# Python ↔ R Data Wrangling Cheat Sheet (pandas ↔ dplyr/tidyr)

Assume: R uses dplyr/tidyr; Python uses pandas. Snippets are minimal + copy-paste friendly.

## Setup (tiny example data)

| R | Python |
|---|---|
| ```
library(dplyr)
library(tidyr)

df <- tibble(
  name = c("A","B","C","C"),
  age  = c(25, 30, 22, 22),
  score= c(88, 92, 95, 95),
  grp  = c("g1","g1","g2","g2")
)
``` | ```
import pandas as pd

df = pd.DataFrame({
  "name": ["A","B","C","C"],
  "age":  [25, 30, 22, 22],
  "score":[88, 92, 95, 95],
  "grp":  ["g1","g1","g2","g2"]
})
``` |

## Select columns

| R | Python |
|---|---|
| ```
df %>% select(name, score)
df %>% select(starts_with("s"))
``` | ```
df[["name", "score"]]
df.loc[:, df.columns.str.startswith("s")]
``` |

## Filter rows

| R | Python |
|---|---|
| `df %>% filter(age > 24, grp == "g1")` | `df[(df["age"] > 24) & (df["grp"] == "g1")]` |

## Create/transform columns (mutate)

| R | Python |
|---|---|
| ```
df %>%
  mutate(
    passed = score > 90,
    score_z = (score - mean(score)) / sd(score))
  )
``` | ```
df.assign(
  passed = df["score"] > 90,
  score_z = (df["score"] - df["score"].mean()) / df["
``` |

## Group + summarize (aggregation)

| R | Python |
|---|---|
| ```
df %>%
  group_by(grp) %>%
  summarise(
    n = n(),
    mean_score = mean(score),
    max_age = max(age),
    .groups = "drop"
  )
``` | ```
df.groupby("grp").agg(
  n=("grp", "size"),
  mean_score=("score", "mean"),
  max_age=("age", "max")
).reset_index()
``` |

## Sort (arrange)

| R | Python |
|---|---|
| `df %>% arrange(desc(score), age)` | `df.sort_values(["score", "age"], ascending=[False, Tr` |

## Rename columns

| R | Python |
|---|---|
| df %>% rename(points = score) | df.rename(columns={"score": "points"}) |

## Distinct / drop duplicates

| R | Python |
|---|---|
| df %>% distinct()<br>df %>% distinct(name, .keep_all = TRUE) | df.drop_duplicates()<br>df.drop_duplicates(subset=["name"], keep="first") |

## Join / merge

| R | Python |
|---|---|
| lookup <- tibble(grp=c("g1","g2"), label=c("Group...<br><br>df %>% left_join(lookup, by = "grp") | lookup = pd.DataFrame({"grp":["g1","g2"], "label":["G...<br><br>df.merge(lookup, on="grp", how="left") |

## Reshape: pivot wider

| R | Python |
|---|---|
| # example: mean score by grp & name, wide by grp<br>df %>%<br>  group_by(name, grp) %>% summarise(mean_score=mean(score), .groups="drop") %>%<br>  pivot_wider(names_from = grp, values_from = mean_score) | tmp = (df.groupby(["name","grp"], as_index=False)<br>                   .agg(mean_score=("score","mean")))<br>tmp.pivot(index="name", columns="grp", values="mean_s... |

## Reshape: pivot longer (wide → long)

| R | Python |
|---|---|
| wide <- tibble(name=c("A","B"), g1=c(1,2), g2=c(3,4))<br><br>wide %>%<br>  pivot_longer(cols = c(g1, g2), names_to="grp", values_to... | wide = pd.DataFrame({"name":["A","B"], "g1":[1,2], "g...<br><br>wide.melt(id_vars=["name"], value_vars=["g1","g2"],<br>  var_name="grp", value_name="value") |

## Handle missing values

| R | Python |
|---|---|
| df2 <- tibble(x=c(1, NA, 3), y=c("a", NA, "c"))<br><br>df2 %>% drop_na()<br>df2 %>% replace_na(list(x = 0, y = "missing")) | import numpy as np<br>df2 = pd.DataFrame({"x":[1, np.nan, 3], "y":["a", np....<br><br>df2.dropna()<br>df2.fillna({"x": 0, "y": "missing"}) |

## Split-apply-combine (within groups)

| R | Python |
|---|---|
| df %>%<br>  group_by(grp) %>%<br>  mutate(rank_in_grp = dense_rank(desc(score))) %>%<br>  ungroup() | df.assign(<br>  rank_in_grp = df.groupby("grp")["score"].rank(metho... |