

# reQTL Model Selection Using QQplot

Genevieve Roberts

2025-07-24

## Load Packages

```
library(lme4)
library(lmerTest)
library(emmeans)
library(data.table)
library(here)
library(dplyr)
library(stringr)
library(tidyr)
library(tibble)
library(purrr)
library(ggplot2)
library(broom)
library(broom.mixed)
```

## Explore the real data

This is a presumably “real” reQTL SNP in melanocytes. The genotype is rs2910686 and the gene expression is ERAP2. It’s presumably real because it is also a known psoriasis GWAS SNP and in high LD to a causal SNP validated in vitro.

```
#load the real data
real_dat <- read.csv(here::here("notebooks/long_form_reQTL_data.csv")) %>%
  arrange(donor)

#make sure the reference group is set to "PBS"
real_dat$condition <- relevel(factor(real_dat$condition), ref = "PBS")
real_dat$genotype.nt <- relevel(factor(real_dat$genotype.nt), ref = "TT")
```

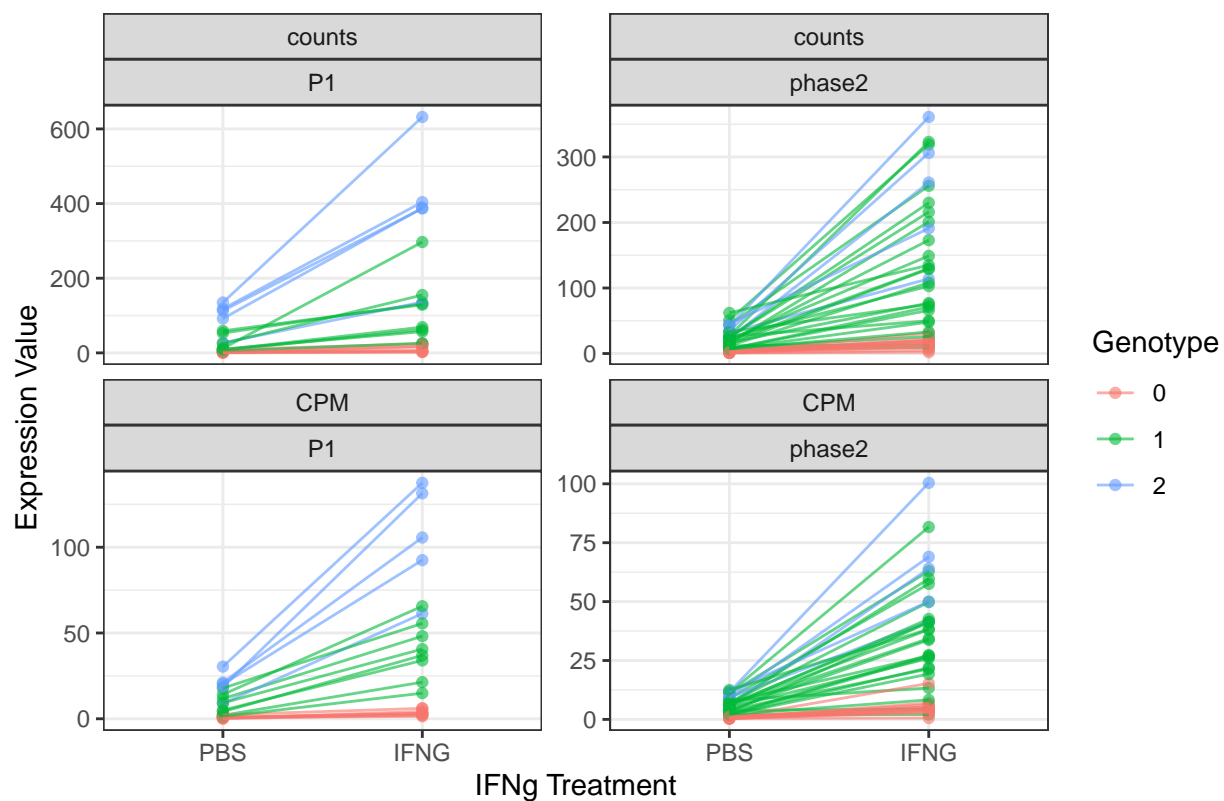
## Plot the real data

```
# Create the plot

# Reshape to long format for faceting
real_dat_plot <- real_dat %>%
  pivot_longer(cols = c(counts, CPM), names_to = "scale_type", values_to = "expression_value")

# Create the plot
ggplot(real_dat_plot, aes(x = factor(condition),
                          y = expression_value,
                          color = factor(genotype.num))) +
  geom_point(alpha = 0.6) +
  geom_line(aes(group = donor), alpha = 0.6) +
  labs(
    title = "Interaction between Genotype and IFNg_treatment on Expression",
    x = "IFNg Treatment",
    y = "Expression Value",
    color = "Genotype"
  ) +
  facet_wrap(scale_type ~ phase, scales = "free_y") +
  theme_bw()
```

## Interaction between Genotype and IFNg\_treatment on Expression



## Compare various model fits

### Linear mixed effects (CPM) with random intercept for donor

```
# Fit the Poisson regression model
real_linear_model <- lmer(CPM ~ condition * genotype.num +
  (1 | donor), #random intercept for donor
  data = real_dat)

# Clean up
real_tidy_linear <- tidy(real_linear_model)
pander(real_tidy_linear)
```

Table 1: Table continues below

effect	group	term	estimate	std.error
fixed	NA	(Intercept)	0.1691	2.876
fixed	NA	conditionIFNG	0.8466	3.287
fixed	NA	genotype.num	6.747	2.656
fixed	NA	conditionIFNG:genotype.num	32	3.047
ran_pars	donor	sd__(Intercept)	8.137	NA
ran_pars	Residual	sd__Observation	11.47	NA

statistic	df	p.value
0.05879	102.7	0.9532
0.2576	59.39	0.7976
2.54	103.7	0.01255
10.5	59.39	3.649e-15
NA	NA	NA
NA	NA	NA

### Linear mixed effects (CPM) with random intercept for donor nested within phase

```
# Fit the Poisson regression model
real_linear_model <- lmer(CPM ~ condition * genotype.num +
  (1 | phase/donor), #random intercept for donor, nested within phase
  data = real_dat)

# Clean up and pull the interaction term
real_tidy_linear <- tidy(real_linear_model)
pander(real_tidy_linear)
```

Table 3: Table continues below

effect	group	term	estimate	std.error
fixed	NA	(Intercept)	1.32	4.402
fixed	NA	conditionIFNG	0.8466	3.289
fixed	NA	genotype.num	6.569	2.584
fixed	NA	conditionIFNG:genotype.num	32	3.048
ran_pars	donor:phase	sd__(Intercept)	7.454	NA
ran_pars	phase	sd__(Intercept)	4.754	NA

effect	group	term	estimate	std.error
ran_pars	Residual	sd__Observation	11.48	NA

statistic	df	p.value
0.2998	2.026	0.7923
0.2574	59.22	0.7977
2.542	104.6	0.01249
10.5	59.22	3.86e-15
NA	NA	NA
NA	NA	NA
NA	NA	NA

### Negative binomial mixed effects (counts) with random intercept for donor

```
# Fit the Poisson regression model
real_nb_model <- glmer.nb(counts ~ condition * genotype.num +
  (1 | donor),
  data = real_dat)

# Clean up and pull the interaction term
real_tidy_nb <- tidy(real_nb_model)
pander(real_tidy_nb)
```

Table 5: Table continues below

effect	group	term	estimate	std.error
fixed	NA	(Intercept)	0.9906	0.1913
fixed	NA	conditionIFNG	1.663	0.1615
fixed	NA	genotype.num	1.676	0.1635
fixed	NA	conditionIFNG:genotype.num	0.01622	0.1341
ran_pars	donor	sd__(Intercept)	0.6431	NA

statistic	p.value
5.178	2.246e-07
10.3	6.949e-25
10.25	1.23e-24
0.1209	0.9038
NA	NA

### LRT comparing nested models

```
# Fit a reduced model
real_poisson_model_red <- glmer.nb(counts ~ condition + genotype.nt + (1 | donor), #no interaction term
  data = real_dat)

# Fit the Poisson regression model
real_poisson_model <- glmer.nb(counts ~ condition * genotype.nt + (1 | donor), #interaction term)
```

```

data = real_dat)

# Do a likelihood ratio test
tidy_lrt <- tidy(anova(real_poisson_model_red, real_poisson_model, test = "LRT"))
tidy_lrt

## # A tibble: 2 x 9
##   term                npar  AIC   BIC logLik minus2logL statistic    df p.value
##   <chr>              <dbl> <dbl> <dbl> <dbl>      <dbl>    <dbl> <dbl> <dbl>
## 1 real_poisson_mode~    6 1043. 1060. -516.      1031.    NA      NA    NA
## 2 real_poisson_model    8 1047. 1070. -516.      1031.    0.246    2    0.884

```

## Generate QQplots for Model Candidates

```
### functions
#function to permute outcome and record results for various model formulas
generate_permuted_null <- function(model_formula,
                                   data,
                                   interaction_term_regex = ":",
                                   permute_var = NULL,
                                   n_perm = 100,
                                   model_type = c("lmer", "glmer", "glmer.nb", "glm"),
                                   family = poisson) {

  model_type <- match.arg(model_type)
  model_formula_str <- paste(deparse(model_formula), collapse = " ")

  # Determine outcome variable for permutation
  if (is.null(permute_var)) {
    outcome_var <- all.vars(model_formula)[1]
  } else {
    outcome_var <- permute_var
  }

  # Fit the model to the observed data
  fit_model <- switch(
    model_type,
    lmer      = lmer(model_formula, data = data),
    glmer     = glmer(model_formula, data = data, family = family),
    glmer.nb  = glmer.nb(model_formula, data = data),
    glm       = glm(model_formula, data = data, family = family)
  )

  # Extract observed interaction term results
  observed_results <- tidy(fit_model) %>%
    filter(str_detect(term, interaction_term_regex)) %>%
    select(term, estimate, se = std.error, statistic = any_of("statistic"), p.value) %>%
    mutate(type = "observed")

  # Generate permuted results
  permuted_results <- map_dfr(1:n_perm, function(i) {
    permuted_data <- data %>%
      mutate(!!sym(outcome_var) := sample(!!sym(outcome_var)))

    perm_model <- tryCatch({
      switch(
        model_type,
        lmer      = lmer(model_formula, data = permuted_data),
        glmer     = glmer(model_formula, data = permuted_data, family = family),
        glmer.nb  = glmer.nb(model_formula, data = permuted_data),
        glm       = glm(model_formula, data = permuted_data, family = family)
      )
    }, error = function(e) return(NULL))

    if (is.null(perm_model)) return(tibble())
  })
```

```

    tidy(perm_model) %>%
      filter(str_detect(term, interaction_term_regex)) %>%
      select(term, estimate, se = std.error, statistic = any_of("statistic"), p.value) %>%
      mutate(type = "null")
  })

  # Combine and annotate
  bind_rows(observed_results, permuted_results) %>%
    mutate(model_formula = model_formula_str)
}

#function to plot multiple model qqplots on one plot together
plot_qq_chisq <- function(sim_df) {
  model_formula <- unique(sim_df$model_formula)

  # Compute chi-squared statistics from estimates and SEs
  sim_df <- sim_df %>%
    mutate(
      chisq = (estimate / se)^2,
      p.value = pchisq(chisq, df = 1, lower.tail = FALSE)
    )

  # Get null and observed
  null_chisq <- sim_df %>% filter(type == "null") %>% pull(chisq)
  observed_chisq <- sim_df %>% filter(type == "observed") %>% pull(chisq)
  empirical_pval <- mean(null_chisq >= observed_chisq, na.rm = TRUE)

  # Prepare data for QQ plot
  qq_df <- data.frame(
    expected = sort(qchisq(ppoints(length(null_chisq)), df = 1)),
    observed = sort(null_chisq),
    type = "null"
  )

  observed_point <- data.frame(
    expected = qchisq(0.5 / length(null_chisq), df = 1),
    observed = observed_chisq,
    label = paste0("Empirical p = ", signif(empirical_pval, 3))
  )

  ggplot(qq_df, aes(x = expected, y = observed)) +
    geom_point(alpha = 0.4) +
    geom_abline(slope = 1, intercept = 0, color = "gray") +
    geom_point(data = observed_point, aes(x = expected, y = observed), color = "red", size = 3) +
    geom_text(data = observed_point, aes(label = label), vjust = -1.5, color = "red") +
    labs(x = "Expected Chi-sq (df=1)", y = "Observed Chi-sq",
         title = "QQ Plot of Permuted vs Observed Interaction Chi-Sq",
         subtitle = paste("Model:", model_formula)) +
    theme_bw()
}

plot_qq_chisq_multi <- function(sim_df) {
  # Compute chi-squared statistics

```



```

sim_df <- sim_df %>%
  mutate(
    chisq = (estimate / se)^2,
    p.value = pchisq(chisq, df = 1, lower.tail = FALSE),
    model_formula_family=str_c(model_formula, " ", paste0(family), " ", paste0(model_type))
  )

# Calculate empirical p-values per model_formula
empirical_pvals <- sim_df %>%
  group_by(model_formula_family) %>%
  summarise(
    observed_chisq = chisq[type == "observed"],
    null_chisq = list(chisq[type == "null"]),
    empirical_pval = mean(null_chisq[[1]] >= observed_chisq, na.rm = TRUE),
    .groups = "drop"
  )

# Prepare data for QQ plot (null points only)
qq_df <- sim_df %>%
  filter(type == "null") %>%
  group_by(model_formula_family) %>%
  arrange(chisq) %>%
  mutate(
    expected = qchisq(ppoints(n()), df = 1)
  ) %>%
  ungroup()

# Prepare observed points (one per model)
observed_points <- empirical_pvals %>%
  mutate(
    expected = qchisq(1 - 0.5 / lengths(null_chisq), df = 1),
    observed_pval = pchisq(observed_chisq, df = 1, lower.tail = FALSE),
    label = paste0("p = ", signif(observed_pval, 3))
  )

# Plot
ggplot(qq_df, aes(x = expected, y = chisq, color = model_formula_family)) +
  geom_point(alpha = 0.4) +
  geom_abline(slope = 1, intercept = 0, color = "gray") +
  geom_point(data = observed_points, aes(x = expected, y = observed_chisq), shape = 4, size = 4, stroke = 2) +
  geom_text(data = observed_points, aes(x = expected, y = observed_chisq, label = label), vjust = -1, size = 8) +
  labs(
    x = "Expected Chi-sq",
    y = "Observed Chi-sq",
    color = "Model Formula",
    title = "QQ Plot of Permuted vs Observed Interaction Chi-Sq (Multiple Models)"
  ) +
  theme_bw()
}

```

## Create a grid of models and plot them together

First plot is all tested models

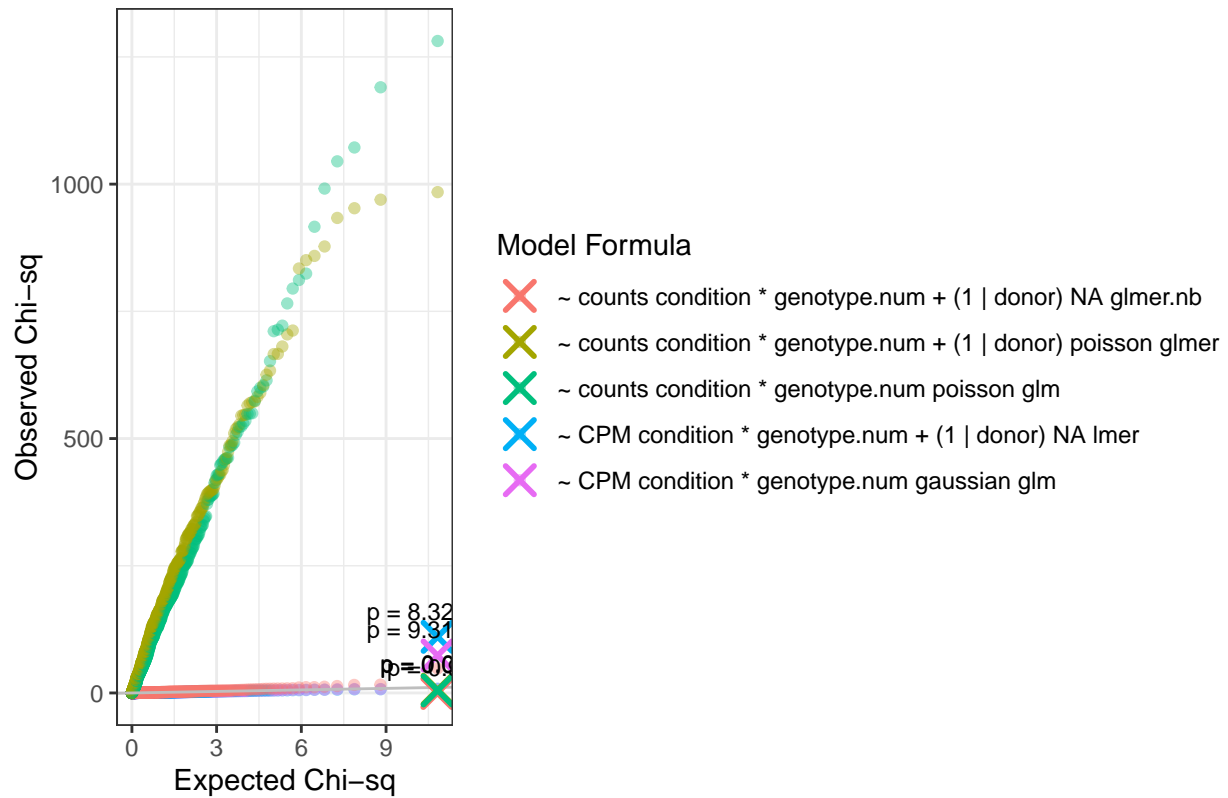
```
# Define models we'd like to try
model_grid <- tibble(
  model_formula = list(
    as.formula("CPM ~ condition * genotype.num + (1 | donor)"),
    as.formula("CPM ~ condition * genotype.num"),
    as.formula("counts ~ condition * genotype.num"),
    as.formula("counts ~ condition * genotype.num + (1 | donor)"),
    as.formula("counts ~ condition * genotype.num + (1 | donor)"),
  ),
  model_type = c("lmer", "glm", "glm", "glmer", "glmer.nb"),
  family = list(NA, gaussian(), poisson(), poisson(), NA)
)

# Run all models and return one flat dataframe
flat_qq_results <- pmap_dfr(
  list(model_grid$model_formula, model_grid$model_type, model_grid$family),
  function(model_formula, model_type, family) {
    result <- generate_permuted_null(
      model_formula = model_formula,
      data = real_dat,
      model_type = model_type,
      family = family,
      n_perm = 500
    )

    tibble(
      model_formula = paste(as.character(model_formula), collapse = " "),
      model_type = model_type,
      family = if (inherits(family, "family")) family$family else NA_character_
    ) %>%
      bind_cols(result) %>%
      select(-"model_formula...10") %>%
      rename(model_formula="model_formula...1")
  }
)

plot_qq_chisq_multi(flat_qq_results)
```

## QQ Plot of Permuted vs Observed Interaction Chi-Sq (Multiple Models)



Second plot is only the non-poisson models to “zoom in”

```
plot_qq_chisq_multi(
  filter(flat_qq_results, is.na(family) | family != "poisson")
)
```

QQ Plot of Permuted vs Observed Interaction Chi-Sq (Multiple Models)

