

A Quick Introduction to Ox

Materials are mostly lifted from [../OxMetrics7/ox/doc/OxIntro.pdf](#)

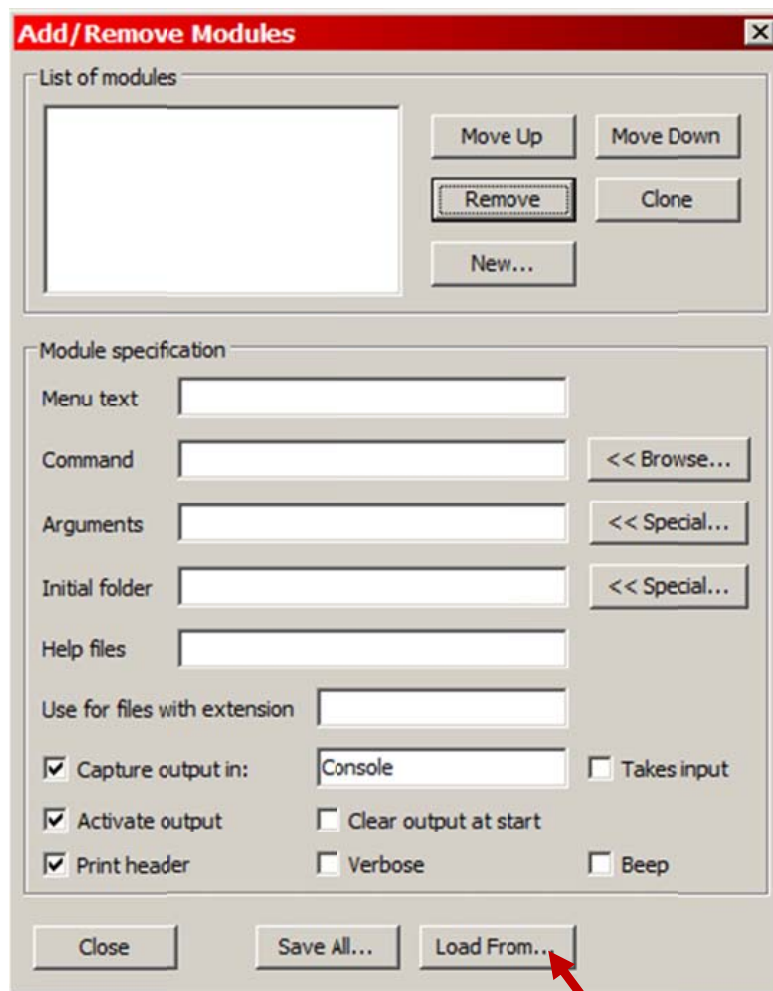
Chapter 1: Ox Environment

Installation

The latest version of Ox console and OxEdit can be downloaded at

www.doornik.com/download/oxmetrics7/Ox_Console/.

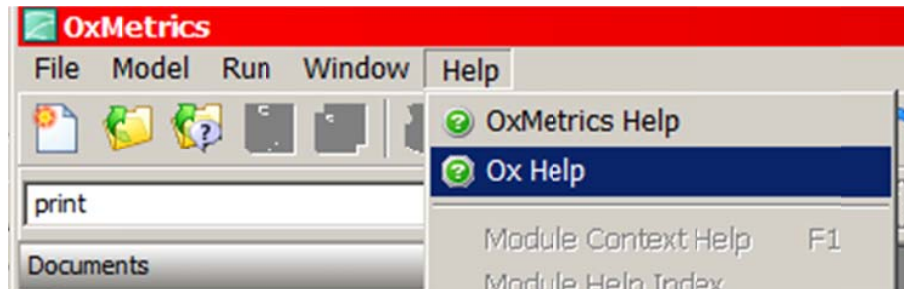
The latest version (i.e., Ox 7 Console) used to require that you specify certain modules before you can run your code. To do so, choose [Add/Remove Modules](#) in the [Tools](#) menu.



Load the file [OxCons.tool](#), which can be found in [../OxMetrics7/ox/bin/OxEdit](#).

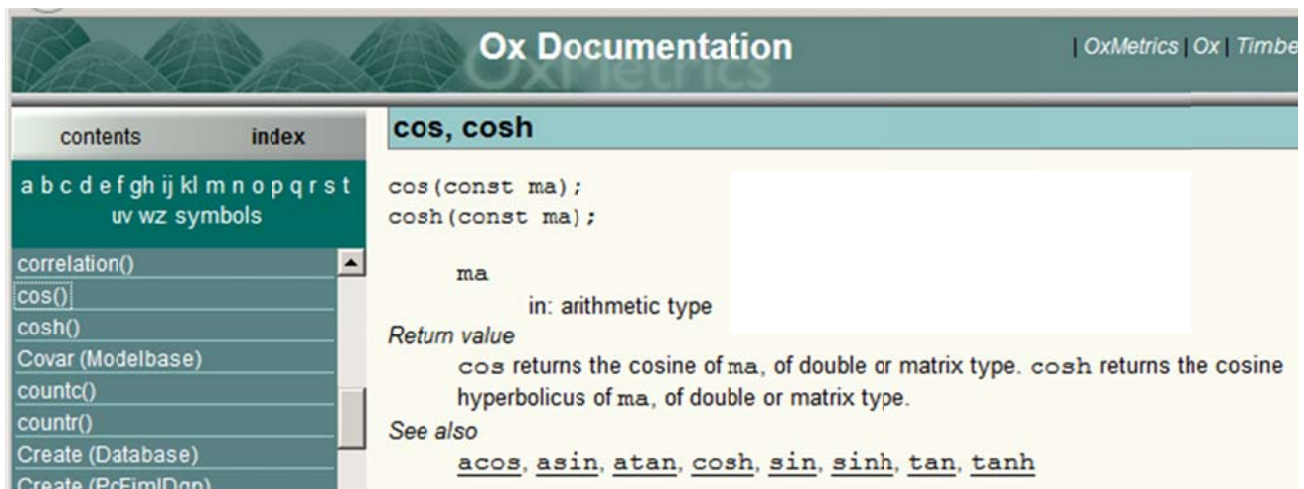
1 *Help and Documentation*

2
3 The help files can be found in the [../ox/doc/index.html](http://ox/doc/index.html). To access the master file, use [Ox Help](#).



4 The entries at the top give access to the table of contents, and to the index. The capture below shows
5 the help on `cos` - the cosine function.

6



7 *Using OxMetrics and OxRun*

8
9 Ox Console works with OxMetrics, and offers some very useful features for developing and running Ox
10 programs. The first thing to note when opening [../OxMetrics/ox/samples/myfirst.ox](http://OxMetrics/ox/samples/myfirst.ox) using Open in
11 OxMetricss [File](#) menu is the syntax highlighting:

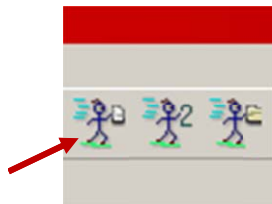
```

1 #include <oxstd.h> // include the Ox standard library header
2
3 main()              // function main is the starting point
4 {
5     decl m1, m2;     // declare two variables, m1 and m2
6
7     m1 = unit(3);    // assign to m1 a 3 x 3 identity matrix
8     m1[0][0] = 2;    // set top-left element to 2
9     m2 = <0,0,0;1,1,1>; //m2 is a 2 x 3 matrix, the first row
10                        // consists of zeros, the second of ones
11
12     print("two matrices", m1, m2);    // print the matrices
13 }

```

1 There are several features to make programming easier. Unmatched parentheses are shown in red
2 (forgetting a closing `)` or `}` is quite a common mistake). You can select of block of lines, and then use
3 **Comment In/Comment Out** to make temporary changes (there is unlimited undo/redo as well).

5 To run a file, you can click on the Run button on the top toolbar:



6 The output will appear in a new window, entitled **+Ox Output**.

```

myfirst.ox  +Ox Output
----- Ox at 10:53:17 on 06-Sep-2013 -----
Ox Professional version 6.00 (Windows/U/MT) (C) J.A. Doornik, 1994-2009
two matrices
2.0000    0.00000    0.00000
0.00000    1.0000    0.00000
0.00000    0.00000    1.0000
0.00000    0.00000    0.00000
1.0000    1.0000    1.0000

```

8 *Compilation and Run-Time Errors*

10 Program statements are processed in two steps, resulting in two potential types of errors:

12 (1) Compilation errors

1 The statements are scanned in and compiled into some kind of internal code. Errors which occur
2 at this stage are compilation errors. No statements are executed when there is a compilation error.
3 Compilation errors could be caused by undeclared variables, wrong number of function arguments,
4 forgetting a semicolon at the end of a statement (among many other reasons). For example, these two
5 messages are caused by one undeclared variable at line 10 of the program:

```
6  
7 myfirst.ox (10): y undeclared identifier  
8 myfirst.ox (10): lvalue expected  
9
```

10 Occasionally, a syntax error leads to a large list of error messages. Then, correcting the first mistake
11 could well solve most of the problem.

12 13 (2) Run-time errors

14
15 When the code which does not have syntax errors is executed, things can still go wrong, resulting in
16 a run-time error. An example is trying to multiply two matrices which have non matching dimension.
17 Here, this happened at line 10 in the main function:

```
18  
19 Runtime error: matrix[3][3] * matrix[2][3] bad operand  
20 Runtime error occurred in main(10), call trace:  
21 myfirst.ox (10): main  
22
```

23 Chapter 2: Syntax

24 25 *Introduction*

26
27 This chapter gives a brief overview of the main syntax elements of the Ox language. The most important
28 features of the Ox language are:

- 29 • The matrix is a standard type in Ox. You can work directly with matrices, for example adding
30 or multiplying two matrices together. Ox also has the standard scalar types for integers (type
31 int), and real numbers (type double). A vector is a matrix with one column or one row, whereas
32 a 1×1 matrix is often treated as a scalar. Ox will keep track of matrix dimensions for you.
- 33 • Variables are implicitly typed. So a variable can start as an integer, then become a 10×1 matrix,
34 then a 2×2 matrix. and so on. As in most languages, variables must be explicitly declared.
- 35 • Ox has strings and arrays as built-in types, to allow for higher dimensional matrices, or arrays
36 of strings.
- 37 • Ox has an extensive numerical and statistical library.
- 38 • Ox allows you to write object-oriented code (this is an optional feature).

39 The syntax of Ox is modeled on C and C++ (and also Java), so if you are familiar with these languages,
40 you will recognize the format for loops, functions, etc. However, prior knowledge of these language is
41 not assumed.

42 43 *Comment*

44

1 Ox has two types of comment: */* . . .*/* for blocks of comment, and *//* for comment up to the end
2 of a line.

3
4 When writing functions, it is useful to add comment to document the function, especially the role of
5 the arguments, and the return value.

6 7 *Program Layout*

8
9 Our first complete program is:

10

```
#include <oxstd.h>

main()
{
    print("Hello world");
}
```

11 This program does only print a line of text, but is worth discussing anyway:

- 12 • The first line includes a header file. The contents of the file `oxstd.h` are literally inserted at
13 the point of the `#include` statement. The name is between `<` and `>` to indicate that it is a
14 standard header file: it actually came with the Ox system. The purpose of that file is to declare
15 all standard library functions, so that they can be used from then onwards.
- 16 • This short program has one function, called `main`. It has no arguments, hence the empty paren-
17 theses (these are compulsory). An Ox program starts execution at the `main` function; without
18 `main`, there might be a lot of code, but nothing will happen.
- 19 • A block of code (here the *function body* of the `main` function), is enclosed in curly braces.

20 Most Ox code presented here uses syntax coloring: comment is shown as italic, and reserved words
21 (also called keywords) are bold. Both OxMetrics and OxEdit use syntax coloring. This is a purely
22 visual aide, the actual Ox code is a plain text file.

23 24 *Statements*

25
26 Statements are commands to do something, for example, some computation. Important ingredients of
27 statements are variables, to store or access a result, and operators (discussed in the next chapter) to
28 combine existing results into new ones. A statement is terminated with a semicolon (`;`). Please note
29 that, when you are copying a code from a paper, Ox makes a distinction between lower case and upper
30 case letters.

31

```

#include <oxstd.h>

main()
{
    decl n, sigma, x, beta, eps;
    n = 4; sigma = 0.25;
    x = 1 ~ ranu(n, 2);
    beta = <1; 2; 3>;

    eps = sigma * rann(n, 1);
    print("x", x, "beta", beta, "epsilon", eps);
}

```

1 Some remarks on this program:

- 2 • `decl` is used to declare the variables of this program.
- 3 • `n = 4` simply assigns the integer value 4 to the variable `n`.
- 4 • `sigma = 0.25` simply assigns a real value.
- 5 • `;` terminates each statement.
- 6 • `ranu` and `rann` are library functions for generating uniform and normal random numbers.
- 7 • `print` is a standard library function used for printing.
- 8 • `*` multiplies two variables.
- 9 • `~` concatenates two variables. Here we concatenate an integer with a 4×2 matrix.

10 Exercise

11
12 As a first exercise run the Ox program listed above. This code, `oxlut2c.ox`, can be found in the
13 `ox/tutorial` folder.

14
15 Use the help system to discover the meaning of `ranu` and the difference between `print` and `println`.

16
17 Add a line for computing $y = X\beta + \epsilon$. Also print the value of y . This requires:

- 18 (1) declaring the variable y ;
- 19 (2) inserting a statement computing $y = X\beta + \epsilon$ and storing it in y ; and
- 20 (3) adding a statement to print the y variable.

21
22 The `rows()` and `sizer()` functions returns the number of rows of a matrix, the `columns()` and `sizec()`
23 functions the number of columns. Add a print statement to report the dimensions of x in the above
24 program.

25
26 Here are some of the things which can go wrong in the previous exercises:

- 27 (1) Forget a comma. For example `decl a, b c;` needs a comma after the `b`.

1 (2) Forget a semicolon, as for example in: `n = 4 sigma = 0.25;`

2
3 (3) Adding a semicolon after the function header, as in:

4
5 `main();`
6 `{`
7 `}`

8
9 (4) Omitting the curly braces, as in:

10
11 `main()`
12 `{`
13 `print("some text");`

14
15 (5) Forget to declare a variable. The new `y` variable must be declared before it can be used.

16
17 (6) Any typo, such as writing `priny` instead of `print`.

18
19 (7) In some parts of the world a comma is used as a decimal separator. Ox uses a dot. Perhaps
20 surprisingly at this stage, this is valid code: `sigma = 0,25;`

21 *Identifiers*

22
23
24 Identifiers (names of variables and functions) may be up to 60 characters long. They may consist of
25 the characters `[A-Z]`, `[a-z]`, `[_]`, and `[0-9]`, but may not start with a digit.

26 *Style*

27
28
29 It is useful to raise the issue of programming style at this early stage. A consistent style makes a
30 program more readable, and easier to extend. Even in the very small programs of these tutorials it
31 helps to be consistent. Often, a program of a few lines grows over time as functionality is added. With
32 experience, a good balance between conciseness and clarity will be found.

33
34 Later, we shall introduce a system of name decoration, which will increase the readability of a computer
35 program. For example, we would prefix all global variables with `g_`, such as `g_dMisval` (but we shall
36 do our best to avoid global variables as much as possible).

37 *Matrix Constants*

38
39
40 The previous code used various types of constants: 4 is an integer constant, 0.25 is a double constant,
41 and "x" is a string constant. Most interesting is the value assigned to `beta`, which is a matrix constant.
42 This is a list of numbers inside `<` and `>`, where a comma separates elements within a row, and a semi-
43 colon separates rows. Remember that you can only use numbers in a matrix constant, no variables:
44 `<1,2,sigma>` is illegal. In that case use `1~2~sigma`.

45 *Creating a Matrix*

46
47
48 There are several ways to create a matrix in Ox, as the following examples show.

1 (1) Using matrix constants

2
3 (2) Reading matrices directly from disk

4
5 (3) Using library functions such as `unit`, `zeros`, `ones`:

6
7 Some of the relevant library functions are:

- 8 • `zeros(r,c)` creates an $r \times c$ matrix of zeros;
- 9 • `ones(r,c)` creates an $r \times c$ matrix of ones;
- 10 • `unit(r)` creates an $r \times r$ identity matrix;
- 11 • `constant(d,r,c)` creates an $r \times c$ matrix filled with ds ;
- 12 • `range(i,j)` creates an $1 \times j - i + 1$ matrix with $i, i + 1, \dots, j$;
- 13 • `diag(x)` creates an square matrix with with the elements of x in the diagonal.

14 (4) Using matrix concatenation

15
16 Example: `x = (0 ~ 1) | (2 ~ 3);`

17
18 Concatenation works as follows: `~` is for horizontal concatenation, and `|` is for vertical concatenation.
19 The above line will create the following matrix:

$$\begin{pmatrix} 0 & 1 \\ 2 & 3 \end{pmatrix}.$$

20 The code uses parentheses to ensure that the two horizontal concatenations are done first. Ox always
21 does horizontal before vertical concatenation, so the parentheses are redundant. When in doubt it is
22 better to write them anyway.

23 24 **Chapter 3: Operators**

25 26 *Introduction*

27
28 A language like Ox needs at least three kinds of operators to work with matrices:

- 29
30 (1) index operators, to access elements in a matrix;
- 31 (2) matrix operators, for standard matrix operations such as matrix multiplication;
- 32 (3) dot operators, for element by element operations.

33
34 There are quite a few additional operators, for example to work with logical expressions these are
35 discussed later in this chapter.

36 37 *Index Operators*

38
39 Previously, we learned various ways to create a matrix. The counterpart is the extraction of single
40 elements or specific rows or columns from the matrix. Ox has a flexible indexing syntax to achieve
41 this. But first:

Initially you might forget this and make a few mistakes, but before too long it will become second nature. Ox has adopted this convention for compatibility with most modern languages, and because it leads to faster programs. There is an option to start at index one, which is explained in the Ox manual (and not really recommended). The available indexing options are:

- Single element indexing

A matrix usually has two indices: `[i][j]` indexes element (i, j) of a matrix, where `[0][0]` is the first (top left) element.

- Range indexing

Either `i` or `j` may be replaced by a range, such as `i1:i2`. If the lower value of a range is missing, zero is assumed; if the upper value is missing, the upper bound is assumed.

- Empty index

The empty index `[]` selects all rows (when it is the first index) or all columns (when it is the second). When there is only one index `[]` the result is a column vector, filled by taking the elements from the index and row by row.

- Indexing a vector (i.e. a matrix with one row or one column)

When a matrix is used with one index, it is treated as a vector. In that case it does not matter whether the vector is a row or a column vector.

- Using a matrix as an index

In this case the matrix specifies which rows (if it is the first index) or columns (for the second index) to select. The elements in the indexing matrix must be integers (if not, they are truncated to integers by removing the fractional part).

Here are some examples:

$$x = \begin{pmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{pmatrix}, y = (0 \quad 1 \quad 2), z = \begin{pmatrix} 0 \\ 3 \\ 6 \end{pmatrix}$$

$$x[0][0] = 0, x[][1:] = \begin{pmatrix} 1 & 2 \\ 4 & 5 \end{pmatrix}, x[1][y] = (3 \quad 4 \quad 5),$$

$$y[:1] = y[0][:1] = (0 \quad 1), z[:1] = z[:1][0] = \begin{pmatrix} 0 \\ 3 \end{pmatrix}.$$

Write a program to verify these examples.

Write a program which creates a 4×4 matrix of random numbers. Then extract the 2×2 leading submatrix using (a) range indexing and (b) matrix indexing.

Matrix Operators

1 All operators + - * / work as expected when both operands are an integer or a double: when both
 2 operands are an integer, the result is an integer, otherwise it will be a double. The exception is divi-
 3 sion of two integers: this will produce a double, so 1/2 equals 0.5 (C and C++ programmers take note!).

4
 5 When matrices are involved, things get more interesting. Obviously, when two matrices have the same
 6 size we can add them element by element (+), or subtract them element by element (-). Although not
 7 a standard matrix algebraic operation, adding a column vector to a row vector works is allowed in Ox
 8 (and at times very useful). It works like a table:

$$(x_0 \quad x_1) + \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} x_0 + y_0 & x_1 + y_0 \\ x_0 + y_1 & x_1 + y_1 \\ x_0 + y_2 & x_1 + y_2 \end{pmatrix}.$$

9 For matrix multiplication use *, then element i, j of the result is the inner product of row i (left
 10 operand) and column j (right operand).

11
 12 Division (/), when the right operand is a matrix, corresponds to post multiplication with the inverse.

13
 14 In some applications, (e.g., creating the constant term for regression), when concatenating an integer
 15 (or double) and a matrix, the scalar is automatically replicated to get the same number of rows (~) or
 16 columns (|). When concatenating two non-matching matrices, the shortest one is padded with zeros
 17 at the end. (So there is a difference between $1 \sim \langle 1; 1 \rangle$ and $\langle 1 \rangle \sim \langle 1; 1 \rangle$; a warning is printed for
 18 the latter.)

19
 20 A square matrix may be raised to an integer power, using ^, for example A2 equals A*A. To summarize:

Operator	Operation
'	transpose, X'y is short for X'*y
^	(matrix) power
*	(matrix) multiplication
/	(matrix) division
+	addition
-	subtraction
~	horizontal concatenation
	vertical concatenation

21 Some operations are illegal, resulting in an error message. Here is an example:

22
 23 Runtime error: matrix[4][1] * matrix[3][4] bad operand
 24 Runtime error occurred in main(16), call trace:
 25 C:\Program Files\Ox\tutorial\oxut3a.ox (16): main

26
 27 The first says that we cannot multiply a 4×1 matrix into a 3×4 matrix. The error occurred in the
 28 main function, at line 16. In OxEdit or OxMetrics you can double click on the line with the error to
 29 jump directly to the problematic code.

30
 31 *Dot Operators*

32

1 Dot operators are element-by-element operators. For adding and subtracting matrices there is only
 2 the dot version, already used in the previous section (written as + and -). Element-by-element mul-
 3 tiplication is denoted by .* and ./ is used for element-by- element division. As with addition and
 4 subtraction, dot conformity implies that either operand may be a row (or column) vector. This is then
 5 swept through the rows (columns) of the other operand. For example:

$$(x_0 \ x_1) .* \begin{pmatrix} y_0 & y_1 \\ y_2 & y_3 \\ y_4 & y_5 \end{pmatrix} = \begin{pmatrix} x_0 y_0 & x_1 y_1 \\ x_0 y_2 & x_1 y_3 \\ x_0 y_4 & x_1 y_5 \end{pmatrix}.$$

6 To summarize:

Operator	Operation
.^	element-by-element power
.*	element-by-element multiplication
./	element-by-element division
+	addition
-	subtraction

7 Relational and Equality Operators

8
 9 Relational operators compare both operands, and exist in matrix version and in element by element
 10 (or dot) version. The first version always returns an integer, even when both arguments are matrices.
 11 The return value 0 stands for FALSE, and 1 for TRUE. When comparing a matrix to a scalar, the
 12 result is only 1 (TRUE) if it holds for each element of the matrix.

Operator	Operation
<	less than
>	greater than
<=	less than or equal to
>=	greater than or equal to
==	is equal
!=	is not equal

13 The second form of relational operator is the dotted version: this does an element by element compar-
 14 ison, and returns a *matrix* of 0s and 1s. The dotted versions are:

Operator	Operation
.<	element-by-element less than
.>	element-by-element greater than
.<=	element-by-element less than or equal to
.>=	element-by-element greater than or equal to
.==	element-by-element is equal
.!=	element-by-element is not equal

15 Often code is more readable when using the predefined constant TRUE and FALSE, instead of the
 16 numbers 1 and 0. These are defined in oxford.h. Relational operators are especially important in con-
 17 ditional expressions and loops, and these are discussed in the next chapter.

18

1 Logical Operators

2
3 These are closely related to the relational operators, and also have non-dotted and dotted versions.
4 The first evaluate to either zero or one:

Operator	Operation
<code>&&</code>	logical-and
<code> </code>	logical-or

5 If an expression involves several logical operators after each other, evaluation will stop as soon as
6 the final result is known. For example in `(1 || checkval(x))` the function `checkval` is never called,
7 because the result will be true regardless of its outcome. This is called a *boolean shortcut*.

8
9 The dotted versions perform the logical expression for each element when matrices are involved (there-
10 fore they cannot have boolean shortcuts):

Operator	Operation
<code>.&&</code>	element-by-element logical-and
<code>. </code>	element-by-element logical-or

11 Some procedures are available for selecting or dropping rows/columns based on a logical decision.
12 These are `selectifr`, `selectifc`, `deleteifr` and `deleteifc`; `vecindex` may be used to translate the
13 0-1s to indices. A very useful, but slightly more complex operator is the dot-conditional operator.

14
15 Here are some examples using these functions:

Expression	Outcome
<code>u</code>	1 0 1 0 2
<code>u .> 0</code>	1 0 1 0 1
<code>vecindex(u)</code>	0 2 4
<code>vecindex(u .> 1)</code>	4
<code>selectifc(u, u .> 0)</code>	1 1 2
<code>selectifc(u, u .> 1)</code>	2

16 Ox has more operators (e.g., conditional operators). See the Ox book for the full list.

17 Chapter 4: Input and Output

18
19 The program below shows a way of reading elements of an external file and creating an output file.

20 The output file `first.out` is a 4×1 matrix containing the squared elements of the first column of the
21 4×2 matrix `x`;

22
23 The output file is read into the vector `y`.

24
25 (To the best of my knowledge) “`%#M`” specifies the matrix to be printed row by row with the default
26 format.
27
28

- 1
- 2 Details about print formatting can be found under `print` of the Ox documentation.
- 3
- 4 Specifying paths can be avoided by keeping the program and the related files in the same folder.
- 5

```
#include <oxstd.h>

main()
{
    decl x=rann(4,2);
    print("x =", x);

    decl file, tmp;

    tmp=sprint("first.out");
    file=fopen(tmp,"w");
    fprintf(file,"%#M",x[][0].^2);
    fclose(file);

    decl y;
    file=fopen(tmp);
    fscan(file,"%#M",4,1,&y);
    fclose(file);

    print("y =", y);
}
```

```
#include <oxstd.h>
main()
{
    decl x=rann(4,2);
    print ("x =x", x);
    decl file, tmp;
    tmp=sprint("first.out");
    file=fopen(tmp,"w");
    fprintf(file,"%#M",x[][0].^2);
    fclose(file);
    decl y;
    file=fopen(tmp);
    fscan(file,"%#M",4,1,&y);
    fclose;
    print("y =",y);
}
```

1 Chapter 5: Program Flow and Program Design

3 *Introduction*

5 Ox is a complete programming language, with **if** statements and **for** loops. However, where you need
6 loops in more traditional languages, you can often use the special matrix statements available in Ox.
7 Try to avoid loops whenever you can: the vectorized version will often be very much faster than using
8 loops. On the other hand, you'll discover that loops cannot be avoided altogether: some code just
9 doesn't vectorize (or the resulting code might get too complex to maintain).

11 *for Loops*

13 The **for** loop has the following syntax:

```
14  
15     for ( initialization ; condition ; incrementation )  
16     {  
17         statements  
18     }
```

19 For example:

```
20  
21  
22     decl i;  
23  
24     for (i = 0; i < 4; i++)  
25     {  
26         print(' ', i);  
27     }
```

28 Prints: 0 1 2 3.

31 At the end of the loop, **i** will have the value 4. Since the condition is checked prior to executing the
32 loop, it is possible that the body is not executed at all (then **i** will have the initial value).

34 It is allowed to have more than one statement in the initialization or incrementation part of the **for**
35 loop. A comma is then required as a separator.

37 A loop within a loop is also allowed.

39 Write a function which multiplies two matrices, $A_{3 \times 2}$ and $B_{2 \times 4}$, using **for** loops. Implement the mul-
40 tiplication using a) a two-tiered loop and b) a three-tiered loop.

42 Compare the results with using the matrix multiplication operator.

44 *while Loops*

46 The example for the **for** loop can also be written using a **while** loop:

```
47  
48     decl i=0;
```

```

1   while (i < 4)
2   {
3       print(' ', i);
4       i++;
5   }

```

In this case, the `for` loop is more readable. But if there is not a clear initialization or incrementation part, the `while` form might be preferred.

Again, the `while` loop is not executed at all when `i` starts at 4 or above. If a loop must be executed at least once, use the `do while` loop:

break and continue

Two special commands are available inside loops:

```

17   break;

```

Terminates the loop in which the command appears, for example:

```

20   for (i = 0; i < 4; i++)
21   {
22       if (i==2) break;
23       print(' ', i);
24   }

```

Prints: 0 1

```

28   continue;

```

Starts with the next iteration of the loop, for example:

```

32   for (i = 0; i < 4; i++)
33   {
34       if (i==2) continue;
35       print(' ', i);
36   }

```

Prints: 0 1 3

Conditional Statements

In the previous section we used `if` statements to illustrate the use of `continue` and `break`. The full syntax is:

```

45   if (condition)
46   {
47       statements
48   }
49   else if (condition)

```

```

1  {
2      statements
3  }
4  else
5  {
6      statements
7  }

```

When using if-else syntax, the conditions are checked sequentially. For example, the condition for **else if** is only examined if the condition for **if** is not satisfied; in the same manner, when the condition for **else if** is not satisfied, **else** statement is automatically implemented.

```

12
13  decl x=<1,-2,-1,0,2>;
14
15  for (decl i = 0; i < 5; i++)
16  {
17      if (x[i]<0) print (" ", "1");
18      else if (x[i]==0) print (" ", "2");
19      else print (" ", "3");
20  }

```

Prints: 3 1 1 2 3

How would you implement the same program using three **if** statements?

Would this implementation be more efficient? Why or why not?

Vectorization

The program below draws n observations from a standard normal distribution, and computes the mean of the positive numbers ONLY.

The mean of positive numbers are computed two ways: using **for** loop and **if** statement, and using **selectifr**.

The time each procedure takes to compute the mean is kept tracked by **timer()**.

The **meanc(x)** function returns a $1 \times c$ matrix holding the means of the columns of x . The "**meanr**" function operates on the rows.

In the above program, which of the two procedures is more efficient? How much faster is it? You might need to increase the size of n to get more stable estimates.

This example demonstrates the advantage of using vector functions over loops.


```

main()
{
  decl n=100000;

  decl x=rann(n,1);
  decl count,total,time;
  count=total=0;

  time=timer();//save starting time

  for(decl i=0;i<n;i++){
    if(x[i]>0){
      count=count+1;
      total=total+x[i];
    }//end if
  }//end i

  decl t1=timespan(time);

  time=timer();//save starting time again

  x=selectifr(x,x.>0);
  decl t2=timespan(time);

  print("Time 1 = ",t1,"\n");
  print("Time 2 = ",t2,"\n");
  print("Mean =: ",total/count~meanc(x));
}

```

```

main()
{
  decl n=100000;
  decl x=rann(n,1);
  decl count,total,time;
  count=total=0;
  time=timer();//save starting time
  for(decl i=0;i<n;i++){
    if(x[i]>0){
      count=count+1;
      total=total+x[i];
    }//end if
  }//end i
  decl t1=timespan(time);

```

```

time=timer();//save starting time again
x=selectifr(x,x.>0);
decl t2=timespan(time);
print("Time 1 = ",t1,"\n");
print("Time 2 = ",t2,"\n");
print("Mean =: ",total/count~meanc(x));
}

```

Additional Exercises

1. Finding \sqrt{S} : The square-root of S can be solved using the so-called Babylonian method. This is an iterative solution where the n^{th} estimate is equal to

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{S}{x_n} \right).$$

a) Write a program that provides an estimate for $n = 10$, where $x_0 = S$. How close is the approximate to the exact value (i.e., \sqrt{S}).

b) Write a program that stops when $|x_n - \sqrt{S}| \leq 0.0001$. How many iterations are needed when $S = 2$? $S = 20$? $S = 200$?. (Note that Ox does not have an absolute value function. However, you can recall that $|x| = \sqrt{x^2}$.)

2. The BoxMuller transformation is a method for generating pairs of independent random numbers from $N(0, 1)$, given random numbers from $U(0, 1)$. Given U_1 and U_2 are independent random variables from $U(0, 1)$, if we let

$$Z_1 = \sqrt{-2 \log(U_1)} \cos(2\pi U_2)$$

and

$$Z_2 = \sqrt{-2 \log(U_1)} \sin(2\pi U_2),$$

then Z_1 and Z_2 are independent random variables from $N(0, 1)$.

a) Generate n random number from $U(0, 1)$, and designate it as U_1 ; do the same for U_2 . Transform U_1 and U_2 into Z_1 and Z_2 . (Note that Ox has a predefined constant for π . This constant requires the inclusion of `oxfloat.h`. Also, `log` is the natural logarithm.)

b) Find the correlation between Z_1 and Z_2 to examine if the two random variables are independent.

c) Let $Z_{2n \times 1}$ be the concatenation of Z_1 and Z_2 . Find its mean and variance. Based on these statistics, is it reasonable to assume that $Z \sim N(0, 1)$?

3. By definition, the statistic $Q = \sum_{i=1}^n Z_i^2$, where Z_i is identically and independently distributed as $N(0, 1)$ is said to be $Q \sim \chi^2(n)$. The mean and variance of Q are $E(Q) = n$ and $V(Q) = 2n$, respectively.

- 1 a) Generate n random number from $N(0, 1)$, and compute Q .
- 2
- 3 b) Repeat the process N times.
- 4
- 5 c) Find the mean and variance of the N Q statistics. Are they close to what you would expect them
- 6 to be? Again, N may need to be large for the sample statistics to be stable.
- 7
- 8 d) Can this be done without using any loop?
- 9

Functions - Part I

Using Functions

The function is a fundamental building block when writing Ox programs. Functions allow for splitting complex tasks up in manageable bits. The best ones are those which only interact with the outside via the arguments (the inputs) and the return value (the outputs, if any). Then, when there are no external variables used inside the function, the function can be treated as an isolated piece of code: the only thing which matters is the documentation of the function.

Up to this point, only one function has been used, the main function. Execution of an Ox program starts at main, from which other functions are called; there is no action outside functions. Ox comes with a vast library of functions for your convenience. These are all documented in the help and the Ox book. Whether a function is written in C, and added to the Ox system (as for the standard library), or written in Ox itself (such as the maximization functions and the Database class), does not make any difference to the user of the function.

Simple Functions

The most simple Ox function has no arguments, and returns no value. The syntax is:

```
function_name ()  
{  
    statements  
}
```

For example:

```
#include <oxstd.h>  
  
sometext()  
{  
    print("Some text\n");  
}  
  
main()  
{  
    sometext();  
}
```

```

#include <oxstd.h>
sometext()
{
    print("Some text\n");
}
main()
{
    sometext();
}

```

1 We've created the `sometext` function, and call it from the `main` function. When the program is run,
 2 it just prints `Some text`. Note that, to call the function, the empty parentheses are required.

3 *Function Argument*

4
 5
 6 A function can take arguments. In the header of the function code, the arguments are listed, separated
 7 by commas. This example takes one argument:

```

#include <oxstd.h>

dimensions(const mX)
{
    println("the argument has ", rows(mX), " rows");
}

main()
{
    dimensions( zeros(40, 5) );
}

```

```

#include <oxstd.h>
dimensions(const mX)
{
    println("the argument has ", rows(mX), " rows");
}
main()
{
    dimensions( zeros(40, 5) );
}

```

8 The `const` which precedes each argument indicates the function is only accessing the value, not chang-
 9 ing it. Although any change made to `mX` is only local to the function (once the function returns, both
 10 will have their old values back), it is very useful to use `const` wherever possible: the compiler can then

1 generate much faster code.

2
3 **Exercise:** Modify the `dimensions` function to give it two arguments, printing the number of rows in
4 both arguments.

5
6 *Returning a Value*

7
8 The **return** statement returns a value from the function, and also *exits the function*. So, when the
9 program flow reaches a **return** statement, control returns to the caller, without executing the remain-
10 der of the function. The syntax of the **return** statement is:

11
12 `return return value ;`

13
14 Or, to exit from a function which does not have a return value:

15
16 `return;`

17
18 For example:

```
MyOls1(const mY, const mX)
{
    return (mX'mX)^-1 * (mX'mY);
}
```

```
MyOls1(const mY, const mX)
{
    return (mXmX)^-1 * (mXmY);
}
```

19 Or, using the library function `olsc` as shown in the next example. This estimates and prints the
20 coefficients of the linear regression model. But first, look up in Ox Help what exactly the `olsc` is
21 supposed to do.

22
23 In the example below, the dependent variable is in the $n \times 1$ vector `mY`, and the regressors in the $n \times k$
24 matrix `mX`. The `&b` part is explained below. Any local variable (here: `b`) must be declared; `b` only exists
25 while the function is active. With **return** the result is returned to the caller, and the function exited.

```

#include <oxstd.h>

MyOls(const mY, const mX)
{
    decl b;

    olsc(mY, mX, &b);
    print("in MyOls(): b=", b);
    return b;
}

main()
{
    decl b;
    // mY argument, mX argument, both just random
    b = MyOls( rann(10, 1), ranu(10, 2) );
    // b now holds the result
    print("in main(): b=", b);
}

```

```

#include <oxstd.h>
MyOls(const mY, const mX)
{
    decl b;
    olsc(mY, mX, &b);
    print("in MyOls(): b=", b);
    return b;
}
main()
{
    decl b;
    // mY argument, mX argument, both just random
    b = MyOls( rann(10, 1), ranu(10, 2) );
    // b now holds the result
    print("in main(): b=", b);
}

```

- 1 In MyOls, move the line with the `return` statement to above the `print` statement, and compare the
- 2 output with the old version.
- 3
- 4 Test the function using the program given underneath: the task is to use `MyOls()` for the regression.
- 5 The data are observations on the weight of chickens (y) versus the amount of feed they were given (X).
- 6

```

#include <oxstd.h>

MyOls(const mY, const mX)
{
    decl b;
    olsc(mY, mX, &b);
    return b;
}

main()
{
    decl y = <0.58; 1.10; 1.20; 1.30; 1.95;
              2.55; 2.60; 2.90; 3.45; 3.50;
              3.60; 4.10; 4.35; 4.40; 4.50>;
    decl x1 = <1 : 15>'; // note transpose!
    decl mx = 1~ x1~ x1 : 2; // design matrix with quadratic effect
    print(y~ mx); // print all data
    // Finish the code: Use MyOls to regress y on mx & print results
}

```

```

#include <oxstd.h>
MyOls(const mY, const mX)
{
    decl b;
    olsc(mY, mX, &b);
    return b;
}
main()
{
    decl y = <0.58; 1.10; 1.20; 1.30; 1.95;
              2.55; 2.60; 2.90; 3.45; 3.50;
              3.60; 4.10; 4.35; 4.40; 4.50>;
    decl x1 = <1 : 15>'; // note transpose!
    decl mx = 1 ~ x1 ~ x1 .^2; // design matrix with quadratic effect
    print(y ~mx); // print all data
    // Finish the code: Use MyOls to regress y on mx & print results
}

```

1 *Function Declaration*

2

3 A function can only be called when the compiler knows about it. In the program below, the `MyOls()`

4 function can be used inside `main`, because the source code is already known at that stage. If `MyOls()`

5 were to be moved below `main` it cannot be used any more: the compiler has not yet encountered

6 `MyOls()`. However, there is a way to inform about the existence of `MyOls()`, without yet giving the

7 source code, namely by declaring the function. This amounts to copying the header only, terminated

8 with a semicolon. To illustrate the principle:

9


```

#include <oxstd.h>

MyOls(const mY, const mX); // forward declaration of MyOls,
                           // so that it can be used in main

main()
{
    // now MyOls may be used here
}

MyOls(const mY, const mX)
{
    // code of MyOls
}

```

- 1 The header files (e.g. `oxstd.h`) mainly list all the function declarations together, whereas the source
- 2 code resides elsewhere.
- 3
- 4 An option for small bits of code is to write the function to an `.ox` file, and just include the whole file
- 5 into the source file which needs it.
- 6
- 7 Save the following code as `MyOls.ox`:

```

MyOls(const mY, const mX)
{
    decl b;
    olsc(mY, mX, &b);
    return b;
}

```

- 8 Then adjust your program as follows:

```

#include <oxstd.h>
#include "myols.ox"

main()
{
    // now MyOls may be used here
}

```

1 *Returning Values in an Argument*

2
3 Often, a function needs to return more than one value. It was pointed out before that a function
4 cannot make a permanent change to an argument. However, this can be changed using the ampersand
(`&`). The following program illustrates the principle.

```
#include <oxstd.h>

test1(x) // no const, because x will be changed
{
    x = 1;
    println("in test1: x=", x);
}

test2(const ax)
{
    // Note: indexing starts at 0 in 0x
    ax[0] = 2;
    println("in test2: x=", ax[0]);
}

main()
{
    decl x = 10;
    println("x = ", x);
    test1(x); // pass x
    println("x = ", x);
    test2(&x); // pass reference to x
    println("x = ", x);
}
```

5
6 The program prints:

```
7     x = 10
8     in test1: x=1
9     x = 10
10    in test2: x=2
11    x = 2
```

12 This is happening:

- 13 • When calling `test2`, it receives in `&x` the address of the variable `x`, not its contents. In other
14 words, we are now working with a reference to `x`, rather than directly with `x`.
- 15 • Inside `test2`, the `ax` argument holds this address. To access the contents at that address, we use
16 subscript 0: `ax[0]` is the contents of the address, which we can now change.
- 17 • `ax[0] = 2` does precisely that: it changes `x` itself, because `x` resides at that address.

```

#include <oxstd.h>
test1(x) // no const, because x will be changed
{
  x = 1;
  println("in test1: x=", x);
}
test2(const ax)
{
  // Note: indexing starts at 0 in Ox
  ax[0] = 2;
  println("in test2: x=", ax[0]);
}
main()
{
  decl x = 10;
  println("x = ", x);
  test1(x); // pass x
  println("x = ", x);
  test2(&x); // pass reference to x
  println("x = ", x);
}

```

1 Consider the variable as a **mailbox**: a location at which a value can be stored. In the first case (**test1**),
 2 we just pass the content of the box: the function can read or change the content, but we (the function
 3 caller) still have the original version.

4
 5 In the second case (**test2**), we pass the key to the mailbox (the ‘address’): this allows the function
 6 to put something new in the box, which will then permanently replace the content once the function
 7 returns.

8
 9 Whether the variable is passed by *value* as in **test1**, or by *reference* as in **test2** is determined by the
 10 author of the function. The function user must follow the conventions adopted by the function author.

11
 12 Finally, Ox makes no distinction between functions which **do** return a value, and those that **do not**.
 13 If you wish, you may ignore the return value of a function altogether. (But, of course, if a function
 14 has no return value, it should not be used in an expression.)

15
 16 For example:

```

decl b0 = MyOls(my, mx)[0];
MyOls(my, mx);
print(MyOls(my, mx));

```

17 The first line directly indexes the returned vector. The second line calls the function, but does not use

1 the returned vector. The third line prints the returned vector.

2

3 Modify `MyOls` to print the following information:

4 Number of observations: `xx`

5 Coefficients:

6 `xx`

7 `xx`

8 Error variance:

9 `xx`

10 Modify `MyOls` to compute the estimated error variance. Return this through an argument.

11

12 Recall: Error variance is $\hat{\sigma}^2 = \sum e_i^2 / (n - p - 1)$, where p is the number of predictor.

13

14 Additional Exercises

15

16 a) Standardizing the Columns of a Matrix

17

18 i) Look up in Ox Help how the function `standardize` works? Use the function to create sX , the standardized version the $n \times p$ matrix X . Print the (column) means and covariance of sX .

19

20 ii) Write your own function called, `mystandardize`, that turns X into sX .

21

22 b) Square Root of a Matrix

23

24 A matrix $\mathbf{A}^{1/2}$ is said to be a square root of \mathbf{A} if $\mathbf{A}^{1/2}\mathbf{A}^{1/2} = \mathbf{A}$. (What does this imply about the dimension of $\mathbf{A}^{1/2}$, and, therefore, \mathbf{A} ?) In general, a matrix can have many square roots (i.e., the square root is not unique).

25

26 i) What is the square root of a diagonal matrix \mathbf{D} ? For real numbers, what constraints are necessary for the elements of \mathbf{D} ?

27

28 One (iterative) method for obtaining $\mathbf{A}^{1/2}$ is the matrix extension of the Babylonian method. Let \mathbf{X}_0 be the identity matrix \mathbf{I} . In the $t + 1$ iteration, the estimate of $\mathbf{A}^{1/2}$ is given by

$$\mathbf{X}_{t+1} = \frac{1}{2} (\mathbf{X}_t + \mathbf{A}\mathbf{X}_t^{-1}).$$

32 ii) Write a function that obtains $\mathbf{A}^{1/2}$ using the Babylonian method. In addition to the solution, the function should also indicate how many iterations were needed to arrive at the solution. Use different dimensions and types of matrix (i.e., symmetric, nonsymmetric) to test your function.

33

34 Note: Check how the function `invert` is used in Ox.

35

36 Another (exact) method is based the eigendecomposition or spectral decomposition. If the $n \times n$ matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^{-1}$, where \mathbf{V} is the matrix of eigenvectors and \mathbf{D} is the diagonalized eigenvalues, the the square root of \mathbf{A} is given by

$$\mathbf{A}^{1/2} = \mathbf{V}\mathbf{D}^{1/2}\mathbf{V}^{-1}.$$

1 iii) Show why this is true.

2

3 iv) Look up in Ox Help how the functions associated with eigendecomposition.

4

5 v) Write a function that obtains $\mathbf{A}^{1/2}$ using the eigendecomposition. Again, use different dimensions
6 and types of matrix (i.e., symmetric, nonsymmetric) to test your function.

7

8 vi) What happens to $\mathbf{A}^{1/2}$ when \mathbf{A} is symmetric? How does $\mathbf{A}^{1/2}\mathbf{A}^{1/2}$ compare with $\mathbf{A}^{1/2}(\mathbf{A}^{1/2})'$?

9

10 c) Decomposition of a Symmetric Matrix

11

In psychometrics, we will need to work a lot with positive-definite matrices (e.g., the covariance matrix $\mathbf{\Sigma}$). By definition, a positive-definite matrix is necessarily symmetric. One important decomposition of a positive-definite matrix \mathbf{A} is the Cholesky (also, Choleski) decomposition, which is of the form

$$\mathbf{A} = \mathbf{L}\mathbf{L}',$$

12 where \mathbf{L} is a (unique) lower triangular matrix.

13

14 i) Look up in Ox Help how the function `choleski` works.

15

16 ii) Decompose a (positive-definite) matrix using the Babylonian, eigendecomposition, and Cholesky
17 methods, and compare the results.

18

19 d) Affine Transformations

20

Suppose that \mathbf{X} is a random variable with mean $\boldsymbol{\mu}_X$ and covariance $\mathbf{\Sigma}_X$, then the affine transformation

$$\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X}$$

21 will result in the random variable \mathbf{Y} , with mean $\boldsymbol{\mu}_Y = \mathbf{a} + \mathbf{B}\boldsymbol{\mu}_X$ and covariance $\mathbf{\Sigma}_Y = \mathbf{B}\mathbf{\Sigma}_X\mathbf{B}'$.

22

23 Let the elements of $\mathbf{X}_{4 \times 1}$ be independently and identically distributed as χ^2_{10} . Also, let

24

$$\mathbf{a} = \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix} \text{ and } \mathbf{B} = \begin{pmatrix} 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}.$$

25 i) What are $\boldsymbol{\mu}_X$ and $\mathbf{\Sigma}_X$?

26

27 ii) Let $\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X}$. What is the dimension of \mathbf{Y} ? Derive $\boldsymbol{\mu}_Y$ and $\mathbf{\Sigma}_Y$. (Hint: You can use Ox to
28 carry out the matrix operations.

29

30 We can also derive the answers to i) and ii) using simulation.

31

32 iii) Generate $\mathbf{X}_{N \times 4}$, a matrix of random numbers from χ^2_{10} , where N is a very large number. Find the
33 column means and the covariance of \mathbf{X} , and compare to i).

34

35 Note that the formula for the affine transformation involves **column** vectors. However, \mathbf{X} and \mathbf{Y} are
36 made up of **row** vectors. Consequently, you also need to transpose the formula.

37

1 iv) Transform \mathbf{X} to \mathbf{Y} . What should be the dimension of \mathbf{Y} ? Find the column means and the covari-
2 ance of \mathbf{Y} , and compare to ii).

3
4 The formulas associated with the affine transformation can be applied to the sample counterparts. As
5 in, $\bar{\mathbf{Y}} = \mathbf{a} + \mathbf{B}\bar{\mathbf{X}}$ and $\mathbf{S}_Y = \mathbf{B}\mathbf{S}_X\mathbf{B}'$.

6
7 v) In the following transformation, $\mathbf{Y} = \mathbf{S}_X^{-1/2}(\mathbf{X} - \bar{\mathbf{X}}) = \mathbf{S}_X^{-1/2}\mathbf{X} - \mathbf{S}_X^{-1/2}\bar{\mathbf{X}}$, what are \mathbf{a} and \mathbf{B} ?

8
9 vi) What would $\bar{\mathbf{Y}}$ and \mathbf{S}_Y be? You can use derivation, simulation, or both to answer this question.

10
11 e) Generating Random Numbers from Multivariate Normal Distributions

12
13 It is important to know that if \mathbf{X} is a normal random variable, the random variable \mathbf{Y} resulting from
14 an affine transformation is also normal. This result will allow us to generate random numbers from
15 a multivariate normal distribution with specific $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ given random numbers from the standard
16 normal distribution.

17
18 i) If \mathbf{X} is from $MVN(\mathbf{0}, \mathbf{I})$, solve for \mathbf{a} and \mathbf{B} in the following equations:

$$\begin{aligned} E(\mathbf{Y}) = \boldsymbol{\mu} &= \mathbf{a} + \mathbf{B}\boldsymbol{\mu}_X \\ Var(\mathbf{Y}) = \boldsymbol{\Sigma} &= \mathbf{B}\boldsymbol{\Sigma}_X\mathbf{B}' \end{aligned}$$

19 ii) What are our options for \mathbf{B} ?

20
21 Let \mathbf{X} , an $N \times p$ matrix of random numbers from $N(0, 1)$, \mathbf{M} , the desired mean vector, and \mathbf{Sigma} , the
22 desired covariance matrix, be variables in `main()`.

23
24 iii) Write a function, `mvn_gen1`, that involves **passing** \mathbf{X} , plus other necessary variables, and **returns**
25 multivariate normal random numbers with the desired mean and covariance.

26
27 iv) Write a function, `mvn_gen2`, that involves **passing** reference to \mathbf{X} , plus other necessary variables,
28 and **changes** the value of \mathbf{X} to multivariate normal random numbers with the desired mean and co-
29 variance.

30

Functions - Part II

Functions as Arguments

A function may be passed as argument to another function, and then called from within that function. To pass a function as argument, just pass the name (without parentheses). The argument is then used as any other function, but there can be no argument checking at compile time, only at run time.

The examples in this section involves maximization of a function of several parameters. Fortunately, maximization code is provided with Ox, and we shall use that to illustrate using functions as arguments.

The library function `MaxBFGS` implements the BFGS ([Broyden-Fletcher-Goldfarb-Shanno](#)) method (other available unconstrained maximization methods are Newton's method and the [Nelder-Mead simplex method](#)). Further information on all these functions is in the Ox manual and online help.

Details of the procedures are beyond our current objectives, but there is a vast literature on non-linear optimization techniques to consult. Note that many texts on optimization focus on *minimization*, rather than *maximization*, but of course that is just a matter of reversing the sign.

Consider the linear regression model with two predictors:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \epsilon_i.$$

Using least-squares, estimating the regression coefficients requires minimizing the sum of squared residuals. As in, the goal is to minimize the objective function

$$Q(\beta) = \sum_{i=1}^N [Y_i - (\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i})]^2.$$

In practice, X_1 , X_2 and Y are observed, and β_0 , β_1 and β_2 are the parameters to be estimated.

For our purposes, we can simulate data that we will analyze using the following code:

```
#include <oxstd.h>
decl X1,X2,mX,Y;
main()
{
  decl beta;
  decl N=10;
  X1 = rann(N, 1);
  X2 = rann(N, 1);
  mX = 1 ~ X1 ~ X2;
  beta = rann(3, 1);
  Y = mX * beta + rann(N, 1);
}
```

```

#include <oxstd.h>

decl X1,X2,mX,Y;

main()
{
    decl beta;
    decl N=10;
    X1 = rann(N, 1);
    X2 = rann(N, 1);
    mX = 1 ~ X1 ~ X2;
    beta = rann(3, 1);
    Y = mX * beta + rann(N, 1);
}

```

1 Once we have the our data, we can obtain estimates of β_0 , β_1 and β_2 by minimizing the objective func-
 2 tion above. To achieve this goal, several functions for the nonlinear optimization *without constraints*
 3 can be considered. Examples of such functions are **MaxBFGS**, and **MaxSimplex**, which implements
 4 the Nelder-Mead simplex method.

5
 6 To use the built-in functions for maximization, the objective function should be defined according to
 7 the following format:

```

8
9 func(const vP, const adFunc, const avScore, const amHess);
10

```

11 obeying the following rules:

- 12 • **vP** is a $p \times 1$ matrix of parameter values at which the function is to be evaluated.
- 13 • **adFunc** must be the address of a variable on input. On output, the function value at the supplied
 14 parameters should be stored at the address.
- 15 • **avScore** holds either 0 on input, or the address of the score variable. If it was not 0 on input,
 16 the first derivatives of the function (the scores, a $p \times 1$ vector) should be stored at the address.
- 17 • We ignore the **amHess** argument.
- 18 • **func** should return 1 if it was successful, and 0 if it failed to evaluate the function at the supplied
 19 parameter values.

20 To minimize the objective function above, the following function can be defined:

```

funcReg(const vP, const adFunc, const avScore, const amHessian)
{
    adFunc[0] = -1 * sumc ( (Y - mX * vP) .^2 );
    return 1;
}

```


1 Note: \mathbf{Y} is an $N \times 1$ matrix, \mathbf{mX} is a $N \times 3$ matrix, and \mathbf{vP} is a 3×1 matrix, and `sumc ((Y - mX * vP) .^2)`
2 is $\sum_i^N [Y_i - (\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i})]^2$.

3
4 Next, in the `main` function, specify the initial value of the parameters and use this function name (i.e.,
5 `funcReg`) as an argument in `MaxBFGS` or `MaxSimplex`, as in the code below. To make it run, the header
6 `#import <maximize>` is required.

```
decl vP, func;  
vP = rann(3, 1); //initial values  
//If using BFGS  
MaxBFGS(funcReg, &vP, &func, 0, TRUE);  
//If using Simplex  
MaxSimplex(funcReg, &vP, &func, 0);
```

7 The function `MaxBFGS`,

```
8  
9     MaxBFGS(const func, const avP, const adFunc, const amInvHess, const fNumDer);
```

10
11 has five arguments, where

12
13 `func`

14 in: a function computing the function value, optionally with derivatives

15
16 `avP`

17 in: address of $p \times 1$ matrix with starting values

18 out: $p \times 1$ matrix with final coefficients

19
20 `adFunc`

21 in: address

22 out: double, final function value

23
24 `amInvHess`

25 in: 0 or address of $p \times p$ matrix, initial (inverse negative) quasi-Hessian \mathbf{K} ; a possible starting value
26 is the identity matrix (which is used when 0 is used as argument)

27 out: if not 0 on input: final \mathbf{K} (not reliable as estimate of actual Hessian)

28
29 `fNumDer`

30 in: 0 - `func` provides analytical first derivatives; 1 - use numerical first derivatives

31
32 Refer to Ox Help to know more about the function `MaxSimplex`.

33

```

#include <oxstd.h>
#import <maximize>

decl X1,X2,mX,Y;

//Optimization for minimizing the residuals of regression
//-----Linear Regression-----//
//y=b0+b1x1+b2x2+e

funcReg(const vP, const adFunc, const avScore, const amHessian)
{
    adFunc[0] = -1 * sumc ( (Y - mX * vP) .^ 2 );
    return 1; // function value - 1 indicates success
}

main()
{
    decl beta; //regression coefs (True values) <b0, b1, b2>
    decl N=1000; //sample size
    decl vP, func, est;

    //-----Data simulation-----//
    X1 = rann(N, 1);
    X2 = rann(N, 1);
    mX = 1 ~ X1 ~ X2; //N x 3
    beta = rann(3, 1); //3 x 1
    Y = mX * beta + rann(N, 1); //N x 1

    //True values
    print("\n", "True parameters:", beta');

    //-----Linear Regression Solution-----
    print("\n", "Estimates using OLS:");
    print((invert(mX'mX)*mX'Y)');

    //-----Optimization using BFGS-----
    vP = rann(3, 1); // starting values
    est = MaxBFGS(funcReg, &vP, &func, 0, TRUE); //MaxBFGS method
    print("\n", "Estimates using BFGS:");
    print(vP');

    //-----Optimization using Simplex-----
    vP = rann(3, 1);
    est = MaxSimplex(funcReg, &vP, &func, 0); //Simplex method
    print("\n", "Estimates using Simplex:");
    print(vP');
}

```

```

#include <oxstd.h>
#import <maximize>
decl X1,X2,mX,Y;
//Optimization for minimizing the residuals of regression
//-----Linear Regression-----//
//y=b0+b1x1+b2x2+e
funcReg(const vP, const adFunc, const avScore, const amHessian)
{
adFunc[0] = -1 * sumc ( (Y - mX * vP) .^ 2 );
return 1; // function value - 1 indicates success
}
main()
{
decl beta; //regression coefs (True values) <b0, b1, b2>'
decl N=1000; //sample size
decl vP, func, est;
//-----Data simulation-----//
X1 = rann(N, 1);
X2 = rann(N, 1);
mX = 1 ~ X1 ~ X2; //N x 3
beta = rann(3, 1); //3 x 1
Y = mX * beta + rann(N, 1); //N x 1
//True values
print("\n","True parameters:",beta');
//-----Linear Regression Solution-----
print("\n","Estimates using OLS:");
print((invert(mX'mX)*mX'Y)');
//-----Optimization using BFGS-----
vP = rann(3, 1); // starting values
est = MaxBFGS(funcReg, &vP, &func, 0, TRUE); //MaxBFGS method
print("\n","Estimates using BFGS:");
print(vP');
//-----Optimization using Simplex-----
vP = rann(3, 1);
est = MaxSimplex(funcReg, &vP, &func, 0); //Simplex method
print("\n","Estimates using Simplex:");
print(vP');
}

```

- 1 In the example above, the optimization uses numerical first derivatives of the objective function. To
- 2 accelerate the speed, analytical first derivatives can be provided with the objective function, and used
- 3 for optimization purpose.

4

The partial derivatives for the objective function

$$Q(\beta) = \sum [Y_i - (\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i})]^2$$

1 are:

$$\begin{aligned}\frac{\partial Q(\beta)}{\partial \beta_0} &= -\sum 2[Y_i - (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2})], \\ \frac{\partial Q(\beta)}{\partial \beta_1} &= -\sum 2X_{i1}[Y_i - (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2})], \\ \frac{\partial Q(\beta)}{\partial \beta_2} &= -\sum 2X_{i2}[Y_i - (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2})].\end{aligned}$$

2 We can modify our objective function to include the partial derivatives.

```
funcReg(const vP, const adFunc, const avScore, const amHessian)
{
    adFunc[0] = -1 * sumc ( (Y - mX * vP) .^ 2 );
    if (avScore) //if !0: compute score
    { //this bit is not needed for numerical derivatives
        (avScore[0])[0] = 2 * sumc (Y - mX * vP);
        (avScore[0])[1] = 2 * sumc (X1 .* (Y - mX * vP));
        (avScore[0])[2] = 2 * sumc (X2 .* (Y - mX * vP));
    }
    return 1; // function value - 1 indicates success
}
```



3 [What happened to the negative signs?](#)

4
5 In the main function, let's make the following modifications:

```
//-----Optimization using BFGS-----
vP = rann(3, 1); // starting values
est = MaxBFGS(funcReg, &vP, &func, 0, TRUE); //MaxBFGS method
print("\nMaxBFGS using numerical analysis: \nFunction value is", func);
print("Estimates are:", vP');
vP = <1;2;3>; //starting values
MaxBFGS(funcReg, &vP, &func, 0, FALSE);
print("\nMaxBFGS using analytical derivative analysis:\nFunction value is", func);
print("Estimates are:", vP');
```

6 It should be noted that the function within which we define the objective function has a return value
7 (i.e., 1 or 0). It can be used to determine whether the function evaluation succeeded (1) or failed (0).

8
9 Furthermore, the optimization functions (e.g., `MaxBFGS`, `MaxSimplex`) have their own return values as
10 well. These return values can be used to determine whether the optimization has converged or not.

11
12 **Include a code that shows the return values of the optimization functions, and explain the output.**

13

1 In some situations, the objective function needs to be maximized subject to certain constraints (e.g.,
2 $0 \leq x \leq 1$).

3
4 The following is an example of a nonlinear optimization with *constraints* using the built-in function
5 MaxSQPF. It is of the form:

```
6  
7     MaxSQPF(const func, const avP, const adFunc, const amHessian,  
8             const fNumDer, const cfunc_gt0, const cfunc_eq0, vLo, vHi);
```

9
10 where

11 **func**

12 in: a function computing the function value, optionally with derivatives

13 **avP**

14 in: address of $p \times 1$ matrix with starting values

15 out: $p \times 1$ matrix with final coefficients

16 **adFunc**

17 in: address

18 out: double, final function value

19 **amHessian**

20 in: 0 or address

21 out: if not 0 on input: final Hessian approximation matrix **B**

22 **fNumDer**

23 in: 0: func provides analytical first derivatives; 1: use numerical first derivatives

24 **vLo**

25 in: $p \times 1$ matrix with lower bounds, or $\langle \rangle$ if there are no lower bounds

26 **vHi**

27 in: $p \times 1$ matrix with upper bounds, or $\langle \rangle$ if there are no upper bounds

28
29 The **cfunc_gt0** argument can be zero, or a function evaluating the nonlinear constraints (which will
30 be constrained to be positive) with the following format:

```
31  
32     cfunc_gt0(const avF, const vP);
```

33 **avF**

34 in: address

35 out: $m \times 1$ matrix with constraints at **vP**

36 **vP**

37 in: $p \times 1$ matrix with coefficients

38 returns

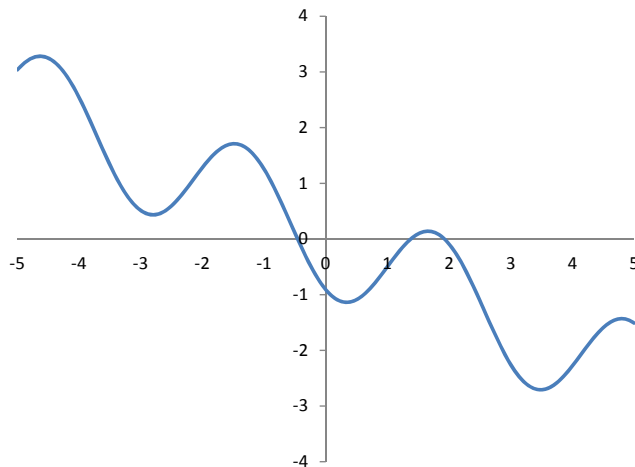
39 1: successful; 0: constraint evaluation failed

1

Assume we want to find the maximum of the objective function

$$Q(x) = \sin(2x - 2) - .5x,$$

with the constrain that $x \in [1, +\infty)$. Below is the graph of $Q(x)$.



2

3 From the graph, we can see that the global maximum of $Q(x)$ “occurs” at $-\infty$ (i.e., it does not exist).

4 However, a few local maxima exist. We want to find global maximum for $x \geq 1$. Let’s see how we can
5 solve this problem using **MaxSQPF**.

6

7 First, define the objective function using the following code:

```
func(const vP, const adFunc, const avScore, const amHessian)
{
  adFunc[0] =sin(2*vP[0]-2)-.5*vP[0];
  return 1;
}
```

8 Then set the constraints $x \geq 1$ using the following code:

```
cfunc_gt0(const avF, const vP)
{
  avF[0] = vP - 1;
  return 1;
}
```

9 Finally, the optimization can be put into action by calling it in **main**.

```

main(){
decl fvalue; //function value
decl vP=2.5; //initial value for x
decl vlow = <>; //lower bound of x, set to null
decl vhigh = <>; //upper bound of x, set to null
MaxSQPF(func, &vP, &fvalue, 0, TRUE, cfunc_gt0, 0, vlow, vhigh);
println("\nMaximum function value is ",fvalue,"\nand it can be found at x = ", vP);
}

```

- 1 To make it work, the following headers are required:

```

#include <oxstd.h>
#import <maxsqp>

```

This is how the Ox final program should look like.

```

#include <oxstd.h>
#import <maxsqp>

func(const vP, const adFunc, const avScore, const amHessian)
{
    adFunc[0] =sin(2*vP[0]-2)-.5*vP[0];
    return 1;
}

cfunc_gt0(const avF, const vP)
{
    avF[0] = vP - 1;
    return 1;
}

main(){
    decl fvalue; //function value
    decl vP=2.5; //initial value for x
    decl vlow = <>; //lower bound of x, set to null
    decl vhigh = <>; //upper bound of x, set to null

    MaxSQPF(func, &vP, &fvalue, 0, TRUE, cfunc_gt0, 0, vlow, vhigh);

    println("\nMaximum function value is ",fvalue);
    println("and it can be found at x = ", vP);
}

```

1 Consider the following questions:

2
3 (1) Can we do this optimization by setting lower bound or upper bound of x instead of the constraint
4 function?

5
6 (2) How do we set the constraint $x^2 \leq 4$?

7
8 (3) How do we use MaxSQPF for unconstrained optimization?

9
10 The following is an example to estimate the coefficients in logistic regression using optimization tech-
11 niques.

12 Let the regression model be

$$\text{logit}P(Y_i = 1|\beta) = \beta_0 + \beta_1 X_{1i},$$

or

$$P(Y_i = 1|\beta) = \frac{\exp(\beta_0 + \beta_1 X_{1i})}{1 + \exp(\beta_0 + \beta_1 X_{1i})}.$$

To simplify the notation, we can denote $P(Y_i = 1|\beta)$ as P_i . Then, the objective function (in this case, the log-likelihood) can be written as

$$Q(\beta) = \sum_{i=1}^N Y_i \log P_i + (1 - Y_i) \log(1 - P_i).$$

13 We can optimize $Q(\beta)$ using numerical first derivative or analytical first derivative. To use the latter,
14 the first derivatives should be provided with the objective function. It can be easily shown that the
15 first derivatives with respect to β_0 , β_1 and β_2 can be shown as

$$\begin{aligned} \frac{\partial Q(\beta)}{\partial \beta_0} &= \sum (Y_i - P_i), \\ \frac{\partial Q(\beta)}{\partial \beta_1} &= \sum X_{1i} \times (Y_i - P_i), \end{aligned}$$

16 Based on the objective function and the first derivatives above, we can define the objective function
as follows:

```
LogistReg(const vP, const adFunc, const avScore, const amHessian)
{
    decl ex=exp(mX * vP);
    decl P=ex./(1 + ex);
    adFunc[0] = sumc(Y.*log(P)+(1-Y).*log(1-P));
    if (avScore)
    {
        (avScore[0])[0] = sumc (Y - P);
        (avScore[0])[1] = sumc (X1 .* (Y - P));
    }
    return 1;
}
```

17


```

LogistReg(const vP, const adFunc, const avScore, const amHessian)
{
decl ex=exp(mX * vP);
decl P=ex./(1 + ex);
adFunc[0] = sumc(Y.*log(P)+(1-Y).*log(1-P));
if (avScore)
{
(avScore[0])[0] = sumc (Y - P);
(avScore[0])[1] = sumc (X1 .* (Y - P));
}
return 1;
}

```

- 1 The following code in the main can be used for data generation.

```

decl beta;
decl N=500;
decl vP, func, est;
//-----Begin Data simulation-----
X1 = rann(N, 1);
mX = 1 ~ X1;
beta = <-.5; 2>;
decl p = exp(mX * beta) ./ ( 1 + exp(mX * beta) );
Y = p .> ranu(N, 1);
//-----Data simulation End-----

```

- 2 The following code in the main can be used to obtain estimates of β with analytical and numerical
- 3 derivatives.

```

//---MaxBFGS using analytical derivatives
vP=<1, 2>';
est=MaxBFGS(LogistReg, &vP, &func, 0, FALSE);
println("\nMaxBFGS with analytical derivatives");
println(" Convergence: ",MaxConvergenceMsg(est));
println(" Function value = ", func," Parameter estimates = ", vP');
//---MaxBFGS using numerical derivatives
vP=<1, 2>';
est=MaxBFGS(LogistReg, &vP, &func, 0, TRUE);
println("\nMaxBFGS with numerical derivatives");
println(" Convergence: ",MaxConvergenceMsg(est));
println(" Function value = ", func," Parameter estimates = ", vP');

```

- 1 The following code in the `main` involves two ways of maximization with constraints - through a constraint functions or through bounds.

```
//---This is one way to maximize with a constraint
// - using constraint function
vP=<1, 2>';
decl vlow = <>;
decl vhigh = <>;
est=MaxSQPF(LogistReg, &vP, &func, 0, TRUE, cfunc_gt0, 0, vlow, vhigh);
println("\nMaximization with a constraint function");
println(" Convergence: ",MaxConvergenceMsg(est));
println(" Function value =\n      ", func,"\n Parameter estimates = ", vP');
//---This is another way to maximize with a constraint
//- using lower and upper bounds
vP=<1, 2>';
vlow = <M_INF_NEG, 0>';
vhigh = <>;
est=MaxSQPF(LogistReg, &vP, &func, 0, TRUE, 0, 0, vlow, vhigh);
println("\nMaximization with bounds");
println(" Convergence: ",MaxConvergenceMsg(est));
println(" Function value =\n      ", func,"\n Parameter estimates = ", vP');
```

- 3 Note that for the maximization with constraint function to work, we need to define a constraint function
- 4 outside `main`. The function below constrain the second parameter to be positive.

```
cfunc_gt0(const avF, const vP)
{
    avF[0] = vP[1];
    return 1;
}
```

- 5 Lastly, make sure to include and import the necessary headers, and declare the appropriate global
- 6 variables:

```
#include <oxstd.h>
#include <oxfloat.h>
#import <maximize>
#import <maxsqp>
decl X1, mX, Y;
```

- 7 1) Create a program from these codes, and run it.

8

- 1 2) Compare the results from the different maximization algorithms.
- 2
- 3 3) Compare the results from what you would get in SPSS.

Model Summary			
Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	483.857 ^a	.312	.422

a. Estimation terminated at iteration number 5 because parameter estimates changed by less than .001.

Variables in the Equation						
	B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 ^a X1	1.761	.170	107.465	1	.000	5.817
Constant	-.611	.116	27.967	1	.000	.543

a. Variable(s) entered on step 1: X1.

In addition to the estimates, we are typically also interested in the standard errors. We can calculate the standard errors using the function `Num2Derivative`. This function is used to numerically calculate the second derivatives for objective function at user-specified points. In our case, the standard errors of the coefficients using numerical second derivatives of the log-likelihood at the converged parameter values can be estimated. Because the variance-covariance matrix is minus the inverse of the second derivatives:

$$\text{Var}(\hat{\beta}) \approx -\mathbf{H}^{-1}(\hat{\beta}) \text{ where } \mathbf{H} = \left\{ \frac{\partial^2 l}{\partial \beta_i \partial \beta_j} \right\}.$$

4 The standard errors are the square root of the diagonal elements of the variance-covariance matrix.

5
6 We can now add the following code in the `main` function to calculate the standard errors of the esti-
7 mates. Note that the standard errors will only be calculated when the return value of the optimization
8 functions are `MAX_CONV` or `MAX_WEAK_CONV`, which means the estimates have converged strongly or
9 weakly.

10

```
// if estimation converges, compute the standard errors
decl mhess;
if (est == MAX_CONV || est == MAX_WEAK_CONV)
{
  if (Num2Derivative(LogistReg, vP, &mhess))
  {
    mhess = -invert(mhess);
    print("\n Standard errors:", sqrt(diagonal(mhess)'));
  }
}
```

1 Add this code to the `main` function, and compare the results to the standard errors provided by SPSS.

2

3 Additional Exercises

4

5 1. Refer to Ox Help to check how the functions `MaxSimplex` and `MaxNewton` work. Modify the pro-
6 gram to obtain the estimates and standard errors of the logistic regression coefficients using these two
7 functions.

8

9 2. The following are responses of five examinees on 10 items. These responses were generated based
10 on a 2PL model.

```
decl X=<
1,1,0,0,0,1,0,0,0,0;
1,1,0,1,0,1,0,1,0,0;
1,0,1,1,0,1,1,0,0,0;
0,0,1,0,0,1,1,1,0,0;
1,1,1,0,0,1,1,0,0,0>;
```

11 The true item parameters a and b are

```
decl b=<-2; -1; 0; 1; 2; -2; -1; 0; 1; 2>;
decl a=<0.81; 1.44; 0.70; 0.57;
0.61; 0.67; 0.81; 0.72; 0.71; 0.95>;
```

Recall that the 2PL is

$$P_{ij} = P(X = 1|\theta_i, a_j, b_j) = \frac{\exp[1.7a_j(\theta_i - b_j)]}{1 + \exp[1.7a_j(\theta_i - b_j)]},$$

and the corresponding likelihood is

$$L = \prod_{i=1}^N \prod_{j=1}^J P_{ij}^{X_{ij}} (1 - P_{ij})^{(1-X_{ij})}$$

where N is the number of examinees, and J is the number of items. Accordingly, the log-likelihood is

$$l = \log L = \sum_{i=1}^N \sum_{j=1}^J [X_{ij} \log P_{ij} + (1 - X_{ij}) \times \log(1 - P_{ij})]$$

12 a) Based on these equations, estimate the abilities using optimization methods with numerical deriva-
13 tives. You may want to constrain the person abilities to between -3 and 3.

14

b) Estimate person abilities using appropriate optimization method with the analytical derivative with respect to θ given below:

$$\frac{\partial l}{\partial \theta_i} = \sum_{j=1}^J [1.7a_j(X_{ij} - P_{ij})].$$

- 1
2 c) Estimate the standard errors of the estimates.
3
4 3. Two attributes are required for a item. Therefore, examinees can be classified into one of the four
5 groups, 00, 10, 01 and 11, depending on the attribute mastery patterns. Assume we have the following
6 information about the expected number of examinees N_g , and the expected number of examinees
7 answering correctly R_g in group g .

α_g	N_g	R_g
00	50.23	10.08
10	69.32	15.47
01	31.01	4.89
11	152.44	124.66

The likelihood function is

$$L = \prod_{g=1}^4 P(X = 1|\alpha_g)^{R_g} (1 - P(X = 1|\alpha_g))^{N_g - R_g},$$

9 where $P(X = 1|\alpha_g)$ is the probability of answering correctly the item given the g^{th} attribute pattern α_g .

10
11 If we want to fit a DINA model, $P(\alpha_1) = P(\alpha_2) = P(\alpha_3) = \varrho$ and $P(\alpha_4) = 1 - \varsigma$, where ϱ and ς are
12 the guessing and slip parameters, respectively.

- 13
a) Show that likelihood in terms of ϱ and ς is

$$L = \varrho^{\sum_{g=1}^3 R_g} \times (1 - \varrho)^{\sum_{g=1}^3 (N_g - R_g)} \times (1 - \varsigma)^{R_4} \times (\varsigma)^{N_4 - R_4}$$

- 14 b) Derive the log-likelihood function, and estimate ϱ and ς using an optimization method without a
15 constraint.

- 16
17 c) This time, include the following in your optimization function as constraints: $s, g \geq 0$ and $1 - s > g$.

- 18
19 d) The following are the first derivatives of the log-likelihood function with respect to ϱ and ς :
20

$$\begin{aligned} \frac{\partial l}{\partial \varrho} &= \frac{\sum_{g=1}^3 R_g}{\varrho} - \frac{\sum_{g=1}^3 (N_g - R_g)}{1 - \varrho} \\ \frac{\partial l}{\partial \varsigma} &= \frac{N_4 - R_4}{\varsigma} - \frac{R_4}{1 - \varsigma}. \end{aligned}$$

21 Estimate ϱ and ς with the the analytical derivatives in your optimization.

- 22
23 e) The maximum likelihood estimates of ϱ and ς can be written in closed form, as in:

$$\begin{aligned} \hat{\varrho} &= \frac{\sum_{g=1}^3 R_g}{\sum_{g=1}^3 N_g} \\ \hat{\varsigma} &= 1 - \frac{R_4}{N_4} \end{aligned}$$

- 1 Compute the MLEs using these formulas, and compare the estimates with the optimization results.
- 2
- 3 f) Find the standard errors of the estimates.
- 4
- 5 g) Lastly, find the DINO model parameter estimates and standard errors for the same data.
- 6

Data Generation and Analysis

Data Generation

In psychometrics, (Monte Carlo) simulation study is one of the primary tools, if not the primary tool, used by researchers to investigate the viability of newly developed models and procedures. Conducting simulation studies requires generating lots and lots of data. The reliability of the findings depends on how well the data have been generated from the underlying models. In addition to standard distributions, simulation studies in psychometric research also involve sampling from nonstandard distributions. We will discuss procedures for sampling from distributions commonly encountered in psychometric literature. We will also discuss general techniques that can be used for sampling from standard and nonstandard distributions.

Generating Continuous Latent Variable

Univariate

Typically, when the latent variable is univariate, it is assumed to have a standard normal distribution. Most programming languages have a function that can generate random numbers from ordinary distributions such as $N(0, 1)$. If not, techniques, such as the Box-Muller transformation, can be used. As you already know, sampling from the more general $N(\mu, \sigma^2)$ involves a simple linear transformation.

Occasionally, we may need to work with nonnormal distributions to manipulate the kurtosis or skewness of the sample. For example, we may want to work with a distribution that has heavier tails. Instead of the normal distribution, we can use the t-distribution. In Ox, this can be done using `rant(r, c, df)`. Alternatively, we can also use `rann(r, c) ./ sqrt(ranchi(r, c, df) / df)`.

For skewed distribution, the skew-normal distribution can be used. Its density function is given by

$$f(x) = \frac{2}{\omega} \phi\left(\frac{x - \xi}{\omega}\right) \Phi\left[\alpha \left(\frac{x - \xi}{\omega}\right)\right],$$

where ξ , ω , and α are the location, scale, and shape parameters, respectively,

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right),$$

and

$$\Phi(x) = \int_{-\infty}^x \phi(t) dt.$$

Its mean and variance are:

$$\mu = \xi + \omega \delta \sqrt{\frac{2}{\pi}}$$

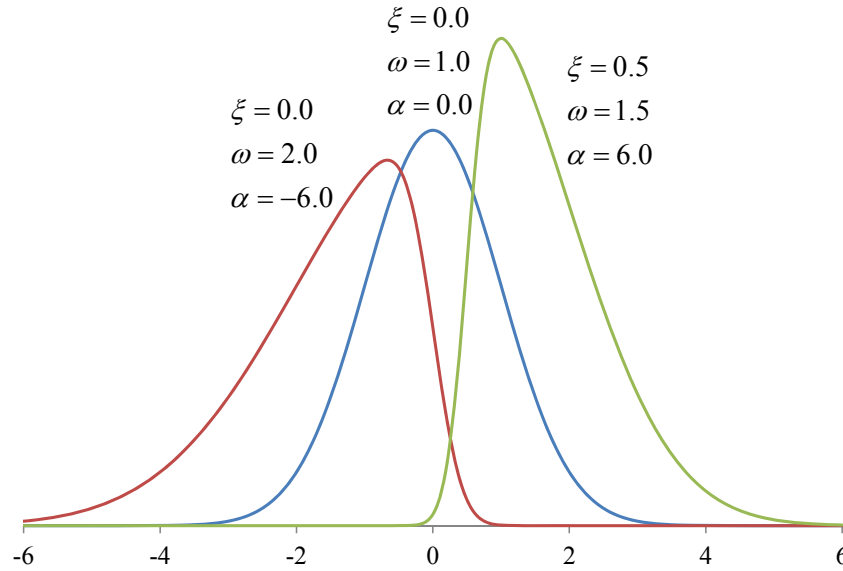
and

$$\sigma^2 = \omega^2 \left(1 - \frac{2\delta^2}{\pi}\right),$$

where

$$\delta = \frac{\alpha}{\sqrt{1 + \alpha^2}}.$$

Below is a graph showing three skew-normal distributions for different value of ξ , ω , and α . When $\xi = 0$, $\omega = 1$, and $\alpha = 0$, we get the standard normal distribution.



A Detour: Rejection Sampling (From *Bayesian Data Analysis*, pages 284-285)

Of the many procedures for sampling from a target distribution, *rejection sampling* or *acceptance-rejection method* is possibly the simplest and more general method of doing so. It can also be used in conjunction with other more complex procedures (i.e., Markov chain Monte Carlo).

Let $f(x)$ be the target density. To perform rejection sampling, we need a positive function $g(x)$ defined for all x for which $f(x) > 0$ that has the following properties:

- We need to be able to sample from the probability density $h(x)$ proportional to $g(x)$, as in, $g(x) = M \times h(x)$, for $M > 1$.
- The importance ratio $f(x)/g(x)$ must have a known bound. That is, there must be a known constant M for which $f(x)/h(x) \leq M$ for all x .

To use rejection sampling to obtain a single draw from $f(x)$, the procedure is carried out in two steps:

1. Sample X at random from $h(x)$, the probability density proportional to $g(x)$.
2. Accept the sample with probability $f(x)/g(x)$. If sample is rejected, repeat the previous step.

A good approximate density $h(x)$ for rejection sampling should be roughly proportional to $f(x)$. This can lead to most of the draws being accepted. If $h(x)$ is quite different from $f(x)$, we may need to chose a large M to ensure that $f(x) \leq g(x)$. This can lead to an inefficient sampling process (i.e., many of the draws will be rejected).

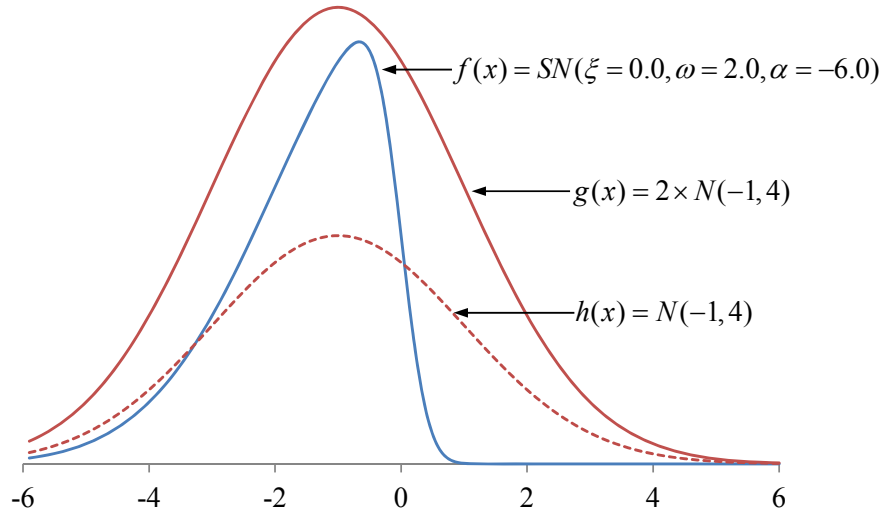
End of Detour

Suppose we are interested in sampling from a skew-normal distribution with parameters $\xi = 0$, $\omega = 2$, and $\alpha = -6$. Because we do not know a *direct* sampling method to obtain observations from this

1 distribution, we can use rejection sampling.

2

3 Below is a graph showing the density curve of $f(x) = SN(\xi = 0, \omega = 2, \alpha = -6)$. We will sample
4 from a known distribution, $h(x) = N(-1, 4)$. It can be verified that $f(x)/h(x) \leq 2$. So, we set
5 $g(x) = 2 \times h(x)$.



6 Below is a function for drawing from a skew-normal distribution. Instead of a single draw, the function
7 samples n draws at a time from $h(x)$ to be more efficient. The M for this function is appropriate for
8 the current $f(x)$ and $h(x)$ parameters. This means that a different M might be needed if $f(x)$ and
9 $h(x)$ are altered.

```
decl N=10000,n=100;
decl mu=-1,sigma=2;
decl xi=0,omega=2,alpha=-6;
decl M=2;
drawSN(decl n)
{
decl h=rann(n,1)*sigma+mu;
decl g=densn((h-mu)/sigma);
decl f=2/omega*densn((h-xi)/omega).*probn(alpha*(h-xi)/omega);
decl prob=f./(M*g);
return selectifr(h,prob.>ranu(n,1));
}
```

10 We will need the code below in the *main*. The **while** loop keeps track of the sampled draws, number
11 of attempts, and number of accepted draws. We also compare μ and σ^2 with \bar{x} and s^2 .

```
decl X=<>;
```

```

decl count0=0,count1=0; //count0->total attempts;count1->number of accepts
decl time=timer();
while(count1<N)
{
decl tmp=drawSN(n);
decl acc=rows(tmp);
count0=count0+n;
count1=count1+acc;
X=X|tmp;
} //end while
print("\n Elapsed Time:",timespan(time),"\n");
print("\n # of Attempts vs. # of Accepts:", count0~count1);
X=X[0:N-1];
decl delta=alpha/sqrt(1+alpha^2);
print("\n Theoretical Mean and Variance:");
print(xi+omega*delta*sqrt(2/M_PI )~omega^2*(1-2*delta^2/M_PI ));
print("\n Sample Mean and Variance:", meanc(X)~varc(X));

```

1. Use different values for n , and examine how it affects the a) time and b) number of attempts needed to reach the sample size N .

3

2. Examine the impact of changing $M = 2$ to $M = 10$ on the number of attempts.

5

Another way of generating skewed (or more generally, nonnormal) distribution is to use the Fleishman's cubic transformation method. The random variable Y can be derived from $X \sim N(0, 1)$ using the following transformation:

$$Y = a + bX + cX^2 + dX^3.$$

The coefficients a , b , c , and d can be chosen so that the distribution has specific mean (μ_Y), variance (σ_Y^2), skewness (γ_{1Y}), and kurtosis (γ_{2Y}). For example, for $\mu_Y = 0$ and $\sigma_Y^2 = 1$, $\gamma_{1Y} = 0$, and $\gamma_{2Y} = 0$ (i.e., standard normal distribution), we set $a = c = d = 0$ and $b = 1$; for $\mu_Y = 0$ and $\sigma_Y^2 = 1$, $\gamma_{1Y} = 0.5$, and $\gamma_{2Y} = 0$ (i.e., a skewed distribution), we set $a = -c \approx -0.0926$, $b \approx 1.0399$, and $d \approx -0.0166$.

10

Although we know the moments of Y , we do not know its exact distribution.

12

It should be noted that for latent variable models to be identified, we typically constrain the latent variable to have a mean of $\mu = 0$ and variance $\sigma^2 = 1$.

15

Multivariate

17

Many of the models in the psychometric literature call for multivariate latent variables. More often than not, they are also multivariate normal.

20

As discussed in Lecture 2, we can generate $Y_{p \times 1} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ by sampling $X_{p \times 1}$ from $MVN(\mathbf{0}, \mathbf{I})$, and performing the affine transformation

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{X}.$$

1 Again, for identifiability purposes, we typically set $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma} = \mathbf{R}$.

3 Incorporating Covariates

5 Increasingly, covariates are being included as part of latent variable models. In addition to the response
6 data, covariates can be a source of ancillary information to improve estimation.

8 Let $\boldsymbol{\theta}_{D \times 1}$ be the latent variable, and $\mathbf{Y}_{p \times 1}$ the vector of covariates. The objective is to generate $\boldsymbol{\theta}$ that
9 has the following structures:

- 10 • (1) the marginal covariance of $\boldsymbol{\theta}$ is a correlation matrix, and
- 11 • (2) the proportion of variance accounted for by the covariates in regressing $\boldsymbol{\theta}$ on \mathbf{Y} is ψ .

Let

$$\begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{Y} \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\theta\theta} & \boldsymbol{\Sigma}_{\theta Y} \\ \boldsymbol{\Sigma}_{Y\theta} & \boldsymbol{\Sigma}_{YY} \end{bmatrix} \right).$$

By design, we can let $\boldsymbol{\Sigma}_{YY}$ be equal to \mathbf{I} . Using the properties of conditional distributions of a MVN distribution (Result 4.6 of Johnson, and Wichern, 2007; Theorem 3.2.4 of Mardia, Kent, & Bibby, 1979), the conditional distribution of $\boldsymbol{\theta}$ given $\mathbf{Y} = \mathbf{y}$ is

$$\boldsymbol{\theta} | \mathbf{Y} = \mathbf{y} \sim \text{MVN}(\boldsymbol{\Sigma}_{\theta Y} \mathbf{y}, \boldsymbol{\Sigma}_{\theta\theta} - \boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}).$$

12 The two structures given above imply the following.

$$\begin{aligned} \{\boldsymbol{\Sigma}_{\theta\theta}\}_{dd} &= 1; \\ \{\boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd} / \{\boldsymbol{\Sigma}_{\theta\theta}\}_{dd} &= \{\boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd} = \psi; \end{aligned}$$

13 where $\{\boldsymbol{\Sigma}\}_{rc}$ is the element in row r and column c , of the matrix $\boldsymbol{\Sigma}$.

15 It can be assumed without loss of generality that the proportions of variance accounted for by the
16 covariates are the same for all dimensions, and that each covariate has the same predictive power. It
17 follows that

$$\begin{aligned} \{\boldsymbol{\Sigma}_{\theta\theta} - \boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd} &= 1 - \psi \\ \{\boldsymbol{\Sigma}_{\theta\theta} - \boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd'} &= \rho(1 - \psi) \\ \{\boldsymbol{\Sigma}_{\theta\theta}\}_{dd'} &= \rho + \psi - \rho\psi \\ \{\boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd} = \{\boldsymbol{\Sigma}_{\theta Y} \boldsymbol{\Sigma}_{Y\theta}\}_{dd'} &= \psi \\ \{\boldsymbol{\Sigma}_{\theta Y}\} &= \sqrt{\psi/p} \end{aligned}$$

18 Following are the steps in generating N $\boldsymbol{\theta}$ observations with a correlation of ρ between the dimensions,
19 and where the proportion variance accounted for by the covariates is equal to ψ .

- 20 1. Generate a $\mathbf{Y}_{N \times p}$ matrix of covariates, as in, N observations from $\text{MVN}_p(\mathbf{0}, \mathbf{I})$.
- 21 2. Compute \mathbf{YB} , the matrix *conditional* mean vector for the N examinees, where \mathbf{B} is $p \times D$ and
22 $\{\mathbf{B}\} = \sqrt{\psi/p}$.
- 23 3. Generate observations from $\text{MVN}_D(\{\mathbf{YB}\}'_i, \boldsymbol{\Sigma}_{\theta|Y})$, $i = 1, 2, \dots, N$, where $\{\mathbf{YB}\}_i$ is the i^{th} row of
24 \mathbf{YB} , and $\boldsymbol{\Sigma}_{\theta|Y} = (1 - \psi)[(1 - \rho)\mathbf{I}_{D \times D} + \rho]$.

Using this algorithm, the *marginal* distribution of $\boldsymbol{\theta}$ is a multivariate normal distribution with the parameters

$$\begin{aligned}\boldsymbol{\mu}_{\boldsymbol{\theta}} &= \mathbf{0} \\ \boldsymbol{\Sigma}_{\boldsymbol{\theta}\boldsymbol{\theta}} &= (1 - \rho - \psi + \rho\psi)\mathbf{I}_{D \times D} + (\rho + \psi - \rho\psi)\end{aligned}$$

Approximate normal distributions can be obtained when using discrete covariates by adding an additional step to the algorithm. For example for a dichotomous variable, the additional step, Step 1.5, is to code the covariate as -1 if the observation is below the median, or 1 otherwise. This modification preserves the mean and covariance structure of the covariates. The rest of the steps still apply.

Higher-Order Models

In some applications, the structure of the covariance $\boldsymbol{\Sigma}$, typically, expressed as \mathbf{R} , of the latent variable $\boldsymbol{\theta}_{D \times 1}$ can be accounted for by positing a higher-order latent variable $\boldsymbol{\omega}_{p \times 1}$. For our purposes, it would suffice to treat $\boldsymbol{\omega}$ as univariate (i.e., $\boldsymbol{\omega} = \omega$).

Specifically, we can express the d^{th} element of $\boldsymbol{\theta}$ as a linear function of ω , as in

$$\theta_{di} = \lambda_d \times \omega_i + \epsilon_{di},$$

where λ_d is the latent coefficient in regressing the θ_d on ω , ϵ_{di} is the error term that is independently and identically distributed as $N(0, 1 - \lambda_d^2)$, and $|\lambda_d| \leq 1$.

This formulation implies the following properties:

$$\begin{aligned}\rho(\theta_d, \theta_{d'} | \omega) &= 0 \\ \rho(\theta_d, \omega) &= \lambda_d \\ \rho(\theta_d, \theta_{d'}) &= \lambda_d \times \lambda_{d'}\end{aligned}$$

Show that

(i) $\theta_{di} | \omega \sim N(\lambda_d \times \omega_i, 1 - \lambda_d^2)$, and

(ii) if $\omega \sim N(0, 1)$, the marginal distribution of θ_d is also $N(0, 1)$.

To generate $\boldsymbol{\theta} \sim \text{MVN}(0, \mathbf{R})$ using the higher-order model, where \mathbf{R} is compound symmetric (i.e., $\{\mathbf{R}\}_{dd'} = \rho$):

1. Generate ω_i from $N(0, 1)$.

2. Then, given ω_i , generate θ_{di} from $N(\lambda_d \times \omega_i, 1 - \lambda_d^2)$, where $\lambda_d = \sqrt{\rho}$.

It should be noted that when $D = 2$, the estimates of the regression parameters will not be unique because different sets of λ_1 and λ_2 resulting in the same ρ can be found. When $D = 3$, the higher-order model perfectly fits the structure of $\boldsymbol{\Sigma}$. **Why?** When $D \geq 4$, the true correlational structure of $\boldsymbol{\theta}$ may be more complex than what a linear model can fit. **Again, why?**

Generating Discrete Latent Variable

1 In this section, we will focus on a specific type of discrete latent variable - (multivariate) binary la-
2 tent variables $\alpha_{K \times 1}$ (i.e., attribute vectors). Of course, α need not be binary/dichotomous, as in,
3 α_k need not assume the value 0 or 1 only. However, only a few works involving polytomous at-
4 tributes are currently available so it may not be worthwhile to delve into this type of attribute in
5 the time being. In addition, extending procedures for dichotomous to polytomous attributes should be
6 relatively straightforward. Below is a discussion of how to generate attributes with different structures.

7 Saturated Structure

10 A saturated attribute structure is the most general formulation of α because no *prior* constraints
11 are imposed on the attributes. The prevalence of an attribute pattern, $P(\alpha_l)$, for $l = 1, \dots, 2^K$, is
12 unrestricted, and is only subject to the constraint that $\sum_{\forall l} P(\alpha_l) = 1$. It should be noted, however,
13 that attributes based on a saturated formulation are not necessarily unstructured - they may still
14 follow a particular structure, albeit implicitly. For example, although $P(\alpha_l)$ may be assigned without
15 specific restrictions, the resulting attributes may still end up following a particular structure, say, a
16 higher-order structure.

18 To sample α from a saturated model or structure, we use the $P(\alpha_l)$. Given the probabilities, we can sample
19 α using a multinomial distribution. In Ox, we can use `ranmultinomial(const n, const vp)`, where
20 `n` is the number of trials (i.e., N), and `vp` is a vector of probabilities of length L , with probabilities
21 summing to one. This function will return $\mathbf{n}_{1 \times L} = (n_1, \dots, n_L)$, a vector of nonnegative integers,
22 where $\sum_{\forall l} n_l = N$.

24 Given \mathbf{n} , we can create $\mathbf{A}_{N \times K}$, which contains n_1 rows of α_1 , n_2 rows of α_2 , ..., and n_L rows of α_L .

26 Below are codes we can use to sample α from a saturated model.

28 The first code starts with the declaration of a few global variables. Note the difference between `alpha`
29 and `alpha0`.

31 The function `pattern` will create `alpha`, a $2^K \times K$ matrix containing all the possible attribute patterns.
32 The attributes patterns are ordered as follows:

- 33 1. First, increasing with respect to $\sum_{\forall k} \alpha_{lk}$, the number of mastered attributes.
- 34 2. Then, for a fixed number of mastered attributes, decreasing with respect to $\sum_{\forall k} \alpha_{lk} 10^{K-k+1}$, the
35 “numerical value” of the attribute pattern, which can be interpreted as “attributes with smaller
36 k are mastered earlier.”

37 In this particular ordering, it will be easy to recreate `alpha` manually. More importantly, this particular
38 ordering makes it easier to locate specific attribute patterns.

```
decl N=20, K=3; //N->sample size; K->number of attributes
decl alpha; //alpha-> 2^KxK possible attribute patterns
decl alpha0; //alpha0-> NxK sampled attribute patterns
```

```

pattern(){ //create a matrix of all the possible attribute patterns
alpha=<>;
for(decl k=1;k<=K;k++){
decl tmp=<>;
decl base=zeros(2^(K-k),1)|ones(2^(K-k),1);
for(decl l=0;l<2^(k-1);l++){
tmp=tmp|base;}
alpha=alpha~tmp;}
alpha=alpha~sumr(alpha);
decl tmp=zeros(2^K,1);
for(decl k=0;k<K;k++){
tmp=tmp-10^(K-k-1)*alpha[] [k];
}
alpha=alpha~tmp;
alpha=sortbyc(alpha,K^(K+1));
alpha=alpha[] [0:K-1];
}

```

- 1 The code below should appear in `main`. It starts by calling the function `pattern`. In lieu of specifying
- 2 the 2^K probabilities associated with the 2^K attribute patterns, we have a general code for making all
- 3 the attribute patterns equally likely (i.e., the attribute patterns are uniformly distributed).
- 4
- 5 Note that, originally, `ex` is used to hold $(P(\alpha_1), \dots, P(\alpha_{2^K}))'$; eventually, the same variable is used to
- 6 hold (n_1, \dots, n_{2^K}) . In doing so, we have one less variable to declare and keep track of.
- 7
- 8 After declaring `alpha0` as a null matrix, we use a loop to convert \mathbf{n} to \mathbf{A} .

```

ranseed(130930);
pattern();
decl ex=ones(2^K,1)./(2^K); //attribute patterns are equiprobable
ex=ranmultinomial(N,ex);
alpha0=<>;
for(decl l=0;l<2^K;l++){
alpha0=alpha0|(ones(ex[l],1)*alpha[l] []);
}
print(ex'~alpha);
print(alpha0);

```

- 9 Assemble the codes, and examine the output. Which attribute pattern has the highest (lowest) fre-
- 10 quency?

11 Higher-Order Structure

13

1 The attribute patterns are said to have a higher-order structure if they follow the model presented by
2 de la Torre and Douglas (2004) where an examinees overall (or higher-order) ability θ determines the
3 probability that the examinee masters each of the attributes. (In general, the higher-order ability can
4 be multivariate. However, estimating θ accurately may be challenging.)

5 Specifically, θ_i and α_{ik} are related as follows:

$$P(\alpha_{ik} = 1|\theta_i) = \frac{\exp[\lambda_{1k}(\theta_i - \lambda_{0k})]}{1 + \exp[\lambda_{1k}(\theta_i - \lambda_{0k})]}.$$

6 where λ_{0k} and λ_{1k} are the difficulty and discrimination parameters of attribute k . It is not unusual to
7 use a common discrimination parameter across the attributes (i.e., $\lambda_{1k} = \lambda_1$).

8 In this particular model, it is assumed that $\theta_i \sim N(0, 1)$. In addition, the attributes are also assumed
to be conditionally independent given θ , as in,

$$P(\boldsymbol{\alpha}_i|\theta_i) = \prod_{\forall k} P(\alpha_{ik}|\theta_i).$$

9 The latter assumption allow us to generate $\boldsymbol{\alpha}_i$ by generating θ_i from $N(0, 1)$, and then α_{ik} given θ_i ,
10 $k = 1, \dots, K$, individually. Compared to sampling from a multinomial distribution, the higher-order
11 model allows us to sample $\boldsymbol{\alpha}$ without needing to know the corresponding $P(\boldsymbol{\alpha})$.

12
13 Sampling $P(\alpha_{ik} = 1|\theta_i)$ is identical to sampling $P(X_{ij} = 1|\theta_i)$, which will be discussed below.

15 Hierarchical Structure

16
17 If some the attributes are related such that the mastery of one attribute is a prerequisite to the mastery
18 of another (i.e., attribute mastery is sequential), the structure of the attribute vectors is said to be
19 hierarchical. An important implication of the model is that some attribute patterns are precluded by
20 the structure. For example, if $K = 3$, and $\alpha_1 \rightarrow \alpha_2, \alpha_3$ (read, mastery of α_1 is a prerequisite to the
21 mastery of either α_2 or α_3 , then the attribute patterns 010, 001, and 011 are deemed impossible. Thus,
22 instead of eight, we only need to deal with five attribute patterns.

23
24 To generate $\boldsymbol{\alpha}$ from a hierarchical structure, we can employ the same procedure used in the saturated
25 attribute structure. However, we need to make sure that impossible attribute patterns are assigned a
26 probability of 0.

27
28 The code below can be used to replace the original definition of `ex`.

```
decl ex=<1,1,0,0,1,1,0,1>; //three attribute patterns are deemed impossible
ex= ex/sumc(ex); // normalization to ensure the probabilities add to 1
```

29 Examine the output for the code, and determine which of the possible attribute patterns has the
30 highest (lowest) frequency.

31
32 Note that, in this code, the possible attribute patterns are equiprobable. However, this need not always
33 be the case. For example, if we have reasons to believe that 000 is twice as likely as the remaining

1 attribute patterns, we can declare `ex=<2,1,0,0,1,1,0,1>'` instead.

2 *Generating Binary Item Response Data with Underlying Continuous Latent Variable*

3
4
5 The most typical models for binary item response data involve a unidimensional continuous latent variable. These models are the 1PL, 2PL and 3PL, which are under the logistic item response model family.

6
7 Recall that the 3PL is of the form:

$$P(X_{ij} = 1|\theta_i) = \gamma_j + (1 - \gamma_j) \frac{\exp[1.7\alpha_j(\theta_i - \beta_j)]}{1 + \exp[1.7\alpha_j(\theta_i - \beta_j)]},$$

8 and constraining the appropriate item parameters will yield the 1PL and the 2PL.

9
10 The Ox function below can be used to generate items from the 3PL. Note that including several
11 variables as arguments allows us to minimize our number of global variables to be declared.

```
decl N=10, J=5;
gen3PL(decl t, a, b, c)
{
  decl tmp=exp(1.7*(t*a'-ones(N,1)*(a.*b)'));
  tmp=c'+ones(N,1)*(1-c)'*.tmp./(1+tmp);
  return tmp.>ranu(N,J);
}
```

12 The function needs to be called from the `main`.

```
ranseed(130930);
decl theta, alpha, beta, gamma;
decl X;
theta=rann(N,1);
theta=sortc(theta); //sorting the abilities from lowest to highest
alpha=.75+ranu(J,1);
beta=.90*rann(J,1);
beta=sortc(beta);
gamma=.25*ranu(J,1);
print("\nItem Parameters:",alpha~beta~gamma);
X=gen3PL(theta, alpha, beta, gamma);
print("\nTheta & Data:",theta~X);
```

13 **Modify the item parameters to generate 1PL and 2PL data.**

14
15 For multidimensional IRT models, we only need to modify how the original `tmp` is computed. Of
16 course, MIRT models also involve α_j and θ_i , and may require using δ_j instead of β_j .

17

Generating Polytomous Item Response Data with Underlying Continuous Latent Variable

In addition to dichotomous response, we also have to deal with polytomous response (i.e., $X = 0, 1, \dots, m - 1$). Polytomous response can be classified as either ordered or unordered. *Generalized partial credit* and *graded response* models are example of IRT models for ordered polytomous response, whereas the *nominal* and *multiple-choice* models are examples of IRT models for unordered polytomous response.

Although the different polytomous models calculate $P(X_{ij} = x|\theta)$ differently, the same process can be used to generate polytomous data regardless of how $P(X_{ij} = x|\theta)$ has been calculated. For this reason, we will focus only on one model, the GPCM.

Recall that, in the GPCM,

$$P(X_{ij} = x|\theta_i) = \frac{\exp[\sum_{k=0}^x 1.7\alpha_j(\theta_i - \beta_{jk})]}{\sum_{l=0}^{m-1} \exp[\sum_{k=0}^l 1.7\alpha_j(\theta_i - \beta_{jk})]}.$$

For identifiability, we can constraint $1.7\alpha_j(\theta_i - \beta_{j0}) = 0$.

Another Detour: The Inverse CDF Method

Inverse CDF (also, *inverse transform*, among many other names) sampling is a basic method for sampling numbers at random from any probability distribution given its cumulative distribution function (CDF). Although this procedure is applicable to both continuous and discrete distributions, we will employ the procedure only for the latter.

The basic idea is to sample u from $U(0, 1)$, and then return the largest number x from the domain of the distribution f_X such that $P(X \leq x) \leq u$.

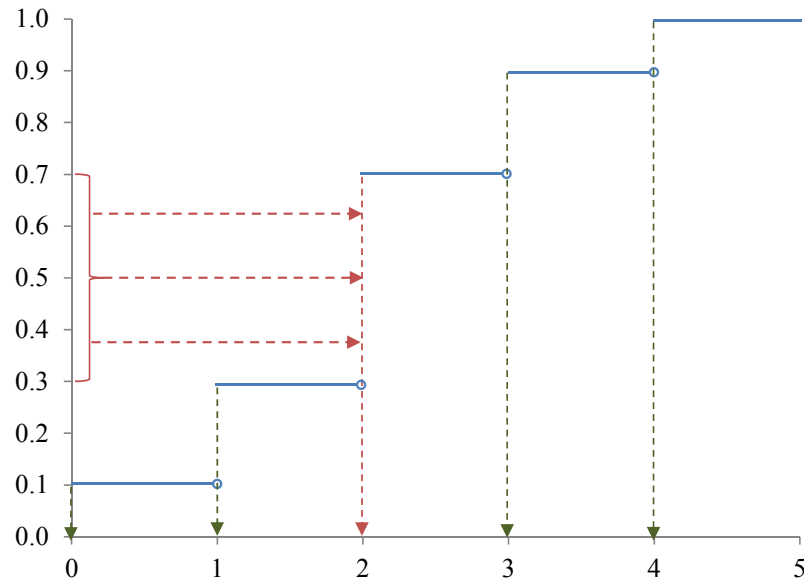
We illustrate this method using the random variable X , with the following probability and cumulative probability distribution functions.

X	$P(X = x)$	$P(X \leq x)$
0	0.1	0.1
1	0.2	0.3
2	0.4	0.7
3	0.2	0.9
4	0.1	1.0

If our sampling procedure is correct, 10% of the time, we should get a “0” or a “4”, 20% of the time, a “1” or a “3”, and 40% of the time, a “2”.

If we sample u from $U(0, 1)$. 40% of the time, $0.3 \leq u \leq 0.7$. The graph below shows that, for $X = 0, 1, 2$, $P(X \leq x) \leq u$. Accordingly, we choose $X = 2$ whenever $0.3 \leq u \leq 0.7$.

A roundabout way to implement this procedure is to use several `if-else` statements. The following code, where $Px = P(X \leq x)$, can sample correctly from this distribution.



```

if(u<P0) X=0;
else if(u<P1) X=1;
else if(u<P2) X=2;
else if(u<P3) X=3;
else X=4;

```

- 1 However, this code is inefficient, particularly if we need to generate an $N \times J$ matrix of responses as it
 2 will require several `for` loops. Below is an equivalent, but more efficient code that can be used with
 3 matrices.

```

X=(u>P0)+(u>P1)+(u>P2)+(u>P3);

```

4 [End of Detour](#)

- 5 We can use the Inverse CDF procedure to efficiently sample polytomous responses. The following code
 6 generates the response category probabilities using the GPCM for $m = 5$.

```

decl N=10, J=5, m=5;
decl p0, p1, p2, p3, p4;

```

```

probGPC(decl t, a, b) //specifically for m=5
{
  decl pp, pp0, pp1, pp2, pp3, pp4;
  pp0=ones(N,J);
  pp1=exp(1.7*(t*a'-ones(N,1)*(a.*b[] [0]))');
  pp2=pp1.*exp(1.7*(t*a'-ones(N,1)*(a.*b[] [1]))');
  pp3=pp2.*exp(1.7*(t*a'-ones(N,1)*(a.*b[] [2]))');
  pp4=pp3.*exp(1.7*(t*a'-ones(N,1)*(a.*b[] [3]))');
  pp=pp0+pp1+pp2+pp3+pp4;
  p0=pp0./pp;
  p1=pp1./pp;
  p2=pp2./pp;
  p3=pp3./pp;
  p4=pp4./pp;
}

```

- 1 Given the computed probabilities, the code below converts the probabilities into item responses.

```

genGPC(decl aX) //specifically for m=5
{
  decl P0, P1, P2, P3;
  P0=p0;
  P1=P0+p1;
  P2=P1+p2;
  P3=P2+p3;
  decl u=ranu(N,J);
  aX[0]=(u.>P0)+(u.>P1)+(u.>P2)+(u.>P3);
}

```

- 2 Finally, to use these functions, we need to call them in **main**. Before we can do that, we also need to
- 3 declare and initialize a few variables.

```

ranseed(130930);
decl theta, alpha, beta;
decl X;
theta=rann(N,1);
theta=sortc(theta); //sorting the abilities from lowest to highest
alpha=.75+ranu(J,1);
beta=.90*rann(J,m-1);

```

```

beta=sortr(beta); // sorting the category difficulties from easiest to hardest
beta=sortbyc(beta,0); // sorting the item from easiest to hardest
print("\nItem Parameters:",alpha~beta);
probGPC(theta, alpha, beta);
genGPC(&X);
print("\nTheta & Data:",theta~X);

```

1 *Generating Binary Item Response Data with Underlying Discrete Latent Variable*

2
3 In this section, we will consider data generation procedures for two types of CDMs: reduced and saturated. For the former, we will use the DINA model, and for the latter, the G-DINA model. We assume
4 that the matrix $\mathbf{A}_{N \times K}$ of attribute vectors has already been generated.
5

6 DINA Model

7
8
9 This function `alphaeta` below will require a few of the global variables. This function simply requires
10 the $N \times K$ matrix `alpha0` as input and converts it to the $N \times J$ matrix `eta` based on the attribute
11 specification in the $J \times K$ matrix `Q`.

```

decl N=10, J=6, K=3;
decl alpha; //all possible attribute patterns
decl Q;
decl alpha0, eta; //alpha0->generated attribute vectors; eta->converted from alpha0
alphaeta(decl a)
{
eta=a*Q';
eta=eta.==(ones(N,1)*sumc(Q'));
}

```

12 In addition, we need the following in `main`: initializing the `Q`-matrix, and the guessing and slip param-
13 eters. Once `alphaeta` has been called, we can create the matrix of probabilities, and then convert it
14 to the matrix of item responses.

```

Q=<1,0,0; 0,1,0; 0,0,1; 1,1,0; 1,0,1; 0,1,1>;
decl g=.05+.25*ranu(J,1);
decl s=.05+.25*ranu(J,1);
alphaeta(alpha0);
decl prob=(1-eta).*(ones(N,1)*g') + eta.*(ones(N,1)*(1-s'));
decl X=prob.>ranu(N,J);
print(alpha0, eta);
print(prob,X);

```

1 **Modify alphaeta to generate responses for the DINO model.**

3 G-DINA Model

5 Because the G-DINA model is more general, we can expect the code to generate the item responses to
6 also be more complex, particularly if the code has to have sufficient generality.

8 In the G-DINA model, Item j can have as many as 2^{K_j} unique parameters, where $K_j = \sum_k q_{jk}$. In
9 this discussion, we will use $P(X_j = 1|\alpha_{lj})$ as our parameters, where α_{lj} denotes the reduced attribute
10 vector for Item j .

12 Let $K_{j(max)} = \max(K_1, \dots, K_J)$. We can collect all the item parameters in the $J \times 2^{K_{j(max)}}$ matrix \mathbf{B} .
13 For $K_j < K_{j(max)}$, we only need entries for the first 2^{K_j} columns.

15 The importance of clarifying the meaning of the columns cannot be overstated. For $K_j = 1$, the
16 relevant columns, from left to right, represent the reduced attribute vectors 0 and 1; for $K_j = 2$, 00,
17 10, 01 and 11; for $K_j = 3$, 000, 100, 010, 001, 110, 101, 011 and 111; and so forth.

19 Determining $P(X_{ij} = 1|\alpha_i)$ is a matter of determining the location of α_{ij} , and assigning the corre-
20 sponding item parameter (i.e., probability). To implement this in Ox, we need we additional functions.

22 First is the function **patternj**. This function is identical to the earlier function **pattern**, except that
23 instead of fixed K , this function uses K_j to create all the possible *reduced* attribute patterns.

```
patternj(decl Kj)
{
  decl eta=<>;
  for(decl k=1;k<=Kj;k++){
    decl tmp=<>;
    decl base=zeros(2^(Kj-k),1)|ones(2^(Kj-k),1);
    for(decl l=0;l<2^(k-1);l++){
      tmp=tmp|base;}
    eta=eta~tmp;}
  eta=eta~sumr(eta);
  decl tmp=zeros(2^Kj,1);
  for(decl k=0;k<Kj;k++){
    tmp=tmp-10^(Kj-k-1)*eta[][k];
  }
  eta=eta~tmp;
  eta=sortbyc(eta,Kj~(Kj+1));
  return(eta[][0:Kj-1]);
}
```

24 The other function is **probj**, which determines the probability of correct response to Item j for all the
25 N examinees.

```

probj(decl jj)
{
  decl alpha_j=selectifc(alpha0,Q[jj][].==1);
  decl Kj=sumr(Q[jj][]);
  decl eta_j=patternj(Kj);
  decl prob_j=zeros(N,1);
  for(decl i=0;i<N;i++){
    decl tmp=vecindex(prodr((ones(2^Kj,1)*alpha_j[i][]).==eta_j).==1);
    prob_j[i]=beta[jj][tmp];
  }//end i
  return(prob_j);
}

```

- 1 In **main**, we need to define the item parameter matrix **beta**, which is declared as a global variable.
- 2 Currently, **beta** contains the parameters for the DINA model. However, by changing the entries of
- 3 **beta**, we can generate item responses for other CDMs.
- 4
- 5 Note that the probabilities are determined one item at a item, so this program is not as efficient as
- 6 the program for a specific reduced model, such as for the DINA model. As alluded to earlier, there is
- 7 always a tradeoff between efficiency and generality.

```

beta=<
.2,.8, 0, 0;
.2,.8, 0, 0;
.2,.8, 0, 0;
.2,.2,.2,.8;
.2,.2,.2,.8;
.2,.2,.2,.8;
.2,.2,.2,.8>;
decl X=<>;
for(decl j=0;j<J;j++){
  X=X~probj(j);
}
print(X);
X=X.>ranu(N,J);

```

- 8 Modify **beta** to generate responses for the DINO model.
- 9

10 Reading and Writing Files

- 11
- 12 Although simulation studies can be done one fell swoop, it is always advisable to break them down
- 13 into smaller steps. For example, for a simulation study that has the following three components, data
- 14 generation, analysis of data, and summarizing the results, it would be a good idea to save the data after
- 15 they have been generated, perform the analysis as a separate step, and carry out the summarization

1 as another step. By doing it step-by-step, it would be easier to redo or extend the simulation study
2 when the situation calls for it.

3
4 In Ox, reading and writing files can be done in many ways. However, knowing a few basic templates
5 would suffice. We can illustrate by solving the following problem.

6
7 We are interested in examining the behavior of the sample mean as a function of sample size ($N =$
8 $25, 100$), and the underlying distribution ($N(0, 1)$ and $U(-3, 3)$). Specifically, we want to examine the
9 mean and variance of the sample mean across $R = 10$ replications.

10
11 Even though the problem is simple enough, we will carry out the data generation as one step, and the
12 remaining procedure as another step.

13
14 The program is short enough so we can declare a few global variables.

```
decl N=<25,100>;
decl dist={"Normal","Uniform"};
decl R=10;
```

15 The following code in `main` will generate $2 \times 2 \times 10 = 40$ data sets that will be written to files. The files
16 can be uniquely identified by the name of the distribution, sample size, and replications number.

17

```
decl X, file, tmp;
for(decl i=0; i<2; i++){
  for(decl j=0; j<2; j++){
    for(decl r=0; r<R; r++){
      if (j==0) X=rann(N[i],1);
      else X=-3+6*ranu(N[i],1);
      tmp=sprintf(dist[j],"_",N[i],"_",r,".dat");
      file=fopen(tmp,"w");
      fprintf(file,"%#M",X);
      fclose(file);
    }}//end r, j, i
```

18 Most of the code for generating the data can be recycled to read the files. As we read each file, we are
19 also computing the corresponding mean and saving it in the right cell of the $R \times 4$ matrix `mean`.

20
21 The variables `r` and `count` allow us to identify the correct rows and columns, respectively.

22
23 Finally, we print the summary as a screen output.

24
25 Incidentally, the default number of columns of a matrix that the screen output can accommodate is
26 six. To increase the number of characters, hence, columns, use `format(int)`, where `int` is the number

1 of characters. For example, `format(1000)`.

2

```
decl mean=zeros(R,4);
decl X, file, tmp;
decl count=0;
for(decl i=0; i<2; i++){
for(decl j=0; j<2; j++){
for(decl r=0; r<R; r++){
tmp=sprint(dist[j],"_",N[i],"_",r,".dat");
file=fopen(tmp);
fscan(file,"%#M",N[i],1,&X);
fclose(file);
mean[r][count]=meanc(X);
} //end r
count=count+1;
} //end j, i
print("\n    Normal25    Uniform25    Normal100    Uniform100");
print("\nMn: ", meanc(mean),"Var",varc(mean));
```

3 Based on results in sampling theory, we know that $E(\bar{x}) = \mu_x$, and $\sigma_{\bar{x}}^2 = \sigma_x^2/N$, where N is the sample
4 size. By the way, in case you do not remember, for the uniform distribution, $\sigma_x^2 = (b-a)^2/12$, where
5 a and b are the lower and upper limits, respectively.

6

7 How do the sample statistics compare with their expected values? Try increasing R to see if it has any
8 impact on the results.

9

10 Additional Exercises

11

12 1. Generate $R = 100$ data sets for $N = 500$ examinees and $J = 25$ items that conform to the 3PL. Fix
13 $\alpha_j = 1.0$ and $\gamma_j = 0.2$, but, for each replication, sample θ_j and β_j from $N(0, 1)$.

14

15 2. Find the MLE ($\hat{\theta}_i$) and EAP ($\tilde{\theta}_i$) of θ_i .

16

You may use last week's optimization function to find the MLE. For the EAP, recall that the expected
a posteriori is

$$E(\theta|\mathbf{X}_i) = \int \theta \cdot p(\theta|\mathbf{X}_i) d\theta,$$

where

$$p(\theta|\mathbf{X}_i) = \frac{L(\mathbf{X}_i|\theta) \cdot p(\theta)}{\int L(\mathbf{X}_i|\theta) \cdot p(\theta) d\theta}.$$

17 More often than not, we assume $p(\theta) \sim N(0, 1)$.

18

But as we have learned in the past, we can approximate any continuous distribution with a discrete
distribution, and replace integration with summation. The posterior distribution can be approximated

as

$$p(A_q|\mathbf{X}_i) = c \cdot L(\mathbf{X}_i|A_q) \cdot W(A_q),$$

where A_q is a quadrature node, $W(A_q)$ is its corresponding weight, and

$$c^{-1} = \sum_{q=1}^Q L(\mathbf{X}_i|A_q) \cdot W(A_q).$$

Once we have the discrete version of the posterior distribution, we can easily approximate $E(\theta|\mathbf{X}_i)$ as

$$\sum_{q=1}^Q A_q \cdot p(A_q|\mathbf{X}_i).$$

3. For each replication, find the correlation $Cor(\check{\theta}_i, \theta_i)$, bias, $\sum_{\forall i}(\check{\theta}_i - \theta_i)/N$, and MSE, $\sum_{\forall i}(\check{\theta}_i - \theta_i)^2/N$, where $\check{\theta} = \hat{\theta}, \tilde{\theta}$.

4. Average the correlation, bias, and MSE across the replications. Determine which is the better estimate, MLE or EAP? Is the result consistent across the different criteria?

Optional Exercises

1. In the higher-order attribute structure model, we can obtain the *marginal* probability of α_l from the *joint* probability $P(\alpha_l, \theta)$ by integrating out θ , as in,

$$P(\alpha_l) = \int P(\alpha_l, \theta) d\theta = \int P(\alpha_l|\theta) p(\theta) d\theta \approx \sum_{\forall q} P(\alpha_l|A_q) W(A_q).$$

Recall that $P(\alpha_l|\theta) = \prod_{\forall k} P(\alpha_{lk}|\theta)$. Note that $P(\alpha_{lk}|\theta)$ is a shorthand of writing $P(\alpha_{lk} = a|\theta)$, $a = 0, 1$. For $a = 1$, $P(\alpha_{lk}|\theta) = P(\alpha_k = 1|\theta)$; for $a = 0$, $P(\alpha_{lk}|\theta) = 1 - P(\alpha_k = 1|\theta)$.

For example, $P(101|\theta)$ can be written as $P(\alpha_1 = 1|\theta) \times [1 - P(\alpha_2 = 1|\theta)] \times P(\alpha_k = 1|\theta)$.

More generally, we can write the $P(\alpha_l|\theta)$ as

$$P(\alpha_1 = 1|\theta)^{\alpha_{l1}} \times [1 - P(\alpha_1 = 1|\theta)]^{1-\alpha_{l1}} \times \dots \times P(\alpha_K = 1|\theta)^{\alpha_{lK}} \times [1 - P(\alpha_K = 1|\theta)]^{1-\alpha_{lK}}.$$

For $K = 5$, $\lambda_{1k} = 2$, $\lambda_{0k} = (-1.0, -0.5, 0.0, 0.5, 1.0)'$, and $\theta \sim N(0, 1)$, find $P(\alpha_l)$, for $l = 1, \dots, 32$.

As an alternative to using the integration formula above, we can find the marginal probabilities computationally by generating an extremely large number of attribute vectors, and recording their relative frequencies.

It can be noted that, by using these probabilities to generate the attribute vectors using the multinomial distribution, the attribute vectors will have the same (more or less, subject to Monte Carlo errors) relative frequencies. Hence, given the relative frequencies alone, we cannot be certain that the true generating model is the saturated or the higher-order model.

For $\lambda_{1k} = 0$, show that $P(\alpha_l) = 1/2^K$. Note that $P(\alpha_k = 1|\theta) = P(\alpha_k = 1) = 0.5$.

1 This mean that, when the attributes are unrelated to the higher-order trait, the higher-order structure
 2 is equivalent to the *uniform attribute structure*. Note that the exact values of the difficulty parameters
 3 are immaterial.

4
 5 2. For greater flexibility, the higher-order model can be written in the slope-intercept form, as in,

$$P(\alpha_{ik} = 1|\theta_i) = \frac{\exp(\gamma_{0k} + \gamma_{1k}\theta_i)}{1 + \exp(\gamma_{0k} + \gamma_{1k}\theta_i)}.$$

5 This allows us to set the value of γ_{1k} independent of γ_{0k} . For example, even if $\gamma_{1k} = 0$, the value of
 6 γ_{0k} still matters.

7
 8 When $\gamma_{1k} = 0$, the higher-order model in the slope-intercept form reduces to

$$P(\alpha_{ik} = 1|\theta_i) = P(\alpha_{ik} = 1) = p_k = \frac{\exp(\gamma_{0k})}{1 + \exp(\gamma_{0k})},$$

8 where p_k is the proportion of individuals in the population who mastered attribute k .

9
 10 Thus, $P(\boldsymbol{\alpha}_l)$ simplifies to

$$p_1^{\alpha_{l1}} \times q_1^{1-\alpha_{l1}} \times \dots \times p_K^{\alpha_{lK}} \times q_K^{1-\alpha_{lK}},$$

10 for $q_k = 1 - p_k$.

11
 12 In addition to being able to compute the marginally probability of an attribute more easily (e.g.,
 13 $P(101) = p_1q_2p_3$, without any integration), we can also generate the attribute vectors more easily. As
 14 in, α_k is sampled independently from $Ber(p_k)$.

15
 16 This is the *independence model* (i.e., the attributes are independent of each other).

17
 18 When both $\gamma_{0k} = \gamma_{1k} = 0$, the independence model produces what attribute structure?

19
 20 For $K = 5$, $\gamma_{1k} = 0$, and $\gamma_{0k} = (-2.0, -1.0, 0.0, 1.0, 2.0)'$, find p_k , for $k = 1, \dots, 5$, and $P(\boldsymbol{\alpha}_l)$, for
 21 $l = 1, \dots, 32$. Solve this problem in two ways: analytically and computationally.

22

MMLE & EM Algorithms - Part I

Some Background: Structural vs. Incidental Parameters

In latent variable modeling, we need to distinguish between structural and incidental parameters. For example, in item response theory and cognitive diagnosis models, the structural parameters are the item (e.g., α_j , β_j , g_j , s_j) parameters; the incidental parameters are the person-specific parameters (e.g., θ_i , α_i).

Another way of viewing the distinction between structural and incidental parameters is to take different samples, and note which parameters should remain the same (up to some random errors), and which should not. The former are the structural parameters, the latter the incidental parameters.

With small enough data and powerful enough optimization tool, we might be tempted to maximize, say, the likelihood function to estimate the structural and incidental parameters simultaneously via joint maximum likelihood estimation. Theoretically, this is not advisable. Neyman & Scott (1948) have shown that in the presence of incidental parameters, the maximum likelihood estimates of structural parameters may not be consistent.

Marginalized Maximum Likelihood Estimation

Recall that $P(A) = \int P(A, B)dB = \int P(A|B) \times P(B)dB$. In the same way, instead of maximizing the likelihood function, which contains both structural and incidental parameters, we can work with the marginalized maximum likelihood function, where the incidental parameters are integrated out. That is, we will maximize the marginalized likelihood

$$L(\mathbf{X}_i) = \int L(\mathbf{X}_i|\theta) \times g(\theta)d\theta.$$

The integral can be approximated by using summation and quadrature nodes, as in,

$$L(\mathbf{X}_i) \approx \sum_{q=1}^Q L(\mathbf{X}_i|A_q) \times w(A_q).$$

For N individuals, the (marginalized) likelihood is

$$L(\mathbf{X}) = \prod_{i=1}^N L(\mathbf{X}_i).$$

Instead of the likelihood, it would be more convenient to work with the loglikelihood

$$l(\mathbf{X}) = \log L(\mathbf{X}) = \sum_{i=1}^N \log L(\mathbf{X}_i).$$

Let ξ_j be the parameter of interest. To obtain an estimate of ξ_j , we need to find

$$\frac{\partial l(\mathbf{X})}{\partial \xi_j} = \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial \xi_j}.$$

1 The derivative on the right-hand side is

$$\frac{\partial L(\mathbf{X}_i)}{\partial \xi_j} = \sum_{q=1}^Q w(A_q) \times \frac{\partial L(\mathbf{X}_i|A_q)}{\partial \xi_j}.$$

2 For dichotomous response (i.e., $X_{ij} = 0, 1$),

$$L(\mathbf{X}_i|A_q) = \prod_{j=1}^J P_{jq}^{X_{ij}} \times (1 - P_{jq})^{1-X_{ij}},$$

3 where $P_{jq} = P(X_j = 1|A_q)$. So

$$\begin{aligned} \frac{\partial L(\mathbf{X}_i|A_q)}{\partial \xi_j} &= \prod_{j' \neq j} \left[P_{j'q}^{X_{ij'}} \times (1 - P_{j'q})^{1-X_{ij'}} \right] \times \frac{\partial \left[P_{jq}^{X_{ij}} \times (1 - P_{jq})^{1-X_{ij}} \right]}{\partial \xi_j} \\ &= \prod_{j' \neq j} \left[P_{j'q}^{X_{ij'}} \times (1 - P_{j'q})^{1-X_{ij'}} \right] \times P_{jq}^{X_{ij}} \times (1 - P_{jq})^{1-X_{ij}} \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \\ &= \prod_{j'=1}^J \left[P_{j'q}^{X_{ij'}} \times (1 - P_{j'q})^{1-X_{ij'}} \right] \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \\ &= L(\mathbf{X}_i|A_q) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \end{aligned}$$

4 Therefore,

$$\frac{\partial L(\mathbf{X}_i)}{\partial \xi_j} = \sum_{q=1}^Q w(A_q) \times L(\mathbf{X}_i|A_q) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j}.$$

5 It follows that

$$\begin{aligned} \frac{\partial l(\mathbf{X})}{\partial \xi_j} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q w(A_q) \times L(\mathbf{X}_i|A_q) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \\ &= \sum_{i=1}^N \sum_{q=1}^Q \frac{1}{L(\mathbf{X}_i)} L(\mathbf{X}_i|A_q) \times w(A_q) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \\ &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q|\mathbf{X}_i) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \xi_j} \end{aligned}$$

6 Let P_{jq} be the 2PL, as in,

$$P_{jq} = \frac{\exp[1.7\alpha_j(A_q - \beta_j)]}{1 + \exp[1.7\alpha_j(A_q - \beta_j)]}.$$

7 This means we have a vector of parameters: $\xi_j = (\alpha_j, \beta_j)$. To maximize the loglikelihood, we need to
8 set the gradient $g(\xi)$ to $\mathbf{0}$, which is equivalent to

$$\begin{pmatrix} \frac{\partial l(\mathbf{X})}{\partial \alpha_j} \\ \frac{\partial l(\mathbf{X})}{\partial \beta_j} \end{pmatrix} = \mathbf{0}.$$

1 To solve these equations, we first show that

$$\begin{aligned}\frac{\partial P_{jq}}{\partial \alpha_j} &= \underline{1.7 \times P_{jq} \times (1 - P_{jq}) \times (A_q - \beta_j)} \\ \frac{\partial P_{jq}}{\partial \beta_j} &= \underline{-1.7 \times P_{jq} \times (1 - P_{jq}) \times \alpha_j}\end{aligned}$$

2 It follows that

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \alpha_j} &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] [1.7 \times P_{jq} \times (1 - P_{jq}) \times (A_q - \beta_j)] \\ &= 1.7 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times (X_{ij} - P_{jq}) \times (A_q - \beta_j) \\ \frac{\partial l(\mathbf{X})}{\partial \beta_j} &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] [-1.7 \times P_{jq} \times (1 - P_{jq}) \times \alpha_j] \\ &= -1.7 \alpha_j \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times (X_{ij} - P_{jq})\end{aligned}$$

3 An aside for now: These derivatives can be simplified as

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \alpha_j} &= 1.7 \sum_{q=1}^Q (A_q - \beta_j) \left[\sum_{i=1}^N P(A_q | \mathbf{X}_i) \times X_{ij} - P_{jq} \sum_{i=1}^N P(A_q | \mathbf{X}_i) \right] \\ &= 1.7 \sum_{q=1}^Q (A_q - \beta_j) \times [R_{jq} - P_{jq} \times N_q] \\ \frac{\partial l(\mathbf{X})}{\partial \beta_j} &= 1.7 \alpha_j \sum_{q=1}^Q \left[\sum_{i=1}^N P(A_q | \mathbf{X}_i) \times X_{ij} - P_{jq} \sum_{i=1}^N P(A_q | \mathbf{X}_i) \right] \\ &= -1.7 \alpha_j \sum_{q=1}^Q [R_{jq} - P_{jq} \times N_q]\end{aligned}$$

4 where R_{jq} is the *expected* number of correct response on item j for quadrature node A_q , and N_q is the
5 expected number of examinees in quadrature node A_q .

6
7 Because we are solving the zero of the function, we can use Newton's method, which is a zero-finder.
8 To implement Newton's method, we need $\nabla g(\boldsymbol{\xi}_j)$, the Jacobian matrix of $g(\boldsymbol{\xi}_j)$ at $\boldsymbol{\xi}_j$, which is, in this
9 case, a 2×2 matrix of second derivatives.

10

1 The following are the entries of $\nabla g(\boldsymbol{\xi}_j)$:

$$\begin{aligned}\frac{\partial^2 l^*(\mathbf{X})}{\partial \alpha_j^2} &= -1.7^2 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times (A_q - \beta_j)^2 \times P_{jq} \times (1 - P_{jq}) \\ \frac{\partial^2 l^*(\mathbf{X})}{\partial \beta_j^2} &= -1.7^2 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times \alpha_j^2 \times P_{jq} \times (1 - P_{jq}) \\ \frac{\partial^2 l^*(\mathbf{X})}{\partial \alpha_j \partial \beta_j} &= -1.7 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times [(X_{ij} - P_{jq}) - 1.7 \alpha_j \times (A_q - \beta_j) \times P_{jq} \times (1 - P_{jq})]\end{aligned}$$

2 At cycle $c + 1$, the updated estimate of $\boldsymbol{\xi}_j$ is given by

$$\boldsymbol{\xi}_j^{c+1} = \boldsymbol{\xi}_j^c - [\nabla g(\boldsymbol{\xi}_j^c)]^{-1} g(\boldsymbol{\xi}_j^c).$$

3 As can be seen from the above equation, Newton's method is an iterative solution. However, it should
4 be noted that the Newton cycles are nested within iteration t . We can denote the nested relationship
5 as follows:

$$\boldsymbol{\xi}_j^{t(c+1)} = \boldsymbol{\xi}_j^{t(c)} - [\nabla g(\boldsymbol{\xi}_j^{t(c)})]^{-1} g(\boldsymbol{\xi}_j^{t(c)}).$$

6 Once the Newton step has converged, we can suppress the superscript for the Newton cycle, and write
7 the item parameter estimate at iteration t as $\boldsymbol{\xi}_j^t$.

8
9 Before we can proceed with iteration $t + 1$, we need to update the posterior distributions using the
10 most current estimates (i.e., $\boldsymbol{\xi}_j^t$).

11
12 This process is repeated until no substantial improvement (i.e., parameter estimates or likelihood re-
13 mains *essentially* the same) between two successive iterations can be observed.

14
15 Below is an rough sketch of how to implement the algorithm.

- 16 1. At iteration $t = 0$, assign initial values to the item parameters, as in, $\boldsymbol{\xi}_1^0, \dots, \boldsymbol{\xi}_J^0$.
- 17 2. Compute the posterior distributions across the N examinees.
- 18 3. Find the updated estimates of the item parameters.
- 19 4. Repeat steps 2 and 3 until convergence.

20 To implement the MMLE algorithm, we will need to define the following $N \times Q$ matrices.

- 21 • $\mathbf{P} = \{P(A_q | \mathbf{X}_i)\}$ contains the posterior probabilities; each entry is potentially unique
- 22 • $\mathbf{R} = \{X_{ij} - P_{jq}\}$ contains the differences between observed and model-based responses; each
23 entry is potentially unique, and computed as $\langle X_{1j}; \dots; X_{Nj} \rangle - \langle P_{j1}; \dots; P_{jQ} \rangle'$
- 24 • $\mathbf{D} = \{A_q - \beta_j\}$ contains the differences between quadrature node and the difficulty parameter;
25 the rows are repeated, as in, $\mathbf{1}_{N \times 1} \times \langle A_1; \dots; A_Q \rangle' - \beta_j$;

1 Approximate Standard Errors of MLEs

2
3 The information matrix of the estimator of $\boldsymbol{\xi}$ is $\mathbf{I}(\boldsymbol{\xi}) = -E\{\partial^2 l(\mathbf{X})/\partial \boldsymbol{\xi}^2\}$. The second derivative of
4 the marginalized loglikelihood for item j with respect to the parameters ξ_{jk} and $\xi_{jk'}$ is given by

$$\frac{\partial^2 l(\mathbf{X})}{\partial \xi_{jk} \partial \xi_{jk'}} = \sum_{i=1}^N \left[\frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L^2(\mathbf{X}_i)}{\partial \xi_{jk} \partial \xi_{jk'}} - \frac{1}{L^2(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial \xi_{jk}} \times \frac{\partial L(\mathbf{X}_i)}{\partial \xi_{jk'}} \right].$$

5 The first term will vanish when the expectation is taken. Hence,

$$\begin{aligned} \frac{\partial^2 l(\mathbf{X})}{\partial \xi_{jk} \partial \xi_{jk'}} &= \sum_{i=1}^N \left[- \left(\sum_{q=1}^Q p(A_q) \frac{\partial L(\mathbf{X}_i | A_q)}{\partial \xi_{jk}} \right) \times \frac{1}{L^2(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial \xi_{jk'}} \right] \\ &= - \sum_{i=1}^N \frac{1}{L^2(\mathbf{X}_i)} \times \left[\sum_{q=1}^Q p(A_q) \frac{\partial L(\mathbf{X}_i | A_q)}{\partial \xi_{jk}} \right] \times \frac{\partial L(\mathbf{X}_i)}{\partial \xi_{jk'}}. \end{aligned}$$

6 Taking the derivatives will give us

$$\begin{aligned} & - \sum_{i=1}^N \frac{1}{L^2(\mathbf{X}_i)} \left\{ \sum_{q=1}^Q p(A_q) \times L(\mathbf{X}_i | A_q) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk}} \right\} \times \\ & \left\{ \sum_{q=1}^Q p(A_q) \times L(\mathbf{X}_i | A_q) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk'}} \right\} \end{aligned}$$

7 We can distribute one $L^{-1}(\mathbf{X}_i)$ to each of the two factors in the braces to simplify them. For example,
8 the first factor becomes

$$\begin{aligned} & \sum_{q=1}^Q \frac{p(A_q) L(\mathbf{X}_i | A_q)}{L(\mathbf{X}_i)} \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk}} \\ &= \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk}}. \end{aligned}$$

9 Therefore,

$$\begin{aligned} \frac{\partial^2 l(\mathbf{X})}{\partial \xi_{jk} \partial \xi_{jk'}} &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk}} \right\} \times \\ & \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \xi_{jk'}} \right\}, \end{aligned}$$

10 which represents the sum of the cross-products of expected values based on the examinees' posterior
11 distributions. Instead of computing the expectation, the information matrix can be approximated by
12 evaluating $\partial^2 l(\mathbf{X})/(\partial \xi_{jk} \partial \xi_{jk'})$ at $\hat{\boldsymbol{\xi}}$ using the observed \mathbf{X} resulting in $\mathbf{I}(\hat{\boldsymbol{\xi}})$. Finally, $\mathbf{I}^{-1}(\hat{\boldsymbol{\xi}})$ provides
13 an approximation of $\text{Cov}(\hat{\boldsymbol{\xi}})$, and the root of its diagonal elements represents the $\text{SE}(\hat{\boldsymbol{\xi}})$.

14
15 Although the two second derivatives appear to be notationally the same, the second derivatives to
16 obtain the standard errors do not invoke the use of the posterior distributions until *after* the necessary

1 derivatives have been taken.

2

3 For the 2PL, we will have the following second derivatives:

$$\begin{aligned}
 \frac{\partial^2 l(\mathbf{X})}{\partial \alpha_j^2} &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \alpha_j} \right\}^2 \\
 &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times [1.7 \times P_{jq}(1 - P_{jq}) \times (A_q - \beta_j)] \right\}^2 \\
 &= -1.7^2 \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times (X_{ij} - P_{jq}) \times (A_q - \beta_j) \right\}^2
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial^2 l(\mathbf{X})}{\partial \beta_j^2} &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \beta_j} \right\}^2 \\
 &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times [-1.7 \times P_{jq}(1 - P_{jq}) \times \alpha_j] \right\}^2 \\
 &= -1.7^2 \alpha_j^2 \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times (X_{ij} - P_{jq}) \right\}^2
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial^2 l(\mathbf{X})}{\partial \alpha_j \partial \beta_j} &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \alpha_j} \right\} \times \\
 &\quad \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \beta_j} \right\} \\
 &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times [1.7 \times P_{jq}(1 - P_{jq}) \times (A_q - \beta_j)] \right\} \times \\
 &\quad \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times [-1.7 \times P_{jq}(1 - P_{jq}) \times \alpha_j] \right\} \\
 &= -1.7^2 \alpha_j \sum_{i=1}^N \left\{ \sum_{q=1}^Q \underbrace{p(A_q | \mathbf{X}_i)}_P \times \underbrace{(X_{ij} - P_{jq})}_R \times \underbrace{(A_q - \beta_j)}_D \right\} \left\{ \sum_{q=1}^Q \underbrace{p(A_q | \mathbf{X}_i)}_P \times \underbrace{(X_{ij} - P_{jq})}_R \right\}
 \end{aligned}$$

1 *An Implementation of MMLE for 2PL in Ox: 2PL_MMLE_I.ox*

2

3 This program requires including only one header file: `#include <oxstd.h>`.

4

```
decl X,P,R,D,N=2000;
decl MaxIt=100;
decl MaxC=50;
decl MaxDiff=.001;
decl A,Q,W;
decl a,b;
decl J=15;
decl Par0,SE;
```

5 X is the $N \times J$ matrix of response data, P is the $N \times Q$ matrix of posterior probabilities, R is the $N \times Q$
6 matrix of $X_{iq} - P(X_j = 1|A_q)$, and D is the $N \times Q$ matrix of $A_q - \beta_j$.

7

8 `MaxIt`, `MaxC`, and `MaxDiff` are the maximum number of iterations, number of Newton cycles, and al-
9 lowable change in item parameter estimates between two consecutive iterations or cycles, respectively.

10

11 `A`, `Q`, and `W` are the quadrature nodes, number of quadrature nodes, and weight of the quadrature nodes,
12 respectively.

13

14 `a`, `b`, and `J` are discrimination parameters, difficulty parameters, and number of items, respectively.

15

16 `Par0` contains the parameter estimates, whereas `SE` the standard errors.

17

```
initiate(){
A=<-3.5:[.25]3.5>'; //quadrature nodes
Q=rows(A);
W=densn(A); //quadrature weights
a=ones(J,1); //specifying initial a parameters
b=.025+.95*meanc(X); //generating initial b parameters based on the responses
b=-quann(b');
Par0=a~b;
SE=zeros(J,2);
} //end initiate
```

18 The first function we will discuss is `initiate`. This function gives initial values to the variables that
19 will be used in the later functions.

20

21 At $t = 0$, `A`, `Q`, and `W` are initiated. All the discrimination parameters are assigned an initial value of
22 1. The difficulty parameter of item j is determined using the proportion correct p_j and its normal

1 quantile equivalent (after a correction to ensure that $0.025 \leq p_j \leq .975$).

2

```
update_post(){
P=<>;          // posterior matrix: empty --> N x Q
for (decl i=0;i<N;i++){
decl ex=exp(1.7*(ones(Q,1)*Par0[][0]')*(A*ones(1,J)-ones(Q,1)*Par0[][1]'));//Q x J
decl p=ex./(1+ex);          //Q by J
decl Xi=ones(Q,1)*X[i][];   //Q by J
decl post=exp(sumr(Xi.*log(p)+(1-Xi).*log(1-p))).*W; //Q x 1
post=post./sumc(post);      // Q by 1
P=P|post';                  // posterior matrix: N by Q
}}// end i; end update_post
```

3 The next function is `update_post`, which is designed to update the posterior distributions at the end
4 of each iteration.

5

6 For this functions, we start with a null `P`. This variable is built up one row at a time.

7

8 To find the posterior distribution for examinee i , we need the \mathbf{X}_i . However, it should be noted that
9 A_q is used in place of θ_i .

10

11 Note also that although it is more stable to work with the loglikelihood, the posterior requires the
12 actual likelihood, hence, the use of `exp()`.

13

14 To ensure that proper posteriors are calculated, a normalization step is required.

15

16 The next function, `newton`, is the most critical function of this program. This function estimates the
17 parameters one item at a time, thus, requires the input `j`.

18

19 The first seven lines declare and initiate the variables that will used inside the function. Note that it
20 is important to **reset** the `diff`, the largest difference in the parameter estimates between two cycles,
21 and `it`, the number of Newton cycles, when a new item is involved.

22

23 It should be noted that although the function is called `newton`, it also contains steps for computing `R`
24 and `D`. The computation of `P` is done in `update_post`.

25

26 As mentioned earlier, once we have `P`, `R` and `D`, finding the values of the gradient and Jacobian is a
27 matter of plugging in the right variables.

28

29 `Parjtmp` represents the updated estimates of the parameters of item j at this particular cycle. The
30 maximization continues for item j until `diff` remains large AND the number of cycles is below `MaxC`.

31

32 The last few lines update the initial values for next cycle, and add 1 to the cycle counter.

33

```

newton(const j){
decl Xj=X[][j]*ones(1,Q);    //N by Q
decl diff=0.1;
decl it=0;
decl aj=Par0[j][0]; // starts at the estimate from prior iteration
decl bj=Par0[j][1]; // starts at the estimate from prior iteration
decl Parjtmp;

while (diff>MaxDiff && it<MaxC){ //Newton cycle
decl d1=ones(2,1); // gradient vector
decl d2=ones(2,2); // Jacobian matrix
decl ex=exp(1.7*aj*(ones(N,1)*A'-bj'*ones(N,Q))); //N x Q
decl Prob=ex./(1+ex);           //N x Q
decl R=Xj-Prob; //N x Q
decl D=ones(N,1)*A'-bj; //N x Q
//Gradient
d1[0][0]=sumr(sumc(1.7*P.*R.*D));
d1[1][0]=sumr(sumc(-1.7*aj*P.*R));
//Jacobian (Hessian)
d2[0][0]=-sumr(sumc(1.7^2*P.*Prob.*(1-Prob).*(D.^2)));
d2[1][1]=-sumr(sumc(1.7^2*aj^2*P.*Prob.*(1-Prob)));
d2[1][0]=sumr(sumc(-1.7*P.*(R-1.7*aj*D.*Prob.*(1-Prob))));
d2[0][1]=d2[1][0];
Parjtmp=(aj|bj)-invert(d2)*d1;
diff=max(sqrt((Parjtmp-(aj|bj)).^2));
aj=Parjtmp[0];
bj=Parjtmp[1];
it++;
} //end Newton cycle
return Parjtmp;
} //end newton

```

1 As the name indicates, `findSE` is the function called to compute the standard errors *analytically* after
2 the item parameters have been estimated.

3
4 It starts by making sure that the posteriors are computed using the most recent item parameters es-
5 timates. It then finds the standard errors one item at a time.

6
7 Except for SE, the variables in this function are identical to those in `newton`.

8

```

findSE(){
update_post();
for (decl j=0; j<J; j++){
decl Xj=X[][j]*ones(1,Q);
decl aj=Par0[j][0];
decl bj=Par0[j][1];

decl d2=ones(2,2);
decl ex=exp(1.7*aj*(ones(N,1)*A'-bj'*ones(N,Q)));
decl Prob=ex./(1+ex);
decl R=Xj-Prob;
decl D=ones(N,1)*A'-bj;
d2[0][0]=-1.7^2*sumc(sumr(P.*R.*D).^2);
d2[1][1]=-1.7^2*sumc(sumr(aj*P.*R).^2);
d2[1][0]=-1.7^2*sumc(sumr(P.*R.*D).*sumr(aj*P.*R));
d2[0][1]=d2[1][0];
SE[j][]=sqrt(diagonal(-invert(d2)));
} //end j
} //end findSE

```

```

iterate(){
decl diff=0.1;
decl it=0;
while(diff>MaxDiff && it<MaxIt){
update_post();
decl Par1=zeros(J,2);
for (decl j=0; j<J; j++){
Par1[j][]=newton(j)';
} // end j
diff=max(sqrt((Par1-Par0).^2));
Par0=Par1;
it++;
} //end while
} //end iterate

```

- 1 The `iterate` function governs the movement from one iteration to another.
- 2
- 3 This function starts with declaring and initiating `diff` and `it`, which determine when the iterative
- 4 process should be terminated.
- 5
- 6 `Par1` temporarily holds updated parameter estimates, which is obtained one item at a time. For each
- 7 item, `newton` is called to do the heavy lifting.
- 8

1 After going through all the items, constraints are examined to determine whether another iteration
 2 is necessary. If it is, the updated estimates serve as the starting values and the iteration counter is
 3 increased by 1.

4

```
main(){
format(10000);
decl file=sprint("twopl.dat");
file=fopen(file);
fscan(file,"%#M",N,J,&X);
fclose(file);

initiate();
iterate();
findSE();
print("\n    The following results are based on MMLE_I algorithm \n");
print("\n          a          b          SE(a)          SE(b)");
print(Par0~SE);
} //end main
```

5 The main function reads the data, calls `initiate`, `iterate`, and `findSE`. Lastly, the estimates and
 6 standard errors are printed as screen output.

7

The following results are based on MMLE_I algorithm

a	b	SE(a)	SE(b)
1.8298	-0.96470	0.12453	0.033606
1.4379	-0.71420	0.085263	0.032320
1.3799	-0.54251	0.076313	0.030351
1.2632	-0.39560	0.070886	0.029683
1.3484	-0.22642	0.070125	0.027642
1.7499	-0.027444	0.10337	0.023287
1.7306	0.053287	0.097310	0.023628
0.77912	0.11880	0.044111	0.039933
1.2035	0.33102	0.064823	0.030640
1.2760	0.37469	0.069572	0.029967
1.1632	0.71873	0.065910	0.037458
1.6739	0.87003	0.10027	0.033419
1.5440	1.0198	0.098429	0.037284
0.79812	1.1151	0.050936	0.061815
1.7701	1.7574	0.19376	0.067353

1 MMLE Estimation when Only the Gradient is Available: 2PL_MMLE_II.ox

2
3 As you witness above, deriving the gradient and Hessian can be time consuming. In addition, if you
4 are not really careful, it is also easy to make minor mistakes that can have major implications. This
5 next program can be used to obtain the (marginalized) maximum likelihood estimate of the item
6 parameters when you only have the time and patience to derive the gradient.

7
8 Without computing ~~the second~~ the second derivatives, we will obtain the standard errors *numerically*
9 in conjunction with the gradient-based item parameter estimates. We will introduce this later.

10
11 For this program, in addition to including `oxstd.h`, we also need to import `maxsqp` and `solvenle` for
12 this program.

13
14 The variables declared globally are the same as those in `2PL_MMLE_I`, except for `j` and `fL`. The former
15 is the item identifier, and the latter will contain the full likelihood (i.e., the likelihoods based all the J
16 items and the estimates from the previous iteration).

17
18 Both need to be declared as global variables for the functions in this program.

19
20 The function below, `TwoPL_NLE`, contains `avF[0]`, which holds the system of nonlinear equations (i.e.,
21 the gradient) for which the zeros need to be found.

22
23 To define the elements of `avF[0]`, we will use the same variables defined earlier. The only difference
24 lies in using the reference to the item parameters of item j , rather than the values of the item param-
25 eters per se.

```
TwoPL_NLE(const avF, const vP){ //vP is <aj;bj>
decl ex=exp(1.7*vP[0].*(ones(N,1)*A'-vP[1]));
decl Prob=ex./(1+ex);
decl Xj=X[][j]*ones(1,Q);
decl R=Xj-Prob;
decl D=ones(N,1)*A'-vP[1]*ones(N,Q);
avF[0]=sumr(sumc(1.7*P.*R.*D))|sumr(sumc(-1.7*vP[0]*P.*R));
return 1;
} //end TwoPL_NLE
```

27 The function `update_post` is the same function in `2PL_MMLE_I`.

28
29 The `iterate` function for this program is largely the same as the previous `iterate` functions. The
30 major difference is in the use of `SolveNLE`, a zero-finder function. As input, it takes the function where
31 the system of equations is defined and the address of the initial values for the parameters of item j . It
32 returns the updated estimates to the same address.

```

iterate(){
decl diff=0.1;
decl it=0;
while(diff>MaxDiff && it<MaxIt){
update_post();
decl Par1=zeros(J,2);
for (j=0;j<J;j++){
decl vP=Par0[j] []';
SolveNLE(TwoPL_NLE,&vP);
Par1[j] []=vP';
}
diff=max(sqrt((Par1-Par0).^2));
Par0=Par1;
it++;
} //end while
} //end iterate

```

1 To estimate standard errors, we need the function TwoPL_DL_obj, fullLike and findSE_DL (not shown
2 here). All of them will be explained later.

3
4 The structure of the main function of this program is identical to that of 2PL_MMLE_I.

5

The following results are based on MMLE_II algorithm

a	b	DL SE(a)	DL SE(b)
1.8298	-0.96470	0.12650	0.033550
1.4379	-0.71420	0.083727	0.032176
1.3799	-0.54251	0.076882	0.030212
1.2632	-0.39560	0.068563	0.030135
1.3484	-0.22642	0.072029	0.027664
1.7499	-0.027444	0.098138	0.023397
1.7306	0.053287	0.096376	0.023590
0.77912	0.11880	0.044347	0.039954
1.2035	0.33102	0.063779	0.030675
1.2760	0.37469	0.068003	0.029953
1.1632	0.71873	0.064365	0.036824
1.6739	0.87003	0.10232	0.032290
1.5440	1.0198	0.096554	0.037094
0.79812	1.1151	0.049487	0.060885
1.7701	1.7574	0.17177	0.064979

The EM Algorithm

It is possible to bypass the calculation of (some of the) derivatives when estimating item parameters. One technique is the EM Algorithm. Generally speaking, the EM is a data-augmentation algorithm that augments the observed data with latent data to simplify calculations, and make the subsequent maximization easier.

In traditional maximum likelihood estimation, the estimate is obtained

$$\hat{\xi} = \arg \max_{\xi} L(\mathbf{X}|\xi) = \arg \max_{\xi} l(\mathbf{X}|\xi)$$

In latent variable modeling context, this is can be difficult to do.

To simplify the problem, \mathbf{X} is treated as *incomplete* data, and $\mathbf{Z} = (\mathbf{X}, \boldsymbol{\theta})$ as the *complete* data, where \mathbf{X} is observed, whereas $\boldsymbol{\theta}$ is missing. We assume that \mathbf{X} and $\boldsymbol{\theta}$ are jointly distributed as:

$$p(\mathbf{Z}|\xi) = p(\mathbf{X}, \boldsymbol{\theta}|\xi) = p(\boldsymbol{\theta}|\mathbf{X}, \xi)p(\mathbf{X}|\xi).$$

$L(\xi|\mathbf{Z}) = L(\xi|\mathbf{X}, \boldsymbol{\theta}) = p(\mathbf{X}, \boldsymbol{\theta}|\xi)$, defined as the *complete-data* likelihood function, is considered a random variable because $\boldsymbol{\theta}$ is *missing*, and assumed to be a random variable.

In contrast, we call $L(\xi|\mathbf{X})$ as the *incomplete-data* likelihood function.

In the EM algorithm, the E-step involves finding the expected value of the complete data log-likelihood with respect to $\boldsymbol{\theta}$ given the observed data \mathbf{X} and the current estimate of ξ , denoted by ξ^t . Specifically, this expectation is defined as:

$$\begin{aligned} Q(\xi, \xi^t) &= E[l(\mathbf{X}, \boldsymbol{\theta}|\xi^t)] \\ &= \int_{\boldsymbol{\theta}} \log p(\mathbf{X}, \boldsymbol{\theta}|\xi) \times p(\boldsymbol{\theta}|\mathbf{X}, \xi^t) d\boldsymbol{\theta} \\ &= \int_{\boldsymbol{\theta}} \sum_i \log p(\mathbf{X}_i, \boldsymbol{\theta}|\xi) \times p(\boldsymbol{\theta}|\mathbf{X}, \xi^t) d\boldsymbol{\theta} \\ &= \sum_i \int_{\boldsymbol{\theta}} \log p(\mathbf{X}_i, \boldsymbol{\theta}|\xi) \times p(\boldsymbol{\theta}|\mathbf{X}, \xi^t) d\boldsymbol{\theta} \\ &= \sum_i \int_{\theta} \log p(\mathbf{X}_i, \theta|\xi) \times p(\theta|\mathbf{X}_i, \xi^t) d\theta \end{aligned}$$

where $p(\theta|\mathbf{X}_i, \xi^t)$ is the posterior distribution of θ given the response vector \mathbf{X}_i and the most current estimates of the item parameters. We can approximate this expectation using quadrature nodes as follows:

$$Q(\xi, \xi^t) \approx \sum_i \sum_q \log p(\mathbf{X}_i, A_q|\xi) \times p(A_q|\mathbf{X}_i, \xi^t).$$

Using our previous notation, we can rewrite this as

$$Q(\xi, \xi^t) \approx \sum_{i=1}^N \sum_{q=1}^Q \log L(\mathbf{X}_i|A_q) \times P(A_q|\mathbf{X}_i).$$

1 The M-step involves maximizing the expectation to find

$$\boldsymbol{\xi}^{t+1} = \arg \max_{\boldsymbol{\xi}} Q(\boldsymbol{\xi}, \boldsymbol{\xi}^t).$$

2 This we can do one item at a time. To implement the M-step, first find:

$$\begin{aligned} \frac{\partial Q(\boldsymbol{\xi}, \boldsymbol{\xi}^t)}{\partial \boldsymbol{\xi}_j} &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \frac{\partial \log L(\mathbf{X}_i | A_q)}{\partial \boldsymbol{\xi}_j} \\ &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \frac{1}{L(\mathbf{X}_i | A_q)} \times \frac{\partial L(\mathbf{X}_i | A_q)}{\partial \boldsymbol{\xi}_j} \\ &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \frac{1}{L(\mathbf{X}_i | A_q)} \times L(\mathbf{X}_i | A_q) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \boldsymbol{\xi}_j} \\ &= \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \frac{\partial P_{jq}}{\partial \boldsymbol{\xi}_j} \end{aligned}$$

3 Note that this derivative is equal to the derivative given in the last equation on page 68, which indicates
4 that the MMLE results can be obtained via EM algorithm.

5
6 The remaining steps of the EM algorithm follows those of MMLE. For example, we can use numerical
7 algorithm to obtain $\boldsymbol{\xi}^{t+1}$.

9 *The Bock-Aitken Solution*

10 Bock and Aitken (1981) showed that the expected marginalized log-likelihood given the most recent
item parameter estimates can be written as

$$E[l(\mathbf{X}) | \boldsymbol{\xi}^t] = \sum_{q=1}^Q \sum_{j=1}^J R_{jq} \log P_{jq} + (N_q - R_{jq}) \times \log(1 - P_{jq}) + \text{constant},$$

11
12 where (again) R_{jq} is the expected number of correct response on item j for quadrature node A_q and
13 N_q is the expected number of examinees in quadrature node A_q based on the posterior distribution.

14
15 With respect to $\boldsymbol{\xi}_j$, the objective function based on the above equation can be written as

$$O_{EM} = \sum_{q=1}^Q [R_{jq} \times \log P_{jq} + (N_q - R_{jq}) \times \log(1 - P_{jq})].$$

To estimate α_j and β_j , we need to find

$$\begin{pmatrix} \frac{\partial O_{EM}}{\partial \alpha_j} \\ \frac{\partial O_{EM}}{\partial \beta_j} \end{pmatrix} = \mathbf{0}.$$

1

$$\begin{aligned}
\frac{\partial O_{EM}}{\partial \alpha_j} &= \sum_{q=1}^Q \left[\frac{R_{jq}}{P_{jq}} \times \frac{\partial P_{jq}}{\partial \alpha_j} - \frac{N_q - R_{jq}}{1 - P_{jq}} \times \frac{\partial P_{jq}}{\partial \alpha_j} \right] \\
&= \sum_{q=1}^Q \left[\frac{R_{jq}}{P_{jq}} - \frac{N_q - R_{jq}}{1 - P_{jq}} \right] \times \frac{\partial P_{jq}}{\partial \alpha_j} \\
&= \sum_{q=1}^Q \left[\frac{R_{jq} - P_{jq} \times N_q}{P_{jq} \times (1 - P_{jq})} \right] \times [1.7 \times P_{jq} \times (1 - P_{jq}) \times (A_q - \beta_j)] \\
&= 1.7 \sum_{q=1}^Q (A_q - \beta_j) \times [R_{jq} - P_{jq} \times N_q]
\end{aligned}$$

2 Similarly,

$$\frac{\partial O_{EM}}{\partial \beta_j} = -1.7 \alpha_j \sum_{q=1}^Q [R_{jq} - P_{jq} \times N_q]$$

Comparing these derivatives to the derivatives we set aside on page 69, we can see that

$$\begin{pmatrix} \frac{\partial l(\mathbf{X})}{\partial \alpha_j} \\ \frac{\partial l(\mathbf{X})}{\partial \beta_j} \end{pmatrix} = \begin{pmatrix} \frac{\partial O_{EM}}{\partial \alpha_j} \\ \frac{\partial O_{EM}}{\partial \beta_j} \end{pmatrix}.$$

3 This shows that optimizing the $l(\mathbf{X})$ is equivalent to optimizing O_{EM} . To avoid involving any deriva-
4 tives, we can optimize O_{EM} numerically.

5
6 It should be noted that the numerical optimization also involves the Hessian matrix, which is also
7 derived numerically. However, this Hessian matrix cannot be used to compute the standard errors -
8 these standard errors are underestimates because the objective function involves a single item rather
9 than the entire likelihood.

10

11 *Estimation using Objective Function based on the Bock-Aitken Soluation: 2PL_MMLE_EM*

12

13 For this program, in addition to including `oxstd.h`, we also need to import `maxsqp`.

14

15 The variables declared globally are the same as those in 2PL_MMLE_II, except for `Ng`, `Rg` and `SE_EM`.
16 The former two contain N_q and R_{jq} , respectively, and the last contains the standard errors estimated
17 using numerical second derivatives of O_{EM} , which, as stated before, underestimates the standard errors.

18

19 The function below, `TwoPL_MMLE_EM`, contains `adFunc[0]`, which holds O_{EM} , the objective function
20 we aim to maximize.

21

```

TwoPL_MMLE_EM(const vP, const adFunc, const avScore, const amHessian)
{
  //vP is <aj;bj>
  decl ex=exp(1.7*vP[0].*(A'-vP[1]));          //1 x Q
  decl Prob=ex./(1+ex);          //1 x Q
  adFunc[0] = sumr(Rg[j][].*log(Prob)+(Ng[j][]-Rg[j][]).*log(1-Prob));
  return 1;
}
//end TwoPL_MMLE_EM

```

1 To define the elements of `adFunc[0]`, we will use the same variables defined earlier, except `Rg` and
2 `Ng`. Note that `Rg` and `Ng` are $J \times Q$ matrices. `Prob` is a $1 \times Q$ matrix of probabilities and the O_{EM} is
3 calculated for item j .

4
5 The function `update_post` is the same as that in `2PL_MMLE_I`.

6
7 Based on the posteriors calculated from `update_post`, N_q and R_{jq} are calculated using the function
8 `NgRg` below. Recall that $N_q = \sum_i^N P(A_q|X_i)$ and $R_{jq} = \sum_i^N (P(A_q|X_i) \times X_i)$.
9

```

NgRg()
{
  Ng=ones(J,1)*sumc(P);          //J x Q
  Rg=zeros(J,Q); //J x Q

  for (decl j=0;j<J;j++){ //j loop: for each item
    decl Xjq=X[][j]*ones(1,Q);  //N x Q
    Rg[j][]=sumc(Xjq.*P);       //1 x Q
  }
  //end NgRg
}

```

```

findSE_MMLE_EM(){
  decl mhess;
  for (j=0;j<J;j++){
    decl vP=Par0[j][]';
    if (Num2Derivative(TwoPL_MMLE_EM, vP, &mhess)){
      SE_EM[j][]=sqrt(diagonal(-invert(mhess)));
    }
  }
  //end findSE_MMLE_EM
}

```

1 The function `findSE_MMLE_EM` can be used to estimate the standard errors based on the Hessian matrix
2 of O_{EM} . However, this is for pedagogical purposes only to show that the estimates of standard errors
3 are inaccurate when comparing with those obtained in `2PL_MMLE_I`.

4
5 The following `iterate` function is largely similar to the previous `iterate` function except two differ-
6 ences. First, in the E step, N_q and R_{jq} are calculated by calling `NgRg` function; and second, `MaxSQPF` is
7 used to maximize O_{EM} . The use of `decl low=<0,-10>'`; and `decl hi=<5,10>'`; specifies the bounds
8 of the item parameters. That is, $\hat{\alpha}_j \in (0, 5)$ and $\hat{\beta}_j \in (-10, 10)$. The declared `f` stores the value of O_{EM} .
9

```
iterate(){
decl diff=0.1;
decl it=0;

while(diff>MaxDiff && it<MaxIt){
update_post();
NgRg();
decl Par1=zeros(J,2);
for (j=0;j<J;j++){
decl f;
decl low=<0,-10>';
decl hi=<5,10>';
decl vP=Par0[j] [];
MaxSQPF(TwoPL_MMLE_EM, &vP, &f,0, 1, 0, 0, low, hi);
Par1[j] []=vP';
}
diff=max(sqrt((Par1-Par0).^2));
Par0=Par1;
it++;
} //end while
} //end iterate
```

10 The other three functions, `TwoPL_DL_obj`, `fullLike`, and `findSE_DL` (not shown here), are the same
11 as those in `2PL_MMLE_II`, and can be used to correctly estimate standard errors, and will be explained
12 later.

13

The following results are based on EM algorithm

a	b	DL SE(a)	DL SE(b)	EM SE(a)	EM SE(b)
1.8299	-0.96467	0.12650	0.033548	0.094765	0.028293
1.4379	-0.71418	0.083729	0.032175	0.068832	0.029378
1.3800	-0.54247	0.076880	0.030210	0.064348	0.028219
1.2633	-0.39563	0.068569	0.030134	0.058327	0.028672

1.3485	-0.22640	0.072031	0.027662	0.061115	0.026457
1.7501	-0.027453	0.098147	0.023396	0.078708	0.021970
1.7307	0.053281	0.096382	0.023589	0.077795	0.022177
0.77913	0.11879	0.044348	0.039953	0.040129	0.039298
1.2035	0.33101	0.063782	0.030674	0.055490	0.029480
1.2761	0.37464	0.068004	0.029951	0.058627	0.028619
1.1632	0.71867	0.064362	0.036822	0.056304	0.034710
1.6740	0.87003	0.10233	0.032292	0.083445	0.028971
1.5440	1.0197	0.096542	0.037088	0.079434	0.033140
0.79815	1.1151	0.049485	0.060877	0.044878	0.057427
1.7701	1.7573	0.17177	0.064982	0.12737	0.050932

1 *Estimation using only the Likelihood as the Objective function - The Jiffy Discounted Likelihood Tech-*
2 *nique: 2PL_MMLE_DL.ox*

4 This is another method for finding parameter estimates without using analytical derivatives. It pro-
5 vides correct item parameter estimates, as well as correct standard errors.

7 Additionally, it offers a way to use the item parameter estimates to numerically find the correct standard
8 errors (Note that we have not introduced how to obtain the correct standard errors for 2PL_MMLE_II
9 and 2PL_MMLE_EM). However, be forewarned: This approach is computationally slow!

11 Recall that the objective function in MMLE is

$$l(\mathbf{X}) = \sum_{i=1}^N \log \sum_{q=1}^Q L(\mathbf{X}_i | A_q) \times w(A_q).$$

12 As it is, it already involves two summations (i.e., across N examinees and across Q quadrature nodes).

14 Hidden in $L(\mathbf{X}_i | A_q)$, the likelihood of the response vector \mathbf{X}_i given A_q , is the product across the J
15 items. Thus, this problem involves an $N \times Q \times J$ array, which can be unwieldy to handle.

17 To be able to implement marginalized maximum likelihood estimation using only the (marginalized)
18 likelihood as objective function, we need to employ a trick. This trick requires us to recognize (i) that
19 we are estimating the parameters one item at a time, and (ii) that

$$L(\mathbf{X}_i | A_q) = \left[\prod_{j' \neq j} P_{j'q}^{X_{ij'}} (1 - P_{j'q})^{1-X_{ij'}} \right] \times P_{jq}^{X_{ij}} (1 - P_{jq})^{1-X_{ij}}.$$

20 Thus, when working with item j ,

$$L^{-j}(\mathbf{X}_i | A_q) = \prod_{j' \neq j} P_{j'q}^{X_{ij'}} (1 - P_{j'q})^{1-X_{ij'}} = \frac{L(\mathbf{X}_i | A_q)}{P_{jq}^{X_{ij}} (1 - P_{jq})^{1-X_{ij}}}$$

21 is considered to be a constant with respect to ξ_j .

22

1 At iteration t , we can compute $L^{(t-1)}(\mathbf{X}_i|A_q)$, the $L(\mathbf{X}_i|A_q)$ based on $\xi_1^{(t-1)}, \dots, \xi_J^{(t-1)}$, and hold this
 2 scalar constant across the J items.

3
 4 At the same time, we also need to discount the impact of item j by dividing $L^{(t-1)}(\mathbf{X}_i|A_q)$ by
 5 $P_{jq}^{(t-1)X_{ij}} (1 - P_{jq}^{(t-1)})^{1-X_{ij}}$, where $P_{jq}^{(t-1)}$ is the probability of success on item j given A_k and $\xi_j^{(t-1)}$.
 6

7 Hence, at iteration t , the item parameters of item j can be estimated by maximizing the following
 8 modified objective function:

$$\begin{aligned} l_j(\mathbf{X}) &= \sum_{i=1}^N \log \sum_{q=1}^Q \frac{L^{(t-1)}(\mathbf{X}_i|A_q)}{P_{jq}^{(t-1)X_{ij}} (1 - P_{jq}^{(t-1)})^{1-X_{ij}}} \times [P_{jq}^{X_{ij}} (1 - P_{jq})^{1-X_{ij}}] \times w(A_q) \\ &= \sum_{i=1}^N \log \sum_{q=1}^Q L^{-j(t-1)}(\mathbf{X}_i|A_q) \times [P_{jq}^{X_{ij}} (1 - P_{jq})^{1-X_{ij}}] \times w(A_q) \\ &= \sum_{i=1}^N \log \sum_{q=1}^Q [P_{jq}^{X_{ij}} (1 - P_{jq})^{1-X_{ij}}] \times \frac{P(A_q|\mathbf{X}_i)}{P_{jq}^{(t-1)X_{ij}} (1 - P_{jq}^{(t-1)})^{1-X_{ij}}} \end{aligned}$$

9 where candidate values for ξ_j affect only P_{jq} .

10
 11 It is worth noting that all the factors involved in this objective function are scalar (because $L^{(t-1)}(\mathbf{X}_i|A_q)$
 12 or $P(A_q|\mathbf{X}_i)$ can be computed elsewhere), allowing us to carry our maximization over a two-dimensional
 13 array (i.e., an $N \times Q$ matrix).

14
 15 Using the discounted likelihood technique (DLT), we can obtain the MMLE estimate using only the
 16 likelihood. In essence, we can implement MMLE in a *jiffy* (i.e., quickly).

17
 18 However, it should be noted that the DLT is only quick with respect our ability to implement MMLE,
 19 not the actual implementation time.

20
 21 Below are some of the functions in this program.

```
TwoPL_DL_obj(const vP, const adFunc, const avScore, const amHessian)//vP is <aj;b>
{
  decl ex0=exp(1.7*Par0[j][0]*(ones(N,1)*A'-Par0[j][1])); //N x Q
  decl Prob0=ex0./(1+ex0); //N x Q
  decl Xj=X[] [j]*ones(1,Q); //N X Q
  decl lj0=(Prob0.^Xj).*(1-Prob0).^(1-Xj);
  decl fL_j=fL./lj0;
  decl ex=exp(1.7*vP[0]*(ones(N,1)*A'-vP[1])); //Q x 1
  decl Prob=ex./(1+ex); //Q x J
  adFunc[0] = sumc(log(sumr(fL_j.*(Prob.^Xj).*((1-Prob).^(1-Xj)).*(ones(N,1)*W'))));
  return 1;
} //end TwoPL_DL
```

- 1 The function `TwoPL_DL_obj` contains the objective function to be maximized.
- 2
- 3 The variable `lj0` is the likelihood for item j than needs to be discounted from the full likelihood `fL`.
- 4 The “discounted likelihood” is given by `fL_j`.
- 5
- 6 Note that two types of probabilities are computed for item j : one based on parameter estimates from
- 7 the the previous iteration, and another based on values to be estimated at this iteration.
- 8
- 9 The objective function in `adFunc[0]` is an implementation of $l_j(\mathbf{X})$.

```

fullLike(){
fL=<>;                                     //empty -> N x Q
for(decl i=0;i<N;i++){
decl ex=exp(1.7*(ones(Q,1)*Par0[] [0] ') .*(A*ones(1,J)-ones(Q,1)*Par0[] [1] ')); //Q x J
decl Prob=ex./(1+ex);                      //Q x J
decl Xi=ones(Q,1)*X[i] [];                 //Q x J
decl tmp=sumr(Xi.*log(Prob)+(1-Xi).*log(1-Prob));
fL=fL|exp(tmp');
} //end i
} //end fullLike

```

- 10 The function `fullLike` finds `fL`. Note that it loops across the N examinees, which represents the N
- 11 rows. `Prob` is a $Q \times J$ matrix of probabilities, but reduces to $Q \times 1$ after taking the row sums. By
- 12 exponentiating, row sum is equivalent to row product; and by transposing, we get a $1 \times Q$ of likelihoods
- 13 for each examinee.

- 14
- 15 The function `findSE_DL` estimates the standard errors using the numerical second derivatives of the
- 16 “discounted likelihood”.
- 17

```

findSE_DL(){
fullLike();
decl mhess;
for (j=0;j<J;j++){
decl vP=Par0[j] [];
if (Num2Derivative(TwoPL_DL_obj, vP, &mhess)){
SE[j] []=sqrt(diagonal(-invert(mhess)));
} //end if
} //end j
} //end findSE_DL

```

1 Again, `iterate` for this program is highly similar to the previous `iterate` functions.

2

```
iterate(){
decl diff=0.1;
decl it=0;
while(diff>MaxDiff && it<MaxIt){
fullLike();
decl Par1=zeros(J,2);
for (j=0;j<J;j++){
decl f;
decl vP=Par0[j][]';
decl low=<0,-10>';
decl hi=<5,10>';
MaxSQPF(TwoPL_DL_obj, &vP, &f,0, 1, 0, 0, low, hi);
Par1[j][]=vP';
}
diff=max(sqrt((Par1-Par0).^2));
Par0=Par1;
it++;
} //end while
} //end iterate
```

3 Lastly, the main function has the same structure as the previous main functions.

4

The following results are based on MMLE_DL algorithm

a	b	DL SE(a)	DL SE(b)
1.8305	-0.96407	0.12653	0.033533
1.4384	-0.71361	0.083746	0.032160
1.3804	-0.54208	0.076912	0.030202
1.2637	-0.39517	0.068585	0.030124
1.3489	-0.22608	0.072056	0.027655
1.7506	-0.027098	0.098183	0.023389
1.7313	0.053604	0.096420	0.023583
0.77936	0.11895	0.044360	0.039941
1.2039	0.33117	0.063802	0.030665
1.2765	0.37484	0.068028	0.029943
1.1635	0.71876	0.064384	0.036812
1.6744	0.87002	0.10235	0.032281
1.5443	1.0197	0.096568	0.037082
0.79836	1.1150	0.049498	0.060860
1.7703	1.7573	0.17177	0.064969

1 *Numerical Calculation of the Standard Error using*

2
3 The DLT finds an estimate of the Hessian matrix in the process of estimating the item parameters.
4 As a result, the DLT finds the item parameters and the standard errors in one step, without requiring
5 any derivatives.

6
7 At the same time, the DLT can also be used to estimate standard errors using estimated item param-
8 eters. Particularly, it can be incorporated in 2PL_MMLE_II and 2PL_MMLE_EM.

9
10 As alluded to earlier, the functions TwoPL_DL_obj, fullLike and findSE_DL were already used in
11 2PL_MMLE_II and 2PL_MMLE_EM to estimate standard errors.

12
13 When using the DLT in conjunction with the estimated item parameters, the estimates are obtained
14 by solving a system of equations using SolveNLE in 2PL_MMLE_II or maximizing O_{EM} in 2PL_MMLE_EM.

15
16 After which, the findSE_DL is called, where the Hessian of the “discounted likelihood” is evaluated at
17 the estimates. With the Hessian matrix, the standard errors can be estimated.

18
19 1. Run the four programs (MMLE_I, MMLE_II, MMLE_EM, MMLE_DL) and compare the different
20 item parameter estimates and standard errors.

21
22 2. Note if there are discernable differences in the computational times.

23
24 3. What happens when you increase the sample size N ?

25
26 **Additional Exercises: Estimating the 1PL**

27
28 1. To derive the MMLE algorithm for 1PLM, we may start by finding the appropriate derivatives. Fill
29 in the missing steps.

30
Recall that the item response function for 1PLM is given by

$$P(X_{ij} = 1|\theta_i, \beta_j) = \frac{\exp(1.7 \times [\theta_i - \beta_j])}{1 + \exp[1.7 \times (\theta_i - \beta_j)]}.$$

31 Then, the marginalized likelihood is given as

$$\begin{aligned} L(\mathbf{X}_i) &= \int_{\theta} L(\mathbf{X}_i|\theta) \times g(\theta) d\theta \\ &\approx \sum_{q=1}^Q L(\mathbf{X}_i|A_q) \times w(A_q). \end{aligned}$$

32 For N individuals, the marginalized loglikelihood can be written as

$$\begin{aligned} l(\mathbf{X}) &= \log \prod_{i=1}^N L(\mathbf{X}_i) \\ &= \sum_{i=1}^N \log L(\mathbf{X}_i). \end{aligned}$$

1 The first derivative of the marginalized loglikelihood is

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \beta_j} &= \frac{\partial \left[\sum_{i=1}^N \log L(\mathbf{X}_i) \right]}{\partial \beta_j} \\ &= \left(\quad \right) \times \frac{\partial l(\mathbf{X}_i)}{\partial \beta_j},\end{aligned}$$

2 where

$$\begin{aligned}\frac{\partial l(\mathbf{X}_i)}{\partial \beta_j} &= \frac{\partial \left[\sum_{q=1}^Q L(\mathbf{X}_i | A_q) \right]}{\partial \beta_j} \\ &= \sum_{q=1}^Q w(A_q) \times \frac{\partial [L(\mathbf{X}_i | A_q)]}{\partial \beta_j},\end{aligned}$$

3 and

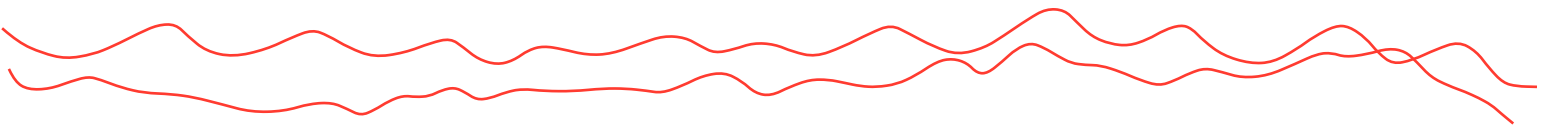
$$\begin{aligned}\frac{\partial L(\mathbf{X}_i | A_q)}{\partial \beta_j} &= \frac{\partial}{\partial \beta_j} \prod_{j'=1}^J \left(\quad \right) \\ &= \prod_{j' \neq j}^J P_{j'q}^{X_{ij}} \times (1 - P_{j'q})^{(1-X_{ij})} \times \frac{\partial \left(\quad \right)}{\partial \beta_j} \\ &= \prod_{j=1}^J P_{jq}^{X_{ij}} \times (1 - P_{jq})^{(1-X_{ij})} \times \frac{\partial P_{jq}}{\partial \beta_j} \times \left(\quad \right) \\ &= \prod_{j=1}^J \frac{\partial P_{jq}}{\partial \beta_j} \times \left(\quad \right)\end{aligned}$$

4 Finally,

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \beta_j} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q w(A_q) \times L(\mathbf{X}_i | A_q) \times \frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \times \frac{\partial P_{jq}}{\partial \beta_j} \\ &= \left(\quad \right) \\ &= -1.7 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times (X_{ij} - P_{jq}).\end{aligned}$$

Show that the second derivative is

$$\frac{\partial^2 l^*(\mathbf{X})}{\partial \beta_j^2} = -1.7^2 \sum_{i=1}^N \sum_{q=1}^Q P(A_q | \mathbf{X}_i) \times P_{jq} \times (1 - P_{jq}).$$



1 The above second derivative is for the purpose of estimating β_j . To obtain $SE(\beta_j)$, we need to find

$$\begin{aligned}\frac{\partial^2 l(\mathbf{X})}{\partial \beta_j^2} &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times \frac{\partial P_{jq}}{\partial \beta_j} \right\}^2 \\ &= - \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times \left[\frac{X_{ij} - P_{jq}}{P_{jq}(1 - P_{jq})} \right] \times (\quad) \right\}^2 \\ &= (\quad) \sum_{i=1}^N \left\{ \sum_{q=1}^Q p(A_q | \mathbf{X}_i) \times (\quad) \right\}^2\end{aligned}$$

2 2. Based on these derivations, write four programs that estimate the item parameters of the 1PL model
3 using

- 4
- 5 (a) a complete MMLE algorithm that involves the gradient and Hessian
- 6
- 7 (b) an MMLE algorithm that involves the gradient only, and the standard errors that are determined
8 *analytically* (same as in (a) above);
- 9 (c) an EM algorithm based on the Bock and Aitken solution; and
- 10
- 11 (d) optimization based on the original objective function only.
- 12

13 3. Write also two program that computes the standard errors *analytically* and *numerically*.

14

15 Generate the data using the same program in the 2PL example.

16

17 4. Compare the estimates, standard errors, and speed of the four programs for different values of N
18 and J .

19

20 Optional Exercises: Numerical Second Derivatives

21

22 Let's start with an example involving a single variable, and the first derivative. Let $f(x) = 3x^2 - x + 5$.
23 The first derivative of this function is $f'(x) = 6x - 1$. The slope at $x = 2$ is $f'(2) = 6(2) - 1 = 11$.

24

25 The second derivative is $f''(x) = 6$, which is a constant.

26

27 We can numerically by examining two points in the vicinity of $(2, f(2)) = (2, 15)$. This could be the
28 points $(x - h, f(x - h))$ and $(x + h, f(x + h))$, where h is a small number, say, 0.0001.

29

The slope (i.e., first derivative) of f at x can be approximated by

$$f'(x) \approx \frac{f(x + h) - f(x - h)}{2h},$$

30 which is simply the $\Delta f(x) / \Delta x$, or the familiar *rise over run*.

31

Applying the first derivative approximation to our example, we get

$$f'(2) \approx \frac{f(2.0001) - f(1.9999)}{2 \times 0.0001} = 11.00000000002,$$

1 which is close enough to 11 for government work.

2
3 Now, obtaining the second derivatives numerically requires applying the same principle, and recognizing that f'' is just the derivative of f' . That is, the second derivative is just the slope of the slopes (i.e., $\Delta f'(x) / \Delta x$).

6
7 For this we need three points: $(x, f(x))$, $(x - h, f(x - h))$ and $(x + h, f(x + h))$.

8
9 The second derivative of f at x can be approximated by

$$\begin{aligned} f''(x) &\approx \frac{\frac{f(x+h)-f(x)}{h} - \frac{f(x)-f(x-h)}{h}}{h} \\ &= \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}. \end{aligned}$$

Applying the second derivative approximation to our example, we get

$$f''(2) \approx \frac{f(2.0001) + f(1.9999) - 2f(2.0000)}{0.0001^2} = 6.0000005,$$

10 which, again, is close enough to 6 for government work.

11
12 We can apply these approximations to multiple variables, but let's limit our discussion to two variables.

13
14 The approximations are

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} &\approx [f(x+h, y) - 2f(x, y) + f(x-h, y)] / h^2 \\ \frac{\partial^2 f}{\partial y^2} &\approx [f(x, y+h) - 2f(x, y) + f(x, y-h)] / h^2 \\ \frac{\partial^2 f}{\partial x \partial y} &\approx [f(x-h, y-h) + f(x+h, y+h) - f(x-h, y+h) - f(x+h, y-h)] / (4h^2) \end{aligned}$$

15 Let $f(x, y) = x^3 + 2y^2 - 5xy^2$ be the function, and we want to find the second derivatives at the point
16 $(1, 2)$.

17
18 Analytically,

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} &= 6x \rightarrow 6 \\ \frac{\partial^2 f}{\partial y^2} &= 4 - 10x \rightarrow -6 \\ \frac{\partial^2 f}{\partial x \partial y} &= -10y \rightarrow -20 \end{aligned}$$

1 Numerically,

$$\frac{\partial^2 f}{\partial x^2} \approx [f(1.0001, 2) - 2f(1, 2) + f(0.9999, 2)] / 0.0001^2 \rightarrow 5.9999998$$

$$\frac{\partial^2 f}{\partial y^2} \approx [f(1, 2.0001) - 2f(1, 2) + f(1, 1.9999)] / 0.0001^2 \rightarrow -6.0000005$$

$$\begin{aligned} \frac{\partial^2 f}{\partial x \partial y} &\approx [f(0.9999, 1.9999) + f(1.0001, 2.0001) - f(0.9999, 2.0001) - f(1.0001, 1.9999)] / (4 \times 0.0001^2) \\ &\rightarrow -19.9999999 \end{aligned}$$

2 1. Write a function that numerically computes the second derivatives of a two-variable objective func-
3 tion.

4
5 2. Compare the approximations from this function to those of `Num2Derivative`.

6
7 3. Modify one of the estimation programs for the 2PL so that it derives the standard errors using your
8 function for the numerical second derivatives.

The Expectation-Maximization Algorithm - Part II

EM Algorithms for Cognitive Diagnosis Models

The EM algorithms for CDMs are straightforward extensions of the EM algorithms for traditional IRT models, at least when the attribute structure has a saturated formulation. However, we need to note a few differences when CDMs are involved. First, using quadrature nodes to approximate continuous θ is no longer necessary - α is already discrete. And second, for a class of CDMs (e.g., DINA, DINO, G-DINA), estimates of the item parameters can be obtained separately from a partitioning of the attribute vectors.

An EM Algorithm the G-DINA Model

Again, let's start with the marginalized likelihood of the response vector of examinee i , but this time involving CDMs:

$$L(\mathbf{X}_i) = \sum_{l=1}^L L(\mathbf{X}_i | \alpha_l) \times p(\alpha_l),$$

where $p(\alpha_l)$ is the prior probability of the attribute vector α_l , and $L = 2^K$.

For a dichotomous response,

$$L(\mathbf{X}_i | \alpha_l) = \prod_{j=1}^J P_{\alpha_{jl}^*}^{X_{ij}} \times (1 - P_{\alpha_{jl}^*})^{1-X_{ij}},$$

only consider the attributes that are measured in item j : reduced attribute vector

where $P_{\alpha_{jl}^*} = P(X_j = 1 | \alpha_{jl}^*)$ is the probability of a correct response to item j given the reduced attribute vector α_{jl}^* , and $jl = 1, 2, \dots, 2^{K_j}$, K_j being $\sum_{k=1}^K q_{jk}$. Note that for the G-DINA model, the number of parameters for item j is equal to 2^{K_j} .

total number of required attributes in item j

To obtain the MLE of \mathcal{P}_j , where $\mathcal{P}_j = P_{\alpha_{jl}^*}$, for $jl = 1, \dots, 2^{K_j}$, we need to maximize

maximizing likelihood is the same as maximizing the loglikelihood

$$l(\mathbf{X}) = \log \prod_{i=1}^N L(\mathbf{X}_i) = \sum_{i=1}^N \log L(\mathbf{X}_i)$$

with respect to \mathcal{P}_j . As in,

$$\begin{aligned} \frac{\partial l(\mathbf{X})}{\partial \mathcal{P}_j} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial \mathcal{P}_j} \\ &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{l=1}^L p(\alpha_l) \times \frac{\partial L(\mathbf{X}_i | \alpha_l)}{\partial \mathcal{P}_j}. \end{aligned}$$

Now, the derivative on the right-hand side can be shown to be

$$\frac{\partial L(\mathbf{X}_i | \alpha_l)}{\partial \mathcal{P}_j} = L(\mathbf{X}_i | \alpha_l) \times \left[\frac{X_{ij} - P_{\alpha_{jl}^*}}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \times \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j}.$$

1 We can use this in $\partial l(\mathbf{X})/\partial P_{\alpha_{jl}^*}$ to get

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \mathcal{P}_j} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{l=1}^L p(\alpha_l) \times L(\mathbf{X}_i|\alpha_l) \times \left[\frac{X_{ij} - P_{\alpha_{jl}^*}}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \times \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j} \\ &= \sum_{i=1}^N \sum_{l=1}^L \frac{L(\mathbf{X}_i|\alpha_l) \times p(\alpha_l)}{L(\mathbf{X}_i)} \times \left[\frac{X_{ij} - P_{\alpha_{jl}^*}}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \times \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j} \\ &= \sum_{i=1}^N \sum_{l=1}^L \left[\frac{X_{ij} - P_{\alpha_{jl}^*}}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \times P(\alpha_l|\mathbf{X}_i) \times \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j},\end{aligned}$$

where $p(\alpha_l|\mathbf{X}_i)$ is the posterior probability that examinee i has the attribute pattern α_l . At this juncture, we are dealing with two types of attribute patterns, α_l and α_{jl}^* . However, we know that for a subset of the attribute patterns, $\alpha_l \rightarrow \alpha_{jl}^*$. Hence,

$$P(\alpha_{jl}^*|\mathbf{X}_i) = \sum_{\forall l \ni \alpha_l \rightarrow \alpha_{jl}^*} P(\alpha_l|\mathbf{X}_i).$$

This will lead us to the following simplification:

$$\frac{\partial l(\mathbf{X})}{\partial \mathcal{P}_j} = \sum_{i=1}^N \sum_{jl=1}^{L_j} \left[\frac{X_{ij} - P_{\alpha_{jl}^*}}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \times P(\alpha_{jl}^*|\mathbf{X}_i) \times \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j},$$

2 where $L_j = 2^{K_j}$. This reduces the number of inside terms from 2^K to just 2^{K_j} .

3
4 By interchanging the summations, and making some rearrangements, we get,

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \mathcal{P}_j} &= \sum_{jl=1}^{L_j} \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \sum_{i=1}^N (X_{ij} - P_{\alpha_{jl}^*}) \times P(\alpha_{jl}^*|\mathbf{X}_i) \\ &= \sum_{jl=1}^{L_j} \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] \left[\sum_{i=1}^N X_{ij} \times P(\alpha_{jl}^*|\mathbf{X}_i) - P_{\alpha_{jl}^*} \sum_{i=1}^N P(\alpha_{jl}^*|\mathbf{X}_i) \right] \\ &= \sum_{jl=1}^{L_j} \frac{\partial P_{\alpha_{jl}^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] [R_{jl}^* - P_{\alpha_{jl}^*} N_{jl}^*],\end{aligned}$$

5 where $N_{jl}^* = \sum_{i=1}^N P(\alpha_{jl}^*|\mathbf{X}_i)$ is the expected number of examinees with *reduced* attribute pattern
6 α_{jl}^* , and $R_{jl}^* = \sum_{i=1}^N X_{ij} \times P(\alpha_{jl}^*|\mathbf{X}_i)$ is the expected number of examinees with attribute pattern α_{jl}^*
7 answering item j correctly.

8
9 By expanding the summation, the derivative can be written as

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \mathcal{P}_j} &= \frac{\partial P_{\alpha_1^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_1^*}(1 - P_{\alpha_1^*})} \right] [R_1^* - P_{\alpha_1^*} N_1^*] + \frac{\partial P_{\alpha_2^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_2^*}(1 - P_{\alpha_2^*})} \right] [R_2^* - P_{\alpha_2^*} N_2^*] \\ &\quad + \dots + \frac{\partial P_{\alpha_{L_j}^*}}{\partial \mathcal{P}_j} \times \left[\frac{1}{P_{\alpha_{L_j}^*}(1 - P_{\alpha_{L_j}^*})} \right] [R_{L_j}^* - P_{\alpha_{L_j}^*} N_{L_j}^*].\end{aligned}$$

We can now start taking the 2^{K_j} derivatives one at a time. For the j^{th} derivative, it should be clear that

$$\frac{\partial P_{\alpha_{jl}^*}}{\partial P_{\alpha_{jl}^*}} = \begin{cases} 0 & \text{if } j' \neq j \\ 1 & \text{if } j' = j \end{cases}.$$

Thus, only one of the 2^{K_j} terms will survive. For $\mathcal{P}_j = P_{\alpha_{jl}^*}$, this can be written as

$$\frac{\partial l(\mathbf{X})}{\partial P_{\alpha_{jl}^*}} = \left[\frac{1}{P_{\alpha_{jl}^*}(1 - P_{\alpha_{jl}^*})} \right] (R_{jl}^* - P_{\alpha_{jl}^*} N_{jl}^*).$$

To find $\hat{P}_{\alpha_{jl}^*}$, we need to find the zero of the derivative. Fortunately, this has a closed-form solution, as in,

$$\hat{P}_{\alpha_{jl}^*} = \frac{R_{jl}^*}{N_{jl}^*}.$$

Because we do not need to use an iterative solution, we also do not need the second derivatives.

It should be noted finding the parameters of item j involves solving the system of equations $\nabla l(\mathbf{X}) = \mathbf{0}$. However, each equation is independent of the other equations so the solution can be derived one equation at a time.

As before:

1. At iteration 0, assign initial values to the item parameters, as in, $P_{\alpha_1^*}^0, \dots, P_{\alpha_{L_j}^*}^0$.
2. At iteration t , find $(R_{jl}^{t*}, N_{jl}^{t*}), \dots, (R_{L_j}^{t*}, N_{L_j}^{t*})$ based on $P_{\alpha_1^*}^{t-1}, \dots, P_{\alpha_{L_j}^*}^{t-1}$.
3. From the expected quantities, find the updated item parameter estimates $P_{\alpha_1^*}^t, \dots, P_{\alpha_{L_j}^*}^t$.
4. Repeat steps 2 and 3 until convergence.

With α_l , it is harder to determine what an appropriate prior should be. Typically, we can start with a uniform prior, but more likely than not, it is not the correct prior. To provide a more accurate prior, we can update the prior distribution using the information available in the posterior distributions of the examinees. We can estimate the prior as

$$p(\alpha_l) = \sum_{i=1}^N P(\alpha_l | \mathbf{X}_i) / N.$$

In determining the prior distribution using information from the data (i.e., posterior distributions), we employ what is called an *empirical Bayes* method.

Computing the Standard Errors of the G-DINA Model Parameter Estimates

At this point, we have different ways of finding the standard errors of the item parameter estimates: one is to use the information matrix, and the other is to use a numerical solution.

The information matrix of the estimator of $\mathcal{P}_j = \{P_{\alpha_{jl}^*}\}$ is $\mathbf{I}(\mathcal{P}_j) = -E\{\partial^2 l(\mathbf{X}) / \partial \mathcal{P}_j^2\}$.

1 Again, we can start with

$$\begin{aligned}\frac{\partial^2 l(\mathbf{X})}{\partial P\alpha_{jl}^* \partial P\alpha_{jl'}^*} &= -\sum_{i=1}^N \frac{\partial l(\mathbf{X})}{\partial P\alpha_{jl}^*} \frac{\partial l(\mathbf{X})}{\partial P\alpha_{jl'}^*} \\ &= -\sum_{i=1}^N \left[\frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial P\alpha_{jl}^*} \right] \times \left[\frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial P\alpha_{jl'}^*} \right].\end{aligned}$$

2 Now,

$$\begin{aligned}\frac{\partial L(\mathbf{X}_i)}{\partial P\alpha_{jl}^*} &= \frac{\partial \left[\sum_{j'l'=1}^{L_j} L(\mathbf{X}_i | P\alpha_{j'l'}^*) \times w(P\alpha_{j'l'}^*) \right]}{\partial P\alpha_{jl}^*} \\ &= \sum_{j'l'=1}^{L_j} \frac{\partial L(\mathbf{X}_i | P\alpha_{j'l'}^*)}{\partial P\alpha_{jl}^*} \times w(P\alpha_{j'l'}^*) \\ &= \sum_{j'l'=1}^{L_j} \left[L(\mathbf{X}_i | P\alpha_{j'l'}^*) \times \frac{X_{ij} - P\alpha_{jl'}^*}{P\alpha_{j'l'}^* (1 - P\alpha_{j'l'}^*)} \times \frac{\partial P\alpha_{j'l'}^*}{\partial P\alpha_{jl}^*} \right] \times w(P\alpha_{j'l'}^*) \\ &= L(\mathbf{X}_i | P\alpha_{jl}^*) \times \frac{X_{ij} - P\alpha_{jl}^*}{P\alpha_{jl}^* (1 - P\alpha_{jl}^*)} \times w(P\alpha_{jl}^*)\end{aligned}$$

3 because only one of the terms will have a nonzero derivative (i.e., the term involving $P\alpha_{jl}^*$).

4

5 Thus,

$$\begin{aligned}\frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial P\alpha_{jl}^*} &= \frac{1}{L(\mathbf{X}_i)} \times \left[L(\mathbf{X}_i | P\alpha_{jl}^*) \times \frac{X_{ij} - P\alpha_{jl}^*}{P\alpha_{jl}^* (1 - P\alpha_{jl}^*)} \times w(P\alpha_{jl}^*) \right] \\ &= P(\alpha_{jl} | \mathbf{X}_i) \times \frac{X_{ij} - P\alpha_{jl}^*}{P\alpha_{jl}^* (1 - P\alpha_{jl}^*)}\end{aligned}$$

6 Evaluated at the solution, the second derivative of the marginalized loglikelihood for item j with respect
7 to the parameters $P\alpha_{jl}^*$ and $P\alpha_{jl'}^*$ is given by

$$\frac{\partial^2 l(\mathbf{X})}{\partial P\alpha_{jl}^* \partial P\alpha_{jl'}^*} = -\sum_{i=1}^N \left\{ P(\alpha_{jl} | \mathbf{X}_i) \times \left[\frac{X_{ij} - \hat{P}\alpha_{jl}^*}{\hat{P}\alpha_{jl}^* (1 - \hat{P}\alpha_{jl}^*)} \right] \right\} \times \left\{ P(\alpha_{jl'} | \mathbf{X}_i) \left[\frac{X_{ij} - \hat{P}\alpha_{jl'}^*}{\hat{P}\alpha_{jl'}^* (1 - \hat{P}\alpha_{jl'}^*)} \right] \right\}.$$

8 The negative of $\mathbf{I}^{-1}(\hat{\mathcal{P}}_j)$ provides an approximation of $\text{Cov}(\hat{\mathcal{P}}_j)$, and the root of its diagonal elements
9 represents the $\text{SE}(\hat{\mathcal{P}}_j)$.

10

11 *Implementation of EM Algorithm for the G-DINA Model in Ox*

12 

13 The program `GDINA_MMLE.ox` will estimate the item parameters of the G-DINA model. It requires
14 both item responses and the Q-matrix as input.

15

16 This program includes the following headers: `oxstd.h`, `oxprob.h`, and `oxfloat.h`.

17

1 We start by declaring a few global variables. The first set are variables **N**, sample size, **J**, test length,
2 **K**, number of attributes, **data**, item response data file name, **qmatrix**, Q-matrix filename, **crit**, maxi-
3 mum tolerable difference between iterations, and **maxit**, maximum number of iterations, are variables
4 that users can define.

5
6 The remaining global variables are the $J \times K$ Q-matrix **Q**, $N \times J$ data matrix **X**, $N \times 2^K$ posterior
7 distributions **post**, $2^K \times K$ matrix all possible attribute patterns **alpha**, $J \times 2^{K_j+1}$ matrix containing
8 the estimates and standard errors **est**, $2^K \times 1$ vector of log-prior probabilities **lprior**, and $2^K \times J$
9 matrix latent group memberships **eta**.

```
decl N=1000,J=30,K=5;
decl data={"sample.dat"};
decl qmatrix={"q_sample.txt"};
decl crit=0.001,maxit=999;
decl Q,X;
decl post,alpha,est;
decl lprior,eta;
```

10 The function pattern requires a number \mathcal{K} as input, and returns a $2^K \times \mathcal{K}$ matrix containing the 2^K
11 unique attribute binary patterns. Recall return matrix is designed such that the patterns are increasing
12 in terms of number of 1s first, and then decreasing in terms of numerical value. This can function will
13 be called to generate **alpha** (requires K), and all the possible latent groups for each item (requires K_j).

14

```
pattern(const KK){
decl patt=<>;
for(decl k=1;k<=KK;k++){//creating all possible patterns
decl tmp=<>;
decl base=zeros(2^(KK-k),1)|ones(2^(KK-k),1);
for(decl l=0;l<2^(k-1);l++){
tmp=tmp|base;}//end l
patt=patt~tmp;}//end k
patt=patt~sumr(patt);
decl tmp=zeros(2^KK,1);
for(decl k=0;k<KK;k++){
tmp=tmp-10^(K-k-1)*patt[] [k];}// end k
patt=patt~tmp;
//order the patterns by column KK than by column KK+1
patt=sortbyc(patt,KK~(KK+1));
return(patt[] [0:KK-1]);
};//end pattern
```

15 The function **alphaeta** creates the $2^K \times J$ matrix **eta**, where the l_j^{th} entry indicates which latent
16 group α_l will be classified with respect to item j . For example, if item 5 requires α_1 and α_3 of $K = 5$

1 attributes, then the attribute patterns 00001, 11000, 00111, and 10110, will be classified in latent
 2 groups 0, 1, 2, and 3, respectively.

3

```
alphaeta(){
eta=zeros(2^K,J);
for(decl i=0;i<2^K;i++){
for(decl j=0;j<J;j++){
decl Kj=sumr(Q[j] []);
decl Lj=vecindex(Q[j] []==1);
decl Qj=pattern(Kj);
decl tmp_alpha=alpha[i] [Lj];
eta[i] [j]=vecindex(prodr(ones(2^Kj)*tmp_alpha==Qj)==1);
}}//end j; end k
}//end alphaeta
```

4 The function `findpost` updates the posterior distributions one examinee at a time. The first step
 5 requires creating a $2^K \times J$ matrix of item parameters based on the locations indicated in `eta`.

6

```
findpost(){
decl par=zeros(2^K,J);
for(decl cl=0; cl<2^K;cl++){
for(decl j=0;j<J;j++){
par[cl] [j]=est[j] [eta[cl] [j]];//assign item par based on eta
}}//end cl; end j
for(decl i=0;i<N;i++){
decl XX=ones(2^K,1)*X[i] [];
decl tmp=sumr(XX.*log(par)+(1-XX).*log(1-par))+lprior;
tmp=exp(tmp);
post[i] []=(tmp./sumc(tmp))';
}}//end i; end findpost
```

7 The function `initiate` makes sure that the variables have the appropriate starting values and dimen-
 8 sions.

9

10 For example, `lprior` is assigned to have the log of a uniform prior; `est` is a $J \times 2^{K_{j(max)}}$ matrix of -1s;
 11 and $P(\mathbf{1}) = .8$ whereas $P(\textcircled{1}) = .2$.

12

13 This is also where `alpha` and `eta` are created.

14

```

initiate(){
alpha=pattern(K);
alphaeta();
lprior=ones(2^K,1)./(2^K);
lprior=log(lprior/sumc(lprior));
est=-1*ones(J,2^(maxc(sumr(Q)))));
for(decl j=0;j<J;j++){
decl Kj=sumr(Q[j] []);
est[j] [0:(2^Kj-1)]=.2;
est[j] [2^Kj-1]=.8;
} // end j
post=ones(N,2^K);
findpost();
} //end initiate

```

```

iterate(){
decl est0, tmp; // est0->estimates from prior iteration
decl maxabs=.1, it=0;
decl n1, r1; // expected # of examinees & correct resp
while (maxabs>crit && it<maxit){
est0=est;
for(decl j=0;j<J;j++){ //loop for items
decl Kj=sumr(Q[j] []);
decl Lj=vecindex(Q[j] [] .==1);
decl Qj=pattern(Kj);
decl tmp_alpha=alpha[] [Lj]; //reduced alpha patterns
for(decl kj=0;kj<2^Kj;kj++){ //loop for latent groups
tmp=ones(2^K,1)*Qj[kj] [] .==tmp_alpha;
tmp=vecindex(prodr(tmp) .==1); //index of att. patt. in group kj
tmp=sumr(post[] [tmp]);
n1=sumc(tmp);
r1=sumc(X[] [j] .*tmp);
est[j] [kj]=r1/n1; //estimate of P_j(kj)
} //end kj; end j
maxabs=max(sqrt((est0-est).^2));
it=it+1;
findpost();
lprior=log(meanc(post))'; //update prior
} //end while; end iterate

```

- 1 The function `iterate` is where the item parameters are estimated. It starts by declaring and initiating
- 2 the variables that will be used within the function.

3

1 Note that the variable `tmp` takes on several meaning within the function.
2
3 The `while` loop is in effect as long as the difference in parameter estimates between two iterations
4 remain large AND the maximum number of iterations has not been reached.
5
6 The item parameters are estimated one item at a time, and within an item, one latent group at a time.
7
8 After an iterations, the variables in the `while` conditions are updated. The posterior distributions of
9 the examinees are updated, and the mean probabilities across examinees serves as the updated prior.
10

```
info_factor(const j,const tmp_alpha,const Qj, const kj){
decl tmp=ones(2^K,1)*Qj[kj][].==tmp_alpha;//Pj
tmp=vecindex(prodr(tmp).==1);
tmp=sumr(post[][tmp]); //prob
decl P=est[j][kj]*ones(N,1);
return(tmp.*(X[][j]-P)./(P.*(1-P)));
}
```

11 The function `info_factor` computes the factors that will used in computing the information matrix.
12

```
stderr(){
decl KKj=2^maxc(sumr(Q));
decl est2=ones(J,2^(maxc(sumr(Q))));
findpost();
for(decl j=0;j<J;j++){
decl Kj=sumr(Q[j][]);
decl Lj=vecindex(Q[j][]==1);
decl tmp_se2=zeros(2^Kj,2^Kj);//information matrix-upper triangle only
decl Qj=pattern(Kj);
decl tmp_alpha=alpha[][Lj]; //reduced alpha pattern
for(decl kj1=0; kj1<2^Kj;kj1++){
for(decl kj2=kj1;kj2<2^Kj;kj2++){
decl tmp_est1=info_factor(j,tmp_alpha,Qj,kj1);//finds the first factor
decl tmp_est2=tmp_est1;
if(kj2>kj1) tmp_est2=info_factor(j,tmp_alpha,Qj,kj2);//finds the second factor
tmp_se2[kj1][kj2]=sumc(tmp_est1.*tmp_est2);//sum of cross products across N
}}//end kj1,kj2
tmp_se2=tmp_se2+tmp_se2'-diag(diagonal(tmp_se2));//completing the matrix
est2[j][0:2^Kj-1]=sqrt(diagonal(invert(tmp_se2)));
} //end j
est=est~est2;
} //end stderr
```

1 The function `stderr` computes standard errors of the item parameter estimates, and updates `est` to
2 include both item parameter estimates and standard errors.
3
4 The function creates `est2`, a $J \times 2^{K_{j(max)}}$ matrix of 1s, to hold the standard errors. Prior to computing
5 the standard errors, the posterior distributions of the examinees are updated.
6
7 The standard errors are computed one item at a item. For greater efficiency, only the upper triangle
8 of the information matrix, `tmp_se2`, is computed, and `info_factor` is called once for the diagonal
9 entries, and twice for off-diagonal entries.
10
11 The information matrix is filled before additional operations are carried out.
12
13 Lastly, the standard errors in `est2` are appended to the parameter estimates in `est`.
14

```

printout(){
decl file0;
file0=fopen("est.out","w");
print("\nParameter estimates and SEs in the file: est.out\n");
fprintf(file0,"%#M",est);
fclose(file0);
decl cl=zeros(2^K,1);
for(decl k=0;k<K;k++){
cl=cl+alpha[][k]*10^(K-k-1);}
decl par=sprintf("    \"%0",K,".0f\"");
decl postp=meanc(post)';
print("\n","Latent Classes and their Posterior Probabilities:");
print("%cf",{par," %.4f ",par," %.4f"},
(cl~postp)[0:(2^K)/2-1] []~(cl~postp)[(2^K)/2:2^K-1] [],"\n");
print("Estimates of Attribute Prevalence:");
decl out=(range(1,K)|sumc(alpha.*(meanc(post)'*ones(1,K))))';
print("%cf",{ "          %0 .0f", "          %.4f "},out);
} //end printout

```

15 In addition to printing the filename where `est` can be found, the function `printout` also performs
16 additional computations from the posterior distributions of the examinees.
17
18 As mentioned earlier, it computes an overall posterior distribution by averaging the posterior distri-
19 butions of the examinees.
20
21 The prevalence of each attribute in the population is estimated by computing the marginal probability
22 from the overall posterior distribution.
23
24 Finally, the `main` function reads the data and Q-matrix, and calls the functions `initiate`, `iterate`,
25 `stderr`, and `printout`.
26

1 Additional Exercises

1. Modify the program so it can provide the EAP and MAP estimates of each examinee's attribute vector.

2. Modify the program so it can estimate the parameters of the DINA model.

Hints for modifying `GDINA_MMLE.ox`:

To be consistent with the G-DINA notation, instead of g_j and s_j , we will estimate the parameters $g_j = P_j(\mathbf{1})$ and $1 - s_j = P_j(\mathbf{1})$. We can get the estimate of s_j and its standard error by applying a simple linear transformation, as in, $\hat{s}_j = 1 - \hat{P}_j(\mathbf{1})$.

The function `pattern` is unchanged.

The function `alphaeta` is dramatically simplified. Recall in the DINA generation code, we first use `alpha*Q'` to create a $2^K \times J$ matrix that counts the number of required attributes mastered by each attribute pattern for each item. We then convert the entry in row l column j to 1 if it is equal to K_j ; otherwise, it is converted to 0.

For the function `findpost`, we do not need the two loops to find the value of `par` - it is either `est[] [1]` if `eta` is 1, and `est[] [0]` if `eta` is 0. You should be able to write this in one line by using `eta` and `1-eta`. The loop for the N examinees is unchanged.

We only need to change `est` to an $J \times 2$ matrix, with the first column entries set to 0.2, and the second column entries to 0.8. No loop necessary.

In the function `iterate`, the changes are in loop for the items. We do not need K_j , L_j , Q_j , and α_j^* when dealing with the DINA model. The loop for estimating the item parameters by latent group can also be taken out. The number of steps is fixed to two - one to find $\hat{P}_j(\mathbf{1})$, and another $\hat{P}_j(\mathbf{1})$.

To estimate $\hat{P}_j(\mathbf{1})$, we need to get the indices of the entries of `eta[] [j]` that are equal to 1. We the proper index, `r1`, `n1`, and $\hat{P}_j(\mathbf{1})$ is unchanged.

To estimate $\hat{P}_j(\mathbf{1})$, we can find the indices of the entries of `eta[] [j]` that are equal to 0. However this is not necessary. We just need to know that $r_0 + r_1 = \sum_{\forall i} X_{ij}$ and $n_0 + n_1 = N$.

The rest of the steps remain unchanged.

For the function `info_factor`, we only need to know two things: which item is needed, `j`, and which group is involved, `kj`. The remaining argument can be dropped.

For the DINA model, there are only two possible values for `kj`, 0 or 1. The indexing line should involved `eta[] [j]` and `kj` to make it more generic. This is the only change for this function.

For the function `stderr`, a few simplifications are needed. One, `est2` and `tmp_se2` have fixed dimensions.

Again, in the `j` loop, we do not need K_j , L_j , Q_j , and α_j^* . In addition, the upper limits of the loop is

1 fixed - this refers to the number of item j parameters.

2
3 Next, `tmp_est1` needs to be created from the simplified `info_factor`. We can assign `tmp_est2` to
4 be equal to `tmp_est1`. This assignment is to save time, and is correct when `kj1` and `kj2` are equal.
5 However, when they are not (i.e., `kj1` is greater than `kj2`), we need to call `info_factor` once more.

6
7 The remaining steps of `stderr` are unchanged.

8
9 The functions `printout` and `main` can be used as they are.

10
11 Those who are interested reporting \hat{s}_j , what should be $SE[1 - \hat{P}_j(\mathbf{1})]$?

12
13 To verify your results, compare the results by running the DINA estimation program you were given
14 last semester.

15 Optional Exercises

16 1. DINA Program from Scratch

17
18 **Warning: This is NOT for the faint of heart!**

19
20 Instead of modifying the G-DINA estimation program, you may want to write the DINA estimation
21 program from scratch, including the finding the first derivatives and the information matrix.

22 2. Alternative Ways of Parameter Estimation

23
24
25 a) Because the estimates of the DINA and G-DINA models can be expressed in closed form, we do not
26 need the second derivatives. However, we still the first derivatives to carry out the E-step, and find
27 N_{lj}^* and R_{lj}^* . As an exercise, use `SolveNLE` in the M-step.

28
29 b) How about estimating the G-DINA model parameters without any derivatives, as in, use `MaxSQPF`?

30
31 c) Still another way of estimating the DINA model is to carry out the G-DINA estimation E-steps,
32 and use `SolveNLE` in the M-step. What would system of equations need to be solved to obtain $\hat{P}_j(\mathbf{1})$,
33 and another $\hat{P}_j(\mathbf{1})$ given N_{lj}^* and R_{lj}^* , $lj = 1, \dots, 2^{K_j}$.

34
35 3. Related to 2b), can we also obtain the standard errors numerically?

36
37
38

The Expectation-Maximization Algorithm - Part III

EM Algorithms for Cognitive Diagnosis Models

Previously, we learned how to estimate the structural parameters of item response models with underlying discrete and continuous latent variables. In this section, we will combine the two models by estimating the structural parameters of higher-order CDMs. In addition to the item parameters (e.g., $P_{\alpha_{jl}^*}$), the structural parameters include the parameters that relate θ_i and α_i (i.e., λ_{0k} , λ_1).

As mentioned previously, the complete specification involved a model for the attribute structure $P(\alpha_l)$, and a model for the item response (i.e., CDM).

Today, we will focus on an EM algorithm where $P(\alpha_l)$ is expressed in terms of a higher-order model (i.e., $P(\alpha_l|\theta)$), and use the DINA model as the CDM. The algorithm can be easily extended for other CDMs.

An EM Algorithm for the Higher-Order DINA Model

Let the marginalized likelihood of the response vector of examinee i be

$$L(\mathbf{X}_i) = \int_{\theta} \sum_{l=1}^L L(\mathbf{X}_i|\alpha_l) \times P(\alpha_l|\theta) \times p(\theta) d\theta,$$

where

$$L(\mathbf{X}_i|\alpha_l) = \prod_{j=1}^J P_{jl}^{X_{ij}} (1 - P_{jl})^{1-X_{ij}},$$

$$P(\alpha_l|\theta) = \prod_{k=1}^K P(\alpha_k = 1|\theta)^{\alpha_{lk}} [1 - P(\alpha_k = 1|\theta)]^{1-\alpha_{lk}},$$

$P_{jl} = P(X_j = 1|\alpha_l)$ is the probability of a correct response to item j given the attribute pattern α_l ; $P(\alpha_l|\theta)$ and $p(\alpha_k = 1|\theta)$ are the probabilities of the attribute pattern α_l and mastering attribute k conditional on θ , respectively; and $p(\theta)$ is assumed to be $N(0, 1)$.

The conditional probability of an attribute mastery given θ is defined as

$$P(\alpha_k = 1|\theta) = \frac{\exp[\lambda_1(\theta - \lambda_{0k})]}{1 + \exp[\lambda_1(\theta - \lambda_{0k})]}, \text{ for } k = 1, \dots, K,$$

where λ_1 is the common attribute discrimination parameter, and λ_{0k} is the location parameter with respect to attribute k .

Because α_{jk} is latent, a simpler higher-order model may be necessary to get stable parameter estimates.

As in the past, we can use quadrature nodes to approximate the marginalized likelihood, as in,

$$L(\mathbf{X}_i) \approx \sum_{q=1}^Q \sum_{l=1}^L L(\mathbf{X}_i|\alpha_l) \times P(\alpha_l|A_q) \times W(A_q),$$

where A_q , $W(A_q)$, and Q are the quadrature node, weight of the node, and number of nodes, respectively.

Recall that for the DINA model, given the j^{th} row of the Q-matrix \mathbf{q}_j ,

$$P_{jl} = \begin{cases} g_j & \text{if } \alpha'_l \mathbf{q}_j < \mathbf{q}'_j \mathbf{q}_j \\ 1 - s_j & \text{if } \alpha'_l \mathbf{q}_j = \mathbf{q}'_j \mathbf{q}_j \end{cases}.$$

Let $\boldsymbol{\beta} = \{\beta_p\} = \{\{\boldsymbol{\beta}_{j\eta}\}, \lambda_{01}, \dots, \lambda_{0K}, \lambda_1\}$, where $\beta_{j0} = g_j$ and $\beta_{j1} = s_j$. For the higher-order DINA model with a common attribute discrimination parameter, the total number of parameters is $(2 \times J) + K + 1$.

To obtain the MLE of β_p , maximize

$$l(\mathbf{X}) = \log \prod_{i=1}^N L(\mathbf{X}_i) = \sum_{i=1}^N \log L(\mathbf{X}_i)$$

with respect to β_p :

$$\begin{aligned} \frac{\partial l(\mathbf{X})}{\partial \beta_p} &= \sum_{i=1}^N \frac{\partial L(\mathbf{X}_i)}{\partial \beta_p} / L(\mathbf{X}_i) \\ &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q \sum_{l=1}^L W(A_q) \frac{\partial L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \times P(\boldsymbol{\alpha}_l | A_q)}{\partial \beta_p}. \end{aligned}$$

The derivations need to be carried out one parameter at a time.

Case I: $\beta_p = \beta_{j\eta}$

The derivative in $\partial l(\mathbf{X}) / \partial \beta_p$ simplifies to

$$\frac{\partial L(\mathbf{X}_i | \boldsymbol{\alpha}_l) P(\boldsymbol{\alpha}_l | A_q)}{\partial \beta_{j\eta}} = P(\boldsymbol{\alpha}_l | A_q) \times \frac{\partial L(\mathbf{X}_i | \boldsymbol{\alpha}_l)}{\partial \beta_{j\eta}}.$$

As before, it can be shown that the derivative on the right-hand side is equal to

$$\frac{\partial L(\mathbf{X}_i | \boldsymbol{\alpha}_l)}{\partial \beta_{j\eta}} = L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \times \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{X_{ij} - P_{jl}}{P_{jl}(1 - P_{jl})} \right].$$

Thus, the original derivative can be written as

$$\frac{\partial l(\mathbf{X})}{\partial \beta_{j\eta}} = \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q \sum_{l=1}^L W(A_q) \times L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \times \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{X_{ij} - P_{jl}}{P_{jl}(1 - P_{jl})} \right].$$

By interchanging the summations, we can have the following simplifications:

$$\begin{aligned} \frac{\partial l(\mathbf{X})}{\partial \beta_{j\eta}} &= \sum_{l=1}^L \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \times \\ &\quad \sum_{i=1}^N \left[\frac{L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \times \sum_{q=1}^Q P(\boldsymbol{\alpha}_l | A_q) \times W(A_q)}{L(\mathbf{X}_i)} \right] \times [X_{ij} - P_{jl}] \end{aligned}$$

$$\begin{aligned}
&= \sum_{l=1}^L \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \sum_{i=1}^N \frac{L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \times P(\boldsymbol{\alpha}_l)}{L(\mathbf{X}_i)} \times [X_{ij} - P_{jl}] \\
&= \sum_{l=1}^L \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \sum_{i=1}^N P(\boldsymbol{\alpha}_l | \mathbf{X}_i) \times [X_{ij} - P_{jl}] \\
&= \sum_{l=1}^L \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \left[\sum_{i=1}^N X_{ij} \times P(\boldsymbol{\alpha}_l | \mathbf{X}_i) - P_{jl} \sum_{i=1}^N P(\boldsymbol{\alpha}_l | \mathbf{X}_i) \right] \\
&= \sum_{l=1}^L \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] [R_{jl} - P_{jl}N_l],
\end{aligned}$$

where $P(\boldsymbol{\alpha}_l | \mathbf{X}_i)$ is the posterior probability that examinee i has the attribute pattern $\boldsymbol{\alpha}_l$, $N_l = \sum_{i=1}^N P(\boldsymbol{\alpha}_l | \mathbf{X}_i)$ is the expected number of examinees with attribute pattern $\boldsymbol{\alpha}_l$, and $R_{jl} = \sum_{i=1}^N P(\boldsymbol{\alpha}_l | \mathbf{X}_i) X_{ij}$ is the expected number of examinees with attribute pattern $\boldsymbol{\alpha}_l$ answering item j correctly.

It should be noted that, in obtaining $P(\boldsymbol{\alpha}_l) = \sum_{q=1}^Q P(\boldsymbol{\alpha}_l | A_q) W(A_q)$, the values of $\lambda_{01}, \dots, \lambda_{0K}, \lambda_1$, are fixed to some values (i.e., estimates from the previous iteration).

For item j , this derivative can be written such that the summation across the L attribute patterns are partitioned into two summations: one for $\eta_j = 0$, and one for $\eta_j = 1$. As in,

$$\begin{aligned}
\frac{\partial l(\mathbf{X})}{\partial \beta_{j\eta}} &= \sum_{\{\boldsymbol{\alpha}_l: \boldsymbol{\alpha}'_l \mathbf{q}_j < \mathbf{q}'_j \mathbf{q}_j\}} \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \times [R_{jl} - P_{jl}N_l] + \\
&\quad \sum_{\{\boldsymbol{\alpha}_l: \boldsymbol{\alpha}'_l \mathbf{q}_j = \mathbf{q}'_j \mathbf{q}_j\}} \frac{\partial P_{jl}}{\partial \beta_{j\eta}} \times \left[\frac{1}{P_{jl}(1 - P_{jl})} \right] \times [R_{jl} - P_{jl}N_l] \\
&= \frac{\partial g_j}{\partial \beta_{j\eta}} \times \left[\frac{1}{g_j[1 - g_j]} \right] \sum_{\{\boldsymbol{\alpha}_l: \boldsymbol{\alpha}'_l \mathbf{q}_j < \mathbf{q}'_j \mathbf{q}_j\}} [R_{jl} - g_j N_l] + \\
&\quad \frac{\partial(1 - s_j)}{\partial \beta_{j\eta}} \times \left[\frac{1}{(1 - s_j)s_j} \right] \sum_{\{\boldsymbol{\alpha}_l: \boldsymbol{\alpha}'_l \mathbf{q}_j = \mathbf{q}'_j \mathbf{q}_j\}} [R_{jl} - (1 - s_j)N_l] \\
&= \frac{\partial g_j}{\partial \beta_{j\eta}} \times \left[\frac{1}{g_j[1 - g_j]} \right] \times [R_{jl}^{(0)} - g_j N_{jl}^{(0)}] + \\
&\quad \frac{\partial(1 - s_j)}{\partial \beta_{j\eta}} \times \left[\frac{1}{(1 - s_j)s_j} \right] \times [R_{jl}^{(1)} - (1 - s_j)N_{jl}^{(1)}],
\end{aligned}$$

where $N_{jl}^{(0)}$ is the expected number of examinees lacking at least one of the required attributes for item j , and $R_{jl}^{(0)}$ is the expected number of examinees among $N_{jl}^{(0)}$ correctly answering item j . $N_{jl}^{(1)}$ and $R_{jl}^{(1)}$ have the same interpretation except that they pertain to the examinees with all the required attributes for item j . $N_{jl}^{(0)} + N_{jl}^{(1)}$ is equal to N_l for all j .

When $\eta = 0$ (i.e., $\beta_{j0} = g$), $\partial P_{jl} / \partial \beta_{j\eta}$ is 1 for the first term, and 0 for the second term. Thus, the maximization of $\partial l(\mathbf{X})$ with respect to β_{j0} simplifies to solving for g_j in the equation

$$\left[\frac{1}{g_j(1 - g_j)} \right] \times [R_{jl}^{(0)} - g_j \times N_{jl}^{(0)}] = 0.$$

1 This gives the estimator $\hat{g}_j = R_{jl}^{(0)}/N_{jl}^{(0)}$.

2
3 Similarly, the maximization of $\partial l(\mathbf{X})$ with respect to β_{j1} is equivalent to solving for s_j in the equation

$$-\left[\frac{1}{(1-s_j)s_j}\right] \times [R_{jl}^{(1)} - (1-s_j) \times N_{jl}^{(1)}] = 0,$$

4 which results in the estimator $\hat{s}_j = [N_{jl}^{(1)} - R_{jl}^{(1)}]/N_{jl}^{(1)}$.

5
6 **Case II:** $\beta_p = \lambda_{0k}, k = 1, \dots, K$

7
With respect to the higher-order parameters, the derivative in $\partial l(\mathbf{X})/\partial \beta_p$ simplifies to

$$\frac{\partial L(\mathbf{X}_i|\boldsymbol{\alpha}_l)P(\boldsymbol{\alpha}_l|A_q)}{\partial \lambda_{0k}} = L(\mathbf{X}_i|\boldsymbol{\alpha}_l) \times \frac{\partial P(\boldsymbol{\alpha}_l|A_q)}{\partial \lambda_{0k}}.$$

8 Treating α_k like X_{ij} , the derivative on the right-hand side can be written as

$$\begin{aligned} \frac{\partial P(\boldsymbol{\alpha}_l|A_q)}{\partial \lambda_{0k}} &= \prod_{k=1}^K P(\alpha_k = 1|A_q)^{\alpha_{lk}} [1 - P(\alpha_k = 1|A_q)]^{1-\alpha_{lk}} \\ &\quad \times \frac{\partial P(\alpha_k = 1|A_q)}{\partial \lambda_{0k}} \times \left\{ \frac{\alpha_{lk} - P(\alpha_k = 1|A_q)}{P(\alpha_k = 1|A_q) \times [1 - P(\alpha_k = 1|A_q)]} \right\} \\ &= P(\boldsymbol{\alpha}_l|A_q) \times \frac{\partial P(\alpha_k = 1|A_q)}{\partial \lambda_{0k}} \times \left\{ \frac{\alpha_{lk} - P(\alpha_k = 1|A_q)}{P(\alpha_k = 1|A_q) \times [1 - P(\alpha_k = 1|A_q)]} \right\} \\ &= P(\boldsymbol{\alpha}_l|A_q) \times (-\lambda_1) \times [\alpha_{lk} - P(\alpha_k = 1|A_q)]. \end{aligned}$$

9 By substituting and interchanging the summations,

$$\begin{aligned} \frac{\partial l(\mathbf{X})}{\partial \lambda_{0k}} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q \sum_{L=1}^L W(A_q) \times L(\mathbf{X}_i|\boldsymbol{\alpha}_l) \times P(\boldsymbol{\alpha}_l|A_q) \times (-\lambda_1) \times [\alpha_{lk} - P(\alpha_k = 1|A_q)] \\ &= \sum_{q=1}^Q -\lambda_1 \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1|A_q)] \sum_{i=1}^N \frac{L(\mathbf{X}_i|\boldsymbol{\alpha}_l) \times P(\boldsymbol{\alpha}_l|A_q) \times W(A_q)}{L(\mathbf{X}_i)} \\ &= \sum_{q=1}^Q -\lambda_1 \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1|A_q)] \sum_{i=1}^N P(\boldsymbol{\alpha}_l, A_q|\mathbf{X}_i) \\ &= \sum_{q=1}^Q -\lambda_1 \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1|A_q)] \times N_{lq} \\ &= \sum_{q=1}^Q -\lambda_1 \sum_{l=1}^L [R_{lkq} - P(\alpha_k = 1|A_q) \times N_{lq}] \\ &= \sum_{q=1}^Q -\lambda_1 \times \left[\sum_{l=1}^L R_{lkq} - \sum_{l=1}^L P(\alpha_k = 1|A_q) \times N_{lq} \right] \\ &= \sum_{q=1}^Q -\lambda_1 \times [R_{kq} - P(\alpha_k = 1|A_q) \times N_q] \\ &\propto -\sum_{q=1}^Q [R_{kq} - P(\alpha_k = 1|A_q) \times N_q], \end{aligned}$$

1 where $P(\boldsymbol{\alpha}_l, A_q | \mathbf{X}_i)$ is the joint posterior probability of $\boldsymbol{\alpha}_l$ and A_q for examinee i , N_{lq} is the expected
2 number of examinees with $\boldsymbol{\alpha}_l$ and A_q , R_{lkq} is the expected number of examinees with $\boldsymbol{\alpha}_l$ and A_q who
3 have mastered attribute k , N_q and R_{kq} are the expected number examinees with A_q and expected
4 number of examinees with A_q who have mastered α_k , respectively.

5
6 Find $\hat{\lambda}_{0k}$ by solving

$$\sum_{q=1}^Q [R_{kq} - P(\alpha_k = 1 | A_q) \times N_q] = 0.$$

7 **Case III:** $\beta_p = \lambda_1$

8
9 As in Case II, the derivative in $\partial l(\mathbf{X}) / \partial \beta_p$ also simplifies to

$$\frac{\partial L(\mathbf{X}_i | \boldsymbol{\alpha}_l) P(\boldsymbol{\alpha}_l | A_q)}{\partial \lambda_1} = L(\mathbf{X}_i | \boldsymbol{\alpha}_l) \frac{\partial P(\boldsymbol{\alpha}_l | A_q)}{\partial \lambda_1}.$$

9 However, it should be noted that α_1 appears across all $P(\alpha_k = 1 | A_q)$. For this reason, it might be
10 instructive to go over how the derivative is obtained.

11

$$\begin{aligned} \frac{\partial P(\boldsymbol{\alpha}_l | A_q)}{\partial \lambda_1} &= \sum_{k=1}^K \prod_{k' \neq k} P(\alpha_{k'} = 1 | A_q)^{\alpha_{lk'}} [1 - P(\alpha_{k'} = 1 | A_q)]^{1 - \alpha_{lk'}} \times \\ &\quad \frac{\partial P(\alpha_k = 1 | A_q)^{\alpha_{lk}} [1 - P(\alpha_k = 1 | A_q)]^{1 - \alpha_{lk}}}{\partial \lambda_1}, \end{aligned}$$

12 and the derivative on the right-hand side is equal to

$$\begin{aligned} &P(\alpha_k = 1 | A_q)^{\alpha_{lk}} [1 - P(\alpha_k = 1 | A_q)]^{1 - \alpha_{lk}} \frac{\partial P(\alpha_k = 1 | A_q)}{\partial \lambda_1} \times \\ &\quad \frac{\alpha_{lk} - P(\alpha_k = 1 | A_q)}{P(\alpha_k = 1 | A_q) [1 - P(\alpha_k = 1 | A_q)]}. \end{aligned}$$

13 This allows us to write

$$\begin{aligned} \frac{\partial P(\boldsymbol{\alpha}_l | A_q)}{\partial \lambda_1} &= \sum_{k=1}^K \left\{ \prod_{k=1}^K P(\alpha_k = 1 | A_q)^{\alpha_{lk}} [1 - P(\alpha_k = 1 | A_q)]^{1 - \alpha_{lk}} \right\} \\ &\quad \times \frac{\partial P(\alpha_k = 1 | A_q)}{\partial \lambda_1} \times \left\{ \frac{\alpha_{lk} - P(\alpha_k = 1 | A_q)}{P(\alpha_k = 1 | A_q) [1 - P(\alpha_k = 1 | A_q)]} \right\} \\ &= \sum_{k=1}^K P(\boldsymbol{\alpha}_l | A_q) \times \frac{\partial P(\alpha_k = 1 | A_q)}{\partial \lambda_1} \times \left\{ \frac{\alpha_{lk} - P(\alpha_k = 1 | A_q)}{P(\alpha_k = 1 | A_q) [1 - P(\alpha_k = 1 | A_q)]} \right\} \\ &= \sum_{k=1}^K P(\boldsymbol{\alpha}_l | A_q) \times (A_q - \lambda_{0k}) \times [\alpha_{lk} - P(\alpha_k = 1 | A_q)], \end{aligned}$$

14 because $\partial P(\alpha_k = 1 | A_q) / \partial \lambda_1 = (A_q - \lambda_{0k}) \times P(\alpha_k = 1 | A_q) \times [1 - P(\alpha_k = 1 | A_q)]$.

15

1 Again, by substituting and interchanging the summations,

$$\begin{aligned}
\frac{\partial l(\mathbf{X})}{\partial \lambda_1} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{q=1}^Q \sum_{L=1}^L W(A_q) L(\mathbf{X}_i | \boldsymbol{\alpha}_L) \times \\
&\quad \times \sum_{k=1}^K P(\boldsymbol{\alpha}_L | A_q) \times (A_q - \lambda_{0k}) \times [\alpha_{Lk} - P(\alpha_k = 1 | A_q)] \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1 | A_q)] \sum_{i=1}^N \frac{L(\mathbf{X}_i | \boldsymbol{\alpha}_L) \times P(\boldsymbol{\alpha}_L | A_q) \times W(A_q)}{L(\mathbf{X}_i)} \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1 | A_q)] \sum_{i=1}^N P(\boldsymbol{\alpha}_L, A_q | \mathbf{X}_i) \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \sum_{l=1}^L [\alpha_{lk} - P(\alpha_k = 1 | A_q)] \times N_{lq} \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \sum_{l=1}^L [R_{lkq} - P(\alpha_k = 1 | A_q) \times N_{lq}] \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \times \left[\sum_{l=1}^L R_{lkq} - \sum_{l=1}^L P(\alpha_k = 1 | A_q) \times N_{lq} \right] \\
&= \sum_{q=1}^Q \sum_{k=1}^K (A_q - \lambda_{0k}) \times [R_{kq} - P(\alpha_k = 1 | A_q) \times N_q].
\end{aligned}$$

2 Find $\hat{\lambda}_1$ by finding the zero of the above equation.

3
4 As previously, to implement the algorithm, Step 1 starts with initial values for \mathbf{g} , \mathbf{s} , and $\boldsymbol{\lambda}$. Step 2
5 involves computing $N_{jl}^{(0)}$, $R_{jl}^{(0)}$, $N_{jl}^{(1)}$ and $R_{jl}^{(1)}$, and finding $\hat{\mathbf{g}}$ and $\hat{\mathbf{s}}$. Step 3 involves computing N_q and
6 R_{kq} , and finding $\hat{\boldsymbol{\lambda}}$. Steps 2 and 3 are repeated until convergence.

7
8 We can easily derive the standard errors of the item parameter estimates using the formula for the
9 information matrix in the previous lecture. We will derive the standard errors of the higher-order
10 parameters numerically.

11
12 *Implementation of EM Algorithm for the HO-DINA model in Ox*

13
14 HODINA_MMLE.ox is an implementation of EM algorithm for the HO-DINA model. As with other pro-
15 grams for CDMs, this program will require the file names of the data and the Q-matrix.

16
17 We start by including a few headers and importing the code `solvenle`.

```

#include <oxstd.h>
#include <oxfloat.h>
#import <solvenle>

```

1 The first set of global variables allows users to easily modify the program to suit their specific needs.

```
decl N=2000,J=30,K=5;
decl data={"X_hodina.dat"};
decl qmatrix={"Q.txt"};
decl crit=0.001,maxit=999;
```

2 The next set of global variables include variables that are used generally with the DINA model, and
3 specifically with its higher-order formulation. Examples of latter are A, W, **nnodes**, **Rkq**, **Nq**, **la**, **la1**,
4 and 10. The HO-specific variable **prior_q** is an $2^K \times 1$ matrix with $P(\alpha_l | A_q)$ as entries.

```
decl Q,X;
decl A,W,nnodes;
decl Rkq,Nq;
decl la,la1,la0;
decl like,post,alpha,g,s,se,lse;
decl r0,r1,n0,n1;
decl lprior,eta;
decl it,maxabs;
decl alpha1,m1;
decl prior_q,theta;
```

5 The functions **pattern** and **alphaeta** are the same functions as before.

6
7 The function **findpriorq** will return a $Q \times 2^K$ matrix where each row represents the probability (like-
8 lihood) of each attribute pattern α_l given a particular A_q .

```
findpriorq(const l1, const l0){//QxL
decl tmp=exp(l1*(A-l0'));
tmp=tmp./(1+tmp);
tmp=log(tmp)*alpha'+log(1-tmp)*(1-alpha');
return(exp(tmp));
} //end findpriorq
```

```
findprior(){//Lx1
decl tmp=findpriorq(la1,la0).*(W*ones(1,2^K));
return(sumc(tmp)');
} //end findprior
```

- 1 The function `findprior` calls `findpriorq`, and summarizes the probability of each α_l by finding the
2 weighted sum of the likelihoods across the Q quadrature nodes.
3
4 The function `findlike` computes the entries of an $N \times 2^K$ matrix where each row represent one
5 examinee, and the column is the likelihood of a response vector for a fixed α_l .

```
findlike(){//NxL
for(decl i=0;i<N;i++){
decl XX=ones(2^K,1)*X[i][];
decl tmp=eta.*(ones(2^K,1)*(1-s)')+(1-eta).*(ones(2^K,1)*g');
tmp=sumr(XX.*log(tmp)+(1-XX).*log(1-tmp));
like[i][]=exp(tmp)';
}}//end i, findlike
```

- 6 The function `findpost` is similar to the previous functions for finding the posterior attribute distri-
7 bution of each examinee except for how the prior (this time involving θ) is computed. It updates the
8 $N \times 2^K$ variable `post`.

```
findpost(){//NxL
for(decl i=0;i<N;i++){
decl XX=ones(2^K,1)*X[i][];
decl tmp=eta.*(ones(2^K,1)*(1-s)')+(1-eta).*(ones(2^K,1)*g');
tmp=sumr(XX.*log(tmp)+(1-XX).*log(1-tmp))+lprior;
tmp=exp(tmp);
post[i][]=(tmp./sumc(tmp))';
}}//end i, findpost
```

- 9 The function `HO_DINA` contains the system of equations that needs to be solved for λ .

```
HO_DINA(const avF, const lam){//d is always vertical
decl p=lam[0]*(A-lam[1:]'); //Q*K
p=exp(p)./(1+exp(p));
decl tmp0=sumc(W*ones(1,K).*(Rkq-N*p))'; //(1*K)'
decl tmp=-sumc(W.*sumr((A-lam[1:]').*(Rkq-Nq*ones(1,K).*p)));
avF[0]= tmp|tmp0;
return 1;
}//end HO_DINA
```

- 10 The function `initiate` gives starting values to previously declared variables.
11
12 The function `est_itepar` finds g_j and s_j . It requires the argument `j`, and will be called from `iterate`.


```

est_itempar(const jj){
decl tmp=Q[jj][]*alpha';
tmp=tmp.==(Q[jj][]*ones(K,2^K));
tmp=sumr(post.*(ones(N,1)*tmp));
n0=N-sumc(tmp);
r0=sumc(X[][jj].*(1-tmp));
n1=sumc(tmp);
r1=sumc(X[][jj].*tmp);
g[jj]=r0/n0;
s[jj]=(n1-r1)/n1;
} //est_itempar

```

- 1 The function `est_hopar` estimate λ , and is also called from `iterate`. After computing N_q and R_{kq} ,
- 2 which are obtained from the first derivatives, `SolveNLE` is called to solved the system of equations.

```

est_hopar(){
like=zeros(N,2^K);
decl post_q;
decl Rlkq=<>, Ilq=<>;
findlike();
prior_q=findpriorq(la1,la0); //p(a1|Aq)=QxL
decl isum=like*prior_q'*W; //NxL x LxQ x Qx1
for(decl q=0;q<nnodes;q++){
post_q=like.*(ones(N,1)*prior_q[q][]); //NxL
post_q=post_q./(isum*ones(1,2^K));
Ilq=Ilq~sumc(post_q)';
Rlkq=Rlkq|sumc(alpha.*sumc(post_q)'.*ones(1,K));
} //end q
Nq=sumc(Ilq)';
Rkq=Rlkq;
la=la1|la0;
decl stat=SolveNLE(HO_DINA, &la);
la1=la[0];
la0=la[1:];
} //est_hopar

```

- 3 In addition to estimating the item parameters, the `iterate` function also estimates the higher-order
- 4 parameters. And as before, it also updates the prior, posteriors, and iteration counter.
- 5
- 6 It should be noted that the stopping rule is based on the improvement on the item parameters only,
- 7 and does not include the higher-order parameters.
- 8
- 9 The function `SE_item` will find the standard errors of g_j and s_j . It does so analytically using the
- 10 cross-products of the first derivatives. The function will be called from `stderr`.

```

SE_item(const jj){
decl tmp,tmp00,tmp11,tmp01;
tmp=Q[jj][]*alpha';
tmp=tmp.==(Q[jj][]*ones(K,2^K));
tmp=sumr(post.*(ones(N,1)*tmp));
tmp00=sumc(((X[] [jj]-g[jj]*ones(N,1)).*(1-tmp)).^2)/((g[jj]*(1-g[jj])).^2);
tmp11=sumc(((X[] [jj]-(1-s[jj])*ones(N,1)).*tmp).^2)/((s[jj]*(1-s[jj])).^2);
tmp01=-sumc(((X[] [jj]-g[jj]*ones(N,1)).*(1-tmp)).*((X[] [jj]-(1-s[jj])*ones(N,1)).*tmp)
/(g[jj]*(1-g[jj])*s[jj]*(1-s[jj])));
se=se|sqrt(diagonal(invert((tmp00~tmp01)|(tmp01~tmp11))));
} //end SE_item

```

- 1 The function `SE_HO_N` is specifically designed to compute the the marginalized loglikelihood with 11
- 2 and 10 as arguments. The arguments are necessary to afford the function flexibility in computing the
- 3 marginalized loglikelihoods around $\hat{\lambda}$, which are need for determining the second derivative.

```

SE_HO_N(const l1,const l0){
decl tmp=0;
for(decl i=0;i<1;i++){
decl XX=ones(2^K,1)*X[i][];
decl tmp1=eta.*(ones(2^K,1)*(1-s)')+(1-eta).*(ones(2^K,1)*g');
tmp1=ones(nnodes,1)*exp(sumr(XX.*log(tmp1)+(1-XX).*log(1-tmp1)))';
tmp1=tmp1.*findpriorq(l1,l0);
tmp1=tmp1.*(W*ones(1,2^K));
tmp=tmp+sumc(sumr(log(tmp1)));
} //end i
return(tmp);
} //end SE_HO

```

- 4 The function `stderr` will compute the standard errors of \hat{g} and \hat{s} , and then $\hat{\lambda}$.
- 5
- 6 Computing the standard errors of \hat{g}_j and \hat{s}_j is done one item at a time, and require simply calling the
- 7 function `SE_item` with `j` as the argument.
- 8
- 9 Computing the standard errors of $\hat{\lambda}$ is done numerically. In addition to $\hat{\lambda}$, the marginalized loglikeli-
- 10 hood needs to also be evaluated at $\hat{\lambda} \pm h$. The evaluations for $\hat{\lambda}_1$ and $\hat{\lambda}_0$ are carried out separately.
- 11
- 12 Note that only the diagonal elements of the information matrix are computed. Most likely, the result-
- 13 ing standard errors are larger than the actual values.
- 14
- 15 Write a code that evaluates all the elements of the information matrix. Compare the standard errors
- 16 computed from this matrix to those computed using the diagonal elements only.

```

stderr(){
se=<>;
lse=<>;
findpost();
for(decl j=0;j<J;j++){// Analytical SE of item parameter estimates
SE_item(j);
}// end j
// Numerical SE of HO parameter estimates; diagonals only
decl h=.0001;
decl lamh=la1-h;
decl laph=la1+h;
decl tmp=(SE_HO_N(laph,la0)+SE_HO_N(lamh,la0)-2*SE_HO_N(la1,la0))/h^2;
lse=lse|tmp;
for(decl k=0;k<K;k++){
laph=la0;
laph[k]=laph[k]+h;
lamh=la0;
lamh[k]=lamh[k]-h;
tmp=(SE_HO_N(la1,laph)+SE_HO_N(la1,lamh)-2*SE_HO_N(la1,la0))/h^2;
lse=lse|tmp;
}//end k
lse=sqrt((-1)./lse);
}//end stderr

```

```

find_theta(){
prior_q=findpriorq(la1,la0); //p(a1|Aq)=QxL
findlike(); //NxL
decl tmp2=zeros(N,nnodes);
for (decl l=0;l<2^K;l++){
decl tmp1=ones(N,1)*(prior_q[l].*W)'; //Nx1 x 1xQ
tmp1=(like[l].*ones(1,nnodes)).*tmp1;
tmp2=tmp2+tmp1;
}//end l
tmp2=tmp2./(sumr(tmp2)*ones(1,nnodes)); //normalizing
theta=sumr(tmp2.*(ones(N,1)*A'));//expected value
}//end find_theta

```

- 1 The `find_theta` function computes the $E(\theta|\mathbf{X}_i)$. The function is based on the following result:

$$\begin{aligned}
P(\theta|\mathbf{X}_i, \boldsymbol{\alpha}, \boldsymbol{\lambda}) &\propto P(\mathbf{X}_i|\theta, \boldsymbol{\alpha}, \boldsymbol{\lambda}) \times P(\theta|\boldsymbol{\alpha}, \boldsymbol{\lambda}) \\
&\propto L(\mathbf{X}_i|\boldsymbol{\alpha}) \times P(\boldsymbol{\alpha}|\theta, \boldsymbol{\lambda}) \times P(\theta|\boldsymbol{\lambda}) \\
&= L(\mathbf{X}_i|\boldsymbol{\alpha}) \times P(\boldsymbol{\alpha}|\theta, \boldsymbol{\lambda}) \times P(\theta).
\end{aligned}$$

- 2 The function `printout` prints two examinee-related output files: `alpha.out` and `theta.out`. The

1 former is the attribute classification, whereas the later the expected a posteriori of θ .

2
3 In addition, the function provides the following screen output: item parameter and higher-order pa-
4 rameter estimates and the associated standard errors, and the joint and marginal probabilities of the
5 attributes.

6
7 Finally, the `main` function reads the data and Q-matrix, and calls the functions `initiate`, `iterate`,
8 `stderr`, and `printout`.

9
10 1. Using `HODINA_gen.ox`, generate data for different values of λ , q , and s , and N .

11
12 2. How are the different estimates affected by the different generating values?

13 14 An E-M Algorithms for the LLM

15
16 Recall that the LLM (linear logistic model) is just the additive model under the logit link. Its item
17 response function is given by

$$P(X_j = 1 | \alpha_{jl}^*) = P_{\alpha_{jl}^*} = \frac{\exp\left(\lambda_{j0} + \sum_{k=1}^{K_j} \lambda_{jk} \alpha_{lk}\right)}{1 + \exp\left(\lambda_{j0} + \sum_{k=1}^{K_j} \lambda_{jk} \alpha_{lk}\right)},$$

18 or

$$\text{logit}(P_{\alpha_{jl}^*}) = \lambda_{j0} + \sum_{k=1}^{K_j} \lambda_{jk} \alpha_{lk}.$$

19 As can be seen from the model, the number of LLM parameters for item j is equal to $K_j + 1$.

20
To obtain the MLE of λ_j , let's again start with the marginalized likelihood of the response vector of
examinee i :

$$L(\mathbf{X}_i) = \sum_{l=1}^L L(\mathbf{X}_i | \alpha_l) p(\alpha_l).$$

where, as usual, $p(\alpha_l)$ is the prior probability of $p(\alpha_l)$, and for dichotomous response,

$$L(\mathbf{X}_i | \alpha_l) = \prod_{j=1}^J P_{\alpha_{jl}^*}^{X_{ij}} \times (1 - P_{\alpha_{jl}^*})^{(1-X_{ij})}.$$

Finding $\hat{\lambda}_j$ requires maximizing

$$l(\mathbf{X}) = \log \prod_{i=1}^N L(\mathbf{X}_i) = \sum_{i=1}^N \log L(\mathbf{X}_i)$$

21 with respect to λ_j .

22

23 Taking the derivative,

$$\begin{aligned}
\frac{\partial l(\mathbf{X})}{\partial \lambda_j} &= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \times \frac{\partial L(\mathbf{X}_i)}{\partial \lambda_j} \\
&= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \times \sum_{l=1}^L p(\boldsymbol{\alpha}_l) \times \frac{\partial L(\mathbf{X}_i|\boldsymbol{\alpha}_l)}{\partial \lambda_j} \\
&= \sum_{i=1}^N \frac{1}{L(\mathbf{X}_i)} \sum_{l=1}^L p(\boldsymbol{\alpha}_l) \times L(\mathbf{X}_i|\boldsymbol{\alpha}_l) \times \left[\frac{X_{ij} - P_{\boldsymbol{\alpha}_{jl}^*}}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \times \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j} \\
&= \sum_{i=1}^N \sum_{l=1}^L P(\boldsymbol{\alpha}_l|\mathbf{X}_i) \times \left[\frac{X_{ij} - P_{\boldsymbol{\alpha}_{jl}^*}}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \times \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j}
\end{aligned}$$

1 where $P(\boldsymbol{\alpha}_l|\mathbf{X}_i)$ is the posterior probability that examinee i has the attribute pattern $\boldsymbol{\alpha}_l$.

2

3 As in the previous lecture, we can group the latent classes, as in, $\boldsymbol{\alpha}_l \rightarrow \boldsymbol{\alpha}_{jl}^*$, to get

$$P(\boldsymbol{\alpha}_l|\mathbf{X}_i) = \sum_{\forall l: \boldsymbol{\alpha}_l \rightarrow \boldsymbol{\alpha}_{jl}^*} P(\boldsymbol{\alpha}_l|\mathbf{X}_i),$$

4 which leads to the following simplification:

$$\frac{\partial l(\mathbf{X})}{\partial \lambda_j} = \sum_{jl=1}^{L_j} P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i) \left[\frac{1}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \sum_{i=1}^N (X_{ij} - P_{\boldsymbol{\alpha}_{jl}^*}) \times \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j},$$

5 where the outer sum involves 2^{K_j} , instead of 2^K terms.

6

7 Again, by interchanging the summations, and making some rearrangements, we get

$$\begin{aligned}
\frac{\partial l(\mathbf{X})}{\partial \lambda_j} &= \sum_{jl=1}^{L_j} \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \sum_{i=1}^N (X_{ij} - P_{\boldsymbol{\alpha}_{jl}^*}) \times P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i) \\
&= \sum_{jl=1}^{L_j} \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \times \left[\sum_{i=1}^N X_{ij} \times P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i) - P_{\boldsymbol{\alpha}_{jl}^*} \times \sum_{i=1}^N P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i) \right] \\
&= \sum_{jl=1}^{L_j} \frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \lambda_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{jl}^*}(1 - P_{\boldsymbol{\alpha}_{jl}^*})} \right] \times \left[R_{jl}^* - P_{\boldsymbol{\alpha}_{jl}^*} \times N_{jl}^* \right],
\end{aligned}$$

8 where $R_{jl}^* = \sum_{i=1}^N X_{ij} \times P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i)$ and $N_{jl}^* = \sum_{i=1}^N P(\boldsymbol{\alpha}_{jl}^*|\mathbf{X}_i)$.

9

10 By expanding the summation, the derivative can be written as

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \boldsymbol{\lambda}_j} &= \frac{\partial P_{\boldsymbol{\alpha}_{j1}^*}}{\partial \boldsymbol{\lambda}_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{j1}^*}(1 - P_{\boldsymbol{\alpha}_{j1}^*})} \right] \times [R_{j1}^* - P_{\boldsymbol{\alpha}_{j1}^*} \times N_{j1}^*] + \\ &\quad \frac{\partial P_{\boldsymbol{\alpha}_{j2}^*}}{\partial \boldsymbol{\lambda}_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{j2}^*}(1 - P_{\boldsymbol{\alpha}_{j2}^*})} \right] \times [R_{j2}^* - P_{\boldsymbol{\alpha}_{j2}^*} \times N_{j2}^*] + \dots + \\ &\quad \frac{\partial P_{\boldsymbol{\alpha}_{jL_j}^*}}{\partial \boldsymbol{\lambda}_j} \times \left[\frac{1}{P_{\boldsymbol{\alpha}_{jL_j}^*}(1 - P_{\boldsymbol{\alpha}_{jL_j}^*})} \right] \times [R_{jL_j}^* - P_{\boldsymbol{\alpha}_{jL_j}^*} \times N_{jL_j}^*].\end{aligned}$$

Because

$$\frac{\partial P_{\boldsymbol{\alpha}_{jl}^*}}{\partial \boldsymbol{\lambda}_j} = P_{\boldsymbol{\alpha}_{jl}^*} \times (1 - P_{\boldsymbol{\alpha}_{jl}^*}) \times \frac{\partial \text{logit}(P_{\boldsymbol{\alpha}_{jl}^*})}{\partial \boldsymbol{\lambda}_j},$$

1 and by denoting $R_{jl}^* - P_{\boldsymbol{\alpha}_{jl}^*} \times N_{jl}^*$ as w_{jl} and $\text{logit}(P_{\boldsymbol{\alpha}_{jl}^*})$ as ℓ_{jl} , we can write $\partial l(\mathbf{X})/\partial \boldsymbol{\lambda}_j$ in matrix
2 notation, as in,

$$\begin{aligned}\frac{\partial l(\mathbf{X})}{\partial \boldsymbol{\lambda}_j} &= \left[\frac{\partial \ell_{jl}}{\partial \boldsymbol{\lambda}_j} \right]_{L_j \times (K_j+1)}' \times \mathbf{w}_j \\ &= \begin{bmatrix} \frac{\partial \ell_1}{\partial \lambda_{j0}} & \frac{\partial \ell_1}{\partial \lambda_{j1}} & \dots & \frac{\partial \ell_1}{\partial \lambda_{jK_j}} \\ \frac{\partial \ell_2}{\partial \lambda_{j0}} & \dots & \dots & \vdots \\ \vdots & \dots & \dots & \vdots \\ \frac{\partial \ell_{L_j}}{\partial \lambda_{j0}} & \dots & \dots & \frac{\partial \ell_{L_j}}{\partial \lambda_{jK_j}} \end{bmatrix}' \times \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{L_j} \end{bmatrix}.\end{aligned}$$

The logit of the item response function of LLM can also be written in matrix notation, as in,

$$\begin{bmatrix} \ell_1 \\ \ell_2 \\ \vdots \\ \ell_{L_j} \end{bmatrix} = [\mathbf{1} \quad \mathbf{A}_j^*] \times \begin{bmatrix} \lambda_{j0} \\ \lambda_{j1} \\ \vdots \\ \lambda_{jL_j} \end{bmatrix},$$

3 where $\mathbf{1}$ is an $L_j \times 1$ matrix of 1s, and \mathbf{A}_j^* is $2^{L_j} \times L_j$ matrix of stacked $\boldsymbol{\alpha}_{jl}^*$. $[\mathbf{1} \quad \mathbf{A}_j^*]$ is also called our
4 *design* matrix.

5

For example, if two attributes are required for item j , we have,

$$\begin{bmatrix} \ell_1 \\ \ell_2 \\ \ell_3 \\ \ell_4 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} \lambda_{j0} \\ \lambda_{j1} \\ \lambda_{j2} \end{bmatrix}.$$

It should be easy to extract what each row means. For instance, for row 2, which represents P_{10} ,

$$\ell_2 = \text{logit}(P_{10}) = \lambda_{j0} + \lambda_{j1}.$$

6 Note that the derivative of ℓ_2 with respect to λ_0 , λ_1 , and λ_2 is simply $[1 \quad 1 \quad 0]$, which corresponds to
7 the second row of the design matrix.

8

1 It can be easily shown that, in general,

$$\begin{bmatrix} \frac{\partial \ell_1}{\partial \lambda_{j0}} & \frac{\partial \ell_1}{\partial \lambda_{j1}} & \cdots & \frac{\partial \ell_1}{\partial \lambda_{jK_j}} \\ \frac{\partial \ell_2}{\partial \lambda_{j0}} & \cdots & \cdots & \vdots \\ \vdots & \cdots & \cdots & \vdots \\ \frac{\partial \ell_{L_j}}{\partial \lambda_{j0}} & \cdots & \cdots & \frac{\partial \ell_{L_j}}{\partial \lambda_{jK_j}} \end{bmatrix} = [\mathbf{1} \quad \mathbf{A}_j^*].$$

Thus,

$$\frac{\partial l(\mathbf{X})}{\partial \boldsymbol{\lambda}_j} = [\mathbf{1} \quad \mathbf{A}_j^*]' \times \mathbf{w}_j.$$

2 We get $\hat{\boldsymbol{\lambda}}_j$ by solving for $\boldsymbol{\lambda}_j$ in

$$[\mathbf{1} \quad \mathbf{A}_j^*]' \times \mathbf{w}_j = \mathbf{0}.$$

3 If we use `solveNLE` in Ox, we do not need the second derivatives to get the MLEs.

4
5 Finding the MLEs and the corresponding standard errors are two distinct steps. We will consider our
6 options for the latter later.

7
8 *Implementation of EM Algorithm for the LLM in Ox (I)*

9
10 The program `LLM_NLE.ox` will estimate the item parameters of the LLM by solving the zero of the first
11 derivative of the marginalized loglikelihood. It requires both item responses and the Q-matrix as input.

12
13 This program includes the following headers:

```
#include <oxstd.h>
#include <oxprob.h>
#include <oxfloat.h>
#import <maxsqp>
#import <solveNLE>
```

```
decl N=1000, J=30,K=5;
decl X,Q,Qj;
decl Kj,Kjmax;
decl Ng,Rg,Ngj,Rgj;
decl alpha,eta,est,post,coef;
decl prior,lprior;
decl lambda,SE,jj,fL,Prob;
```

14 We first declare a few basic global variables. Some of which are the same as the variables that we have
15 defined previously.

16

1 In the program we declare additional global variables: **Ng** and **Rg** are the matrix for the number of
2 examinees in each latent group and the number of examinees in each latent group answering item
3 correctly; **Ngj** and **Rgj** are the j^{th} row of **Ng** and **Rg**; and **coef** is a matrix for item parameters, which
4 in our case is λ_j .
5
6 The last set of variables will be used in conjunction with the objective function, which has a specific
7 form. The objective function will be needed for the numerical second derivatives.
8
9 The functions **pattern**, **alphaeta** and **findpost** are the same functions as before.
10
11 The function **NgRg** will calculate **Ng** and **Rg** for item j .
12

```
NgRg(){
for (decl j=0;j<J;j++){ //j loop: for each item
decl Kj=sumr(Q[j] []);
decl Qj=pattern(Kj);
decl Lj=vecindex(Q[j] [] ==1);
decl Rjtmp=<>;
decl Ijtmp=<>;
//re-group the 2^K groups in to 2^Kj groups
for(decl kj=0;kj<2^Kj;kj++){
decl tmp_loc=vecindex(prodr(ones(2^K,1)*Qj[kj] [] ==alpha[] [Lj]) ==1);
decl tmp=sumr(post[] [tmp_loc]);
Ijtmp=Ijtmp|sumc(tmp); //2^Kj col-vector
Rjtmp=Rjtmp|sumc(X[] [j].*tmp); //2^Kj col-vector
} //end kj
Ng[j] [0:2^Kj-1]=Ijtmp'; // # of examinees expected in latent class
Rg[j] [0:2^Kj-1]=Rjtmp'; // # of examinees expected to answer correctly
} //end j
} //end NgRg
```

13 To solve the system of first derivatives, we will use the **SolveNLE** function in Ox, rather than writing
14 our own zero-finder function.
15
16 The function **LLM_NLE**, which will be called in **SolveNLE**, can be used for this purpose. In this code,
17 **M** is the design matrix $[\mathbf{1} \quad \alpha_j^*]$, **Qj** is the reduced attribute pattern α_{jl} obtained from the function
18 **pattern** with the argument of **Kj**, and **w** is w_{jl} .
19


```

LLM_NLE(const avF, const d){
decl M=ones(2^Kj,1)~Qj;
decl P=M*d;
P=exp(P)./(1+exp(P));
decl w=(Rgj-Ngj.*P);
avF[0]=M'*w;
return 1;
} //end LLM_NLE

```

- 1 As with the 2PL, fullLike provides an $N \times 2^K$ matrix where the l^{th} element in row i is $L(\mathbf{X}_i|\boldsymbol{\alpha}_l)$.

```

fullLike(){//returns N x 2^K
fL=<>;
Prob=zeros(2^K,J);
for(decl j=0;j<J;j++){
decl tmp=vecindex(Q[j] [] .==1);
decl astar=(1~alpha[] [tmp])*coef[j] [0:sumr(Q[j] [])]';
Prob[] [j]=exp(astar)./(1+exp(astar));
} //end j
for(decl i=0;i<N;i++){
decl Xi=ones(2^K,1)*X[i] [];
decl tmp=sumr(Xi.*log(Prob)+(1-Xi).*log(1-Prob));
fL=fL|exp(tmp');
} //end i, fullLike

```

- 2 LLM_obj_func is the objective function that we will call in finding the standard errors *numerically*.
3 Note that this function uses the discounted likelihood to carry out the computations more efficiently.

```

LLM_obj_func(const vP, const adFunc, const avScore, const amHessian){
decl Prob0=ones(N,1)*Prob[] [jj]'; //N x 2^K
decl Xj=X[] [jj]*ones(1,2^K); //N x 2^K
decl lj0=(Prob0.^Xj).*(1-Prob0).^(1-Xj);
decl fL_j=fL./lj0;
decl tmp=vecindex(Q[jj] [] .==1);
decl astar=(1~alpha[] [tmp])*vP;
decl Prob=exp(astar)./(1+exp(astar)); //2^K x 1
Prob=ones(N,1)*Prob';
adFunc[0]=sumc(log(sumr(fL_j.*(Prob.^Xj).*((1-Prob).^(1-Xj)).*(ones(N,1)*prior'))));
return 1;
} //end LLM_obj_func

```

1 The function `findSE` is used to find the standard errors numerically once the item parameters have
2 been estimated.

3

```
findSE(){
fullLike();
decl mhess;
for (jj=0;jj<J;jj++){
decl vP=coef[jj][0:sumr(Q[jj][])];
if (Num2Derivative(LLM_obj_func, vP, &mhess)){
SE[jj][0:sumr(Q[jj][])]=sqrt(diagonal(-invert(mhess)));
}}//end if, j, findSE
```

4 The function `initiate` initializes the global variables. In addition to assigning specific values, it also
5 called the functions `pattern`, `alphaeta`, `NgRg` and `findpost` functions. The initial values for the
6 probabilities, λ , and the prior distribution are also specified here.

```
initiate(){
alpha=pattern(K);
alphaeta();
Kjmax=max(sumr(Q));
est=-1*ones(J,2^(maxc(sumr(Q))));
coef=-1*ones(J,Kjmax+1);
Ng=zeros(J,2^Kjmax);
Rg=zeros(J,2^Kjmax);
for(decl j=0;j<J;j++){
Kj=sumr(Q[j]);
est[j][0:(2^Kj-1)]=.2;
est[j][2^Kj-1]=.8;
coef[j][0:Kj]=0.3;
}// end j
lprior=ones(2^K,1)./(2^K);
lprior=log(lprior/sumc(lprior));
post=ones(N,2^K);
findpost();
NgRg();
}//end initiate
```

7 The function `iterate` is where the item parameters are estimated. The `while` loop is in effect as
8 long as the difference in parameter estimates between two iterations remain large AND the maximum
9 number of iterations has not been reached.

10

11 Within the `while` loop, `est`, the success probabilities of the different latent groups for each item, will
12 be calculated after finishing the estimation of λ s.

13

- 1 The posterior distributions, E-step, iteration number, and maximum item difference between two
- 2 consecutive iterations will be updated after each iteration.

```

iterate(){
decl crit=0.001;
decl maxit=999;
decl it=0;
decl diff=<>;
decl maxabs=1;
decl f;
while (maxabs>crit&&it<maxit){
diff=coef;
for (decl j=0;j<J;j++) {
Kj=sumr(Q[j][]);
Qj=pattern(Kj);
Rgj=Rg[j][0:2^Kj-1]';
Ngj=Ng[j][0:2^Kj-1]';
decl coefj=coef[j][0:Kj]';
SolveNLE(LLM_NLE, &coefj);
coef[j][0:Kj]=coefj';
lambda=(ones(2^Kj,1)~Qj)*coef[j][0:Kj]';
est[j][0:(2^Kj-1)]=(exp(lambda)./(1+exp(lambda)))';
} //end j
prior=meanc(post)';
lprior=log(prior/sumc(prior));
findpost();
NgRg();
it++;
maxabs=max(sqrt((diff-coef).^2));
} //end while
} //end iterate

```

- 3 Finally, the main function will be used to read responses and Q-matrix, call the *initiate* and *iterate*
- 4 function.
- 5
- 6 The main function also provides a few screen output:
- 7
- 8 The first set of output is $\hat{\lambda}_j$ from `coef`. As before, -1 means no value.
- 9
- 10 The second set of output is $SE(\hat{\lambda}_j)$ from `SE`. As before, 1 means no value.
- 11
- 12 The last set of output is $\hat{\lambda}_j \rightarrow P\alpha_{j_l}^*$ from `est`.

```

main(){
format(10000);
decl file0, tmp0;
file0=fopen("sample_data.dat","r");
fscan(file0,"%#M",N,J,&X);
fclose(file0);
file0=fopen("sample_Q.dat","r");
fscan(file0,"%#M",J,K,&Q);
fclose(file0);
initiate();
iterate();
print ("\nEstimates of Lambdas:",coef);
decl SE=findSE();
print("\nStandard Errors of Lambda Estimates:",SE);
print("\n Latent Group Success Probabilities:",est);
} //end main

```

1 *An Implementation of EM Algorithm for the LLM in Ox (II)*

2
3 Recall that given a provisional posterior distribution and the response vector \mathbf{X}_j , we can compute
4 $R_{jl}^* = \sum_{i=1}^N X_{ij} \times P(\boldsymbol{\alpha}_{jl}^* | \mathbf{X}_i)$ and $N_{jl}^* = \sum_{i=1}^N P(\boldsymbol{\alpha}_{jl}^* | \mathbf{X}_i)$.
5
6 We can use this information as another way of estimating $\boldsymbol{\lambda}_j$. We can to maximize the loglikelihood
7 associated with item j given R_{jl}^* and N_{jl}^* , and this loglikelihood is given by

$$\begin{aligned}
l(\mathbf{X}_j) &= \log \prod_{l=1}^{L_j} P_{\boldsymbol{\alpha}_{jl}^*}^{R_{jl}^*} (1 - P_{\boldsymbol{\alpha}_{jl}^*})^{N_{jl}^* - R_{jl}^*} \\
&= \sum_{l=1}^{L_j} \left[R_{jl}^* \times \log(P_{\boldsymbol{\alpha}_{jl}^*}) + (N_{jl}^* - R_{jl}^*) \times \log(1 - P_{\boldsymbol{\alpha}_{jl}^*}) \right].
\end{aligned}$$

8 In LLM_SQP.ox, the function associated with $l(\mathbf{X}_j)$ is given below:

```

LLM_obj_func0( const vP, const avF, const avS, const avM){
Qj=pattern(Kj);
decl M=ones(2^Kj,1)~Qj;
decl P=M*vP;
P=exp(P)./(1+exp(P));
avF[0]=sumc(Rgj.*log(P)+(Ngj-Rgj).*log(1-P));
return 1;
} // end LLM_obj_func0

```

1 This function differentiates LLM_SQP.ox from LLM_NLE.ox. Of course, we need to change how we obtain
 2 the estimate in `iterate`: from `SolveNLE` to `MaxSQPF`.

3
 4 Note that with `MaxSQPF`, we don't need to know the derivatives. Thus, this can be used with any CDM!

6 Additional Exercises

7
 8 1. Comparing the saturated DINA and HO-DINA models

9
 10 a) Generate DINA model data using attributes that have i) uniformly distributed structure, and ii)
 11 higher-order structure.

12
 13 b) Estimate the data using the estimation programs for the saturated DINA and HO-DINA models.

14
 15 c) Compare the item parameter estimates, attribute distribution, attribute prevalence, and attribute
 16 classification of the two models.

17
 18 2. Estimating the parameters of the *A*-CDM by maximizing the log-likelihood.

19
 20 What we need to do first is to modify the following code in `iterate` function in `LLM_SQP.ox`, to it to
 21 be applicable to the *A*-CDM.

```
est[j][0:(2^Kj-1)]=(exp(lambda)./(1+exp(lambda)))';
```

22 Note that the above code is based on the item response function of the LLM IRF. Recall, this item
 23 response function is

$$P(X_j = 1|\alpha_{jl}^*) = \frac{\exp\left(\lambda_{j0} + \sum_{k=1}^{K_j^*} \lambda_{jk} \alpha_{lk}\right)}{1 + \exp\left(\lambda_{j0} + \sum_{k=1}^{K_j^*} \lambda_{jk} \alpha_{lk}\right)}.$$

24 For the *A*-CDM, that above code should be replaced by:

```
est[j][0:(2^Kj-1)]=lambda';
```

25 Also, because of the formulation in the identity link, we need to add constraints to the parameters for
 26 the estimates to be acceptable. Specifically, we use the constraints:

$$0 \leq P\alpha_{jl}^* \leq 1, \quad jl = 1, \dots, L_j.$$

27 The code below can be used as the constraints, and needs to be called in `MaxSQPF`. Modify `MaxSQPF` to
 28 include these constraints.

29

```

cfunc(const avF, const vP){
decl M=ones(2^Kj,1)~Qj;
decl P=M*vP;
avF[0]=P/(1-P);
return 1;
} //end cfunc

```

Optional Exercises

1. Normality of θ Not Assumed (not sure this can be done)

To relax the normality assumption about θ , $W(A_q)$ can be estimated instead of being derived from $N(0, 1)$. To do this, we can find $W(A_q|\mathbf{X})$, the posterior weight of the quadrature node A_q given the data.

The posterior weight can be written as

$$\begin{aligned}
W(A_q|\mathbf{X}) &\propto p(\mathbf{X}|A_q) \times W(A_q) \\
&= \left[\sum_{l=1}^L p(\mathbf{X}, \boldsymbol{\alpha}_l | A_q) \right] \times W(A_q) \\
&= \left[\sum_{l=1}^L p(\mathbf{X} | \boldsymbol{\alpha}_l, A_q) \times p(\boldsymbol{\alpha}_l | A_q) \right] \times W(A_q) \\
&= \left[\sum_{l=1}^L p(\mathbf{X} | \boldsymbol{\alpha}_l) \times p(\boldsymbol{\alpha}_l | A_q) \right] \times W(A_q).
\end{aligned}$$

The last equality follows from the conditional independence of \mathbf{X} given $\boldsymbol{\alpha}_l$. The posterior weights need to be normalized to ensure that $\sum_{q=1}^Q W(A_q|\mathbf{X}) = 1$, and the mean and variance of the latent trait are always equal to 0 and 1, respectively.

This step can be viewed as Step 4 of the HO-IRT EM iteration process, and needs to be repeated together with Steps 2 and 3.

2. As mentioned before, the logit of the higher-order structure written in two ways.

$$\lambda_1(\theta - \lambda_{0k}) \quad \text{or} \quad \lambda_1 \times \theta + \lambda_{0k}^*.$$

The former is called the *discrimination-difficulty* form; the latter, the *slope-intercept* form. The latter has at least two advantages. First, it is computationally more stable. That is, even if the $\lambda_1 = 0$, we can still estimate λ_{0k}^* .

Why is not possible with the former formulation?

Second, it makes higher-order CDMs more general. Recall that when $\lambda_1 = 0$, we have the independence model, as in,

$$P(\boldsymbol{\alpha}_l | \theta) = \prod_{\forall k} P(\alpha_{lk} | \theta) = \prod_{\forall k} P(\alpha_{lk}).$$

1 However, $P(\alpha_{lk})$, the probability of mastering attribute k , can vary in size (i.e., some attributes are
2 more difficult to master than others.
3
4 A special case is when $\lambda_{0k^*} = 0$, $P(\alpha_{lk}) = 0.5$ for all k . This will produce a uniform attribute structure.
5
6 Note that, whenever $\lambda_1 = 0$, $\lambda_{0k^*}^* = 0$ is implied in the discrimination-difficulty formulation.
7
8 If we do not want $\lambda_1 = 0$ to imply $\lambda_{0k^*} = 0$, we need to estimate the higher-order parameters using
9 the slope-intercept form.
10
11 Write an EM algorithm for the higher-order DINA model that estimates λ_1 and λ_{0k^*} .
12

Markov Chain Monte Carlo - Part I

Introduction

Background

Statistical models are formal mathematical expressions of processes that provide description of the assumed structure of the observed data.

Notwithstanding the inability of simple conventional models to describe accurately the intricate structures of most real-life phenomena, these models remain attractive and in use because of their computational ease.

However in many situations, adequate understanding of these phenomena can only be facilitated by using more complex models.

In these settings, traditional methods of estimation (e.g., numerical evaluation, analytic approximation) cannot address the computational requisites inherent in these models.

Simulation-based analyses such as MCMC are particularly suited for **high-dimensional, complex problems**. The application of this method had been limited in the past by its computational requirements.

However, the advent of affordable and fast computing machines has rendered this impediment a non-issue, and has renewed the interest of more than a few researchers in MCMC. As a result, the breadth of its application has increased tremendously in recent years.

Methods of Analysis

A researcher's statistical persuasion - frequentist or Bayesian - determines the choice of techniques and inferences.

From a frequentist perspective, random quantities are comprised of the data, and perhaps, random effects for some applications.

From a Bayesian perspective, no distinction exists between the observable data and the unobservable parameters - both are considered random variables.

The frequentist inference is based on the likelihood function, whereas the Bayesian inference is based on the posterior distribution.

The mathematical computations in each perspective seek solutions to different problems. The frequentist problem is finding an explicit definition of estimators as solutions to maximization problems, whereas the Bayesian problem is finding an implicit representation of estimators as an integral.

Although MCMC lends itself naturally to Bayesian inference, it can also be applied to optimization problems.

A Bayesian approach to parameter estimation in IRT has been done in the past. Although the estimation procedures have a Bayesian flavor, they differ from the current line of Bayesian estimation in that numerical procedures and approximation are used instead of stochastic integration.

Let \mathbf{X} denote the response data, and $\boldsymbol{\zeta}$ denote the collection of parameters (e.g., ability and item parameters).

In a Bayesian context, inference about $\boldsymbol{\zeta}$ given \mathbf{X} requires setting up the joint distribution $p(\boldsymbol{\zeta}, \mathbf{X})$.

The joint distribution can be written as the product of the prior distribution $p(\boldsymbol{\zeta})$ and the sampling distribution or likelihood function $p(\mathbf{X}|\boldsymbol{\zeta})$:

$$p(\boldsymbol{\zeta}, \mathbf{X}) = p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta}).$$

Using Bayes theorem, the posterior distribution of $\boldsymbol{\zeta}$ given \mathbf{X} is

$$\begin{aligned} p(\boldsymbol{\zeta}|\mathbf{X}) &= \frac{p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta})}{\int p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta})d\boldsymbol{\zeta}} \\ &= \frac{p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta})}{p(\mathbf{X})}. \end{aligned}$$

For given data, \mathbf{X} is considered constant and $p(\mathbf{X})$ does not depend on $\boldsymbol{\zeta}$ so it can be omitted from the posterior. Hence, it can be written in an unnormalized form as

$$p(\boldsymbol{\zeta}|\mathbf{X}) \propto p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta}).$$

It is not uncommon in Bayesian applications for the posterior to be known only up a proportionality constant.

Inference about any the features of posterior (e.g., moments, quantiles, regions of highest posterior density) can be made. Specifically, finding the posterior expectations of a function of $\boldsymbol{\zeta}$ is of particular interest.

Denoting the function of $\boldsymbol{\zeta}$ as $f(\boldsymbol{\zeta})$, its posterior expectation is given by

$$\begin{aligned} E[f(\boldsymbol{\zeta})|\mathbf{X}] &= \int f(\boldsymbol{\zeta})p(\boldsymbol{\zeta}|\mathbf{X})d\boldsymbol{\zeta} \\ &= \frac{\int f(\boldsymbol{\zeta})p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta})d\boldsymbol{\zeta}}{\int p(\boldsymbol{\zeta})p(\mathbf{X}|\boldsymbol{\zeta})d\boldsymbol{\zeta}}. \end{aligned}$$

For high-dimensional, non-standard and complex models, analytical or numerical solutions or approximations of $E[f(\boldsymbol{\zeta})]$ are not typically available.

One possible solution to this problem without constraining the models (e.g., by using conjugate priors) is to perform Monte Carlo integration based on Markov chains that reduces the computational involvedness of $E[f(\boldsymbol{\zeta})]$ into a sequence of simpler calculations.

1 Monte Carlo Integration

2
3 For notational convenience and greater generality, denote the distribution of interest as $\pi(\zeta)$, which
4 in this case is the posterior distribution. The expectation problem can then be reexpressed as

$$E[f(\zeta)] = \int f(\zeta)\pi(\zeta)d\zeta.$$

5 If a sample $\zeta^{(t)}$, $t = 1, 2, \dots, T$, can be drawn from $\pi(\zeta)$, the expected value of any transformation of
6 ζ can be estimated by the average of the corresponding transformation of $\zeta^{(t)}$. Hence, the estimator
7 of $E[f(\zeta)]$ is

$$\tilde{E}[f(\zeta)] = \frac{1}{T} \sum_{i=1}^T f(\zeta^{(t)}).$$

8 This estimator is a method of moments estimator of $E[f(\zeta)]$, and by the strong law of large numbers
9 is a strongly consistent estimator of $E[f(\zeta)]$.

10
11 The approximation error of $E[f(\zeta)]$ can be probabilistically determined, and because the sample size T
12 is determined by the researcher, the approximation error can be made as small as desired by increasing
13 the number of draws from $\pi(\zeta)$.

14
15 In Bayesian inference, however, sampling directly from $\pi(\zeta)$ is not always possible.

16
17 One solution to this problem is to indirectly generate *independent* samples from $\pi(\zeta)$ using the
18 acceptance-rejection (A-R) method.

19 **Acceptance-Rejection Sampling**

20
21
22 From the equation of the posterior distribution, it is easy to see that the it can be written as

$$\pi(\zeta) = g(\zeta)/c$$

23 where $g(\zeta) = p(\zeta)p(\zeta|\mathbf{X})$ and c is the normalizing constant.

24
25 Suppose that for all values of ζ , a distribution $q(\zeta)$ which can be easily simulated and a known constant
26 k are available such that the blanketing condition $g(\zeta) \leq kq(\zeta)$ is satisfied, then a sample from $\pi(\zeta)$
27 can be drawn as follows:

- 28 1. Draw a candidate ζ^* from $q(\zeta)$;
- 29 2. Draw u from $\mathcal{U}(0,1)$;
- 30 3. Accept ζ^* if $u \leq g(\zeta^*)/kq(\zeta^*)$; otherwise, repeat steps 1 through 3.

31 Recall that this procedure is known as the *acceptance-rejection* (A-R) sampling. For the method to be
32 efficient, the k must be as small as possible.

33
34 However, finding the optimum k is in itself computationally formidable.

35

Hence, there is still a need for a straightforward method of drawing samples from any arbitrary distribution (i.e., MCMC).

Markov Chains

A sequence of random variable $\{\zeta^{(0)}, \zeta^{(1)}, \dots, \zeta^{(t)}, \dots\}$ is said to be a Markov chain if the transition kernel $p(\mathbf{s}, A)$ satisfies the condition

$$\begin{aligned} p(\mathbf{s}, A) &= p(\zeta^{(t+1)} \in A | \zeta^{(0)} = \mathbf{s}_0, \dots, \zeta^{(t-1)} = \mathbf{s}_{t-1}, \zeta^{(t)} = \mathbf{s}) \\ &= p(\zeta^{(t+1)} \in A | \zeta^{(t)} = \mathbf{s}), \end{aligned}$$

Stated differently, this equation means that, given the current state, the future and the past states are independent.

The transition kernel is the conditional probability that a chain will move from \mathbf{s} to a point in A .

The ergodic theorem states that if the expected value of a function of ζ , $E_\pi[f(\zeta)]$, is finite, then the ergodic average $\bar{f}(\zeta^{(T)}) = 1/T \sum_{t=1}^T f(\zeta^{(t)})$ converges in probability to $E_\pi[f(\zeta)]$ as $T \rightarrow \infty$. The central limit theorem for ergodic Markov chains provides the asymptotic distribution of the ergodic average as

$$\lim_{T \rightarrow \infty} \sqrt{T} \frac{\bar{f}(\zeta^{(T)}) - E_\pi[f(\zeta)]}{\sigma(f)} = N(0, 1),$$

where $\sigma(f)$ is real-valued but may be hard to obtain because of the dependence in the samples. The results above apply to most Markov chains used in simulations at present.

The Metropolis-Hastings Algorithm

This method was originally developed by *Metropolis*, Rosenbluth, Rosenbluth, Teller, and Teller (1953); a generalization of the method is provided by *Hastings* (1970).

As in the A-R method, M-H starts by sampling \mathbf{z} from a candidate-generating distribution q . However, unlike in A-R, the Markovian nature of the chain allows \mathbf{z} to depend on the current state \mathbf{s} .

For MCMC, the candidate-generating distribution is defined as $q(\mathbf{s}, \mathbf{z})$. For most applications, the use of $q(\mathbf{s}, \mathbf{z})$ alone will not satisfy an important condition (i.e., *reversibility*) for the samples to be deemed acceptable.

For instance, $q(\mathbf{s}, \mathbf{z})$ can be such that the process moves more often from \mathbf{s} to \mathbf{z} than the other way around. To achieve equilibrium, the number of moves from \mathbf{s} to \mathbf{z} can be restricted by introducing $\alpha(\mathbf{s}, \mathbf{z}) < 1$, the probability of move from \mathbf{s} to \mathbf{z} .

The M-H transition kernel from \mathbf{s} to \mathbf{z} is given by

$$p_{MH}(\mathbf{s}, \mathbf{z}) = q(\mathbf{s}, \mathbf{z})\alpha(\mathbf{s}, \mathbf{z}).$$

Whereas the A-R algorithm samples a new candidate value when \mathbf{z} is rejected, the M-H algorithm simply returns the value \mathbf{s} .

Hence, reversibility is satisfied if the transition kernel $p_{MH}(\mathbf{s}, \mathbf{z})$ is defined with

$$\alpha(\mathbf{s}, \mathbf{z}) = \min \left\{ \frac{\pi(\mathbf{z})q(\mathbf{z}, \mathbf{s})}{\pi(\mathbf{s})q(\mathbf{s}, \mathbf{z})}, 1 \right\}.$$

An important feature of the M-H algorithm, particularly for Bayesian applications, is that **the normalizing constant is not required** in computing $\alpha(\mathbf{s}, \mathbf{z})$ because it appears in both the numerator and denominator of $\alpha(\mathbf{s}, \mathbf{z})$.

Below are the steps in drawing T draws from $\pi(\boldsymbol{\zeta})$ using the M-H algorithm:

1. Start from an arbitrary value $\mathbf{s}^{(0)}$;
2. Draw a candidate \mathbf{z}^* from $q(\mathbf{s}^{(t)}, \mathbf{z})$;
3. Draw u from $\mathcal{U}(0,1)$;
4. Return $\mathbf{s}^{(t+1)} = \mathbf{z}^*$ if $u \leq \alpha(\mathbf{s}^{(t)}, \mathbf{z}^*)$; otherwise, return $\mathbf{s}^{(t+1)} = \mathbf{s}^{(t)}$.
5. Repeat steps 2 to 4 for $t = 1, 2, \dots, T$;

The Gibbs Sampler

Suppose the elements of $\boldsymbol{\zeta}$ can be written as $\boldsymbol{\zeta} = \{\zeta_1, \zeta_2, \dots, \zeta_d, \dots, \zeta_D\}$.

The joint distribution $\pi(\boldsymbol{\zeta})$ can be written as $\pi(\zeta_d, \boldsymbol{\zeta}'_d)$ where $\boldsymbol{\zeta}'_d = \{\zeta_1, \dots, \zeta_{d-1}, \zeta_{d+1}, \dots, \zeta_D\}$ for $d = 1, \dots, D$.

One particular goal in a multivariate set-up is to be able to describe the *marginal* distribution $\pi(\zeta_d)$ defined to be

$$\pi(\zeta_d) = \int \pi(\zeta_d, \boldsymbol{\zeta}'_d) d\boldsymbol{\zeta}'_d.$$

For example, one may be interested in the expected value of a function of ζ_d , $E[f(\zeta_d)]$. If a sample $\{\zeta_d^{(1)}, \zeta_d^{(2)}, \dots, \zeta_d^{(T)}\}$ can be drawn directly from $\pi(\zeta_d)$, then, as before, the expected value can be estimated by

$$\tilde{E}[f(\zeta_d)] = \frac{\sum_{t=1}^T f(\zeta_d^{(t)})}{T}.$$

The discussion in the previous section on M-H algorithm assumes that $\boldsymbol{\zeta}$ can be sampled from a multivariate distribution, and hence allows the algorithm to be carried out by treating $\boldsymbol{\zeta}$ en bloc. This may not always be easy or efficient to implement.

An alternative technique that often simplifies the calculation is to apply the algorithm to the elements of $\boldsymbol{\zeta}$ one at a time. Suppose the transition kernel $q_d(s_d, z_d)$ is defined as

$$q_d(s_d, z_d) = \pi(\zeta_d | \boldsymbol{\zeta}'_d) = \frac{\pi(\zeta_d, \boldsymbol{\zeta}'_d)}{\int \pi(\zeta_d, \boldsymbol{\zeta}'_d) d\zeta_d},$$

and is called the *full conditional* distribution of ζ_d .

Suppose further that samples can be drawn directly from these conditional distributions. If a sequence is drawn as follows:

1. Start from an arbitrary value $\zeta^{(0)} = \{\zeta_1^{(0)}, \zeta_2^{(0)}, \dots, \zeta_D^{(0)}\}$;
2. At iteration $t + 1$, draw a candidate ζ_d^* from $\pi(\zeta_d^{(t)} | \zeta_{<d}^{(t+1)}, \zeta_{>d}^{(t)})$;
3. Repeat for $d = 1, 2, \dots, D$;
4. Repeat steps 2 and 3 for $t = 1, 2, \dots, T$;

then the process that yields the sequence

$$\{\zeta_1^{(0)}, \zeta_2^{(0)}, \dots, \zeta_D^{(0)}, \zeta_1^{(1)}, \zeta_2^{(1)}, \dots, \zeta_D^{(1)}, \dots, \zeta_1^{(T)}, \zeta_2^{(T)}, \dots, \zeta_D^{(T)}\}$$

is said to be the *Gibbs sampler*.

In Gibbs sampling, all candidate values are accepted (i.e., $\alpha(s_d, z_d) = 1$).

For distribution with multiple components, there are many ways of obtaining draws from the marginal distribution which can be derived from the joint and conditional distributions.

However, the defining characteristic of the Gibbs sampler is its use of the full set of the univariate conditional distributions to define the iterations.

In practice, Gibbs sampling has referred to any algorithm that recursively samples from the full conditional distributions, univariate or multivariate, directly or indirectly.

Under mild conditions

$$\pi(\zeta_1^{(t)}, \zeta_2^{(t)}, \dots, \zeta_D^{(t)}) \xrightarrow{\mathcal{D}} \pi(\zeta_1, \zeta_2, \dots, \zeta_D)$$

as $t \rightarrow \infty$, where $\xrightarrow{\mathcal{D}}$ denotes convergence in distribution.

This implies that for all d , $\pi(\zeta_d^{(t)}) \xrightarrow{\mathcal{D}} \pi(\zeta_d)$ holds.

Moreover, provided that each of the components of ζ is visited infinitely often, convergence still holds even if the components are updated in no particular order, or only some of components are updated in each iteration.

From the above discussion, it can be seen that the Gibbs sampler is a special case of the more general M-H algorithm.

Because the Gibbs sampler is a special case of M-H algorithm, the ergodic theorem and convergence also apply to samples obtained via the Gibbs sampling.

However, for many researchers, this process is more appropriately called the *M-H-within-Gibbs algorithm*. That is, the sampling is from full conditional distributions (i.e., Gibbs sampling), and M-H is employed *within* Gibbs whenever the full condition distributions cannot be sampled directly.

That is, in situations where the full conditional distributions are unknown and cannot be sampled from directly, sampling from $q_d(s_d, z_d)$, which is typically different from $\pi(\zeta_d | \zeta_d')$, and then computing the corresponding $\alpha(s_d, z_d)$ can be used in conjunction with direct sampling from the known full conditional distributions.

1 **Implementation Issues**

3 *Candidate-Generating Distribution*

5 Described below are some of the families of candidate-generating distributions $q(\mathbf{s}, \mathbf{z})$ that are used in
6 the M-H algorithm.

8 Random Walk Chain

9 One family of distribution for simulating \mathbf{z} uses a random perturbation ϵ from a distribution $q(\epsilon)$ such
that

$$\mathbf{z} = \mathbf{s} + \epsilon.$$

10 It follows that $q(\mathbf{s}, \mathbf{z})$ is of the form $q(\mathbf{z} - \mathbf{s})$. Because the candidate value is just the current values
11 plus some noise, the Markov chain associated with this process is called a random walk.

12
13 The multivariate normal, t , and uniform distributions have been used as possible options for $q(\epsilon)$.
14 The advantage of the above distributions, and any symmetric distributions for that matter (i.e.,
15 $q(\epsilon)=q(-\epsilon)$), is that the probability of move simplifies to

$$\alpha(\mathbf{s}, \mathbf{z}) = \min \left\{ \frac{\pi(\mathbf{z})q(\mathbf{z}, \mathbf{s})}{\pi(\mathbf{s})q(\mathbf{s}, \mathbf{z})}, 1 \right\} = \min \left\{ \frac{\pi(\mathbf{z})}{\pi(\mathbf{s})}, 1 \right\}.$$

16 Independence Chain

17
18 Another family of distributions for simulating \mathbf{z} are those that draw the candidate value independent
19 of \mathbf{s} (i.e., $q(\mathbf{s}, \mathbf{z}) = q(\mathbf{z})$).

20
21 This gives a sequence that is called an *independence chain*.

22
23 Whereas the scale parameters are sufficient for generating the random walk chains, both location and
24 scale parameters are needed to be specified in generating an independence chain.

25
26 It should be noted that although the process of generating candidate values are independent of the
27 current values, the resulting chain is still a Markov chain because the $\alpha(\mathbf{s}, \mathbf{z})$ depends on \mathbf{s} .

29 M-H Acceptance-Rejection Sampling

30
31 The M-H acceptance-rejection method circumvents the problem of having to find the constant k that
32 satisfies the blanketing condition required in the acceptance-rejection sampling.

33
34 For this method, the choice of k does not necessarily result in the target distribution being blanketed
35 or dominated for all values of \mathbf{s} .

36
37 The method proceeds by dividing the support into two regions: dominated and non-dominated. The
38 probability of move $\alpha(\mathbf{s}, \mathbf{z})$ is then adjusted depending on whether \mathbf{s} or \mathbf{z} or both fall in the dominated
39 region.

40
41 Other families of candidate-generating distributions exist including those that give autoregressive
42 chains and grid-based chains.

43

1 Regardless of how the candidate values are sampled, the whole support must be traversed for the chain
2 to converge to the desired distribution. This is why the choice of $q(\mathbf{s}, \mathbf{z})$, which determines $\alpha(\mathbf{s}, \mathbf{z})$, is
3 of particular importance.

4
5 The size of the move or jump needs to be carefully managed. On one hand, candidate-generating dis-
6 tributions that propose small moves will have high acceptance rate. On the other hand, distributions
7 that propose big moves will invariably produce values that will have low probabilities of acceptance.

8
9 In both cases, it will take a long time for the chain to traverse the entire support, particularly in the
10 tails.

11
12 Hence, a large number of iterations is needed before the chain can converge. This problem can be
13 diagnosed by large autocorrelations among the sampled values.

14
15 To improve the rate of convergence, candidate-generating distributions similar to the target distribu-
16 tion can be used. Sampling from this distribution can produce high acceptance rate without making
17 moves that are too small.

18
19 Another way of achieving faster convergence is to fine-tune the scale parameter so that an optimal
20 acceptance rate can be obtained.

21
22 Although the results are not conclusive, some heuristic rules exist regarding the optimal acceptance
23 rate. It is recommended that for one or two dimensions, the optimal rate is around 0.50, whereas for
24 higher dimensions, the optimal rate is around 0.25.

25
26 Equivalently, in rare occasions where the asymptotic variance of the $\zeta^{(t)}$ is available, the scale param-
27 eter of $2.28/\sqrt{D}$ times the asymptotic variance can be used.

28
29 It is possible that even with these optimal acceptance rates, the autocorrelation may remain high.
30 Such situations warrant the use of other families of candidate-generating distributions.

31
32 High autocorrelation does not mean the chain will not converge. It simply means we need more itera-
33 tions to get there.

34 *Blocking*

35
36
37 The updating in M-H and Gibbs sampling sections can be viewed as two types of blocking that are
38 used in MCMC.

39
40 However, aside from updating the elements of ζ one at a time or simultaneously, the algorithm can
41 also be applied to blocks of ζ .

42
43 Using blocks avoids the problem of low acceptance rates that are often encountered when updating is
44 done simultaneously.

45
46 Furthermore, in some applications the elements of ζ can have natural and meaningful groupings, up-
47 dating the elements one at a time may be inefficient.

48
49 A recommendation is to block elements that are correlated and belong to the same context in the

1 model formulation.

2
3 Updating the chain is done by sampling directly from the conditional distributions of the blocks when-
4 ever possible, or using the M-H algorithm otherwise.

5 *Number of Chains*

6
7
8 Several options as to the number of chains that need to be run have been proposed.

9
10 The first is the use of many short chains as suggested. However, unless independent samples from
11 the stationary distribution are required, this approach may be wasteful because independence is not
12 required for the ergodic theorem to hold.

13
14 The more serious debate is choosing between several long chains and one very long chain.

15
16 This issue is closely related to the problem of convergence diagnosis (i.e., determining which is the
17 better procedure in assessing convergence).

18
19 One advantage of using multiple chains is that if the chains have not reached the stationary distribu-
20 tion, the use of multiple chains from dispersed starting points may detect the lack of convergence.

21
22 There are cases where each of the chains appear to have converged when examined individually but
23 may exhibit lack of convergence when compared to other chains.

24
25 Another benefit of using multiple chains is that by starting from dispersed points, the dependence of
26 the chains on the starting values is reduced.

27
28 However, one drawback of using multiple chains in that convergence is dictated by the slower chains.
29 Also, the MCMC algorithm needs to be altered to process samples obtained from different chains.

30
31 Finally, because a single chain that gives misleading results do not happen too often, the extra cost in
32 implementing this procedure may not be worthwhile.

33
34 The major shortcoming of using a single chain is that it severely suffers from the “you’ve only seen
35 where you’ve been” defect (i.e., parts of the support that have not been visited by the chain may not
36 be detected).

37
38 The primary strength of using one long chain, especially for applications where the algorithm mixes
39 slowly, is that it has a higher chance of reaching the stationary distribution compared to comparative
40 number of samples from multiple but shorter chains.

41 *Starting Values*

42
43
44 Since the MCMC algorithm is based on Markov chains, it is assured to converge to the stationary
45 distribution $\pi(\zeta)$ regardless of the starting value - as long as the chain is run long enough, it will forget
46 its initial value.

47
48 However, care should still be exercised in selecting the initial values because extreme values can in-
49 crease the number of iterations before convergence.

1
2 In severe cases, because of the numerical instability at the extreme tails of the distribution, the chain
3 may fail to converge to the main support of the distribution.

4 5 *Assessing Convergence*

6
7 In practical applications of the MCMC methods, two of the important decisions to be made pertain
8 to the length of burn-in and the length of the chain. These two questions are part of a bigger issue,
9 that of convergence.

10
11 Burn-in refers to the sufficiently long number of iterations m before a sample $\zeta^{(m+1)}$ and values there-
12 after, are said to be drawn from the stationary distribution $\pi(\zeta)$.

13
14 The samples drawn during the burn-in period, presumably still affected by the starting values, are
15 usually discarded, and inference is based on the remaining $T - m$ samples only.

16
17 The most obvious and commonly used method of determining burn-in is visual inspection of the MCMC
18 output.

19
20 Formal methods of determining the length of burn-in are available, which is subsumed under conver-
21 gence diagnostics.

22
23 The more practical of these methods are those that provide assessment of convergence based of the
24 output of the chain alone without resorting to the transition kernel.

25
26 However, it has been suggested that if extreme starting values are avoided, computing m may not be
27 necessary for a sufficiently long chain because it is usually negligible, between 1% to 2% of the total
28 length of the chain.

29
30 Determining the length of the chain is motivated by obtaining a desired precision of the estimate of
31 $E[f(\zeta)]$.

32
33 The variance of $\tilde{E}[f(\zeta)]$ is called the *Monte Carlo error*. One obvious way of computing the Monte
34 Carlo error is to run several chains and then compute the variance of the estimates.

35
36 We want to minimize this variance relative to $\tilde{V}[f(\zeta)]$ by using a sufficiently long enough chain.

37
38 The components of ζ can be assessed for convergence one at a time using the Gelman-Rubin poten-
39 tial scale reduction factor statistics, or en bloc using the Brooks-Gelman multivariate potential scale
40 reduction factor.

41 42 *Output Analysis*

43
44 Aside monitoring convergence, the output of the chain is primarily used for inference.

45
46 The earlier discussion of ergodic theorem specifically addressed the asymptotic behavior of the mean
47 of a function of the sample.

1 However, as stated earlier, in practice chains are not run infinitely, and for a multi-component ζ , the
2 interest lies in the characteristics of the marginal distributions.

3
4 The characteristics of interest of a particular component of ζ can be computed by using the draws
5 after the burn-in $\{\zeta_d^{(m+1)}, \dots, \zeta_d^{(T)}\}$, and at the same time ignoring the output for other components.

6
7 For example, the mean and variance of a function of ζ_d can be estimated as

$$\tilde{E}[f(\zeta_d)] = \frac{1}{T-m} \sum_{t=m+1}^T f(\zeta_d^{(t)}),$$

8 and

$$\tilde{V}[f(\zeta_d)] = \frac{1}{T-m-1} \sum_{t=m+1}^T \left(f(\zeta_d^{(t)}) - \tilde{E}[f(\zeta_d)] \right)^2,$$

9 respectively.

10
11 The same procedure can be followed in computing the quantiles and credibility intervals of interest.

12
13 **Important Note:** $f(\zeta_d) = \zeta_d$ is just a special case. For example, instead of ζ_d , we might be interested
14 in $\log(\zeta_d)$.

15
16 In addition, marginal densities can also be obtained from the output as the mean of the conditional
17 densities.

18
19 For some purposes, saving all the MCMC draws may require a large storage. As a result, a strategy
20 of saving only every r^{th} draws after the burn-in, called *thinning*, has been implemented to deal with
21 this problem.

22
23 At the same time, thinning the chain, which produces approximately independent draws, has been
24 used to handle the dependency problem in the samples.

25
26 However, aside from advantages is computer storage, this method produces no gain in efficiency, and
27 estimation based on the thinned chain is, in fact, less precise compared to estimation based on the
28 entire chain.

29 30 Preliminary Examples

31 32 Univariate Normal Distribution

33
34 Suppose $X|\theta \sim N(\theta, \sigma^2)$, where σ^2 is known. We also assume the prior distribution $\theta \sim N(\mu_0, \tau_0^2)$.

35
36 We can write the prior and likelihood as:

$$\begin{aligned} p(\theta) &= \frac{1}{\sqrt{2\pi\tau}} \exp \left[-\frac{1}{2} \frac{(\theta - \mu_0)^2}{\tau_0^2} \right] \\ &\propto \exp \left[-\frac{1}{2} \frac{(\theta - \mu_0)^2}{\tau_0^2} \right] \end{aligned}$$

$$\begin{aligned}
p(X|\theta) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{(x-\theta)^2}{\sigma^2}\right] \\
&\propto \exp\left[-\frac{1}{2} \frac{(x-\theta)^2}{\sigma^2}\right].
\end{aligned}$$

1 What we are interested in is the posterior distribution of θ given X . As in, what is $p(\theta|X)$?

2
3 Note that a normal distribution is considered a *conjugate prior* for a normal likelihood because the
4 resulting posterior is also normal.

5
6 This mean that we are assured that $p(\theta|X)$ is normal. We just need to find the posterior mean μ_1 and
7 variance τ_1^2 .

8
9 By definition of a posterior,

$$\begin{aligned}
p(\theta|X) &\propto p(X|\theta) \times p(\theta) \\
&\propto \exp\left[-\frac{1}{2} \frac{(X-\theta)^2}{\sigma^2}\right] \times \exp\left[-\frac{1}{2} \frac{(\theta-\mu_0)^2}{\tau_0^2}\right] \\
&= \exp\left\{-\frac{1}{2} \left[\frac{(X-\theta)^2}{\sigma^2} + \frac{(\theta-\mu_0)^2}{\tau_0^2}\right]\right\}.
\end{aligned}$$

We need to simplify the exponent so the posterior will be of the form

$$p(\theta|X) \propto \exp\left[-\frac{1}{2} \frac{(\theta-\mu_1)^2}{\tau_1^2}\right].$$

10 Note that for this problem the only random variable of interest is θ because once the data (i.e., x) are
11 observed, they are no longer random (i.e., they can be considered fixed).

12
13 To get the desired form, we first need to expand the terms in the exponent.

$$\begin{aligned}
p(\theta|X) &\propto \exp\left[-\frac{1}{2} \left(\frac{X^2}{\sigma^2} - 2\frac{X\theta}{\sigma^2} + \frac{\theta^2}{\sigma^2} + \frac{\theta^2}{\tau_0^2} - 2\frac{\mu_0\theta}{\tau_0^2} + \frac{\mu_0^2}{\tau_0^2}\right)\right] \\
&\propto \exp\left[-\frac{1}{2} \left(-2\frac{X\theta}{\sigma^2} + \frac{\theta^2}{\sigma^2} + \frac{\theta^2}{\tau_0^2} - 2\frac{\mu_0\theta}{\tau_0^2}\right)\right] \\
&= \exp\left\{-\frac{1}{2} \left[\theta^2 \left(\frac{1}{\tau_0^2} + \frac{1}{\sigma^2}\right) - 2\theta \left(\frac{\mu_0}{\tau_0^2} + \frac{x}{\sigma^2}\right)\right]\right\}
\end{aligned}$$

14 To complete the square, we can let

$$\begin{aligned}
\frac{1}{\tau_1^2} &= \frac{1}{\tau_0^2} + \frac{1}{\sigma^2} \\
\mu_1 &= \frac{\frac{1}{\tau_0^2}\mu_0 + \frac{1}{\sigma^2}X}{\frac{1}{\tau_0^2} + \frac{1}{\sigma^2}} \\
&= \left(\frac{\mu_0}{\tau_0^2} + \frac{X}{\sigma^2}\right) \times \tau_1^2.
\end{aligned}$$

1 Then

$$\begin{aligned} p(\theta|X) &\propto \exp \left[-\frac{1}{2} \left(\frac{\theta^2}{\tau_1^2} - 2\theta \frac{\mu_1}{\tau_1^2} \right) \right] \\ &= \exp \left[-\frac{1}{2} \left(\frac{\theta^2 - 2\theta\mu_1}{\tau_1^2} \right) \right] \\ &= \exp \left[-\frac{1}{2} \left(\frac{\theta^2 - 2\theta\mu_1 + \mu_1^2}{\tau_1^2} \right) + c \right] \\ &\propto \exp \left[-\frac{1}{2} \frac{(\theta - \mu_1)^2}{\tau_1^2} \right]. \end{aligned}$$

2 Therefore,

$$\begin{aligned} \theta|X &\sim N(\mu_1, \tau_1^2) \\ &= N \left(\frac{\frac{1}{\tau_0^2}\mu_0 + \frac{1}{\sigma^2}X}{\frac{1}{\tau_0^2} + \frac{1}{\sigma^2}}, \left(\frac{1}{\tau_0^2} + \frac{1}{\sigma^2} \right)^{-1} \right). \end{aligned}$$

3 What happens to the influence of the prior on the posterior as τ_0^2 gets larger?

4
5 If we have N independent draws from $x_i|\theta \sim N(\theta, \sigma^2)$, then the posterior can be written as

$$\begin{aligned} \theta|\mathbf{X} &\sim N(\mu_N, \tau_N^2) \\ &= N \left(\frac{\frac{1}{\tau_0^2}\mu_0 + \frac{N}{\sigma^2}\bar{X}}{\frac{1}{\tau_0^2} + \frac{N}{\sigma^2}}, \left(\frac{1}{\tau_0^2} + \frac{N}{\sigma^2} \right)^{-1} \right) \end{aligned}$$

6 1. Show that the posterior distribution of $\theta|\mathbf{X}$ is indeed $N(\mu_N, \tau_N^2)$.

7
8 2. What happens to the posterior when the prior becomes less and less informative (i.e., $\tau_0^2 \rightarrow \infty$)?

9
10 Once we know the posterior distribution, we can now make inferences about the distribution. For
11 example, we can now find its mean $E(\theta|\mathbf{X})$, variance $V(\theta|\mathbf{X})$, and the 95% CI (credibility interval)
12 for $\theta|\mathbf{X}$.

13
14 At this point, we have a few ways of doing it.

15
16 1. Since we know the exact form of the posterior distribution, we can derive μ_N , τ_N^2 , and 95% CI for
17 $\theta|\mathbf{X}$ analytically.

18
19 2. Again since we know the form of the distribution, can sample *directly* from the posterior.

20
21 3. For pedagogical purposes, we can also sample from the posterior indirectly using the M-H algorithm.

22
23 For 1, the solution is:

$$\begin{aligned} \tilde{E}(\theta|\mathbf{X}) &= \mu_N \\ \tilde{V}(\theta|\mathbf{X}) &= \tau_N^2 \\ \text{95\% CI for } \theta|\mathbf{X} &: [\mu_N - 1.96 \times \tau_N, \mu_N + 1.96 \times \tau_N] \end{aligned}$$

1 For 2, the solution is:

2

3 First, sample T draws from $N(\theta_N, \tau_N^2)$. Call the the sample $\theta^* = \{\theta^{*(t)}\}$.

$$\begin{aligned}\tilde{E}(\theta|\mathbf{X}) &= \bar{\theta}^* = \frac{\sum_{t=1}^T \theta^{*(t)}}{T} \\ \tilde{V}(\theta|\mathbf{X}) &= \frac{\sum_{t=1}^T (\theta^{*(t)} - \bar{\theta}^*)^2}{T-1} \\ 95\% \text{ CI for } \theta|\mathbf{X} &: [q_{.025}(\theta^*), q_{.975}(\theta^*)]\end{aligned}$$

4 where $q_\alpha(\mathbf{x})$ represents the α^{th} quantile of the draws in \mathbf{x} .

5

6 For 3, we need to first set up the M-H algorithm.

7

8 We will use a random walk chain, where $\epsilon \sim N(0, \sigma_\epsilon^2)$, as in, $\theta^* = \theta^{*(t)} + \epsilon$.

9

10 Next we need to define the transition probability $\alpha(\theta^{*(t)}, \theta^*)$ at iteration $t+1$. Because of the symmet-
11 ric nature of the candidate-generating distribution, we do not need $q(\theta^{*(t)}, \theta^*)$ and $q(\theta^*, \theta^{*(t)})$ because
12 they will cancel each other.

13

14 In this example,

$$q(\theta^{*(t)}, \theta^*) \propto \exp\left[-\frac{1}{2} \frac{(\theta^* - \theta^{*(t)})^2}{\sigma_\epsilon}\right]$$

15 and

$$q(\theta^*, \theta^{*(t)}) \propto \exp\left[-\frac{1}{2} \frac{(\theta^{*(t)} - \theta^*)^2}{\sigma_\epsilon}\right],$$

16 which are equal to each other.

17

18 Thus,

$$\begin{aligned}\alpha(\theta^{*(t)}, \theta^*) &= \min\left\{\frac{p(\theta^*|\mathbf{X})}{p(\theta^{*(t)}|\mathbf{X})}, 1\right\} \\ &= \min\left\{\frac{\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \sum_{i=1}^N \frac{(X_i - \theta^*)^2}{\sigma^2}\right] \times \frac{1}{\sqrt{2\pi}\tau_0} \exp\left[-\frac{1}{2} \frac{(\theta^* - \mu_0)^2}{\tau^2}\right]}{\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \sum_{i=1}^N \frac{(X_i - \theta^{*(t)})^2}{\sigma^2}\right] \times \frac{1}{\sqrt{2\pi}\tau_0} \exp\left[-\frac{1}{2} \frac{(\theta^{*(t)} - \mu_0)^2}{\tau^2}\right]}, 1\right\} \\ &= \min\left\{\frac{\exp\left[-\frac{1}{2} \sum_{i=1}^N \frac{(X_i - \theta^*)^2}{\sigma^2}\right] \times \exp\left[-\frac{1}{2} \frac{(\theta^* - \mu_0)^2}{\tau^2}\right]}{\exp\left[-\frac{1}{2} \sum_{i=1}^N \frac{(X_i - \theta^{*(t)})^2}{\sigma^2}\right] \times \exp\left[-\frac{1}{2} \frac{(\theta^{*(t)} - \mu_0)^2}{\tau^2}\right]}, 1\right\}\end{aligned}$$

19 As can be seen from above, when using M-H, we only need the posterior distribution to up a proportion-
20 ality constant. In addition, with symmetric candidate-generating distribution $q(\theta^{(t)}, \theta^*) = q(\theta^*, \theta^{(t)})$,
21 we only need to compute the $p(\mathbf{X}|\theta) \times p(\theta)$.

22

23 We will accept θ^* , as in, set $\theta^{*(t+1)} = \theta^*$ if $\alpha(\theta^{*(t)}, \theta^*) > u$; otherwise, $\theta^{*(t+1)} = \theta^{*(t)}$.

24

1 Let m be the burn-in. We need to sample $T + m$ draws using the M-H algorithm. Call the sample
2 $\theta^{**} = \{\theta^{**(t)}\}$, where $t = 1, \dots, T + m$.

3
4 From this sample, discard the first m draws. Call the remaining sampled draws $\theta^* = \{\theta^{*(t+m)}\} =$
5 $\{\theta^{*(t)}\}$, where $t = 1, \dots, T + m$.

6
7 Once we have θ^* , proceed in the same manner as in 2.

8
9 *MCMC Implementation in Ox for a Univariate Normal with Unknown Mean*

10
11 The program `uniN.ox` will makes inferences about the unknown μ in $N(\mu, \sigma^2)$. Inferences are made in
12 three different ways: analytically, by direct sampling, and by M-H sampling.

```
#include <oxstd.h>
#include <oxprob.h>
#include <oxfloat.h>
```

13 We need to declare a few global variables. Except for `X`, these variables can be manipulated by the
14 user.

```
decl N=1000; //sample size
decl mu0=0,tau20=1; //prior parameters
decl mu=0.5,sigma2=2; //parameters for X|mu
decl T=10000; //size of sample from the posterior
decl m=5000; //burn-in
decl X; //data
```

15 The function `MH_N` samples from the posterior using the M-H algorithm. The argument `s0` is the
16 current value $\theta^{*(t)}$ of the chain.

```
MH_N(decl s0){//s0 is the current value
decl s1=.1*rann(1,1)+s0;//sampling candidate value from N(s0,.1^2)
decl ll0=sumc(-.5*(X-s0).^2/sigma2)+(-.5*(s0-mu0)^2/tau20);//posterior when par=s0
decl ll1=sumc(-.5*(X-s1).^2/sigma2)+(-.5*(s1-mu0)^2/tau20);//posterior when par=s1
decl alpha=exp(ll1-ll0);//acceptance probability
if(alpha>ranu(1,1)) s0=s1;//decide to accept sampled value s1 or keep s0
return s0;
} //end MH_N
```

17 In the `main` function, we set the random seed, and declare a few local variables.

18

```

ranseed(120);
decl smp; //sample from the posterior
decl muN,tau2N; //posterior parameters
decl m_est,t2_est; //est of mu_N and tau2N
decl LL,UL; //lower and upper limits

```

1 The next code generates the data X , and solves μ_N and τ_N^2 analytically.

2

```

X=rann(N,1)*sqrt(sigma2)+mu;//generating data
tau2N=1/(1/tau20+N/sigma2);//analytic posterior variance
muN=(mu0/tau20+N/sigma2*meanc(X))*tau2N;//analytic posterior mean

```

3 For the **analytic** solution, we use μ_N , τ_N^2 , and $\mu_N \pm 1.96\tau_N$ as our estimates.

4

```

print("\nAnalytic Solution");
m_est=muN;
t2_est=tau2N;
LL=muN-1.96*sqrt(t2_est);
UL=muN+1.96*sqrt(t2_est);
print("\nPosterior Mean and Variance: ",m_est~t2_est);
print("95% Credibility Interval: ",LL~UL);

```

5 For the *direct-sampling* solution, we first sample **smp** from $N(\mu_N, \tau_N^2)$. We then estimate the parameters
6 of interest using their sample counterparts.

```

print("\nSolution Based on Direct Sampling");
smp=rann(T,1)*sqrt(tau2N)+muN;
m_est=meanc(smp);
t2_est=varc(smp);
LL=quantilec(smp,.025);
UL=quantilec(smp,.975);
print("\nPosterior Mean and Variance: ",m_est~t2_est);
print("95% Credibility Interval: ",LL~UL);

```

7 For the *indirect-sampling* solution, we use \bar{x} as the initial value; **smp** is set to a null vector. We then
8 write a loop to sample indirectly from the posterior using the function **MH_N**.

9

1 We then drop the burn-in (i.e., the first m draws) from `smp`. The remaining steps are identical to the
 2 solution above.

```
print("\nSolution Based on Indirect Sampling");
decl smp0=meanc(X);
smp=<>;
for(decl t=0;t<T+m;t++){//M-H step
smp0=MH_N(smp0);
smp=smp|smp0;
}//end t
smp=smp[m:];
m_est=meanc(smp);
t2_est=varc(smp);
LL=quantilec(smp,.025);
UL=quantilec(smp,.975);
print("\nPosterior Mean and Variance: ",m_est~t2_est);
print("95% Credibility Interval: ",LL~UL);
```

3 Change the values of N , μ_0 , τ_0^2 , μ , σ^2 , T , and m to see how it affects the results.

4

5 *Multivariate Normal Distribution*

6

7 Suppose $\mathbf{X}|\boldsymbol{\theta} \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is known. We also assume the prior distribution $\boldsymbol{\theta} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0)$.
 8 Both $\boldsymbol{\Sigma}$ and $\boldsymbol{\Lambda}_0$ need to be positive definite.

9

10 In the multivariate normal, the prior and likelihood are written as:

$$\begin{aligned} p(\boldsymbol{\theta}) &\propto \exp \left[-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)' \boldsymbol{\Lambda}_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_0) \right] \\ p(\mathbf{X}|\boldsymbol{\theta}) &\propto \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\theta}) \right]. \end{aligned}$$

11 What we need is the posterior distribution $p(\boldsymbol{\theta}|\mathbf{X})$.

12

13 Again, because a multivariate normal distribution is considered a conjugate prior for a normal likeli-
 14 hood, we are assured that $p(\boldsymbol{\theta}|\mathbf{X})$ is also normal.

15

16 This means that we just need to find the posterior mean $\boldsymbol{\mu}_1$ and variance $\boldsymbol{\Lambda}_1$.

17

18 Again,

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{X}) &\propto p(\mathbf{X}|\boldsymbol{\theta}) \times p(\boldsymbol{\theta}) \\ &= \exp \left\{ -\frac{1}{2} [(\mathbf{X} - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\theta}) + (\boldsymbol{\theta} - \boldsymbol{\mu}_0)' \boldsymbol{\Lambda}_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_0)] \right\}. \end{aligned}$$

We need to simplify the exponent so the posterior will be of the form

$$p(\boldsymbol{\theta}|\mathbf{X}) \propto \exp \left[-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu}_1)' \boldsymbol{\Lambda}_1^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_1) \right].$$

1 To get the desired form, we need to expand the terms in the exponent.

$$\begin{aligned}
 p(\boldsymbol{\theta}|\mathbf{X}) &\propto \exp \left[-\frac{1}{2} (\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X} - 2\mathbf{X}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} + \boldsymbol{\theta}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} + \boldsymbol{\theta}'\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\theta} - 2\boldsymbol{\theta}'\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + \boldsymbol{\mu}_0'\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0) \right] \\
 &\propto \exp \left[-\frac{1}{2} (-2\mathbf{X}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} + \boldsymbol{\theta}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} + \boldsymbol{\theta}'\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\theta} - 2\boldsymbol{\theta}'\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0) \right] \\
 &= \exp \left\{ -\frac{1}{2} [\boldsymbol{\theta}'(\boldsymbol{\Lambda}_0^{-1} + \boldsymbol{\Sigma}^{-1})\boldsymbol{\theta} - 2\boldsymbol{\theta}'(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + \boldsymbol{\Sigma}^{-1}\mathbf{X})] \right\}
 \end{aligned}$$

2 Why is $\mathbf{X}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\theta} = \boldsymbol{\theta}'\boldsymbol{\Sigma}^{-1}\mathbf{X}$?

3

4 To complete the square, let

$$\begin{aligned}
 \boldsymbol{\Lambda}_1^{-1} &= \boldsymbol{\Lambda}_0^{-1} + \boldsymbol{\Sigma}^{-1} \\
 \boldsymbol{\mu}_1 &= \boldsymbol{\Lambda}_1(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + \boldsymbol{\Sigma}^{-1}\mathbf{X}).
 \end{aligned}$$

5 Then

$$\begin{aligned}
 p(\boldsymbol{\theta}|\mathbf{X}) &\propto \exp \left[-\frac{1}{2} (\boldsymbol{\theta}'\boldsymbol{\Lambda}_1^{-1}\boldsymbol{\theta} - 2\boldsymbol{\theta}'\boldsymbol{\Lambda}_1^{-1}\boldsymbol{\mu}_1) \right] \\
 &= \exp \left[-\frac{1}{2} (\boldsymbol{\theta}'\boldsymbol{\Lambda}_1^{-1}\boldsymbol{\theta} - 2\boldsymbol{\theta}'\boldsymbol{\Lambda}_1^{-1}\boldsymbol{\mu}_1 + \boldsymbol{\mu}_1'\boldsymbol{\Lambda}_1^{-1}\boldsymbol{\mu}_1) + c \right] \\
 &\propto \exp \left[-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu}_1)' \boldsymbol{\Lambda}_1^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_1) \right].
 \end{aligned}$$

6 Therefore,

$$\begin{aligned}
 \boldsymbol{\theta}|\mathbf{X} &\sim N(\boldsymbol{\mu}_1, \boldsymbol{\Lambda}_1) \\
 &= N\left((\boldsymbol{\Lambda}_0^{-1} + \boldsymbol{\Sigma}^{-1})^{-1}(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + \boldsymbol{\Sigma}^{-1}\mathbf{X}), (\boldsymbol{\Lambda}_0^{-1} + \boldsymbol{\Sigma}^{-1})^{-1}\right)
 \end{aligned}$$

7 If we have N independent draws from $\mathbf{x}_i|\boldsymbol{\theta} \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma})$, show that the posterior can be written as

$$\begin{aligned}
 \boldsymbol{\theta}|\mathbf{X} &\sim N(\boldsymbol{\mu}_N, \boldsymbol{\Lambda}_N) \\
 &= N\left((\boldsymbol{\Lambda}_0^{-1} + N\boldsymbol{\Sigma}^{-1})^{-1}(\boldsymbol{\Lambda}_0^{-1}\boldsymbol{\mu}_0 + N\boldsymbol{\Sigma}^{-1}\bar{\mathbf{X}}), (\boldsymbol{\Lambda}_0^{-1} + N\boldsymbol{\Sigma}^{-1})^{-1}\right)
 \end{aligned}$$

8 What does it mean for the prior to be uninformative (i.e., what should $\boldsymbol{\Lambda}_0$ be)?

9

10 Let $\mathbf{X}_{2 \times 1} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is known. Suppose $\mathbf{X}_{N \times 2}$, and we are interested in estimating $E(\boldsymbol{\theta}|\mathbf{X})$,
 11 $V(\boldsymbol{\theta}|\mathbf{X})$, and 95% CIs for $\theta_1|\mathbf{X}$ and $\theta_2|\mathbf{X}$. (Note that these are *marginal* (i.e., individual) CIs, not
 12 *joint* (i.e., simultaneous) credibility ellipse.)

13

14 As in the univariate case, we have at least three ways of solving the problem: analytically, sampling
 15 directly from the posterior, and sampling indirectly from the posterior.

16

17 For the analytic solution:

$$\begin{aligned}
 \tilde{E}(\boldsymbol{\theta}|\mathbf{X}) &= \boldsymbol{\mu}_N \\
 \tilde{V}(\boldsymbol{\theta}|\mathbf{X}) &= \boldsymbol{\Lambda}_N \\
 \text{95\% CI for } \theta_d|\mathbf{X} &: \mu_{Nd} \pm 1.96 \times \sqrt{\lambda_{Ndd}}
 \end{aligned}$$

1 for $d = 1, 2$, and where λ_{Ndd} is the d^{th} diagonal element of $\mathbf{\Lambda}_N$.

2
3 For the direct sampling solution:

4
5 We sample T draws from $N(\boldsymbol{\mu}_N, \mathbf{\Lambda}_N)$, and call it $\boldsymbol{\Theta}^* = \{\boldsymbol{\theta}^{*(t)}\}$.

$$\begin{aligned}\tilde{E}(\boldsymbol{\theta}|\mathbf{X}) &= \bar{\boldsymbol{\theta}}^* = \frac{\sum_{t=1}^T \boldsymbol{\theta}^{*(t)}}{T} \\ \tilde{V}(\boldsymbol{\theta}|\mathbf{X}) &= \frac{\sum_{t=1}^T \left(\boldsymbol{\theta}^{*(t)} - \bar{\boldsymbol{\theta}}^* \right) \left(\boldsymbol{\theta}^{*(t)} - \bar{\boldsymbol{\theta}}^* \right)'}{T-1} \\ 95\% \text{ CI for } \boldsymbol{\theta}_d|\mathbf{X} &: [q_{0.025}(\boldsymbol{\theta}_d^*), q_{0.975}(\boldsymbol{\theta}_d^*)]\end{aligned}$$

6 for $d = 1, 2$, and where $\boldsymbol{\theta}_d^*$ is the vector of sampled draws for dimension d .

7
8 For the indirect sampling solution, we need to first set up the M-H algorithm.

9
10 We will again use a random walk chain and symmetric candidate generating distribution, as in
11 $\boldsymbol{\theta}_d^* = \boldsymbol{\theta}_d^{*(t)} + \epsilon_d$, where $\epsilon_d \sim N(0, \sigma_\epsilon^2)$. Even though $\boldsymbol{\theta}$ is multivariate, and possibly correlated, it is
12 acceptable to use univariate and independent candidate generating distributions.

13
14 Next we need to define the transition probability $\alpha(\boldsymbol{\theta}^{*(t)}, \boldsymbol{\theta}^*)$ at iteration $t + 1$.

15
16 Thus,

$$\begin{aligned}\alpha(\boldsymbol{\theta}^{*(t)}, \boldsymbol{\theta}^*) &= \min \left\{ \frac{p(\boldsymbol{\theta}^*|\mathbf{X})}{p(\boldsymbol{\theta}^{*(t)}|\mathbf{X})}, 1 \right\} \\ &= \min \left\{ \frac{\exp \left[-\frac{1}{2} \sum_{i=1}^N (\mathbf{X}_i - \boldsymbol{\theta}^*)' \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\theta}^*) - \frac{1}{2} (\boldsymbol{\theta}^* - \boldsymbol{\mu}_0)' \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\theta}^* - \boldsymbol{\mu}_0) \right]}{\exp \left[-\frac{1}{2} \sum_{i=1}^N (\mathbf{X}_i - \boldsymbol{\theta}^{*(t)})' \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\theta}^{*(t)}) - \frac{1}{2} (\boldsymbol{\theta}^{*(t)} - \boldsymbol{\mu}_0)' \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\theta}^{*(t)} - \boldsymbol{\mu}_0) \right]}, 1 \right\}\end{aligned}$$

17 The ratio involves the term of the form $\sum_{i=1}^N \mathbf{y}_i' \mathbf{V} \mathbf{y}_i$ in the exponent. Computationally, this can be
18 expensive to implement as it would require writing a `for` loop that computes $\mathbf{y}' \mathbf{V} \mathbf{y}$ one observation at
19 a time because of the two matrix multiplications involved.

20
21 As an alternative, we can use Cholesky decomposition to get $\mathbf{V} = \mathbf{C}\mathbf{C}'$, where \mathbf{C} is a lower triangular
22 matrix.

23
24 This allows us to write

$$\begin{aligned}\mathbf{y}' \mathbf{V} \mathbf{y} &= \mathbf{y}' \mathbf{C} \mathbf{C}' \mathbf{y} \\ &= (\mathbf{C}' \mathbf{y})' \mathbf{C}' \mathbf{y} \\ &= \mathbf{y}^{*'} \mathbf{y}^* \\ &= \sum_{d=1}^D y_d^{*2}\end{aligned}$$

25 This means that, instead of two the matrix multiplications in $\mathbf{y}' \mathbf{V} \mathbf{y}$, we only need to carry out one
26 matrix multiplication (i.e., $\mathbf{C}' \mathbf{y}$), and then sum the square of the resulting terms.

1 In Ox, this means you we can now use `sumr` and `sumc` without using a `for` loop.

2
3 Note that $\alpha(\boldsymbol{\theta}^{*(t)}, \boldsymbol{\theta}^*)$ is a scalar so we will accept or reject $\boldsymbol{\theta}^*$ en bloc.

4
5 Set $\boldsymbol{\theta}^{*(t+1)} = \boldsymbol{\theta}^*$ if $\alpha(\boldsymbol{\theta}^{*(t)}, \boldsymbol{\theta}^*) > u$; otherwise, $\boldsymbol{\theta}^{*(t+1)} = \boldsymbol{\theta}^{*(t)}$.

6
7 Again, let m be the burn-in. We need to sample $T + m$ draws using the M-H algorithm. Call the
8 sample $\boldsymbol{\Theta}^{**} = \{\boldsymbol{\theta}^{***(t)}\}$, where $t = 1, \dots, T + m$.

9
10 From this sample, discard the first m draws. Call the remaining sampled draws $\boldsymbol{\Theta}^* = \{\boldsymbol{\theta}^{***(t+m)}\} =$
11 $\{\boldsymbol{\theta}^{*(t)}\}$, where $t = 1, \dots, T + m$.

12
13 As in the univariate case, once we have $\boldsymbol{\Theta}^*$, we can proceed as above.

14
15 Because $\boldsymbol{\theta}$ is multivariate, we actually have a **fourth** option - to use Gibbs sampling.

16
17 Note that Gibbs sampling is particularly useful if we do not know how to sample $\boldsymbol{\theta}|\mathbf{X}$, which is defi-
18 nitely not the case here. We still do it just because it's fun!

19
20 To implement Gibbs sampling, we need to find the *full conditional distributions* $p(\theta_1|\mathbf{X}, \theta_2)$ and
21 $p(\theta_2|\mathbf{X}, \theta_1)$.

22
23 We already know the posterior distribution $p(\boldsymbol{\theta}|\mathbf{X})$, which is $N(\boldsymbol{\mu}_N, \boldsymbol{\Lambda}_N)$. By partitioning $\boldsymbol{\theta}$, we get
24 the following (full) conditional distributions:

$$\begin{aligned} p(\theta_1|\mathbf{X}, \theta_2) &= N\left(\mu_{N1} + \frac{\lambda_{N12}}{\lambda_{N22}}(\theta_2 - \mu_{N2}), \lambda_{N11} - \frac{\lambda_{N12}^2}{\lambda_{N22}}\right) \\ p(\theta_2|\mathbf{X}, \theta_1) &= N\left(\mu_{N2} + \frac{\lambda_{N12}}{\lambda_{N11}}(\theta_1 - \mu_{N1}), \lambda_{N22} - \frac{\lambda_{N12}^2}{\lambda_{N11}}\right) \end{aligned}$$

25 This gives us two univariate distributions, which we know how to sample from directly. If we do not
26 know how to sample directly from one or more of these distributions, we can use A-R or M-H.

27
28 For this example, we will use direct sampling. Because we are using a Markov chain, we need to sample
29 $T + m$ draws from the posterior. We accomplish this by sampling iteratively from the full conditional
30 distributions.

31
32 At iteration $t + 1$:

- 33 • Draw $\theta_1^{(t+1)}$ from $N\left(\mu_{N1} + \frac{\lambda_{N12}}{\lambda_{N22}}(\theta_2^{(t)} - \mu_{N2}), \lambda_{N11} - \frac{\lambda_{N12}^2}{\lambda_{N22}}\right)$
- 34 • Then, draw $\theta_2^{(t+1)}$ from $N\left(\mu_{N2} + \frac{\lambda_{N12}}{\lambda_{N11}}(\theta_1^{(t+1)} - \mu_{N1}), \lambda_{N22} - \frac{\lambda_{N12}^2}{\lambda_{N11}}\right)$
- 35

36 Because we employ direct sampling, all draws are accepted. However, we still need to discard the
37 burn-in from $\boldsymbol{\Theta}^{**}$, and based our inferences on $\boldsymbol{\Theta}^*$ only.

38
39 *MCMC Implementation in Ox for a Multivariate Normal with Unknown Mean Vector*

40

1 The program `multiN.ox` provides a few inferences about unknown $\boldsymbol{\mu}_{2 \times 1}$ in $N(\boldsymbol{\mu}, \boldsymbol{\Lambda})$. The solutions
2 are derived analytically, and using direct sampling of the posterior, M-H algorithm, and Gibbs sampler.
3
4 We use the same headers as `uniN.ox`
5
6 The global variables the multivariate counterparts of the univariate global variables.
7

```
decl N=1000; //sample size
decl mu0=<0;0>,Lambda0=<1,0;0,1>; //prior parameters
decl mu=<.5;-1>,Sigma=<1,.5;.5,1>; //parameters for X|mu
decl T=10000; //size of sample from the posterior
decl m=5000; //burn-in
decl X; //data
```

8 The function `MH_N` generalizes to accommodate the multivariate nature of $\boldsymbol{\mu}$. For example, we get
9 the initial draw from `rann(2,1)`, and `sumr(((X-s0')*choleski(invert(Sigma))).^2)` computes
10 $(\mathbf{x}_i - \boldsymbol{\theta}^{*(t)})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}^{*(t)})$.

```
MH_N(decl s0){//s0 is the current value
decl s1=.1*rann(2,1)+s0;//sampling candidate value from N(s0,.1^2)
decl ll0=sumr(((X-s0')*choleski(invert(Sigma))).^2);
ll0=-.5*sumc(ll0)+(-.5*(s0-mu0)'invert(Lambda0)*(s0-mu0));//post. when par=s0
decl ll1=sumr(((X-s1')*choleski(invert(Sigma))).^2);
ll1=-.5*sumc(ll1)+(-.5*(s1-mu0)'invert(Lambda0)*(s1-mu0));//post. when par=s1
decl alpha=exp(ll1-ll0);//acceptance probability
if(alpha>ranu(1,1)) s0=s1;//decide to accept sampled value s1 or keep s0
return s0;
} //end MH_N
```

11 The function `Gibbs` is new to this program. It first computes the mean and variance of the full
12 conditional $p(\theta_1 | \mathbf{X}, \theta_2)$, before sampling from this distribution. The same this is done for $\theta_2 | \mathbf{X}, \theta_1$.

```
Gibbs(decl s0,M,L){//s0=value, M=post mean, L=post var
decl m=M[0]+L[0][1]/L[1][1]*(s0[1]-M[1]); //mean of par1
decl v=L[0][0]-L[0][1]^2/L[1][1]; //variance of par1
s0[0]=rann(1,1)*sqrt(v)+m; //sampling par1 given par 2
m=M[1]+L[0][1]/L[0][0]*(s0[0]-M[0]); //mean of par2
v=L[1][1]-L[0][1]^2/L[0][0]; //variance of par2
s0[1]=rann(1,1)*sqrt(v)+m; //sampling par2 given par 1
return s0;
} //end Gibbs
```

1 The `main` function is largely similar to the original `main` function except for this time we are making
 2 sure that the dimensions permit matrix operations, and a code is added for the Gibbs sampler solution.

3
 4 Change the values of N , μ_0 , Λ_0 , μ , Σ , T , and m to see how it affects the results.

6 MCMC Estimation of Ability

7
 8 In this example, let N and J be the number of examinees, and test length, respectively, and $P_{ij} =$
 9 $P(X_j = 1|\theta_i)$ be the 2PL. Also, let $p(\theta_i) \sim N(0, 1)$. The goal is to find $p(\theta_i|\mathbf{X}_i)$, for $i = 1, \dots, N$.

10
 11 Specifically, we are interested in $E(\theta_i|\mathbf{X}_i)$ and $V(\theta_i|\mathbf{X}_i)$.

12
 13 To start, we need to find the full conditional distribution for each examinee.

14
 15 The joint posterior is $p(\mathbf{X}, \boldsymbol{\theta})$, and the posterior distribution is

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{X}) &\propto p(\mathbf{X}|\boldsymbol{\theta}) \times p(\boldsymbol{\theta}) \\ &= \prod_{i=1}^N p(\mathbf{X}_i|\theta_i) \times p(\theta_i) \\ &= \prod_{i=1}^N \left[\prod_{j=1}^J P_{ij}^{X_{ij}} (1 - P_{ij})^{1-X_{ij}} \right] \times p(\theta_i) \end{aligned}$$

16 Let

$$\begin{aligned} \mathbf{X}^{(-i)} &= (\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \mathbf{X}_{i+1}, \dots, \mathbf{X}_N) \\ \boldsymbol{\theta}^{(-i)} &= (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_N) \end{aligned}$$

17 The full conditional distribution for θ_i is

$$\begin{aligned} p(\theta_i|\mathbf{X}, \boldsymbol{\theta}^{(-i)}) &\propto p(\mathbf{X}|\theta_i, \boldsymbol{\theta}^{(-i)}) \times p(\theta_i|\boldsymbol{\theta}^{(-i)}) \\ &= p(\mathbf{X}_i|\theta_i) \times p(\mathbf{X}^{(-i)}|\boldsymbol{\theta}^{(-i)}) \times p(\theta_i) \\ &= p(\mathbf{X}_i|\theta_i) \times p(\theta_i) \times c \\ &\propto p(\mathbf{X}_i|\theta_i) \times p(\theta_i) \end{aligned}$$

18 Why can we let $p(\mathbf{X}^{(-i)}|\boldsymbol{\theta}^{(-i)}) = c$?

19
 20 The $p(\theta_i|\mathbf{X}, \boldsymbol{\theta}^{(-i)})$ indicates that, as far as examinee i , is concerned, we only need to deal with \mathbf{X}_i .
 21 This dramatically simplifies the problem because the original N -dimensional problem can be solved
 22 one dimension at a time.

23
 24 However, although $p(\theta_i|\mathbf{X}, \boldsymbol{\theta}^{(-i)})$ is unidimensional, we actually do not know how to sample from this
 25 distribution. Therefore, we will use M-H.

26
 27 To start, at iteration $t + 1$, the candidate value will be drawn as $\theta_i^* \sim N(\theta_i^{(t)}, \sigma_\epsilon^2)$.

28
 29 The acceptance probability is

$$\alpha(\theta_i^{(t)}, \theta_i^*) = \min \left\{ \frac{p(\mathbf{X}_i | \theta_i^*) \times p(\theta_i^*)}{p(\mathbf{X}_i | \theta_i^{(t)}) \times p(\theta_i^{(t)})}, 1 \right\}.$$

1 We set $\theta_i^{(t+1)} = \theta_i^*$ if $\alpha(\theta_i^{(t)}, \theta_i^*) > u$; otherwise $\theta_i^{(t+1)} = \theta_i^{(t)}$.

2
3 **An Important Note:** Although we deal with the examinees one at a time, it does not mean we need
4 to also do the different steps of the M-H algorithm *separately* (i.e., write a loop for the N examinees).
5 For greater efficiency, they can be done *simultaneously*.

6
7 In this example, the M-H code can be written such that we should be able to *simultaneously* sample
8 θ^* , *simultaneously* evaluate the acceptance probabilities, and *simultaneously* determine $\theta^{(t+1)}$, and still
9 ensure that the elements of $\theta^{(t+1)}$ are updated *separately*, rather than *en bloc*.

10
11 Also, for this problem, saving all the draws means we have to store a $T + m \times N$ matrix of draws.
12 This may not be an issue at this point, but other applications that deal with multidimensional data
13 and require structural parameter estimation, it can dramatically slow down the algorithm.

14
15 If we already know what inferences we are interested in, it would suffice to just save the sufficient
16 statistics.

17
18 For the $E(\theta|\mathbf{X})$ and $V(\theta|\mathbf{X})$, the sufficient statistics are

$$\begin{aligned} S_\theta &= \sum_{t=m+1}^{T+m} \theta^{(t)} \\ S_{\theta^2} &= \sum_{t=m+1}^{T+m} \theta^{(t)^2} \end{aligned}$$

19 Given the sufficient statistics,

$$\begin{aligned} \tilde{E}(\theta|\mathbf{X}) &= S_\theta / T \\ \tilde{V}(\theta|\mathbf{X}) &= \frac{S_{\theta^2} - T \times \tilde{E}^2(\theta|\mathbf{X})}{T - 1} \end{aligned}$$

20 In practice, we do not need to get all the $T + m$ draws before computing S_θ and S_{θ^2} . Instead, we can
21 keep a running total of S_θ and S_{θ^2} at each iteration after the burn-in period.

22
23 This mean that instead of a $T + m \times N$ matrix, we need to only keep track of an $N \times 2$ matrix.

24
25 Note that the latter is not a function of T , so the number of iterations does not affect the size of the
26 matrix that needs to be stored in the memory. In some application, this can dramatically affect the
27 efficiency of the algorithm.

28
29 *MCMC Implementation in Ox for a Finding the EAP of the Ability*

30
31 The program `eap_theta.ox` estimates $E(\theta_i|\mathbf{X}_i)$ and $V(\theta_i|\mathbf{X}_i)$ using MCMC, specifically M-H algo-
32 rithm, for the 2PL. In addition to the estimates, it also evaluates how well the parameters are estimated.

33
34 We start with the basic headers.

35

```

#include <oxstd.h>
#include <oxprob.h>
#include <oxfloat.h>

```

1 Next, we define the global variables. These variables are defined globally to make the implementation
2 of the M-H algorithm more convenient.

3

```

decl N=500,J=10; //N=sample size, J=test length
decl T=5000; //size of sample from the posterior
decl m=1000; //burn-in
decl theta,alpha,beta,X; //ability, item parameters,data
decl theta0; //current value of theta
decl s_theta,s_theta2; //sufficient statistics

```

4 The function TwoPL computes the 2PL $N \times J$ probabilities given the θ , α , and β .

5

```

TwoPL(decl t,a,b){
decl tmp=exp(1.7*t*a'-ones(N,1)*(a.*b)');
return(tmp./(1+tmp));
} //TwoPL

```

6 The function MH_draw() updates theta0. Note that drawing θ^* , evaluating the acceptance probabili-
7 ties, and determining which elements of theta0 need to be updated is done simultaneously for all the
8 N examinees. The use of vecindex allows us to selectively update the elements theta0.

```

MH_draw(){
decl theta1=theta0+.25*rann(N,1);
decl tmp=TwoPL(theta0,alpha,beta);
decl ll0=sumr(X.*log(tmp)+(1-X).*log(1-tmp))+log(densn(theta0));
tmp=TwoPL(theta1,alpha,beta);
decl ll1=sumr(X.*log(tmp)+(1-X).*log(1-tmp))+log(densn(theta1));
decl acc=exp(ll1-ll0); //acceptance probability
decl ind=vecindex(acc.>ranu(N,1));
theta0[ind][]=theta1[ind][];
} //end MH_draw

```

1 The function `initiate` makes sure that the variables are set up properly. When multiple chains are
 2 involved, this becomes more crucial. Of course, we may need to add some noise to the initial value of
 3 `theta0`. Otherwise, the different chains will all have the same starting values.

```
initiate(){
  theta0=quann(.025+.95*meanr(X));
  s_theta=0,s_theta2=0;
} //end initiate
```

4 The function `update` keeps a running total of S_θ and S_{θ^2} . It is called in the `main` for iterations
 5 $t = m + 1, \dots, T + m$.

```
update(){
  s_theta=s_theta+theta0;
  s_theta2=s_theta2+theta0.^2;
} //end update
```

7 The first part of the `main` function generates θ , α and β , and simulates the $N \times J$ data matrix \mathbf{X} .

8
 9 It then calls the `initiate` function. The `i` loop draws from the posterior using `MH_draw`, and calls
 10 `update` after the burn-in.

```
initiate();
for(decl i=0;i<T+m;i++){
  MH_draw();
  if(i>=m) update();
} //end i
```

11 The last few lines estimate $\tilde{\theta} = \tilde{E}(\theta|\mathbf{X})$ and $\tilde{V}(\theta|\mathbf{X})$, and evaluates the quality of the estimates by
 12 computing correlation between the true and estimate θ , mean posterior standard deviation, and root
 13 mean squared error:

$$\begin{aligned}\tilde{\rho}_{\theta, \tilde{\theta}} &= \text{Corr}(\theta, \tilde{\theta}) \\ \text{MPSD} &= \sqrt{\sum_{i=1}^N \frac{V(\theta_i|\mathbf{X}_i)}{N}} \\ \text{RMSE} &= \sqrt{\sum_{i=1}^N \frac{(\tilde{\theta}_i - \theta_i)^2}{N}}\end{aligned}$$


```

decl est=s_theta/T;
decl tmp=correlation(theta~est);
print("\nCorr(True,Est):      ",tmp[0][1]);
tmp=(s_theta2-T*est.^2)/(T-1);
print("\nMean Posterior SD: ",sqrt(meanc(tmp))[0]);
tmp=sqrt(meanc((theta-est).^2));
print("\nRMSE:                ",tmp[0],"\n");

```

1 Change N , J , α , and β to see how they affect the quality of the estimates.

3 Additional Exercises

5 1. Monte Carlo Variance

7 Whenever our inferences are based on draws from the posterior, Markov chain or otherwise, *Monte Carlo* (MC) *errors* are involved. This simply means that we can expect the estimates to vary from one sample to another.

11 The amount of variability is called the *MC variance*, which we can denote as σ_{MC}^2 . We can approximate σ_{MC}^2 by drawing S samples, and computing,

$$\sigma_{MC}^2 = \frac{\sum_{s=1}^S \left[\tilde{E}^{(s)}(\theta|\mathbf{X}) - \overline{\tilde{E}(\theta|\mathbf{X})} \right]^2}{S},$$

13 if the variable of interest is $E(\theta|\mathbf{X})$, which is usually the case.

15 Thus, the total uncertainty associated with $\theta|\mathbf{X}$ is the posterior variance plus the Monte Carlo variance (i.e., $V(\theta|\mathbf{X}) + \sigma_{MC}^2$).

18 The Monte Carlo variance can represent a substantial proportion of the total uncertainty if the number of draws in a sample is not sufficiently large.

21 a) Modify `uniN.ox` to allow you to obtain S samples for the same data.

23 b) Find the σ_{MC}^2 associated with the estimates based on direct and indirect samplings procedures.

25 c) Manipulate T and m to see how they affect σ_{MC}^2 for the different sampling procedures.

27 2. Extend `multiN.ox` to make it applicable to the case $D = 3$.

29 3. Comparing Different Methods for Estimating the Ability θ

31 Previously, we learned how to find the MLE $\hat{\theta}$ and EAP $\tilde{\theta}$ through analytic maximization, and the use of quadrature nodes, respectively.

- 1 The current method of finding $\tilde{\theta}$ is based on a Monte Carlo method, specifically, MCMC.
2
3 The standard errors of $\hat{\theta}$ and $\tilde{\theta}$ are computed from the second derivative and the posterior variance,
4 respectively.
5
6 a) Using the correlation between the true and estimated θ , mean standard error, and root mean squared
7 error, compare the three methods of estimating θ and their standard errors.
8
9 b) What happens to the estimates and the standard errors as J increases.
10

Markov Chain Monte Carlo - Part II

An MCMC Algorithm for the 2PL

One advantage of MCMC is the ease by which existing algorithms can be extended to accommodate more complex models.

In this example, we will extend the algorithm for obtaining the EAP of θ in the context of the 2PL to include item parameter estimation.

Again, let N and J be the number of examinees, and test length, respectively, and $P_{ij} = P(X_j = 1|\theta_i)$ be the 2PL.

The goal this time is to make inferences about $p(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}|\mathbf{X})$, where $\boldsymbol{\alpha} = \{\alpha_j\}$, $\boldsymbol{\beta} = \{\beta_j\}$, and $\boldsymbol{\theta} = \{\theta_i\}$.

Prior, Posterior, and Full Conditional Distributions

Let the priors be

$$\begin{aligned} p(\theta_i) &\sim N(0, 1) \\ p(\alpha_j) &\sim \log N(0, 1) \\ p(\beta_j) &\sim N(0, 1) \end{aligned}$$

where $x \sim \log N(\mu, \sigma^2)$ means $\log x \sim N(\mu, \sigma^2)$.

The likelihood of the data given $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$, and $\boldsymbol{\theta}$ is

$$\begin{aligned} L(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}) &= \prod_{i=1}^N L(\mathbf{X}_i|\boldsymbol{\alpha}, \boldsymbol{\beta}, \theta_i) \\ &= \prod_{j=1}^J L(\mathbf{X}_j|\alpha_j, \beta_j, \boldsymbol{\theta}) \\ &= \prod_{i=1}^N \prod_{j=1}^J P_{ij}^{X_{ij}} (1 - P_{ij})^{1-X_{ij}}, \end{aligned}$$

where

$$P_{ij} = P(X = 1|\alpha_j, \beta_j, \theta_i) = \frac{\exp[1.7\alpha_j(\theta_i - \beta_j)]}{1 + \exp[1.7\alpha_j(\theta_i - \beta_j)]}.$$

The joint posterior distribution is

$$\begin{aligned} p(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}|\mathbf{X}) &\propto L(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}) \times p(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}) \\ &= L(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}) \times p(\boldsymbol{\alpha}) \times p(\boldsymbol{\beta}) \times p(\boldsymbol{\theta}) \\ &= L(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}) \times \prod_{j=1}^J p(\alpha_j) \times \prod_{j=1}^J p(\beta_j) \times \prod_{i=1}^N p(\theta_i). \end{aligned}$$

1 To implement Gibbs sampling, we need the full conditional distributions, and for this example, we will
 2 define the following full distributions: $p(\alpha_j, \beta_j | \mathbf{X}, \boldsymbol{\alpha}^{(-j)}, \boldsymbol{\beta}^{(-j)}, \boldsymbol{\theta})$ for $j = 1, \dots, J$, and $p(\theta_i | \mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}^{(-i)})$
 3 for $i = 1, \dots, N$, where

$$\begin{aligned}\boldsymbol{\alpha}^{(-j)} &= \{\alpha_1, \dots, \alpha_{j-1}, \alpha_{j+1}, \dots, \alpha_J\} \\ \boldsymbol{\beta}^{(-j)} &= \{\beta_1, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_J\} \\ \boldsymbol{\theta}^{(-i)} &= \{\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_N\}\end{aligned}$$

4 Based on these formulations, we will update the parameters of item j *en bloc*.

5
 6 The full conditional distribution of $\{\alpha_j, \beta_j\}$ is

$$\begin{aligned}p(\alpha_j, \beta_j | \mathbf{X}, \boldsymbol{\alpha}^{(-j)}, \boldsymbol{\beta}^{(-j)}, \boldsymbol{\theta}) &\propto L(\mathbf{X} | \alpha_j, \beta_j, \boldsymbol{\alpha}^{(-j)}, \boldsymbol{\beta}^{(-j)}, \boldsymbol{\theta}) \times p(\alpha_j, \beta_j | \boldsymbol{\alpha}^{(-j)}, \boldsymbol{\beta}^{(-j)}, \boldsymbol{\theta}) \\ &\propto L(\mathbf{X}_j | \alpha_j, \beta_j, \boldsymbol{\theta}) \times p(\alpha_j, \beta_j) \\ &= L(\mathbf{X}_j | \alpha_j, \beta_j, \boldsymbol{\theta}) \times p(\alpha_j) \times p(\beta_j)\end{aligned}$$

7 What is the second proportionality due to, the first or the second factor?

8
 9 The full conditional distribution of θ_i is

$$\begin{aligned}p(\theta_i | \mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}^{(-i)}) &\propto L(\mathbf{X} | \boldsymbol{\alpha}, \boldsymbol{\beta}, \theta_i, \boldsymbol{\theta}^{(-i)}) \times p(\theta_i | \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}^{(-i)}) \\ &\propto L(\mathbf{X}_i | \boldsymbol{\alpha}, \boldsymbol{\beta}, \theta_i) \times p(\theta_i)\end{aligned}$$

10 Note that neither $p(\alpha_j, \beta_j | \mathbf{X}, \boldsymbol{\alpha}^{(-j)}, \boldsymbol{\beta}^{(-j)}, \boldsymbol{\theta})$ nor $p(\theta_i | \mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\theta}^{(-i)})$ is a known distribution so we
 11 would have to use M-H to obtain draws from the full conditional distributions.

12
 13 At iteration $t + 1$:

14
 15 For $j = 1, \dots, J$, draw α_j^* from $N(\alpha_j^{(t)}, \sigma_\alpha^2)$ and β_j^* from $N(\beta_j^{(t)}, \sigma_\beta^2)$.

16
 The acceptance probability \mathcal{A} associated with the parameters of item j is

$$\mathcal{A}(\{\alpha_j^{(t)}, \beta_j^{(t)}\}, \{\alpha_j^*, \beta_j^*\}) = \min \left\{ \frac{L(\mathbf{X}_j | \alpha_j^*, \beta_j^*, \boldsymbol{\theta}^{(t)}) \times p(\alpha_j^*) \times p(\beta_j^*)}{L(\mathbf{X}_j | \alpha_j^{(t)}, \beta_j^{(t)}, \boldsymbol{\theta}^{(t)}) \times p(\alpha_j^{(t)}) \times p(\beta_j^{(t)})}, 1 \right\}.$$

17 Set $\{\alpha_j^{(t+1)}, \beta_j^{(t+1)}\} = \{\alpha_j^*, \beta_j^*\}$ if $\mathcal{A}(\{\alpha_j^{(t)}, \beta_j^{(t)}\}, \{\alpha_j^*, \beta_j^*\}) > u$; otherwise, set $\{\alpha_j^{(t+1)}, \beta_j^{(t+1)}\} = \{\alpha_j^{(t)}, \beta_j^{(t)}\}$.

18
 19 Then for $i = 1, \dots, N$, draw θ_i^* from $N(\theta_i^{(t)}, \sigma_\theta^2)$.

20
 The acceptance probability \mathcal{A} associated with ability i is

$$\mathcal{A}(\theta_i^{(t)}, \theta_i^*) = \min \left\{ \frac{L(\mathbf{X}_i | \boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\beta}^{(t+1)}, \theta_i^*) \times p(\theta_i^*)}{L(\mathbf{X}_i | \boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\beta}^{(t+1)}, \theta_i^{(t)}) \times p(\theta_i^{(t)})}, 1 \right\}.$$

21 Set $\theta_i^{(t+1)} = \theta_i^*$ if $\mathcal{A}(\theta_i^{(t)}, \theta_i^*) > u$; otherwise, set $\theta_i^{(t+1)} = \theta_i^{(t)}$.

22
 23 Again, although updating of the item parameters will be done one item at a time, we will sample the
 24 candidate values simultaneously to make the algorithm more efficient.

1 Whenever possible, we will sample the draws en bloc even though updating will be done one item or
2 person at a time.

3
4 Note that the sampling process for θ_i is essentially identical to the sampling process when the item
5 parameters are fixed.

6
7 The only difference is that we are now using $\alpha^{(t+1)}$ and $\beta^{(t+1)}$ instead of α and β .
8

9 [A Detour: The Four-Parameter Beta Distribution](#)

The four-parameter beta random variable is a linear transformation of a beta random variable. If $X \sim \text{Beta}(v, \omega)$, for $0 < x < 1$, its density function is

$$f(x) = \frac{x^{v-1}(1-x)^{\omega-1}}{\beta(v, \omega)},$$

10 where $\beta(v, \omega) = \int_0^1 u^{v-1}(1-u)^{\omega-1} \partial u$. The mean and variance of the distribution are

$$\begin{aligned} E(X) &= v/(v + \omega); \\ V(X) &= v\omega/\{(v + \omega)^2(v + \omega + 1)\}. \end{aligned}$$

A random variable $Y = (b - a)X + a$ is distributed as 4-Beta(v, ω, a, b) for $a < b$, and has a density function

$$g(y) = \frac{(y - a)^{v-1}(b - y)^{\omega-1}}{\beta(v, \omega)(b - a)^{v+\omega-1}} \quad \text{for } a < y < b.$$

11 The corresponding mean and variance are

$$\begin{aligned} E(Y) &= (b - a)v/(v + \omega) + a; \\ V(Y) &= (b - a)^2 v\omega/\{(v + \omega)^2(v + \omega + 1)\}. \end{aligned}$$

12 Below are some features of the distribution that makes it attractive.

13
14 First, the support of the four-parameter beta distribution can be altered by adjusting a and b . Hence,
15 the beta distribution is a special case of the 4-parameter beta distribution (i.e., when $a = 0$ and $b = 1$).
16

17 Second, a symmetric or asymmetric distribution can be obtained by varying v and ω . A symmetric
18 distribution is produced when $v = \omega$.
19

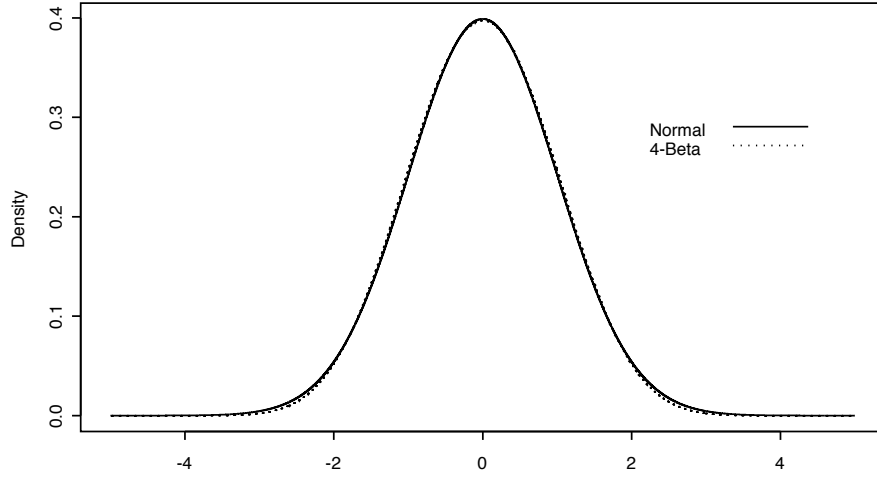
20 A special case is when $v = \omega = 1$; it is the uniform distribution over the interval (a, b) .
21

22 When $v = \omega \approx 12.66$ and $a = -5$ and $b = 5$, the maximum difference between the standard normal
23 and the 4-parameter beta densities over the interval $(-5, 5)$ is about 0.005.
24

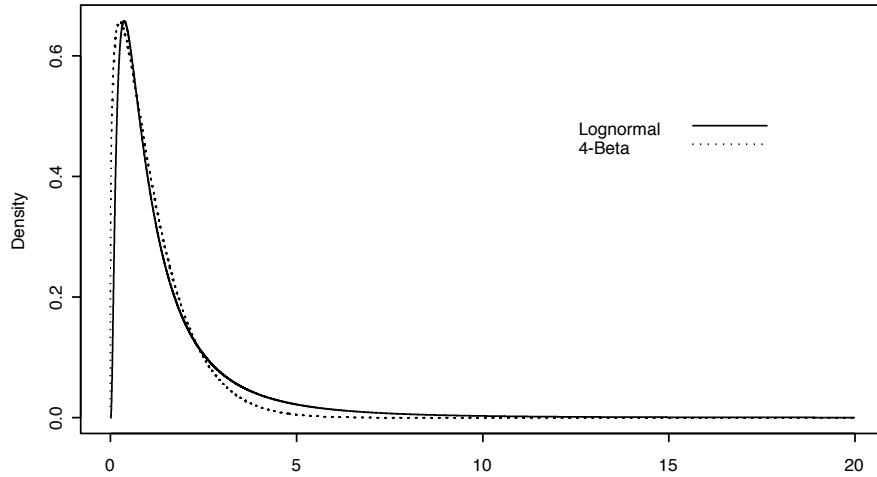
25 When $v < \omega$, the distribution is skewed to the right, and vice versa.
26

27 Finally, the variance of the distribution can be modified by changing the magnitudes of v and ω : larger
28 values correspond to a smaller variances.
29

Normal(0,1) and 4-Beta(-5,5,12.66,12.66) Densities



Lognormal(0,1) and 4-Beta(0,20,1.30,22.00) Densities



1 The figure below shows the 4-Beta approximations to some of the commonly used prior distributions
 2 associated with the item parameters.

3
 4 In using priors with constrained support in the interval (a, b) , we are assured that the parameter esti-
 5 mates will not run away.

6
 7 For example, if instead of $N(0, 1)$, we use the prior $4\text{-Beta}(-5, 5, v, \omega)$ for the difficulty parameter, we
 8 are guaranteed that $|\hat{\beta}_j| < 5$.

9

10 [End of Detour](#)

11 *Implementation in Ox of the MCMC Algorithm for the 2PL*

12

13 The program `mcmc_2PL.ox` is Bayesian estimation program for estimating the item parameters of the
 14 2PL. It is a straightforward extension of `eap_theta.ox`.

15

1 The headers are the same.

2
3 The global variables are mostly the same except we added a few variable related to the item parameters.

4
5 The variables `alpha0` and `beta0` will hold the current values of the item parameters; the variable `draws`
6 will be $T \times 2J$ matrix holding *all* the α and β draws after the burn-in period.

7

```
decl N=1000,J=10; //N=sample size, J=test length
decl T=2000; //size of sample from the posterior
decl m=1000; //burn-in
decl theta,alpha,beta,X; //ability, item parameters, data
decl theta0,alpha0,beta0; //current value of theta, alpha, beta
decl draws,s_theta,s_theta2;//itempar draws,sufficient statistics
```

8 The function `TwoPL` is the same.

9

```
draw_ab(){
decl alpha1=alpha0+.05*rann(J,1);
decl beta1= beta0+.10*rann(J,1);
decl tmp=TwoPL(theta0,alpha0,beta0);
decl ll0=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+
  log(densbeta(alpha0/20,1.3,22))+
  log(densbeta((beta0+5)/10,12.66,12.66));
tmp=TwoPL(theta0,alpha1,beta1);
decl ll1=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+
  log(densbeta(alpha1/20,1.3,22))+
  log(densbeta((beta1+5)/10,12.66,12.66));
decl acc=exp(ll1-ll0);//acceptance probability
decl ind=vecindex(acc.>ranu(J,1));
alpha0[ind]=alpha1[ind];
beta0[ind]= beta1[ind];
} //end draw_ab
```

10 The function `draw_ab` draws the item parameters. Again, we will use a symmetric candidate-generating
11 functions (i.e., normal) to make it easier for us to compute the acceptance probabilities.

12

13 Note that σ_α and σ_β need not be the same.

14

15 In fact, in some cases where convergence rates vary extremely across items, we may actually need to
16 use a different candidate variance for each parameter, as in, we may need σ_{α_j} and σ_{β_j} .

17

1 For the priors of α_j and β_j , we use four-parameter beta distributions instead of the lognormal and
2 normal distributions.
3
4 Note that `acc` is $J \times 1$ so α_j^* and β_j^* will be accepted or rejected at the same time.
5
6 The function `draw_theta` is almost the same to `MH_draw` in `eap_theta.ox`. But instead of `alpha`
7 and `beta`, the function involves `alpha0` and `beta0` to reflect the fact that the item parameters may
8 be different from one iteration to another.
9
10 In addition to the original variable, the function `initiate` also provides starting values for `alpha0`
11 and `beta0`, and sets `draws` as a null matrix.
12
13 In the `update` function will also have `draws=draws|(alpha0'~beta0')`; to keep track of all the item
14 parameter draws after the burn-in.
15

```

evaluate(){
decl est=meanc(draws);
est=shape(est,J,2);
decl se=sqrt(varc(draws));
se=shape(se,J,2);
print("\n      True Alpha   True Beta");
print(alpha~beta);
print("\n      Est. Alpha   Est. Beta   SE(Alpha)   SE(Beta)");
print(est~se);
est=s_theta/T;
decl tmp=correlation(theta~est);
println("\n      Ability Estimates  ");
print("Corr(True,Est):      ",tmp[0][1]);
tmp=(s_theta2-T*est.^2)/(T-1);
print("\nMean Posterior SD: ",sqrt(meanc(tmp))[0]);
tmp=sqrt(meanc((theta-est).^2));
print("\nRMSE:              ",tmp[0],"\n");
} //end evaluate

```

16 Because this program needs to performs more analyses, we need to write a separate function, `evaluate`,
17 to examine the quality of the item parameter and ability estimates.
18
19 In addition to the measures of the quality of the ability estimates, this function also gives the item
20 parameter estimates (i.e., posterior means) and the standard errors (i.e., posterior standard deviations).
21
22 Note that the original dimension of the posterior means and standard deviations is $1 \times 2J$. We need
23 to use the function `shape` to change them to $J \times 2$.
24
25 As before, the `main` function generates the data, calls `initiate`, uses a loop to call `draw_ab` and
26 `draw_theta` $T + m$ times, and updates the variables of interest within the loop after m draws.
27

1 After the draws, the `main` function also calls `evaluate`, and then prints the time it takes to carry out
2 the MCMC algorithm.

3

4 1. Modify T and m to see how they affect the results.

5

6 2. How do N and J affect the time it takes to implement the algorithm?

7

8 An MCMC Algorithm for the HO-DINA Model

9

10 Let N , J , and K be the number of examinees, test length, and number of attributes, respectively.

11

12 We will assume a 1PL the $P(\alpha_k = 1|\theta_i)$, and the DINA model for $P_{ij} = P(X_j = 1|\alpha_i)$.

13 In addition, we also assume $P(\alpha_i|\lambda, \theta_i) = \prod_{k=1}^K P(\alpha_{ik}|\lambda, \theta_i)$.

14

15 The goal this time is to make inferences about $p(\lambda, g, s, \theta, \alpha|\mathbf{X})$, where $\lambda = \{\lambda_1, \lambda_0\} = \{\lambda_1, \lambda_{01}, \dots, \lambda_{0K}\}$,
16 $g = \{g_j\}$, $s = \{s_j\}$, $\theta = \{\theta_i\}$, and $\alpha = \{\alpha_i\}$.

17

18 *Prior, Posterior, and Full Conditional Distributions*

19

20 Let the priors be λ , θ , α , g and s :

$$\begin{aligned} p(\lambda_1) &\sim 4 - \text{Beta}(0, 16, 1.5, 21.5) \\ p(\lambda_{0k}) &\sim N(-1, 1) \\ p(\theta_i) &\sim 4 - \text{Beta}(-4, 4, 8.21, 8.21) \\ p(\alpha_{ik} = 1|\lambda, \theta_i) &\sim \text{Ber}\left(\{1 + \exp[-1.7\lambda_1(\theta_i - \lambda_{0k})]\}^{-1}\right) \\ p(g_j) &\sim 4 - \text{Beta}(0, .9, 1.5, 2) \\ p(1 - s_j) &\sim 4 - \text{Beta}(.1, 1, 2, 1.5) \end{aligned}$$

21 The likelihood of the data given, g , s , and α is

$$\begin{aligned} L(\mathbf{X}|g, s, \alpha) &= \prod_{i=1}^N L(\mathbf{X}_i|g, s, \alpha_i) \\ &= \prod_{j=1}^J L(\mathbf{X}_j|g_j, s_j, \alpha) \\ &= \prod_{i=1}^N \prod_{j=1}^J \left[s_j^{1-X_{ij}} (1 - s_j)^{X_{ij}} \right]^{\eta_{ij}} \left[g_j^{X_{ij}} (1 - g_j)^{1-X_{ij}} \right]^{1-\eta_{ij}}, \end{aligned}$$

22 where $\eta_{ij} = I(\alpha'_i q_j = q'_j q_j)$.

23

24 The joint posterior distribution is

$$p(\lambda, g, s, \theta, \alpha|\mathbf{X}) \propto L(\mathbf{X}|g, s, \alpha) \times p(\alpha|\lambda, \theta) \times p(\theta) \times p(\lambda) \times p(g) \times p(s).$$

1 To implement Gibbs sampling, we need the following full conditional distributions:

$$\begin{aligned}
p(\lambda_1|\mathbf{X}, \boldsymbol{\lambda}_0, \boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{g}, \mathbf{s}) &\propto p(\boldsymbol{\alpha}|\boldsymbol{\lambda}, \boldsymbol{\theta}) \times p(\lambda_1) \\
p(\lambda_{0k}|\mathbf{X}, \lambda_1, \boldsymbol{\lambda}_0^{(-k)}, \boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{g}, \mathbf{s}) &\propto p(\boldsymbol{\alpha}_k|\lambda_1, \lambda_{0k}, \boldsymbol{\theta}) \times p(\lambda_{0k}) \\
p(\theta_i|\mathbf{X}, \boldsymbol{\lambda}, \boldsymbol{\theta}^{(-i)}, \boldsymbol{\alpha}, \mathbf{g}, \mathbf{s}) &\propto p(\boldsymbol{\alpha}_i|\boldsymbol{\lambda}, \boldsymbol{\theta}_i) \times p(\theta_i) \\
p(\boldsymbol{\alpha}_i|\mathbf{X}, \boldsymbol{\lambda}, \boldsymbol{\theta}, \boldsymbol{\alpha}^{(-i)}, \mathbf{g}, \mathbf{s}) &\propto L(\mathbf{X}_i|\mathbf{g}, \mathbf{s}, \boldsymbol{\alpha}_i) \times p(\boldsymbol{\alpha}_i|\boldsymbol{\lambda}, \boldsymbol{\theta}_i) \\
p(g_j|\mathbf{X}, \boldsymbol{\lambda}, \boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{g}^{(-j)}, \mathbf{s}) &\propto L(\mathbf{X}_j|g_j, s_j, \boldsymbol{\alpha}) \times p(g_j) \\
p(s_j|\mathbf{X}, \boldsymbol{\lambda}, \boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{g}, \mathbf{s}^{(-j)}) &\propto L(\mathbf{X}_j|g_j, s_j, \boldsymbol{\alpha}) \times p(s_j)
\end{aligned}$$

2 **Verify that the above full conditional distributions are correct.**

3

4 Based on these formulations, we will update all the parameters *one at a time*, including the item
5 parameters.

6

7 None of the full conditional distribution is known so we have to use the M-H algorithm.

8

9 At iteration $t + 1$:

10

1. Draw λ_1^* from $N(\lambda_1^{(t)}, \sigma_{\lambda_1}^2)$, and accept λ_1^* with probability

$$\mathcal{A}(\lambda_1^{(t)}, \lambda_1^*) = \min \left\{ \frac{p(\boldsymbol{\alpha}^{(t)}|\lambda_1^*, \boldsymbol{\lambda}_0^{(t)}, \boldsymbol{\theta}^{(t)}) \times p(\lambda_1^*)}{p(\boldsymbol{\alpha}^{(t)}|\lambda_1^{(t)}, \boldsymbol{\lambda}_0^{(t)}, \boldsymbol{\theta}^{(t)}) \times p(\lambda_1^{(t)})}, 1 \right\}.$$

2. For $k = 1, \dots, K$, draw λ_{0k}^* from $N(\lambda_{0k}^{(t)}, \sigma_{\lambda_0}^2)$, and accept λ_{0k}^* with probability

$$\mathcal{A}(\lambda_{0k}^{(t)}, \lambda_{0k}^*) = \min \left\{ \frac{p(\boldsymbol{\alpha}^{(t)}|\lambda_1^{(t+1)}, \lambda_{0k}^*, \boldsymbol{\lambda}_0^{(-j)(t)}, \boldsymbol{\theta}^{(t)}) \times p(\lambda_{0k}^*)}{p(\boldsymbol{\alpha}^{(t)}|\lambda_1^{(t+1)}, \lambda_{0k}^{(t)}, \boldsymbol{\lambda}_0^{(-j)(t)}, \boldsymbol{\theta}^{(t)}) \times p(\lambda_{0k}^{(t)})}, 1 \right\}.$$

3. For $i = 1, \dots, N$, draw θ_i^* from $N(\theta_i^{(t)}, \sigma_{\theta}^2)$, and accept θ_i^* with probability

$$\mathcal{A}(\theta_i^{(t)}, \theta_i^*) = \min \left\{ \frac{p(\boldsymbol{\alpha}_i^{(t)}|\boldsymbol{\lambda}^{(t+1)}, \boldsymbol{\theta}^*) \times p(\theta_i^*)}{p(\boldsymbol{\alpha}_i^{(t)}|\boldsymbol{\lambda}^{(t+1)}, \boldsymbol{\theta}^{(t)}) \times p(\theta_i^{(t)})}, 1 \right\}.$$

4. For $i = 1, \dots, N$ and $k = 1, \dots, K$, draw α_{ik}^* from $\text{Ber}(.5)$, and accept α_{ik}^* with probability

$$\mathcal{A}(\alpha_{ik}^{(t)}, \alpha_{ik}^*) = \min \left\{ \frac{L(\mathbf{X}_i|\mathbf{g}^{(t)}, \mathbf{s}^{(t)}, \boldsymbol{\alpha}^*) \times p(\alpha_{ik}^*|\boldsymbol{\lambda}^{(t+1)}, \boldsymbol{\theta}^{(t+1)})}{L(\mathbf{X}_i|\mathbf{g}^{(t)}, \mathbf{s}^{(t)}, \boldsymbol{\alpha}^{(t)}) \times p(\alpha_{ik}^{(t)}|\boldsymbol{\lambda}^{(t+1)}, \boldsymbol{\theta}^{(t+1)})}, 1 \right\}.$$

5. For $j = 1, \dots, J$, draw g_j^* from $N(g_j^{(t)}, \sigma_g^2)$, and accept g_j^* with probability

$$\mathcal{A}(g_j^{(t)}, g_j^*) = \min \left\{ \frac{L(\mathbf{X}_j|g_j^*, s_j^{(t)}, \boldsymbol{\alpha}^{(t+1)}) \times p(g_j^*)}{L(\mathbf{X}_j|g_j^{(t)}, s_j^{(t)}, \boldsymbol{\alpha}^{(t+1)}) \times p(g_j^{(t)})}, 1 \right\}.$$

6. For $j = 1, \dots, J$, draw s_j^* from $N(s_j^{(t)}, \sigma_s^2)$ and accept s_j^* with probability

$$\mathcal{A}(s_j^{(t)}, s_j^*) = \min \left\{ \frac{L(\mathbf{X}_j|g_j^{(t+1)}, s_j^*, \boldsymbol{\alpha}^{(t+1)}) \times p(s_j^*)}{L(\mathbf{X}_j|g_j^{(t+1)}, s_j^{(t)}, \boldsymbol{\alpha}^{(t+1)}) \times p(s_j^{(t+1)})}, 1 \right\}.$$

A Detour: The Gelman-Rubin Convergence Statistic R

In addition to graphical method, which can be unreliable, we can also compute the Gelman-Rubin statistic to diagnose the convergence of the chain.

The statistic, which involves multiple parallel chains with dispersed starting values to estimate the portion of the posterior mean estimator that is due to the Monte Carlo error, is used to assess the convergence of scalar parameters.

This means that even if θ is multivariate, we assess the convergence of the chain element by element.

Let θ be the scalar parameter, and let $\theta_c^{(t)}$ denote the t^{th} draw of the c^{th} chain, where $c = 1, \dots, C$.

We define the between- and within-chain variances as

$$B = \frac{T}{C-1} \sum_{c=1}^C (\bar{\theta}_c - \bar{\theta})^2$$

$$W = \frac{1}{C} \sum_{c=1}^C S_c^2,$$

where

$$\bar{\theta}_c = \frac{1}{T} \sum_{t=1}^T \theta_c^{(t)}$$

$$\bar{\theta} = \frac{1}{C} \sum_{c=1}^C \bar{\theta}_c$$

$$S_c^2 = \frac{1}{T-1} \sum_{t=1}^T (\theta_c^{(t)} - \bar{\theta}_c)^2$$

Note that multiple chains are needed to be able to compute B .

Note also that $\sqrt{\hat{R}}$ is computed using means and variances of the draws across the different chains.

This means that we can monitor convergence by simply computing the sufficient statistics $s_c = \sum_{t=1}^T \theta_c^{(t)}$ and $ss_c = \sum_{t=1}^T \theta_c^{(t)2}$, without having to save to the draws.

The marginal posterior variance of the parameter, $V(\theta|\mathbf{X})$, can be estimated by computing the weighted average of W and B as follows:

$$\hat{V}^*(\theta|\mathbf{X}) = \frac{T-1}{T} W + \frac{1}{T} B.$$

If the chain has not reached the target distribution, $\bar{\theta}_c$ from the different chains will be different from each other resulting in $\hat{V}^*(\theta|\mathbf{X})$ overestimating $V(\theta|\mathbf{X})$.

On the other hand, if the chain has reached its target distribution, $\hat{V}^*(\theta|\mathbf{X}) \approx W$ and $B \approx 0$.

This means that $\hat{V}^*(\theta|\mathbf{X}) \rightarrow V(\theta|\mathbf{X})$ as $n \rightarrow \infty$.

1 The *potential scale reduction factor* is estimated by

$$\sqrt{\hat{R}} = \sqrt{\frac{\hat{V}^*(\theta|\mathbf{X})}{W}},$$

2 which approaches 1 when $T \rightarrow \infty$.

3
4 As a rule of thumb, $\sqrt{\hat{R}} < 1.2$ is considered acceptable. However, in some applications, something
5 smaller (e.g., 1.1), might be preferable.

6
7 If the chains have converged, all the draws after the burn-in are from the same target distribution.

8
9 Therefore, we can combine all the draws, and base our inferences on the pooled draws, instead of the
10 draws from a single chain.

11
12 With C chains each with T draws after burn-in, we will have $T \times C$ instead of T draws to work with.
13 With more draws, the Monte Carlo error can be made smaller.

14
15 For the us to be able to conclude that convergence is reached for entire distribution, potential scale
16 reduction factor must be small for all the parameters of interest.

17
18 In latent variable modeling, we typically limit this to the structural parameters.

19
20 For example, in the HO-DINA model, we monitor the convergence of $\boldsymbol{\lambda}$, \mathbf{g} , and \mathbf{s} .

21
22 For us to conclude that convergence has been reached for entire distribution, we only need to make
23 sure that $\sqrt{\hat{R}_{\max}}$ is sufficiently small.

24
25 It should be noted that, for some researchers, convergence of single parameters is not sufficient to
26 establish *joint* convergence of the parameters.

27
28 To monitor joint convergence, the Brooks-Gelman *multivariate potential scale reduction factor* may be
29 used.

30
31 However, in some applications, the number of parameters involved in latent variable models can be
32 extremely large. This can make computing of the multivariate potential scale reduction factor difficult.

33
34 An alternative is to monitor the convergence of the parameters by block. For example, the multivariate
35 potential scale reduction factor can be computed for $\boldsymbol{\lambda}$, and $\{g_j, s_j\}$ separately.

36
37 **End of Detour**

38 *Implementation in Ox of the MCMC Algorithm for the HO-DINA Model*

39
40 The program `mcmc_hodina.ox` implements the MCMC algorithm for the HO-DINA model outlined
41 above. In addition to parameter estimation, the program also runs multiple (i.e., `Nchain`) chains to
42 assess convergence. The program stops when `max_sqrtR` (i.e., $\sqrt{\hat{R}_{\max}}$) is less than `crit_sqrtR`, or the
43 number of draws (per chain) is equal to `max_it`.

44

1 The minimum number of draws is set to `min_it`, of which `m` draws are discarded as burn-in.
2
3 After the initial `min_it` draws, convergence is assessed. If convergence is achieved, sampling is terminated; if not, additional `inc` draws are obtained before convergence is assessed again.
4
5
6 We declare the sufficient statistics needed for monitoring convergence. The current number of variables is for `Nchain=4`. We also declare `ndraws` to keep track of the total number of iterations after burn-in - the program can have different stopping points. This number can reflect the thinned value.
7
8
9
10 We only declare two types of sufficient statistics: one that is chain-specific, and another that cuts across chain.
11
12 For the former type, we have the sufficient statistics for the structural parameters λ , g , and s , which are also used to monitor convergence.
13
14 For the latter type, we have the sufficient statistics for the incidental variables θ and α .

```
decl N=536,J=15,K=5;//sample size, test length, # of attributes
decl files={"qmatrix.txt","X536.dat"};//data and Q-matrix filenames
//min and max chain lengths, burn-in, increment; # of it->not fixed
decl min_it=2000,max_it=50000,m=1000,inc=1000;
decl crit_sqrtR=1.2,tot_it;//conv. criterion, tot_it per chain
decl Q,X;//Q-matrix and data variables
decl theta0,lambda0,gs0,alpha0;//variables for current values
//Candidate variances
decl cand_theta=.25,cgs0=.02,cgs1=.025,clambda0=.1,clambda1=.05;
decl Nchain=4,max_sqrtR;//number of chains, max of the sqrtR
//Sufficient stats for monitoring convergence
decl c1s,c1s2,c2s,c2s2,c3s,c3s2,c4s,c4s2,ndraws;
//Sufficient stats for incidental parameters
decl stheta,sstheta,salpha,ssalpha;
```

15 The function `OnePL` computes the $N \times K$ probabilities based on the arguments `t` and `l`.
16

```
OnePL(decl t,l){
return ones(N,K)./(1+exp(-1.7*1[0][0]*(t*ones(1,K)-ones(N,1)*1[] [1]')));
}end OnePL
```

17 The function `draw_theta` updates the higher-order latent trait θ . Note that for this program, by using
18 a 4-Beta($-4, 4, 8.21, 8.21$), we are effectively constraining θ_i to be in the interval $[-4, 4]$.
19

20 Again, θ is sampled en bloc, but updated individually.
21

22 Note that `Ox` does not have a function for the four-parameter beta distribution so Y needs to be
23 transformed $(Y - a)/(b - a)$ so it can be evaluated using the usual beta distribution function.
24

```

draw_theta(){
decl theta1=cand_theta*rann(N,1)+theta0;
decl p=OnePL(theta0,lambda0);
decl d0=sumr(alpha0.*log(p)+(1-alpha0).*log(1-p))+
    log(densbeta((theta0+4)/8,8.21,8.21));
p=OnePL(theta1,lambda0);
decl d1=sumr(alpha0.*log(p)+(1-alpha0).*log(1-p))+
    log(densbeta((theta1+4)/8,8.21,8.21));
decl acc=d1-d0;
decl accind=vecindex(ranu(N,1).<exp(acc));
theta0[accind]=theta1[accind];
} //end draw_theta

```

- 1 The function `draw_lambda` updates the higher-order structural parameter λ . It consists of two steps -
- 2 one for updating λ_0 , and one for updating λ_1 .
- 3
- 4 For λ_0 , we need to sum across the rows; for λ_1 , we need to sum across the rows and the columns.

```

draw_lambda(){
//Step for the difficulty parameters
decl lambda1=lambda0[] [0]~(clambda0*rann(K,1)+lambda0[] [1]);
decl d0=sumc(alpha0.*log(OnePL(theta0,lambda0)))+
    (1-alpha0).*log(1-OnePL(theta0,lambda0)))+log(densn(lambda0[] [1]'+1));
decl d1=sumc(alpha0.*log(OnePL(theta0,lambda1)))+
    (1-alpha0).*log(1-OnePL(theta0,lambda1)))+log(densn(lambda1[] [1]'+1));
decl acc=d1-d0;
decl accind=vecindex(ranu(1,K).<exp(acc));
lambda0[accind]=lambda1[accind];
//Step for the common discrimination parameter
lambda1=(clambda1*rann(1,1)+lambda0[] [0])~lambda0[] [1];
d0=sumr(sumc(alpha0.*log(OnePL(theta0,lambda0)))+
    (1-alpha0).*log(1-OnePL(theta0,lambda0)))+
    +log(densbeta(lambda0[0] [0]/16,1.5,21.5));
d1=sumr(sumc(alpha0.*log(OnePL(theta0,lambda1)))+
    (1-alpha0).*log(1-OnePL(theta0,lambda1)))+
    +log(densbeta(lambda1[0] [0]/16,1.5,21.5));
acc=d1-d0;
if(ranu(1,1)<exp(acc)) lambda0[] [0]=lambda1[] [0];
} //end draw_lambda

```

- 5 Note that because we are using a common discrimination parameter, we need to make sure that, at
- 6 each iteration, only one value is sampled, and its prior density is counted only once.
- 7

- 1 The function `draw_alpha` updates α . The candidate values are sampled as α_{ik} , but accepted or rejected
2 as α_i .

```
draw_alpha(){
decl alpha1=ranbinomial(N,K,1,.5);
decl tmp=alpha0*Q';
tmp=tmp.*(ones(N,1)*sumc(Q'));
tmp=tmp.*(ones(N,1)*gs0[][1]')+(1-tmp).*(ones(N,1)*gs0[][0]');
decl d0=sumr(X.*log(tmp)+(1-X).*log(1-tmp))
      +sumr(alpha0.*log(OnePL(theta0,lambda0))+
      (1-alpha0).*log(1-OnePL(theta0,lambda0)));
tmp=alpha1*Q';
tmp=tmp.*(ones(N,1)*sumc(Q'));
tmp=tmp.*(ones(N,1)*gs0[][1]')+(1-tmp).*(ones(N,1)*gs0[][0]');
decl d1=sumr(X.*log(tmp)+(1-X).*log(1-tmp))
      +sumr(alpha1.*log(OnePL(theta0,lambda0))+
      (1-alpha1).*log(1-OnePL(theta0,lambda0)));
decl acc=d1-d0;
decl accind=vecindex(ranu(N,1).<exp(acc));
alpha0[accind][]=alpha1[accind][];
} //end draw_alpha
```

```
draw_gs(){
decl tmp0=alpha0*Q';
tmp0=tmp0.*(ones(N,1)*sumc(Q'));
decl gs1=(cgs0*rann(J,1)+gs0[][0])~gs0[][1];
decl tmp=tmp0.*(ones(N,1)*gs0[][1]')+(1-tmp0).*(ones(N,1)*gs0[][0]');
decl d0=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+log(densbeta(gs0[][0]/.9,1.5,2));
tmp=tmp0.*(ones(N,1)*gs1[][1]')+(1-tmp0).*(ones(N,1)*gs1[][0]');
decl d1=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+log(densbeta(gs1[][0]/.9,1.5,2));
decl acc=d1-d0;
decl accind=vecindex(ranu(J,1).<exp(acc));
gs0[accind][]=gs1[accind][];
gs1=gs0[][0]~(cgs1*rann(J,1)+gs0[][1]);
tmp=tmp0.*(ones(N,1)*gs0[][1]')+(1-tmp0).*(ones(N,1)*gs0[][0]');
d0=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+log(densbeta((gs0[][1]-.1)/.9,2,1.5));
tmp=tmp0.*(ones(N,1)*gs1[][1]')+(1-tmp0).*(ones(N,1)*gs1[][0]');
d1=sumc(X.*log(tmp)+(1-X).*log(1-tmp))'+log(densbeta((gs1[][1]-.1)/.9,2,1.5));
acc=d1-d0;
accind=vecindex(ranu(J,1).<exp(acc));
gs0[accind][]=gs1[accind][];
} //end draw_gs
```

- 3 The function `draw_gs` updates the item parameters. Note that the algorithm estimates $1-s_j$, rather s_j .

- 1
2 The function `initiate` is for initializing the current states and sufficient statistics. This function is
3 called C times, once for each chain.

```
initiate(decl nc){
  if(nc==0) {c1s=c1s2=c2s=c2s2=c3s=c3s2=c4s=c4s2=ndraws=0;
             max_sqrtR=2;}
  theta0=rann(N,1)*.05+standardize(meanr(X));
  decl tmp=(-3*ranu(K,1)+1);
  lambda0=(1+2*ranu(1,1)*ones(K,1))~tmp;
  alpha0=ranbinomial(N,K,1,.5);
  gs0=ranu(J,1)*.5~.5+ranu(J,1)*.5;
  stheta=sstheta=0;
  salpha=ssalpha=0;
} //end initiate
```

- 4 The function `draw_dina` is called to carry out the different steps of the M-H algorithm.
5
6 The function `save` updates the sufficient statistics needed to estimate \sqrt{R} , α , and θ .

```
save(decl nc){
  if(nc==0){
    c1s=c1s+(lambda0[0][0]~lambda0[][1]~vec(gs0)');
    c1s2=c1s2+(lambda0[0][0]~lambda0[][1]~vec(gs0)')^2;
    ndraws=ndraws+1;
  } //end nc==0
  if(nc==1){
    c2s=c2s+(lambda0[0][0]~lambda0[][1]~vec(gs0)');
    c2s2=c2s2+(lambda0[0][0]~lambda0[][1]~vec(gs0)')^2;
  } //end nc==1
  if(nc==2){
    c3s=c3s+(lambda0[0][0]~lambda0[][1]~vec(gs0)');
    c3s2=c3s2+(lambda0[0][0]~lambda0[][1]~vec(gs0)')^2;
  } //end nc==2
  if(nc==3){
    c4s=c4s+(lambda0[0][0]~lambda0[][1]~vec(gs0)');
    c4s2=c4s2+(lambda0[0][0]~lambda0[][1]~vec(gs0)')^2;
  } //end nc==3
  stheta=stheta+theta0; sstheta=sstheta+theta0.^2;
  salpha=salpha+alpha0; ssalpha=ssalpha+alpha0.^2;
} //end save
```

- 7 When the first chain is involved, `ndraws` is also incremented.
8

- 1 The function `compute_sqrtR` computes $\sqrt{\hat{R}_{\max}}$. It first computes the $\sqrt{\hat{R}}$ for each of the structural
 2 parameters, and picks out the maximum of these convergence statistics.

```
compute_sqrtR(){
  decl B,W,V;
  decl m=(c1s|c2s|c3s|c4s)/ndraws;
  B=sumc(m.^2)-sumc(m).^2/Nchain;
  B=ndraws/(Nchain-1)*B;
  W=c1s2-c1s.^2/ndraws;
  W=W+c2s2-c2s.^2/ndraws;
  W=W+c3s2-c3s.^2/ndraws;
  W=W+c4s2-c4s.^2/ndraws;
  W=W/(Nchain*(ndraws-1));
  V=(1-1/ndraws)*W+1/ndraws*B;
  max_sqrtR=maxc(sqrt(V./W)');
} //end compute_sqrtR
```

```
print_est(const et){
  format(1000);
  decl tmp=sqrt((sstheta-(stheta.^2)/(Nchain*ndraws))/(Nchain*ndraws-1));
  decl par=sprint("theta.out");
  decl file=fopen(par,"w");
  fprintf(file,"%#M",stheta/(Nchain*ndraws)~tmp);
  fclose(file);
  par=sprint("alpha.out");
  file=fopen(par,"w");
  fprintf(file,"%#M",salp/(Nchain*ndraws).>.5);
  fclose(file);
  decl m=c1s+c2s+c3s+c4s;
  decl se=c1s2+c2s2+c3s2+c4s2;
  se=sqrt((se-m.^2/(Nchain*ndraws))/(Nchain*ndraws-1));;
  m=m/(Nchain*ndraws);
  print("\n","***** Higher-Order Structural Parameters *****","\n");
  print("Common Discrimination", "\n", "      Est.      SE", "%10.4f", m[0]~se[0]);
  print("Intercepts", "\n", "      Att.      Est.      SE", "%10.4f",
        range(1,K)'~m[1:K]'~se[1:K]');
  print("\n","***** Item Parameters *****","\n");
  print("      Item      g      1-s      SE(g)      SE(1-s)", "%10.4f",
        range(1,J)'~shape(m[] [K+1:2*J+K], J, 2)~shape(se[] [K+1:2*J+K], J, 2));
  print("\n", "Elapsed Time: ", timespan(et), "\n");
} //print_est
```

- 3 The function `print_est` computes the estimates of θ , α , λ , g , and $1 - s$.

4

1 The estimates of the incidental parameters are printed to file, whereas the estimates of the structural
 2 parameters are printed as screen output.
 3
 4 The main function has a few components.
 5
 6 The first component records the starting time, and counts the total number of structural and incidental
 7 parameters.

```
ranseed(32671);
print("Starting Time: ",time(),"\n");
decl et=timer();
decl npar=(N+1)*(K+1)+2*J;//total number of parameters
```

8 The second component reads the Q-matrix and data files.

9
 10 The next component initiates the variables and carries out the initial number of draws for each of the
 11 C chains.

12
 13 After going through all the chains, $\sqrt{\hat{R}}$ is computed, and total number of iteration is updated.

```
for(decl nc=0;nc<Nchain;nc++){
  initiate(nc);
  for(decl j=0;j<min_it;j++){
    draw_dina();
    if(j>=m && imod(j,10)==0) save(nc);//saving and thinning suff. stats
  }//end sampling
  decl par=sprint("tmp_",nc,".out");//saving most recent values
  decl file=fopen(par,"w");
  fprintf(file,"%#M",lambda0[0][0]~lambda0[][1]~vec(gs0)~theta0~vec(alpha0)');
  fclose(file);
}//end nc

compute_sqrtR();
tot_it=min_it;
print("\n","      Iteration      Max SqrtR");
print(tot_it~max_sqrtR);
```

14 When $\sqrt{\hat{R}}$ is greater than `crit_sqrtR` after the initial number of draws, additional draws are sampled.

15
 16 After `inc` additional draws, `crit_sqrtR` is computed again to examine convergence.

17
 18 The sampling process continues until $\sqrt{\hat{R}}$ is less than `crit_sqrtR`, or the maximum number of itera-
 19 tions is reached.

20

1 Note that we are thinning the chains by saving only every tenth draw. The function `imod` helps us to
2 this.

```
while (max_sqrtR>crit_sqrtR && tot_it<=max_it){
  for(decl nc=0;nc<Nchain;nc++){
    decl tmp1;
    tmp0=sprint("tmp_",nc,".out");
    file0=fopen(tmp0);
    fscan(file0,"%#M",1,npar,&tmp1);
    fclose(file0);

    lambda0[][0]=tmp1[0];
    lambda0[][1]=tmp1[1:K]';
    gs0=shape(tmp1[K+1:2*J+K],J,2);
    theta0=tmp1[2*J+K+1:2*J+K+N]';
    alpha0=shape(tmp1[2*J+K+N+1:npar-1],N,K);

    for(decl j=0;j<inc;j++){
      draw_dina();
      if(imod(j,10)==0) save(nc);
    }//end j

    decl par=sprint("tmp_",nc,".out");
    decl file=fopen(par,"w");
    fprintf(file,"%#M",lambda0[0][0]~lambda0[][1]~vec(gs0)'~theta0'~vec(alpha0)');
    fclose(file);
  }//end additional draws

  compute_sqrtR();
  tot_it=tot_it+inc;
  print(tot_it~max_sqrtR);
} //end while
```

3 Note that the chains are restarted right where they left off.

4
5 By saving the most recent values after each sampling chunk, we do not need to restart the chains if we
6 they have not converged about a particular number of iterations.

7
8 Finally, when the chains have converged, or maximum number of iterations is reach, the function
9 `print_est` is called.

10
11 Modify `min_it`, `max_it`, `m`, `inc`, and `crit_sqrtR` to see how they affect the time needed to get the
12 estimates.

13
14 How do these changes affect the estimates and standard errors?

15

1 **Additional Exercises**

- 2
- 3 1. Compare the MML and MCMC estimates of the 2PL in terms of computational speed and quality
- 4 of the estimates.
- 5
- 6 2. Updating α_j and β_j Separately
- 7
- 8 Modify `mcmc_2PL.ox` so that α_j and β_j are updated separately. This will require one step for sampling
- 9 and updating α_j , and another step for sampling and updating β_j . However, the two steps can be
- 10 carried out in the same function.
- 11
- 12 How does updating α_j and β_j separately affect the computational efficiency of the algorithm and qual-
- 13 ity of the parameter estimates?
- 14
- 15 3. Compare the MML and MCMC estimates of the HO-DINA model.
- 16
- 17 4. Modify `mcmc_hodina.ox` to use five parallel chains to evaluate convergence.
- 18

Markov Chain Monte Carlo - Part III

An MCMC Algorithm for Multi-Unidimensional Ability Estimation

Most assessments involve batteries of tests that measure highly correlated abilities. However, it is also not uncommon in practice for these tests to be scored *one test at a time*.

In so doing, only actual responses from a specific test are taken into consideration in determining examinees' scores.

In doing so, the correlation structure of the abilities are ignored.

Because available information relevant to examinees' standings on the tests is not incorporated, this method of ability estimation is deemed suboptimal.

To improve efficiency, we can estimate the abilities *simultaneously*.

For this, we would need a multidimensional model.

The extension of the 3PL model to the multidimensional context as follows:

$$P(X_{ij} = 1 | \boldsymbol{\theta}_i, \boldsymbol{\alpha}_j, \beta_j, \gamma_j) = \gamma_j + (1 - \gamma_j) \frac{\exp(\boldsymbol{\alpha}'_j \boldsymbol{\theta}_i + \beta_j)}{1 + \exp(\boldsymbol{\alpha}'_j \boldsymbol{\theta}_i + \beta_j)}$$

where

- $P(X_{ij} = 1 | \boldsymbol{\theta}_i, \boldsymbol{\alpha}_j, \beta_j, \gamma_j)$ is the probability of examinee i responding to item j correctly;
- X_{ij} is the response of examinee i to item j (0 = incorrect, 1 = correct);
- $\boldsymbol{\theta}_i$ is the ability vector of the examinee;
- $\boldsymbol{\alpha}_j$ is the vector of item parameters related to the discrimination of the item;
- β_j is the parameter related to the difficulty of the item;
- γ_j is the pseudo-guessing parameter of the item;
- $i = 1, \dots, N$ (the total number of examinees); and
- $j = 1, \dots, J$ (the total number of items).

For the purposes of this example, we will assume a simple structure is (i.e., each item measures one dimension of ability, and thus, $\boldsymbol{\alpha}_j$ contains only one non-zero element).

Individually, the tests are unidimensional; *collectively*, they are multidimensional. Thus, what we have is a *multi-unidimensional* tests.

Under this assumption, the multidimensional 3PL model can be reexpressed as:

$$P_{ij(d)} = P(X_{j(d)} = 1 | \theta_{i(d)}, \alpha_{j(d)}, \beta_{j(d)}, \gamma_{j(d)}) = \gamma_{j(d)} + (1 - \gamma_{j(d)}) \frac{\exp(\alpha_{j(d)}\theta_{i(d)} + \beta_{j(d)})}{1 + \exp(\alpha_{j(d)}\theta_{i(d)} + \beta_{j(d)})}$$

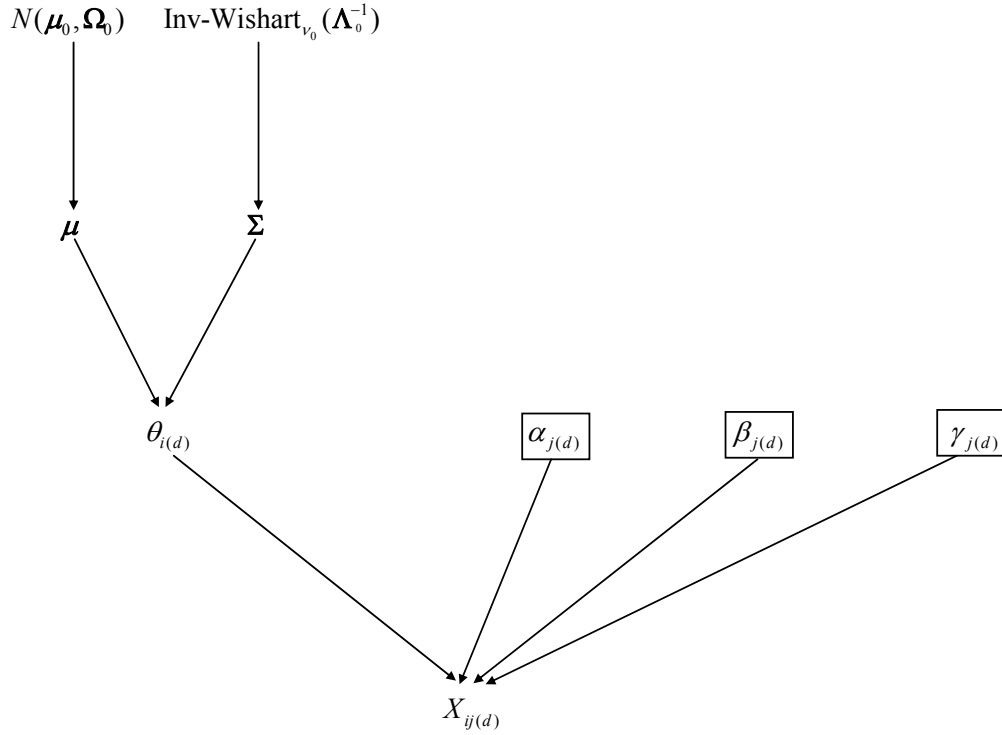
1 where

- 2 • $X_{ij(d)}$ is the response of examinee i to the j^{th} item of dimension d ;
- 3 • $\theta_{i(d)}$ is the d^{th} component of the vector $\boldsymbol{\theta}_i$ (i.e., $\boldsymbol{\theta}_i = \{\theta_{i(d)}\}$);
- 4 • $\alpha_{j(d)}$, $\beta_{j(d)}$, and $\gamma_{j(d)}$ are the parameters of the j^{th} item of dimension d ;
- 5 • $d = 1, \dots, D$ (the number of dimensions);
- 6 • $j(d) = 1(d), \dots, J(d)$; and
- 7 • $\sum_{d=1}^D J(d) = J$.

8 For this example, we will assume that the item parameters are known (i.e., we are only interested in
9 scoring the responses from multiple tests),

10 In addition to $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_i\}$, we will also estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, which represent the *hyperparameters*.

11 Below is a directed acyclic graph of the of the model. The variables in boxes are assumed to be known.
12
13
14



15 The likelihood of the data matrix \mathbf{X} is given by

$$\begin{aligned}
L(\mathbf{X}|\Theta) &= \prod_{i=1}^N \prod_{d=1}^D \prod_{j(d)=1}^{J(d)} P_{ij(d)}^{X_{ij(d)}} (1 - P_{ij(d)})^{1-X_{ij(d)}} \\
&= \prod_{i=1}^N L(\mathbf{X}_i|\theta_i).
\end{aligned}$$

1 *Prior, Posterior, and Full Conditional Distributions*

2

3 The the prior distribution of θ_i is parameterized as

$$\theta_i|\mu, \Sigma \sim N(\mu, \Sigma),$$

4 where

$$\begin{aligned}
\mu &\sim N(\mu_0, \Omega_0) \\
\Sigma &\sim \text{Inv-Wishart}_{\nu_0}(\Lambda_0^{-1}).
\end{aligned}$$

5

A Detour: The Inverse-Wishart Distribution

6 The Wishart distribution is the D -dimensional generalization of the univariate χ^2 distribution.

7

8 For a $D \times D$ symmetric, positive definite matrix Λ , and $\nu \geq D + 2$, the $D \times D$ symmetric, positive
9 definite matrix $\Sigma^{-1} \sim \text{Wishart}_{\nu}(\Lambda)$ with probability density function

$$p(\Sigma^{-1}|\nu, \Lambda) = c \cdot |\Sigma|^{-\frac{\nu-D-1}{2}} \exp\left[-\frac{1}{2}\text{tr}(\Lambda^{-1}\Sigma^{-1})\right],$$

10 where $c = \left[2^{\nu D/2} \pi^{D(D-1)/4} \prod_{d=1}^D \Gamma\left(\frac{\nu+1-d}{2}\right)\right]^{-1} |\Lambda|^{-\nu/2}$ and $E[\Sigma^{-1}] = \nu\Lambda$.

11

If $\Sigma^{-1} \sim \text{Wishart}_{\nu}(\Lambda)$ then $\Sigma \sim \text{Inv-Wishart}_{\nu}(\Lambda^{-1})$ with probability density function

$$p(\Sigma|\nu, \Lambda^{-1}) \propto |\Sigma|^{-\frac{\nu+D+1}{2}} \exp\left[-\frac{1}{2}\text{tr}(\Lambda\Sigma^{-1})\right],$$

12 and the expected value is $E(\Sigma) = (\nu - D - 1)^{-1}\Lambda$. The constant of proportionality is also c .

13

14 Let

$$\begin{aligned}
\Sigma &\sim \text{Inv-Wishart}_{\nu_0}(\Lambda_0^{-1}) \\
\theta_i|\mu_i, \Sigma &\sim N(\mu_i, \Sigma)
\end{aligned}$$

15 The posterior distribution of Σ is

$$\Sigma|\Theta, \mathbf{M} \sim \text{Inv-Wishart}_{\nu_N}(\Lambda_N^{-1})$$

16 where $\mathbf{M} = \{\mu_i\}$, $\nu_N = \nu_0 + N$, and $\Lambda_N = \Lambda_0 + \sum_{i=1}^N (\theta_i - \mu_i)(\theta_i - \mu_i)'$.

17

18 The derivation of this result is shown below.

$$\begin{aligned}
p(\Sigma) &\propto |\Sigma|^{-\frac{\nu+D+1}{2}} \exp \left[-\frac{1}{2} \text{tr}(\Lambda \Sigma^{-1}) \right], \\
p(\theta_i | \mu_i, \Sigma) &\propto |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\theta_i - \mu_i)' \Sigma^{-1} (\theta_i - \mu_i) \right].
\end{aligned}$$

1 Using the properties of the trace of a matrix, the likelihood of θ_i can be written as

$$p(\theta_i | \mu_i, \Sigma) \propto |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \text{tr} \left[(\theta_i - \mu_i)(\theta_i - \mu_i)' \Sigma^{-1} \right] \right\}.$$

2 The likelihood of the Θ is

$$\begin{aligned}
p(\Theta | M, \Sigma) &\propto \prod_{i=1}^N p(\theta_i | \mu_i, \Sigma) \\
&= |\Sigma|^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \text{tr} \left[(\theta_i - \mu_i)(\theta_i - \mu_i)' \Sigma^{-1} \right] \right\}.
\end{aligned}$$

3 Combining the prior density and the likelihood

$$\begin{aligned}
p(\Sigma | \Theta, M) &\propto p(\Sigma) \times p(\Theta | M, \Sigma) \\
&= |\Sigma|^{-\frac{\nu+D+1}{2}} \exp \left[-\frac{1}{2} \text{tr}(\Lambda \Sigma^{-1}) \right] \\
&\quad |\Sigma|^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \text{tr} \left[(\theta_i - \mu_i)(\theta_i - \mu_i)' \Sigma^{-1} \right] \right\} \\
&= |\Sigma|^{-\frac{\nu+N+D+1}{2}} \exp \left\{ -\frac{1}{2} \text{tr} \left[\left(\Lambda + \sum_{i=1}^N (\theta_i - \mu_i)(\theta_i - \mu_i)' \right) \Sigma^{-1} \right] \right\} \\
&\sim \text{Inv-Wishart}_{\nu_N}(\Lambda_N^{-1}).
\end{aligned}$$

4 Below are the steps in drawing an observation from $\text{Inv-Wishart}_{\nu}(\Lambda^{-1})$.

- 5 1. Draw θ_i , $i = 1, \dots, \nu$, from $N(\mathbf{0}, \Lambda)$;
- 6 2. Find $\mathbf{S} = \sum_{i=1}^{\nu} \theta_i \theta_i'$;
- 7 3. \mathbf{S}^{-1} is a draw from $\text{Inv-Wishart}_{\nu}(\Lambda^{-1})$.

8 End of Detour

9 The joint posterior distribution can be expressed as

$$\begin{aligned}
p(\Theta, \mu, \Sigma | \mathbf{X}) &\propto p(\mathbf{X} | \Theta) \times p(\Theta | \mu, \Sigma) \times p(\mu) \times p(\Sigma) \\
&= \prod_{i=1}^N p(\mathbf{X}_i | \theta_i) \times p(\theta_i | \mu, \Sigma) \times p(\mu) \times p(\Sigma).
\end{aligned}$$

10 Again, this is not a known distribution so we need to sample from the full conditional distributions.

11

1 We need the full conditional distributions for Σ , μ , and θ_i .

2

3 The full conditional distribution of μ is

$$\begin{aligned} p(\mu|\mathbf{X}, \Sigma, \Theta) &\propto p(\mathbf{X}|\mu, \Sigma, \Theta) \times p(\mu|\Sigma, \Theta) \\ &\propto p(\mathbf{X}|\Theta) \times p(\Theta|\mu, \Sigma) \times p(\mu|\Sigma) \\ &\propto p(\Theta|\mu, \Sigma) \times p(\mu) \end{aligned}$$

4 Given the prior distribution and the hyperdistributions above, the full conditional posterior distribution
5 of μ is

$$\begin{aligned} \mu|\mathbf{X}, \Sigma, \Theta &\sim N(\mu_N, \Omega_N) \\ &= N\left((\Omega_0^{-1} + N\Sigma^{-1})^{-1}(\Omega_0^{-1}\mu_0 + N\Sigma^{-1}\bar{\theta}), (\Omega_0^{-1} + N\Sigma^{-1})^{-1}\right). \end{aligned}$$

6 This full conditional distribution is a known distribution, and can be sampled directly.

7

8 The full conditional distribution of Σ is

$$\begin{aligned} p(\Sigma|\mathbf{X}, \mu, \Theta) &\propto p(\mathbf{X}|\mu, \Sigma, \Theta) \times p(\Sigma|\mu, \Theta) \\ &\propto p(\mathbf{X}|\Theta) \times p(\Theta|\mu, \Sigma) \times p(\Sigma|\mu) \\ &\propto p(\Theta|\mu, \Sigma) \times p(\Sigma) \end{aligned}$$

9 Again, given the prior distribution and the hyperdistributions above, the full conditional posterior
10 distribution of Σ is

$$\begin{aligned} \Sigma|\mathbf{X}, \mu, \Theta &\sim \text{Inv-Wishart}_{\nu_N}(\Lambda_N^{-1}) \\ &= \text{Inv-Wishart}_{\nu_0+N}\left(\Lambda_0 + \sum_{i=1}^N(\theta_i - \mu)(\theta_i - \mu)'\right). \end{aligned}$$

11 This full conditional distribution is also a known distribution, and can be sampled directly.

12

13 The full conditional distribution of θ_i is

$$\begin{aligned} p(\theta_i|\mathbf{X}, \mu, \Sigma, \Theta^{(-1)}) &\propto p(\mathbf{X}|\mu, \Sigma, \Theta) \times p(\theta_i|\mu, \Sigma, \Theta^{(-1)}) \\ &\propto p(\mathbf{X}^{(-1)}|\Theta^{(-1)}) \times p(\mathbf{X}_i|\theta_i) \times p(\theta_i|\mu, \Sigma) \\ &\propto p(\mathbf{X}_i|\theta_i) \times p(\theta_i|\mu, \Sigma) \end{aligned}$$

14 This full conditional distribution is NOT a known distribution so we have to use the M-H algorithm
15 to sample from this distribution *indirectly*.

16

17 At iteration $t + 1$, for $i = 1, \dots, N$, sample θ_{id}^* from $N(\theta_{id}^{(t)}, \sigma_\epsilon^2)$, and accept θ_i^* with probability

$$\alpha(\theta_i^{(t)}, \theta_i^*) = \min\left(\frac{p(\mathbf{X}_i|\theta_i^*) \times p(\theta_i^*|\mu, \Sigma)}{p(\mathbf{X}_i|\theta_i^{(t)}) \times p(\theta_i^{(t)}|\mu, \Sigma)}, 1\right).$$

1 *Implementation in Ox of the MCMC Algorithm for Estimating Σ , μ , and Θ*

2
3 The program `simultaneous.ox` implements the above MCMC algorithm. This program is specifically
4 for $D = 2$, and $J(1) = J(2) = 10$.

5
6 It starts with the basic headers.

7
8 The first set of global variables describes the characteristics of the data to be analyzed, and variables
9 that will be holding the current values of the parameters to be estimated.

```
decl N=1000,D=2,J=20;//sample size, dimension, total test length
decl Jd=<10;10>//test length per dimension
decl T=5000,m=2000;//iterations,burn-in
decl X;//data
decl mu0,Sigma0,theta0;//variables holding the current values
```

10 The next set of global variables are variables needed to generate the data. μ , σ_1^2 , σ_2^2 , and ρ can be
11 specified here; `theta` will be generated given the hyperparameters, and `alpha`, `beta`, `gamma` will be
12 read from a file.

```
decl mu=<0;-.2>;
decl Sigma;
decl S11=1,rho=.8,S22=1.2;
decl theta,alpha,beta,gamma;
```

13 The last set of variables are the parameters of the hyperdistributions, σ_ϵ , and the sufficient statistics.

14
15 The variable `design` is a $D \times J$ matrix that converts $\Theta_{N \times D}$ into an $N \times J$ matrix $\Theta_{N \times D}^*$ such that
16 θ_{id} is in the appropriately dimension (i.e., d).

```
decl nu_0,Lambda_0;
decl mu_0,Omega_0;
decl cand_sd;
decl smu,sSigma,stheta;
decl ssmu,ssSigma,sstheta;
decl design;
```

17 The function `ThreePL` computes the 3PL probabilities and requires $\Theta_{N \times D}^*$, and the item parameters.

```

ThreePL(const t,a,b,c){
decl one=ones(N,1),gam,prob;
gam=one*c';
prob=gam+(1-gam)./(1+exp(-1.7*((one*a').*t-one*(a.*b)')));
return prob;
}//ThreePL

```

- 1 The function `draw_mu` samples μ directly from its full conditional distribution given the current values
2 of Σ and Θ .

```

draw_mu(){
decl OmegaN=invert(invert(Omega_0)+N*invert(Sigma0));
decl muN=OmegaN*(invert(Omega_0)*mu_0+N*invert(Sigma0)*(meanc(theta0)'));
mu0=choleski(OmegaN)*rann(D,1)+muN;
}//end draw_mu

```

- 3 The function `draw_Sigma` samples Σ directly from its full conditional distribution given the current
4 values of μ and Θ .

```

draw_Sigma(){
decl nun=nu_0+N;
decl Lambdan=Lambda_0+((theta0-mu0')'*(theta0-mu0'));
decl the=rann(nun,D)*(choleski(invert(Lambdan)))';
Sigma0=invert(the'*the);
}//draw_Sigma

```

- 5 The function `logmvd` computes $-\frac{1}{2}(\theta_i - \mu)' \Sigma^{-1}(\theta_i - \mu)$ given $\theta_i - \mu$.

```

logmvd(decl m){
m=m*choleski(invert(Sigma0));//';
m=-.5*sumr(m.^2);
return m;
}//end logmvd

```

- 6 The function `draw_theta` updates θ_i . Note that in the sampling stage, we sample θ_{id} individually,
7 although in the decision stage, θ_i is accepted as a vector.

```

draw_theta(){
decl theta1=cand_sd*rann(N,D)+theta0;
decl irt0=ThreePL(theta0*design,alpha,beta,gamma);
irt0=sumr(X.*log(irt0)+(1-X).*log(1-irt0))+logmvd(theta0-mu0');
decl irt1=ThreePL(theta1*design,alpha,beta,gamma);
irt1=sumr(X.*log(irt1)+(1-X).*log(1-irt1))+logmvd(theta1-mu0');
decl acc=irt1-irt0;
acc=vecindex(ranu(N,1).<exp(acc));
theta0[acc][]=theta1[acc][];
} //end draw_theta

```

- 1 The function `initiate` assigns values to the prior parameters, and and initial values to the parameters
- 2 to be estimates.
- 3
- 4 It also sets $\sigma_\epsilon = .1$.
- 5
- 6 Finally, the function also sets the sufficient statistics to zero.

```

initiate(){
nu_0=D+2;
Lambda_0=unit(D);
mu_0=zeros(D,1);
Omega_0=unit(D);
mu0=zeros(D,1);
Sigma0=unit(D);
theta0=rann(N,D);
cand_sd=.1;
smu=zeros(D,1),sSigma=zeros(D,D),stheta=zeros(N,D);
ssmu=zeros(D,1),ssSigma=zeros(D,D),sstheta=zeros(N,D);
} //end initiate

```

- 7 The function `update` updates the sufficient statistics after the burn-in. Note that running sum of
- 8 squared draws are only for Θ . Thus, we won't estimate the posterior standard deviation of μ and Σ .

```

update(){
smu=smu+mu0;
sSigma=sSigma+Sigma0;
stheta=stheta+theta0;
sstheta=sstheta+theta0.^2;
} // end update

```

1 The function `print_results` prints the estimates of μ , Σ , and \mathbf{R} , which is the correlation matrix.
2
3 In addition, it also print $Cor(\theta, \tilde{\theta})$, $\sqrt{\sum \tilde{V}(\theta_i | \mathbf{X}_i)/N}$, and $\sqrt{\sum (\theta_i - \tilde{\theta}_i)^2/N}$ to evaluate the quality
4 of the ability estimates.
5

```
print_results(){
println("\n*****Hyperparameter Estimates*****");
print("\nEstimate of Mu:",smu'/T);
sSigma=sSigma/T;
print("\nEstimate of Covariance Matrix:",sSigma);
print("\nEstimate of Correlatiopn Matrix:");
decl tmp=sqrt(invert(diag(diagonal(sSigma))));
print(tmp*sSigma*tmp);
tmp=stheta/T;
println("\n*****Ability Estimate Results*****");
print("\nCorr(True,Est):");
print(diagonal(correlation(theta~tmp)[:D-1][D:]));
decl tmp2=sqrt(meanc((sstheta-T*tmp.^2)/(T-1)));
print("Mean Posterior SD:",tmp2);
print("RMSE:",sqrt(meanc((theta-tmp).^2)));
} //print_results
```

6 The `main` function has several components. The first component reads the ten 3PL item parameters.
7 These parameters are replicated D times to correspond to the total test length J before they are des-
8 ignated as `alpha`, `beta`, and `gamma`.
9
10 The code creates a design matrix appropriate for the different test lengths across the D dimensions.
11

```
tmp=0;
design=zeros(D,J);
for(decl i=0;i<D;++i){
tmp=tmp~sumc(Jd[:i]);
design[i][tmp[i]:tmp[i+1]-1]=ones(1,Jd[i]);
} //end i
```

12 In general, the design matrix is of the following form:

$$\begin{pmatrix} \mathbf{1}_{1 \times J(1)} & \mathbf{0}_{1 \times J(2)} & \cdots & \mathbf{0}_{1 \times J(D)} \\ \mathbf{0}_{1 \times J(1)} & \mathbf{1}_{1 \times J(2)} & \cdots & \mathbf{0}_{1 \times J(D)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{1 \times J(1)} & \mathbf{0}_{1 \times J(2)} & \cdots & \mathbf{1}_{1 \times J(D)} \end{pmatrix}_{D \times J}$$

1 The next few lines create **Sigma** from the variances and covariances, and generates **theta**, and the item
2 response matrix **X**.

3
4 The function **initiate** is then called, after which the functions **draw_mu**, **draw_Sigma**, and **draw_theta**
5 are called $T + m$ times. The function **update** is called only when $t > m$.

6
7 Lastly, the draws are summarized using the function **print_results**.

8
9 1. Change N , μ , σ^2 , and ρ to see how they affect the estimates.

10
11 2. Double $J(d)$ to examine the effect of test length on the estimates. Which of the estimates are most
12 affected by a longer test?.

13
14 3. What happens to the estimates if more discriminating items are involved? Use $\alpha^* = 1.5 \times \alpha$ to
15 generate and analyze the data.

16 17 **An MCMC Algorithm for Multi-Unidimensional Ability Estimation with Ancillary Vari-** 18 **ables**

19
20 In addition to examinees' responses to the test questions, other ancillary information about the exam-
21 inees is also available in many testing situations.

22
23 Examples of these ancillary information include demographic variables such as sex, age, and race, and
24 educational variables such as grade level, courses taken, and previous test scores.

25
26 To achieve greater efficiency, in addition to the responses to other tests, we can also incorporate the
27 ancillary information in the ability estimation.

28
29 Specifically, we will regress θ_i on the covariates (i.e., ancillary variables) \mathbf{Y}_i .

30
31 As in the previous example, we will assume a multi-unidimensional structure.

32
33 Below is a directed acyclic graph of the model. As before the variables in boxes (i.e., item parameters
34 and covariates) are assumed to be known.

35
36 From this model, we can see that the response $X_{ij(d)}$ is only affected by a single ability (i.e., $\theta_{i(d)}$).

37
38 However, our inference about $\theta_{i(d)}$ can be improved by taking into consideration Σ , the covariance
39 between abilities.

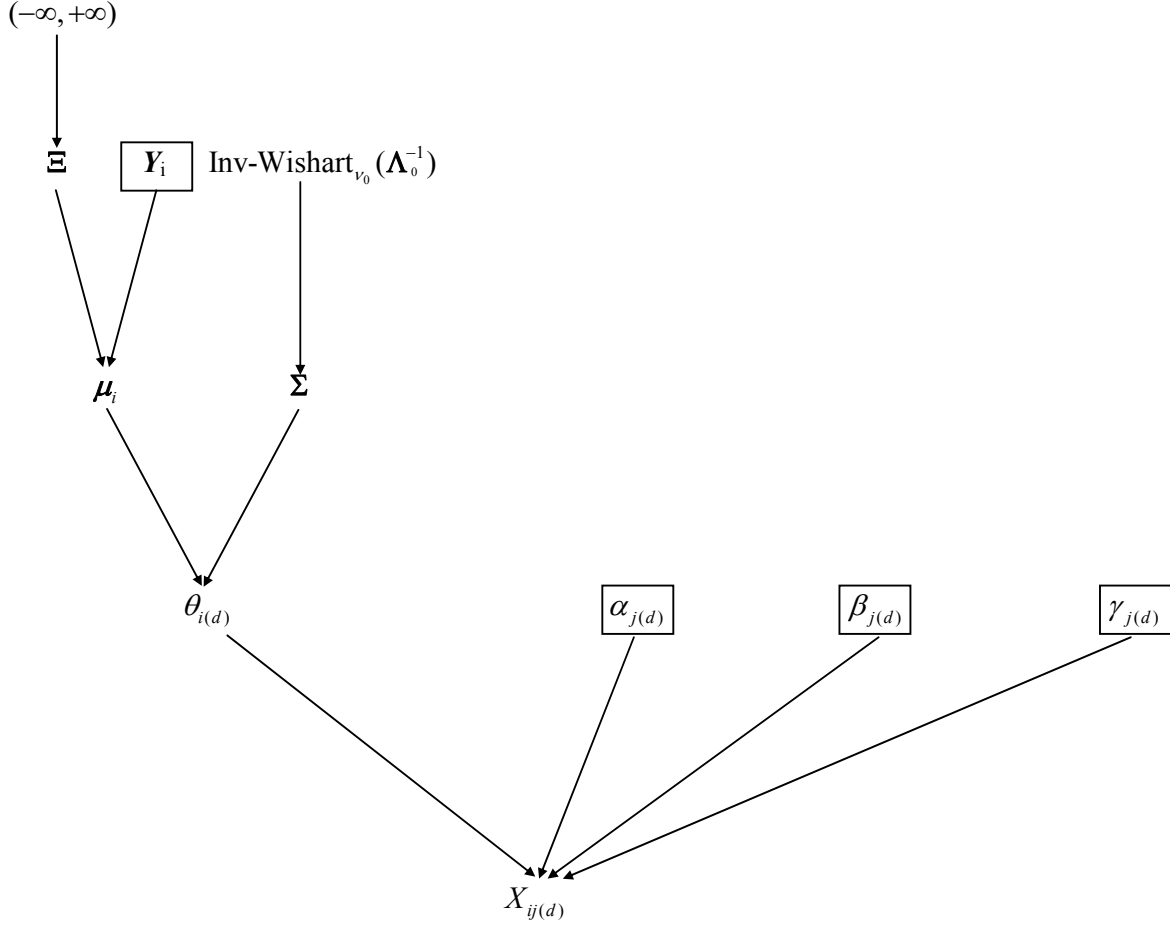
40
41 Additionally, by knowing how the abilities and covariates are related, we can further improve our in-
42 ference about θ from knowing \mathbf{Y} .

43
44 As we have done in the previous example, the correlation between the abilities are estimated, rather
45 than assumed.

46
47 In the same way, we will also empirically determine, rather than assume, the relationship between θ
48 and \mathbf{Y} .

49

- 1 Note that, based on the graph, the covariate \mathbf{Y}_i determines the mean of $\boldsymbol{\theta}_i$, not the ability itself.
2
3 This allows us to have individuals with identical \mathbf{Y} to have identical $\boldsymbol{\mu}$, but actually have different $\boldsymbol{\theta}$.
4 The only exception is when \mathbf{Y} perfectly determines $\boldsymbol{\theta}$, which is next to impossible. In such a situation,
5 no test needs to be administered.
6



- 7 As before, the likelihood of the data matrix \mathbf{X} is given by

$$\begin{aligned}
 L(\mathbf{X}|\boldsymbol{\Theta}) &= \prod_{i=1}^N \prod_{d=1}^D \prod_{j(d)=1}^{J(d)} P_{ij(d)}^{X_{ij(d)}} (1 - P_{ij(d)})^{1-X_{ij(d)}} \\
 &= \prod_{i=1}^N L(\mathbf{X}_i|\boldsymbol{\theta}_i).
 \end{aligned}$$

- 8 *Prior, Posterior, and Full Conditional Distributions*

9

- 10 For examinee i with ability $\boldsymbol{\theta}_i$, the prior distribution is given by

$$\boldsymbol{\theta}_i | \mathbf{Y}_i, \boldsymbol{\Xi}, \boldsymbol{\Sigma} \sim N(\boldsymbol{\Xi}' \mathbf{Y}_i, \boldsymbol{\Sigma}),$$

1 with

$$\begin{aligned}\Sigma &\sim \text{Inv-Wishart}_{\nu_0}(\Lambda_0^{-1}), \\ \Xi &\sim (-\infty, +\infty),\end{aligned}$$

2 where

- 3 • $\mathbf{Y}_i = \{Y_{i1}, Y_{i2}, \dots, Y_{ip}\}$ is the vector the p observable covariates of examinee i ;
- 4 • $\Xi = \{\xi_1, \xi_2, \dots, \xi_d, \dots, \xi_D\}$ is the $p \times D$ matrix of regression parameters;
- 5 • $\mathbf{u}_i = \Xi' \mathbf{Y}_i$ and Σ are the mean vector and the covariance matrix of the multivariate normal
- 6 distribution, respectively;
- 7 • ν_0 and Λ_0 are defined as before; and
- 8 • $(-\infty, +\infty)$ represents the improper prior of Ξ to facilitate sampling from the posterior distribu-
- 9 tion.

10 The joint distribution of Ξ , Σ , and Θ given \mathbf{X} and \mathbf{Y} is

$$\begin{aligned}p(\Xi, \Sigma, \Theta | \mathbf{X}, \mathbf{Y}) &\propto p(\mathbf{X} | \mathbf{Y}, \Xi, \Sigma, \Theta) \times p(\Theta | \mathbf{Y}, \Xi, \Sigma) \times p(\Sigma | \mathbf{Y}, \Xi) \times p(\Xi | \mathbf{Y}) \\ &= p(\mathbf{X} | \Theta) \times p(\Theta | \mathbf{Y}, \Xi, \Sigma) \times p(\Sigma) \times p(\Xi).\end{aligned}$$

11 Below are the full conditional distributions of Ξ , Σ , and θ_i .

12

13 For the regression parameters - $\Xi | \Sigma, \Theta, \mathbf{X}, \mathbf{Y}$:

$$\begin{aligned}p(\Xi | \Sigma, \Theta, \mathbf{X}, \mathbf{Y}) &= p(\Xi | \Sigma, \Theta, \mathbf{Y}) \\ &\propto p(\Theta | \Xi, \Sigma, \mathbf{Y}) \times p(\Xi | \Sigma, \mathbf{Y}) \\ &\propto p(\Theta | \Xi, \Sigma, \mathbf{Y}) \times p(\Xi).\end{aligned}$$

14 For the regression parameters, the full conditional distribution of Ξ , $p(\Xi | \Theta, \Sigma, \mathbf{Y})$, given an improper

15 prior is a known distribution called the *matrix-normal*.

16

17 Specifically, this full conditional distribution is given by

$$N((\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\Theta, \Sigma \otimes (\mathbf{Y}'\mathbf{Y})^{-1}),$$

18 where \otimes is the Kronecker product. We can sample *directly* from this distribution.

19

20 To sample from $p(\Xi | \mathbf{Y}, \Sigma, \Theta)$ in *vector* format, one can use the following algorithm:

$$\text{Vec}(\Xi) = \mathbf{Z}_{I \times D} [\Sigma \otimes (\mathbf{Y}'\mathbf{Y})^{-1}]^{1/2} + \text{Vec}((\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\Theta),$$

21 where $\mathbf{Z} \sim \text{MVN}(\mathbf{0}, \mathbf{I})$, and $\text{Vec}(\cdot)$ stacks the vectors of the argument.

22

23 As an alternative, we can sample from this full conditional distribution without using the Kronecker

24 product by, first, recasting the matrices as follows:

$$\Xi_{(D \times P) \times 1}^* = \text{Vec}(\Xi),$$

$$\Theta_{(N \times D) \times 1}^* = \text{Vec}(\Theta),$$

$$\Sigma_{(N \times D) \times (N \times D)}^* = \begin{pmatrix} \Sigma & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma \end{pmatrix},$$

$$\mathbf{Y}_{(N \times D) \times (D \times P)}^* = \begin{pmatrix} \mathbf{Y}'_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{Y}'_1 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{Y}'_1 \\ \mathbf{Y}'_2 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{Y}'_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{Y}'_2 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{Y}'_N & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{Y}'_N & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{Y}'_N \end{pmatrix},$$

1 where \mathbf{Y}_i is the covariate vector of examinee i (i.e., the transpose of the i th row of the design matrix \mathbf{Y}).

2

Then, we can sample from

$$\Xi^* \sim \text{MVN}((\mathbf{Y}^{*'} \Sigma^{*-1} \mathbf{Y}^*)^{-1} \mathbf{Y}^{*'} \Sigma^{*-1} \Theta^*, (\mathbf{Y}^{*'} \Sigma^{*-1} \mathbf{Y}^*)^{-1}).$$

3 For the covariance matrix - $\Sigma | \mathbf{X}, \mathbf{Y}, \Xi, \Theta$:

$$\begin{aligned} p(\Sigma | \mathbf{X}, \mathbf{Y}, \Xi, \Theta) &= p(\Sigma | \mathbf{Y}, \Xi, \Theta) \\ &= p(\Theta | \mathbf{Y}, \Xi, \Sigma) \times p(\Sigma | \mathbf{Y}, \Xi) \\ &\propto p(\Theta | \mathbf{Y}, \Xi, \Sigma) \times p(\Sigma). \end{aligned}$$

4 The full conditional distribution of Σ , $p(\Sigma | \Xi, \Theta, \mathbf{Y})$, is an Inv-Wishart $_{\nu_N}(\Lambda_N^{-1})$, where $\nu_N = \nu_0 + N$,
5 and $\Lambda_N = \Lambda_0 + \sum (\theta_i - \Xi' \mathbf{Y})(\theta_i - \Xi' \mathbf{Y})'$.

6

7 We already know how to sample from this distribution *directly*.

8

9 For the ability parameters - $\theta_i | \mathbf{X}, \mathbf{Y}, \Xi, \Sigma, \Theta^{(-i)}$:

$$\begin{aligned} p(\theta_i | \mathbf{X}, \mathbf{Y}, \Xi, \Sigma, \Theta^{(-i)}) &\propto p(\mathbf{X} | \mathbf{Y}, \Xi, \Sigma, \Theta) \times p(\theta_i | \mathbf{Y}, \Xi, \Sigma, \Theta^{(-i)}) \\ &\propto p(\mathbf{X}_i | \theta_i) \times p(\theta_i | \mathbf{Y}_i, \Xi, \Sigma). \end{aligned}$$

10 Again, this full conditional distribution is NOT a known distribution so we have to use the M-H algo-
11 rithm to sample from this distribution *indirectly*.

12

1 At iteration $t + 1$, for $i = 1, \dots, N$, sample θ_{id}^* from $N(\theta_{id}^{(t)}, \sigma_\epsilon^2)$, and accept θ_i^* with probability

$$\alpha(\theta_i^{(t-1)}, \theta_i^*) = \min \left\{ \frac{p(\mathbf{X}_i | \theta_i^*) p(\theta_i^* | \mathbf{Y}_i, \Xi^{(t)}, \Sigma^{(t)})}{p(\mathbf{X}_i | \theta_i^{(t-1)}) p(\theta_i^{(t-1)} | \mathbf{Y}_i, \Xi^{(t)}, \Sigma^{(t)})}, 1 \right\}.$$

2 The above model specification simultaneously accounts for responses to other tests and ancillary vari-
3 ables in scoring examinees' responses.

4
5 In the traditional method of scoring, where tests scored one at a time and no covariates are used,
6 $\Xi = \mathbf{0}$ and $\{\Sigma\}_{dd'} = 0$ for $d \neq d'$.

7
8 When tests are scored simultaneously without the benefit of ancillary variables, only Σ is estimated
9 whereas $\Xi = \mathbf{0}$.

10
11 Finally, when ancillary variables are used to score one test at a time, both Ξ and Σ are estimated,
12 with the constraint that $\{\Sigma\}_{dd'} = 0$ for $d \neq d'$.

13
14 When Σ is treated as a diagonal matrix, the algorithm may need to be substantially altered to estimate
15 the diagonal elements of Σ individually.

16
17 It should be noted that the Σ in this model is the covariance of θ *conditional* on \mathbf{Y} , as in $\Sigma_{\theta|\mathbf{Y}}$.

18
19 However, we might also be interested in making inference about $\Sigma_{\theta\theta}$, which is the *marginal* covariance
20 of θ .

21
22 To do so, we just need to express $\Sigma_{\theta\theta}|\mathbf{X}, \mathbf{Y}$ as $f(\mathbf{X}, \mathbf{Y}, \Xi, \Sigma, \Theta)$, where, at iteration $t + 1$:

$$\Sigma_{\theta\theta}^{(t+1)} = \frac{1}{N-1} \sum_{i=1}^N (\theta_i^{(t+1)} - \mu^{(t+1)})(\theta_i^{(t+1)} - \mu^{(t+1)})'.$$

23 From the draws, we can also make inference about $\Sigma_{\theta\theta}|\mathbf{X}, \mathbf{Y}$.

24

25 A Detour: A More General Method of Generating $\mathbf{Y}_{P \times 1}$ and $\theta_{D \times 1}$ Jointly

26 Previously, we discussed generating \mathbf{Y}_i and θ_i jointly such that they follow a particular structure.

27

28 One structure is that $\theta \sim N(\mathbf{0}, \mathbf{R})$. This structure is needed if the item parameters are also to be
29 estimated because, without such a constraint, the model will not be identified, hence, not estimable.

30

31 Another structure is the $\Sigma_{\mathbf{Y}\theta} = \{\sqrt{\psi/P}\}$, which is due to the setting $\psi_d = \psi$.

32

33 Lastly, we fixed the *conditional* covariance $\Sigma_{\theta|\mathbf{Y}}$, and allowed the *marginal* covariance $\Sigma_{\theta\theta}$ to vary.

34

35 In the current method, we will allow Σ to not necessarily be \mathbf{R} , ψ_d to not necessarily be ψ , and $\Sigma_{\theta|\mathbf{Y}}$
36 to vary while fixing $\Sigma_{\theta\theta}$.

37

38 As before, we will let $\mathbf{Y} \sim N(\mathbf{0}, \mathbf{I})$. In practice, we can use standardization or principal component
39 analysis to transform \mathbf{Y} .

40

1 In addition, we will let the matrix of regression coefficients $\Xi = \{\xi_{dp}\} = \{\sigma_d \times \sqrt{\psi_d/P}\}$. This is but
2 one way of splitting ψ_p into the regression weights ξ_{dp} such that

$$\frac{\sum_{d=1}^D \xi_{dp}^2}{\sigma_d^2} = \psi_d.$$

3 Following are the steps in generating N θ observations with a *marginal* covariance of θ is $\Sigma_{\theta\theta}$, and
4 where the proportion variance accounted for by the covariates with respect to θ_d is equal to ψ_d .

- 5 1. Generate a $Y_{N \times P}$ matrix of covariates, as in, N observations from $N(\mathbf{0}, \mathbf{I})$. Standardize \mathbf{Y} .
- 6 2. Compute $\mathbf{M} = \mathbf{Y}\Xi$, the matrix of *conditional* mean vectors for the N examinees.
- 7 3. Compute $\Sigma_{\theta|\mathbf{Y}} = \Sigma_{\theta\theta} - \Sigma_{\theta\mathbf{Y}}\Sigma_{\mathbf{Y}\theta}$.
- 8 4. Generate θ_i from $N(\{\mathbf{M}\}_i', \Sigma_{\theta|\mathbf{Y}})$, $i = 1, 2, \dots, N$, where $\{\mathbf{M}\}_i$ is the i^{th} row of \mathbf{M} .

9 Additional generalizations can be made. For example μ_{θ} need not be equal to $\mathbf{0}$.

10 *Implementation in Ox of the Method for Generating \mathbf{Y} and θ Jointly*

11 The program `ancillary_gen.ox` can be used to generate \mathbf{Y} and θ that follow specific structure. In
12 addition, it also generates the item response data that can be analyzed using the `ancillary.ox`.

13 The program needs to read the ten 3PL item parameters currently contained in the file `best10.txt`.

14 This program is specifically written for $D = 2$, and $P = 2$, and $J(1) = J(2) = 10$. The number of
15 examinees N can be changed.

```
decl N=1000,D=2,P=2,J=20;//sample size, dimension, total test length
decl Jd=<10;10>;//test length per dimension
```

21 The proportion of variance account for by \mathbf{Y} in each dimension (i.e., ψ_1, ψ_2) can be modified.

22 Similarly the variances of θ , σ_1^2 and σ_2^2 , can be specified. Instead of σ_{12} , users can specify ρ .

```
decl psi1=.8,psi2=.5;
decl S11=1,rho=.8,S22=1.2;
```

24 In addition to the screen output, which provides statistics describing the generated \mathbf{Y} and θ , the pro-
25 gram will also create an output file called "ancillary.dat".

26 This file will contain $\mathbf{X}_{N \times 20}$, $\mathbf{Y}_{N \times 2}$ and $\Theta_{N \times 2}$, in this order.

28 **End Detour**

1 *Implementation in Ox of the MCMC Algorithm for Estimating Ξ , Σ , Θ and $\Sigma_{\theta\theta}$*

2
3 The program `ancillary.ox` implements the above MCMC algorithm. This program is designed to be
4 used with `ancillary_gen.ox`, and will estimate Ξ , Σ , and θ_i , and $\Sigma_{\theta\theta}$.

5
6 The program includes the usual headers.

7
8 The first set of global variables is very similar to that in `simultaneous.ox`. However, this time, we
9 need a variable P to denote the number of covariates, and Y to denote the vector of covariates.

10
11 We are also declaring Xi0 to hold the current values of Ξ , and iYtY which is equal to $(Y'Y)^{-1}$. By
12 declaring iYtY, we can compute this variable only once in the entire estimation process even though
13 it will use it each iteration.

```
decl N=1000,D=2,P=2,J=20;//sample size, dimension, total test length
decl Jd=<10;10>//test length per dimension
decl T=5000,m=2000;//iterations,burn-in
decl X,Y;//data, covariates
decl Xi0,Sigma0,mSigma0,theta0;//variables holding the current values
decl theta,alpha,beta,gamma;// abilities, item parameters
decl design,iYtY;//design matrix, inverse of Y'Y
```

14 The next set of global variables represent the priors of the hyperparameters, σ_θ^2 , and the sufficient
15 statistics. We are using an improper prior for Ξ so we do not have parameters to declare for this
16 variable.

```
decl nu_0,Lambda_0;
decl cand_sd;
decl sXi,sSigma,smSigma,stheta;
decl ssXi,ssSigma,ssmSigma,ssttheta;
```

17 The next few lines pertain to the function `ThreePL`, which is unaltered from before.

18
19 The function `draw_Xi` updates Ξ .

20
21 We are using the first of the methods for sampling Ξ discussed above. In Ox, `**` represents \otimes .

```
draw_Xi(){
decl b0;
b0=(vec(rann(P,D)))'*(choleski(Sigma0**iYtY))'+(vec(iYtY*(Y'theta0)))';
Xi0=shape(b0,P,D);
};//end draw_Xi
```

1 The function `draw_Sigma` updates the conditional covariance $\Sigma_{\theta|Y}$.

2
3 That this is identical to the `draw_Sigma` function earlier except that instead of a common μ , each
4 individual gets its own $\mu_i = \Xi'Y_i$.

```
draw_Sigma(){
decl nun=nu_0+N;
decl Lambdan=Lambda_0+(theta0-Y*Xi0)'(theta0-Y*Xi0);
decl the=rann(nun,D)*(choleski(invert(Lambdan)))';
Sigma0=invert(the'the);
} //draw_Sigma
```

5 The function `logmvd` is unchanged.

6
7 The function `draw_theta` is different from `draw_theta` function only in one respect: μ_i is used in place
8 of μ .

```
draw_theta(){
decl theta1=cand_sd*rann(N,D)+theta0;
decl irt0=ThreePL(theta0*design,alpha,beta,gamma);
irt0=sumr(X.*log(irt0)+(1-X).*log(1-irt0))+logmvd(theta0-Y*Xi0);
decl irt1=ThreePL(theta1*design,alpha,beta,gamma);
irt1=sumr(X.*log(irt1)+(1-X).*log(1-irt1))+logmvd(theta1-Y*Xi0);
decl acc=irt1-irt0;
acc=vecindex(ranu(N,1).<exp(acc));
theta0[acc][]=theta1[acc][];
} //end draw_theta
```

9 The function `initiate` remains largely the same, except that it has been adapted to accommodate Ξ .

10
11 The function `update` updates a subset of sufficient statistics (i.e., those necessary for estimating Ξ , Σ ,
12 Θ , and $\Sigma_{\theta\theta}$, and the posterior variance of θ_i .

13
14 Note that the `smSigma`, the sufficient statistic for $\Sigma_{\theta\theta}$, is not sampled, only updated.

```
update(){
sXi=sXi+Xi0;
sSigma=sSigma+Sigma0;
smSigma=smSigma+theta0'theta0/(N-1);
stheta=stheta+theta0;
sstheta=sstheta+theta0.^2;
} // end update
```

1 The function `print_results` is largely similar to its previous definition. However, this time, we print
2 separately for the *conditional* and *marginal* covariance and correlation matrices.

3
4 In the `main` function, we read in the ten item parameters before they are replicated, which is exactly
5 the same process used in generating the data.

6
7 We also read in the data, which this time is composed of \mathbf{X} , \mathbf{Y} , and $\mathbf{\Theta}$. We need to ensure that the
8 right columns are assigned to the right variables, as in, `X`, `Y` and `theta`.

9
10 The last three steps involve initiating the variables, sampling from the posterior, which include updat-
11 ing the sufficient statistics after the burn-in, and printing the results.

12
13 1. For $\rho = .8$ and $\sigma_1^2 = \sigma_2^2 = 1$, what values of $\psi_1 = \psi_2$ will lead to *conditional* correlation of \mathbf{I} ? What
14 does this tell us about the relationship between θ_1 and θ_2 vis-à-vis \mathbf{Y} ?

15
16 2. Change N , ψ_d , σ^2 , and ρ to see how they affect the estimates.

17
18 By comparing $\psi_d \neq 0$ or $\rho \neq 0$ to zero against $\psi_d = 0$ or $\rho = 0$, we can evaluate the impact of
19 incorporating the covariates and responses to other tests in the estimation process.

20
21 3. Double $J(d)$ to examine the effect of test length on the estimates. Which of the estimates are most
22 affected by a longer test?.

23
24 4. What happens to the estimates if more discriminating items are involved? Again, use $\alpha^* = 1.5 \times \alpha$
25 to generate and analyze the data.

26 27 Additional Exercises

28
29 1. Change `simultaneous.ox` to accommodate $D = 4$ dimensions.

30
31 Let $MSE_{\rho=0}$ and MSE_{ρ} be the mean squared error of the ability estimates from scoring the tests *one*
32 *at a time* and scoring the tests, respectively.

33
34 The ratio $RE = MSE_{\rho=0}/MSE_{\rho}$ is the relative efficiency due to simultaneous scoring.

35
36 For example, an $RE = 1.5$ can be interpreted as gaining 50% in efficiency by taking into account
37 responses from other tests. That is, the tests now function as if they are $1.5 \times J(d)$ long.

38
39 Find the MSE of θ when $D = 4$, $J(d) = 10$, and $\rho = 0, .5, .9$. Average the MSE across the D
40 dimensions.

41
42 What is the RE under this set-up?

43
44 What happens when D is reduced to 2 or when $J(d)$ is increased to 20?

45
46 2. Impact of Incorporating Ancillary Information in Ability Estimation

47
48 2ai) For a particular data set, find the EAP of θ *one ability at a time*.

49

1 2aii) For the same data, use `simultaneous.ox`. You need to modify the code so it can to read in the
2 same data used as above.

3

4 2aiii) Finally, use `ancillary.ox` to analyze the same data.

5

6 2b) Compare the three estimates based on $\text{Corr}(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}})$, mean posterior variance, and root-mean squared
7 error.

8

9 Comparing 2ai and 2aii allows us to examine the advantage of *multidimensional* over *unidimensional*
10 scoring

11

12 Comparing 2aii and 2aiii allows us to examine the advantage of using covariates on top of *multidimen-*
13 *sional* scoring.

14

15 Comparing 2ai and 2aiii allows us to examine the impact of *multidimensional* scoring and use of co-
16 variates over *unidimensional* scoring.

17

18 3. Manipulate a few of the variables and repeat the analysis to determine under which conditions the
19 different procedure yield similar/different results.

20

21 Optional Exercises

22

23 1. Modify `ancillary.ox` so that Ξ can be sampled without using the Kronecker product \otimes .

24

25 2. Simultaneous Calibration: Sampling \mathbf{R}

26

27 The presence of multiple tests can also be used to improve item parameter estimation. The impact of
28 multiple abilities on the item parameters is *indirect*, but can be substantial if the sample size is not
29 sufficiently large.

30

31 One way of ensuring that the model is determined is to constraint the ability distribution, as in,
32 $\boldsymbol{\theta} \sim N(\mathbf{0}, \mathbf{R})$.

33

34 The full conditional distribution of \mathbf{R} can be shown to be

$$p(\mathbf{R}|\mathbf{X}, \boldsymbol{\Theta}, \mathbf{B}) \propto p(\boldsymbol{\Theta}|\mathbf{R}) \times p(\mathbf{R}),$$

35 where \mathbf{B} is the matrix of item parameters.

36

37 Because \mathbf{R} , we can define $p(\mathbf{R})$ so its full conditional distribution is an Inverse-Wishart.

38

39 Although we know how to sample directly from this full conditional distribution, the candidate value
40 $\boldsymbol{\Sigma}^*$ is not guaranteed to be \mathbf{R}

41

42 To ensure that we only sample correlation matrices, we need to first convert $\boldsymbol{\Sigma}^* \rightarrow \mathbf{R}^*$, and then accept
43 it with probability

44

$$\alpha(\mathbf{R}^{(t)}, \mathbf{R}^*) = \min \left(\exp \left[\frac{D+1}{2} \left(\log \mathbf{R}^* - \log \mathbf{R}^{(t)} \right) \right], 1 \right).$$

1 Write a program that can simultaneously calibrate D tests. Manipulate N , D , and ρ to see their
 2 impact on the item parameter estimates.

3

4 3. Higher-Order IRT

5

6 When multiple abilities are involved, we can also posit that higher-order ability, say, ω , such that

$$p(\boldsymbol{\theta}_i|\omega_i) = \prod_{d=1}^D p(\theta_{id}|\omega_i),$$

7 where $\theta_{id} = \lambda_d \times \omega_i + \epsilon_{id}$, $|\lambda_d| \leq 1$, and $\epsilon_{id} \sim N(0, 1 - \lambda_d^2)$.

8

9 If we assume $\omega \sim N(0, 1)$, then

$$\begin{aligned}\theta_{id}|\omega_i &\sim N(\lambda_d \times \omega_i, 1 - \lambda_d^2) \\ \theta_{id} &\sim N(0, 1)\end{aligned}$$

10 It can be shown that ω , θ_d , and $\theta_{d'}$ are related as follows:

$$\begin{aligned}\text{Corr}(\omega, \theta_d) &= \lambda_d \\ \text{Corr}(\theta_d, \theta_{d'}) &= \lambda_d \times \lambda_{d'}\end{aligned}$$

11 The full conditional distribution of ω_i^* can be shown to be

$$N\left(\sigma_\omega^{2(t+1)} \times \sum_{d=1}^D \frac{\lambda_d^{(t)} \times \theta_{id}^{(t)}}{1 - \lambda_d^{2(t)}}, \sigma_\omega^{2(t+1)}\right),$$

12 where $\sigma_\omega^{2(t+1)} = \lambda_d^{2(t)} / [1 - \lambda_d^{2(t)}]$. We can sample from this distribution directly.

13

14 However, the remaining parameters (i.e., $\boldsymbol{\theta}_i$ and λ_d) need to be sampled using M-H.

15

16 Note that because the marginal distribution of θ_d is already $N(0, 1)$, we can use this model to simul-
 17 taneously calibrate multiple tests.

18

19 Again, we need M-H to sample \mathbf{B} .

20

21 Write a program that can simultaneously calibrate D tests using the HO-IRT model. Manipulate N ,
 22 D , and ρ to see their impact on the item parameter estimates.

23

24 In addition, for $D = 4$, create a correlation matrix \mathbf{R} such that the correlations cannot be fully written
 25 as $\lambda_d \times \lambda_{d'}$.

26

27 Compare the quality of the estimates obtained using the HO-IRT model and the model for simultaneous
 28 estimation.

29