

# Conversion Guide between R and Python: Data manipulation

Afshine AMIDI and Shervine AMIDI

August 21, 2020

## Main concepts

□ **File management** – The table below summarizes the useful commands to make sure the working directory is correctly set:

Category	R Command	Python Command
Paths	<code>setwd(path)</code>	<code>os.chdir(path)</code>
	<code>getwd()</code>	<code>os.getcwd()</code>
	<code>file.path(path_1, ..., path_n)</code>	<code>os.path.join(path_1, ..., path_n)</code>
Files	<code>list.files(path, include.dirs = TRUE)</code>	<code>os.listdir(path)</code>
	<code>file_test('-f', path)</code>	<code>os.path.isfile(path)</code>
	<code>file_test('-d', path)</code>	<code>os.path.isdir(path)</code>
	<code>read.csv(path_to_csv_file)</code>	<code>pd.read_csv(path_to_csv_file)</code>
	<code>write.csv(df, path_to_csv_file)</code>	<code>df.to_csv(path_to_csv_file)</code>

□ **Exploring the data** – The table below summarizes the main functions used to get a complete overview of the data:

Category	R Command	Python Command
Look at data	<code>df %&gt;% select(col_list)</code>	<code>df[col_list]</code>
	<code>df %&gt;% head(n) / df %&gt;% tail(n)</code>	<code>df.head(n) / df.tail(n)</code>
	<code>df %&gt;% summary()</code>	<code>df.describe()</code>
Data types	<code>df %&gt;% str()</code>	<code>df.dtypes / df.info()</code>
	<code>df %&gt;% NROW() / df %&gt;% NCOL()</code>	<code>df.shape</code>

□ **Data types** – The table below sums up the main data types that can be contained in columns:

R Data type	Python Data type	Description
character	object	String-related data
factor		String-related data that can be put in bucket, or ordered
numeric	float64	Numerical data
int	int64	Numeric data that are integer
POSIXct	datetime64	Timestamps

## Data preprocessing

□ **Filtering** – We can filter rows according to some conditions as follows:

R

```
df %>%
  filter(some_col some_operation some_value_or_list_or_col)
```

where some\_operation is one of the following:

Category	R Command	Python Command
Basic	<code>== / !=</code>	<code>== / !=</code>
	<code>&lt;, &lt;=, &gt;=, &gt;</code>	<code>&lt;, &lt;=, &gt;=, &gt;</code>
	<code>&amp; /  </code>	<code>&amp; /  </code>
Advanced	<code>is.na()</code>	<code>pd.isnull()</code>
	<code>%in% (val_1, ..., val_n)</code> <code>%like% 'val'</code>	<code>.isin([val_1, ..., val_n])</code> <code>.str.contains('val')</code>

□ **Mathematical operations** – The table below sums up the main mathematical operations that can be performed on columns:

Operation	R Command	Python Command
$\sqrt{x}$	<code>sqrt(x)</code>	<code>np.sqrt(x)</code>
$\lfloor x \rfloor$	<code>floor(x)</code>	<code>np.floor(x)</code>
$\lceil x \rceil$	<code>ceiling(x)</code>	<code>np.ceil(x)</code>

## Data frame transformation

□ **Common transformations** – The common data frame transformations are summarized in the table below:

Category	R Command	Python Command
Concatenation	<code>rbind(df_1, ..., df_n)</code> <code>cbind(df_1, ..., df_n)</code>	<code>pd.concat([df_1, ..., df_n], axis=0)</code> <code>pd.concat([df_1, ..., df_n], axis=1)</code>
Dimension change	<code>spread(df, key, value)</code>  <code>gather(df, key, value)</code>	<pre>pd.pivot_table(     df, values='some_values',     index='some_index',     columns='some_column',     aggfunc=np.sum )  pd.melt(     df, id_vars='variable',     value_vars='other_variable' )</pre>