

Clinical Regression in R

Choose the model, encode the estimand, diagnose the fit, report the effect

Rverse Analytics

The outcome distribution chooses the model; the scientific question chooses the contrast. Write the estimand before fitting anything: adjusted mean difference, odds ratio, rate ratio or hazard ratio?

Model map

Outcome	Starting model	Core R call	Effect scale
Continuous	Gaussian linear	<code>lm(y ~ x + z, data = d)</code>	adjusted mean difference
Binary	Logistic	<code>glm(y ~ x + z, family = binomial, data = d)</code>	OR after <code>exp(coef)</code>
Common binary risk	Modified Poisson + robust SE	<code>glm(y ~ x + z, family = poisson("log"), data = d)</code>	risk ratio
Count / person-time	Poisson	<code>glm(events ~ x + z, family = poisson, data = d)</code> <code>offset(log(time))</code>	incidence-rate ratio
Overdispersed count	Negative binomial	<code>MASS::glm.nb(events ~ x + z, family = poisson("log"), data = d)</code> <code>offset(log(time))</code>	incidence-rate ratio
Ordered category	Proportional odds	<code>ordinal::clm(y ~ x + z, data = d)</code>	common OR
Repeated / clustered	Mixed or GEE	<code>lme4::glmer(...)</code> <code>geepack::geeglm(...)</code>	/ conditional / marginal
Time-to-event	Cox PH	<code>survival::coxph(Surv(t, event) ~ x + z, data = d)</code>	hazard ratio

A reproducible fit-to-table pattern

```
library(dplyr)
library(broom)
library(gtsummary)

d <- trial |>
  mutate(
    response = factor(response, levels = c(0, 1), labels = c("No", "Yes")),
    stage = relevel(stage, ref = "T1")
  )

fit <- glm(response ~ scale(age) + stage + grade,
           family = binomial, data = d, na.action = na.exclude)

# Machine-readable estimates
tidy(fit, exponentiate = TRUE, conf.int = TRUE)

# Publication table
tbl_regression(fit, exponentiate = TRUE,
              label = list(`scale(age)` ~ "Age, per 1 SD")) |>
  add_global_p() |>
  add_n() |>
  bold_labels()
```

Encode predictors deliberately

```
# Reference category
d$group <- relevel(factor(d$group), ref = "Control")

# Non-linearity: restricted cubic spline
fit_spline <- glm(y ~ splines::ns(age, df = 4) + sex,
                 family = binomial, data = d)

# Prespecified interaction: effect of treatment differs by sex
fit_int <- glm(y ~ treatment * sex + age,
              family = binomial, data = d)

# Likelihood-ratio test for the interaction block
anova(update(fit_int, . ~ . - treatment:sex), fit_int, test = "LRT")
```

Interpretation: With an interaction, the main treatment coefficient is the treatment effect at the reference level of `sex`; it is not the overall treatment effect.

Diagnostics by model

Question	Useful check	R recipe	Red flag / response
Linearity	residuals vs fitted	<code>plot(fit, which = 1)</code>	curve → transform/spline
Constant variance	scale-location	<code>plot(fit, which = 3)</code>	fan shape → robust SE/model variance
Influential cases	Cook's distance	<code>plot(fit, which = 4)</code>	investigate data and sensitivity
Collinearity	VIF	<code>performance::check_collinearity(fit)</code>	unstable CI → revise predictors
Logistic calibration	observed vs predicted	<code>performance::check_model(fit)</code>	mis-calibration → respecify/validate
Discrimination	ROC/AUC	<code>pROC::roc(d\$y, fitted(fit))</code>	report CI; AUC is not calibration
Count dispersion	Pearson χ^2 / df	<code>sum(residuals(fit, "pearson")^2) / df.residual(fit)</code>	>1 → NB or robust approach

Robust uncertainty and model comparison

```
# Heteroskedasticity-consistent covariance
V <- sandwich::vcovHC(fit, type = "HC3")
lmtest::coefTest(fit, vcov. = V)

# Nested models: likelihood-ratio test
anova(fit_small, fit_large, test = "LRT")

# Predict on the response scale with uncertainty
pred <- predict(fit, newdata = new_patients, type = "link", se.fit = TRUE)
new_patients |>
  mutate(prob = plogis(pred$fit),
         lo = plogis(pred$fit - 1.96 * pred$se.fit),
         hi = plogis(pred$fit + 1.96 * pred$se.fit))
```

Report, do not merely print

- State population, outcome coding, predictors, functional forms and missing-data strategy.
- Report effect estimate + 95% CI + exact *p*; include units and reference categories.
- Separate prespecified adjustment from data-driven selection; avoid univariable *p*-value screening.
- Give events and candidate parameter counts; quantify optimism with internal validation when predicting.
- For prediction models, report calibration and discrimination; a significant coefficient is not validation.