

Python for R Users - Cheat Sheet

Quick correspondences for common data, wrangling, stats, viz, and modeling tasks. Updated Feb 18, 2026.

1) Core data structures

Concept	Python	R
Numeric vector	<code>list, numpy.array</code>	<code>c()</code>
Matrix	<code>numpy.array</code> (2D)	<code>matrix()</code>
Data frame	<code>pandas.DataFrame</code>	<code>data.frame() / tibble()</code>
1D labeled	<code>pandas.Series</code>	named vector
Categorical	<code>pandas.Categorical</code>	<code>factor()</code>
Missing value	<code>numpy.nan / pd.NA</code>	<code>NA</code>

2) Data wrangling: pandas \leftrightarrow dplyr

Task	Python (pandas)	R (dplyr/tidyr)
Select columns	<code>df[['a','b']]</code>	<code>select(df, a, b)</code>
Filter rows	<code>df[df.a > 5]</code>	<code>filter(df, a > 5)</code>
Add/transform column	<code>df['c'] = ...</code>	<code>mutate(df, c = ...)</code>
Group + summarize	<code>df.groupby('a').agg(...)</code>	<code>df %>% group_by(a) %>% summarise(...)</code>
Sort	<code>df.sort_values('a')</code>	<code>arrange(df, a)</code>
Join	<code>pd.merge(x, y, on='id', how='left')</code>	<code>left_join(x, y, by = 'id')</code>
Distinct	<code>df.drop_duplicates()</code>	<code>distinct(df)</code>
Rename	<code>df.rename(columns={'old':'new'})</code>	<code>rename(df, new = old)</code>
Pivot wider	<code>df.pivot(index=..., columns=..., values=...)</code>	<code>pivot_wider(df, names_from=..., values_from=...)</code>
Pivot longer	<code>pd.melt(df, id_vars=[...])</code>	<code>pivot_longer(df, cols=..., names_to=..., values_to=...)</code>

3) Stats: common tests & summaries

Thing	Python	R
Mean / SD	<code>np.mean(x) / np.std(x, ddof=1)</code>	<code>mean(x) / sd(x)</code>
Median / quantile	<code>np.median(x) / np.quantile(x, q)</code>	<code>median(x) / quantile(x, probs=q)</code>
Correlation	<code>scipy.stats.pearsonr(x,y)</code>	<code>cor(x,y)</code>
t-test	<code>scipy.stats.ttest_ind(x,y)</code>	<code>t.test(x,y)</code>
Chi-square	<code>scipy.stats.chi2_contingency(tbl)</code>	<code>chisq.test(tbl)</code>
Linear regression	<code>statsmodels.formula.api.ols('y ~ x1 + x2', data=df).fit()</code>	<code>lm(y ~ x1 + x2, data=df)</code>
Logistic regression	<code>statsmodels.formula.api.logit('y ~ x', data=df).fit()</code>	<code>glm(y ~ x, data=df, family=binomial())</code>

4) Visualization: matplotlib/seaborn ↔ ggplot2

Plot	Python	R
Scatter	sns.scatterplot(data=df, x='x', y='y')	ggplot(df, aes(x, y)) + geom_point()
Line	plt.plot(x, y)	ggplot(df, aes(x, y)) + geom_line()
Histogram	plt.hist(x, bins=30)	ggplot(df, aes(x)) + geom_histogram(bins=30)
Boxplot	sns.boxplot(data=df, x='g', y='y')	ggplot(df, aes(g, y)) + geom_boxplot()

5) Machine learning: scikit-learn ↔ caret/tidymodels

Task	Python	R
Train/test split	from sklearn.model_selection import train_test_split	rsample::initial_split() / caret::createDataPartition()
Fit model	model.fit(X_train, y_train)	fit(model, data=train) / train(...)
Predict	model.predict(X_test)	predict(model, newdata=test)
Random forest	RandomForestClassifier / RandomForestRegressor	ranger / randomForest
Cross-validation	cross_val_score(model, X, y, cv=5)	vfold_cv(data, v=5)

6) File I/O

Task	Python	R
Read CSV	pd.read_csv('file.csv')	read.csv('file.csv') / readr::read_csv()
Write CSV	df.to_csv('file.csv', index=False)	write.csv(df, 'file.csv', row.names=FALSE)
Read Excel	pd.read_excel('file.xlsx')	readxl::read_excel('file.xlsx')
Read JSON	json.load(open('file.json'))	jsonlite::fromJSON('file.json')
Save object	pickle.dump(obj, open('x.pkl', 'wb'))	saveRDS(obj, 'x.rds')
Load object	pickle.load(open('x.pkl', 'rb'))	readRDS('x.rds')

Handy mental mapping

- %>% pipe ≈ method chaining (`df.assign(...).query(...)`) or `.pipe()`
- **tibble** ≈ **DataFrame** (print-friendly, column-first)
- **formula** ($y \sim x$) ≈ statsmodels formula strings
- **NA** handling: in pandas watch out for **NaN** vs **pd.NA** and **dtype** casting

Tip: If you tell me what packages you use most in R (tidyverse, data.table, base, etc.), I can tailor a version that matches your exact workflow.