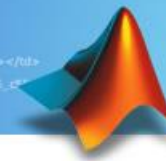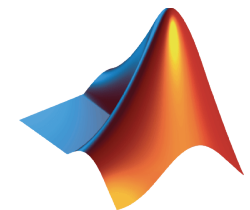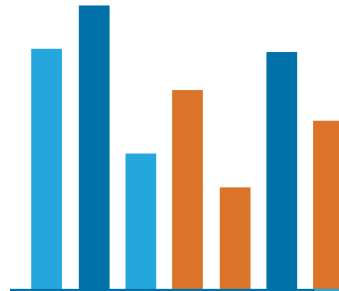# MATLAB for Data Analysis
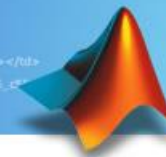
Application Engineer

Jeffrey Liu

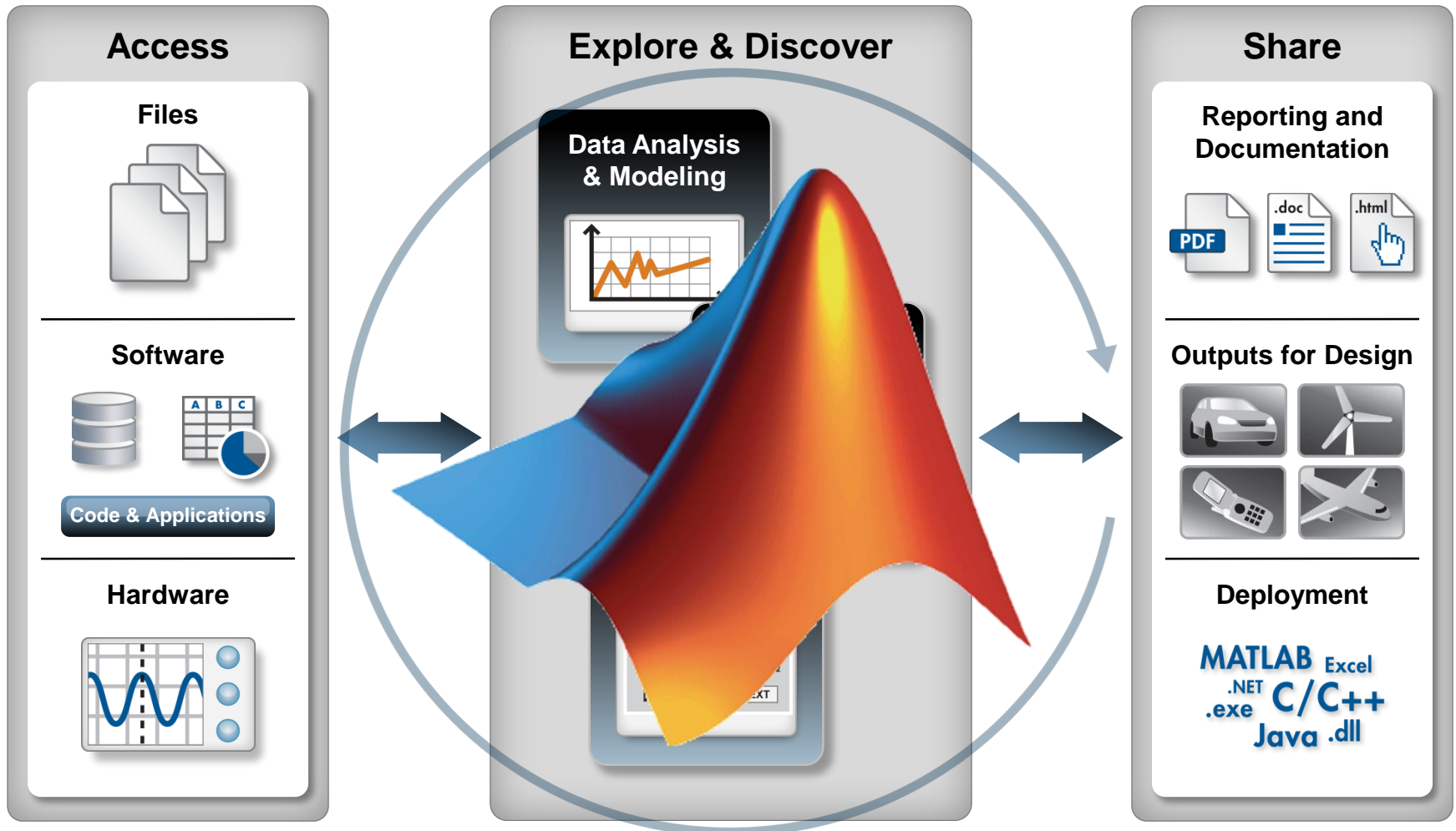# Outline

- Importing Data to MATLAB
  - Create data in MATLAB
  - Working with data files

- Data Types
  - Dataset, table
  - data types

- Data pre-processing
  - Dealing with missing value
  - Locating data
  - Merging data

- Basic data analysis
  - Descriptive statistics
  - Split-apply workflow

- Fit data to models
  - Curve fitting toolbox
  - Statistics models

- Publish your code

# Random Number generation

MATLAB®

- **MATLAB**
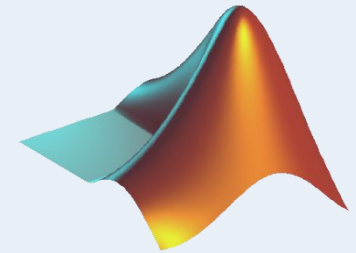  - rand
    - uniform distribution generation between 0~1
  - randi
    - Uniformly distributed pseudorandom integers
  - randn
    - Standard normal distribution generator
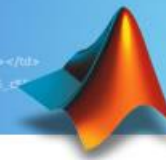- **Statistics Toolbox**
- **Random seed**
  - rng
  - rng shuffle

```
rand
randi
randn
```

Statistics Toolbox™

```
betarnd
binornd
chi2rnd
ncx2rnd
exprnd
evrnd
frnd
ncfrnd
gamrnd
...
```

`random`

MATLAB&SIMULINK

# Interactively Importing Text Files and Spreadsheets

▪ The Import Tool attempts to interpret the contents of spreadsheets and delimited text files, and automatically set several options.

# Reading Fixed structure files

- Files
  - Excel, text, csv, or binary
  - .xpt in SAS
  - .mat
  - Multimedia, scientific
  - Web, XML
- Command
  - load
  - csvread
  - xlsread
  - textread
  - textscan

# Low-Level File I/O

## 1. Open the file

>> fid = fopen('gasprices.csv','r');

## 2. Perform the file operations

Functions for reading and writing text and binary files are shown in the diagram to the right.

## 3. Close the file

>> fid = fclose(fid);

# Vectors, Matrices, and Arrays

- A table of numbers with *m* rows and *n* columns is referred to as an *m*-by-*n* matrix.

- A *vector* is a single row or a single column of numbers. It is therefore a special case of a matrix, where either *m* or *n* is equal to 1.

- A single number is referred to as a *scalar*. It is a 1-by-1 matrix and, equivalently, a 1-element vector.

**Array**:
multiple values in one variable

abc

**Matrix**:
regular 2-D array (*m*-by-*n*)

**Vector**:
1-D array (1-by-*n* or *m*-by-1)

**Scalar**:
single value

MATLAB&SIMULINK

# Cell array

- Cell arrays are an easy way to assemble data of dissimilar types and sizes into a single "container of containers."

- They are typically used to store strings of different length.

- To concatenate elements into a cell array,

  use curly braces ({}) instead of square brackets

# What Is a Table ?

- Many data sets conform to a particular tabular arrangement with the following properties:

- Each column of data has the same type (text, numeric, logical (T/F), etc.).
- Different columns, however, can have different types.
- Each column typically has a unique name.
- Each column has the same number of rows

MATLAB&SIMULINK

# Categorical array

- When text labels intended to represent finite set of possibilities, cell array in unnecessary, Instead, you can use *categorical* array

>> x = {'C', 'B', 'C', 'A', 'B', 'A', 'C'};

>> y = categorical(x);

x

| 'C' | 'B' | 'C' | 'A' | 'B' | 'A' | 'C' |

{ }

1-by-7

y

| C | B | C | A | B | A | C |

1-by-7

MATLAB&SIMULINK

# Extracting Data from a Table

- Index into variables in a table using dot (**.**) notation to reference the variables by name:

  ```
  >> dates = stockPrices.Date;
  >> admPrices = stockPrices.ADM;
  ```

- You can also use curly braces (**{}**) to index into variables in a table. You can index numerically

  ```
  >> twoFin = stockPrices{:,[4 5]};
  ```

  or by name
  ```
  >> twoFin = stockPrices{:,{'ADM','ADN'}};
  ```

| Date | AAL | ABF | ADM | ADN | AGK | ... |
|------|------|-------|--------|--------|--------|-----|
| '1/3/2006' | 2241.76 | 840.50 | 466.75 | 132.75 | 295.74 | |
| '1/4/2006' | 2251.65 | 857.50 | 481.00 | 133.00 | 291.10 | |
| '1/5/2006' | 2193.41 | 855.00 | 477.25 | 135.00 | 292.90 | |
| '1/6/2006' | 2235.17 | 857.00 | 477.25 | 135.50 | 289.03 | |
| '1/9/2006' | 2212.09 | 862.00 | 478.25 | 138.00 | 283.10 | |

**505-by-1**        **505-by-1**        **505-by-93**

**{}**

# Merging Data from two tables

- Use "join" to merge "dataset" variables
  - Import files with dataset array

  `>> Variables = dataset('XLSFile', 'filename');`
  - Merge two datasets

  `>> New = join(A,B, 'Keys', 'date', 'Type', 'outer', 'MergeKeys',true);`

- Use "innerjoin", "outerjoin" to combine "table" variables
  - Import files with table array

  `>> Variables = readrable(date, data, 'var_name');`
  - Merge two tables
    - ✓ innerjoin, outerjoin

MATLAB&SIMULINK

# Exercise 1

- Merging two tables by time in excel file.

# Dealing with missing value

- Possible strategies



Ignore          Delete          Replace

# Avoiding NaNs in Calculations

- Several functions are designed to ignore NaNs in calculations

```
nancov              nanmedian           nansum
nanmax              nanmin              nanvar
nanmean             nanstd
```

- if a column contains all NaNs, ignoring them will result in applying the desired function to an empty array.

| 3 | 2 | 4 | NaN |
|---|---|---|-----|
| 1 | NaN | 2 | 1 |
| NaN | NaN | 6 | 1 |

↓ nanmean

| 2 | 2 | 4 | 1 |
|---|---|---|---|

| 3 | NaN | 4 | NaN |
|---|-----|---|-----|
| 1 | NaN | 2 | 1 |
| NaN | NaN | 6 | 1 |

↓ nanmean

| 2 | NaN | 4 | 1 |
|---|-----|---|---|

MATLAB&SIMULINK

# Locating Missing(Other) Values

- MATLAB provides numerous "**is\***" functions that take an array as input and return a logical output that signifies if the input has a certain characteristic.



| false |
|-------|
| true |
| false |
| true |
| true |

**isfinite** ←→

| Inf |
|-----|
| 3 |
| NaN |
| 1 |
| 7 |

**isnan** →

| false |
|-------|
| false |
| true |
| false |
| false |

- You can use logical indexing to remove elements from an array.

```
>> idx = isnan(x);
>> x(idx) = [];
```

**Numerical comparison**

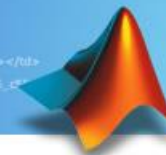| 3 | NaN | 4 | NaN |
|-----|-----|---|-----|
| 1 | NaN | 2 | 1 |
| NaN | NaN | 6 | 1 |

`>> x == NaN`

**Not a Number!**

| F | F | F | F |
|---|---|---|---|
| F | F | F | F |
| F | F | F | F |

✖

| 3 | NaN | 4 | NaN |
|-----|-----|---|-----|
| 1 | NaN | 2 | 1 |
| NaN | NaN | 6 | 1 |

`>> isnan(x)`

| F | T | F | T |
|---|---|---|---|
| F | T | F | F |
| T | T | F | F |

✔

`>> all(isnan(x))`

| F | T | F | F |
|---|---|---|---|

| 3 |
|---|
| 1 |
| 7 |

**x**

NaN

MATLAB&SIMULINK

# Replacing Missing Values in Matrices

# Test Your Knowledge

1. (Select all that apply.)  If x is a (numeric) vector, which commands will remove all the NaNs from x?

    A. x = x(isnan(x));

    B. x = x(x ~= NaN);

    C. x(isnan(x)) = [];

    D. x(x == NaN) = [];

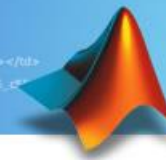    E. x = x(isfinite(x));

2. Suppose x is a 6-by-3 (numeric) matrix of ones.  Three values in the first column are missing (NaN), as are six values in the third column. What will >> nanmean(x) return?
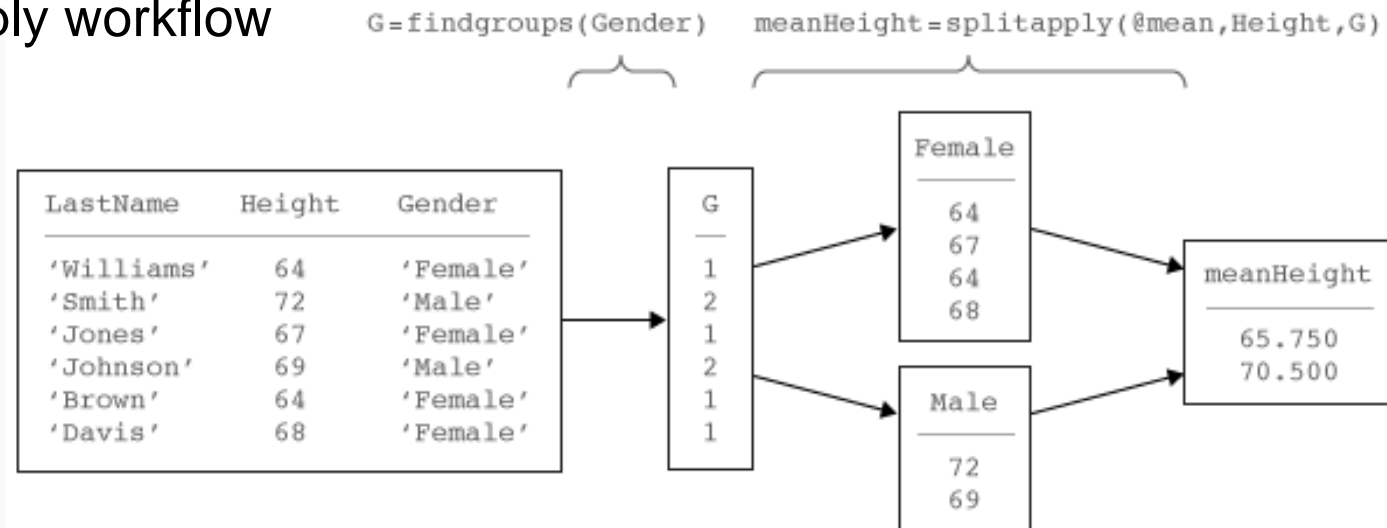
    A. 1 B. [0.5 1 0]

    C. [0.5 1 NaN]

    D. [1 1 NaN]

    E. [1 1 0]

    F. An error message.

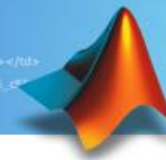# Calculate the data by group

- Function 'grpstats' : Summary statistics organized by group
  - Syntax : statarray = grpstats(tbl,groupvar,whichstats)
    tbl : data in table or dataset array
    groupvar : Column name for grouping in tbl
    whichstats : Types of summary statistics ( numel, std, max, sum…)
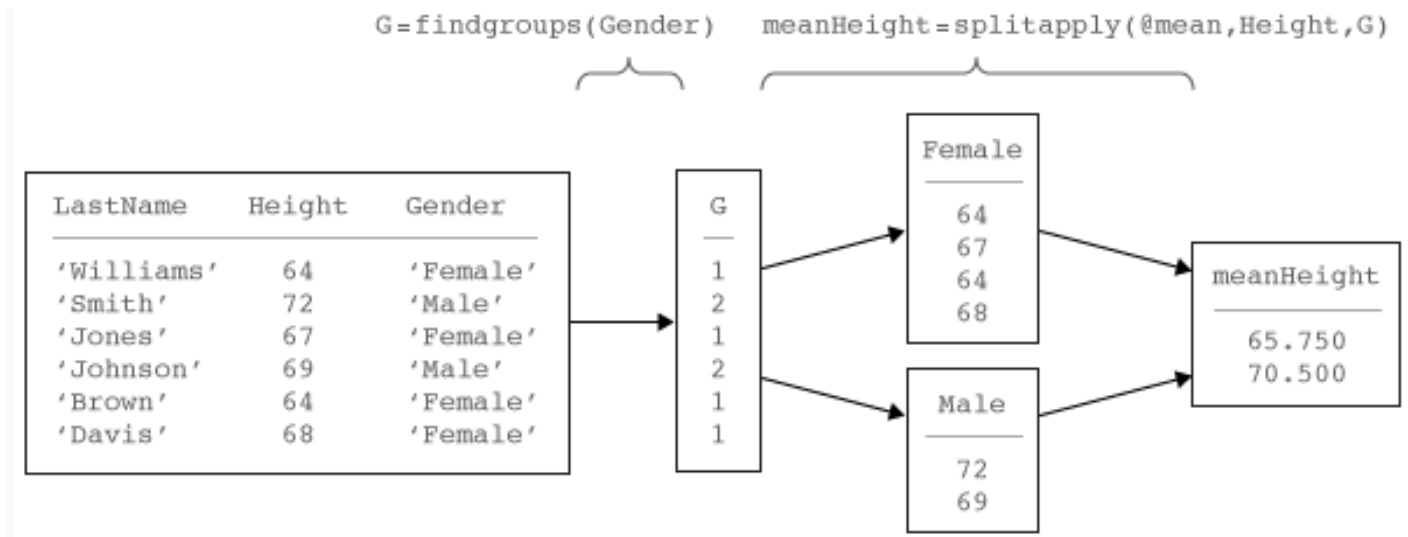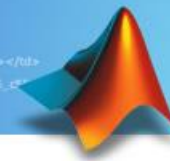
- New Split-apply workflow

# Exercise 2

- 在bank-full.csv中依照三種類別(工作、婚姻、購買與否)區分客戶,並計算類別中的平均收入

# LinearModel Variables

- Given data in vectors `x` (predictor) and `y` (response), you can perform a least-squares linear regression using `fitlm`:

  ```
  >> linmodel = fitlm(x,y)
  ```

- You can visually inspect the fit using the `plot` method:
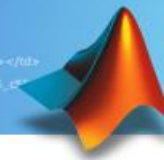
  ```
  >> plot(linmodel)
  ```

- Evaluate the fitted model at chosen predictor values using the **predict** method:

  ```
  >> predict(linmodel,xnew)
  ```
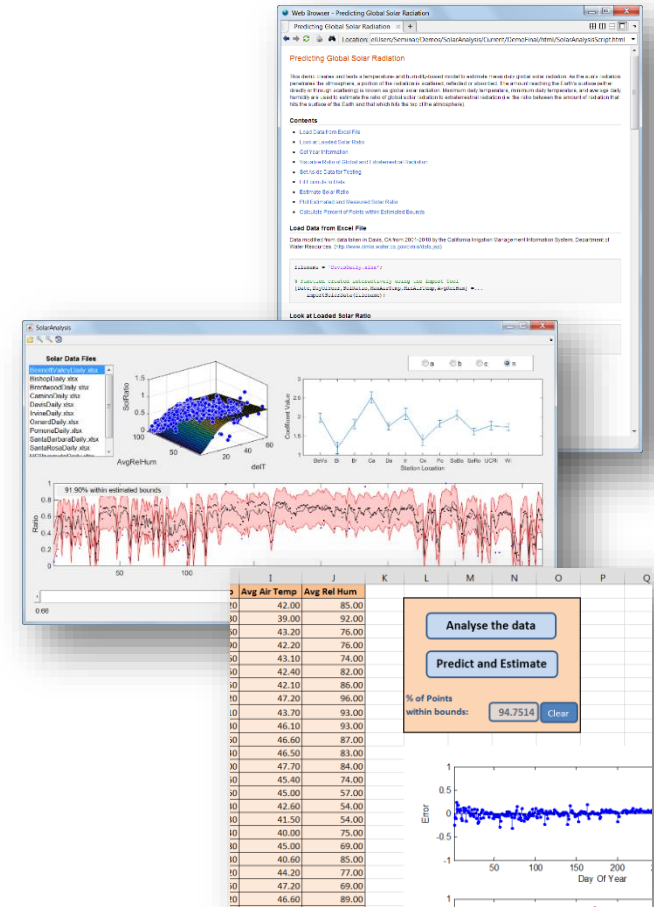
- This value is one of the properties of a **LinearModel** variable:
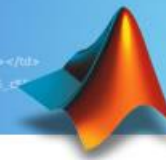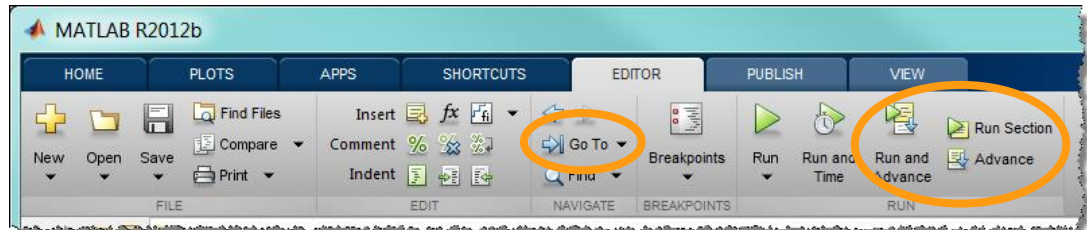
  ```
  >> linmodel.Rsquared
  ```

MATLAB&SIMULINK
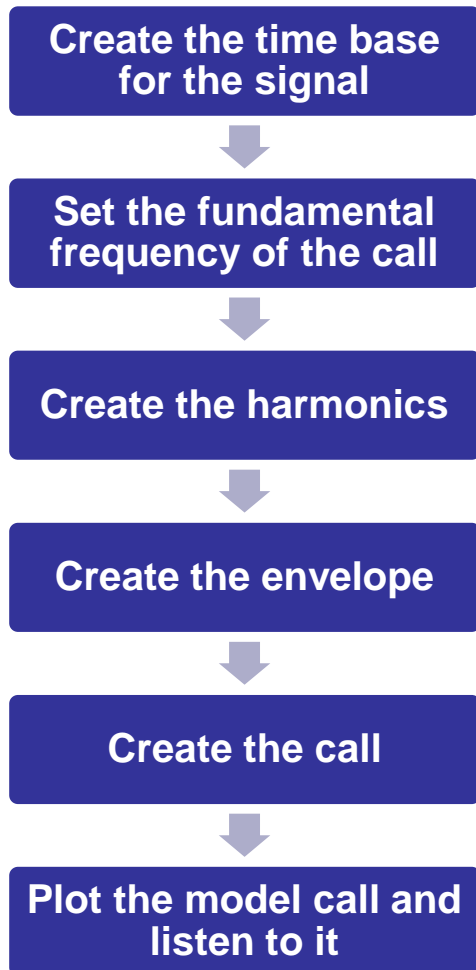
# Sharing Results from MATLAB

- Automatically generate reports

- Create and package applications

- Deploy to other environments

# Code Sections

**Create the time base for the signal**

⬇

**Set the fundamental frequency of the call**

⬇

**Create the harmonics**

⬇

**Create the envelope**

⬇

**Create the call**

⬇

**Plot the model call and listen to it**

## %%

```
%% Set the fundamental frequency of the call.
f0 = 175;

%% Create the harmonics.
y0 = sin(2*pi*f0*t) + sin(2*pi*2*f0*t) + sin(2*pi*3

%% Create the envelope
% Set the additional parameters in the model.
A0 = 2; % Initial amplitude.
B  = 1.5; % Amplitude decay rate.
fm = 0.65; % Frequency of the modulating envelope.
% Create the envelope
A  = A0*exp(-B*t).*sin(2*pi*fm*t);

%% Create the call.
call = A.*y0;
```

# Publishing Code

```matlab
%% Set the fundamental frequency of the call.
f0 = 175;

%% Create the harmonics.
y0 = sin(2*pi*f0*t) + sin(2*pi*2*f0*t) + sin(2*pi*3*

%% Create the envelope
% Set the additional parameters in the model.
A0 = 2; % Initial amplitude.
B  = 1.5; % Amplitude decay rate.
fm = 0.65; % Frequency of the modulating envelope.
% Create the envelope
A  = A0*exp(-B*t).*sin(2*pi*fm*t);

%% Create the call.
call = A.*y0;
```
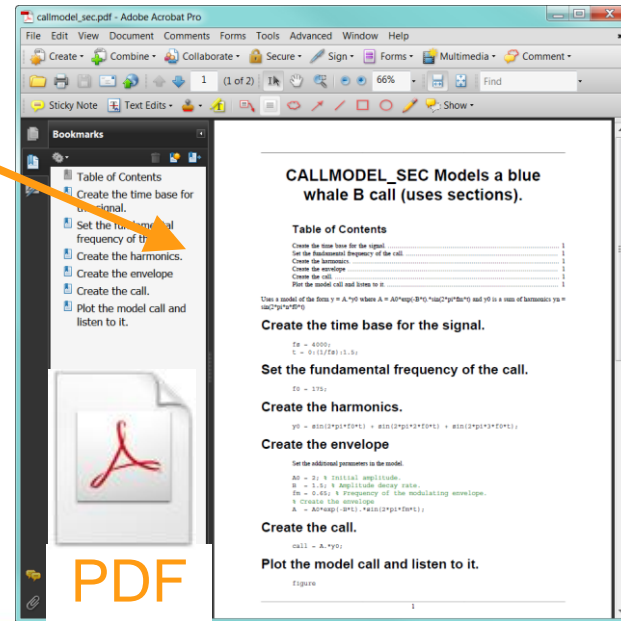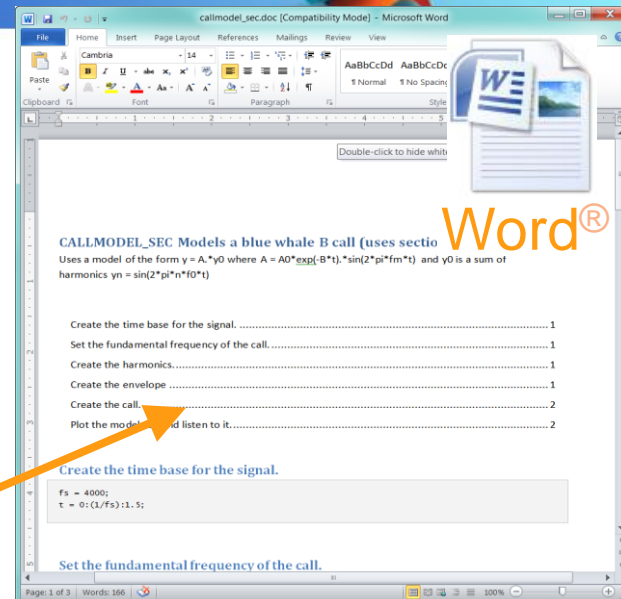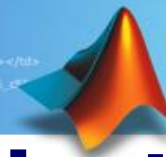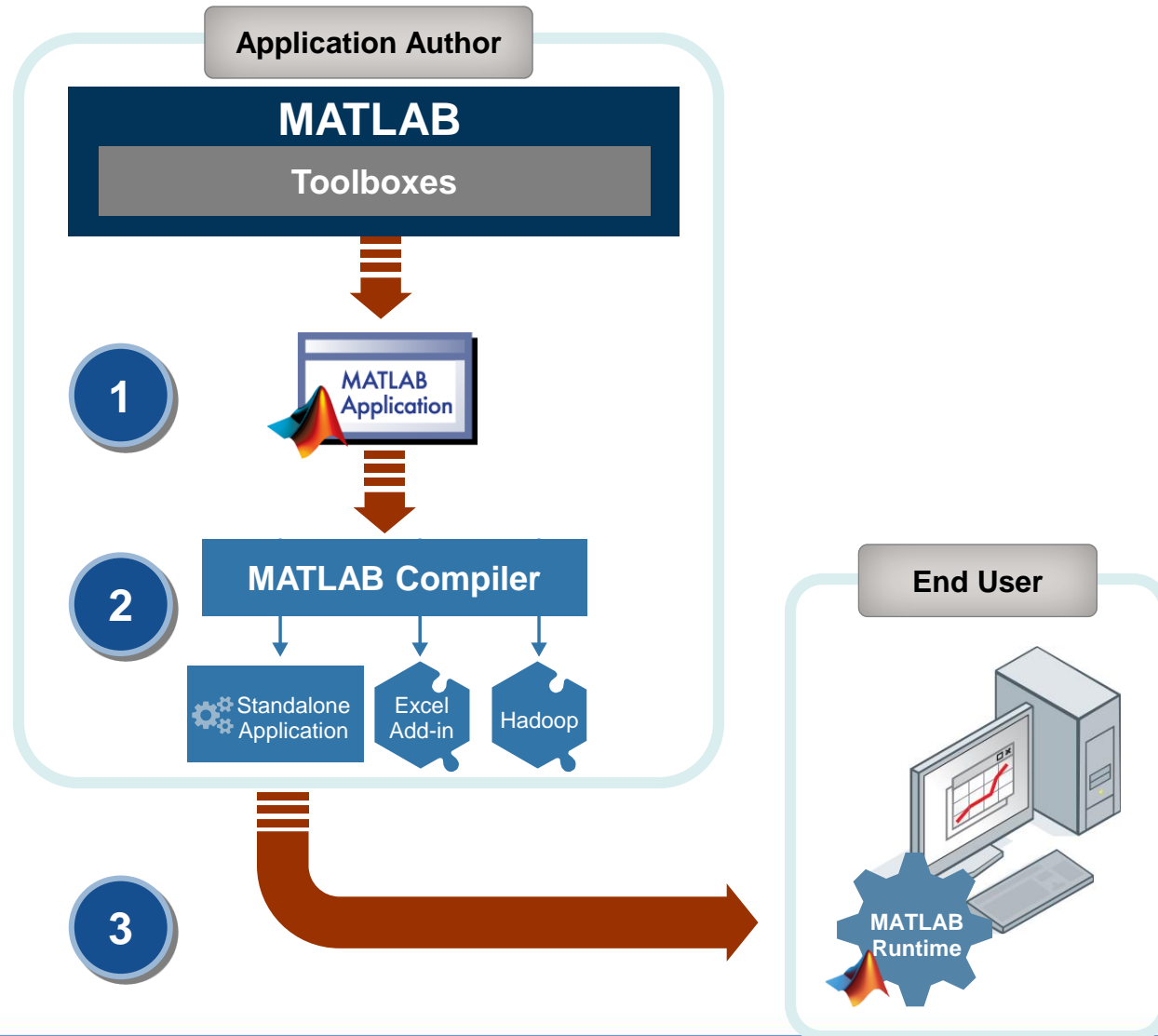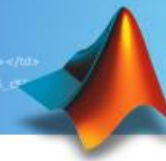
Word®

PDF

# Sharing Standalone Applications

# Thanks for your attention !

MATLAB&SIMULINK