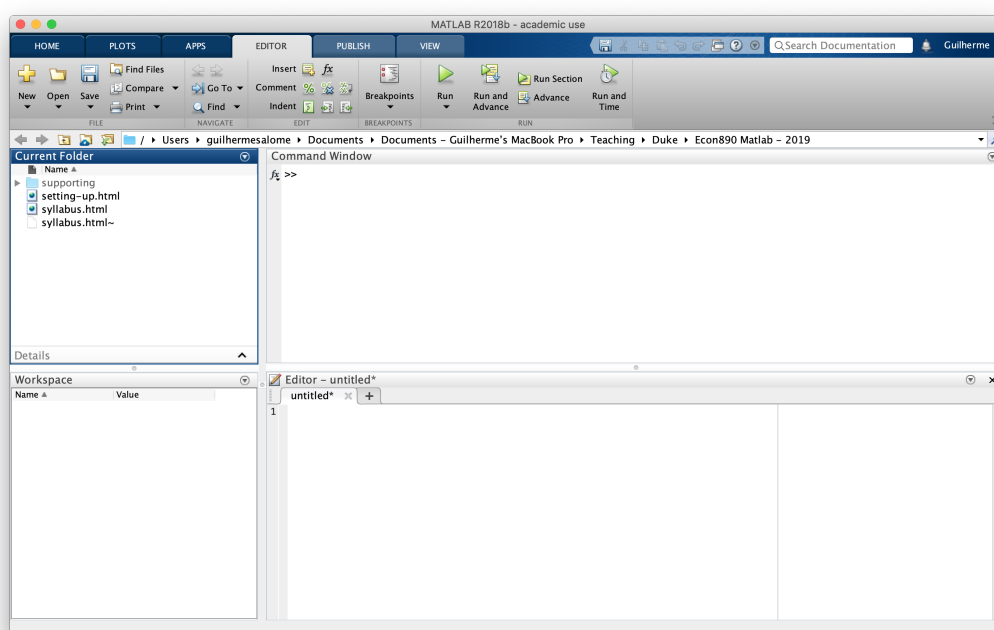# Matlab Basics

# 1   The Matlab Integrated Development Environment

When Matlab is first opened you are presented with an integrated development environment (IDE), which is depicted in Figure 1 below.



**Figure 1:**  The Matlab IDE.

The figure shows four windows from left-to-right: Current Folder, Command Window, Workspace, Editor. The `Current Folder` window displays all files in the working directory Matlab is currently at. This working directory is where Matlab will search for user-defined functions and other files (more on this later). You can change the working directory by using the navigation bar above the window.

The `Command Window` is known as a repeat-eval-print-loop (REPL), and is where we can interactively type and execute Matlab commands. As you type commands in, Matlab will execute them and then wait for more. For example, type the following in the `Command Window` and hit enter to execute:

```
1  cd ~/
```

Matlab will change the working directory to the home folder of your computer. Here, we evaluated a special function `cd` with the argument `~/`.
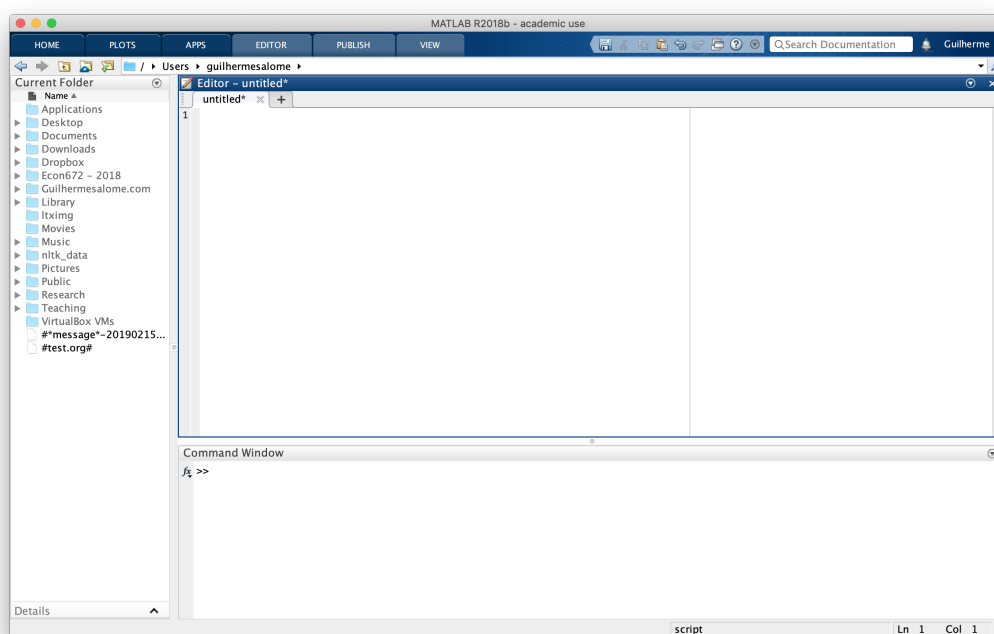
The `Workspace` window shows a list of all variables created so far. The list is empty because we have not created any variable yet. Create a new variable by typing:

```
1  x = 1
```

Now the `Workspace` window should show the name and value of the variable.

Last, the `Editor` window is where we edit `.m` files. These files are used to create Matlab scripts and to define functions. Notice that the `Editor` window can open multiple files by using separate tabs.

All four windows can be rearranged and even removed from view. Feel free to customize the IDE to your needs. I often remove the `Workspace` window, and rearrange the `Editor` window for it be on top, leaving a smaller space for the `Command Window` on the bottom of the IDE. Figure 2 shows the IDE with the windows rearranged.



**Figure 2:** A possible way to organize the windows in the IDE.
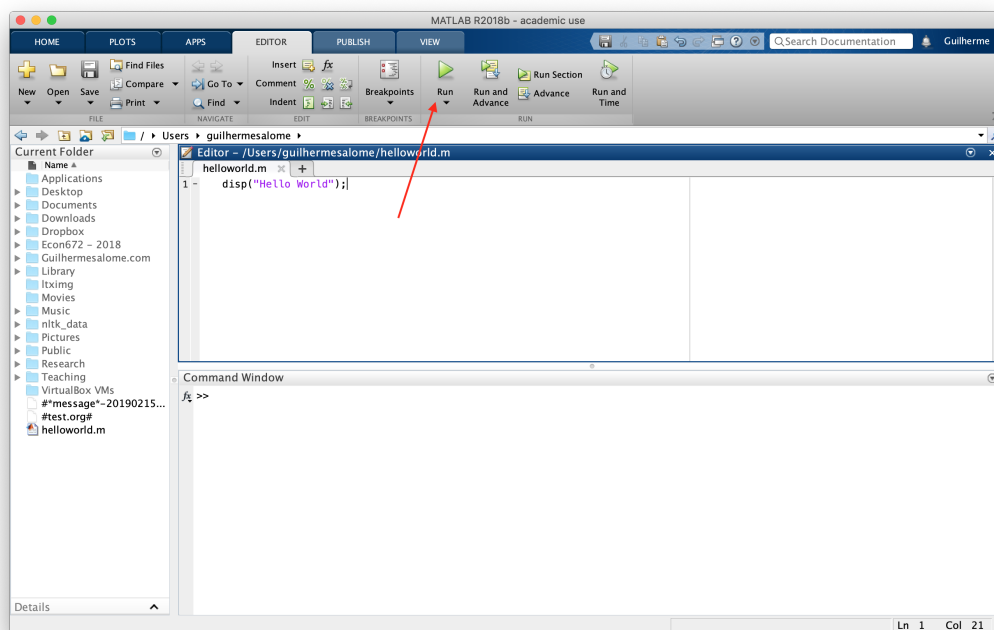
## 2   Hello World!

We will now create a "hello world" program in Matlab. Go on the `Editor` window and type the following:

```
1  % helloworld.m
2  disp("Hello World!");
```

Save the file with a descriptive name, like `helloworld.m`. Now, we would like to run our `helloworld.m` program. There are two ways of doing so:

- Click on `Editor` and then on the green button `Run` (see Figure 3);

- Or, on the `Command Window` type the name of the file without the `.m` extension and hit enter.

**Figure 3:** Running a Matlab script.

When we run the file, Matlab displays the result in the `Command Window`. The first line of the file starts with a percentage symbol, which indicates a comment line. In the second line, the function `disp` is called. This function takes a string as an argument and displays it on the screen.

# 3  Linear Regression with Ordinary Least Squares

To learn the fundamentals of the Matlab programming language we will implement the ordinary least squares (OLS) estimator for the parameters of a linear regression under the classical assumptions.

Let's assume that the economic model of interest is given by:

$$\underbrace{y}_{n\times 1} = \underbrace{X}_{n\times K}\underbrace{\beta}_{K\times 1} + \underbrace{\varepsilon}_{n\times 1}$$

where $y$ is a $n \times 1$ vector of observed dependent variables, $X$ is an $n \times K$ matrix of observed explanatory variables, $\beta$ is a $K \times 1$ vector of unknown parameters and $\varepsilon$ is a $n \times 1$ vector of unobserved explanatory variables.

We know that the OLS estimator of $\beta$ (see Equation 1.2.5 in Hayashi (2000)) is given by:

$$\hat{\beta} = (X'X)^{-1}X'y \tag{1}$$

Our first objective is to implement this estimator.

Using a bottom-up approach, we need to understand:

1. The data types for numbers in Matlab;

2. How to create vectors and matrices;

3. How to do matrix transposition, inversion and multiplication;

4. How to create a function that given $X$ and $y$ produces $\hat{\beta}$;

5. How to generate data to test out our code.

## 3.1   Numeric Data Types for Numbers

Numeric values are stored in the double-precision floating point format, by default. The class name of this type of numeric value is `double`. We can create doubles by simply assigning them to variables:

```
1  x = 10;
2  y = 5
```

Whenever you see code in a box as above, you should run it yourself by typing the commands line by line, in the `Command Window`.

Notice that when the semicolon is omitted, Matlab displays the result of the assignment on the screen. If you want to suppress this display you need the semicolon. The use of semicolon is important when constructing scripts, so that the `Command Window` is not filled with information we do not need to see. However, omitting the semicolon when typing in the `Command Window` is useful to see the result of operations.

```
1  % The two lines below
2  y = 5;
3  disp(y);
4  % are equivalent to:
5  y = 5
```

The usual mathematical operations can be performed with numbers:

```
1   x + y
2   x - y
3   % multiplication
4   x*y
5   % floating point division
6   y/x
7   x/y
8   % exponentiation
9   x^y
10  % remainder of integer division
11  mod(13, 5)
12  % integer division
13  floor(13/5)
14  % absolute value
15  abs(-1)
16  % maximum
17  max(x, y)
18  % minimum
19  min(x, y)
```

```matlab
20  % exponential
21  exp(4)
22  % narutal logarithm
23  log(exp(1))
24  % square root
25  sqrt(4)
26  4^0.5
27  % factorial
28  factorial(x)
29  % pi
30  pi
31  sqrt(pi)
32  % trigonometry
33  sin(0)
34  cos(0)
35  tan(0)
```

The result of the operations above are still `double`.

You can verify the type and size in memory of a variable with the special[1] function `whos`:

```matlab
1  whos x;
2  whos x y;
```

It is important to notice that the number displayed in the `Command Window` is a formatted version of the actual number stored in a variable. For example, $\pi$ is displayed as 3.1416. By default, Matlab displays numbers using at most 5 digits, and it may use scientific notation for numbers that are too big or too small.

```matlab
1  pi
2  0.0000009809809
3  213124412123
```

We can change the display format for numeric values between 5 digits and 15 digits using the special command `format`:

```matlab
1  pi
2  % Use 15-digits to display numbers
3  format long;
4  pi
5  % Use 5-digits to display numbers
6  format short;
7  pi
8  % Force scientific notation with 5-digits
9  format short e;
10 pi
```

---

[1]The function `whos` is being called with a special syntax called the `command syntax`. Usually, in a function call the parameters are passed inside parentheses, as in `mod(x, y)`. However, the `command syntax` allows us to pass parameters without using parentheses, as in `whos x y`. It is not often used, except with a few Matlab commands. More details about this `command syntax` are available at this page.

```matlab
11  % Force scientific notation with 15-digits
12  format long e;
13  pi
```

Observe that, as with other implementations of floating-point arithmetic (as in C or FORTRAN), Matlab is also subject to the usual rounding errors.

```matlab
1  3*(4/3 - 1) - 1
```

However, this is not a bug, but an issue resulting from not being possible to write $\frac{4}{3}$ as a binary number using a finite number of bits.

There are other numeric types available, such as integers, complex numbers, infinity, not a number, among others. It is worth mentioning that infinity and not a number can result from some arithmetic operations:

```matlab
1  % Finding infinity in Matlab
2  1/0
3  -1/0
4  % Not a number
5  1/0 - 1/0
```

Using `double` is sufficient for most applications, and for this reason we will skip the other numeric types. However, applications that are memory bound or that require a lot of optimization may benefit from the other numeric types. The official documentation on all numeric types is available on this page.

## 3.2   Vectors and Matrices

Vectors and matrices are structured collections of numbers of type `double` (by default). We can create vectors and matrices by:

```matlab
1  % A 5x1 vector
2  y = [10; -2; 3.4; exp(3); log(4.5)]
3  % A 5x3 matrix
4  X = [1, 2, 3;
5       1,-3, 5;
6       1, 4, 0;
7       1, 0,-1;
8       1, 4.5, 3.3]
```

Notice that commas separate columns and semicolons separate rows. The use of commas to separate columns is optional.

We can very obtain the dimension of a vector with the function `length`, and the dimensions of a matrix with the function `size`:

```matlab
1  length(y)
2  size(X)
```

What would be the result of calling `length` on a matrix? We can figure that out by reading the `help` documentation on the `length` function.

```matlab
1  help length
```

In the case of a matrix, it will return the value of the highest dimension.

Notice that the function `size` applied on a matrix returns two numbers. We can access these two numbers separately by: saving them on two separate variables, or asking the function `size` to return only one of the numbers.

```
1  A = ones(2, 3)
2  help size
3  % save the dimensions to two variables
4  [d1, d2] = size(A)
5  % ask size to return only the 1st dimension
6  d1 = size(A, 1)
7  % ask size to return only the 2nd dimension
8  d2 = size(A, 2)
```

To access the elements of vectors and matrices we use the following syntax:

```
1   % The first element of a vector has the index 1
2   y(1)
3   % The second element of a vector has the index 2
4   y(2)
5   % And so on until the last element
6   y(5)
7   % For a matrix we need to define both dimensions
8   % The very first value at the first row and first column
9   X(1, 1)
10  % The value on the first row and second column
11  X(1, 2)
12  % The value on the last row and the last column
13  X(5, 3)
14  % There is a special syntax to get all rows of a column
15  X(:, 2) % all rows of 2nd column
16  X(:, 1) % all rows of 1st column
17  % The same notation applies to get all columns of a row
18  X(4, :) % all columns of the 4th row
19  X(1, :) % all columns of the 1st row
20  % We can use a similar notation to access slices of the
       matrix
21  X(2:4, 4) % rows 2 to 4 of the 4th column
22  X(3, 1:2) % columns 1 to 2 of the 3rd row
23  X(2:4, 1:2) % columns 1 to 2 of rows 2 to 4
```

The values stored in arrays (vectors, matrices and higher-dimensional matrices) can be modified in place:

```
1  a = [0 1; 2 3]
2  a(2, 1) = 500
3  % multiple values can be updated at once
4  a(:, 1) = [30; 30] % modifying an entire column
5  a(2, :) = [20 20] % modifying an entire row
```

There are a few functions that assist in the creation of matrices:

```
 1  % Matrix full of zeros
 2  help zeros
 3  zeros(2)
 4  zeros(4, 1)
 5  % Matrix full of ones
 6  ones(3)
 7  ones(2, 3)
 8  % Identity matrix
 9  eye(5)
10  eye(2, 3)
11  % Sparse Identity Matrix
12  A = eye(1000)
13  B = speye(1000) % optimized
14  whos A B
15  % Diagonal matrix
16  diag([1 2 3 4 5])
```

We can also generate a vector by using a range of numbers:

```
 1  % start:stop
 2  A = 1:4
 3  % the values are incremented by 1 until the value is greater
      than
 4  % the stop, which is 4
 5  B = -5: 3.5
 6  % we can change the increment using the syntax start:
      increment:stop
 7  C = -5:0.5:3.5
```

We can combine other matrices to create new ones:

```
 1  A = [X(:, 1), X(:, 3)]
```

## 3.3   Matrix Operations

We can now create matrices and modify them.  Now we will cover the basic matrix operations. Let's start with matrix and scalar operations:

```
 1  X = [1 2 3; 4 5 6]
 2  % operations with a scalar are elementwise
 3  X + 3
 4  X - 2
 5  % special notation for elementwise multiplication, division
      and exponentiation
 6  X.*3
 7  X./3
 8  X.^3
 9  % log, exp and sqrt apply elementwise
10  log(X)
```

```
11  exp(X)
12  sqrt(X)
```

Matrix and matrix operations:

```
1   A = [10 20 30; 0 -5 -10];
2   B = [2 2 2; 5 10 20];
3   % elementwise operations (matrices have the same dimensions)
4   A + B
5   A - B
6   A./B
7   A.*B
8   % concatenation
9   % concatenate rows (vertical concatenation)
10  C = vertcat(A, B)
11  C = [A; B]
12  % concatenate columns (horizontal concatenation)
13  D = horzcat(A, B)
14  D = [A B]
15  % repetition
16  repmat(A, 2, 3) % repeats matrix A: 2 times in row, 3 times
       in column
```

Now, for linear algebra:

```
1   A = [10 2 30;
2        0 -5 -10]
3   B = ones(size(A, 2), 1)
4   % Matrix multiplication
5   C = A*B
6   % Matrix tranposition
7   transpose(A)
8   A'
9   % Matrix inversion
10  C = [A; B']
11  C^(-1)
12  inv(C)
13  % Solving linear system Cx=B
14  solution = C\B
15  % Eigenvalues and eigenvectors
16  eigenvalues = eig(C)
17  [eigenvectors, eigenvalues] = eig(C)
18  % Singular value decomposition
19  [U, S, V] = svd(C)
20  % Kronecker Product
21  kron(A, B)
```

There are many other useful linear algebra functions. The full documentation is available on the linear algebra reference page.

## 3.4   Functions (and Conditionals)

We now have the main ingredients to implement the OLS estimator from Equation 1:

```matlab
y = [10; -2; 3.4; exp(3); log(4.5)]
X = [1, 2, 3;
     1,-3, 5;
     1, 4, 0;
     1, 0,-1;
     1, 4.5, 3.3]
% The OLS estimator for the regression y = X*beta + epsilon
beta_hat = inv(X'*X)*X'*y
```

Now that we have the estimator, we want to be able to use it for whatever matrices $X$ and $y$ we have. To do that, we create a function.

Functions in Matlab are created in separate files ending in `.m`. Each function has its own `.m` file[2]. The syntax for a function is `function [output1, ..., outputN] = foo(input1, ..., inputM)`, then the actual function code, and last we finish with the keyword `end`. Let's create a simple function for the OLS estimator:

```matlab
% linreg_ols.m
function [beta] = linreg_ols(y, X)
beta = inv(X'*X)*X'*y;
end
```

Whenever you see a box as above where the first line starts with a comment of the type `% somename.m`, it means that you should type the code in a new file and save it with the name given in the comment. Save the file on the current folder Matlab is using.

Notice that the name of the function <u>must</u> be the same name of the file. After saving the file, we can use the function by just calling its name:

```matlab
beta = linreg_ols(y, X)
```

Matlab can access all `.m` files in the current folder and on sub-folders within it.[3]

There are a couple of modifications we need to do to the `linreg_ols` function. First, we need to add documentation about it, so that our future selves know what it actually does. We do so by adding comments to the function right after the line where the function is declared. These comments should accomplish the following:

- A description of the function does in one line (or very few lines);

- A description of the function inputs;

- A description of the function outputs;

- Examples of usage.

---

[2]Strictly speaking there can be more than one function in a single `.m` file, but there must be one function with the same name as the file name, and the other functions in the file are known as local functions. These local functions are not available outside of the `.m` file, and their use is to break down the code in smaller pieces.

[3]It is possible to use functions and other files in folders outside the current folder. To do so, you need to add the path to that specific folder to the search path Matlab uses. This can be done with the function `addpath`. See the `addpath` reference for details.

After finishing core functions it is absolutely necessary to add documentation. If you stop using your code, even for a few days, it is very easy to forget what you did and what different functions do. Having the documentation is fundamental to quickly remember how to operate your code. We can add comments in Matlab with a percentage symbol. Let's add comments so that the usage of the function is clear:

```matlab
% linreg_ols.m
function [beta] = linreg_ols(y, X)
% LINREG_OLS estimates beta in the linear regression y = X*
    beta + epsilon
% via ordinary least squares
%
% Args:
%     y: A nx1 vector of dependent variables
%     X: A nxK matrix of independent variables
%
% Returns:
%     beta: A Kx1 vector of the estimated beta coefficients
%
% Examples:
%     y = [1;2;3];
%     X = rand(3);
%     beta = linreg_ols(y, X)
%
% Reference:
% Equation 1.2.5 in Hayashi, F. (2010).
% Econometrics. Princeton University Press. isbn:
    9780691010182.

beta = inv(X'*X)*X'*y;

end
```

You can see the documentation of the function with the special command `help`.

```matlab
help linreg_ols
```

Often we want to add an intercept to the linear regression. Let's extend our `linreg_ols` function to automatically add an intercept. To conserve space in the lecture notes, I will remove the documentation for our function, but you should not.

```matlab
% linreg_ols.m
function [beta] = linreg_ols(y, X)
% add intercept to X
X = horzcat(ones(size(X, 1), 1), X); % X is a local variable
    to this function
beta = inv(X'*X)*X'*y;
end
```

If we want the intercept to be optional, then we need to add a new input to the function that will represent whether or not we want the intercept to be added to $X$. We

can do so by using the conditional statement `if`.

### 3.4.1   Conditionals

The syntax for conditional statements is:

```
1  if condition
2      % code to run if condition is true
3  elseif another_condition
4      % code to run if condition is false and
5      % another_condition is true
6  else
7      % code to run if condition and another_condition are
           false
8  end
```

In this case the `condition` is pretty simple, we can add an input with the name `add_intercept` and if its value is `true` we add the intercept, otherwise we do not. In Matlab, the possible logical (boolean) values are `true` and `false`, but the numbers `1` and `0` can also be used to the same extent. The basic relational operators are:

```
1  % equal to
2  true == 1
3  % not equal to
4  true ~= false
5  % greater than
6  3 > 2
7  % greater or equal than
8  2 >= 2
9  % less than
10 3 < 2
11 % less or equal than
12 -1 <= 0
```

It is also possible to combine different conditions with these logical operators:

```
1  % and
2  (1 > 0) && (2 > 1)
3  % or
4  (1 < 0) || (2 < 1)
5  % not
6  ~(1 > 0)
```

### 3.4.2   Optional Intercept

Let's modify `linreg_ols` to make the intercept optional:

```
1  % linreg_ols.m
2  function [beta] = linreg_ols(y, X, add_intercept)
3  if add_intercept:
4      X = horzcat(ones(size(X, 1), 1), X);
```

```matlab
5  end
6  beta = inv(X'*X)*X'*y;
7  end
```

Finally, we update the documentation:

```matlab
1  % linreg_ols.m
2  function [beta] = linreg_ols(y, X, add_intercept)
3  % LINREG_OLS estimates beta in the linear regression y = X*
      beta + epsilon
4  % via ordinary least squares
5  %
6  % Args:
7  %     y: A nx1 vector of dependent variables
8  %     X: A nxK matrix of independent variables
9  %     add_intercept: A boolean variable. If true, then a
      column of
10 %         ones is added to X, so that an intercept for the
      linear
11 %         regression is estimated.
12 %
13 % Returns:
14 %     beta: A Kx1 vector of the estimated beta coefficients
15 %
16 % Examples:
17 %     y = [1;2;3];
18 %     X = rand(3);
19 %     beta = linreg_ols(y, X)
20 %
21 % Reference:
22 % Equation 1.2.5 in Hayashi, F. (2010).
23 % Econometrics. Princeton University Press. isbn:
      9780691010182.
24 if add_intercept:
25    X = horzcat(ones(size(X, 1), 1), X);
26 end
27 beta = inv(X'*X)*X'*y;
28 end
```

## 3.5   Generating Data (and Loops and Text)

To test if our function is actually doing what we think it is doing, we need to test it. Testing is imperative when programming, since generating code that is 100\

In this case, we are going to simulate data for $X$, $\varepsilon$ and $y$, while fixing true values for the parameters $\beta$. First, to simulate the data we will assume that In Matlab, the basic random number generation functions can generate values from the uniform and normal distributions. However, in the "Statistics and Machine Learning Toolbox" there are random number generators for many other distributions. Let's assume $X$ comes from

the uniform distribution on the interval $(0, 1)$, while $\varepsilon$ comes from the standard normal distribution. We will generate these values in a loop.

### 3.5.1   Loops

In Matlab, we can create `for` loops and `while` loops. The syntax for `for` loops is:

```matlab
% syntax for for loops
for variable = range
    % variable takes the value in the range
    % run some code
end
% example
for i = 1:10
    disp(i)
end
% example with nested loops
X = zeros(3)
for i = 1:size(X, 1)
    for j = 1:size(X, 2)
        X(i, j) = i + j;
    end
end
```

For loops will terminate when all values in the range are exhausted.

The syntax for `while` loops is similar:

```matlab
% syntax for while loops
while condition
    % if condition is true
    % run some code
end
% the condition must be a boolean value

% example
i = 0
while i < 10
    disp(i)
    i = i + 1 % no i++ or += on Matlab
end
```

While loops will terminate when the condition is no longer true.

If you write an infinite loop, you can break out of it with `Ctrl + C`. In addition, if you need a loop to terminate earlier you can use the keyword `break`. For example:

```matlab
numb_iter = 0
max_iter = 10
distance = 1000
while distance > 0.01
    distance = distance/2
    if numb_iter >= max_iter
```

```
 7          disp("I will break now")
 8          break
 9      else
10          numb_iter = numb_iter + 1
11      end
12 end
```

### 3.5.2   Generating Data in a Loop

Let's now create our test data. We will generate a matrix $X$ with 100 rows and 5 columns, and $\varepsilon$ will be a vector with 100 rows. We will create this data in a script, so that we can run it multiple times with different number of rows and columns.

```
 1 % script: test_linreg_ols.m
 2 % Generates data from the linear model y = X*beta + epsilon,
 3 % given the true beta values. Estimates beta using the
     function
 4 % linreg_ols and displays a table comparing true values to
 5 % parameters.
 6
 7 % number of rows and columns
 8 nrows = 100;
 9 ncols = 5;
10 % true betas
11 beta = [-1; 2; 0.5; 3; 0.1];
12 % generate data for X, epsilon and y
13 X = zeros(nrows, ncols);
14 epsilon = zeros(nrows, 1);
15 y = zeros(nrows, 1);
16 for i = 1:nrows
17     X(i, :) = rand(1, 5);
18     epsilon(i) = randn();
19     y(i) = X(i,:)*beta + epsilon(i);
20 end
21 % estimate the betas from X and y
22 beta_hat = linreg_ols(y, X, true);
```

Whenever we run the script above, we will generate data and estimate the values of $\beta$. Notice that we added an intercept in the estimation, even though the true model does not have an intercept. To make the comparison between true parameters values and their estimates, we would like to display the values side by side, with some text indicating what are the true parameters and what are the estimates.

### 3.5.3   Characters and Strings

In Matlab, there are two data types for text: `char` and `string`. A `char` holds a sequence of characters. It works like a vector of numbers, but holds characters instead. To create a `char` we use single quotes ':

```matlab
1  name = 'Guilherme'
2  % name is an array of characters, the individual letters of
       my name
3  length(name)
4  name(1)
5  name(length(name))
6  % we can also concatenate character arrays
7  disp(['My name is ' name '!'])
```

Characters are meant to store short pieces of text, and are used to specify file names, labels and titles for plots, and inputs for functions.

A `string` array is meant to hold and manipulate larger pieces of text. Matlab provides various functions for manipulating text stored in `string` arrays. String arrays are created in a fashion similar to matrices, but with brackets and text in double quotes ":

```matlab
1  % creating a simple string array
2  full_name = ["Guilherme", "Salome"];
3  % full_name is a string array
4  size(full_name)
5  full_name(1, 1)
6  full_name(1, 2)
7  % a string array with many rows and columns
8  students_names = ["First Name", "Last Name";
9                    "A", "B";
10                   "C", "D"];
11 size(students_names)
12 students_names(2, 2)
13 % a string array can also contain a single element
14 % in this case, we do not need the brackets
15 my_name = "Guilherme Salome"
16 size(my_name) % it is still an array
17 % we can find the length of a string with the strlength
       function
18 strlength(my_name)
19 % creating a string array that also has numbers will
20 % implicitely convert numbers to strings
21 data = ["Parameters", 1, 2, 3];
```

Let's explore some of the common operations we might perform with strings. A full list of functions that operate on strings are available on the reference page for characters and strings.

```matlab
1  % we can append a string to string arrays
2  first_names = ["Mary", "John", "Paul"]
3  full_names = first_names + " Smith"
4  % split strings
5  split_sentence = split("This is a sentence about to be split
       ")
6  % split can be used to clean csv data
```

```matlab
 7  csv_data = "1.2, 3.2, -0.5, 5.3"
 8  split_data = split(csv_data, ',')
 9  % join strings
10  a_burguer = join(["Bread", "Tomato", "Lettuce", "Burguer", "
       Cheese", ...
11                    "Bread"])
12  % if a command is too big for a single line, we can use
13  % the ... notation to continue it in the next line
14
15  % joining strings can be useful to store data in the csv file
16  data = [1.2; 3.2; -0.5; 5.3];
17  % the text manipulation functions apply to string arrays
18  % so we need to convert the numbers to a string array
19  data_string = string(data)
20  csv_data = join(data_string, ',')
21  % sorting strings
22  sort(full_names)
```

We can now display our parameter estimates using string arrays:

```matlab
1  % compare the true betas to the estimates
2  disp(horzcat(["True Values"; 0; beta], ["Estimates"; beta_hat
      ]))
```

### 3.5.4   Table

We can create a better display by using a `table`. A `table` contains rectangular data that is column-oriented. We can create a table to display the true parameter values and their estimates. We create a table by passing it vectors with the values, each vector will be a column in the table. Additionally, we can pass a string array containing the column names.

```matlab
 1  % script: test_linreg_ols.m
 2  % Generates data from the linear model y = X*beta + epsilon,
 3  % given the true beta values. Estimates beta using the
       function
 4  % linreg_ols and displays a table comparing true values to
 5  % parameters.
 6
 7  % number of rows and columns
 8  nrows = 100;
 9  ncols = 5;
10  % true betas
11  beta = [-1; 2; 0.5; 3; 0.1];
12  % generate data for X, epsilon and y
13  X = zeros(nrows, ncols);
14  epsilon = zeros(nrows, 1);
15  y = zeros(nrows, 1);
16  for i = 1:nrows
```

```
17        X(i, :) = rand(1, 5);
18        epsilon(i) = randn();
19        y(i) = X(i,:)*beta + epsilon(i);
20  end
21  % estimate the betas from X and y
22  beta_hat = linreg_ols(y, X, true);
23  % display a table with the true parameters and estimates
24  estimates = table([0; beta], beta_hat, 'VariableNames', ["
        Parameters", ...
25                          "Estimates"]);
26  disp(estimates)
```

We can now modify this script to see what happens to the estimates when the number of data points increases, or when the number of parameters increases and so on. We have covered most of the basics of Matlab by coding the `linreg_ols` function and the `test_linreg_ols.m` script.

# 4   Standard Errors for OLS Estimates

We will now extended the `linreg_ols` function to compute standard errors for the parameter estimates and the t-statistics.

Under the classical assumptions, the standard error of the OLS estimates (Equation 1.4.4 in Hayashi (2000)) is given by:

$$SE(\hat{\beta}_i) = \sqrt{s^2 \cdot [(X'X)^{-1}]_{[i,i]}}$$

where $s^2$ (Equation 1.2.13 in Hayashi (2000)) is computed from the residuals of the estimation:

$$e \equiv y - X\hat{\beta}$$

$$s^2 \equiv \frac{e'e}{n - K}$$

Let's implement this estimator in the `linreg_ols` function:

```
1  % linreg_ols.m
2  function [beta, s2, stderr] = linreg_ols(y, X, add_intercept)
3  % Returns:
4  %     beta: A Kx1 vector of the estimated beta coefficients
5  %     s2: Estimate of the residual variance
6  %     stderr: A Kx1 vector of the estimators of the standard
        errors
7  %         of the beta estimates
8  if add_intercept
9      X = horzcat(ones(size(X, 1), 1), X);
10 end
11 % compute inv(X'X) just once (used multiple times)
12 XX = X'*X;
13 % equivalent to inv(X'X) but faster and more numerically
        stable
```

```matlab
14  XXinverse = XX\eye(length(XX));
15  beta = XXinverse*X'*y;
16  % estimate Var(beta_hat|X)
17  e = y - X*beta;
18  s2 = (e'*e)/(size(X, 1) - size(X, 2));
19  stderr = diag((s2*XXinverse)^0.5);
20  end
```

Now, the t-statistic for the null-hypothesis of insignificant coefficients is:

$$t_i \equiv \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)} \sim t_{n-K}$$

To compute the p-value for the test (see Pages 38 and 39 of Hayashi (2000)) we need the cumulative density function (cdf) of the student's T distribution, which is available in the function `tcdf`.

```matlab
1  % linreg_ols.m
2  function [beta, s2, stderr, t_stat, p_value] = linreg_ols(y,
      X, add_intercept)
3  % LINREG_OLS estimates beta in the linear regression y = X*
      beta + epsilon
4  % via ordinary least squares
5  %
6  % Args:
7  %     y: A nx1 vector of dependent variables
8  %     X: A nxK matrix of independent variables
9  %     add_intercept: A boolean variable. If true, then a
      column of
10 %         ones is added to X, so that an intercept for the
      linear
11 %         regression is estimated.
12 %
13 % Returns:
14 %     beta: A Kx1 vector of the estimated beta coefficients
15 %     s2: Estimate of the residual variance
16 %     stderr: A Kx1 vector of the estimators of the standard
      errors
17 %         of the beta estimates
18 %     t_stat: A Kx1 vector with the t-statistics for the null
19 %         hypothesis that each beta is zero (separately)
20 %     p_value: A Kx1 vector with the p-values for the t-test
21 %
22 % Examples:
23 %     y = [1;2;3];
24 %     X = rand(3);
25 %     beta = linreg_ols(y, X, true)
26 %
27 % Reference:
28 % Equation 1.2.5 in Hayashi, F. (2010).
```

```matlab
29  % Econometrics. Princeton University Press. isbn:
        9780691010182.
30  if add_intercept
31      X = horzcat(ones(size(X, 1), 1), X);
32  end
33  % compute inv(X'X) just once (used multiple times)
34  XX = X'*X;
35  % equivalent to inv(X'X) but faster and more numerically
        stable
36  XXinverse = XX\eye(length(XX));
37  % number of observations and explanatory variables
38  n = size(X, 1);
39  K = size(X, 2);
40  beta = XXinverse*X'*y;
41  % estimate Var(beta_hat|X)
42  e = y - X*beta;
43  s2 = (e'*e)/(n - K);
44  stderr = diag((s2*XXinverse)^0.5);
45  % t-statistics
46  t_stat = beta./stderr;
47  p_value = (1 - tcdf(abs(t_stat), n - K))*2;
48  end
```

# 5   Organizing the Regression Results in an Object

We could further extend the `linreg_ols` function to compute all the other statistics commonly found in other specialized statistics software. Every time we add a new computation to the function, we would have to extend the number of outputs. When there are too many outputs in a function, we need a better way to organize them. To organize multiple outputs of a function we can use a `struct`.

A `struct` (structure array) is a data type that groups other data by fields. Each field has a name and can contain any type of data. We access fields in a structure array with the `.` notation. Let's create a simple `struct`:

```matlab
1   % data to store
2   beta_hat = [1;2;3];
3   stderr = [0.5, 0.3, 0.43];
4   % create a structure named results
5   results.b = beta_hat;
6   results.se = stderr;
7   % access the values in the struct
8   results.b
9   results.se
10  % creating a struct using the struct keyword
11  results = struct('b', beta_hat, 'se', stderr);
12  results.b
13  results.se
```

```
14  % creating a struct with empty fields (possibly to be filled
        later)
15  results = struct('b', [], 'se', []);
16  results.b
17  results.b = beta_hat;
18  results.b
```

We can use a `struct` to easily organize the output of `linreg_ols`:

```
 1  % linreg_ols.m
 2  function results = linreg_ols(y, X, add_intercept)
 3  % LINREG_OLS estimates beta in the linear regression y = X*
        beta + epsilon
 4  % via ordinary least squares
 5  %
 6  % Args:
 7  %     y: A nx1 vector of dependent variables
 8  %     X: A nxK matrix of independent variables
 9  %     add_intercept: A boolean variable. If true, then a
        column of
10  %        ones is added to X, so that an intercept for the
        linear
11  %        regression is estimated.
12  %
13  % Returns:
14  %     results: A struct with the following fields:
15  %        beta: A Kx1 vector of the estimated beta coefficients
16  %        s2: Estimate of the residual variance
17  %        stderr: A Kx1 vector of the estimators of the
        standard
18  %        errors of the beta estimates
19  %        t_stat: A Kx1 vector with the t-statistics for the
        null
20  %           hypothesis that each beta is zero (separately)
21  %        p_value: A Kx1 vector with the p-values for the t-
        test
22  %
23  % Examples:
24  %     y = [1;2;3];
25  %     X = rand(3);
26  %     beta = linreg_ols(y, X, true)
27  %
28  % Reference:
29  % Equation 1.2.5 in Hayashi, F. (2010).
30  % Econometrics. Princeton University Press. isbn:
        9780691010182.
31  if add_intercept
32      X = horzcat(ones(size(X, 1), 1), X);
33  end
```

```matlab
34  % compute inv(X'X) just once (used multiple times)
35  XX = X'*X;
36  % equivalent to inv(X'X) but faster and more numerically
        stable
37  XXinverse = XX\eye(length(XX));
38  % number of observations and explanatory variables
39  n = size(X, 1);
40  K = size(X, 2);
41  beta = XXinverse*X'*y;
42  % estimate Var(beta_hat|X)
43  e = y - X*beta;
44  s2 = (e'*e)/(n - K);
45  stderr = diag((s2*XXinverse)^0.5);
46  % t-statistics
47  t_stat = beta./stderr;
48  p_value = (1 - tcdf(abs(t_stat), n - K))*2;
49  % create results struct
50  results = struct('b', beta, 's2', s2, 'stderr', stderr, ...
51                  't_stat', t_stat, 'p_value', p_value);
52  end
```

While a `struct` can be very useful to organize data, it is not the fastest data type available in Matlab and should not be overused.

# 6 Validating Input and Debugging

The `linreg_ols` function is well documented and provides the bare bones for linear regressions. However, it is missing input validation. What should `linreg_ols` return if the matrix $X$ has more columns than lines? Or if $y$ and $X$ do not have the same number of rows?

## 6.1 Validating Input

We can validate data with conditionals, and if something fails the validation, we throw an error. In Matlab, we can throw an error with the function `error`, which takes a string as input and displays it to the user while interrupting the execution of the code. We can add the following lines to the beginning of the `linreg_ols` function:

```matlab
1   % linreg_ols.m
2   function results = linreg_ols(y, X, add_intercept)
3   %% Input validation
4   if size(y, 1) ~= size(X, 1)
5       error("y and X have a different number of rows")
6   end
7   if size(X, 1) < size(X, 2)
8       error("X has less rows than columns")
9   end
10  %% Code to run after inputs are validated ...
11  end
```

If the inputs fail the validation, the code is interrupted and a message is displayed in red so that the user can take action to fix the issues. It is important that the message is clear and identifies the issue as directly as possible.
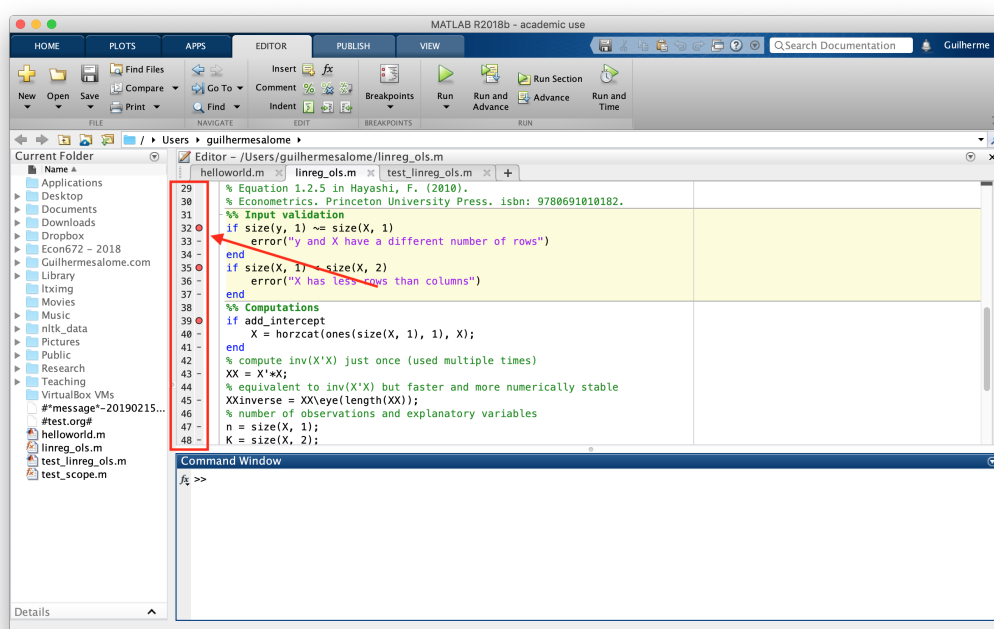
Observe the use of two percentage symbols instead of one. Lines that start with `%%` are also comments, but Matlab consider these lines as the start of a section. A section is just a tool to organize code in `.m` files and has no effect on the code.

Now, what happens if we forget to specify whether we want an intercept? That is, we execute `linreg_ols(y, X)`. In this case, `add_intercept` is a required input, so Matlab will warn you that there are not enough input arguments. We can see what happens in detail by using the `debugging` features Matlab provides.

## 6.2   Debugging

Debugging in the Matlab IDE is done by setting breakpoints, and then executing your code normally. When the code reaches a breakpoint, the execution is paused and you have access to the variables available at that point in your code. You can then choose to continue executing the code in different ways. You can continue the execution line by line, continue until where the cursor is at, or continue until the code finishes or finds another breakpoint.
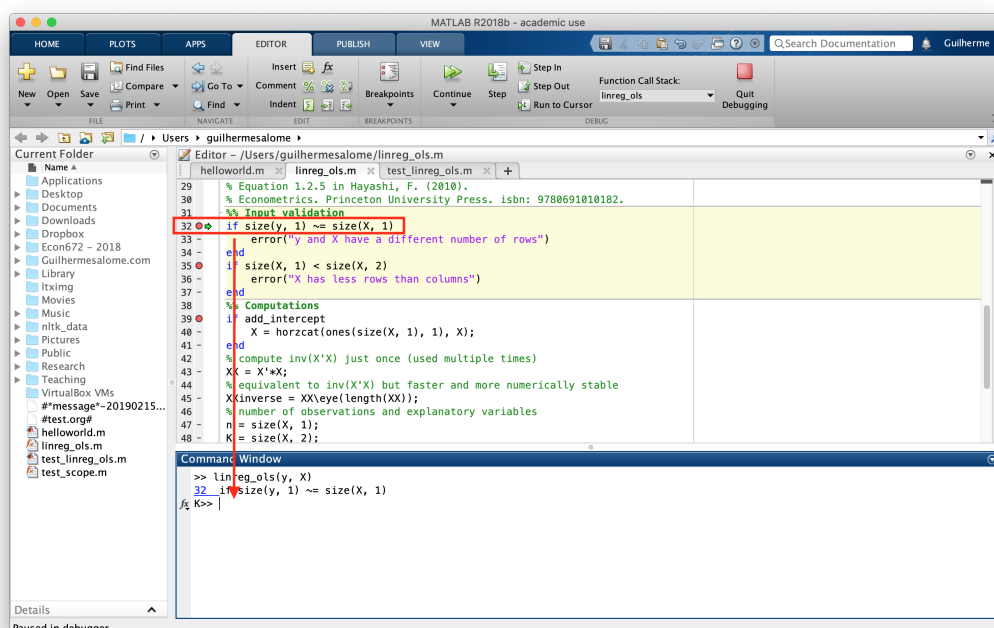
You can add a breakpoint by clicking right next to the number line in the code editor. See Figure 4 on how to add a breakpoint. When you click a red circle will show up indicating that a breakpoint was added. You can remove breakpoints by clicking on the red circles.



**Figure 4:**   Debugging in Matlab by Adding Breakpoints.

Now, in the `Command Window` execute `linreg_ols(y, X)`. The code will start executing and it will stop at the first breakpoint. A green arrow will indicate where the

execution paused, and at the `Command Window` you will be able to run any code and inspect the values of the variables. See Figure 5 where the code was paused at the first breakpoint.



**Figure 5:**  Debugging: Execution Paused at Breakpoint.

Notice that even though we called `linreg_ols` without supplying `add_intercept`, the function still executes. You can continue the execution by clicking on the button "Step", see Figure 6.

Keep clicking on "Step" until the code is paused at the third breakpoint. At this time, the variable `add_intercept` is not defined, and if you continue the execution you will get an error message.

Debugging is useful when dealing with large amounts of code and provides a way to ascertain the value of variables in your program. The code for `linreg_ols` will run even if the `add_intercept` variable is not passed, but will raise an error as soon as this variable is required. We can modify our code to allow `add_intercept` to be truly an optional variable, so that if it is omitted some default behavior will occur.

## 6.3   Optional Variables and Default Values

In Matlab, we can add optional variables to a function by passing the keyword `varargin` as an input. Then, inside the function, a variable named `varargin` is available and contains all optional arguments passed to it. Since optional arguments could be of any type, `varargin` must be able to hold different types of data.

In Matlab, a `cell array` is a data type that can hold data of various data types, and is the data type used for `varargin`. A `cell array` is similar to a number array and is indexed in the same way, but each element of a `cell array` is a `cell`, and a `cell` can have any data type. That is, a single `cell` can hold a number, a character, a string, or even other arrays. We can create a `cell array` with braces `{ }`:
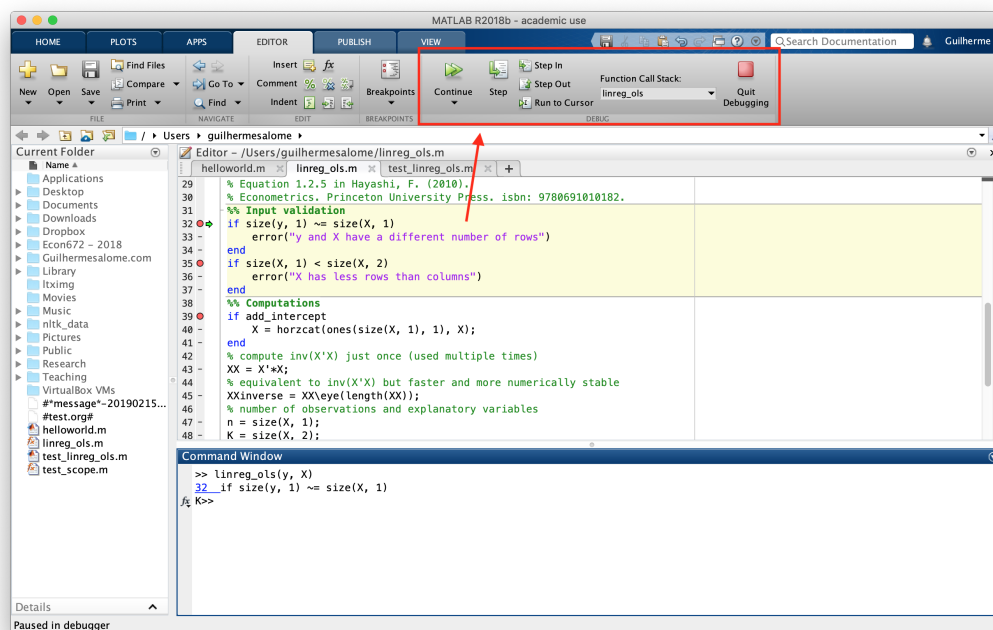
**Figure 6:** Debugging: Continue Execution.

```
1  % create a cell array
2  data = {"Parameters", [1, 2, 3]}
3  % use () to refer to cells, not their contents
4  data(1, 1)
5  data(1, 2)
6  % to access the actual contents of a cell index with {}
7  data{1, 1}
8  data{1, 2}
```

We can modify `linreg_ols` to accept optional arguments using `varargin`. Inside `linreg_ols`, the variable `varargin` is a cell array holding any extra inputs passed to the function. If no extra input was passed, then `varargin` is an empty cell array. We can check whether `varargin` is empty with the function `isempty`. To let the default behavior of the function be adding the intercept, we can verify if `varargin` is empty, and if it is we add the intercept.

```
1  % linreg_ols.m
2  function results = linreg_ols(y, X, varargin)
3  %% Input validation
4  if size(y, 1) ~= size(X, 1)
5      error("y and X have a different number of rows")
6  end
7  if size(X, 1) < size(X, 2)
8      error("X has less rows than columns")
9  end
10 %% Add Intercept if Required
11 if isempty(varargin) || (varargin{1} == true):
```

```matlab
12        X = horzcat(ones(size(X, 1), 1), X);
13 end
14 %% Computations
15 % compute inv(X'X) just once (used multiple times)
16 XX = X'*X;
17 % equivalent to inv(X'X) but faster and more numerically
      stable
18 XXinverse = XX\eye(length(XX));
19 % number of observations and explanatory variables
20 n = size(X, 1);
21 K = size(X, 2);
22 beta = XXinverse*X'*y;
23 % estimate Var(beta_hat|X)
24 e = y - X*beta;
25 s2 = (e'*e)/(n - K);
26 stderr = diag((s2*XXinverse)^0.5);
27 % t-statistics
28 t_stat = beta./stderr;
29 p_value = (1 - tcdf(abs(t_stat), n - K))*2;
30 % create results struct
31 results = struct('b', beta, 's2', s2, 'stderr', stderr, ...
32                  't_stat', t_stat, 'p_value', p_value);
33 end
```

It is possible to tie together optional inputs and input validation.

## 6.4   Parsing Inputs

Matlab has a built-in way to do input parsing and allowing optional inputs: `inputParser`. The `inputParser` is an object created inside a function, which takes care of parsing inputs.

We use the `inputParser` by first creating a `parser object`. Then, we tell this `parser object` what are the inputs the function should receive. For example, we tell the `parser object` it <u>must</u> receive a variable named $y$, it <u>must</u> receive a variable named $X$, and that it <u>could</u> receive a variable named `intercept`. We can also tell the `parser object` what conditions each variable must satisfy for it to be accepted and the function to run normally. We then tell the `parser object` to parse the inputs, and the parsed inputs will be available in the `parser object` as a `struct`. Let's do an example:

```matlab
1 % example_inputparser.m
2 function example_inputparser(arg1, arg2, varargin)
3 % create the parser object
4 parser = inputParser;
5 % add a required argument to the parser: arg1
6 addRequired(parser, 'arg1');
7 % add a required argument to the parser: arg2
8 addRequired(parser, 'arg2');
9 % add an optional argument to the parser
10 addOptional(parser, 'intercept', true);
```

```matlab
11  % actually parse the inputs
12  parse(parser, arg1, arg2, varargin{:});
13  inputs = parser.Results;
14  % display the inputs
15  disp({inputs.arg1, inputs.arg2, inputs.intercept});
16  end
```

The `inputParser` code above is quite more extensive than what we had before, but the code itself is also very clear on what it is doing (it is self-documenting). We create the `parser` object, and add two required arguments with the `addRequired` function. The `addRequired` function informs the `parser` object that it must receive two variables, the first variable should receive the name `arg1` and the second variable should receive the name `arg2`. Then, we call `addOptional` to inform the `parser` object that it can receive a third variable, and if it does, this variable should be named `intercept`. If no third variable is received, then the variable `intercept` should be created and its value should be `true`. Test the function by calling:

```matlab
1  example_inputparser(1, 2)
2  example_inputparser(1, 2, false)
3  example_inputparser()
```

The functions `addRequired` and `addOptional` define inputs that must appear in the correct order (positional arguments). While `addRequired` defines required inputs, `addOptional` defines optional inputs, but both inputs are positional arguments. However, we can also define inputs that can come in different order, but to do so these inputs must also be named (name-value pair arguments). We can add these name-value pair inputs with the `addParameter` function:

```matlab
1  % example_inputparser.m
2  function example_inputparser(arg1, arg2, varargin)
3  % create the parser object
4  parser = inputParser;
5  % add a required argument to the parser: arg1
6  addRequired(parser, 'arg1');
7  % add a required argument to the parser: arg2
8  addRequired(parser, 'arg2');
9  % add an optional argument to the parser
10 addOptional(parser, 'intercept', true);
11 % add a name-value pair input to the parser
12 addParameter(parser, 'cov_type', 'standard');
13 % actually parse the inputs
14 parse(parser, arg1, arg2, varargin{:});
15 inputs = parser.Results;
16 % display the inputs
17 disp({inputs.arg1, inputs.arg2, inputs.intercept, inputs.
       cov_type});
18 end
```

The order for inputs is always the following: required arguments, optional arguments, name-value pairs.

The advantage of using `inputParser` is that we can also do input validation at the same time. Each of the functions `addRequired`, `addOptional` and `addParameter` can take a validation function that is used to check whether the inputs are valid. Let's add input validation to `example_inputparser`:

```matlab
% example_inputparser.m
function example_inputparser(arg1, arg2, varargin)
% create the parser object
parser = inputParser;
% add a required argument to the parser: arg1
addRequired(parser, 'arg1', @validate_arg1);
% add a required argument to the parser: arg2
addRequired(parser, 'arg2', @validate_arg2);
% add an optional argument to the parser
addOptional(parser, 'intercept', true, @validate_intercept);
% add a name-value pair input to the parser
addParameter(parser, 'cov_type', 'standard', @
    validate_cov_type);
% actually parse the inputs
parse(parser, arg1, arg2, varargin{:});
inputs = parser.Results;
% display the inputs
disp({inputs.arg1, inputs.arg2, inputs.intercept, inputs.
    cov_type});
end

%% Functions for input validation
% these functions must return a boolean value or throw an
    error
function bool = validate_arg1(arg1)
    bool = isnumeric(arg1);
end

function bool = validate_arg2(arg2)
    bool = isnumeric(arg2) && size(arg2, 1) >= size(arg2, 2);
end

function bool = validate_intercept(intercept)
    bool = islogical(intercept) && length(intercept) == 1;
end

function bool = validate_cov_type(cov_type)
    covs = ["standard", "White", "HAC"];
    bool = ischar(cov_type) && ismember(cov_type, covs);
end
```

The functions declared in the section "Functions for input validation" are passed as arguments to the functions `addRequired`, `addOptional` and `addParameter`. To pass a function as an argument to another function we use the `@` operator in front of the func-

tion name. When we use the `@` syntax with the function name we are creating a function handle. A function handle allows a function to be passed around as any other variable, and calling a function handle is the same as calling the function itself.

Notice that most of the functions above have just one line of code. We can simplify the code by using anonymous functions.

## 6.5   Anonymous Functions

An anonymous function is a function that is written in a single line and is itself a function handle. We can create an anonymous function using the syntax `@(inputs) code;`. For example:

```matlab
% create an anonymous function and assign to a variable
square_root = @(x) x.^0.5;
% call the anonymous function
square_root(4)
% since the anonymous function is a function handle
% we can pass it as inputs to other functions
integral(square_root, 0, 1)
% we can create an anonymous within an anonymous function
poly_integral = @(a) integral(@(x) (a*x.^2 + x + 2), 0, 1);
```

The validating functions in `example_inputparser` can be substituted by anonymous functions, which makes the code clearer. Let's modify `example_inputparser` to make use of anonymous functions:

```matlab
% example_inputparser.m
function example_inputparser(arg1, arg2, varargin)
% create the parser object
parser = inputParser;
% add a required argument to the parser: arg1
addRequired(parser, 'arg1', @isnumeric);
% add a required argument to the parser: arg2
addRequired(parser, 'arg2', @(x) (isnumeric(x) && size(x, 1)
    >= size(x, 2)));
% add an optional argument to the parser
addOptional(parser, 'intercept', true, @(x) (islogical(x) &&
    length(x) == 1));
% add a name-value pair input to the parser
addParameter(parser, 'cov_type', 'standard', @
    validate_cov_type);
% actually parse the inputs
parse(parser, arg1, arg2, varargin{:});
inputs = parser.Results;
% display the inputs
disp({inputs.arg1, inputs.arg2, inputs.intercept, inputs.
    cov_type});
end

%% Functions for input validation
```

```matlab
21  % these functions must return a boolean value or throw an
        error
22  function bool = validate_cov_type(cov_type)
23      covs = ["standard", "White", "HAC"];
24      bool = ischar(cov_type) && ismember(cov_type, covs);
25  end
```

We can make use of these features to validate the input for `linreg_ols`:

```matlab
 1  % linreg_ols.m
 2  function results = linreg_ols(y, X, varargin)
 3  %% Input validation
 4  parser = inputParser;
 5  addRequired(parser, 'y', @isnumeric)
 6  addRequired(parser, 'X', @validate_X);
 7  addParameter(parser, 'intercept', true, @islogical);
 8  parse(parser, y, X, varargin{:});
 9  y = parser.Results.y;
10  X = parser.Results.X;
11  if size(y, 1) ~= size(X, 1)
12      error("y and X have a different number of rows")
13  end
14  %% Add Intercept if Required
15  if parser.Results.intercept
16      X = horzcat(ones(size(X, 1), 1), X);
17  end
18  %% Computations
19  % compute inv(X'X) just once (used multiple times)
20  XX = X'*X;
21  % equivalent to inv(X'X) but faster and more numerically
        stable
22  XXinverse = XX\eye(length(XX));
23  % number of observations and explanatory variables
24  n = size(X, 1);
25  K = size(X, 2);
26  beta = XXinverse*X'*y;
27  % estimate Var(beta_hat|X)
28  e = y - X*beta;
29  s2 = (e'*e)/(n - K);
30  stderr = diag((s2*XXinverse)^0.5);
31  % t-statistics
32  t_stat = beta./stderr;
33  p_value = (1 - tcdf(abs(t_stat), n - K))*2;
34  % create results struct
35  results = struct('b', beta, 's2', s2, 'stderr', stderr, ...
36                  't_stat', t_stat, 'p_value', p_value);
37  end
38
39  %% Functions for Input Validation
```

```matlab
40  function bool = validate_X(X)
41      bool = isnumeric(X);
42      if size(X, 1) < size(X, 2)
43          error("X has less rows than columns");
44      end
45  end
```

We can now call `linreg_ols` with or without a name-value pair specifying the inclusion of an intercept:

```matlab
1  linreg_ols(y, X)
2  linreg_ols(y, X, 'intercept', true)
```

# 7  Importing Data and Saving Results

We will now use our `linreg_ols` function with a real data set. Download this data set on housing prices in California, and save the file in your working folder. This data set was first used in Pace and Barry (1997), and was later modified for use in Géron (2017). The data contains information on household prices in California in 1990 (based on the census at that time). The first line of the file contains the name of the explanatory variables.

## 7.1  Importing Data

Matlab has various built-in functions for importing data of different types. We can import `.csv` files to Matlab with the `readmatrix` function[4].

```matlab
1  data = readmatrix('housing.csv');
```

The default behavior of `readmatrix` is to infer the type of the file from its extension, in this case `.csv`. Knowing the file extension, the function knows that the values are comma-separated, and the function also attempts to detect if there is a header line and the number of variables (columns) to be imported. The `housing.csv` data is mostly numeric, with the exception of the last column. The function `readmatrix` imports the data into a numeric matrix, which works well for the first 9 columns of `housing.csv`. However, the last column of the file contains categorical data, which is not automatically converted to numbers, and so Matlab treats all its values as `NaN`. If we are not interested in using the last column we could work with `readmatrix`.

The resulting `data` matrix has all the information from `housing.csv`, with the exception of the names of the variables and the last column of the file. We can keep all the columns of the data and the names of the variables if we import the data to a `table` instead of a numeric array. To do so, we use the function `readtable`. As with `readmatrix`, the function `readtable` can figure out the delimiter, header and number of variables from the file extension, but instead of loading the data into a numeric array, it loads the data into a `table`.

---

[4]There is also a GUI interface for importing data under the menu tab `Home`, via the button `Import Data`. Older versions of Matlab have different functions for importing data. Before the R2019a release, for example, `.csv` files could be imported with the function `csvread`. You can still use `csvread`, but `readmatrix` should be preferred since it is an improved version.

```matlab
1  % import housing data
2  data = readtable('housing.csv');
3  % verify that everything was imported correctly
4  data(1:5, :)
5  % notice that the names of the variables are also
6  % displayed, and that the last is correctly imported
7
8  % data in a table can be accessed by normal indexing as in a
9  % numeric array
10 data(1:10, 1)
11 data(1:10, 2)
12 data(1:10, 10)
13 % but we can also access the data like in a struct, where the
       field
14 % names are the variable names
15 data.longitude
16 data.longitude(1:5)
17 data.ocean_proximity(1:5)
18 % notice that the ocean_proximity data is stored as character
19 % arrays in cells
20 % we can convert it to a string array:
21 data.ocean_proximity = string(data.ocean_proximity);
22 data.ocean_proximity(1:10)
23 % we can see all the variable names in a table via its
       properties
24 data.Properties.VariableNames
25 % we can access the columns of the table by index with the
       names
26 data(1:5, 'longitude')
27 % for multiple columns we pass a string array
28 data(1:5, ["population", "households", "median_income"])
29 % similar to cells, when we index with parentheses we do not
       get the
30 % actual values stored in a table, but we get a subtable
31 % to get the actual values we use the curly brackets notation
32 data{1:5, ["population", "households", "median_income"]}
```

The data has a few missing values (added on purpose). We can count how many lines have missing values with the functions isnan and sum. The function isnan takes a numeric array and finds all NaN. The function sum can be used to sum across values in rows or columns (or both) of a matrix.

```matlab
1  % dimensions of table
2  size(data)
3  % find missing data no the first column
4  sum(isnan(data{:, 1}))
5  % find missing data on all columns, except the last one
6  sum(isnan(data{:, 1:9}))
```

```matlab
7  % we see that there are 207 missing data points for the
8  % total_bedrooms variable
9  % indeed:
10 sum(isnan(data.total_bedrooms))
11 % let's remove these rows from the table
12 % first we need to find the indices for the rows with missing
       values
13 missing_values = isnan(data.total_bedrooms);
14 % missing_values is a boolean array, it has 1's where there
       are
15 % NaNs in the data.total_bedrooms array
16 % we can use this logical (boolean) array to index the
       original
17 % table and select only the rows without NaN's
18 data_no_missing = data(~missing_values, :);
19 % this gives us a new table, from which we can extract data
       for the regression
```

We can also remove rows with missing values (represented by `NaN`) automatically using the function `rmmissing`:

```matlab
1  size(data)
2  data_no_missing = rmmissing(data);
3  size(data_no_missing)
```

Let's run a regression of the logarithm of the median house value on the logarithm of the number of bedrooms in the house.

```matlab
1  results = linreg_ols(log(data.median_house_value), ...
2                       log(data.total_bedrooms));
3  disp(results.b);
```

## 7.2   Saving Results, Cleaning the Workspace and Loading Variables

Now that we can generate results from data, we need to save these results. Saving and loading variables in Matlab can be done with the functions `save` and `load`. There are several ways of using the function `save`:

```matlab
1  % Save ALL variables in the workspace
2  save('all_workspace_variables');
3  % The above creates a file named all_workspace_variables.mat
4  % .mat files can be easily imported with the load function
5
6  % Save only ONE variable
7  save('housing_clean', 'data_no_missing');
8  % saves the variable data_no_missing on the file
       housing_clean.mat
9
10 % Save MULTIPLE variables
```

```
11  save('data_and_results', 'data_no_missing', 'results');
12  % saves the variables data_no_missing and results in
13  % the file data_and_results.mat
```

The files created with the function `save` have the `.mat` extension and are binary files.

Let's now clear our workspace to and reload the variables. First, clear the command window by executing `clc`. This command simply clears the `Command Window` so that you have a clear screen to type commands. You can still use the arrow keys to execute previous commands. Second, use the special command `clear` to delete all variables from the workspace. After executing this command the variables previously created are gone. It is also possible to clear specific variables by executing `clear variable_name`. Third, use `load` to load the `.mat` files, which will take the variables stored in those files and make them available in the workspace.

```
1   % Load ALL variables from before
2   load('all_workspace_variables.mat');
3   % Load only the HOUSING variable from before
4   clear;
5   load('all_workspace_variables.mat', 'housing');
6   % Notice we pass the name of the variable to load so that
        only that
7   % variable is loaded
8   % Load the cleaned housing data
9   clear;
10  load('housing_clean.mat');
11  % Load the cleaned housing data and the regression results
12  clear;
13  load('data_and_results.mat');
14  % Load only the regression results
15  clear;
16  load('data_and_results.mat', 'results');
```

Notice that the variables are loaded with their original names.

# 8   Summary

We have covered the basics of the Matlab programming language. You should be able to:

- Create `.m` script files to run code

- Work with vectors and matrices

- Do basic computations and linear algebra with numeric arrays

- Create function files

- Use conditionals

- Use for-loops and while-loops

- Generate data from statistical distributions

- Create arrays for text data

- Present data in a neat manner using `table`

- Organize function results in `struct`

- Validate input variables with conditionals

- Debug code with the built-in debugger

- Let functions accept any number of inputs with `varargin`

- Use `inputParser` to parse input variables and specify default values

- Pass functions as arguments with function handles

- Simplify small functions with anonymous functions

- Import `.csv` data

- Clean rows from a table with missing values

- Save variables to `.mat` files

- Clear the workspace

- Load variables stored in `.mat` files

# 9 Assignment

All assignments should be submitted to the Github repository you have been assigned to. The deadline is July 26th by midnight. You should write a report in Latex with the solutions to the problems below. If the problem requires you to code, then the code should also be included in the report. For a quick guide on how to add Matlab code to your Latex files, refer to the Section "Adding Matlab Code" on this tutorial.

**Problem 1** *Update the code for generating data to function without loops (vectorization). To do so, you will need to read the documentation for the random number generator functions we have used: `rand` and `randn`. After vectorizing the problem, verify you get the same estimates for the parameters as before. For this part, you will need to fix the generation of random numbers, so that you generate the same random numbers when you run your code without vectorization and with vectorization. Read the documentation for the `rng` function, and if you need more details read about managing the global stream.*

**Problem 2** *Extend `linreg_ols` to also output the $R^2$ of the regression (see Equation 1.2.18 in Hayashi (2000).*

**Problem 3** *Run the linear regression suggested in Equation (8) of Pace and Barry (1997).*

**Problem 4** *Matlab has been around since 1984 and has many packages for common problems. Figure out if Matlab has a function for estimating the parameters of a linear regression using OLS. If it does, what is the name of the function? How does the function actually estimate the coefficients? To answer this, you need to look at the code of the function. You can do so with the special command* **type**. *For example, if the name of the function is* **ols**, *then you can execute* **type ols** *to see the source code.*

**Problem 5** *When the classical assumption of homoskedasticity fails, we need a different estimators for the standard errors of the OLS estimates. White (1980) proposes a heteroskedasticity-robust estimator for the standard errors of the OLS estimates, which is now known as White's standard error. Equation 2.4.1 in Hayashi (2000) shows White's standard error:*

$$\widehat{SE(\hat{\beta}_i)} \equiv \sqrt{\frac{1}{n}[S_{xx}^{-1}\hat{S}S_{xx}^{-1}]_{[i,i]}}$$

*There is a slight change of notation in this part of the Hayashi (2000), and $x_i'$ is the ith row of $X$ ($x_i$ is a column vector with the explanatory variables in the ith row of $X$). The term $S_{xx}$ is the sample mean of $x_i x_i'$: $S_{xx} = \frac{1}{n}\sum_{i=1}^{n} x_i x_i'$ (Equation 2.3.6 in Hayashi (2000)). The term $\hat{S}$ is an estimator for a matrix of fourth moments, and it is defined as $\hat{S} = \frac{1}{n}\sum_{i=1}^{n} e_i^2 x_i x_i'$ (Equation 2.5.1 in Hayashi (2000)), and $e_i$ is as before (residual for the ith observation). Implement this estimator in the* **linreg_ols** *function. You could write a local function to estimate White's standard error in the* **linreg_ols.m** *file.*

**Problem 6** *(Optional) Implement the t-test based on White's standard error. Also, compute the p-value from given the statistic.*

**Problem 7** *Newey and West (1987) proposes another estimator to the standard errors for the OLS estimates under weaker assumptions. The estimator the authors propose is robust not only to heteroskedasticity, but also autocorrelation, and is known as the HAC standard errors. Does Matlab have a function to compute the HAC standard errors? If so, extend* **linreg_ols** *to also compute the HAC standard errors.*

**Problem 8** *(Optional) Implement the t-test based on the HAC standard error. Also, compute the p-value from given the statistic.*

**Problem 9** *Extend* **linreg_ols** *to accept an optional input named* **cov_type**, *which specifies the type of standard errors estimator to report. Use* **inputParser** *to do so.*

# References

Géron, Aurélien (2017). *Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems*. " O'Reilly Media, Inc.". URL: `https://isbnsearch.org/isbn/9781491962299`.

Hayashi, F. (2000). *Econometrics*. Princeton University Press. ISBN: 9780691010182. URL: `https://books.google.com/books?id=QyIW8WUIyzcC`.

Newey, Whitney K. and Kenneth D. West (1987). "A simple, positive semi-definite, heteroskedasticity and autocorrelation consistent covariance matrix". In: *Econometrica* 55.3, p. 703. URL: `https://login.proxy.lib.duke.edu/login?url=https://search-proquest-com.proxy.lib.duke.edu/docview/214867593?accountid=10598`.

Pace, R Kelley and Ronald Barry (1997). "Sparse spatial autoregressions". In: *Statistics & Probability Letters* 33.3, pp. 291–297. URL: `https://doi.org/10.1016/S0167-7152(96)00140-X`.

White, Halbert (1980). "A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity". In: *Econometrica* 48.4, pp. 817–838. URL: `https://www.jstor.org/stable/1912934`.