# **R For Data Science** *Cheat Sheet* Tidyverse for Beginners

## **Tidyverse**

The **tidyverse** is a powerful collection of R packages that are actually data tools for transforming and visualizing data. All packages of the tidyverse share an underlying philosophy and common APIs.

The core packages are:



• ggplot2, which implements the grammar of graphics. You can use it to visualize your data.



 dplyr is a grammar of data manipulation. You can use it to solve the most common data manipulation challenges.



 tidyr helps you to create tidy data or data where each variable is in a column, each observation is a row end each value is a cell.



• readr is a fast and friendly way to read rectangular data.



purrr enhances R's functional programming (FP) toolkit by providing a complete and consistent set of tools for working with functions and vectors.



• tibble is a modern re-imaginging of the data frame.



• stringr provides a cohesive set of functions designed to make working with strings as easy as posssible



• forcats provide a suite of useful tools that solve common problems with factors.

You can install the complete tidyverse with:

> install.packages("tidyverse")

Then, load the core tidyverse and make it available in your current R session by running:

> library(tidyverse)

Note: there are many other tidyverse packages with more specialised usage. They are not loaded automatically with library(tidyverse), so you'll need to load each one with its own call to library().

## **Useful Functions**

> t	tidyverse_conflicts()	Conflicts between tidyverse and other packages
	cidyverse_deps()	List all tidyverse dependencies
> t	cidyverse_logo()	Get tidyverse logo, using ASCII or unicode characters
> t > t	tidyverse_packages() tidyverse_update()	List all tidyverse packages Update tidyverse packages

## Loading in the data

library(datasets)	Load the datasets package
library(gapminder)	Load the gapminder package
attach(iris)	Attach iris data to the R search pat

# dplyr

#### Filter

filter() allows you to select a subset of rows in a data frame.

#### Arrange

arrange () sorts the observations in a dataset in ascending or descending order based on one of its variables.

```
> iris %>% Sort in sepal le se
```

Sort in ascending order of sepal length Sort in descending order of sepal length

Combine multiple dplyr verbs in a row with the pipe operator %>%:

```
> iris %>%
filter(Species=="virginica") %>%
filter for species "virginica"
then arrange in descending
order of sepal length
```

#### Mutate

mutate () allows you to update or create new columns of a data frame.

```
> iris %>%
    mutate(Sepal.Length=Sepal.Length*10)
> iris %>%
    mutate(SLMm=Sepal.Length*10)
Change Sepal.Length to be in millimeters
Create a new column called SLMm
```

Combine the verbs filter(), arrange(), and mutate():

```
> iris %>%
  filter(Species=="Virginica") %>%
  mutate(SLMm=Sepal.Length*10) %>%
  arrange(desc(SLMm))
```

#### Summarize

summarize() allows you to turn many observations into a single data point.

>	iris %>%	Summarize to find the
	<pre>summarize(medianSL=median(Sepal.Length))</pre>	median sepal length
>		Filter for virginica then
	filter(Species=="virginica") %>%	summarize the median
	<pre>summarize(medianSL=median(Sepal.Length))</pre>	sepal length

You can also summarize multiple variables at once:

 ${\tt group\_by}$  ( ) allows you to summarize within groups instead of summarizing the entire dataset:

## ggplot2

#### Scatter plot

Scatter plots allow you to compare two variables within your data. To do this with ggplot2, you use  $\texttt{geom}\ \texttt{point}$  ()

#### **Additional Aesthetics**

#### Color



#### Size



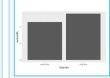
#### Faceting



### Line Plots



#### **Bar Plots**



## Histograms



#### **Box Plots**

Find median and max

