# 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 |

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

```python
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 |

```python
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'])))
```
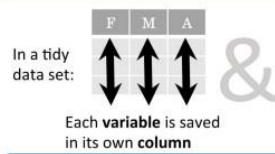Create DataFrame with a MultiIndex

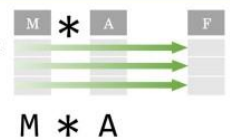## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```python
df = (pd.melt(df)
        .rename(columns={
            'variable' : 'var',
            'value' : 'val'})
        .query('val >= 200')
    )
```
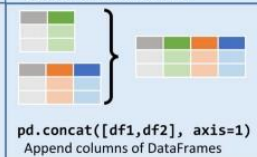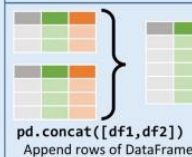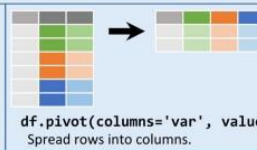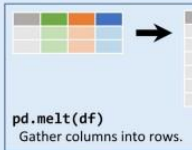
## Reshaping Data – Change the layout of a data set

**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)**
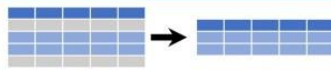Order rows by values of a column (high to low).

**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(['Length','Height'], axis=1)**
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)**
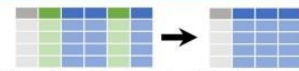Randomly select fraction of rows.

**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.

### 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 |

## 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']]**
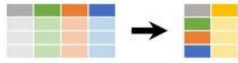Select rows meeting logical condition, and only the specific columns .

# 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()` | `min()` |
| Sum values of each object. | Minimum value in each object. |
| `count()` | `max()` |
| Count non-NA/null values of each object. | Maximum value in each object. |
| `median()` | `mean()` |
| Median value of each object. | Mean value of each object. |
| `quantile([0.25,0.75])` | `var()` |
| Quantiles of each object. | Variance of each object. |
| `apply(function)` | `std()` |
| Apply function to each object. | 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)`
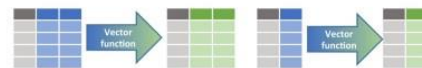Replace all NA/null data with value.

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

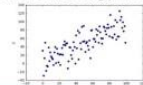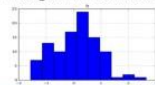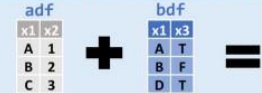| | |
|---|---|
| `shift(1)` | `shift(-1)` |
| Copy with values shifted by 1. | Copy with values lagged by 1. |
| `rank(method='dense')` | `cumsum()` |
| Ranks with no gaps. | Cumulative sum. |
| `rank(method='min')` | `cummax()` |
| Ranks. Ties get min rank. | Cumulative max. |
| `rank(pct=True)` | `cummin()` |
| Ranks rescaled to interval [0, 1]. | Cumulative min. |
| `rank(method='first')` | `cumprod()` |
| Ranks. Ties go to first value. | 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(['_merge'],axis=1)`
Rows that appear in ydf but not zdf (Setdiff).

# Data Wrangling
## with dplyr and tidyr
### Cheat Sheet
**R**Studio

---

## Tidy Data - A foundation for wrangling in R

In a tidy data set:

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

&

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

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

---

## Syntax - Helpful conventions for wrangling

**dplyr::tbl_df(iris)**

Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen:

```
Source: local data frame [150 x 5]

  Sepal.Length Sepal.Width Petal.Length
1          5.1         3.5          1.4
2          4.9         3.0          1.4
3          4.7         3.2          1.3
4          4.6         3.1          1.5
5          5.0         3.6          1.4
..         ...         ...          ...
Variables not shown: Petal.Width (dbl),
  Species (fctr)
```

**dplyr::glimpse(iris)**

Information dense summary of tbl data.

**utils::View(iris)**

View data set in spreadsheet-like display (note capital V).

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|---|
| 1 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 2 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 3 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 5 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |
| 6 | 5.4 | 3.9 | 1.7 | 0.4 | setosa |
| 7 | 4.6 | 3.4 | 1.4 | 0.3 | setosa |
| 8 | 5.0 | 3.4 | 1.5 | 0.2 | setosa |

**dplyr::%>%**

Passes object on left hand side as first argument (or . argument) of function on righthand side.

```
x %>% f(y)      is the same as  f(x, y)
y %>% f(x, ., z) is the same as  f(x, y, z )
```
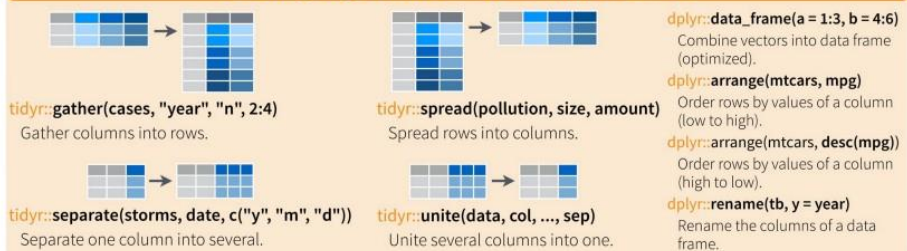
"Piping" with %>% makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```

---

## Reshaping Data - Change the layout of a data set

**tidyr::gather(cases, "year", "n", 2:4)**
Gather columns into rows.

**tidyr::spread(pollution, size, amount)**
Spread rows into columns.

**tidyr::separate(storms, date, c("y", "m", "d"))**
Separate one column into several.

**tidyr::unite(data, col, ..., sep)**
Unite several columns into one.

**dplyr::data_frame(a = 1:3, b = 4:6)**
Combine vectors into data frame (optimized).

**dplyr::arrange(mtcars, mpg)**
Order rows by values of a column (low to high).

**dplyr::arrange(mtcars, desc(mpg))**
Order rows by values of a column (high to low).

**dplyr::rename(tb, y = year)**
Rename the columns of a data frame.

---

## Subset Observations (Rows)

**dplyr::filter(iris, Sepal.Length > 7)**
Extract rows that meet logical criteria.

**dplyr::distinct(iris)**
Remove duplicate rows.

**dplyr::sample_frac(iris, 0.5, replace = TRUE)**
Randomly select fraction of rows.

**dplyr::sample_n(iris, 10, replace = TRUE)**
Randomly select n rows.

**dplyr::slice(iris, 10:15)**
Select rows by position.

**dplyr::top_n(storms, 2, date)**
Select and order top n entries (by group if grouped data).

### Logic in R - ?Comparison, ?base::Logic

| | | | | |
|---|---|---|---|---|
| < | Less than | != | Not equal to | |
| > | Greater than | %in% | Group membership | |
| == | Equal to | is.na | Is NA | |
| <= | Less than or equal to | !is.na | Is not NA | |
| >= | Greater than or equal to | &,\|,!,xor,any,all | Boolean operators | |

---

## Subset Variables (Columns)

**dplyr::select(iris, Sepal.Width, Petal.Length, Species)**
Select columns by name or helper function.

### Helper functions for select - ?select

**select(iris, contains("."))**
Select columns whose name contains a character string.

**select(iris, ends_with("Length"))**
Select columns whose name ends with a character string.

**select(iris, everything())**
Select every column.

**select(iris, matches(".t."))**
Select columns whose name matches a regular expression.

**select(iris, num_range("x", 1:5))**
Select columns named x1, x2, x3, x4, x5.

**select(iris, one_of(c("Species", "Genus")))**
Select columns whose names are in a group of names.

**select(iris, starts_with("Sepal"))**
Select columns whose name starts with a character string.

**select(iris, Sepal.Length:Petal.Width)**
Select all columns between Sepal.Length and Petal.Width (inclusive).

**select(iris, -Species)**
Select all columns except Species.

---

devtools::install_github("rstudio/EDAWR") for data sets     Learn more with **browseVignettes(package = c("dplyr", "tidyr"))** • dplyr 0.4.0• tidyr 0.2.0 • Updated: 1/15

## Summarise Data



**dplyr::summarise(iris, avg = mean(Sepal.Length))**
Summarise data into single row of values.
**dplyr::summarise_each(iris, funs(mean))**
Apply summary function to each column.
**dplyr::count(iris, Species, wt = Sepal.Length)**
Count number of rows with each unique value of variable (with or without weights).



Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

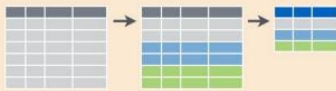| | |
|---|---|
| **dplyr::first** | **min** |
| First value of a vector. | Minimum value in a vector. |
| **dplyr::last** | **max** |
| Last value of a vector. | Maximum value in a vector. |
| **dplyr::nth** | **mean** |
| Nth value of a vector. | Mean value of a vector. |
| **dplyr::n** | **median** |
| # of values in a vector. | Median value of a vector. |
| **dplyr::n_distinct** | **var** |
| # of distinct values in a vector. | Variance of a vector. |
| **IQR** | **sd** |
| IQR of a vector. | Standard deviation of a vector. |

## Group Data

**dplyr::group_by(iris, Species)**
Group data into rows with the same value of Species.
**dplyr::ungroup(iris)**
Remove grouping information from data frame.

**iris %>% group_by(Species) %>% summarise(...)**
Compute separate summary row for each group.



## Make New Variables



**dplyr::mutate(iris, sepal = Sepal.Length + Sepal. Width)**
Compute and append one or more new columns.
**dplyr::mutate_each(iris, funs(min_rank))**
Apply window function to each column.
**dplyr::transmute(iris, sepal = Sepal.Length + Sepal. Width)**
Compute one or more new columns. Drop original columns.



Mutate uses **window functions**, functions that take a vector of values and return another vector of values, such as:

| | |
|---|---|
| **dplyr::lead** | **dplyr::cumall** |
| Copy with values shifted by 1. | Cumulative `all` |
| **dplyr::lag** | **dplyr::cumany** |
| Copy with values lagged by 1. | Cumulative `any` |
| **dplyr::dense_rank** | **dplyr::cummean** |
| Ranks with no gaps. | Cumulative `mean` |
| **dplyr::min_rank** | **cumsum** |
| Ranks. Ties get min rank. | Cumulative `sum` |
| **dplyr::percent_rank** | **cummax** |
| Ranks rescaled to [0, 1]. | Cumulative `max` |
| **dplyr::row_number** | **cummin** |
| Ranks. Ties got to first value. | Cumulative `min` |
| **dplyr::ntile** | **cumprod** |
| Bin vector into n buckets. | Cumulative `prod` |
| **dplyr::between** | **pmax** |
| Are values between a and b? | Element-wise `max` |
| **dplyr::cume_dist** | **pmin** |
| Cumulative distribution. | Element-wise `min` |

**iris %>% group_by(Species) %>% mutate(...)**
Compute new variables by group.

## Combine Data Sets



**Mutating Joins**

 **dplyr::left_join(a, b, by = "x1")**
Join matching rows from b to a.

 **dplyr::right_join(a, b, by = "x1")**
Join matching rows from a to b.

 **dplyr::inner_join(a, b, by = "x1")**
Join data. Retain only rows in both sets.

 **dplyr::full_join(a, b, by = "x1")**
Join data. Retain all values, all rows.

**Filtering Joins**

 **dplyr::semi_join(a, b, by = "x1")**
All rows in a that have a match in b.

 **dplyr::anti_join(a, b, by = "x1")**
All rows in a that do not have a match in b.



**Set Operations**

 **dplyr::intersect(y, z)**
Rows that appear in both y and z.

 **dplyr::union(y, z)**
Rows that appear in either or both y and z.

 **dplyr::setdiff(y, z)**
Rows that appear in y but not z.

**Binding**

 **dplyr::bind_rows(y, z)**
Append z to y as new rows.

 **dplyr::bind_cols(y, z)**
Append z to y as new columns.
Caution: matches rows by position.