

## 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**

First value of a vector.

**dplyr::last**

Last value of a vector.

**dplyr::nth**

Nth value of a vector.

**dplyr::n**

# of values in a vector.

**dplyr::n\_distinct**

# of distinct values in a vector.

**IQR**

IQR of a vector.

**min**

Minimum value in a vector.

**max**

Maximum value in a vector.

**mean**

Mean value of a vector.

**median**

Median value of a vector.

**var**

Variance of a vector.

**sd**

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**

Copy with values shifted by 1.

**dplyr::lag**

Copy with values lagged by 1.

**dplyr::dense\_rank**

Ranks with no gaps.

**dplyr::min\_rank**

Ranks. Ties get min rank.

**dplyr::percent\_rank**

Ranks rescaled to [0, 1].

**dplyr::row\_number**

Ranks. Ties got to first value.

**dplyr::ntile**

Bin vector into n buckets.

**dplyr::between**

Are values between a and b?

**dplyr::cume\_dist**

Cumulative distribution.

**dplyr::cumall**

Cumulative all

**dplyr::cumany**

Cumulative any

**dplyr::cummean**

Cumulative mean

**cumsum**

Cumulative sum

**cummax**

Cumulative max

**cummin**

Cumulative min

**cumprod**

Cumulative prod

**pmax**

Element-wise max

**pmin**

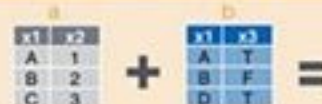
Element-wise min

**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.

**dplyr::anti\_join(a, b, by = "x1")**

All rows in a that do not have a match in b.

**dplyr::semi\_join(a, b, by = "x1")**

All rows in a that have a match in b.

**dplyr::setdiff(y, z)**

Rows that appear in y but not 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.

**dplyr::intersect(y, z)**

Rows that appear in both y and z.

**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.

**dplyr::setdiff(y, z)**

Rows that appear in y but not 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.

**dplyr::intersect(y, z)**

Rows that appear in both y and z.

**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.

**dplyr::setdiff(y, z)**

Rows that appear in y but not 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.

**dplyr::intersect(y, z)**

Rows that appear in both y and z.

**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.

**dplyr::setdiff(y, z)**

Rows that appear in y but not 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.

**dplyr::intersect(y, z)**

Rows that appear in both y and z.

**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.

**dplyr::setdiff(y, z)**

Rows that appear in y but not 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.

**dplyr::intersect(y, z)**

Rows that appear in both y and z.

**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.



# Data Wrangling with dplyr and tidyr

Cheat Sheet



## 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	
..	...	...	...	..

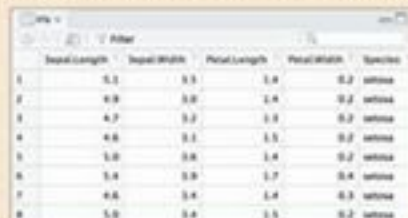
Variables not shown: Petal.Width (dbl), Species (fctr)

**dplyr::glimpse(iris)**

Information dense summary of tbl data.

**utils::View(iris)**

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



**dplyr::%>%**

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

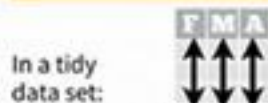
$x \rightsquigarrow f(y)$  is the same as  $f(x, y)$

$y \rightsquigarrow f(x, \dots, z)$  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)
```

## Tidy Data - A foundation for wrangling in R



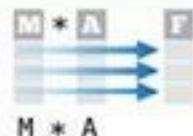
In a tidy data set:

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.



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



**tidyr::gather(cases, "year", "n", 2:4)**

Gather columns into rows.



**tidyr::separate(storms, date, c("y", "m", "d"))**

Separate one column into several.



**tidyr::spread(pollution, size, amount)**

Spread rows into columns.



**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).

## 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("x"))  
  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("x.*"))  
  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.
```

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

<	Less than	!=	Not equal to
>	Greater than	is.na()	Group membership
==	Equal to	is.na()	Is NA
<=	Less than or equal to	is.na()	Is not NA
>=	Greater than or equal to	!is.na(), xor, any, all	Boolean operators



# Python For Data Science Cheat Sheet

## PySpark Basics

Learn Python for data science interactively at [www.datacamp.com](https://www.datacamp.com)



### Spark

PySpark is the Spark Python API that exposes the Spark programming model to Python



### Initializing Spark

#### SparkContext

```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master = "local[2]")
```

#### Inspect SparkContext

```
>>> sc.version
Retrieve SparkContext version
>>> sc.pythonVer
Retrieve Python version
>>> sc.master
Master URL, to connect to
>>> sc.sparkHome
Path where Spark is installed on worker nodes
>>> sc.sparkUser
Retrieve name of the Spark User running SparkContext
>>> sc.appName
Return application name
>>> sc.applicationId
Retrieve application ID
>>> sc.defaultParallelism
Return default level of parallelism
>>> sc.defaultMinPartitions
Default minimum number of partitions for RDDs
```

#### Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = SparkConf()
>>> conf.setMaster("local")
>>> conf.setAppName("My app")
>>> conf.set("spark.executor.memory", "1g")
>>> sc = SparkContext(conf=conf)
```

#### Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[8] --py-file code.py
```

Set which master the context connects to with the `--master` argument, and add Python `.zip`, `.egg` or `.py` files to the runtime path by passing a comma-separated list to `--py-files`.

#### Loading Data

##### Parallelized Collections

```
>>> rdd = sc.parallelize([(1, 'a', 7), (1, 'a', 2), (1, 'b', 2)])
>>> rdd2 = sc.parallelize([(1, 'a', 2), (1, 'b', 3), (1, 'b', 1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize([(1, 'a', 1), (1, 'a', 2), (1, 'a', 3), (1, 'b', 1), (1, 'b', 2), (1, 'b', 3)])
```

##### External Data

Read either one text file from HDFS, a local file system or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`.

```
>>> textFile = sc.textFile("mydir/text/*", 1)
>>> textFile2 = sc.wholeTextFiles("mydir/text/*")
```

### Retrieving RDD Information

#### Basic Information

```
>>> rdd.getNumPartitions()
List the number of partitions
>>> rdd.count()
Count RDD instances
>>> rdd.countByKey()
Count RDD instances by key
>>> rdd.countByValue()
Count RDD instances by value
>>> rdd.collectAsMap()
Return (key,value) pairs as a dictionary
>>> rdd3.sum()
Sum of RDD elements
>>> sc.parallelize([1]).isEmpty()
Check whether RDD is empty
```

#### Summary

```
>>> rdd3.max()
Maximum value of RDD elements
>>> rdd3.min()
Minimum value of RDD elements
>>> rdd3.mean()
Mean value of RDD elements
>>> rdd3.stdev()
Standard deviation of RDD elements
>>> rdd3.variance()
Compute variance of RDD elements
>>> rdd3.histogram()
Compute histogram by bins
>>> rdd3.stats()
Summary statistics (count, mean, stdev, max & min)
```

### Applying Functions

```
>>> rdd.map(lambda x: x[1], x[0])
Apply a function to each RDD element
>>> rdd2 = rdd.flatMap(lambda x: x[1], x[0])
Apply a function to each RDD element and flatten the result
>>> rdd5.collect()
>>> rdd4.flatMapValues(lambda x: x)
Apply a flatMap function to each (key,value) pair of RDD without changing the keys
```

### Selecting Data

```
>>> rdd.collect()
Return a list with all RDD elements
>>> rdd.take(2)
Take first 2 RDD elements
>>> rdd.first()
Take first RDD element
>>> rdd.top(2)
Take top 2 RDD elements
>>> rdd1.sample(False, 0.15, 81).collect()
Return sampled subset of RDD
>>> rdd.filter(lambda x: "a" in x)
Filter the RDD
>>> rdd.distinct().collect()
Return distinct RDD values
>>> rdd.keys().collect()
Return (key,value) RDD's keys
```

### Iterating

```
>>> def g(x): print(x)
>>> rdd.foreach(g)
Apply a function to all RDD elements
```

### Reshaping Data

#### Reducing

```
>>> rdd.reduceByKey(lambda x, y: x+y)
Merge the RDD values for each key
>>> rdd.reduce(lambda a, b: a + b)
Merge the RDD values
```

#### Grouping by

```
>>> rdd3.groupBy(lambda x: x[1])
Return RDD of grouped values
>>> rdd3.groupByKey()
Group RDD by key
```

#### Aggregating

```
>>> seqOp = (lambda x, y: x[0]+y, x[1]+1)
>>> combOp = (lambda x, y: x[0]+y[0], x[1]+y[1])
>>> rdd3.aggregate((0,0), seqOp, combOp)
Aggregate RDD elements of each partition, and then the results
>>> rdd3.aggregateByKey(0, seqOp, combOp)
Aggregate values of each RDD key
>>> rdd3.fold(0, seqOp)
Aggregate the elements of each partition, and then the results
>>> rdd3.foldByKey(0, seqOp)
Merge the values for each key
>>> rdd3.keyBy(lambda x: x[1])
Create tuples of RDD elements by applying a function
```

### Mathematical Operations

```
>>> rdd.subtract(rdd2)
Return each RDD value not contained in RDD2
>>> rdd2.subtractByKey(rdd)
Return each (key,value) pair of RDD2 with no matching key in RDD
>>> rdd.cartesian(rdd2).collect()
Return the Cartesian product of RDD and RDD2
```

### Sort

```
>>> rdd2.sortBy(lambda x: x[1])
Sort RDD by given function
>>> rdd2.sortByKey()
Sort (key, value) RDD by key
```

### Repartitioning

```
>>> rdd.repartition(4)
New RDD with 4 partitions
>>> rdd.coalesce(1)
Decrease the number of partitions in the RDD to 1
```

### Saving

```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://name:node:port/parent/path",
    "org.apache.hadoop.mapreduce.lib.output.TextOutputFormat")
```

### Stopping SparkContext

```
>>> sc.stop()
```

### Execution

```
$ ./bin/spark-submit example/src/main/python/pi.py
```





## Stats - An alternative way to build a layer

Some plots visualize a **transformation** of the original data set. Use a **stat** to choose a common transformation to visualize, e.g. `a + geom_bar(stat = "bin")`



Each stat creates additional variables to map aesthetics to. These variables use a common `..name..` syntax.

stat functions and geom functions both combine a stat with a geom to make a layer, i.e. `stat_bin(geom="bar")` does the same as `geom_bar(stat="bin")`



`a + stat_bin(binswidth = 1, origin = 10)` 1D distributions  
`x, y | count, ..ncount.., density, ..ndensity..`  
`a + stat_bin2d(binswidth = 1, binsize = "x")`  
`x, y | count, ..ncount..`  
`a + stat_density2d(adjust = 1, kernel = "gaussian")`  
`x, y | count, ..density.., scaled`

`f + stat_bin2d(bins = 30, drop = TRUE)` 2D distributions  
`x, y | count, ..density..`  
`f + stat_binhex(bins = 30)`  
`x, y | count, ..density..`  
`f + stat_density2d(contour = TRUE, n = 100)`  
`x, y | color, size | ..level..`

`m + stat_contour(aes(z = z))` 3 Variables  
`x, y, z | order | ..level..`  
`m + stat_spoke(aes(radius = z, angle = z))`  
`angle, radius, x, yend, yend | ..xend.., ..yend..`  
`m + stat_summary_hex(aes(z = z, bins = 30, fun = mean))`  
`x, y, z | fill | ..value..`  
`m + stat_summary2d(aes(z = z, bins = 30, fun = mean))`  
`x, y, z | fill | ..value..`

`g + stat_boxplot(coef = 1.5)` Comparisons  
`x, y | lower, middle, upper, outliers`  
`g + stat_density2d(adjust = 1, kernel = "gaussian", scale = "area")`  
`x, y | density, scaled, count, ..n.., width`

`f + stat_ecdf(n = 40)` Functions  
`x, y | ..x.., ..y..`  
`f + stat_quantile(quantiles = c(0.25, 0.5, 0.75), formula = y ~ log(y), method = "log")`  
`x, y | quantile, ..x.., ..y..`  
`f + stat_smooth(method = "auto", formula = y ~ x, se = TRUE, n = 80, fullrange = FALSE, level = 0.9)`  
`x, y | ..x.., ..y.., ymin, ymax`

`ggplot() + stat_function(aes(x = 3:10), fun = dnorm, n = 101, args = list(0, 0.5))` General Purpose  
`x | y`

`f + stat_identity()`  
`ggplot() + stat_qq(aes(sample = 1:100), distribution = qt, dparams = list(0.5))`  
`sample, x, y | ..x.., ..y..`  
`f + stat_sum()`  
`x, y | size | ..sum..`  
`f + stat_summary(fun.data = "mean_ci_boot")`  
`f + stat_unique()`

## Scales

Scales control how a plot maps data values to the visual values of an aesthetic. To change the mapping, add a custom scale.



### General Purpose scales

Use with any aesthetic:

alpha, color, fill, linetype, shape, size

`scale_*_continuous()` - map cont' values to visual values  
`scale_*_discrete()` - map discrete values to visual values  
`scale_*_identity()` - use data values as visual values  
`scale_*_manual(values = c())` - map discrete values to manually chosen visual values

### X and Y location scales

Use with x or y aesthetics (x shown here)

`scale_x_date(labels = date_format("%m/%d"), breaks = date_breaks("2 weeks"))` - treat x values as dates. See `hrtptime` for label formats.  
`scale_x_datetime()` - treat x values as date times. Use same arguments as `scale_x_date()`.  
`scale_x_log10()` - Plot x on log10 scale  
`scale_x_reverse()` - Reverse direction of x axis  
`scale_x_sqrt()` - Plot x on square root scale

### Color and fill scales

Discrete

`n + b + geom_bar(aes(fill = f))`  
`n + scale_fill_brewer(palette = "blues")`  
For palette choices: `library(RColorBrewer)` `display.brewer.all()`  
`n + scale_fill_grey(start = 0.2, end = 0.8, na.value = "red")`

Continuous

`n + b + geom_dotplot(aes(fill = x))`  
`n + scale_fill_gradient(low = "red", high = "yellow")`  
`n + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 20)`  
`n + scale_fill_gradientn(colors = terrain.colors(50), breaks = rainbow(50), heat.colors(), topo.colors(), cm.colors(), RColorBrewer::brewer.pal(1))`

### Shape scales

`n + f + geom_point(aes(shape = f))`  
`n + scale_shape(solid = FALSE)`  
`n + scale_shape_manual(values = c(1, 2))`  
Shape values shown in chart on right

Manual shape values  
`n + f + geom_point(aes(shape = f))`  
`n + scale_shape(solid = FALSE)`  
`n + scale_shape_manual(values = c(1, 2))`  
Shape values shown in chart on right

### Size scales

`n + f + geom_point(aes(size = c(1)))`

`n + scale_size_area(max = 4)`  
Value mapped to area of circle, not radius

## Coordinate Systems

`r + b + geom_bar()`  
`xlim, ylim`  
The default cartesian coordinate system  
`r + coord_fixed(ratio = 1/2)`  
ratio, xlim, ylim  
Cartesian coordinates with fixed aspect ratio between x and y units  
`r + coord_flip()`  
xlim, ylim  
Flipped Cartesian coordinates  
`r + coord_polar(theta = "x", direction = 1)`  
theta, start, direction  
Polar coordinates  
`r + coord_trans(ytrans = "sqrt")`  
xtrans, ytrans, limx, limy  
Transformed cartesian coordinates. Set extras and strains to the name of a window function.

`z + coord_map(projection = "ortho", orientation = c(41, -74, 0))`  
projection, orientation, xlim, ylim  
Map projections from the `mapproj` package (mercator (default), aequalarea, lagrange, etc.)

## Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s + ggplot(mpg, aes(f, fill = drv))`  
`s + geom_bar(position = "dodge")`  
Arrange elements side by side  
`s + geom_bar(position = "fill")`  
Stack elements on top of one another, normalize height  
`s + geom_bar(position = "stack")`  
Stack elements on top of one another  
`f + geom_point(position = "jitter")`  
Add random noise to X and Y position of each element to avoid overplotting

Each position adjustment can be recast as a function with manual `width` and `height` arguments

`s + geom_bar(position = position_dodge(width = 1))`

## Themes

`y + theme_bw()`  
White background with grid lines  
`y + theme_classic()`  
White background no gridlines  
`y + theme_grey()`  
Grey background (default theme)  
`y + theme_minimal()`  
Minimal theme

`ggthemes` - Package with additional ggplot2 themes

## Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t + ggplot(mpg, aes(cty, hwy)) + geom_point()`  
`t + facet_grid(~ fl)`  
facet into columns based on fl  
`t + facet_grid(year ~ .)`  
facet into rows based on year  
`t + facet_grid(year ~ fl)`  
facet into both rows and columns  
`t + facet_wrap(~ fl)`  
wrap facets into a rectangular layout

Set scales to let axis limits vary across facets

`t + facet_grid(y ~ x, scales = "free")`  
x and y axis limits adjust to individual facets  
• `"free_x"` - x axis limits adjust  
• `"free_y"` - y axis limits adjust

Set labeller to adjust facet labels

`t + facet_grid(~ fl, labeller = label_both)`  
`t + facet_grid(~ fl, labeller = label_bquote(alpha ~ .x))`  
`t + facet_grid(~ fl, labeller = label_parsed)`

## Labels

`t + ggtitle("New Plot Title")`  
Add a main title above the plot  
`t + xlab("New X label")`  
Change the label on the X axis  
`t + ylab("New Y label")`  
Change the label on the Y axis  
`t + labs(title = "New title", x = "New x", y = "New y")`  
All of the above

## Legends

`t + theme(legend.position = "bottom")`  
Place legend at "bottom", "top", "left", or "right"  
`t + guides(color = "none")`  
Set legend type for each aesthetic: colorbar, legend, or none (no legend)  
`t + scale_fill_discrete(name = "Title", labels = c("A", "B", "C"))`  
Set legend title and labels with a scale function.

## Zooming

Without clipping (preferred)  
`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`  
With clipping (removes unseen data points)  
`t + xlim(0, 100) + ylim(10, 20)`  
`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(10, 20))`



# Data Visualization with ggplot2 Cheat Sheet

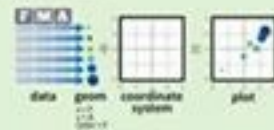


## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with **qplot()** or **ggplot()**

**qplot(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")**  
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

**ggplot(data = mpg, aes(x = cty, y = hwy))**  
Begins a plot that you finish by adding layers to. No defaults, but provides more control than qplot().

**ggplot(mpg, aes(hwy, cty)) +  
geom\_point(aes(color = cyl)) +  
geom\_smooth(method = "lm") +  
coord\_cartesian() +  
scale\_color\_gradient() +  
theme\_bw()**

Add a new layer to a plot with a **geom\_\***() or **stat\_\***() function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

**last\_plot()**

Returns the last plot

**ggsave("plot.png", width = 5, height = 5)**

Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

**Geoms** - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

## One Variable

### Continuous

**a <- ggplot(mpg, aes(hwy))**

**a + geom\_area(stat = "bin")**  
x, y, alpha, color, fill, linetype, size  
**b + geom\_area(aes(y = ..density..), stat = "bin")**  
**a + geom\_density(kernel = "gaussian")**  
x, y, alpha, color, fill, linetype, size, weight  
**b + geom\_density(aes(y = ..density..))**  
**a + geom\_dotplot()**  
x, y, alpha, color, fill  
**a + geom\_freqpoly()**  
x, y, alpha, color, linetype, size  
**b + geom\_freqpoly(aes(y = ..density..))**  
**a + geom\_histogram(binwidth = 5)**  
x, y, alpha, color, fill, linetype, size, weight  
**b + geom\_histogram(aes(y = ..density..))**

### Discrete

**b <- ggplot(mpg, aes(fill))**

**b + geom\_bar()**  
x, alpha, color, fill, linetype, size, weight

## Graphical Primitives

**c <- ggplot(map, aes(long, lat))**

**c + geom\_polygon(aes(group = group))**  
x, y, alpha, color, fill, linetype, size

**d <- ggplot(economics, aes(date, unemploy))**

**d + geom\_path(lineend = "butt",  
linejoin = "round", linewidth = 1)**  
x, y, alpha, color, linetype, size  
**d + geom\_ribbon(aes(ymin = unemploy - 900,  
ymax = unemploy + 900))**  
x, y, alpha, color, fill, linetype, size

**e <- ggplot(seals, aes(x = long, y = lat))**

**e + geom\_segment(aes(xend = long + delta\_long,  
yend = lat + delta\_lat))**  
x, y, alpha, color, linetype, size  
**e + geom\_rect(aes(xmin = long, ymin = lat,  
xmax = long + delta\_long,  
ymax = lat + delta\_lat))**  
x, y, alpha, color, fill, linetype, size

## Two Variables

### Continuous X, Continuous Y

**f <- ggplot(mpg, aes(cty, hwy))**

**f + geom\_blank()**  
**f + geom\_jitter()**  
x, y, alpha, color, fill, shape, size  
**f + geom\_point()**  
x, y, alpha, color, fill, shape, size  
**f + geom\_quantile()**  
x, y, alpha, color, linetype, size, weight  
**f + geom\_rug(sides = "bl")**  
alpha, color, linetype, size  
**f + geom\_smooth(model = lm)**  
x, y, alpha, color, fill, linetype, size, weight  
**f + geom\_text(aes(label = cty))**  
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

### Discrete X, Continuous Y

**g <- ggplot(mpg, aes(class, hwy))**

**g + geom\_bar(stat = "identity")**  
x, y, alpha, color, fill, linetype, size, weight  
**g + geom\_boxplot()**  
lower, middle, upper, x, y, alpha, color, fill, linetype, shape, size, weight  
**g + geom\_dotplot(binaxis = "y",  
stackdir = "center")**  
x, y, alpha, color, fill  
**g + geom\_violin(scale = "area")**  
x, y, alpha, color, fill, linetype, size, weight

### Discrete X, Discrete Y

**h <- ggplot(diamonds, aes(cut, color))**

**h + geom\_jitter()**  
x, y, alpha, color, fill, shape, size

### Continuous Bivariate Distribution

**i <- ggplot(movies, aes(year, rating))**

**i + geom\_bin2d(binwidth = c(5, 0.5))**  
x, y, alpha, color, fill, linetype, size, weight  
**i + geom\_density2d()**  
x, y, alpha, color, linetype, size  
**i + geom\_hex()**  
x, y, alpha, color, fill, size

### Continuous Function

**j <- ggplot(economics, aes(date, unemploy))**

**j + geom\_area()**  
x, y, alpha, color, fill, linetype, size  
**j + geom\_line()**  
x, y, alpha, color, linetype, size  
**j + geom\_step(direction = "hv")**  
x, y, alpha, color, linetype, size

### Visualizing error

**df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)**  
**k <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))**

**k + geom\_crossbar(latten = 2)**  
x, y, ymax, ymin, alpha, color, fill, linetype, size  
**k + geom\_errorbar()**  
x, y, ymax, ymin, alpha, color, linetype, size, width (also **geom\_errorbarh()**)  
**k + geom\_linerange()**  
x, y, ymax, ymin, alpha, color, linetype, size  
**k + geom\_pointrange()**  
x, y, ymax, ymin, alpha, color, fill, linetype, shape, size

### Maps

**data <- data.frame(murder = USArrests\$Murder,  
state = tolower(row.names(USArrests)))**  
**map <- map\_data("state")**  
**i <- ggplot(data, aes(fill = murder))**  
**i + geom\_map(aes(map\_id = state), map = map) +  
expand\_limits(x = map\$long, y = map\$lat)**  
map\_id, alpha, color, fill, linetype, size

## Three Variables

**seals2 <- with(seals, sqrt(delta\_long^2 + delta\_lat^2))**  
**m <- ggplot(seals, aes(long, lat))**

**m + geom\_raster(aes(z = z))**  
x, y, z, alpha, color, linetype, size, weight  
**m + geom\_tile(aes(fill = z))**  
x, y, alpha, color, fill, linetype, size



# A mostly complete chart of Neural Networks

©2016 Fjodor van Veen - [asimovinstitute.org](http://asimovinstitute.org)

Backfed Input Cell

Input Cell

Noisy Input Cell

Hidden Cell

Probabilistic Hidden Cell

Spiking Hidden Cell

Output Cell

Match Input Output Cell

Recurrent Cell

Memory Cell

Different Memory Cell

Kernel

Convolution or Pool

Perceptron (P)



Feed Forward (FF)



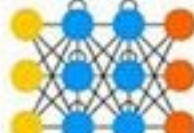
Radial Basis Network (RBF)



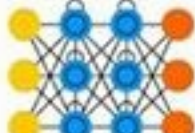
Deep Feed Forward (DFF)



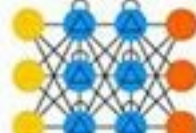
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



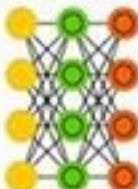
Gated Recurrent Unit (GRU)



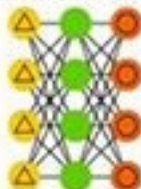
Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)



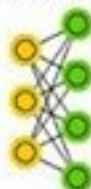
Hopfield Network (HN)



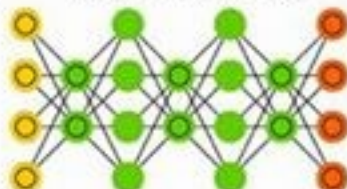
Boltzmann Machine (BM)



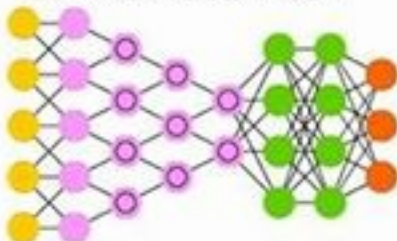
Restricted BM (RBM)



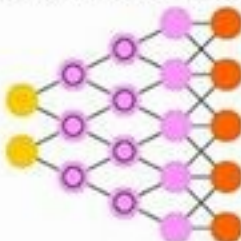
Deep Belief Network (DBN)



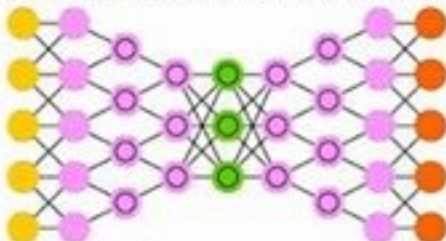
Deep Convolutional Network (DCN)



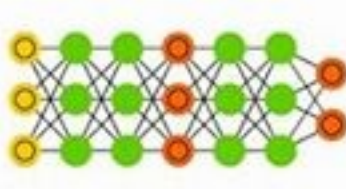
Deconvolutional Network (DN)



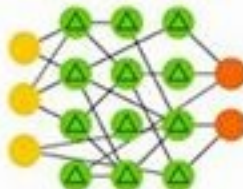
Deep Convolutional Inverse Graphics Network (DCIGN)



Generative Adversarial Network (GAN)



Liquid State Machine (LSM)



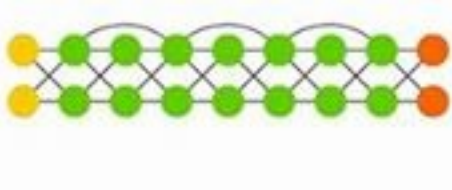
Extreme Learning Machine (ELM)



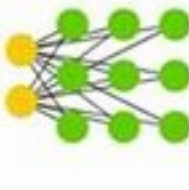
Echo State Network (ESN)



Deep Residual Network (DRN)



Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)





# Python For Data Science Cheat Sheet

## Scikit-Learn

Learn Python for data science interactively at [www.DataCamp.com](https://www.datacamp.com)



### Scikit-learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.



#### A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.cross_validation import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

### Loading The Data

Also see NumPy & Pandas

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10,5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'M', 'F', 'F'])
>>> X[X < 0.7] = 0
```

### Training And Test Data

```
>>> from sklearn.cross_validation import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X,
                                                         y,
                                                         random_state=0)
```

### Preprocessing The Data

#### Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

#### Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

#### Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

### Create Your Model

#### Supervised Learning Estimators

##### Linear Regression

```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

##### Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

##### Naive Bayes

```
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
```

##### KNN

```
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

#### Unsupervised Learning Estimators

##### Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

##### K Means

```
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=3, random_state=0)
```

### Model Fitting

#### Supervised learning

```
>>> lr.fit(X, y)
>>> knn.fit(X_train, y_train)
>>> svc.fit(X_train, y_train)
```

Fit the model to the data

#### Unsupervised Learning

```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```

Fit the model to the data

Fit to data, then transform it

### Prediction

#### Supervised Estimators

```
>>> y_pred = svc.predict(np.random.random((2,5)))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```

Predict labels

Predict labels

Estimate probability of a label

#### Unsupervised Estimators

```
>>> y_pred = k_means.predict(X_test)
```

Predict labels in clustering algos

### Evaluate Your Model's Performance

#### Classification Metrics

##### Accuracy Score

```
>>> knn.score(X_test, y_test)
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)
```

Estimator score method  
Metric scoring functions

##### Classification Report

```
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))
```

Precision, recall, f1-score  
and support

##### Confusion Matrix

```
>>> from sklearn.metrics import confusion_matrix
>>> print(confusion_matrix(y_test, y_pred))
```

#### Regression Metrics

##### Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
>>> y_true = [3, -0.5, 2]
>>> mean_absolute_error(y_true, y_pred)
```

##### Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
>>> mean_squared_error(y_test, y_pred)
```

##### R<sup>2</sup> Score

```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_true, y_pred)
```

#### Clustering Metrics

##### Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

##### Homogeneity

```
>>> from sklearn.metrics import homogeneity_score
>>> homogeneity_score(y_true, y_pred)
```

##### V-measure

```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_true, y_pred)
```

#### Cross-Validation

```
>>> from sklearn.cross_validation import cross_val_score
>>> print(cross_val_score(knn, X_train, y_train, cv=4))
>>> print(cross_val_score(lr, X, y, cv=2))
```

### Tune Your Model

#### Grid Search

```
>>> from sklearn.grid_search import GridSearchCV
>>> params = {"n_neighbors": np.arange(1,3),
             "metric": ["euclidean", "cityblock"]}
>>> grid = GridSearchCV(estimator=knn,
                       param_grid=params)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

#### Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {"n_neighbors": range(1,5),
             "weights": ["uniform", "distance"]}
>>> rsearch = RandomizedSearchCV(estimator=knn,
                                param_distributions=params,
                                cv=4,
                                n_iter=8,
                                random_state=5)
>>> rsearch.fit(X_train, y_train)
>>> print(rsearch.best_score_)
```





# Python For Data Science Cheat Sheet

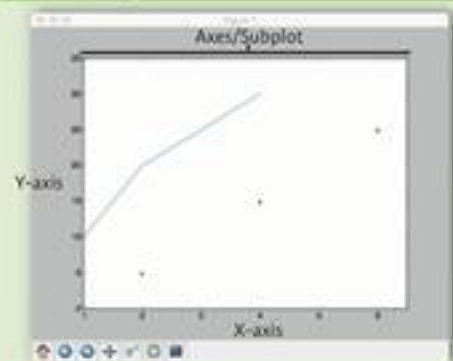
## Matplotlib

Learn Python Interactively at [www.datacamp.com](https://www.datacamp.com)



## Plot Anatomy & Workflow

### Plot Anatomy



### Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4]
>>> y = [10,20,25,30]
>>> fig = plt.figure()
>>> ax = fig.add_subplot(111)
>>> ax.plot(x, y, color='lightblue', linewidth=3)
>>> ax.scatter([2,4,6],
>>>            [5,15,25],
>>>            color='darkgreen',
>>>            marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show()
```

## 4 Customize Plot

### Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha = 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img,
>>>                cmap='seismic')
```

### Markers

```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x, y, marker="s")
>>> ax.plot(x, y, marker="o")
```

### Linestyles

```
>>> plt.plot(x, y, linewidth=4.0)
>>> plt.plot(x, y, ls='solid')
>>> plt.plot(x, y, ls='--')
>>> plt.plot(x, y, '--', x**2, y**2, '-.')
>>> plt.setp(lines, color='r', linewidth=4.0)
```

### Text & Annotations

```
>>> ax.text(1,
>>>         2.1,
>>>         'Example Graph',
>>>         style='italic')
>>> ax.annotate("Sine",
>>>             xy=(8, 0),
>>>             xycoords='data',
>>>             xytext=(10.5, 0),
>>>             textcoords='data',
>>>             arrowprops=dict(arrowstyle="->",
>>>                             connectionstyle="arc3"),)
```

### Mattext

```
>>> plt.title(r'$\sigma_i=155$', fontsize=20)
```

### Limits, Legends & Layouts

#### Limits & Autoscaling

```
>>> ax.margins(x=0.0, y=0.1)
>>> ax.axis('equal')
>>> ax.set(xlim=[0,10.5], ylim=[-1.5,1.5])
>>> ax.set_xlim(0,10.5)
```

#### Legends

```
>>> ax.set(title='An Example Axes',
>>>         ylabel='Y-Axis',
>>>         xlabel='X-Axis')
>>> ax.legend(loc='best')
```

#### Ticks

```
>>> ax.xaxis.set(ticks=range(1,5),
>>>              ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',
>>>                 direction='inout',
>>>                 length=10)
```

#### Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
>>>                       hspace=0.3,
>>>                       left=0.125,
>>>                       right=0.9,
>>>                       top=0.9,
>>>                       bottom=0.1)
```

#### Axis Spines

```
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_position(('outward',10))
```

Add padding to a plot  
Set the aspect ratio of the plot to 1  
Set limits for x-and y-axis  
Set limits for x-axis

Set a title and x-and y-axis labels

No overlapping plot elements

Manually set x-ticks

Make y-ticks longer and go in and out

Adjust the spacing between subplots

Fit subplot(s) in to the figure area

Make the top axis line for a plot invisible  
Move the bottom axis line outward

## 1 Prepare The Data

Also see Lists & NumPy

### 1D Data

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
```

### 2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]
>>> U = -1 - X**2 + Y
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

## 2 Create Plot

```
>>> import matplotlib.pyplot as plt
```

### Figure

```
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

### Axes

All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax3 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2,ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

## 3 Plotting Routines

### 1D Data

```
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x, y)
>>> ax.scatter(x, y)
>>> axes[0,0].bar([1,2,3],[3,4,5])
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2])
>>> axes[1,1].axhline(0.45)
>>> axes[0,1].axvline(0.65)
>>> ax.fill(x, y, color='blue')
>>> ax.fill_between(x, y, color='yellow')
```

Draw points with lines or markers connecting them  
Draw unconnected points, scaled or colored  
Plot vertical rectangles (constant width)  
Plot horizontal rectangles (constant height)  
Draw a horizontal line across axes  
Draw a vertical line across axes  
Draw filled polygons  
Fill between y-values and 0

### 2D Data or Images

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img,
>>>                cmap='gist_earth',
>>>                interpolation='nearest',
>>>                vmin=-2,
>>>                vmax=2)
```

Colormapped or RGB arrays

### Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)
>>> axes[1,1].quiver(y, z)
>>> axes[0,1].streamplot(X, Y, U, V)
```

Add an arrow to the axes  
Plot a 2D field of arrows  
Plot a 2D field of arrows

### Data Distributions

```
>>> ax1.hist(y)
>>> ax3.boxplot(y)
>>> ax3.violinplot(z)
```

Plot a histogram  
Make a box and whisker plot  
Make a violin plot

```
>>> axes2[0].pcolor(data2)
>>> axes2[0].pcolormesh(data2)
>>> CS = plt.contour(Y,X,U)
>>> axes2[2].contourf(data1)
>>> axes2[2] = ax.clabel(CS)
```

Pseudocolor plot of 2D array  
Pseudocolor plot of 2D array  
Plot contours  
Plot filled contours  
Label a contour plot

## 5 Save Plot

### Save figures

```
>>> plt.savefig('foo.png')
```

### Save transparent figures

```
>>> plt.savefig('foo.png', transparent=True)
```

## 6 Show Plot

```
>>> plt.show()
```

## Close & Clear

```
>>> plt.cla()
>>> plt.clf()
>>> plt.close()
```

Clear an axis  
Clear the entire figure  
Close a window

DataCamp

Learn Python for Data Science Interactively





# Python For Data Science Cheat Sheet

## SciPy - Linear Algebra

Learn More Python for Data Science [Interactively](#) at [www.datacamp.com](#)



### SciPy

The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.



### Interacting With NumPy

Also see NumPy

```
>>> import numpy as np
>>> a = np.array([1,2,3])
>>> b = np.array([(1+5j, 2j, 3j), (4j, 5j, 6j)])
>>> c = np.array([(1.5, 2, 3), (4, 5, 6)], [(3, 2, 1), (4, 5, 6)])
```

### Index Tricks

```
>>> np.mgrid[0:5,0:5]
>>> np.ogrid[0:2,0:2]
>>> np.r_[3, [0]*5, -1:1:10j]
>>> np.c_[b, c]
```

Create a dense meshgrid  
Create an open meshgrid  
Stack arrays vertically (row-wise)  
Create stacked column-wise arrays

### Shape Manipulation

```
>>> np.transpose(b)
>>> b.flatten()
>>> np.hstack((b, c))
>>> np.vstack((a, b))
>>> np.hsplit(c, 2)
>>> np.vsplit(d, 2)
```

Permute array dimensions  
Flatten the array  
Stack arrays horizontally (column-wise)  
Stack arrays vertically (row-wise)  
Split the array horizontally at the 2nd index  
Split the array vertically at the 2nd index

### Polynomials

```
>>> from numpy import poly1d
>>> p = poly1d([3, 4, 5])
```

Create a polynomial object

### Vectorizing Functions

```
>>> def myfunc(a):
    if a < 0:
        return a*2
    else:
        return a/2
>>> np.vectorize(myfunc)
```

Vectorize functions

### Type Handling

```
>>> np.real(b)
>>> np.imag(b)
>>> np.real_if_close(c, tol=1000)
>>> np.cast['f'](np.pi)
```

Return the real part of the array elements  
Return the imaginary part of the array elements  
Return a real array if complex parts close to 0  
Cast object to a data type

### Other Useful Functions

```
>>> np.angle(b, deg=True)
>>> g = np.linspace(0, np.pi, num=5)
>>> g[3:] += np.pi
>>> np.unwrap(g)
>>> np.logspace(0, 10, 3)
>>> np.select([c<4], [c*2])
>>> misc.factorial(a)
>>> misc.comb(10, 3, exact=True)
>>> misc.central_diff_weights(3)
>>> misc.derivative(myfunc, 1.0)
```

Return the angle of the complex argument  
Create an array of evenly spaced values (number of samples)  
Unwrap  
Create an array of evenly spaced values (log scale)  
Return values from a list of arrays depending on conditions  
Factorial  
Combine N things taken at k time  
Weights for N-point central derivative  
Find the n-th derivative of a function at a point

## Linear Algebra

You'll use the `linalg` and `sparse` modules. Note that `scipy.linalg` contains and expands on `numpy.linalg`.

```
>>> from scipy import linalg, sparse
```

### Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.mat(np.random.random((10,5)))
>>> D = np.mat([(3,4), (5,6)])
```

### Basic Matrix Routines

```
>>> A.I
>>> linalg.inv(A)
>>> A.T
>>> A.H
```

Inverse  
Inverse  
Transpose matrix  
Conjugate transposition

```
>>> np.trace(A)
>>> linalg.norm(A)
>>> linalg.norm(A, 1)
>>> linalg.norm(A, np.inf)
```

Trace  
Trace  
Frobenius norm  
L1 norm (max column sum)  
L inf norm (max row sum)

```
>>> np.linalg.matrix_rank(C)
>>> linalg.det(A)
>>> linalg.solve(A, b)
>>> E = np.mat(a).T
>>> linalg.lstsq(F, E)
>>> linalg.pinv(C)
>>> linalg.pinv2(C)
```

Rank  
Determinant  
Solving linear problems  
Solver for dense matrices  
Solver for dense matrices  
Least-squares solution to linear matrix equation  
Matrix rank  
Determinant  
Generalized inverse  
Compute the pseudo-inverse of a matrix (least-squares solver)  
Compute the pseudo-inverse of a matrix (SVD)

### Creating Sparse Matrices

```
>>> F = np.eye(3, k=1)
>>> G = np.mat(np.identity(2))
>>> C[C > 0.5] = 0
>>> H = sparse.csr_matrix(C)
>>> I = sparse.csc_matrix(D)
>>> J = sparse.dok_matrix(A)
>>> E.todense()
>>> sparse.ispmatrix_csc(A)
```

Create a 2X2 identity matrix  
Create a 2x2 identity matrix  
Compressed Sparse Row matrix  
Compressed Sparse Column matrix  
Dictionary Of Keys matrix  
Sparse matrix to full matrix  
Identify sparse matrix

### Sparse Matrix Routines

```
>>> sparse.linalg.inv(I)
>>> sparse.linalg.norm(I)
>>> sparse.linalg.solve(H, I)
```

Inverse  
Inverse  
Norm  
Norm  
Solver for sparse matrices

### Sparse Matrix Functions

```
>>> sparse.linalg.expm(I)
```

Sparse matrix exponential

### Asking For Help

```
>>> help(scipy.linalg.diagvd)
>>> np.info(np.matrix)
```

Also see NumPy

### Matrix Functions

```
>>> np.add(A, D)
>>> np.subtract(A, D)
>>> np.divide(A, D)
>>> A @ D
>>> np.multiply(D, A)
>>> np.dot(A, D)
>>> np.vdot(A, D)
>>> np.inner(A, D)
>>> np.outer(A, D)
>>> np.tensordot(A, D)
>>> np.kron(A, D)
```

Addition  
Subtraction  
Division  
Multiplication operator (Python 3)  
Multiplication  
Dot product  
Vector dot product  
Inner product  
Outer product  
Tensor dot product  
Kronecker product

```
>>> linalg.expm(A)
>>> linalg.expm2(A)
>>> linalg.expm3(D)
```

Exponential Functions  
Matrix exponential  
Matrix exponential (Taylor Series)  
Matrix exponential (eigenvalue decomposition)

```
>>> linalg.logm(A)
```

Logarithm Function  
Matrix logarithm

```
>>> linalg.sirm(D)
>>> linalg.cosm(D)
>>> linalg.tanm(A)
```

Trigonometric Functions  
Matrix sine  
Matrix cosine  
Matrix tangent

```
>>> linalg.sinhm(D)
>>> linalg.coshm(D)
>>> linalg.tanhm(A)
```

Hyperbolic Trigonometric Functions  
Hyperbolic matrix sine  
Hyperbolic matrix cosine  
Hyperbolic matrix tangent

```
>>> np.signm(A)
```

Matrix Sign Function  
Matrix sign function

```
>>> linalg.sqrtm(A)
```

Matrix Square Root  
Matrix square root

```
>>> linalg.furm(A, lambda x: x*x)
```

Arbitrary Functions  
Evaluate matrix function

### Decompositions

```
>>> la, v = linalg.eig(A)
>>> l1, l2 = la
>>> v[:,0]
>>> v[:,1]
>>> linalg.eigvals(A)
```

Eigenvalues and Eigenvectors  
Solve ordinary or generalized eigenvalue problem for square matrix  
Unpack eigenvalues  
First eigenvector  
Second eigenvector  
Unpack eigenvalues

```
>>> U, s, Vh = linalg.svd(B)
>>> M, N = B.shape
>>> Sig = linalg.diagsvd(s, M, N)
```

Singular Value Decomposition  
Singular Value Decomposition (SVD)  
Construct sigma matrix in SVD

```
>>> P, L, U = linalg.lu(C)
```

LU Decomposition  
LU Decomposition

### Sparse Matrix Decompositions

```
>>> la, v = sparse.linalg.eigs(F, 1)
>>> sparse.linalg.svds(H, 2)
```

Eigenvalues and eigenvectors  
SVD

DataCamp

Learn Python for Data Science [Interactively](#)





# Python For Data Science Cheat Sheet

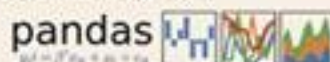
## Pandas Basics

Learn Python for Data Science interactively at [www.datacamp.com](https://www.datacamp.com)



### Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.



Use the following import convention:

```
>>> import pandas as pd
```

### Pandas Data Structures

#### Series

A one-dimensional labeled array capable of holding any data type

A	3
B	-5
C	7
D	4

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

#### DataFrame

Columns			
	Country	Capital	Population
1	Belgium	Brussels	11190846
2	India	New Delhi	1303171035
3	Brazil	Brasilia	207847528

A two-dimensional labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
            'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
            'Population': [11190846, 1303171035, 207847528]}
```

```
>>> df = pd.DataFrame(data,
                      columns=['Country', 'Capital', 'Population'])
```

### I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> pd.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')

Read multiple sheets from the same file
>>> xlsx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

### Asking For Help

```
>>> help(pd.Series.loc)
```

### Selection

Also see NumPy Arrays

#### Getting

```
>>> s['b']
-5
Get one element

>>> df[1:]
   Country  Capital  Population
1  India  New Delhi  1303171035
2  Brazil  Brasilia  207847528
Get subset of a DataFrame
```

### Selecting, Boolean Indexing & Setting

#### By Position

```
>>> df.iloc[[0], [0]]
Select single value by row & column

'Belgium'

>>> df.iat[[0], [0]]
'Belgium'
```

#### By Label

```
>>> df.loc[[0], ['Country']]
'Belgium'

>>> df.at[[0], ['Country']]
'Belgium'
```

#### By Label/Position

```
>>> df.ix[2]
Country      Brazil
Capital    Brasilia
Population  207847528
Select single row of subset of rows
```

```
>>> df.ix[:, 'Capital']
0    Brussels
1    New Delhi
2    Brasilia
Select a single column of subset of columns
```

```
>>> df.ix[1, 'Capital']
'New Delhi'
Select rows and columns
```

#### Boolean Indexing

```
>>> s[~(s > 1)]
Series s where value is not >1
>>> s[(s < -1) | (s > 2)]
s where value is <-1 or >2
>>> df[df['Population'] > 1200000000]
Use filter to adjust DataFrame
```

#### Setting

```
>>> s['a'] = 6
Set index a of Series s to 6
```

### Dropping

```
>>> s.drop(['a', 'c'])
Drop values from rows (axis=0)
>>> df.drop('Country', axis=1)
Drop values from columns (axis=1)
```

### Sort & Rank

```
>>> df.sort_index(by='Country')
Sort by row or column index
>>> s.order()
Sort a series by its values
>>> df.rank()
Assign ranks to entries
```

### Retrieving Series/DataFrame Information

#### Basic Information

```
>>> df.shape
(rows, columns)
>>> df.index
Describe index
>>> df.columns
Describe DataFrame columns
>>> df.info()
Info on DataFrame
>>> df.count()
Number of non-NA values
```

#### Summary

```
>>> df.sum()
Sum of values
>>> df.cumsum()
Cumulative sum of values
>>> df.min()/df.max()
Minimum/maximum values
>>> df.idxmin()/df.idxmax()
Minimum/Maximum index value
>>> df.describe()
Summary statistics
>>> df.mean()
Mean of values
>>> df.median()
Median of values
```

### Applying Functions

```
>>> f = lambda x: x*2
>>> df.apply(f)
Apply function
>>> df.applymap(f)
Apply function element-wise
```

### Data Alignment

#### Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
a    10.0
b     NaN
c     5.0
d     7.0
```

#### Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
a    10.0
b    -5.0
c     5.0
d     7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

#### Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)

read_sql() is a convenience wrapper around read_sql_table() and
read_sql_query()

>>> pd.to_sql('myDf', engine)
```

DataCamp

Learn Python for Data Science interactively





## Summarize Data

**df['w'].value\_counts()**

Count number of rows with each unique value of variable  
**len(df)**

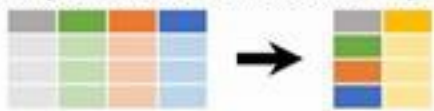
# of rows in DataFrame.

**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)**

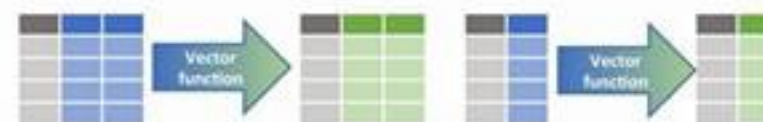
Compute and append one or more new columns.

**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.

**min(axis=1)**

Element-wise min.

**clip(lower=-10,upper=10)**

Trim values at input thresholds

**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.

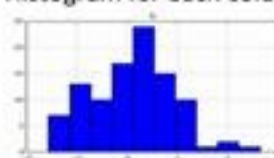
**cumprod()**

Cumulative product.

## Plotting

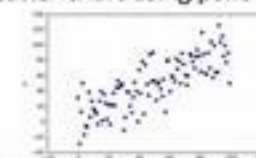
**df.plot.hist()**

Histogram for each column



**df.plot.scatter(x='w',y='h')**

Scatter chart using pairs of points



## Combine Data Sets



Standard Joins

x1	x2	x3
A	1	T
B	2	F
C	3	NaN

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

x1	x2	x3
A	1.0	T
B	2.0	F
D	NaN	T

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

x1	x2	x3
A	1	T
B	2	F

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

x1	x2	x3
A	1	T
B	2	F
C	3	NaN
D	NaN	T

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

Filtering Joins

x1	x2
A	1
B	2

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

x1	x2
C	3

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



Set-like Operations

x1	x2
B	2
C	3

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

x1	x2
A	1
B	2
C	3
D	4

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

x1	x2
A	1

**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 pandas

Cheat Sheet

<http://pandas.pydata.org>

## 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']))
```

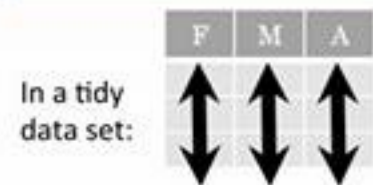
Create DataFrame with a MultiIndex

## Method Chaining

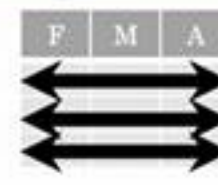
Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)  
     .rename(columns={  
         'variable': 'var',  
         'value': 'val'})  
     .query('val >= 200'))
```

## Tidy Data – A foundation for wrangling in pandas

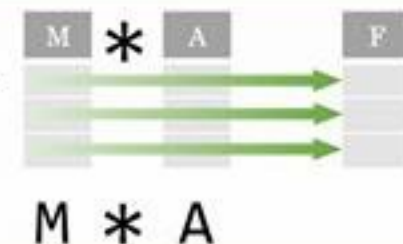


In a tidy data set:  
Each **variable** is saved in its own **column**



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.



## 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.



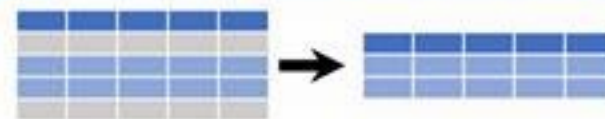
`pd.concat([df1, df2])`  
Append rows of DataFrames



`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
```

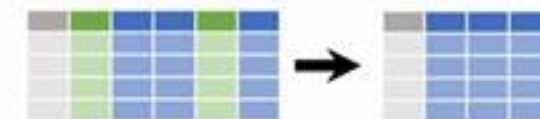
## Subset Observations (Rows)



```
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.
```

## 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

regex	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







# Python For Data Science Cheat Sheet

## Keras

Learn Python for data science interactively at [www.datacamp.com](https://www.datacamp.com)



### Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

#### A Basic Example

```
>>> import numpy as np
>>> from keras.models import Sequential
>>> from keras.layers import Dense
>>> data = np.random.random((1000,100))
>>> labels = np.random.randint(2,size=(1000,1))
>>> model = Sequential()
>>> model.add(Dense(12,
                    activation='relu',
                    input_dim=100))
>>> model.add(Dense(1, activation='sigmoid'))
>>> model.compile(optimizer='rmsprop',
                 loss='binary_crossentropy',
                 metrics=['accuracy'])
>>> model.fit(data, labels, epochs=10, batch_size=32)
>>> predictions = model.predict(data)
```

### Data

(Also see [NumPy](#), [Pandas](#), & [Scikit-Learn](#))

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you split the data in training and test sets, for which you can also resort to the `train_test_split` module of [sklearn.cross\\_validation](#).

#### Keras Data Sets

```
>>> from keras.datasets import boston_housing,
    mnist, cifar10, imdb
>>> (x_train,y_train), (x_test,y_test) = mnist.load_data()
>>> (x_train,y_train2), (x_test2,y_test2) = boston_housing.load_data()
>>> (x_train,y_train3), (x_test3,y_test3) = cifar10.load_data()
>>> (x_train,y_train4), (x_test4,y_test4) = imdb.load_data(word=20000)
>>> num_classes = 10
```

### Other

```
>>> from urllib.request import urlopen
>>> data = np.loadtxt(urlopen("https://archive.ics.uci.edu/ml/machine-learning-databases/pima-indian-diabetes/pima-indian-diabetes.data"),delimiter=",")
>>> X = data[:,0:8]
>>> y = data[:,9]
```

### Preprocessing

#### Sequence Padding

```
>>> from keras.preprocessing import sequence
>>> x_train4 = sequence.pad_sequences(x_train4,maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4,maxlen=80)
```

#### One-Hot Encoding

```
>>> from keras.utils import to_categorical
>>> y_train = to_categorical(y_train, num_classes)
>>> y_test = to_categorical(y_test, num_classes)
>>> y_train3 = to_categorical(y_train3, num_classes)
>>> y_test3 = to_categorical(y_test3, num_classes)
```

### Model Architecture

#### Sequential Model

```
>>> from keras.models import Sequential
>>> model = Sequential()
>>> model2 = Sequential()
>>> model3 = Sequential()
```

#### Multilayer Perceptron (MLP)

##### Binary Classification

```
>>> from keras.layers import Dense
>>> model.add(Dense(12,
                    input_dim=8,
                    kernel_initializer='uniform',
                    activation='relu'))
>>> model.add(Dense(8, kernel_initializer='uniform', activation='relu'))
>>> model.add(Dense(1, kernel_initializer='uniform', activation='sigmoid'))
```

##### Multi-Class Classification

```
>>> from keras.layers import Dropout
>>> model.add(Dense(512, activation='relu', input_shape=(784,)))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(512, activation='relu'))
>>> model.add(Dropout(0.2))
>>> model.add(Dense(10, activation='softmax'))
```

##### Regression

```
>>> model.add(Dense(64, activation='relu', input_dim=train_data.shape[1]))
>>> model.add(Dense(1))
```

#### Convolutional Neural Network (CNN)

```
>>> from keras.layers import Activation, Conv2D, MaxPooling2D, Flatten
>>> model2.add(Conv2D(32, (3,3), padding='same', input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(32, (3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Conv2D(64, (3,3), padding='same'))
>>> model2.add(Activation('relu'))
>>> model2.add(Conv2D(64, (3,3)))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2,2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Flatten())
>>> model2.add(Dense(512))
>>> model2.add(Activation('relu'))
>>> model2.add(Dropout(0.5))
>>> model2.add(Dense(num_classes))
>>> model2.add(Activation('softmax'))
```

#### Recurrent Neural Network (RNN)

```
>>> from keras.layers import Embedding, LSTM
>>> model3.add(Embedding(20000,128))
>>> model3.add(LSTM(128, dropout=0.2, recurrent_dropout=0.2))
>>> model3.add(Dense(1, activation='sigmoid'))
```

#### Train and Test Sets

```
>>> from sklearn.model_selection import train_test_split
>>> X_train,X_test,y_train,y_test = train_test_split(X,
                                                    y,
                                                    test_size=0.33,
                                                    random_state=42)
```

#### Standardization/Normalization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(x_train2)
>>> standardized_X = scaler.transform(x_train2)
>>> standardized_X_test = scaler.transform(x_test2)
```

(Also see [NumPy](#) & [Scikit-Learn](#))

### Inspect Model

```
>>> model.output_shape
>>> model.summary()
>>> model.get_config()
>>> model.get_weights()
```

Model output shape  
Model summary representation  
Model configuration  
List all weight tensors in the model

### Compile Model

#### MLP: Binary Classification

```
>>> model.compile(optimizer='adam',
                 loss='binary_crossentropy',
                 metrics=['accuracy'])
```

#### MLP: Multi-Class Classification

```
>>> model.compile(optimizer='rmsprop',
                 loss='categorical_crossentropy',
                 metrics=['accuracy'])
```

#### MLP: Regression

```
>>> model.compile(optimizer='rmsprop',
                 loss='mse',
                 metrics=['mse'])
```

#### Recurrent Neural Network

```
>>> model3.compile(loss='binary_crossentropy',
                  optimizer='adam',
                  metrics=['accuracy'])
```

### Model Training

```
>>> model3.fit(x_train4,
             y_train4,
             batch_size=32,
             epochs=15,
             verbose=1,
             validation_data=(x_test4,y_test4))
```

### Evaluate Your Model's Performance

```
>>> score = model3.evaluate(x_test,
                           y_test,
                           batch_size=32)
```

### Prediction

```
>>> model3.predict(x_test4, batch_size=32)
>>> model3.predict_classes(x_test4, batch_size=32)
```

### Save/Reload Models

```
>>> from keras.models import load_model
>>> model3.save('model_cifar10.h5')
>>> my_model = load_model('my_model.h5')
```

### Model Fine-tuning

#### Optimization Parameters

```
>>> from keras.optimizers import RMSprop
>>> opt = RMSprop(lr=0.0001, decay=0.1)
>>> model2.compile(optimizer='categorical_crossentropy',
                  optimizer=opt,
                  metrics=['accuracy'])
```

#### Early Stopping

```
>>> from keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=20)
>>> model3.fit(x_train4,
             y_train4,
             batch_size=32,
             epochs=15,
             validation_data=(x_test4,y_test4),
             callbacks=[early_stopping_monitor])
```

