



Syntax for Specifying Linear Models in R

Stephen Parry and Erika Mudrak

Suppose y , x , x_0 , x_1 , x_2 , ... are numeric variables and A , B , C , ... are factors/categorical variables.

Linear Models (including regression, anova, ancova, etc..)

The following syntax works generally, including in `lm`, `glm`, `aov`, `lme4::lmer`, `lme4::glmer`, `nLme::lme`, `MASS::glmPQL`, `glmTMB::glmTMB` functions.

source: Venables, Smith and the R Core Team. An Introduction to R: Notes on R: Programming Environment for Data Analysis and Graphics Version 4.4.1 (2024-06-14), Chapter 11

<http://cran.r-project.org/doc/manuals/R-intro.pdf>

$y \sim x$ $y \sim 1 + x$	Both imply the same simple linear regression model of y on x . The first has an implicit intercept term, and the second an explicit one.
$y \sim 0 + x$ $y \sim -1 + x$ $y \sim x - 1$	Simple linear regression of y on x through the origin (that is, without an intercept term).
$\log(y) \sim x_1 + x_2$	Multiple regression of the transformed variable, $\log(y)$, on x_1 and x_2 (with an implicit intercept term).
$y \sim \text{poly}(x, 2)$ $y \sim 1 + x + I(x^2)$	Polynomial regression of y on x of degree 2. The first form uses orthogonal polynomials, and the second uses explicit powers, as basis.
$y \sim I(x/100)$	Insulate the expression. Inside parentheses all operators have their normal arithmetic meaning, rather than the formula syntax. Here, x is divided by 100 before putting in the model.
$y \sim A$	Analysis of variance model of y , with classes determined by A .
$y \sim A + x$	Analysis of covariance model of y , with classes determined by A , and with covariate x .

Cornell Statistical Consulting Unit

$y \sim A*B$ $y \sim A + B + A:B$	Two-way analysis of variance model of y on A, B and their interaction
$y \sim A/B$ $y \sim A + A:B$	Two-way analysis of variance model of y on B nested in A
$y \sim (A + B + C)^2$ $y \sim A*B*C - A:B:C$	Three factor experiment but with a model containing main effects and two factor interactions only. Both formulae specify the same model.
$y \sim (A + B + C + D \dots)^n$	All terms in parentheses together with interactions up to order n
$y \sim A * x$ $y \sim A/x$ $y \sim A/(1 + x) - 1$	Separate simple linear regression models of y on x within the levels of A, with different codings. The last form produces explicit estimates of as many different intercepts and slopes as there are levels in A.

This syntax works in the aov function only. Note that aov is designed for balanced designs, and the results can be hard to interpret without balance: beware that missing values in the response(s) will likely lose the balance.

$y \sim A*B + \text{Error}(C)$	An experiment with two treatment factors, A and B, and error strata determined by factor C. For example a split plot experiment, with whole plots (and hence also subplots), determined by factor C.
--------------------------------	--

Mixed-Effects Models (models including random effects)

With terms as above including categorical `group` and `block` variables.

The following syntax works to specify the random effects for the `lmer` and `glmer` functions in the `lme4` package, as well as for `glmmTMB::glmmTMB` functions.

From <http://bbolker.github.io/mixedmodels-misc/glmmFAQ.html#model-specification>

Also a good resource: <http://lme4.r-forge.r-project.org/book/Ch2.pdf>

<code>(1 group)</code>	Random group intercept
<code>(x group)</code> <code>(1+x group)</code>	Random slope of x within group with correlated intercept
<code>(0+x group)</code> <code>(-1+x group)</code>	Random slope of x within group: no variation in intercept
<code>(1 group) + (0+x group)</code>	Uncorrelated random intercept and random slope within group
<code>(1 site/block)</code> <code>(1 site) + (1 site:block)</code>	Intercept varying among sites and among blocks within sites (nested random effects)
<code>(x site/block)</code> <code>(x site) + (x site:block)</code>	Slope and intercept varying among sites and among blocks within sites
<code>(1 group1) + (1 group2)</code>	intercept varying among crossed random effects (e.g. site, year)
<code>(x1 site) + (x2 block)</code>	two different effects, varying at different levels
<code>x*site + (x site:block)</code>	fixed effect variation of slope and intercept varying among sites and random variation of slope and intercept among blocks within sites

The following syntax specifies the random effects in the model for the `nlme::lme` (normal distribution only) and `MASS::glmmPQL` (for probability distributions in the exponential family).

<code>random = ~1 group</code>	Random group intercept
<code>random = ~x group</code>	Random slope of x within group with correlated intercept
<code>random = ~1 site/block</code>	Intercept varying among sites and among blocks within sites (nested random effects)
<code>random = ~x site/block</code>	Slope and intercept varying among sites and among blocks within sites