# Common statistical tests are linear models

*Last updated: 27 June, 2019*

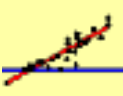See worked examples and more details at the accompanying notebook: https://github.com/eigenfoo/tests-as-linear

| | Common name | Function in scipy.stats | Equivalent linear model in smf.ols | Exact? | The linear model in words | Icon |
|---|---|---|---|---|---|---|
| **Simple Regression: (y ~ 1 + x)** | **y is independent of x**<br>P: One-sample t-test<br>N: Wilcoxon signed-rank | scipy.stats.ttest_1samp(y)<br>scipy.stats.wilcoxon(y) | smf.ols("y ~ 1", data)<br><br>smf.ols("y ~ 1", signed_rank(data)) | ✓<br>for N >14 | One number (intercept, i.e., the mean) predicts **y**.<br> - (Same, but it predicts the *signed rank* of **y**.) |  |
| | P: Paired-sample t-test<br>N: Wilcoxon matched pairs | scipy.stats.ttest_rel(y1, y2)<br>scipy.stats.wilcoxon(y1, y2) | smf.ols("y2_sub_y1 ~ 1", data)<br>smf.ols("y2_sub_y1 ~ 1",<br>signed_rank(data)) | ✓<br>for N >14 | One intercept predicts the pairwise **y₂-y₁** differences.<br> - (Same, but it predicts the *signed rank* of **y₂-y₁**.) |  |
| | **y ~ continuous x**<br>P: Pearson correlation<br>N: Spearman correlation | scipy.stats.pearsonr(x, y)<br>scipy.stats.spearmanr(x, y) | smf.ols("y ~ 1 + x", data)<br>smf.ols("y ~ 1 + x", rank(data)) | ✓<br>for N >10 | One intercept plus **x** multiplied by a number (slope) predicts **y**.<br> - (Same, but with *ranked* **x** and **y**) |  |
| | **y ~ discrete x**<br>P: Two-sample t-test<br>P: Welch's t-test<br>N: Mann-Whitney U | scipy.stats.ttest_ind(y1, y2)<br>N/A<br>scipy.stats.mannwhitneyu(y1, y1) | smf.ols("y ~ 1 + group", data)$^A$<br>N/A<br>smf.ols("y ~ 1 + group",<br>signed_rank(data))$^A$ | ✓<br>✓<br>for N >11 | An intercept for **group 1** (plus a difference if **group 2**) predicts **y**.<br> - (Same, but with one variance *per group* instead of one common.)<br> - (Same, but it predicts the *signed rank* of **y**.) |  |
| **Multiple regression: (y ~ 1 + x₁ + x₂ +...)** | P: One-way ANOVA<br>N: Kruskal-Wallis | scipy.stats.f_oneway(a, b, c)<br>scipy.stats.kruskal(a, b, c) | smf.ols(y ~ 1 + G₂ + G₃ +...+ G_N)$^A$<br>smf.ols(rank(y) ~ 1 + G₂ + G₃ +...+ G_N)$^A$ | ✓<br>for N >11 | An intercept for **group 1** (plus a difference if group ≠ 1) predicts **y**.<br> - (Same, but it predicts the *rank* of **y**.) |  |
| | P: One-way ANCOVA | N/A | smf.ols("y ~ 1 + G₂ + G₃ +...+ G_N + x",<br>data)$^A$ | ✓ | - (Same, but plus a slope on **x**.)<br> *Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.* |  |
| | P: Two-way ANOVA | N/A | smf.ols("y ~ 1 + G₂ + G₃ + ... + G_N +<br>S₂ + S₃ + ... + S_K +<br>G₂*S₂+G₃*S₃+...+G_N*S_K", data) | ✓ | Interaction term: changing **sex** changes the **y ~ group** parameters.<br>*Note: G₂ to N is an indicator (0 or 1) for each non-intercept levels of the **group** variable. Similarly for S₂ to K for sex. The first line (with Gᵢ) is main effect of group, the second (with Sᵢ) for sex and the third is the **group × sex** interaction. For two levels (e.g. male/female), line 2 would just be "S₂" and line 3 would be S₂ multiplied with each Gᵢ.* | [Coming] |
| | **Counts ~ discrete x**<br>N: Chi-square test | scipy.stats.chisquare(data) | **Equivalent log-linear model**<br>sm.GLM(y ~ 1 + G₂ + G₃ + ... + G_N +<br>S₂ + S₃ + ... + S_K +<br>G₂*S₂+G₃*S₃+...+G_N*S_K, family=...)$^A$ | ✓ | Interaction term: (Same as Two-way ANOVA.)<br>*Note: Run glm using the following arguments: glm(model, family=poisson())*<br>*As linear-model, the Chi-square test is log(yᵢ) = log(N) + log(αᵢ) + log(βⱼ) + log(αᵢβⱼ) where αᵢ and βⱼ are proportions. See more info in the accompanying notebook.* | *Same as Two-way ANOVA* |
| | N: Goodness of fit | scipy.stats.chi2_contingency(data) | sm.GLM(y ~ 1 + G₂ + G₃ +...+ G_N,<br>family=...)$^A$ | ✓ | (Same as One-way ANOVA and see Chi-Square note.) | *1W-ANOVA* |

List of common parametric (P) non-parametric (N) tests and equivalent linear models. The notation y ~ 1 + x is R shorthand for y = 1·b + a·x which most of us learned in school. Models in similar colors are highly similar, but really, notice how similar they *all* are across colors! For non-parametric models, the linear models are reasonable approximations for non-small sample sizes (see "Exact" column and click links to see simulations). Other less accurate approximations exist, e.g., Wilcoxon for the sign test and Goodness-of-fit for the binomial test. The signed rank function is `signed_rank(df) = np.sign(df) * df.rank()`. The variables Gᵢ and Sᵢ are "dummy coded" indicator variables (either 0 or 1) exploiting the fact that when Δx = 1 between categories the difference equals the slope. Subscripts (e.g., G₂ or y₁) indicate different columns in data. lm requires long-format data for all non-continuous models. All of this is exposed in greater detail and worked examples at https://eigenfoo.xyz/tests-as-linear/.

$^A$ See the note to the two-way ANOVA for explanation of the notation.

Jonas Kristoffer Lindeløv, George Ho<br>https://lindeloev.net, https://eigenfoo.xyz