \mid Z \right] \right) \right]
\end{aligned}
$$

The normalizing constants are identified almost immediately since

$$
\frac{(2\pi)^{m/2}}{(2\pi)^{m_2/2}} = \frac{(2\pi)^{(m_1 + m_2)/2}}{(2\pi)^{m_2/2}} = (2\pi)^{m_1/2}
$$

for the leading term and by theorem 1 in section A.6.2 we have

$$
|\Sigma| = |\Sigma_{ZZ}| \left| \Sigma_{XX} - \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX} \right|
$$

since $\Sigma_{XX}, \Sigma_{ZZ}$, and $\Sigma$ are positive definite, their determinants are positive and their square roots are real. Hence,

$$
\begin{aligned}
\frac{|\Sigma|^{\frac{1}{2}}}{|\Sigma_{ZZ}|^{\frac{1}{2}}} &= \frac{|\Sigma_{ZZ}|^{\frac{1}{2}} \left| \Sigma_{XX} - \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX} \right|^{\frac{1}{2}}}{|\Sigma_{ZZ}|^{\frac{1}{2}}} \\[2mm]
&= \left| \Sigma_{XX} - \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX} \right|^{\frac{1}{2}} \\[2mm]
&= |Var\left[ X \mid Z \right]|^{\frac{1}{2}}
\end{aligned}
$$

This leaves the exponential terms

$$
\begin{aligned}
&\frac{\exp\left[ -\frac{1}{2} \left( w - \mu \right)^T \Sigma^{-1} \left( w - \mu \right) \right]}{\exp\left[ -\frac{1}{2} \left( z - \mu_Z \right)^T \Sigma_{ZZ}^{-1} \left( z - \mu_z \right) \right]} \\[2mm]
&= \exp\left[ -\frac{1}{2} \left( w - \mu \right)^T \Sigma^{-1} \left( w - \mu \right) + \frac{1}{2} \left( z - \mu_Z \right)^T \Sigma_{ZZ}^{-1} \left( z - \mu_z \right) \right]
\end{aligned}
$$

which require a bit more foundation. We begin with a lemma for the inverse of a partitioned matrix.

**Lemma 1** *Let a symmetric, positive definite matrix $H$ be partitioned as* $\begin{bmatrix} A & B^T \\ B & C \end{bmatrix}$ *where $A$ and $C$ are square $n_1 \times n_1$ and $n_2 \times n_2$ positive definite matrices (their inverses exist). Then,*

$$H^{-1} = \begin{bmatrix} \left(A - B^T C^{-1} B\right)^{-1} & -A^{-1} B^T \left(C - BA^{-1} B^T\right)^{-1} \\ -C^{-1} B \left(A - B^T C^{-1} B\right)^{-1} & \left(C - BA^{-1} B^T\right)^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} \left(A - B^T C^{-1} B\right)^{-1} & -\left(A - B^T C^{-1} B\right)^{-1} B^T C^{-1} \\ -C^{-1} B \left(A - B^T C^{-1} B\right)^{-1} & \left(C - BA^{-1} B^T\right)^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} \left(A - B^T C^{-1} B\right)^{-1} & -A^{-1} B^T \left(C - BA^{-1} B^T\right)^{-1} \\ -\left(C - BA^{-1} B^T\right)^{-1} BA^{-1} & \left(C - BA^{-1} B^T\right)^{-1} \end{bmatrix}$$

**Proof.** $H$ is symmetric and the inverse of a symmetric matrix is also symmetric. Hence, the second and third lines follow from symmetry and the first line. Since $H$ is symmetric, positive definite,

$$\begin{aligned} H &= LDL^T \\ &= \begin{bmatrix} I & 0 \\ BA^{-1} & I \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & C - BA^{-1} B^T \end{bmatrix} \begin{bmatrix} I & A^{-1} B^T \\ 0 & I \end{bmatrix} \end{aligned}$$

and

$$\begin{aligned} H^{-1} &= \left(L^T\right)^{-1} D^{-1} L^{-1} \\ &= \begin{bmatrix} I & -A^{-1} B^T \\ 0 & I \end{bmatrix} \begin{bmatrix} A^{-1} & 0 \\ 0 & \left(C - BA^{-1} B^T\right)^{-1} \end{bmatrix} \begin{bmatrix} I & 0 \\ -BA^{-1} & I \end{bmatrix} \end{aligned}$$

Expanding gives

$$H^{-1} = \begin{bmatrix} X & -A^{-1} B^T \left(C - BA^{-1} B^T\right)^{-1} \\ -\left(C - BA^{-1} B^T\right)^{-1} BA^{-1} & \left(C - BA^{-1} B^T\right)^{-1} \end{bmatrix}$$

where

$$X = A^{-1} + A^{-1} B^T \left(C - BA^{-1} B^T\right)^{-1} BA^{-1} = \left(A - B^T C^{-1} B\right)^{-1}$$

The latter equality follows from some linear algebra. Suppose it's  true

$$\left(A - B^T C^{-1} B\right)^{-1} = A^{-1} + A^{-1} B^T \left(C - BA^{-1} B^T\right)^{-1} BA^{-1}$$

pre- and post-multiply both sides by $A$

$$A \left(A - B^T C^{-1} B\right)^{-1} A = A + B^T \left(C - BA^{-1} B^T\right)^{-1} B$$

post multiply both sides by $A^{-1} \left(A - B^T C^{-1} B\right)$

$$\begin{aligned} A &= \left(A - B^T C^{-1} B\right) \\ &\quad + B^T \left(C - BA^{-1} B^T\right)^{-1} BA^{-1} \left(A - B^T C^{-1} B\right) \\ 0 &= -B^T C^{-1} B + B^T \left(C - BA^{-1} B^T\right)^{-1} BA^{-1} \left(A - B^T C^{-1} B\right) \end{aligned}$$

Expanding the right hand side gives

$$
\begin{aligned}
0 \;=\;& -B^T C^{-1} B + B^T \left( C - B A^{-1} B^T \right)^{-1} B \\
& - B^T \left( C - B A^{-1} B^T \right)^{-1} B A^{-1} B^T C^{-1} B
\end{aligned}
$$

Collecting terms gives

$$
0 = -B^T C^{-1} B + B^T \left( C - B A^{-1} B^T \right)^{-1} \left( I - B A^{-1} B^T C^{-1} \right) B
$$

Rewrite $I$ as $C C^{-1}$ and substitute

$$
0 = -B^T C^{-1} B + B^T \left( C - B A^{-1} B^T \right)^{-1} \left( C C^{-1} - B A^{-1} B^T C^{-1} \right) B
$$

Factor

$$
\begin{aligned}
0 \;=\;& -B^T C^{-1} B + B^T \left( C - B A^{-1} B^T \right)^{-1} \left( C - B A^{-1} B^T \right) C^{-1} B \\
0 \;=\;& -B^T C^{-1} B + B^T C^{-1} B = 0
\end{aligned}
$$

This completes the lemma. $\blacksquare$

Now, we write out the exponential terms and utilize the lemma to simplify.

$$
\exp \left[ -\frac{1}{2} (w - \mu)^T \Sigma^{-1} (w - \mu) + \frac{1}{2} (z - \mu_Z)^T \Sigma_{ZZ}^{-1} (z - \mu_z) \right]
$$

$$
= \exp \left\{
\times
\begin{array}{c}
-\frac{1}{2} \left[ \begin{array}{cc} (x - \mu_X)^T & (z - \mu_Z)^T \end{array} \right] \\
\left[ \begin{array}{cc} \Sigma_{XX \cdot Z}^{-1} & -\Sigma_{XX \cdot Z}^{-1} \Sigma_{XZ} \Sigma_{ZZ}^{-1} \\ -\Sigma_{ZZ}^{-1} \Sigma_{ZX} \Sigma_{XX \cdot Z}^{-1} & \Sigma_{ZZ \cdot X}^{-1} \end{array} \right] \left[ \begin{array}{c} x - \mu_X \\ z - \mu_Z \end{array} \right] \\
+\frac{1}{2} (z - \mu_Z)^T \Sigma_{ZZ}^{-1} (z - \mu_z)
\end{array}
\right\}
$$

$$
= \exp \left\{
-\frac{1}{2}
\left[
\begin{array}{c}
(x - \mu_X)^T \Sigma_{XX \cdot Z}^{-1} (x - \mu_X) \\
- (x - \mu_X)^T \Sigma_{XX \cdot Z}^{-1} \Sigma_{XZ} \Sigma_{ZZ}^{-1} (z - \mu_Z) \\
- (z - \mu_Z)^T \Sigma_{ZZ}^{-1} \Sigma_{ZX} \Sigma_{XX \cdot Z}^{-1} (x - \mu_X) \\
+ (z - \mu_Z)^T \Sigma_{ZZ \cdot X}^{-1} (z - \mu_Z) \\
+ \frac{1}{2} (z - \mu_Z)^T \Sigma_{ZZ}^{-1} (z - \mu_z)
\end{array}
\right]
\right\}
$$

$$
= \exp \left\{
-\frac{1}{2}
\left[
\begin{array}{c}
(x - \mu_X)^T \Sigma_{XX \cdot Z}^{-1} (x - \mu_X) \\
- (x - \mu_X)^T \Sigma_{XX \cdot Z}^{-1} \Sigma_{XZ} \Sigma_{ZZ}^{-1} (z - \mu_Z) \\
- (z - \mu_Z)^T \Sigma_{ZZ}^{-1} \Sigma_{ZX} \Sigma_{XX \cdot Z}^{-1} (x - \mu_X) \\
+ (z - \mu_Z)^T \left( \Sigma_{ZZ \cdot X}^{-1} - \Sigma_{ZZ}^{-1} \right) (z - \mu_Z)
\end{array}
\right]
\right\}
$$

where $\Sigma_{XX \cdot Z} = \Sigma_{XX} - \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX}$ and $\Sigma_{ZZ \cdot X} = \Sigma_{ZZ} - \Sigma_{ZX} \Sigma_{XX}^{-1} \Sigma_{XZ}$. From the last term, write

$$
\Sigma_{ZZ \cdot X}^{-1} - \Sigma_{ZZ}^{-1} = \Sigma_{ZZ}^{-1} \Sigma_{ZX} \Sigma_{XX \cdot Z}^{-1} \Sigma_{XZ} \Sigma_{ZZ}^{-1}
$$

To see this utilize the lemma for the inverse of a partitioned matrix. By symmetry

$$\Sigma_{ZZ}^{-1}\Sigma_{ZX}\Sigma_{XX\cdot Z}^{-1} = \Sigma_{ZZ\cdot X}^{-1}\Sigma_{ZX}\Sigma_{XX}^{-1}$$

Post multiply both sides by $\Sigma_{XZ}\Sigma_{ZZ}^{-1}$ and simplify

$$
\begin{aligned}
\Sigma_{ZZ}^{-1}\Sigma_{ZX}\Sigma_{XX\cdot Z}^{-1}\Sigma_{XZ}\Sigma_{ZZ}^{-1} &= \Sigma_{ZZ\cdot X}^{-1}\Sigma_{ZX}\Sigma_{XX}^{-1}\Sigma_{XZ}\Sigma_{ZZ}^{-1} \\
&= \Sigma_{ZZ\cdot X}^{-1}\left(\Sigma_{ZZ} - \Sigma_{ZZ\cdot X}\right)\Sigma_{ZZ}^{-1} \\
&= \Sigma_{ZZ\cdot X}^{-1} - \Sigma_{ZZ}^{-1}
\end{aligned}
$$

Now, substitute this into the exponential component

$$\exp\left\{-\frac{1}{2}\left[
\begin{array}{c}
(x - \mu_X)^T \Sigma_{XX\cdot Z}^{-1}(x - \mu_X) \\
-(x - \mu_X)^T \Sigma_{XX\cdot Z}^{-1}\Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z) \\
-(z - \mu_Z)^T \Sigma_{ZZ}^{-1}\Sigma_{ZX}\Sigma_{XX\cdot Z}^{-1}(x - \mu_X) \\
+(z - \mu_Z)^T \Sigma_{ZZ}^{-1}\Sigma_{ZX}\Sigma_{XX\cdot Z}^{-1}\Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z)
\end{array}
\right]\right\}$$

Combining the first and second terms and combine the third and fourth terms gives

$$\exp\left\{-\frac{1}{2}\left[
\begin{array}{c}
(x - \mu_X)^T \Sigma_{XX\cdot Z}^{-1}\left(x - \mu_X - \Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z)\right) \\
-(z - \mu_Z)^T \Sigma_{ZZ}^{-1}\Sigma_{ZX}\Sigma_{XX\cdot Z}^{-1}\left(x - \mu_X - \Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z)\right)
\end{array}
\right]\right\}$$

Then, since

$$(x - \mu_X)^T - (z - \mu_Z)^T \Sigma_{ZZ}^{-1}\Sigma_{ZX} = \left(x - \mu_X - \Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z)\right)^T$$

combining these two terms simplifies as

$$\exp\left[-\frac{1}{2}(x - E[x \mid Z = z])^T \Sigma_{XX\cdot Z}^{-1}(x - E[x \mid Z = z])\right]$$

where $E[x \mid Z = z] = \mu_X - \Sigma_{XZ}\Sigma_{ZZ}^{-1}(z - \mu_Z)$. Therefore, the result matches the claim for $f_X(x \mid Z = z)$, the conditional distribution of $X$ given $Z = z$ is normally distributed with mean $E[x \mid Z = z]$ and variance $Var[x \mid Z]$.

## C.2   Special case of precision

Now, we consider a special case of Bayesian normal updating expressed in terms of precision of variables (inverse variance) along with variance representation above. Suppose a variable of interest $x$ is observed with error

$$Y = x + \varepsilon$$

where

$$x \sim N\left(\mu_x, \sigma_x^2 = \frac{1}{\tau_x}\right)$$

and

$$\varepsilon \sim N\left(0, \sigma_\varepsilon^2 = \frac{1}{\tau_\varepsilon}\right)$$

$\varepsilon$ independent of $x$, $\sigma_j^2$ refers to variance, and $\tau_j$ refers to precision of variable $j$. This implies

$$
Var\begin{bmatrix} Y \\ x \end{bmatrix} = \begin{bmatrix} E\left[(Y - \mu_x)^2\right] & E\left[(Y - \mu_x)(x - \mu_x)\right] \\ E\left[(x - \mu_x)(Y - \mu_x)\right] & E\left[(x - \mu_x)^2\right] \end{bmatrix}
$$

$$
= \begin{bmatrix} \sigma_x^2 + \sigma_\varepsilon^2 & \sigma_x^2 \\ \sigma_x^2 & \sigma_x^2 \end{bmatrix}.
$$

Then, the posterior or updated distribution for $x$ given $Y = y$ is normal.

$$(x \mid Y = y) \sim N\left(E\left[x \mid Y = y\right], Var\left[x \mid Y\right]\right)$$

where

$$
\begin{aligned}
E\left[x \mid Y = y\right] &= \mu_x + \frac{\sigma_x^2}{\sigma_x^2 + \sigma_\varepsilon^2}(y - \mu_x) \\
&= \frac{\sigma_\varepsilon^2 \mu_x + \sigma_x^2 y}{\sigma_x^2 + \sigma_\varepsilon^2} \\
&= \frac{\tau_x \mu_x + \tau_\varepsilon y}{\tau_x + \tau_\varepsilon}
\end{aligned}
$$

and

$$
\begin{aligned}
Var\left[x \mid Y\right] &= \sigma_x^2 - \frac{\left(\sigma_x^2\right)^2}{\sigma_x^2 + \sigma_\varepsilon^2} \\
&= \frac{\sigma_x^2\left(\sigma_x^2 + \sigma_\varepsilon^2\right) - \left(\sigma_x^2\right)^2}{\sigma_x^2 + \sigma_\varepsilon^2} \\
&= \frac{\sigma_x^2 \sigma_\varepsilon^2}{\sigma_x^2 + \sigma_\varepsilon^2} \\
&= \frac{1}{\tau_x + \tau_\varepsilon}
\end{aligned}
$$

For both the conditional expectation and variance, the penultimate line expresses the quantity in terms of variance and the last line expresses the same quantity in terms of precision. The precision of $x$ given $Y$ is $\tau_{x|Y} = \tau_x + \tau_\varepsilon$.

## C.3   Truncated normal distribution

Suppose we have a continuum of states that map one-to-one into an unbounded random variable, $x$, with mean $\mu$ and variance $\sigma^2$. Our natural (maximum entropy) probability assignment for $x$ is a normal distribution with mean $\mu$ and variance $\sigma^2$. The density function for $x$ is

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty$$

Suppose we have information that partitions the states, and therefore $x$, into two regions around $t$ creating two truncated distributions for $x$. The density functions are

$$f(x \mid x < t) = \frac{1}{\sqrt{2\pi}\sigma F\left(\frac{t-\mu}{\sigma}\right)} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < t$$

and

$$f(x \mid t < x) = \frac{1}{\sqrt{2\pi}\sigma\left[1-F\left(\frac{t-\mu}{\sigma}\right)\right]} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad t < x < \infty$$

where $F(\cdot)$ is the cumulative standard normal distribution. Of course, the rescaling by $F(\cdot)$ normalizes each distribution such that it integrates to one over the region of support.

Often we're interested in the expected value and, possibly, variance of the truncated outcome random variable $x$. First, we state the result then provide brief derivations followed by a numerical example.

Let

$$\ell(t) = -\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)}, \quad -\infty < x < t$$

and

$$u(t) = \frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1-F\left(\frac{t-\mu}{\sigma}\right)}, \quad t < x < \infty$$

where $\phi(\cdot)$ is the standard normal density function with mean zero and variance one. Then,

$$\begin{aligned} E\left[x \mid x < t\right] &= \mu + \sigma\ell(t) \\ &= \mu - \sigma\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} \end{aligned}$$

and

$$\begin{aligned} E\left[x \mid x > t\right] &= \mu + \sigma u(t) \\ &= \mu + \sigma\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)} \end{aligned}$$

Notice, iterated expectations produces the mean of the untruncated random variable.

$$
\begin{aligned}
E_t\left[E\left[x \mid t\right]\right] &= F\left(\frac{t-\mu}{\sigma}\right) E\left[x \mid x < t\right] + \left[1 - F\left(\frac{t-\mu}{\sigma}\right)\right] E\left[x \mid x > t\right] \\
&= F\left(\frac{t-\mu}{\sigma}\right)\left[\mu - \sigma\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)}\right] \\
&\quad + \left[1 - F\left(\frac{t-\mu}{\sigma}\right)\right]\left[\mu + \sigma\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\right] \\
&= \mu
\end{aligned}
$$

Variances for the truncated distributions are

$$
\begin{aligned}
Var\left[x \mid x < t\right] &= \sigma^2\left[1 - \ell\left(t\right)\left(\ell\left(t\right) - \frac{t-\mu}{\sigma}\right)\right] \\
&= \sigma^2\left[1 + \frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)}\left(-\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} - \frac{t-\mu}{\sigma}\right)\right]
\end{aligned}
$$

and

$$
\begin{aligned}
Var\left[x \mid x > t\right] &= \sigma^2\left[1 - u\left(t\right)\left(u\left(t\right) - \frac{t-\mu}{\sigma}\right)\right] \\
&= \sigma^2\left[1 - \frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\left(\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)} - \frac{t-\mu}{\sigma}\right)\right]
\end{aligned}
$$

To derive these results it's convenient to transform variables. Let $z = \frac{x-\mu}{\sigma}$, or $x = \sigma z + \mu$ so that $dx = \sigma dz$ and $f\left(x\right) dx = \sigma f\left(z\right) dz \equiv \phi\left(z\right) dz = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{z^2}{2}\right] dz$.

$$
E\left[x \mid x < t\right] = \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{t} x f\left(x\right) dx
$$

Now, transform from $x$ to $z$ and utilize $\int z \exp\left[-\frac{z^2}{2}\right] dz = -\exp\left[-\frac{z^2}{2}\right]$.

$$
\begin{aligned}
E\left[x \mid x < t\right] &= \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{\frac{t-\mu}{\sigma}} \left(\sigma z + \mu\right) \phi\left(z\right) dz \\
&= \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \left\{\mu \int_{-\infty}^{\frac{t-\mu}{\sigma}} \phi\left(z\right) dz + \sigma \int_{-\infty}^{\frac{t-\mu}{\sigma}} z\phi\left(z\right) dz\right\} \\
&= \frac{F\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} \mu - \frac{\sigma\phi\left(z\right)\big|_{-\infty}^{\frac{t-\mu}{\sigma}}}{F\left(\frac{t-\mu}{\sigma}\right)} \\
&= \mu - \sigma\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} \\
&= \mu - \sigma\ell\left(t\right)
\end{aligned}
$$

Similarly, the upper support expectation is

$$
\begin{aligned}
E\left[x \mid x>t\right] &= \frac{1}{1-F\left(\frac{t-\mu}{\sigma}\right)} \int_{\frac{t-\mu}{\sigma}}^{\infty} \left(\sigma z+\mu\right) \phi\left(z\right) dz \\
&= \frac{1}{1-F\left(\frac{t-\mu}{\sigma}\right)} \left\{\mu \int_{\frac{t-\mu}{\sigma}}^{\infty} \phi\left(z\right) dz + \sigma \int_{\frac{t-\mu}{\sigma}}^{\infty} z\phi\left(z\right) dz\right\} \\
&= \frac{1-F\left(\frac{t-\mu}{\sigma}\right)}{1-F\left(\frac{t-\mu}{\sigma}\right)} \mu - \frac{\sigma\phi\left(z\right)\mid_{\frac{t-\mu}{\sigma}}^{\infty}}{1-F\left(\frac{t-\mu}{\sigma}\right)} \\
&= \mu + \sigma \frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1-F\left(\frac{t-\mu}{\sigma}\right)} \\
&= \mu + \sigma u\left(t\right)
\end{aligned}
$$

Variances of the truncated distributions involve

$$
Var\left[x \mid x<t\right] = \int_{-\infty}^{t} x^2 f\left(x\right) dx - E\left[x \mid x<t\right]^2
$$

and

$$
Var\left[x \mid x>t\right] = \int_{t}^{\infty} x^2 f\left(x\right) dx - E\left[x \mid x>t\right]^2
$$

As we have expressions for the truncated means, we focus on the second moments and then combine the results.

$$
\begin{aligned}
E\left[x^2 \mid x<t\right] &= \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{t} x^2 f\left(x\right) dx \\
&= \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{\frac{t-\mu}{\sigma}} \left(\sigma z+\mu\right)^2 \phi\left(z\right) dz \\
&= \frac{1}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{\frac{t-\mu}{\sigma}} \left[\sigma^2 z^2 + 2\sigma\mu z + \mu^2\right] \phi\left(z\right) dz \\
&= \mu^2 - 2\sigma\mu \frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} + \frac{\sigma^2}{F\left(\frac{t-\mu}{\sigma}\right)} \int_{-\infty}^{\frac{t-\mu}{\sigma}} z^2 \phi\left(z\right) dz
\end{aligned}
$$

Focusing on the last term, integration by parts produces

$$
\begin{aligned}
\int_{-\infty}^{\frac{t-\mu}{\sigma}} z^2 \phi\left(z\right) dz &= \int_{-\infty}^{\frac{t-\mu}{\sigma}} z\left[z\phi\left(z\right)\right] dz \\
&= -z\phi\left(z\right)\mid_{-\infty}^{\frac{t-\mu}{\sigma}} - \int_{-\infty}^{\frac{t-\mu}{\sigma}} -\phi\left(z\right) dz \\
&= -\frac{t-\mu}{\sigma} \phi\left(\frac{t-\mu}{\sigma}\right) + F\left(\frac{t-\mu}{\sigma}\right)
\end{aligned}
$$

Hence,

$$
\begin{aligned}
E\left[x^2 \mid x < t\right] &= \mu^2 - 2\sigma\mu\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)} \\
&\quad + \frac{\sigma^2}{F\left(\frac{t-\mu}{\sigma}\right)}\left[F\left(\frac{t-\mu}{\sigma}\right) - \frac{t-\mu}{\sigma}\phi\left(\frac{t-\mu}{\sigma}\right)\right] \\
&= \mu^2 + \sigma^2 - \sigma^2\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{F\left(\frac{t-\mu}{\sigma}\right)}\frac{t+\mu}{\sigma} \\
&= \mu^2 + \sigma^2 + \sigma^2\ell\left(t\right)\frac{t+\mu}{\sigma}
\end{aligned}
$$

and

$$
\begin{aligned}
Var\left[x \mid x < t\right] &= \mu^2 + \sigma^2 + \sigma^2\ell\left(t\right)\frac{t+\mu}{\sigma} - \left(\mu + \sigma\ell\left(t\right)\right)^2 \\
&= \mu^2 + \sigma^2 + \sigma^2\ell\left(t\right)\frac{t+\mu}{\sigma} \\
&\quad - \left(\mu^2 + 2\mu\sigma\ell\left(t\right) + \sigma^2\ell\left(t\right)^2\right) \\
&= \sigma^2 + \sigma^2\ell\left(t\right)\frac{t+\mu}{\sigma} - \left(2\mu\sigma\ell\left(t\right) + \sigma^2\ell\left(t\right)^2\right) \\
&= \sigma^2\left\{1 - \ell\left(t\right)\left[\ell\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\}
\end{aligned}
$$

Variance for upper support is analogous.

$$
\begin{aligned}
E\left[x^2 \mid x > t\right] &= \frac{1}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\int_t^\infty x^2 f\left(x\right) dx \\
&= \frac{1}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\int_{\frac{t-\mu}{\sigma}}^\infty \left(\sigma z + \mu\right)^2 \phi\left(z\right) dz \\
&= \frac{1}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\int_{\frac{t-\mu}{\sigma}}^\infty \left[\sigma^2 z^2 + 2\sigma\mu z + \mu^2\right] \phi\left(z\right) dz \\
&= \mu^2 + 2\sigma\mu\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)} + \frac{\sigma^2}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\int_{\frac{t-\mu}{\sigma}}^\infty z^2\phi\left(z\right) dz
\end{aligned}
$$

Focusing on the last term, integration by parts produces

$$
\begin{aligned}
\int_{\frac{t-\mu}{\sigma}}^\infty z^2\phi\left(z\right) dz &= \int_{\frac{t-\mu}{\sigma}}^\infty z\left[z\phi\left(z\right)\right] dz \\
&= -z\phi\left(z\right)\Big|_{\frac{t-\mu}{\sigma}}^\infty - \int_{\frac{t-\mu}{\sigma}}^\infty -\phi\left(z\right) dz \\
&= \frac{t-\mu}{\sigma}\phi\left(\frac{t-\mu}{\sigma}\right) + \left[1 - F\left(\frac{t-\mu}{\sigma}\right)\right]
\end{aligned}
$$

Hence,

$$
\begin{aligned}
E\left[x^2 \mid x > t\right] &= \mu^2 + 2\sigma\mu\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)} \\
&\quad + \frac{\sigma^2}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\left\{\frac{t-\mu}{\sigma}\phi\left(\frac{t-\mu}{\sigma}\right) + \left[1 - F\left(\frac{t-\mu}{\sigma}\right)\right]\right\} \\
&= \mu^2 + \sigma^2 + \sigma^2\frac{\phi\left(\frac{t-\mu}{\sigma}\right)}{1 - F\left(\frac{t-\mu}{\sigma}\right)}\frac{t+\mu}{\sigma} \\
&= \mu^2 + \sigma^2 + \sigma^2 u\left(t\right)\frac{t+\mu}{\sigma}
\end{aligned}
$$

and

$$
\begin{aligned}
Var\left[x \mid x > t\right] &= \mu^2 + \sigma^2 + \sigma^2 u\left(t\right)\frac{t+\mu}{\sigma} - \left(\mu + \sigma u\left(t\right)\right)^2 \\
&= \mu^2 + \sigma^2 + \sigma^2 u\left(t\right)\frac{t+\mu}{\sigma} \\
&\quad - \left(\mu^2 + 2\mu\sigma u\left(t\right) + \sigma^2 u\left(t\right)^2\right) \\
&= \sigma^2 + \sigma^2 u\left(t\right)\frac{t+\mu}{\sigma} - \left(2\mu\sigma u\left(t\right) + \sigma^2 u\left(t\right)^2\right) \\
&= \sigma^2 \left\{1 - u\left(t\right)\left[u\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\}
\end{aligned}
$$

The various components are connected via variance decomposition.

$$
Var\left[x\right] = E_t\left[Var\left[x \mid t\right]\right] + Var_t\left[E\left[x \mid t\right]\right]
$$

where

$$
\begin{aligned}
E_t\left[Var\left[x \mid t\right]\right] &= F\left(\frac{t-\mu}{\sigma}\right)\sigma^2\left\{1 - \ell\left(t\right)\left[\ell\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&\quad + \left(1 - F\left(\frac{t-\mu}{\sigma}\right)\right)\sigma^2\left\{1 - u\left(t\right)\left[u\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&= \sigma^2 + \sigma^2\left\{\begin{array}{c}\phi\left(\frac{t-\mu}{\sigma}\right)\left[\ell\left(t\right) - \frac{t-\mu}{\sigma}\right] \\ -\phi\left(\frac{t-\mu}{\sigma}\right)\left[u\left(t\right) - \frac{t-\mu}{\sigma}\right]\end{array}\right\} \\
&= \sigma^2 + \sigma^2\phi\left(\frac{t-\mu}{\sigma}\right)\left(\ell\left(t\right) - u\left(t\right)\right)
\end{aligned}
$$

and

$$
\begin{aligned}
Var_t\left[E\left[x\mid t\right]\right] &= F\left(\frac{t-\mu}{\sigma}\right)\left[\mu+\sigma\ell\left(t\right)-\mu\right]^2 \\
&\quad + \left(1-F\left(\frac{t-\mu}{\sigma}\right)\right)\left[\mu+\sigma u\left(t\right)-\mu\right]^2 \\
&= F\left(\frac{t-\mu}{\sigma}\right)\left[\sigma\ell\left(t\right)\right]^2 + \left(1-F\left(\frac{t-\mu}{\sigma}\right)\right)\left[\sigma u\left(t\right)\right]^2 \\
&= \sigma^2\left\{F\left(\frac{t-\mu}{\sigma}\right)\ell\left(t\right)^2 + \left(1-F\left(\frac{t-\mu}{\sigma}\right)\right)u\left(t\right)^2\right\} \\
&= \sigma^2\phi\left(\frac{t-\mu}{\sigma}\right)\left\{-\ell\left(t\right)+u\left(t\right)\right\}
\end{aligned}
$$

Then,

$$
\begin{aligned}
Var\left[x\right] &= E_t\left[Var\left[x\mid t\right]\right] + Var_t\left[E\left[x\mid t\right]\right] \\
&= \sigma^2 + \sigma^2\phi\left(\frac{t-\mu}{\sigma}\right)\left(\ell\left(t\right)-u\left(t\right)\right) \\
&\quad +\sigma^2\phi\left(\frac{t-\mu}{\sigma}\right)\left\{-\ell\left(t\right)+u\left(t\right)\right\} \\
&= \sigma^2
\end{aligned}
$$

**Example 13** *Suppose $x \sim N\left(\mu=10,\sigma=2\right)$ and the distribution is truncated at $t=5$. The density function at lower support is*

$$
f\left(x\mid x<5\right)=\frac{1}{2\sqrt{2\pi}\left(0.00621\right)}\exp\left[-\frac{\left(x-10\right)^2}{8}\right],\quad -\infty<x<5
$$

*and at upper support is*

$$
f\left(x\mid x>5\right)=\frac{1}{2\sqrt{2\pi}\left(0.99379\right)}\exp\left[-\frac{\left(x-10\right)^2}{8}\right],\quad 5<x<\infty
$$

*Means of the truncated random variable are*

$$
\begin{aligned}
E\left[x\mid x<5\right] &= \mu+\sigma\ell\left(t\right) \\
&= 10-2\frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)} \\
&= 4.35451
\end{aligned}
$$

*and*

$$
\begin{aligned}
E\left[x\mid x>5\right] &= \mu+\sigma u\left(t\right) \\
&= 10+2\frac{\phi\left(\frac{5-10}{2}\right)}{1-F\left(\frac{5-10}{2}\right)} \\
&= 10.03528
\end{aligned}
$$

*Iterated expectations provides a consistency check.*

$$
\begin{aligned}
E_t\left[E\left[x \mid t\right]\right] &= F\left(\frac{5-10}{2}\right) E\left[x \mid x < 5\right] + \left[1 - F\left(\frac{5-10}{2}\right)\right] E\left[x \mid x > 5\right] \\
&= (0.00621)\,4.35451 + (0.99379)\,10.03528 \\
&= 10
\end{aligned}
$$

*Variances of the truncated random variable are*

$$
\begin{aligned}
Var\left[x \mid x < 5\right] &= \sigma^2 \left\{1 - \ell\left(t\right)\left[\ell\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&= 4\left\{1 + \frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)}\left[-\frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)} - \frac{5-10}{2}\right]\right\} \\
&= 0.3558952
\end{aligned}
$$

*and*

$$
\begin{aligned}
Var\left[x \mid x > 5\right] &= \sigma^2 \left\{1 - u\left(t\right)\left[u\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&= 4\left\{1 - \frac{\phi\left(\frac{5-10}{2}\right)}{1 - F\left(\frac{5-10}{2}\right)}\left[\frac{\phi\left(\frac{5-10}{2}\right)}{1 - F\left(\frac{5-10}{2}\right)} - \frac{5-10}{2}\right]\right\} \\
&= 3.822377
\end{aligned}
$$

*Variance decomposition provides a consistency check.*

$$
Var\left[x\right] = E_t\left[Var\left[x \mid t\right]\right] + Var_t\left[E\left[x \mid t\right]\right]
$$

$$
\begin{aligned}
E_t\left[Var\left[x \mid t\right]\right] &= F\left(\frac{t-\mu}{\sigma}\right)\sigma^2\left\{1 - \ell\left(t\right)\left[\ell\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&\quad + \left(1 - F\left(\frac{t-\mu}{\sigma}\right)\right)\sigma^2\left\{1 - u\left(t\right)\left[u\left(t\right) - \frac{t-\mu}{\sigma}\right]\right\} \\
&= F\left(\frac{5-10}{2}\right)4\left\{1 + \frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)}\left[-\frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)} - \frac{5-10}{2}\right]\right\} \\
&\quad + \left(1 - F\left(\frac{5-10}{2}\right)\right)4 \\
&\quad \times \left\{1 - \frac{\phi\left(\frac{5-10}{2}\right)}{1 - F\left(\frac{5-10}{2}\right)}\left[\frac{\phi\left(\frac{5-10}{2}\right)}{1 - F\left(\frac{5-10}{2}\right)} - \frac{5-10}{2}\right]\right\} \\
&= 3.800852
\end{aligned}
$$

*and*

$$
\begin{aligned}
Var_t\left[E\left[x\mid t\right]\right] &= F\left(\frac{t-\mu}{\sigma}\right)\left[\sigma\ell\left(t\right)\right]^2 + \left(1 - F\left(\frac{t-\mu}{\sigma}\right)\right)\left[\sigma u\left(t\right)\right]^2 \\
&= F\left(\frac{5-10}{2}\right)\left[2\left(-\frac{\phi\left(\frac{5-10}{2}\right)}{F\left(\frac{5-10}{2}\right)}\right)\right]^2 \\
&\quad + \left(1 - F\left(\frac{5-10}{2}\right)\right)\left[2\frac{\phi\left(\frac{5-10}{2}\right)}{1 - F\left(\frac{5-10}{2}\right)}\right]^2 \\
&= 0.199148
\end{aligned}
$$

*Thus, we have*

$$
\begin{aligned}
Var\left[x\right] &= E_t\left[Var\left[x\mid t\right]\right] + Var_t\left[E\left[x\mid t\right]\right] \\
&= 3.800852 + 0.199148 \\
&= 4
\end{aligned}
$$

# Appendix D

## Projections and conditional expectations

## D.1   Gauss-Markov theorem

Consider the data generating process ($DGP$):

$$Y = X\beta + \varepsilon$$

where $\varepsilon \sim \left(0, \sigma^2 I\right)$, $X$ is $n \times p$ (with rank $p$), and $E\left[X^T \varepsilon\right] = 0$, or more generally $E\left[\varepsilon \mid X\right] = 0$.

The *Gauss-Markov theorem* states that $b = \left(X^T X\right)^{-1} X^T Y$ is the minimum variance estimator of $\beta$ amongst linear unbiased estimators. Gauss' insight follows from a simple idea. Construct $b$ (or equivalently, the residuals or estimated errors, $e$) such that the residuals are orthogonal to every column of $X$ (recall the objective is to extract all information in $X$ useful for explaining $Y$ — whatever is left over from $Y$ should be unrelated to $X$).

$$X^T e = 0$$

where $e = Y - Xb$. Rewriting the orthogonality condition yields

$$X^T \left(Y - Xb\right) = 0$$

or the normal equations

$$X^T X b = X^T Y$$

Provided $X$ is full column rank, this yields the usual $OLS$ estimator

$$b = \left(X^T X\right)^{-1} X^T Y$$

It is straightforward to show that $b$ is unbiased (conditional on the data $X$).

$$
\begin{aligned}
E\left[b \mid X\right] &= E\left[\left(X^T X\right)^{-1} X^T Y \mid X\right] \\
&= E\left[\left(X^T X\right)^{-1} X^T \left(X\beta + \varepsilon\right) \mid X\right] \\
&= \beta + \left(X^T X\right)^{-1} X^T E\left[\varepsilon \mid X\right] = \beta + 0 = \beta
\end{aligned}
$$

Iterated expectations yields $E\left[b\right] = E_X\left[E\left[b \mid X\right]\right] = E_X\left[\beta\right] = \beta$. Hence, unbiasedness applies unconditionally as well.

$$
\begin{aligned}
Var\left[b \mid X\right] &= Var\left[\left(X^T X\right)^{-1} X^T Y \mid X\right] \\
&= Var\left[\left(X^T X\right)^{-1} X^T \left(X\beta + \varepsilon\right) \mid X\right] \\
&= E\left[\left\{\beta + \left(X^T X\right)^{-1} X^T \varepsilon - \beta\right\}\left\{\left(X^T X\right)^{-1} X^T \varepsilon\right\}^T \mid X\right] \\
&= \left(X^T X\right)^{-1} X^T E\left[\varepsilon\varepsilon^T\right] X \left(X^T X\right)^{-1} \\
&= \sigma^2 \left(X^T X\right)^{-1} X^T I X \left(X^T X\right)^{-1} \\
&= \sigma^2 \left(X^T X\right)^{-1}
\end{aligned}
$$

Now, consider the stochastic regressors case,

$$
Var\left[b\right] = Var_X\left[E\left[b \mid X\right]\right] + E_X\left[Var\left[b \mid X\right]\right]
$$

The first term is zero since $E\left[b \mid X\right] = \beta$ for all $X$. Hence,

$$
Var\left[b\right] = E_X\left[Var\left[b \mid X\right]\right] = \sigma^2 E\left[\left(X^T X\right)^{-1}\right]
$$

the unconditional variance of $b$ can only be described in terms of the average behavior of $X$.

To show that $OLS$ yields the minimum variance linear unbiased estimator consider another linear unbiased estimator $b_0 = LY$ ($L$ replaces $\left(X^T X\right)^{-1} X^T$). Since $E\left[LY\right] = E\left[LX\beta + L\varepsilon\right] = \beta$, $LX = I$.

Let $D = L - \left(X^T X\right)^{-1} X^T$ so that $DY = b_0 - b$.

$$
\begin{aligned}
Var\left[b_0 \mid X\right] &= \sigma^2 \left[D + \left(X^T X\right)^{-1} X^T\right]\left[D + \left(X^T X\right)^{-1} X^T\right]^T \\
&= \sigma^2 \left(\begin{array}{c} DD^T + \left(X^T X\right)^{-1} X^T D^T + DX \left(X^T X\right)^{-1} \\ + \left(X^T X\right)^{-1} X^T X \left(X^T X\right)^{-1} \end{array}\right)
\end{aligned}
$$

Since

$$
LX = I = DX + \left(X^T X\right)^{-1} X^T X, DX = 0
$$

and

$$Var\left[b_0 \mid X\right] = \sigma^2 \left(DD^T + \left(X^T X\right)^{-1}\right)$$

As $DD^T$ is positive semidefinite, $Var\left[b\right]$ (and $Var\left[b \mid X\right]$) is at least as small as any other $Var\left[b_0\right]$ ($Var\left[b_0 \mid X\right]$). Hence, the Gauss-Markov theorem applies to both nonstochastic and stochastic regressors.

**Theorem 14** *Rao-Blackwell theorem. If $\varepsilon \sim N\left(0, \sigma^2 I\right)$ for the above DGP, b has minimum variance of all unbiased estimators.*

Finite sample inferences typically derive from normally distributed errors and $t$ (individual parameters) and $F$ (joint parameters) statistics. Some asymptotic results related to the Rao-Blackwell theorem are as follows. For the Rao-Blackwell $DGP$, $OLS$ is consistent and asymptotic normally ($CAN$) distributed. Since $MLE$ yields $b$ for the above $DGP$ with normally distributed errors, $OLS$ is asymptotically efficient amongst all $CAN$ estimators. Asymptotic inferences allow relaxation of the error distribution and rely on variations of the laws of large numbers and central limit theorems.

## D.2   Generalized least squares (GLS)

For the *DGP*,

$$Y = X\beta + \varepsilon$$

where $\varepsilon \sim N(0, \Sigma)$ and $\Sigma$ is some general $n \times n$ variance-covariance matrix, then the linear least squares estimator or generalized least squares (*GLS*) estimator is

$$b_{GLS} = \left(X^T \Sigma^{-1} X\right)^{-1} X^T \Sigma^{-1} Y$$

with

$$Var\left[b_{GLS} \mid X\right] = \left(X^T \Sigma^{-1} X\right)^{-1}$$

If $\Sigma$ is known, this can be computed by ordinary least squares (*OLS*) following transformation of the variables $Y$ and $X$ via $\Gamma^{-1}$ where $\Gamma$ is a triangular matrix such that $\Sigma = \Gamma\Gamma^T$, say via Cholesky decomposition. Then, the transformed *DGP* is

$$\begin{aligned} \Gamma^{-1}Y &= \Gamma^{-1}X\beta + \Gamma^{-1}\varepsilon \\ y &= x\beta + \epsilon \end{aligned}$$

where $y = \Gamma^{-1}Y$, $x = \Gamma^{-1}X$, and $\epsilon = \Gamma^{-1}\varepsilon \sim N(0, I)$. To see where the identity variance matrix comes from, consider

$$\begin{aligned} Var\left[\Gamma^{-1}\varepsilon\right] &= \Gamma^{-1}Var\left[\varepsilon\right]\left(\Gamma^{-1}\right)^T \\ &= \Gamma^{-1}\Sigma\left(\Gamma^{-1}\right)^T \\ &= \Gamma^{-1}\Gamma\Gamma^T\left(\Gamma^{-1}\right)^T \\ &= \Gamma^{-1}\Gamma\Gamma^T\left(\Gamma^T\right)^{-1} \\ &= I \end{aligned}$$

Hence, estimation involves projection of $y$ onto $x$

$$E\left[y \mid x\right] = xb$$

where

$$\begin{aligned} b &= \left(x^T x\right)^{-1} x^T y \\ &= \left(X^T\left(\Gamma^{-1}\right)^T \Gamma^{-1} X\right)^{-1} X^T\left(\Gamma^{-1}\right)^T \Gamma^{-1} Y \end{aligned}$$

Since $\Sigma^{-1} = \left(\Gamma\Gamma^T\right)^{-1} = \left(\Gamma^{-1}\right)^T \Gamma^{-1}$, we can rewrite

$$b = \left(X^T \Sigma^{-1} X\right)^{-1} X^T \Sigma^{-1} Y$$

which is the *GLS* estimator for $\beta$. Further, $Var\,[b\mid X] = E\left[(b-\beta)\,(b-\beta)^T\right]$.
Since

$$
\begin{aligned}
b-\beta &= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}Y-\beta\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\left(X\beta+\varepsilon\right)-\beta\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}X\beta-\beta\\
&\quad + \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\varepsilon-\beta\\
&= \beta+\left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\varepsilon-\beta\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\varepsilon
\end{aligned}
$$

$$
\begin{aligned}
Var\,[b\mid X] &= E\left[(b-\beta)\,(b-\beta)^T\right]\\
&= E\left[\left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\varepsilon\varepsilon^T\Sigma^{-1}X\left(X^T\Sigma^{-1}X\right)^{-1}\mid X\right]\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}E\left[\varepsilon\varepsilon^T\mid X\right]\Sigma^{-1}X\left(X^T\Sigma^{-1}X\right)^{-1}\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}\Sigma\Sigma^{-1}X\left(X^T\Sigma^{-1}X\right)\\
&= \left(X^T\Sigma^{-1}X\right)^{-1}
\end{aligned}
$$

as indicated above.

## D.3   Recursive least squares

Suppose the analyst expects a series of noisy signals that require filtering to uncover the signals of interest, $\beta$, where the *DGP* is

$$Y = X\beta + \varepsilon, \quad \varepsilon \mid X \sim (0, V)$$

Instead of recalculating everything with each sample, the analyst can employ recursive least squares to achieve the same results. The idea revolves around the design matrix, $X_t$, Fisher's information matrix, $\Im_t$, and weights from the variance-covariance matrix, $V_t$, for the period $t$ draw.

$$\Im_t = \Im_{t-1} + X_t^T V_t^{-1} X_t$$

where the components may be augmented with zeroes so that the matrices conform.

$$\Im_t = \begin{bmatrix} \Im_{t-1} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & X_t^T V_t^{-1} X_t \end{bmatrix}$$

Let $K_t = \Im_t^{-1} X_t^T V_t^{-1}$ represent the gain, $Y_t - X_t b_{t-1}$ be the innovation, and $b_1 = \Im_1^{-1} X_1^T V_1^{-1} Y_1 = \left( X_1^T V_1^{-1} X_1 \right)^{-1} X_1^T V_1^{-1} Y_1$ be the first sample least squares estimate.[1] Then, the recursive least squares estimate is

$$b_t = b_{t-1} + K_t \left( Y_t - X_t b_{t-1} \right)$$

where again $b_{t-1}$ may be augmented with zeroes to conform.

**Example 15 (smooth accruals)** *Suppose the DGP for cash flows is*

$$\begin{aligned} cf_t &= m_t + e_t \\ m_t &= g\, m_{t-1} + \varepsilon_t \end{aligned}$$

*the variance-covariance matrix, $V$, is diagonal, and $\nu = \frac{\sigma_e}{\sigma_\epsilon}$, $m_0$ and $g$ are known. Then, accruals$_{t-1}$ and $cf_t$ are, collectively, sufficient statistics for the mean of cash flows $m_t$ based on the history of cash flows and $g^{t-1}$accruals$_t$ is an efficient statistic for $m_t$*

$$\begin{aligned} [\widehat{m}_t | cf_1, ..., cf_t] &= g^{t-1} accruals_t \\ &= \frac{1}{den_t} \left\{ \frac{num_t}{g^2} cf_t + g^{t-1}\nu^2 den_{t-1} accruals_{t-1} \right\} \end{aligned}$$

*where accruals$_0 = m_0$,* $\begin{bmatrix} den_t \\ num_t \end{bmatrix} = B^t \begin{bmatrix} den_0 \\ num_0 \end{bmatrix} = S\Lambda^t S^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$*, $B =$* $\begin{bmatrix} 1+\nu^2 & \nu^2 \\ g^2 & g^2\nu^2 \end{bmatrix}$*, $\Lambda$ is diagonal matrix with the eigenvalues of $B$, and $S$ is*

---

[1] If prior beliefs regarding $\beta$ are informed then $b_0$ representing priors is included to construct $b_1$ in analogous fashion to other samples.

*a matrix of the corresponding eigenvectors for $B$. The variance of accruals is equal to the variance of the estimate of the mean of cash flows multiplied by $g^{2(t-1)}$; the variance of the estimate of the mean of cash flows equals the coefficient on current cash flow multiplied by $\sigma_e^2$, $Var\left[\widehat{m}_t\right] = \frac{num_t}{den_t g^2}\sigma_e^2$. The development employs recursive least squares. Let $X_1 = \begin{bmatrix} -\nu \\ 1 \end{bmatrix}$ (a $2{\times}1$ matrix), $X_2 = \begin{bmatrix} g\nu & -\nu \\ 0 & 1 \end{bmatrix}$ (a $2 \times 2$ matrix), $X_t = \begin{bmatrix} 0 & \cdots & 0 & g\nu & -\nu \\ 0 & \cdots & 0 & 0 & 1 \end{bmatrix}$ (a $2 \times t$ matrix with $t-2$ leading columns of zeroes), $Y_1 = \begin{bmatrix} -g\nu m_0 \\ cf_1 \end{bmatrix}$, $Y_2 = \begin{bmatrix} 0 \\ cf_2 \end{bmatrix}$, and $Y_t = \begin{bmatrix} 0 \\ cf_t \end{bmatrix}$. The information matrix for a $t$-period cash flow history is*

$$\Im_t = \Im_{t-1}^a + X_t^T X_t$$

$$= \begin{bmatrix} 1+\nu^2+g^2\nu^2 & -g\nu^2 & 0 & \cdots & 0 \\ -g\nu^2 & 1+\nu^2+g^2\nu^2 & -g\nu^2 & \ddots & \vdots \\ 0 & -g\nu^2 & \ddots & -g\nu^2 & 0 \\ \vdots & \ddots & -g\nu^2 & 1+\nu^2+g^2\nu^2 & -g\nu^2 \\ 0 & \cdots & 0 & -g\nu^2 & 1+\nu^2 \end{bmatrix},$$

*a symmetric tri-diagonal matrix, where $\Im_{t-1}^a$ is $\Im_{t-1}$ augmented with a row and column of zeroes to conform with $\Im_t$. For instance, $\Im_1 = \begin{bmatrix} 1+\nu^2 \end{bmatrix}$ and $\Im_1^a = \begin{bmatrix} 1+\nu^2 & 0 \\ 0 & 0 \end{bmatrix}$. The estimate of the mean of cash flows is derived recursively as*

$$\widehat{m}_t = \widehat{m}_{t-1}^a + K_t\left(z_t - X_t^a \widehat{m}_{t-1}^a\right)$$

*for $t > 1$ where $K_t = \Im_t^{-1} X_t^T$, the gain matrix, and $\widehat{m}_{t-1}^a$ is $\widehat{m}_{t-1}$ augmented with a zero to conform with $\widehat{m}_t$. The best linear unbiased estimate of the current mean is the last element in the vector $\widehat{m}_t$ and its variance is the last row-column element of $\Im_t^{-1}$ multiplied by $\sigma_e^2$.*

**Example 16 (special case)** *Suppose $g = \nu = 1$ for the above DGP. Then,*

$$[\widehat{m}_t | cf_1, ..., cf_t] = accruals_t$$

$$= \frac{1}{F_{2t+1}}\left\{F_{2t}cf_t + F_{2t-1}accruals_{t-1}\right\}$$

*and variance for the most recent mean estimate (the $t$th element) is*

$$Var\left[\widehat{m}_t | cf_1, ..., cf_t\right]_{tt} = \sigma^2 \frac{F_{2t}}{F_{2t+1}}$$

*where $F_t = F_{t-1} + F_{t-2}$, the Fibonacci series.*

# Appendix E

## Two stage least squares IV (2SLS-IV)

In this appendix we develop two stage least squares instrumental variable (*2SLS-IV*) estimation more generally. Instrumental variables are attractive strategies whenever the fundamental condition for regression, $E\left[\varepsilon \mid X\right] = 0$, is violated but we can identify instruments such that $E\left[\varepsilon \mid Z\right] = 0$.

## E.1   General case

For the data generating process

$$Y = X\beta + \varepsilon$$

the *IV* estimator projects $X$ onto $Z$ (if $X$ includes an intercept so does $Z$), $\widehat{X} = Z\left(Z^T Z\right)^{-1} Z^T X$, where $X$ is a matrix of regressors, and $Z$ is a matrix of instruments. Then, the *IV* estimator for $\beta$ is

$$b^{IV} = \left(\widehat{X}^T \widehat{X}\right)^{-1} \widehat{X} Y$$

Provided $E\left[\varepsilon \mid Z\right] = 0$ and $Var\left[\varepsilon\right] = \sigma^2 I$, $b^{IV}$ is unbiased, $E\left[b^{IV}\right] = \beta$, and has variance $Var\left[b^{IV} \mid X, Z\right] = \sigma^2 \left(\widehat{X}^T \widehat{X}\right)^{-1}$.

Unbiasedness is demonstrated as follows.

$$
\begin{aligned}
b^{IV} &= \left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}Y \\
&= \left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\left(X\beta+\varepsilon\right) \\
&= \left(X^TZ\left(Z^TZ\right)^{-1}Z^TZ\left(Z^TZ\right)^{-1}Z^TX\right)^{-1}X^TZ\left(Z^TZ\right)^{-1}Z^T\left(X\beta+\varepsilon\right) \\
&= \left(X^TZ\left(Z^TZ\right)^{-1}Z^TX\right)^{-1}\left(X^TZ\left(Z^TZ\right)^{-1}Z^TX\right)\beta \\
&\quad + \left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\varepsilon \\
&= \beta+\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\varepsilon
\end{aligned}
$$

Then,

$$
\begin{aligned}
E\left[b^{IV}\mid X,Z\right] &= E\left[\beta+\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\varepsilon\mid X,Z\right] \\
&= \beta+\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^TE\left[\varepsilon\mid X,Z\right] \\
&= \beta+0=\beta
\end{aligned}
$$

By iterated expectations,

$$
E\left[b^{IV}\right]=E_{X,Z}\left[E\left[b^{IV}\mid X,Z\right]\right]=\beta
$$

Variance of the $IV$ estimator follows similarly.

$$
Var\left[b^{IV}\mid X,Z\right]=E\left[\left(b^{IV}-\beta\right)\left(b^{IV}-\beta\right)^T\mid X,Z\right]
$$

From the above development,

$$
b^{IV}-\beta=\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\varepsilon
$$

Hence,

$$
\begin{aligned}
Var\left[b^{IV}\mid X,Z\right] &= E\left[\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\varepsilon\varepsilon^T\widehat{X}\left(\widehat{X}^T\widehat{X}\right)^{-1}\mid X,Z\right] \\
&= \sigma^2\left(\widehat{X}^T\widehat{X}\right)^{-1}\widehat{X}^T\widehat{X}\left(\widehat{X}^T\widehat{X}\right)^{-1} \\
&= \sigma^2\left(\widehat{X}^T\widehat{X}\right)^{-1}
\end{aligned}
$$

## E.2   Special case

In the special case where both $X$ and $Z$ have the same number of columns the estimator can be further simplified,

$$b^{IV} = \left(Z^T X\right)^{-1} Z^T Y$$

To see this, write out the estimator

$$
\begin{aligned}
b^{IV} &= \left(\widehat{X}^T \widehat{X}\right)^{-1} \widehat{X} Y \\
&= \left(X^T Z \left(Z^T Z\right)^{-1} Z^T X\right)^{-1} X^T Z \left(Z^T Z\right)^{-1} Z^T Y
\end{aligned}
$$

Since $X^T Z$ and $Z^T X$ have linearly independent columns, we can invert each square term

$$
\begin{aligned}
b^{IV} &= \left(Z^T X\right)^{-1} Z^T Z \left(X^T Z\right)^{-1} X^T Z \left(Z^T Z\right)^{-1} Z^T Y \\
&= \left(Z^T X\right)^{-1} Z^T Y
\end{aligned}
$$

Of course, the estimator is unbiased, $E\left[b^{IV}\right] = \beta$, and the variance of the estimator is

$$Var\left[b^{IV} \mid X, Z\right] = \sigma^2 \left(\widehat{X}^T \widehat{X}\right)^{-1}$$

which can be written

$$Var\left[b^{IV} \mid X, Z\right] = \sigma^2 \left(Z^T X\right)^{-1} Z^T Z \left(X^T Z\right)^{-1}$$

in this special case where $X$ and $Z$ have the same number of columns.

# Appendix F
## Seemingly unrelated regression (SUR)

First, we describe the seemingly unrelated regression $(SUR)$ model. Second, we remind ourselves Bayesian regression works as if we have two samples: one representative of our priors and another from the new evidence. Then, we connect to seemingly unrelated regression $(SUR)$ — both classical and Bayesian strategies are summarized.

We describe the $SUR$ model in terms of a stacked regression as if the latent variables in a binary selection setting are observable.

$$r = X\beta + \varepsilon$$

where

$$r = \begin{bmatrix} U^* \\ Y_1 \\ Y_0 \end{bmatrix}, \ X = \begin{bmatrix} W & 0 & 0 \\ 0 & X_1 & 0 \\ 0 & 0 & X_2 \end{bmatrix}, \ \beta = \begin{bmatrix} \theta \\ \beta_1 \\ \beta_2 \end{bmatrix}, \ \varepsilon = \begin{bmatrix} V_D \\ V_1 \\ V_0 \end{bmatrix},$$

and

$$\varepsilon \sim N\left(0, V = \Sigma \bigotimes I_n\right)$$

Let $\Sigma = \begin{bmatrix} 1 & \sigma_{D1} & \sigma_{D0} \\ \sigma_{D1} & \sigma_{11} & \sigma_{10} \\ \sigma_{D0} & \sigma_{10} & \sigma_{00} \end{bmatrix}$, a $3 \times 3$ matrix, then the Kronecker product, denoted by $\otimes$, is

$$
\begin{aligned}
V &= \Sigma \otimes I_n = \begin{bmatrix} I_n & \sigma_{D1} I_n & \sigma_{D0} I_n \\ \sigma_{D1} I_n & \sigma_{11} I_n & \sigma_{10} I_n \\ \sigma_{D0} I_n & \sigma_{10} I_n & \sigma_{00} I_n \end{bmatrix} \\
&= \begin{bmatrix}
1 & \cdots & 0 & \sigma_{D1} & \cdots & 0 & \sigma_{D0} & \cdots & 0 \\
\vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
0 & \cdots & 1 & 0 & \cdots & \sigma_{D1} & 0 & \cdots & \sigma_{D0} \\
\sigma_{D1} & \cdots & 0 & \sigma_{11} & \cdots & 0 & \sigma_{10} & \cdots & 0 \\
\vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
0 & \cdots & \sigma_{D1} & 0 & \cdots & \sigma_{11} & 0 & \cdots & \sigma_{10} \\
\sigma_{D0} & \cdots & 0 & \sigma_{10} & \cdots & 0 & \sigma_{00} & \cdots & 0 \\
\vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
0 & \cdots & \sigma_{D0} & 0 & \cdots & \sigma_{10} & 0 & \cdots & \sigma_{00}
\end{bmatrix}
\end{aligned}
$$

a $3n \times 3n$ matrix and $V^{-1} = \Sigma^{-1} \otimes I_n$.

## F.1   Classical

Classical estimation of $SUR$ follows generalized least square $(GLS)$.

$$
\widehat{\beta} = \left( X^T \left( \Sigma^{-1} \otimes I_n \right) X \right)^{-1} X^T \left( \Sigma^{-1} \otimes I_n \right) r
$$

and

$$
Var \left[ \widehat{\beta} \right] = \left( X^T \left( \Sigma^{-1} \otimes I_n \right) X \right)^{-1}
$$

## F.2   Bayesian

On the other hand, Bayesian analysis employs a Gibbs sampler based on the conditional posteriors as, apparently, the $SUR$ error structure prevents identification of conjugate priors. Recall, from the discussion of Bayesian linear regression with general error structure the conditional posterior for $\beta$ is $p \left( \beta \mid \Sigma, y; \beta_0, \Sigma_\beta \right) \sim N \left( \overline{\beta}, V_\beta \right)$ where

$$
\begin{aligned}
\overline{\beta} &= \left( X_0^T \Sigma_0^{-1} X_0 + X^T \Sigma^{-1} X \right)^{-1} \left( X_0^T \Sigma_0^{-1} X_0 \beta_0 + X^T \Sigma^{-1} X \widehat{\beta} \right) \\
&= \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right)^{-1} \left( \Sigma_\beta^{-1} \beta_0 + X^T \Sigma^{-1} X \widehat{\beta} \right)
\end{aligned}
$$

$$\widehat{\beta} = \left(X^T \Sigma^{-1} X\right)^{-1} X^T \Sigma^{-1} y$$

and

$$
\begin{aligned}
V_\beta &= \left(X_0^T \Sigma_0^{-1} X_0 + X^T \Sigma^{-1} X\right)^{-1} \\
&= \left(\Sigma_\beta^{-1} + X^T \Sigma^{-1} X\right)^{-1}
\end{aligned}
$$

## F.3    Bayesian treatment effect application

For the application of $SUR$ to the treatment effect setting, we replace $y$ with $r = \begin{bmatrix} U & y_1 & y_0 \end{bmatrix}^T$ and $\Sigma^{-1}$ with $\left(\Sigma^{-1} \otimes I_n\right)$ yielding the conditional posterior for $\beta$, $p\left(\beta \mid \Sigma, y; \beta_0, \Sigma_\beta\right) \sim N\left(\overline{\beta}^{SUR}, V_\beta^{SUR}\right)$ where

$$
\begin{aligned}
\overline{\beta}^{SUR} &= \left(X_0^T \Sigma_0^{-1} X_0 + X^T \left(\Sigma^{-1} \otimes I_n\right) X\right)^{-1} \\
&\quad \times \left(X_0^T \Sigma_0^{-1} X_0 \beta_0 + X^T \left(\Sigma^{-1} \otimes I_n\right) X \widehat{\beta}^{SUR}\right) \\
&= \left(\Sigma_\beta^{-1} + X^T \left(\Sigma^{-1} \otimes I_n\right) X\right)^{-1} \left(\Sigma_\beta^{-1} \beta_0 + X^T \left(\Sigma^{-1} \otimes I_n\right) X \widehat{\beta}^{SUR}\right)
\end{aligned}
$$

$$\widehat{\beta}^{SUR} = \left(X^T \left(\Sigma^{-1} \otimes I_n\right) X\right)^{-1} X^T \left(\Sigma^{-1} \otimes I_n\right) r$$

and

$$
\begin{aligned}
V_\beta^{SUR} &= \left(X_0^T \Sigma_0^{-1} X_0 + X^T \left(\Sigma^{-1} \otimes I_n\right) X\right)^{-1} \\
&= \left(\Sigma_\beta^{-1} + X \left(\Sigma^{-1} \otimes I_n\right) X\right)^{-1}
\end{aligned}
$$

# Appendix G
## Maximum likelihood estimation of discrete choice models

The most common method for estimating the parameters of discrete choice models is maximum likelihood. The likelihood is defined as the joint density for the parameters of interest $\theta$ conditional on the data $X_t$. For binary choice models and $D_t = 1$ the contribution to the likelihood is $F\left(X_t^T\theta\right)$, and for $D_t = 0$ the contribution to the likelihood is $1 - F\left(X_t^T\theta\right)$ where these are combined as binomial draws and $F\left(X_t^T\theta\right)$ is the cumulative distribution function evaluated at $X_t^T\theta$. Hence, the likelihood is

$$L\left(\theta|X\right) = \prod_{t=1}^{n} F\left(X_t^T\theta\right)^{D_t}\left[1 - F\left(X_t^T\theta\right)\right]^{1-D_t}$$

The log-likelihood is

$$\ell\left(\theta|X\right) \equiv logL\left(\theta|X\right) = \sum_{t=1}^{n} D_t log\left(F\left(X_t^T\theta\right)\right) + \left(1 - D_t\right) log\left(1 - F\left(X_t^T\theta\right)\right)$$

Since this function for binary response models like probit and logit is globally concave, numerical maximization is straightforward. The first order conditions for a maximum, $\max_{\theta} \ell\left(\theta|X\right)$, are

$$\sum_{t=1}^{n} \frac{D_t f\left(X_t^T\theta\right)X_{it}}{F\left(X_t^T\theta\right)} - \frac{\left(1-D_t\right)f\left(X_t^T\theta\right)X_{ti}}{1-F\left(X_t^T\theta\right)} = 0 \quad i = 1,\ldots,k$$

where $f\left(\cdot\right)$ is the density function. Simplifying yields

$$\sum_{t=1}^{n} \frac{\left[D_t - F\left(X_t^T\theta\right)\right]f\left(X_t^T\theta\right)X_{ti}}{F\left(X_t^T\theta\right)\left[1-F\left(X_t^T\theta\right)\right]} = 0 \quad i = 1,\ldots,k$$

Estimates of $\theta$ are found by solving these first order conditions iteratively or, in other words, numerically.

A common estimator for the variance of $\hat{\theta}_{MLE}$ is the negative inverse of the Hessian matrix evaluated at $\hat{\theta}_{MLE}$, $\left[ -H\left(D, \hat{\theta}\right) \right]^{-1}$. Let $H\left(D, \theta\right)$ be the Hessian matrix for the log-likelihood with typical element $H_{ij}\left(D, \theta\right) \equiv \frac{\partial^2 \ell_t(D, \theta)}{\partial \theta_i \partial \theta_j}$.[1]

---

[1] Details can be found in numerous econometrics references and chapter 4 of *Accounting and Causal Effects: Econometric Challenges.*

# Appendix H
## Optimization

In this appendix, we briefly review optimization. First, we'll take up linear programming then we'll review nonlinear programming.[1]

## H.1   Linear programming

A linear program $(LP)$ is any optimization frame which can be described by a linear objective function and linear constraints. Linear refers to choice variables, say $x$, of no greater than first degree (affine transformations which allow for parallel lines are included). Prototypical examples are

$$\max_{x \geq 0} \quad \pi^T x$$
$$s.t. \quad Ax \leq r$$

or

$$\min_{x \geq 0} \quad \pi^T x$$
$$s.t. \quad Ax \geq r$$

---

[1] For additional details, consult, for instance, Luenberger and Ye, 2010, *Linear and Nonlinear Programming*, Springer, or Luenberger, 1997 *Optimization by Vector Space Methods*, Wiley.

### H.1.1    basic solutions or extreme points

Basic solutions are typically determined from the standard form for an *LP*. Standard form involves equality constraints except non-negative choice variables, $x \geq 0$. That is, $Ax \leq r$ is rewritten in terms of slack variables, $s$, such that $Ax + s = r$. The solution to this program is the same as the solution to the inequality program.

A basic solution or extreme point is determined from an $m \times m$ submatrix of $A$ composed of $m$ linearly independent columns of $A$. The set of basic feasible solutions then is the collection of all basic solutions involving $x \geq 0$.

Consider an example. Suppose $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$, $r = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$, $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ and $s = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$. Then, $Ax + s = r$ can be written $Bx_s = r$ where $B = \begin{bmatrix} A & I_2 \end{bmatrix}$, $I_2$ is a $2 \times 2$ identity matrix, and $x_s = \begin{bmatrix} x \\ s \end{bmatrix}$. The matrix $B$ has two linearly independent columns so each basic solution works with two columns of $B$, say $B_{ij}$, and the elements other than $i$ and $j$ of $x_s$ are set to zero. For instance, $B_{12}$ leads to basic solution $x_1 = \frac{2r_1 - r_2}{3}$ and $x_2 = \frac{2r_2 - r_1}{3}$. The basic solutions are tabulated below.

| $B_{ij}$ | $x_1$ | $x_2$ | $s_1$ | $s_2$ |
|---|---|---|---|---|
| $B_{12}$ | $\frac{2r_1 - r_2}{3}$ | $\frac{2r_2 - r_1}{3}$ | $0$ | $0$ |
| $B_{13}$ | $r_2$ | $0$ | $r_1 - 2r_2$ | $0$ |
| $B_{14}$ | $\frac{r_1}{2}$ | $0$ | $0$ | $\frac{2r_2 - r_1}{2}$ |
| $B_{23}$ | $0$ | $\frac{r_2}{2}$ | $\frac{2r_1 - r_2}{2}$ | $0$ |
| $B_{24}$ | $0$ | $r_1$ | $0$ | $r_2 - 2r_1$ |

To test feasibility consider specific values for $r$. Suppose $r_1 = r_2 = 10$. The table with a feasibility indicator $(1 \, (x_s \geq 0))$ becomes

| $B_{ij}$ | $x_1$ | $x_2$ | $s_1$ | $s_2$ | feasible |
|---|---|---|---|---|---|
| $B_{12}$ | $\frac{10}{3}$ | $\frac{10}{3}$ | $0$ | $0$ | yes |
| $B_{13}$ | $10$ | $0$ | $-10$ | $0$ | no |
| $B_{14}$ | $5$ | $0$ | $0$ | $5$ | yes |
| $B_{23}$ | $0$ | $5$ | $5$ | $0$ | yes |
| $B_{24}$ | $0$ | $10$ | $0$ | $-10$ | no |

Notice, when $x_2 = 0$ there is slack in the second constraint $(s_2 > 0)$ and similarly when $x_1 = 0$ there is slack in the first constraint $(s_1 > 0)$. Basic feasible solutions, an algebraic concept, are also referred to by their geometric counterpart, extreme points.

Identification of basic feasible solutions or extreme points combined with the fundamental theorem of linear programming substantially reduce the search for an optimal solution.

## *H.1.2   fundamental theorem of linear programming*

For a linear program in standard form where $A$ is an $m \times n$ matrix of rank $m$,

i) if there is a feasible solution, there is a basic feasible solution;

ii) if there is an optimal solution, there is a basic feasible optimal solution.

Further, if more than one basic feasible solution is optimal, the edge between the basic feasible optimal solutions is also optimal. The theorem means the search for the optimal solution can be restricted to basic feasible solutions — a finite number of points.

## *H.1.3   duality theorems*

Optimality programs come in pairs. That is, there is a complementary or dual program to the primary (primal) program. For instance, the dual to the maximization program is a minimization program, and vice versa.

$$
\begin{array}{cc}
\text{primal program} & \text{dual program} \\
\max_{x \geq 0} \quad \pi^T x & \min_{\lambda \geq 0} \quad r^T \lambda \\
s.t. \quad Ax \leq r & s.t. \quad A^T \lambda \geq \pi
\end{array}
$$

or

$$
\begin{array}{cc}
\text{primal program} & \text{dual program} \\
\min_{x \geq 0} \quad \pi^T x & \max_{\lambda \geq 0} \quad r^T \lambda \\
s.t. \quad Ax \geq r & s.t. \quad A^T \lambda \leq \pi
\end{array}
$$

where $\lambda$ is a vector of shadow prices or dual variable values. The dual of the dual program is the primal program.

strong duality theorem

If either the primal or dual has an optimal solution so does the other and their optimal objective function values are equal. If one of the programs is unbounded the other has no feasible solution.

weak duality theorem

For feasible solutions, the objective function value of the minimization program (say, dual) is greater than or equal to the maximization program (say, primal).

The intuition for the duality theorems is straightforward. Begin with the constraints

$$Ax \leq r \quad A^T \lambda \geq \pi$$

Transposing both sides of the first constraint leaves the inequality unchanged.

$$x^T A^T \leq r^T \quad A^T \lambda \geq \pi$$

Now, post-multiply both sides of the first constraint by $\lambda$ and pre-multiply both sides of the second constraint by $x^T$, since both $\lambda$ and $x$ are nonnegative the inequality is preserved.

$$x^T A^T \lambda \leq r^T \lambda \quad x^T A^T \lambda \geq x^T \pi$$

Since $x^T \pi$ is a scalar, $x^T \pi = \pi^T x$. Now, combine the results and we have the relation we were after.

$$\pi^T x \leq x^T A^T \lambda \leq r^T \lambda$$

The solution to the dual lies above that for the primal except when they both reside at the optimum solution, in which case their objective function values are equal.

### H.1.4   example

Suppose we wish to solve

$$\max_{x \geq 0} \quad 10x + 12y$$
$$s.t \quad \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \leq \begin{bmatrix} 10 \\ 10 \end{bmatrix}$$

We only need evaluate the objective function at each of the basic feasible solutions we earlier identified: $10\left(\frac{10}{3}\right) + 12\left(\frac{10}{3}\right) = \frac{220}{3}$, $10(5) + 12(0) = 50 = \frac{150}{3}$, and $10(0) + 12(5) = 60 = \frac{180}{3}$. The optimal solution is $x = y = \frac{10}{3}$.

### H.1.5   complementary slackness

Suppose $x \geq 0$ is an $n$ element vector containing a feasible primal solution, $\lambda \geq 0$ is an $m$ element vector containing a feasible dual solution, $s \geq 0$ is an $m$ element vector containing primal slack variables, and $t \geq 0$ is an $n$ element vector containing dual slack variables. Then, $x$ and $\lambda$ are optimal if and only if (element-by-element)

$$xt = 0$$

and

$$\lambda s = 0$$

These conditions are economically sensible as either the scarce resource is exhausted ($s = 0$) or if the resource is plentiful it has no value ($\lambda = 0$).

# H.2   Nonlinear programming

## H.2.1   unconstrained

Nonlinear programs involve nonlinear objective functions. For instance,

$$\max_{x} \quad f(x)$$

If the function is continuously differentiable, then a local optimum can be found by the first order approach. That is, equate the gradient (a vector of partial derivatives composed of terms, $\frac{\partial f(x)}{\partial x_i}$, $i = 1, \ldots, n$ where there are $n$ choice variables, $x$).

$$\nabla f(x^*) = 0$$

$$\begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

Second order (sufficient) conditions involve the Hessian, a matrix of second partial derivatives.

$$H(x^*) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_n} \end{bmatrix}$$

For a local minimum, the Hessian is positive definite (the eigenvalues of $H$ are positive). While for a local maximum, the Hessian is negative definite (the eigenvalues of $H$ are negative).

## H.2.2   convexity and global minima

If $f$ is a convex function (defined below), then the set where $f$ achieves its local minimum is convex and any local minimum is a global minimum. A function $f$ is convex if for every $x_1$, $x_2$, and $\alpha$ , $0 \le \alpha \le 1$,

$$f(\alpha x_1 + (1 - \alpha) x_2) \le \alpha f(x_1) + (1 - \alpha) f(x_2)$$

If $x_1 \ne x_2$, and $0 < \alpha < 1$,

$$f(\alpha x_1 + (1 - \alpha) x_2) < \alpha f(x_1) + (1 - \alpha) f(x_2)$$

then $f$ is strictly convex. If $g = -f$ and $f$ is (strictly) convex, then $g$ is (strictly) concave.

### H.2.3   example

Suppose we face the problem

$$\min_{x,y} \quad f(x,y) = x^2 - 10x + y^2 - 10y + xy$$

The first order conditions are

$$\nabla f(x,y) = 0$$

or

$$\frac{\partial f}{\partial x} = 2x - 10 + y = 0$$

$$\frac{\partial f}{\partial y} = 2y - 10 + x = 0$$

Since the problem is quadratic and the gradient is composed of linearly independent equations, a unique solution is immediately identifiable.

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 10 \\ 10 \end{bmatrix}$$

or $x = y = \frac{10}{3}$ with objective function value $-\frac{100}{3}$. As the Hessian is positive definite, this solution is a minimum.

$$H = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

Positive definiteness of the Hessian follows as the eigenvalues of $H$ are positive. To see this, recall the sum of the eigenvalues equals the trace of the matrix and the product of the eigenvalues equals the determinant of the matrix. The eigenvalues of $H$ are 1 and 3, both positive.

### H.2.4   constrained — the Lagrangian

Nonlinear programs involve either nonlinear objective functions, constraints, or both. For instance,

$$\max_{x \geq 0} \quad f(x)$$
$$s.t. \quad G(x) \leq r$$

Suppose the objective function and constraints are continuously differentiable concave and an optimal solution exists, then the optimal solution can be found via the Lagrangian. The Lagrangian writes the objective function less a Lagrange multiplier times each of the constraints. As either the multiplier is zero or the constraint is binding, each constraint term equals zero.

$$\mathcal{L} = f(x) - \lambda_1 [g_1(x) - r_1] - \cdots - \lambda_n [g_n(x) - r_n]$$

where $G(x)$ involves $n$ functions, $g_i(x)$, $i = 1, \ldots, n$. Suppose $x$ involves $m$ choice variables. Then, there are $m$ Lagrange equations plus the $n$ constraints that determine the optimal solution.

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial x_1} &= 0 \\
&\vdots \\
\frac{\partial \mathcal{L}}{\partial x_m} &= 0 \\
\lambda_1 \left[ g_1(x) - r_1 \right] &= 0 \\
&\vdots \\
\lambda_n \left[ g_n(x) - r_n \right] &= 0
\end{aligned}
$$

The Lagrange multipliers (shadow prices or dual variable values) represent the rate of change in the optimal objective function for each of the constraints.

$$
\lambda_i = \frac{\partial f(r^*)}{\partial r_i}
$$

where $r^*$ refers to rewriting the optimal solution $x^*$ in terms of the constraint values, $r$. If a constraint is not binding, it's multiplier is zero as it has no impact on the optimal objective function value.[2]

## H.2.5   Karash-Kuhn-Tucker conditions

Originally, the Lagrangian only allowed for equality constraints. This was generalized to include inequality constraints by Karash and separately Kuhn and Tucker. Of course, some regularity conditions are needed to ensure optimality. Various necessary and sufficient conditions have been proposed to deal with the most general settings. The Karash-Kuhn-Tucker theorem supplies first order necessary conditions for a local optimum (gradient of the Lagrangian and the Lagrange multiplier times the inequality constraint equal zero when the Lagrange multipliers on the inequality constraints are non-negative evaluated at $x^*$). Second order necessary (positive semi-definite Hessian for the Lagrangian at $x^*$) and sufficient (positive definite Hessian for the Lagrangian at $x^*$) conditions are roughly akin to those for unconstrained local minima.

---

[2] Of course, these ideas regarding the multipliers apply to the shadow prices of linear programs as well.

### H.2.6   example

Continue with the unconstrained example above with an added constraint.

$$\min_{x,y} \quad f(x,y) = x^2 - 10x + y^2 - 10y + xy$$
$$s.t. \qquad\qquad xy \geq 10$$

Since the unconstrained solution satisfies this constraint, $xy = \frac{100}{9} > 10$, the solution remains $x = y = \frac{10}{3}$.

However, suppose the problem is

$$\min_{x,y} \quad f(x,y) = x^2 - 10x + y^2 - 10y + xy$$
$$s.t. \qquad\qquad xy \geq 20$$

The constraint is now active. The Lagrangian is

$$\mathcal{L} = x^2 - 10x + y^2 - 10y + xy - \lambda(xy - 20)$$

and the first order conditions are

$$\nabla \mathcal{L} = 0$$

or

$$\frac{\partial \mathcal{L}}{\partial x} = 2x - 10 + y - \lambda y = 0$$
$$\frac{\partial \mathcal{L}}{\partial y} = 2y - 10 + x - \lambda x = 0$$

and constraint equation

$$\lambda(xy - 20) = 0$$

A solution to these nonlinear equations is

$$x = y = 2\sqrt{5}$$
$$\lambda = 3 - \sqrt{5}$$

The objective function value is $-29.4427$ which, of course, is greater than the objective function value for the unconstrained problem, $-33.3333$. The Hessian is

$$H = \begin{bmatrix} 2 & 1-\lambda \\ 1-\lambda & 2 \end{bmatrix}$$

with eigenvalues evaluated at the solution, $3 - \lambda = \sqrt{5}$ and $1 + \lambda = 4 - \sqrt{5}$, both positive. Hence, the solution is a minimum.

## H.3    Theorem of the separating hyperplane

The theorem of the separating hyperplane states either there exists a non-negative $y$ such that $Ay = x$ or there exists $\lambda$ such that $A^T\lambda \geq 0$ and $\lambda^T x < 0$. The theorem is about mutual exclusivity — one or the other is true not both. This is similar to the way in which orthogonal complements are mutually exclusive. If one subspace contains all positive vectors the orthogonal complement cannot contain positive vectors. Otherwise, their inner products would be positive and inner products of orthogonal subspaces are zero.

The intuition follows from the idea that vector inner products are proportional to the cosine of the angle between them; if the angle is less (greater) than 90 degrees the cosine is positive (negative). $A^T\lambda \geq 0$ means the columns of A are less than or equal to 90 degrees relative a fixed vector $\lambda$ while $\lambda^T x < 0$ implies the angle between $x$ and $\lambda$ exceeds 90 degrees. The separating hyperplane (hyper simply refers to high dimension) is composed of all vectors orthogonal to a fixed vector $\lambda$.

Consider a simple example.

**Example 17 (simple example)** *Suppose* $A = I$ *and* $x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$, *then* $y = 2\begin{bmatrix} 1 \\ 0 \end{bmatrix} + 3\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ *and there exists no* $\lambda$ *from which to form a plane separating the positive quadrant from* $x$. *On the other hand, suppose* $x = \begin{bmatrix} -2 \\ 3 \end{bmatrix}$ — $x$ *lies outside the positive quadrant and* $\lambda = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ *satisfies the theorem's alternative,* $A^T\lambda = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \geq 0$ *and* $\lambda^T x = -2 < 0$.

Next, consider a simple accounting (incidence matrix) example. That is, a case in which $A$ has a nullspace and a left nullspace.

**Example 18 (simple accounting example)** *If* $A = \begin{bmatrix} -1 & 0 & 1 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \end{bmatrix}$ *and* $x = \begin{bmatrix} 1 \\ 2 \\ -3 \end{bmatrix}$, *then* $y = \begin{bmatrix} 2 \\ 0 \\ 3 \end{bmatrix} + k\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ *and any* $k \geq 0$ *satisfies* $Ay = x$ *and* $y \geq 0$. *Hence, there exists no separating plane based on* $\lambda$. *On the other hand, suppose* $A = \begin{bmatrix} -1 & 0 & -1 \\ 1 & -1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$. *Now,* $y = \begin{bmatrix} 2 \\ 0 \\ -3 \end{bmatrix} +$ $k\begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$ *and no* $k$ *satisfies* $Ay = x$ *and* $y \geq 0$. *Any number of* $\lambda$*s exist.*

For example, $\lambda = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ produces $A^T \lambda = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \geq 0$ and $\lambda^T x = -3 < 0$.

Hence, any $\lambda = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} + k \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ (drawing on the left nullspace of A) where $k > -1$ satisfies the theorem's alternative.

# Appendix I
## Quantum information

Quantum information follows from physical theory and experiments. Unlike classical information which is framed in the language and algebra of set theory, quantum information is framed in the language and algebra of vector spaces. To begin to appreciate the difference, consider a set versus a tuple (or vector). For example, the tuple $\begin{bmatrix} 1 & 1 \end{bmatrix}$ is different than the tuple $\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$ but the sets, $\{1, 1\}$ and $\{1, 1, 1\}$ are the same as set $\{1\}$.

Quantum information is axiomatic. Its richness and elegance is demonstrated in that only four axioms make it complete.[1]

## I.1  Quantum information axioms

### I.1.1  The superposition axiom

A quantum unit (qubit) is specified by a two element vector, say $\begin{bmatrix} \alpha \\ \beta \end{bmatrix}$, with $|\alpha|^2 + |\beta|^2 = 1$.

---

[1]Some argue that the measurement axiom is contained in the transformation axiom, hence requiring only three axioms. Even though we appreciate the merits of the argument, we'll proceed with four axioms. Feel free to count them as only three if you prefer.

Let $|\psi\rangle \equiv \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \alpha\,|0\rangle + \beta\,|1\rangle,^{2}\ \langle\psi| = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}^{\dagger}$ where $\dagger$ is the adjoint (conjugate transpose) operation.

### I.1.2   The transformation axiom

A transformation of a quantum unit is accomplished by unitary (length-preserving) matrix multiplication. The Pauli matrices provide a basis of unitary operators.

$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \qquad X = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$Y = \begin{bmatrix} 0 & -i \\ i & 0 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

where $i = \sqrt{-1}$. The operations work as follows: $I\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}, X\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \beta \\ \alpha \end{bmatrix}, Y\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = -\begin{bmatrix} \beta i \\ \alpha i \end{bmatrix}$, and $Z\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = -\begin{bmatrix} \alpha \\ \beta \end{bmatrix}$. Other useful single qubit transformations are $H = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ and $\Theta = \begin{bmatrix} e^{i\theta} & 0 \\ 0 & 1 \end{bmatrix}$. Examples of these transformations in Dirac notation are[3]

$$H\,|0\rangle = \frac{|0\rangle + |1\rangle}{\sqrt{2}}; \ H\,|1\rangle = \frac{|0\rangle - |1\rangle}{\sqrt{2}}$$

$$\Theta\,|0\rangle = e^{i\theta}\,|0\rangle\,;\ \Theta\,|1\rangle = |1\rangle$$

### I.1.3   The measurement axiom

Measurement occurs via interaction with the quantum state as if a linear projection is applied to the quantum state.[4] The set of projection matrices

---

[2] Dirac notation is a useful descriptor, as $|0\rangle = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ and $|1\rangle = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$.

[3] A summary table for common quantum operators expressed in Dirac notation is provided in section I.2.

[4] This is a compact way of describing measurement. However, conservation of information (a principle of quantum information) demands that operations be reversible. In other words, all transformations (including interactions to measure) be unitary — projections are not unitary. However, there always exist unitary operators that produce the same post-measurement state as that indicated via projection. Hence, we treat projections as an expedient for describing measurement.

are complete as they add to the identity matrix.

$$\sum_m M_m^\dagger M_m = I$$

where $M_m^\dagger$ is the adjoint (conjugate transpose) of projection matrix $M_m$ The probability of a particular measurement occurring is the squared absolute value of the projection. (An implication of the axiom not explicitly used here is that the post-measurement state is the projection appropriately normalized; this effectively rules out multiple measurement.)

For example, let the projection matrices be $M_0 = |0\rangle \langle 0| = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ and $M_1 = |1\rangle \langle 1| = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$. Note that $M_0$ projects onto the $|0\rangle$ vector and $M_1$ projects onto the $|1\rangle$ vector. Also note that $M_0^\dagger M_0 + M_1^\dagger M_1 = M_0 + M_1 = I$. For $|\psi\rangle = \alpha |0\rangle + \beta |1\rangle$, the projection of $|\psi\rangle$ onto $|0\rangle$ is $M_0 |\psi\rangle$. The probability of $|0\rangle$ being the result of the measurement is $\langle \psi| M_0 |\psi\rangle = |\alpha|^2$.

### I.1.4    The combination axiom

Qubits are combined by tensor multiplication. For example, two $|0\rangle$ qubits are combined as $|0\rangle \otimes |0\rangle = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ denoted $|00\rangle$. It is often useful to transform one qubit in a combination and leave another unchanged; this can also be accomplished by tensor multiplication. Let $H_1$ denote a Hadamard transformation on the first qubit. Then applied to a two qubit system,

$$H_1 = H \otimes I = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \text{ and } H_1 |00\rangle = \frac{|00\rangle + |10\rangle}{\sqrt{2}}.$$

Another important two qubit transformation is the controlled not operator,

$$Cnot = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

The controlled not operator flips the target, second qubit, if the control, first qubit, equals $|1\rangle$ and otherwise leaves the target unchanged: $Cnot |00\rangle = |00\rangle$, $Cnot |01\rangle = |01\rangle$, $Cnot |10\rangle = |11\rangle$, and $Cnot |11\rangle = |10\rangle$,

Entangled two qubit states or Bell states are defined as follows,

$$|\beta_{00}\rangle = Cnot \ H_1 \ |00\rangle = \frac{|00\rangle + |11\rangle}{\sqrt{2}}$$

and more generally,

$$\left|\beta_{ij}\right\rangle = Cnot \ H_1 \ |ij\rangle \ \text{ for } i,j = 0,1$$

The four two qubit Bell states form an orthonormal basis.

## I.2   Summary of quantum "rules"

Below we tabulate two qubit quantum operator rules. The column heading indicates the initial state while the row value corresponding the operator identifies the transformed state. Of course, the same rules apply to one qubit (except $Cnot$) or many qubits, we simply have to exercise care to identify which qubit is transformed by the operator (we continue to denote the target qubit via the subscript on the operator).

| operator | $\lvert 00 \rangle$ | $\lvert 01 \rangle$ | $\lvert 10 \rangle$ | $\lvert 11 \rangle$ |
|---|---|---|---|---|
| $I_1$ or $I_2$ | $\lvert 00 \rangle$ | $\lvert 01 \rangle$ | $\lvert 10 \rangle$ | $\lvert 11 \rangle$ |
| $X_1$ | $\lvert 10 \rangle$ | $\lvert 11 \rangle$ | $\lvert 00 \rangle$ | $\lvert 01 \rangle$ |
| $X_2$ | $\lvert 01 \rangle$ | $\lvert 00 \rangle$ | $\lvert 11 \rangle$ | $\lvert 10 \rangle$ |
| $Z_1$ | $\lvert 00 \rangle$ | $\lvert 01 \rangle$ | $-\lvert 10 \rangle$ | $-\lvert 11 \rangle$ |
| $Z_2$ | $\lvert 00 \rangle$ | $-\lvert 01 \rangle$ | $\lvert 10 \rangle$ | $-\lvert 11 \rangle$ |
| $Y_1$ | $i\lvert 10 \rangle$ | $i\lvert 11 \rangle$ | $-i\lvert 00 \rangle$ | $-i\lvert 01 \rangle$ |
| $Y_2$ | $i\lvert 01 \rangle$ | $-i\lvert 00 \rangle$ | $i\lvert 11 \rangle$ | $-i\lvert 10 \rangle$ |
| $H_1$ | $\frac{\lvert 00 \rangle+\lvert 10 \rangle}{\sqrt{2}}$ | $\frac{\lvert 01 \rangle+\lvert 11 \rangle}{\sqrt{2}}$ | $\frac{\lvert 00 \rangle-\lvert 10 \rangle}{\sqrt{2}}$ | $\frac{\lvert 01 \rangle-\lvert 11 \rangle}{\sqrt{2}}$ |
| $H_2$ | $\frac{\lvert 00 \rangle+\lvert 01 \rangle}{\sqrt{2}}$ | $\frac{\lvert 00 \rangle-\lvert 01 \rangle}{\sqrt{2}}$ | $\frac{\lvert 10 \rangle+\lvert 11 \rangle}{\sqrt{2}}$ | $\frac{\lvert 10 \rangle-\lvert 11 \rangle}{\sqrt{2}}$ |
| $\Theta_1$ | $e^{i\theta}\lvert 00 \rangle$ | $e^{i\theta}\lvert 01 \rangle$ | $\lvert 10 \rangle$ | $\lvert 11 \rangle$ |
| $\Theta_2$ | $e^{i\theta}\lvert 00 \rangle$ | $\lvert 01 \rangle$ | $e^{i\theta}\lvert 10 \rangle$ | $\lvert 11 \rangle$ |
| $Cnot$ | $\lvert 00 \rangle$ | $\lvert 01 \rangle$ | $\lvert 11 \rangle$ | $\lvert 10 \rangle$ |

common quantum operator rules

## I.3    Observables and expected payoffs

Measurements involve real values drawn from observables. To ensure measurements lead to real values, observables are also represented by Hermitian matrices. A Hermitian matrix is one in which the complex conjugate of the matrix equals the original matrix, $M^\dagger = M$. Hermitian matrices have real eigenvalues and eigenvalues are the values realized via measurement. Suppose we are working with observable $M$ in state $|\psi\rangle$ where

$$
\begin{aligned}
M \;&=\; \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 \\ 0 & 0 & 0 & \lambda_4 \end{bmatrix} \\
&=\; \lambda_1 \,|00\rangle\,\langle 00| + \lambda_2 \,|01\rangle\,\langle 01| + \lambda_3 \,|10\rangle\,\langle 10| + \lambda_4 \,|11\rangle\,\langle 11|
\end{aligned}
$$

The expected payoff is

$$
\begin{aligned}
\langle M \rangle \;&=\; \langle \psi | \, M \, |\psi\rangle \\
&=\; \lambda_1 \,\langle \psi| 00\rangle\,\langle 00| \psi\rangle + \lambda_2 \,\langle \psi| 01\rangle\,\langle 01| \psi\rangle \\
&\quad + \lambda_3 \,\langle \psi| 10\rangle\,\langle 10| \psi\rangle + \lambda_4 \,\langle \psi| 11\rangle\,\langle 11| \psi\rangle
\end{aligned}
$$

In other words, $\lambda_1$ is observed with probability $\langle \psi| 00\rangle\,\langle 00| \psi\rangle$, $\lambda_2$ is observed with probability $\langle \psi| 01\rangle\,\langle 01| \psi\rangle$, $\lambda_3$ is observed with probability $\langle \psi| 10\rangle\,\langle 10| \psi\rangle$, and $\lambda_4$ is observed with probability $\langle \psi| 11\rangle\,\langle 11| \psi\rangle$.

## I.4   Density operators and quantum entropy

To this point we've focused on pure states. How do we proceed if our state of knowledge indicates a mixture of states (that is, a probability weighted average of states)? Density operators supply the frame for mixed as well as pure states. A density operator is defined as

$$\rho = \sum_{i=1}^{n} p_i \left|\psi_i\right\rangle \left\langle\psi_i\right|$$

where the trace (sum of the diagonal elements) equals one, $tr(\rho) = 1$. For density operator $\rho$ the expected value associated with observable $M = \sum_{k=1}^{n} \lambda_k \left|k\right\rangle \left\langle k\right|$ (via spectral decomposition where $\left|k\right\rangle$ is an orthonormal basis) is $\langle M \rangle = tr(\rho M) = tr(M\rho)$.

This follows as the probability of observing $\lambda_k$ equals

$$\Pr(\lambda_k) = \left\langle k\right| \rho \left|k\right\rangle = \sum_{i=1}^{n} p_i \left\langle k\right| \psi_i \rangle \left\langle\psi_i\right| k \rangle$$

and useful properties of the trace.

$$tr(BA) = tr(AB)$$

To see this, recognize $(AB)_{ik} = \sum_j a_{ij} b_{jk}$ then $tr(AB) = \sum_{i,j} a_{ij} b_{ji}$ and $tr(BA) = \sum_{i,j} b_{ji} a_{ij}$ which is the same as $\sum_{i,j} a_{ij} b_{ji}$. Now, let $B \equiv \left|\psi\right\rangle \left\langle\psi\right|$

$$
\begin{aligned}
tr(AB) &= tr(A \left|\psi\right\rangle \left\langle\psi\right|) \\
&= \sum_i \left\langle i\right| A \left|\psi\right\rangle \left\langle\psi\right| i \rangle \\
&= \left\langle\psi\right| A \left|\psi\right\rangle
\end{aligned}
$$

where $\left|i\right\rangle$ is an orthonormal basis for $\left|\psi\right\rangle$ with first element $\left|\psi\right\rangle$.[5] The second line implements $tr(AB) = \sum_{i,j} a_{ij} b_{ji}$ while the third line follows from orthogonality of the basis for $\left|\psi\right\rangle$, that is, $\langle\psi| i \rangle$ is either 0 or 1.

---

[5] Notice for a pure state $\left|\psi\right\rangle$, $\left\langle\psi\right| A \left|\psi\right\rangle = \langle A \rangle_{\left|\psi\right\rangle}$, the expected value of the observable when the system is in state $\left|\psi\right\rangle$ equals $tr(A \left|\psi\right\rangle \left\langle\psi\right|)$.

Then, applying the first result followed by the third result in reverse and repeating based on the second result we have

$$
\begin{aligned}
\langle M \rangle &= \sum_{k=1}^{n} \Pr\left(\lambda_k\right) \lambda_k \\
&= \sum_{k=1}^{n}\sum_{i=1}^{n} p_i \left\langle k | \psi_i \right\rangle \left\langle \psi_i | k \right\rangle \lambda_k \\
&= \sum_{k=1}^{n} \left\langle k | \rho | k \right\rangle \lambda_k \\
&= tr\left(M\rho\right) \\
&= \sum_{k=1}^{n}\sum_{i=1}^{n} p_i \left\langle \psi_i | k \right\rangle \left\langle k | \psi_i \right\rangle \lambda_k \\
&= \sum_{i=1}^{n} p_i \left\langle \psi_i | M | \psi_i \right\rangle \\
&= tr\left(\rho M\right)
\end{aligned}
$$

Consider an example. Suppose

$$
\rho = \begin{bmatrix} 0.4 & 0 \\ 0 & 0.6 \end{bmatrix} = 0.4 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} + 0.6 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix}
$$

and

$$
M = \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} = 3 \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} - 1 \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}
$$

Then,

$$
\begin{aligned}
\Pr\left(\lambda = 3\right) &= \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0.4 & 0 \\ 0 & 0.6 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \\
&= 0.4 \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \\
&\quad + 0.6 \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \\
&= \frac{1}{2}
\end{aligned}
$$

and

$$\Pr\left(\lambda = -1\right) = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0.4 & 0 \\ 0 & 0.6 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$= 0.4 \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$+0.6 \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$= \frac{1}{2}$$

Hence, the expected value for the observable $M$ is

$$\langle M \rangle = \sum_{k=1}^{2} \Pr\left(\lambda_k\right) \lambda_k$$

$$= \frac{1}{2} 3 + \frac{1}{2}\left(-1\right) = 1$$

$$= tr\left(\rho M\right) = tr\left(M\rho\right)$$

$$= \begin{bmatrix} 0.4 & -0.8 \\ -1.2 & 0.6 \end{bmatrix} = \begin{bmatrix} 0.4 & -1.2 \\ -0.8 & 0.6 \end{bmatrix} = 1$$

or

$$\langle M \rangle = \sum_{i,j} p_j \langle i | M | \psi \rangle \langle \psi | i \rangle$$

$$= 0.4 \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$+0.4 \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$0.6 \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$+0.6 \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= \sum_{j} p_j \langle \psi | M | \psi \rangle$$

$$= 0.4 \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$+0.6 \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= 0.4 + 0.6 = 1$$

## I.4.1   Quantum entropy

von Neumann defines quantum entropy as

$$S = -tr\left(\rho \log \rho\right)$$

For pure states, $\rho$ is a Hermitian matrix whose logarithm (discussed earlier) involves its spectral decomposition $\rho = Q\Lambda Q^T$ where $\Lambda$ has one eigenvalue equal to unity and the remainder equal to zero. Hence, $\log \rho = Q\log\Lambda Q^T = 0$ (by convention, $0\log 0 = 0$) so that von Neumann entropy is zero for pure states.

On the other hand, for mixed states spectral decomposition of $\rho$ is $\sum_{j=1}^{n} \lambda_j \left|j\right\rangle \left\langle j\right|$. This involves nonzero eigenvalues so that

$$
\begin{aligned}
S &= -tr\left(\rho \log \rho\right) \\
&= -tr\left(\sum_{j=1}^{n} \lambda_j \left|j\right\rangle \left\langle j\right| \sum_{j=1}^{n} \log\left(\lambda_j\right) \left|j\right\rangle \left\langle j\right|\right) \\
&= -tr\left(\sum_{j=1}^{n} \lambda_j \log\left(\lambda_j\right) \left|j\right\rangle \left\langle j\right| \left|j\right\rangle \left\langle j\right|\right) \\
&= -tr\left(\sum_{j=1}^{n} \lambda_j \log\left(\lambda_j\right) \left|j\right\rangle \left\langle j\right|\right) \\
&= -\sum_{i=1}^{n} \lambda_i \log \lambda_i
\end{aligned}
$$

The latter result follows from equality of the sum of the eigenvalues and the trace of a matrix.

# I.5   Some trigonometric identities

Most trigonometric identities follow directly from Euler's equation:

$$e^{\pm i\theta} = \cos\theta \pm i\sin\theta$$

**Remark 1**  $\cos^2\theta + \sin^2\theta = 1$

**Proof.**

$$
\begin{aligned}
e^{i\theta}e^{-i\theta} &= e^0 = 1 \\
\left(\cos\theta + i\sin\theta\right)\left(\cos\theta - i\sin\theta\right) &= \cos^2\theta + i\cos\theta\sin\theta - i\cos\theta\sin\theta + \sin^2\theta \\
&= \cos^2\theta + \sin^2\theta
\end{aligned}
$$

∎

**Remark 2** $\frac{1-\cos\theta}{2} = \sin^2\frac{\theta}{2}$

**Proof.**

$$\frac{2-\left(e^{i\theta}+e^{-i\theta}\right)}{4} = \frac{2-2\cos\theta}{4} = \frac{1-\cos\theta}{2}$$

$$\frac{2-\left[\left(e^{\frac{i\theta}{2}}\right)^2+\left(e^{-\frac{i\theta}{2}}\right)^2\right]}{4} = \frac{2-\left[2-4\sin^2\frac{\theta}{2}\right]}{4}$$

$$= \sin^2\frac{\theta}{2}$$

Details related to the second line are below.

$$\left(e^{\frac{i\theta}{2}}\right)^2 = \left(\cos\frac{\theta}{2}+i\sin\frac{\theta}{2}\right)^2$$

$$= \cos^2\frac{\theta}{2}-\sin^2\frac{\theta}{2}+2i\cos\frac{\theta}{2}\sin\frac{\theta}{2}$$

$$\left(e^{-\frac{i\theta}{2}}\right)^2 = \left(\cos\frac{\theta}{2}-i\sin\frac{\theta}{2}\right)^2$$

$$= \cos^2\frac{\theta}{2}-\sin^2\frac{\theta}{2}-2i\cos\frac{\theta}{2}\sin\frac{\theta}{2}$$

$$\left(e^{\frac{i\theta}{2}}\right)^2+\left(e^{-\frac{i\theta}{2}}\right)^2 = 2\cos^2\frac{\theta}{2}-2\sin^2\frac{\theta}{2}$$

$$= 2-4\sin^2\frac{\theta}{2}$$

**Remark 3** $\frac{1-\cos\theta_1\cos\theta_2}{2} = \cos^2\frac{\theta_1}{2}\sin^2\frac{\theta_2}{2}+\sin^2\frac{\theta_1}{2}\cos^2\frac{\theta_2}{2}$

∎

**Proof.** From the two preceding identities

$$\frac{1-\cos\theta}{2} = \sin^2\frac{\theta}{2} = 1-\cos^2\frac{\theta}{2}$$

and

$$\cos\theta = 1-2\sin^2\frac{\theta}{2}$$

$$= 2\cos^2\frac{\theta}{2}-1$$

Therefore,

$$
\begin{aligned}
\frac{1 - \cos\theta_1 \cos\theta_2}{2} &= \frac{1 - \left(1 - 2\sin^2\frac{\theta_1}{2}\right)\left(2\cos^2\frac{\theta_2}{2} - 1\right)}{2} \\
&= \frac{1 - \left(2\cos^2\frac{\theta_2}{2} - 2\sin^2\frac{\theta_1}{2}2\cos^2\frac{\theta_2}{2} - 1 + 2\sin^2\frac{\theta_1}{2}\right)}{2} \\
&= -\cos^2\frac{\theta_2}{2} + 2\sin^2\frac{\theta_1}{2}\cos^2\frac{\theta_2}{2} + 1 - \sin^2\frac{\theta_1}{2} \\
&= \cos^2\frac{\theta_2}{2}\left(2\sin^2\frac{\theta_1}{2} - 1\right) + \cos^2\frac{\theta_1}{2} \\
&= \cos^2\frac{\theta_2}{2}\left(2\sin^2\frac{\theta_1}{2} - \cos^2\frac{\theta_1}{2} - \sin^2\frac{\theta_1}{2}\right) + \cos^2\frac{\theta_1}{2} \\
&= \cos^2\frac{\theta_2}{2}\left(\sin^2\frac{\theta_1}{2} - \cos^2\frac{\theta_1}{2}\right) + \cos^2\frac{\theta_1}{2} \\
&= \cos^2\frac{\theta_2}{2}\sin^2\frac{\theta_1}{2} - \cos^2\frac{\theta_2}{2}\cos^2\frac{\theta_1}{2} + \cos^2\frac{\theta_1}{2} \\
&= \sin^2\frac{\theta_1}{2}\cos^2\frac{\theta_2}{2} + \cos^2\frac{\theta_1}{2}\left(1 - \cos^2\frac{\theta_2}{2}\right) \\
&= \cos^2\frac{\theta_1}{2}\sin^2\frac{\theta_2}{2} + \sin^2\frac{\theta_1}{2}\cos^2\frac{\theta_2}{2}
\end{aligned}
$$

∎

**Remark 4** $\sin 2\theta = 2\cos\theta \sin\theta$

**Proof.**

$$
\begin{aligned}
\frac{e^{i2\theta} - e^{-i2\theta}}{2i} &= \frac{\cos 2\theta + i\sin 2\theta - (\cos 2\theta - i\sin 2\theta)}{2i} \\
&= \frac{2i\sin 2\theta}{2i} \\
&= \sin 2\theta \\
\frac{\left(e^{i\theta}\right)^2 - \left(e^{-i\theta}\right)^2}{2i} &= \frac{(\cos\theta + i\sin\theta)^2 - (\cos\theta - i\sin\theta)^2}{2i} \\
&= \frac{\cos^2\theta + 2i\cos\theta\sin\theta - \sin^2\theta - \left(\cos^2\theta - 2i\cos\sin\theta - \sin^2\theta\right)}{2i} \\
&= \frac{4i\cos\theta\sin\theta}{2i} \\
&= 2\cos\theta\sin\theta
\end{aligned}
$$

∎

# Appendix J
## Common distributions

To try to avoid confusion, we list our descriptions of common multivariate and univariate distributions and their kernels. Others may employ variations on these definitions.

| Multivariate distributions and their support | Density $f(\cdot)$ functions and their kernels | conjugacy |
|---|---|---|
| Gaussian (normal)<br>$x \in R^k$<br>$\mu \in R^k$<br>$\Sigma \in R^{k \times k}$<br>positive definite | $f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{k/2}\|\Sigma\|^{1/2}}$<br>$\times e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$<br>$\propto e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$ | conjugate prior for mean of multi-normal distribution |
| Student t<br>$x \in R^k$<br>$\mu \in R^k$<br>$\Sigma \in R^{k \times k}$<br>positive definite | $f(x; \nu, \mu, \Sigma)$<br>$= \frac{\Gamma[(\nu+k)/2]}{\Gamma\left(\frac{\nu}{2}\right)(\nu\pi)^{k/2}\|\Sigma\|^{1/2}}$<br>$\times \left[1 + \frac{(x-\mu)^T \Sigma^{-1}(x-\mu)}{\nu}\right]^{-\frac{\nu+k}{2}}$<br>$\propto \left[1 + \frac{(x-\mu)^T \Sigma^{-1}(x-\mu)}{\nu}\right]^{-\frac{\nu+k}{2}}$ | marginal posterior for multi-normal with unknown mean and variance |
| Wishart<br>$W \in R^{k \times k}$<br>positive definite<br>$S \in R^{k \times k}$<br>positive definite<br>$\nu > k - 1$ | $f(W; \nu, S) = \frac{1}{2^{\nu k/2}\|S\|^{\nu/2}\Gamma_k\left(\frac{\nu}{2}\right)}$<br>$\times \|W\|^{\frac{\nu-k-1}{2}} e^{-\frac{1}{2}Tr(S^{-1}W)}$<br>$\propto \|W\|^{\frac{\nu-k-1}{2}} e^{-\frac{1}{2}Tr(S^{-1}W)}$ | see inverse Wishart |
| Inverse Wishart<br>$W \in R^{k \times k}$<br>positive definite<br>$S \in R^{k \times k}$<br>positive definite<br>$\nu > k - 1$ | $f(W; \nu, S^{-1}) = \frac{\|S\|^{\nu/2}}{2^{\nu k/2}\Gamma_k\left(\frac{\nu}{2}\right)}$<br>$\times \|W\|^{-\frac{\nu+k+1}{2}} e^{-\frac{1}{2}Tr(SW^{-1})}$<br>$\propto \|W\|^{-\frac{\nu+k+1}{2}} e^{-\frac{1}{2}Tr(SW^{-1})}$ | conjugate prior for variance of multi-normal distribution |

$$\Gamma(n) = (n-1)!, \quad \text{for } n \text{ a positive integer}$$
$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$$
$$\Gamma_k\left(\tfrac{\nu}{2}\right) = \pi^{k(k-1)/4} \prod_{j=1}^{k} \Gamma\left(\tfrac{\nu+1-j}{2}\right)$$

## Multivariate distributions

| Univariate distributions and their support | Density $f(\cdot)$ functions and their kernels | conjugacy |
|---|---|---|
| Beta<br>$x \in (0,1)$<br>$\alpha, \beta > 0$ | $f(x; \alpha, \beta)$<br>$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$<br>$\propto x^{\alpha-1}(1-x)^{\beta-1}$ | beta is conjugate prior to binomial |
| Binomial<br>$s = 1.2....$<br>$\theta \in (0,1)$ | $F(s; \theta)$<br>$= \binom{n}{s} \theta^s (1-\theta)^{n-s}$<br>$\propto \theta^s (1-\theta)^{n-s}$ | beta is conjugate prior to binomial |
| Chi-square<br>$x \in [0,\infty)$<br>$\nu > 0$ | $f(x; \nu)$<br>$= \frac{1}{2^{\nu/2}\Gamma(\nu/2)} \frac{x^{\nu/2-1}}{\exp[x/2]}$<br>$\propto x^{\nu/2-1} e^{-x/2}$ | see scaled, inverse chi-square |
| Inverse chi-square<br>$x \in (0,\infty)$<br>$\nu > 0$ | $f(x; \nu)$<br>$= \frac{1}{2^{\nu/2}\Gamma(\nu/2)} \frac{\exp[-1/(2x)]}{x^{\nu/2+1}}$<br>$\propto x^{-\nu/2-1} e^{-1/(2x)}$ | see scaled, inverse chi-square |
| Scaled, inverse chi-square<br>$x \in (0,\infty)$<br>$\nu, \sigma^2 > 0$ | $f(x; \nu, \sigma^2)$<br>$= \frac{(\sigma^2\nu)^{\nu/2}}{2^{\nu/2}\Gamma(\nu/2)} \frac{\exp[-\nu\sigma^2/(2x)]}{x^{\nu/2+1}}$<br>$\propto x^{-\nu/2-1} e^{-\nu\sigma^2/(2x)}$ | conjugate prior for variance of a normal distribution |
| Exponential<br>$x \in [0,\infty)$<br>$\lambda > 0$ | $f(x; \lambda)$<br>$= \lambda \exp[-\lambda x]$<br>$\propto \exp[-\lambda x]$ | gamma is conjugate prior to exponential |
| Extreme value (logistic)<br>$x \in (-\infty, \infty)$<br>$-\infty < \mu < \infty, s > 0$ | $f(x; \mu, s)$<br>$= \frac{\exp[-(x-\mu)/s]}{s(1+\exp[-(x-\mu)/s])^2}$<br>$\propto \frac{\exp[-(x-\mu)/s]}{(1+\exp[-(x-\mu)/s])^2}$ | posterior for Bernoulli prior and normal likelihood |
| Gamma<br>$x \in [0,\infty)$<br>$\alpha, \gamma > 0$ | $f(x; \alpha, \gamma)$<br>$= x^{\alpha-1} \frac{\exp[-x/\gamma]}{\Gamma(\alpha)\gamma^\alpha}$<br>$\propto x^{\alpha-1} \exp[-x/\gamma]$ | gamma is conjugate prior to exponential and others |
| Inverse gamma<br>$x \in [0,\infty)$<br>$\alpha, \gamma > 0$ | $f(x; \alpha, \gamma)$<br>$= x^{-\alpha-1} \frac{\exp[-\gamma/x]}{\Gamma(\alpha)\gamma^{-\alpha}}$<br>$\propto x^{-\alpha-1} \exp[-\gamma/x]$ | conjugate prior for variance of a normal distribution |
| Gaussian (normal)<br>$x \in (-\infty, \infty)$<br>$-\infty < \mu < \infty, \sigma > 0$ | $f(x; \mu, \sigma)$<br>$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$<br>$\propto \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$ | conjugate prior for mean of a normal distribution |
| Student t<br>$x \in (-\infty, \infty)$<br>$\mu \in (-\infty, \infty)$<br>$\nu, \sigma > 0$ | $f(x; \nu, \mu, \sigma) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)}$<br>$\times \left(1 + \frac{1}{\nu}\frac{(x-\mu)^2}{\sigma^2}\right)^{-\left(\frac{\nu+1}{2}\right)}$<br>$\propto \left(1 + \frac{1}{\nu}\frac{(x-\mu)^2}{\sigma^2}\right)^{-\left(\frac{\nu+1}{2}\right)}$ | marginal posterior for a normal with unknown mean and variance |
| Pareto<br>$x \geq x_0$<br>$\alpha, x_0 > 0$ | $f(x; \alpha, x_0) = \frac{\alpha x_0^\alpha}{x^{\alpha+1}}$<br>$\propto \frac{1}{x^{\alpha+1}}$ | conjugate prior for unknown bound of a uniform |
| Univariate distributions | | |