# Data Manupulate Cheat Sheet -- R vs Python

| Global Setting | | Load & Write Data | |
|---|---|---|---|
| install.packages() | import() | read.csv() | pd.read_csv() |
| getwd() | os.getcwd() | write.csv(, row.name=F) | df.to_csv(, index=False) |
| setwd() | os.chdir() | read.table() | pd.read.table() |
| set.seed() | np.random.seed() | write.xlsx() | df.to_excel() |
| ls() | os.listdir() | **Data Type** | |
| rm() | os.remove() | as.numeric() | as.numeric() |
| **Data Slicing** | | as.character() | as.character() |
| df[1:10, ] | df.iloc[0:10, ] | as.factor() | as.factor() |
| df[, 1:3] | df.iloc[:, 1:3 ] | as.data.frame() | as.data.frame() |
| df[, col] | df.loc[:, col] | as.Date() | as.Date() |
| df$col | df.col | **Basic Function** | |
| df[df$col in c(), ] | df.loc[ df.col.isin([]) ] | seq() | range() |
| df[df$col == value ] | df.loc[ df.col == value ] | rep() | np.repeat() |
| df[df$col == value, col] = value | df.loc[df.col == value, ] = value | length() | len() |
| df[, -1] | df.drop(df.columns[1], 1) | table() | pd.crosstab() |
| df[, c(col1, col2)] = NULL | df.drop([col1, col2], 1) | unique() | set() |
| **Data Wrangling** | | class() | type() |
| df[order(col), ] | df.sort_values([col]) | strsplit() | str.split() |
| df[order(-col), ] | df.sort_values([col], ascending=[0]) | is.null() | pd.isnull() |
| colnames() | df.columns.values() | is.na() | pd.isnan() |
| row.names() | df.index() | paste(A, B, sep='_') | %s_%s' %( A, B) |
| rbind() | pd.concat([, ], axis=0) | setdiff(A, B) | [x for x in A if x not in B] |
| cbind() | pd.concat([, ], axis=1) | grep() | re.findall() |
| str() | df.info() | **Basic Statistics** | |
| dim() | df.shape() | mean() | np.mean() |
| head() | df.head() | rowSums() | df.sum(1) |
| summary() | df.describe() | colSums() | df.sum(0) |
| merge(, by) | pd.merge(, on) | colMeans() | df.mean(0) |
| merge(, all.x=T) | pd.merge(, how='left') | rowMeans() | df.mean(1) |
| merge(, all.y=T) | pd.merge(, how='right') | var() | np.var() |
| merge(, all.x=T, all.y=T) | pd.merge(, how='outer') | sd() | np.std() |
| apply(, MARGIN = 1, function(x) ) | df.apply( lambda x: , axis=1) | t.test() | sp.stats.ttest_ind() |
| apply(, MARGIN = 2, function(x) ) | df.apply( lambda x: , axis=0) | rnorm() | sp.stats.norm.rvs() |
| tidyr::drop_na(df) | df.dropna() | dnorm() | sp.stats.norm.pdf() |
| df[is.na(df)] = value | df.fillna(value) | lm() | linear_model.LinearRegression() |
| df[, mean(col3), .(col1, col2)] | df.groupby([col1, col2])[col3].mean() | | |
| df[, .N, .(col1, col2)] | df.groupby([col1, col2).size() | | |
| unique() | df.drop_duplicates() | | |
| sample_n() | df.sample() | | |
| dcast(, ~, value.var, fun ) | pd.pivot_table(, index, columns, values, aggfunc) | | |
| melt(, id.vars = c(col)) | pd.melt(, id_vars) | | |