# Data Wrangling
## with pandas
## Cheat Sheet
### http://pandas.pydata.org

## Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

Each **variable** is saved in its own **column**

&

Each **observation** is saved in its own **row**

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

## Syntax – Creating DataFrames

|   | a | b | c |
|---|---|---|---|
| 1 | 4 | 7 | 10 |
| 2 | 5 | 8 | 11 |
| 3 | 6 | 9 | 12 |

```
df = pd.DataFrame(
        {"a" : [4 ,5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
      index = [1, 2, 3])
```
Specify values for each column.

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
```
Specify values for each row.

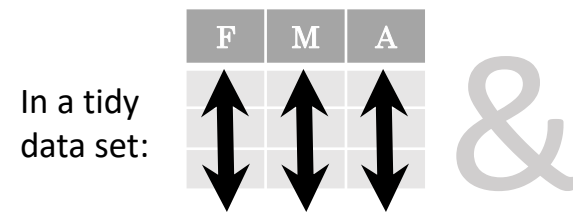|   |   | a | b | c |
|---|---|---|---|---|
| n | v |   |   |   |
| d | 1 | 4 | 7 | 10 |
|   | 2 | 5 | 8 | 11 |
| e | 2 | 6 | 9 | 12 |

```
df = pd.DataFrame(
        {"a" : [4 ,5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
            names=['n','v']))
```
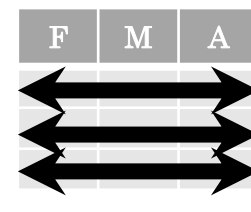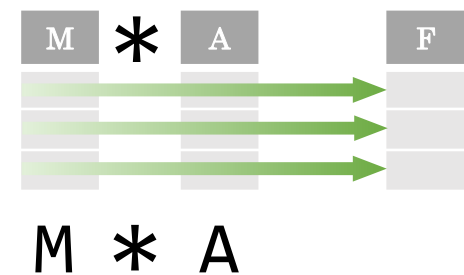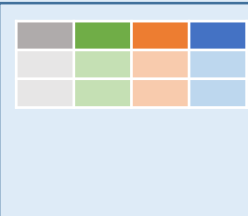Create DataFrame with a MultiIndex

## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.
```
df = (pd.melt(df)
        .rename(columns={
                'variable' : 'var',
                'value' : 'val'})
        .query('val >= 200')
)
```
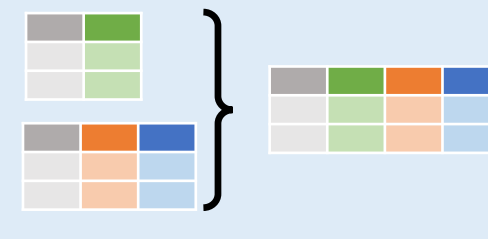
## Reshaping Data – Change the layout of a data set

**pd.melt(df)**
Gather columns into rows.

**df.pivot(columns='var', values='val')**
Spread rows into columns.

**pd.concat([df1,df2])**
Append rows of DataFrames

**pd.concat([df1,df2], axis=1)**
Append columns of DataFrames

**df.sort_values('mpg')**
Order rows by values of a column (low to high).

**df.sort_values('mpg',ascending=False)**
Order rows by values of a column (high to low).

**df.rename(columns = {'y':'year'})**
Rename the columns of a DataFrame

**df.sort_index()**
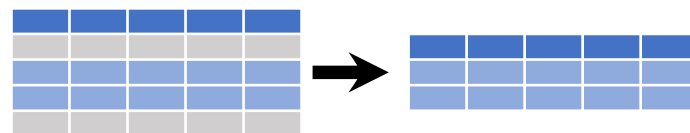Sort the index of a DataFrame

**df.reset_index()**
Reset index of DataFrame to row numbers, moving index to columns.

**df.drop(columns=['Length','Height'])**
Drop columns from DataFrame

## Subset Observations (Rows)

**df[df.Length > 7]**
Extract rows that meet logical criteria.

**df.drop_duplicates()**
Remove duplicate rows (only considers columns).

**df.head(n)**
Select first n rows.

**df.tail(n)**
Select last n rows.

**df.sample(frac=0.5)**
Randomly select fraction of rows.

**df.sample(n=10)**
Randomly select n rows.

**df.iloc[10:20]**
Select rows by position.

**df.nlargest(n, 'value')**
Select and order top n entries.

**df.nsmallest(n, 'value')**
Select and order bottom n entries.

## Subset Variables (Columns)

**df[['width','length','species']]**
Select multiple columns with specific names.

**df['width']** *or* **df.width**
Select single column with specific name.

**df.filter(regex='regex')**
Select columns whose name matches regular expression *regex*.

| regex (Regular Expressions) Examples | |
|---|---|
| `'\.'` | Matches strings containing a period '.' |
| `'Length$'` | Matches strings ending with word 'Length' |
| `'^Sepal'` | Matches strings beginning with the word 'Sepal' |
| `'^x[1-5]$'` | Matches strings beginning with 'x' and ending with 1,2,3,4,5 |
| `'^(?!Species$).*'` | Matches strings except the string 'Species' |

**df.loc[:,'x2':'x4']**
Select all columns between x2 and x4 (inclusive).

**df.iloc[:,[1,2,5]]**
Select columns in positions 1, 2 and 5 (first column is 0).

**df.loc[df['a'] > 10, ['a','c']]**
Select rows meeting logical condition, and only the specific columns .

| Logic in Python (and pandas) | | | | |
|---|---|---|---|---|
| < | Less than | != | | Not equal to |
| > | Greater than | `df.column.isin(values)` | | Group membership |
| == | Equals | `pd.isnull(obj)` | | Is NaN |
| <= | Less than or equals | `pd.notnull(obj)` | | Is not NaN |
| >= | Greater than or equals | `&,\|,~,^,df.any(),df.all()` | | Logical and, or, not, xor, any, all |

# Summarize Data

`df['w'].value_counts()`
    Count number of rows with each unique value of variable
`len(df)`
    # of rows in DataFrame.
`df['w'].nunique()`
    # of distinct values in a column.
`df.describe()`
    Basic descriptive statistics for each column (or GroupBy)



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

`sum()`
    Sum values of each object.
`count()`
    Count non-NA/null values of each object.
`median()`
    Median value of each object.
`quantile([0.25,0.75])`
    Quantiles of each object.
`apply(function)`
    Apply function to each object.

`min()`
    Minimum value in each object.
`max()`
    Maximum value in each object.
`mean()`
    Mean value of each object.
`var()`
    Variance of each object.
`std()`
    Standard deviation of each object.

# Group Data



`df.groupby(by="col")`
    Return a GroupBy object, grouped by values in column named "col".

`df.groupby(level="ind")`
    Return a GroupBy object, grouped by values in index level named "ind".

All of the summary functions listed above can be applied to a group.
Additional GroupBy functions:
`size()`
    Size of each group.
`agg(function)`
    Aggregate group using function.

# Windows

`df.expanding()`
    Return an Expanding object allowing summary functions to be applied cumulatively.
`df.rolling(n)`
    Return a Rolling object allowing summary functions to be applied to windows of length n.

# Handling Missing Data

`df.dropna()`
    Drop rows with any column having NA/null data.
`df.fillna(value)`
    Replace all NA/null data with value.

# Make New Columns



`df.assign(Area=lambda df: df.Length*df.Height)`
    Compute and append one or more new columns.
`df['Volume'] = df.Length*df.Height*df.Depth`
    Add single column.
`pd.qcut(df.col, n, labels=False)`
    Bin column into n buckets.



pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

`max(axis=1)`
    Element-wise max.
`clip(lower=-10,upper=10)`
    Trim values at input thresholds
`min(axis=1)`
    Element-wise min.
`abs()`
    Absolute value.

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

`shift(1)`
    Copy with values shifted by 1.
`rank(method='dense')`
    Ranks with no gaps.
`rank(method='min')`
    Ranks. Ties get min rank.
`rank(pct=True)`
    Ranks rescaled to interval [0, 1].
`rank(method='first')`
    Ranks. Ties go to first value.

`shift(-1)`
    Copy with values lagged by 1.
`cumsum()`
    Cumulative sum.
`cummax()`
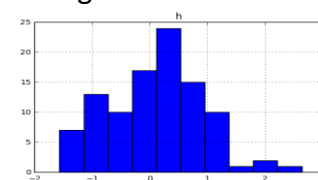    Cumulative max.
`cummin()`
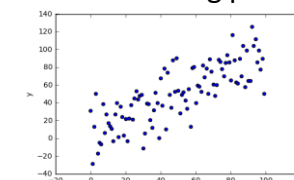    Cumulative min.
`cumprod()`
    Cumulative product.

# Plotting

`df.plot.hist()`
Histogram for each column

`df.plot.scatter(x='w',y='h')`
Scatter chart using pairs of points



# Combine Data Sets



### Standard Joins



`pd.merge(adf, bdf,`
`        how='left', on='x1')`
    Join matching rows from bdf to adf.

`pd.merge(adf, bdf,`
`        how='right', on='x1')`
    Join matching rows from adf to bdf.

`pd.merge(adf, bdf,`
`        how='inner', on='x1')`
    Join data. Retain only rows in both sets.

`pd.merge(adf, bdf,`
`        how='outer', on='x1')`
    Join data. Retain all values, all rows.

### Filtering Joins

`adf[adf.x1.isin(bdf.x1)]`
    All rows in adf that have a match in bdf.

`adf[~adf.x1.isin(bdf.x1)]`
    All rows in adf that do not have a match in bdf.



### Set-like Operations

`pd.merge(ydf, zdf)`
    Rows that appear in both ydf and zdf (Intersection).

`pd.merge(ydf, zdf, how='outer')`
    Rows that appear in either or both ydf and zdf (Union).

`pd.merge(ydf, zdf, how='outer',`
`        indicator=True)`
`.query('_merge == "left_only"')`
`.drop(columns=['_merge'])`
    Rows that appear in ydf but not zdf (Setdiff).

# Python For Data Science *Cheat Sheet*
## Python Basics

Learn More Python for Data Science *Interactively* at www.datacamp.com

## Variables and Data Types

### Variable Assignment

```
>>> x=5
>>> x
 5
```

### Calculations With Variables

| | | |
|---|---|---|
| `>>> x+2`<br>`7` | | Sum of two variables |
| `>>> x-2`<br>`3` | | Subtraction of two variables |
| `>>> x*2`<br>`10` | | Multiplication of two variables |
| `>>> x**2`<br>`25` | | Exponentiation of a variable |
| `>>> x%2`<br>`1` | | Remainder of a variable |
| `>>> x/float(2)`<br>`2.5` | | Division of a variable |

### Types and Type Conversion

| | | |
|---|---|---|
| `str()` | `'5', '3.45', 'True'` | Variables to strings |
| `int()` | `5, 3, 1` | Variables to integers |
| `float()` | `5.0, 1.0` | Variables to floats |
| `bool()` | `True, True, True` | Variables to booleans |

## Asking For Help

```
>>> help(str)
```

## Strings

```
>>> my_string = 'thisStringIsAwesome'
>>> my_string
'thisStringIsAwesome'
```

### String Operations

```
>>> my_string * 2
 'thisStringIsAwesomethisStringIsAwesome'
>>> my_string + 'Innit'
 'thisStringIsAwesomeInnit'
>>> 'm' in my_string
 True
```

## Lists

**Also see NumPy Arrays**

```
>>> a = 'is'
>>> b = 'nice'
>>> my_list = ['my', 'list', a, b]
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

### Selecting List Elements

**Index starts at 0**

| Subset | |
|---|---|
| `>>> my_list[1]` | Select item at index 1 |
| `>>> my_list[-3]` | Select 3rd last item |
| **Slice** | |
| `>>> my_list[1:3]` | Select items at index 1 and 2 |
| `>>> my_list[1:]` | Select items after index 0 |
| `>>> my_list[:3]` | Select items before index 3 |
| `>>> my_list[:]` | Copy my_list |
| **Subset Lists of Lists** | |
| `>>> my_list2[1][0]` | `my_list[list][itemOfList]` |
| `>>> my_list2[1][:2]` | |

### List Operations

```
>>> my_list + my_list
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list * 2
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list2 > 4
True
```

### List Methods

| | |
|---|---|
| `>>> my_list.index(a)` | Get the index of an item |
| `>>> my_list.count(a)` | Count an item |
| `>>> my_list.append('!')` | Append an item at a time |
| `>>> my_list.remove('!')` | Remove an item |
| `>>> del(my_list[0:1])` | Remove an item |
| `>>> my_list.reverse()` | Reverse the list |
| `>>> my_list.extend('!')` | Append an item |
| `>>> my_list.pop(-1)` | Remove an item |
| `>>> my_list.insert(0,'!')` | Insert an item |
| `>>> my_list.sort()` | Sort the list |

### String Operations

**Index starts at 0**

```
>>> my_string[3]
>>> my_string[4:9]
```

### String Methods

| | |
|---|---|
| `>>> my_string.upper()` | String to uppercase |
| `>>> my_string.lower()` | String to lowercase |
| `>>> my_string.count('w')` | Count String elements |
| `>>> my_string.replace('e', 'i')` | Replace String elements |
| `>>> my_string.strip()` | Strip whitespaces |

## Libraries

### Import libraries

```
>>> import numpy
>>> import numpy as np
```
**Selective import**
```
>>> from math import pi
```

pandas — Data analysis
Machine learning
NumPy — Scientific computing
matplotlib — 2D plotting

## Install Python

ANACONDA
Leading open data science platform powered by Python

spyder
Free IDE that is included with Anaconda

jupyter
Create and share documents with live code, visualizations, text, …

## Numpy Arrays

**Also see Lists**

```
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3],[4,5,6]])
```

### Selecting Numpy Array Elements

**Index starts at 0**

| Subset | |
|---|---|
| `>>> my_array[1]`<br>`  2` | Select item at index 1 |
| **Slice** | |
| `>>> my_array[0:2]`<br>`  array([1, 2])` | Select items at index 0 and 1 |
| **Subset 2D Numpy arrays** | |
| `>>> my_2darray[:,0]`<br>`  array([1, 4])` | `my_2darray[rows, columns]` |

### Numpy Array Operations

```
>>> my_array > 3
 array([False, False, False,  True], dtype=bool)
>>> my_array * 2
 array([2, 4, 6, 8])
>>> my_array + np.array([5, 6, 7, 8])
 array([6, 8, 10, 12])
```

### Numpy Array Functions

| | |
|---|---|
| `>>> my_array.shape` | Get the dimensions of the array |
| `>>> np.append(other_array)` | Append items to an array |
| `>>> np.insert(my_array, 1, 5)` | Insert items in an array |
| `>>> np.delete(my_array,[1])` | Delete items in an array |
| `>>> np.mean(my_array)` | Mean of the array |
| `>>> np.median(my_array)` | Median of the array |
| `>>> my_array.corrcoef()` | Correlation coefficient |
| `>>> np.std(my_array)` | Standard deviation |

# Visual Studio Code

Keyboard shortcuts for Windows

## General

| | |
|---|---|
| Ctrl+Shift+P, F1 | Show Command Palette |
| Ctrl+P | Quick Open, Go to File... |
| Ctrl+Shift+N | New window/instance |
| Ctrl+Shift+W | Close window/instance |
| Ctrl+, | User Settings |
| Ctrl+K Ctrl+S | Keyboard Shortcuts |

## Basic editing

| | |
|---|---|
| Ctrl+X | Cut line (empty selection) |
| Ctrl+C | Copy line (empty selection) |
| Alt+ ↑ / ↓ | Move line up/down |
| Shift+Alt + ↓ / ↑ | Copy line up/down |
| Ctrl+Shift+K | Delete line |
| Ctrl+Enter | Insert line below |
| Ctrl+Shift+Enter | Insert line above |
| Ctrl+Shift+\ | Jump to matching bracket |
| Ctrl+] / [ | Indent/outdent line |
| Home / End | Go to beginning/end of line |
| Ctrl+Home | Go to beginning of file |
| Ctrl+End | Go to end of file |
| Ctrl+↑ / ↓ | Scroll line up/down |
| Alt+PgUp / PgDn | Scroll page up/down |
| Ctrl+Shift+[ | Fold (collapse) region |
| Ctrl+Shift+] | Unfold (uncollapse) region |
| Ctrl+K Ctrl+[ | Fold (collapse) all subregions |
| Ctrl+K Ctrl+] | Unfold (uncollapse) all subregions |
| Ctrl+K Ctrl+0 | Fold (collapse) all regions |
| Ctrl+K Ctrl+J | Unfold (uncollapse) all regions |
| Ctrl+K Ctrl+C | Add line comment |
| Ctrl+K Ctrl+U | Remove line comment |
| Ctrl+/ | Toggle line comment |
| Shift+Alt+A | Toggle block comment |
| Alt+Z | Toggle word wrap |

## Navigation

| | |
|---|---|
| Ctrl+T | Show all Symbols |
| Ctrl+G | Go to Line... |
| Ctrl+P | Go to File... |
| Ctrl+Shift+O | Go to Symbol... |
| Ctrl+Shift+M | Show Problems panel |
| F8 | Go to next error or warning |
| Shift+F8 | Go to previous error or warning |
| Ctrl+Shift+Tab | Navigate editor group history |
| Alt+ ← / → | Go back / forward |

| | |
|---|---|
| Ctrl+M | Toggle Tab moves focus |

## Search and replace

| | |
|---|---|
| Ctrl+F | Find |
| Ctrl+H | Replace |
| F3 / Shift+F3 | Find next/previous |
| Alt+Enter | Select all occurences of Find match |
| Ctrl+D | Add selection to next Find match |
| Ctrl+K Ctrl+D | Move last selection to next Find match |
| Alt+C / R / W | Toggle case-sensitive / regex / whole word |

## Multi-cursor and selection

| | |
|---|---|
| Alt+Click | Insert cursor |
| Ctrl+Alt+ ↑ / ↓ | Insert cursor above / below |
| Ctrl+U | Undo last cursor operation |
| Shift+Alt+I | Insert cursor at end of each line selected |
| Ctrl+L | Select current line |
| Ctrl+Shift+L | Select all occurrences of current selection |
| Ctrl+F2 | Select all occurrences of current word |
| Shift+Alt+→ | Expand selection |
| Shift+Alt+← | Shrink selection |
| Shift+Alt + (drag mouse) | Column (box) selection |
| Ctrl+Shift+Alt + (arrow key) | Column (box) selection |
| Ctrl+Shift+Alt +PgUp/PgDn | Column (box) selection page up/down |

## Rich languages editing

| | |
|---|---|
| Ctrl+Space | Trigger suggestion |
| Ctrl+Shift+Space | Trigger parameter hints |
| Shift+Alt+F | Format document |
| Ctrl+K Ctrl+F | Format selection |
| F12 | Go to Definition |
| Alt+F12 | Peek Definition |
| Ctrl+K F12 | Open Definition to the side |
| Ctrl+. | Quick Fix |
| Shift+F12 | Show References |
| F2 | Rename Symbol |
| Ctrl+K Ctrl+X | Trim trailing whitespace |
| Ctrl+K M | Change file language |

## Editor management

| | |
|---|---|
| Ctrl+F4, Ctrl+W | Close editor |
| Ctrl+K F | Close folder |
| Ctrl+\ | Split editor |
| Ctrl+ 1 / 2 / 3 | Focus into 1st, 2nd or 3rd editor group |
| Ctrl+K Ctrl+ ←/→ | Focus into previous/next editor group |
| Ctrl+Shift+PgUp / PgDn | Move editor left/right |
| Ctrl+K ← / → | Move active editor group |

## File management

| | |
|---|---|
| Ctrl+N | New File |
| Ctrl+O | Open File... |
| Ctrl+S | Save |
| Ctrl+Shift+S | Save As... |
| Ctrl+K S | Save All |
| Ctrl+F4 | Close |
| Ctrl+K Ctrl+W | Close All |
| Ctrl+Shift+T | Reopen closed editor |
| Ctrl+K Enter | Keep preview mode editor open |
| Ctrl+Tab | Open next |
| Ctrl+Shift+Tab | Open previous |
| Ctrl+K P | Copy path of active file |
| Ctrl+K R | Reveal active file in Explorer |
| Ctrl+K O | Show active file in new window/instance |

## Display

| | |
|---|---|
| F11 | Toggle full screen |
| Shift+Alt+0 | Toggle editor layout (horizontal/vertical) |
| Ctrl+ = / - | Zoom in/out |
| Ctrl+B | Toggle Sidebar visibility |
| Ctrl+Shift+E | Show Explorer / Toggle focus |
| Ctrl+Shift+F | Show Search |
| Ctrl+Shift+G | Show Source Control |
| Ctrl+Shift+D | Show Debug |
| Ctrl+Shift+X | Show Extensions |
| Ctrl+Shift+H | Replace in files |
| Ctrl+Shift+J | Toggle Search details |
| Ctrl+Shift+U | Show Output panel |
| Ctrl+Shift+V | Open Markdown preview |
| Ctrl+K V | Open Markdown preview to the side |
| Ctrl+K Z | Zen Mode (Esc Esc to exit) |

## Debug

| | |
|---|---|
| F9 | Toggle breakpoint |
| F5 | Start/Continue |
| Shift+F5 | Stop |
| F11 / Shift+F11 | Step into/out |
| F10 | Step over |
| Ctrl+K Ctrl+I | Show hover |

## Integrated terminal

| | |
|---|---|
| Ctrl+` | Show integrated terminal |
| Ctrl+Shift+` | Create new terminal |
| Ctrl+C | Copy selection |
| Ctrl+V | Paste into active terminal |
| Ctrl+↑ / ↓ | Scroll up/down |
| Shift+PgUp / PgDn | Scroll page up/down |
| Ctrl+Home / End | Scroll to top/bottom |

Other operating systems' keyboard shortcuts and additional unassigned shortcuts available at [aka.ms/vscodekeybindings](aka.ms/vscodekeybindings)

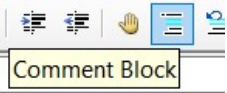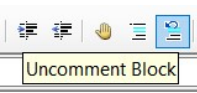# VBA Essentials Cheat Sheet

## VBA Data Types

**String**
Used to hold text

**Long**
Long integer (whole numbers). -2,147,483,648 to 2,147,483,647

**Integer**
Short integer (whole number). -32,768 to 32,767

**Boolean**
True or False

**Boolean**
True or False

**Date**
Holds date data types. 1/1/100 to 12/31/9999

**Single/Double**
Used to hold values with decimals

**Variant**
Catch all data type. When an explicit data type is not declared, variant type is assigned

## VBA Common Operations (Required syntax in bold)

**If Statement**
```
If numGrade > 90 Then
     letterGrade = "A"
ElseIf numGrade > 80 Then
     letterGrade = "B"
Else
     letterGrade = "F"
End If
```

**For … Next Loop**
```
For x=0 to 49
     'Loop Over Code
Next x
```

**For Each … Next Loop**
```
For Each Item In Selection
     Item.Offset(0, 1) = Item * 2
Next
```

**Do … Loop While**
```
Do
     .Range("A1").Offset(Item,0) = Item
Loop While myBool = True
```

**Do While … Loop**
```
Do While myBool=True
     .Range("A1").Offset(Item,0) = Item
Loop
```

## VBA Interacting With User

**Message Box**
```
Msgbox "Hello world"
```

**User Input**
```
usrInput = InputBox("Please Enter Your Name")
```

## Comparison Operators

**Greater Than / Greater Than or Equal**
```
Greater Than : >
Greater Than or Equal: >=
```

**Less Than / Less Than or Equal**
```
Less Than : <
Less Than or Equal: <=
```

**Equal / Not Equal**
```
Equal : =
Not Equal: <>
```

## Logical Operators

**Or**
```
True Or True = True
True Or False = True
False Or False = False
```

**And**
```
True And True = True
True And False = False
False And False = False
```

**Not**
```
Not True = False
Not False = True
```

## Commenting Code

**Single Line Comment**
```
Single line comments are created by
using an apostorpher (') at the
beginning of a line

Msgbox "This line of code will execute"
'Msgbox "This line of code will execute"
```

**Multi Line Comments**
```
View -> Toolbars -> Edit
```


Comment Block    Uncomment Block

## Referencing Workbooks/Worksheets/Ranges

**Workbooks**

**Workbook that contains code:**
ThisWorkbook

**Using the Active Workbook:**
Active Workbook

**Using Numbered Index:**
Workbooks(1)

**Using Workbook Name:**
Workbooks("myWkbk")

**Worksheet**

**Using the Active Worksheet:**
ActiveSheet

**Using the Selected Worksheet:**
Windows.SelectedSheets

**Using Numbered Index:**
Worksheets(1)

**Using Worksheet Name:**
Worksheets("myWksht")

**Range**

**Reference Single Cell:**
Range("A1")

**Refernce Multiple Adjacent Cells:**
Range("A1:C5")

**Reference Multiple Non Adjacent Cells**
Range("A1:A5, C1:C5")

**Using a Named Range**
Range("myRange")

**Cells**

**Refernce All Cells**
Worksheet.Cells

**Rererence Cells with one Parameter**
Cells(3) = "C1"

**Reference Cells With Two Parameters**
Cells(3,3) = "C3"
Cells(3, "E") = "E3"

## Useful Tips

**With … End With**
```
With ThisWorkbook.Worksheets(1)
     .Range("A1")=Month
End With
```

**OffSet**
```
For x=0 to 100
        .Range("A1").Offset(x,0) = Rnd
Next x
```