

PPG Paint Colors: Final Project

Example: read data, save, and reload model object

Grishma Palkar

Overview

This RMarkdown shows how to read in the final project data. It also shows how to calculate the logit-transformed response and setup the binary outcome for use with `caret` or `tidymodels`. It also demonstrates how to fit a simple model (with `lm()`), save that model, and load it back into the workspace. You may find these actions helpful as you work through the project.

You must download the data from Canvas and save the data in the same directory as this RMarkdown file.

Load packages

This example uses the `tidyverse` suite of packages.

```
library(tidyverse)

## --- Attaching packages --- tidyverse 1.3.2 ---
## ✓ ggplot2 3.4.0 ✓ purrr 1.0.1
## ✓ tibble 3.2.1 ✓ dplyr 1.1.2
## ✓ tidyr 1.3.0 ✓ stringr 1.5.0
## ✓ readr 2.1.3 ✓ forcats 0.5.2
## --- Conflicts --- tidyverse_conflicts() ---
## × dplyr::filter() masks stats::filter()
## × dplyr::lag() masks stats::lag()

library(caret)

## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
## lift

library(corrplot)

## corrplot 0.92 loaded

library(coefplot)
```

Read data

Please download the final project data from Canvas. If this Rmarkdown file is located in the same directory as the downloaded CSV file, it will be able to load in the data for you. It is **highly** recommended that you use an RStudio RProject to easily manage the working directory and file paths of the code and objects associated with the final project.

The code chunk below reads in the final project data.

```
df <- readr::read_csv("paint_project_train_data.csv", col_names = TRUE)

## Rows: 835 Columns: 8
## --- Column specification ---
## Delimiter: ","
## chr (2): Lightness, Saturation
## dbl (6): R, G, B, Hue, response, outcome
##
## I use `spec()` to retrieve the full column specification for this data.
## I Specify the column types or set `show_col_types = FALSE` to quiet this message.

The readr::read_csv() function displays the data types and column names associated with the data. However, a glimpse is shown below that reveals the number of rows and also shows some of the representative values for the columns.

df %>% glimpse()

## Rows: 835
## Columns: 8
## $ R          <dbl> 172, 26, 172, 28, 170, 175, 90, 194, 171, 122, 0, 88, 144, ...
## $ G          <dbl> 58, 88, 94, 87, 66, 89, 78, 106, 68, 151, 121, 140, 82, 163...
## $ B          <dbl> 62, 151, 58, 152, 58, 65, 136, 53, 107, 59, 88, 58, 132, 50...
## $ Lightness  <chr> "dark", "dark", "dark", "dark", "dark", "dark", "dark", "da...
## $ Saturation <chr> "bright", "bright", "bright", "bright", "bright", "bright", "br...
## $ Hue        <dbl> 4, 31, 8, 32, 5, 6, 34, 10, 1, 21, 24, 22, 36, 16, 26, 12, ...
## $ response   <dbl> 12, 10, 16, 10, 11, 16, 10, 19, 14, 25, 14, 19, 14, 38, 15...
## $ outcome    <dbl> 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1...
```

The data consist of continuous and categorical inputs. The `glimpse()` shown above reveals the data type for each variable which state to you whether the input is continuous or categorical. The RGB color model inputs, `R`, `G`, and `B` are continuous (`dbl`) inputs. The HSL color model inputs consist of 2 categorical inputs, `Lightness` and `Saturation`, and a continuous input, `Hue`. Two outputs are provided. The continuous output, `response`, and the Binary output, `outcome`. However, the data type of the Binary outcome is numeric because the Binary outcome is **encoded** as `outcome = 1` for the EVENT and `outcome = 0` for the NON-EVENT.

Binary classification task

The Binary output variable, `outcome`, is a numeric variable.

```
df %>% pull(outcome) %>% class()

## [1] "numeric"
```

However, there are **only** two unique values for `outcome`.

```
df %>% count(outcome)

## # A tibble: 2 × 2
##   outcome     n
##   <dbl> <int>
## 1         0  644
## 2         1  191
```

As stated previously, `outcome = 1` denotes the **EVENT** while `outcome = 0` denotes the **NON-EVENT**. Thus, the `outcome` variable uses the 0/1 encoding! This encoding is appropriate for `glm()` and the functions we create in homework assignments, and lecture examples. However, `caret` and `tidymodels` prefer a different encoding. For those reasons, two different binary classification data sets are defined. The first should be used for Parts iiiA) and iiiB) while the second should be used for iiiD).

The data set associated with iiiD) is created for you below. It removes the `response` variable so that way you can focus on the inputs and binary outcome.

```
dfliaA <- df %>%
  select(-response)

dfliaA %>% glimpse()

## Rows: 835
## Columns: 7
## $ R          <dbl> 172, 26, 172, 28, 170, 175, 90, 194, 171, 122, 0, 88, 144, ...
## $ G          <dbl> 58, 88, 94, 87, 66, 89, 78, 106, 68, 151, 121, 140, 82, 163...
## $ B          <dbl> 62, 151, 58, 152, 58, 65, 136, 53, 107, 59, 88, 58, 132, 50...
## $ Lightness  <chr> "dark", "dark", "dark", "dark", "dark", "dark", "dark", "da...
## $ Saturation <chr> "bright", "bright", "bright", "bright", "bright", "bright", "br...
## $ Hue        <dbl> 4, 31, 8, 32, 5, 6, 34, 10, 1, 21, 24, 22, 36, 16, 26, 12, ...
## $ outcome    <dbl> 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1...
```

The data set associated with iiiD) changes the data type of the `outcome` variable. The `ifelse()` function is used to convert `outcome` to a character data type. The value of `outcome = 1` is converted to the string 'event' and the value of `outcome = 0` is converted to 'non_event'. The `outcome` data type is then converted to a factor (R's categorical variable data type) with 'event' forced as the first level.

```
dfliaD <- df %>%
  select(-response) %>%
  mutate(outcome = ifelse(outcome == 1, 'event', 'non_event'),
         outcome = factor(outcome, levels = c('event', 'non_event')))

dfliaD %>% glimpse()

## Rows: 835
## Columns: 7
## $ R          <dbl> 172, 26, 172, 28, 170, 175, 90, 194, 171, 122, 0, 88, 144, ...
## $ G          <dbl> 58, 88, 94, 87, 66, 89, 78, 106, 68, 151, 121, 140, 82, 163...
## $ B          <dbl> 62, 151, 58, 152, 58, 65, 136, 53, 107, 59, 88, 58, 132, 50...
## $ Lightness  <chr> "dark", "dark", "dark", "dark", "dark", "dark", "dark", "da...
## $ Saturation <chr> "bright", "bright", "bright", "bright", "bright", "bright", "br...
## $ Hue        <dbl> 4, 31, 8, 32, 5, 6, 34, 10, 1, 21, 24, 22, 36, 16, 26, 12, ...
## $ outcome    <fct> event, event, event, non_event, non_event, non_event, non_e...
```

By converting `outcome` to a factor, the unique values of the variables are "always known".

```
dfliaD %>% pull(outcome) %>% levels()

## [1] "event"      "non_event"
```

However, the value counts are the same as the original encoding.

```
dfliaD %>% count(outcome)

## # A tibble: 2 × 2
##   outcome     n
##   <fct>   <int>
## 1 event    191
## 2 non_event 644
```

```
dfliaA_train <- dfliaA %>%
  select(R,G,B,Hue,outcome) %>%
  scale() %>% as.data.frame() %>%
  bind_cols(dfliaA %>% select(Lightness,Saturation))

dfliaA_train %>% glimpse()

## Rows: 835
## Columns: 7
## $ R          <dbl> -0.197990120, -2.74419521, -0.197990120, -2.70931447, -0.2327...
## $ G          <dbl> -2.3619736, -1.7631189, -1.6433480, -1.7830807, -2.2022790...
## $ B          <dbl> -1.7994266, -0.1706872, -1.8726283, -0.1523868, -1.8726283...
## $ Hue        <dbl> -1.3548215, 1.3198239, -0.9585777, 1.4188848, -1.2557605, ...
## $ outcome    <dbl> 1.8351266, 1.8351266, 1.8351266, -0.5442689, -0.5442689, ...
## $ Lightness  <chr> "dark", "dark", "dark", "dark", "dark", "dark", "dark", "da...
## $ Saturation <chr> "bright", "bright", "bright", "bright", "bright", "bright", "br...
```

(1) You must therefore train 10 different models!

```
# Model 1: Intercept-only model
mod01_glm <- glm(outcome ~ 1, data = dfliaA_train)
mod01_glm %>% readr::write_rds("model01_glm.rds")

# Model 2:Categorical variables only
mod02_glm <- glm(outcome~(Lightness + Saturation ), data = dfliaA_train)
mod02_glm %>% readr::write_rds("model02_glm.rds")

# Model 3:Continuous variables only
mod03_glm <- glm(outcome~(R+G+B+Hue), data = dfliaA_train)
mod03_glm %>% readr::write_rds("model03_glm.rds")

# Model 4:All categorical and continuous variables – linear additive
mod04_glm <- glm(outcome~(R+G+B+Hue) + Saturation + Hue), data = dfliaA_train)
mod04_glm %>% readr::write_rds("model04_glm.rds")

# Model 5:Interaction of the categorical inputs with all continuous inputs main effects
mod05_glm <- glm(outcome~(R+G+B+Hue)*(Lightness + Saturation), data = dfliaA_train)
mod05_glm %>% readr::write_rds("model05_glm.rds")

# Model 6: Add categorical inputs to all main effect and all pairwise interactions of continuous inputs
mod06_glm <- glm(outcome ~ (R+G+B+Hue)*(R+G+B+Hue)+(Lightness + Saturation ), data = dfliaA_train)
mod06_glm %>% readr::write_rds("model06_glm.rds")

# Model 7 : Interaction of the categorical inputs with all main effect and all pairwise interactions of co
ntinuous inputs
mod07_glm <- glm(outcome ~ (R+G+B+Hue)*(R+G+B+Hue)*(Lightness+Saturation) , data = dfliaA_train)
mod07_glm %>% readr::write_rds("model07_glm.rds")

# 3 models with basis functions of your choice
mod08_glm <- glm(non-linear basis functions based on your EDA.
mod08_glm <- glm(outcome ~ I(R^2) + I(G^2) + I(B^2) + I(Hue^2) + Lightness + Saturation, data = dfliaA_train)
mod08_glm %>% readr::write_rds("model08_glm.rds")

# Model 9:Can consider interactions of basis functions with other basis functions!
mod09_glm <- lm(outcome ~ (R+G+B+Lightness + Saturation + Hue)*(R+G+B+Lightness + Saturation + Hue), data = dfliaA_train)
mod09_glm %>% readr::write_rds("model09_glm.rds")

# Model 10:Can consider interactions of basis functions with the categorical inputs!
mod10_glm <- glm(outcome ~ (R+G+B+Lightness + Saturation + Hue)*(Lightness+Saturation), data = dfliaA_train)
mod10_glm %>% readr::write_rds("model10_glm.rds")
```

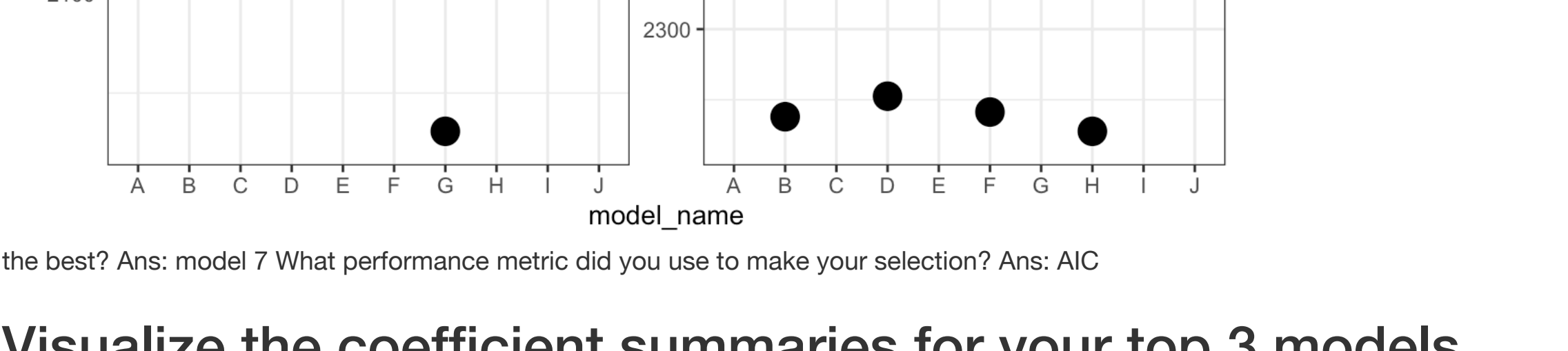
Did you experience any issues or warnings while fitting the generalized linear models? Ans: No.

```
extract_metrics <- function(mod_object, mod_name)
{
  broom::glance(mod_object) %>%
    mutate(model_name = mod_name)
}

glm_mle_results <- purrr::map2_dfr(list(mod01_glm,mod02_glm,mod03_glm,mod04_glm,mod05_glm,mod06_glm,mod07_glm,mod08_glm,mod09_glm,mod10_glm),
                                LETTERS[1:10],extract_metrics)

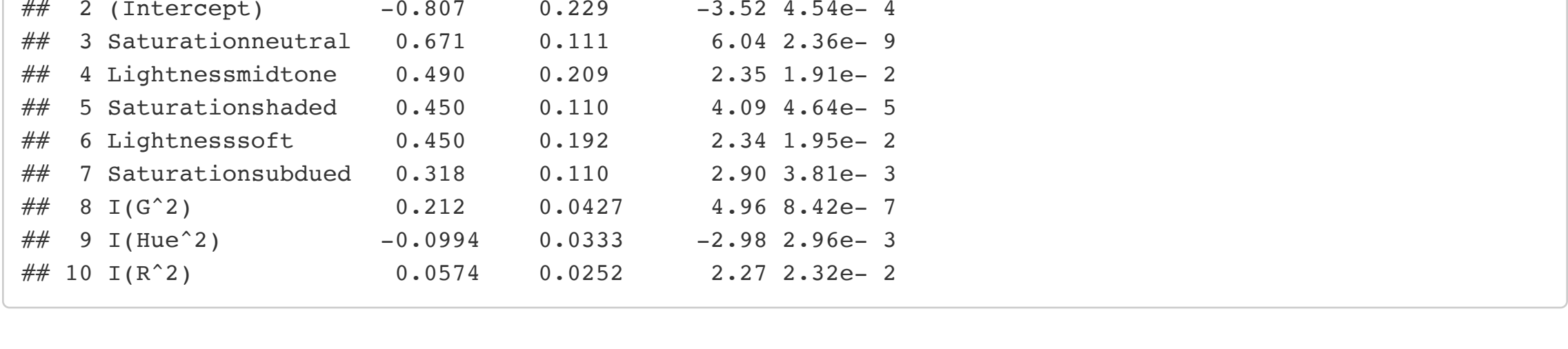
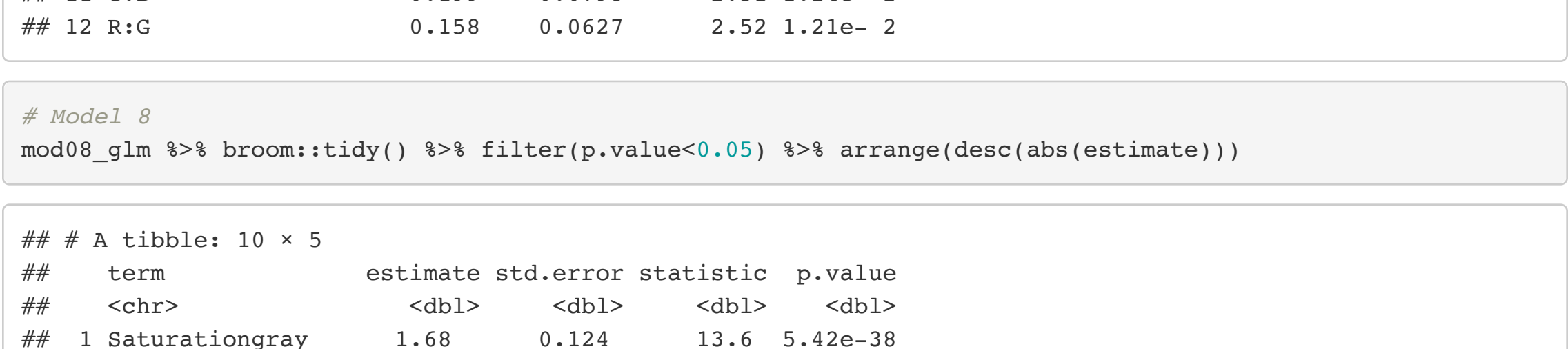
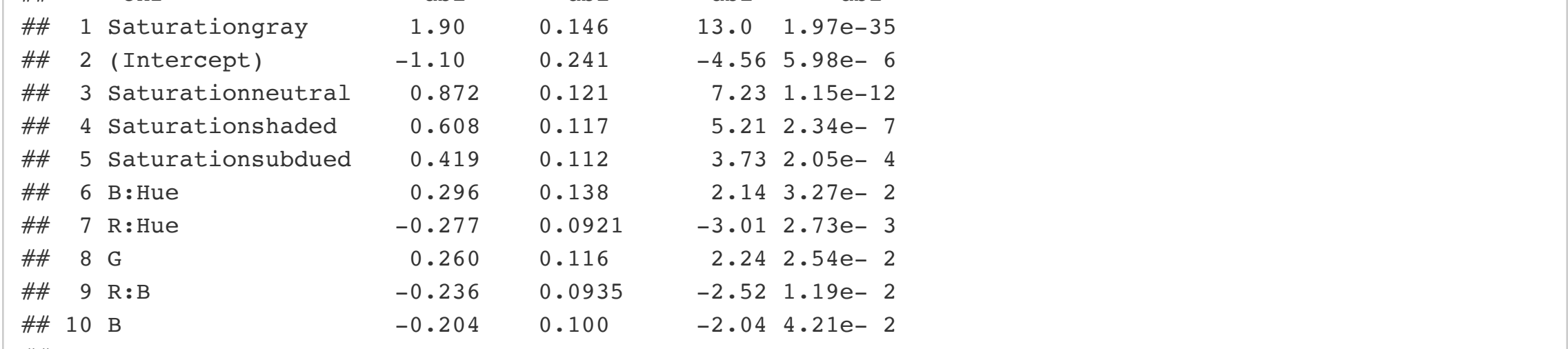
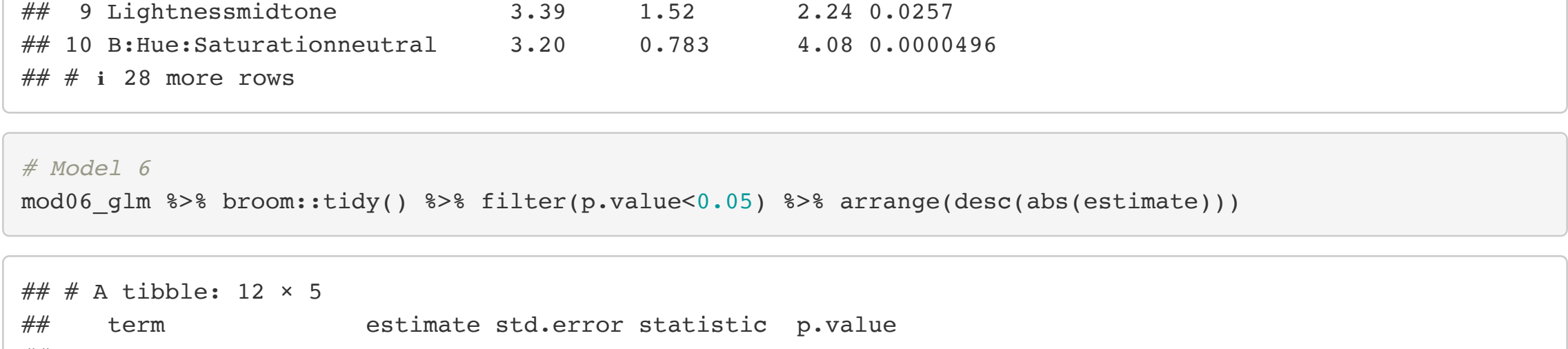
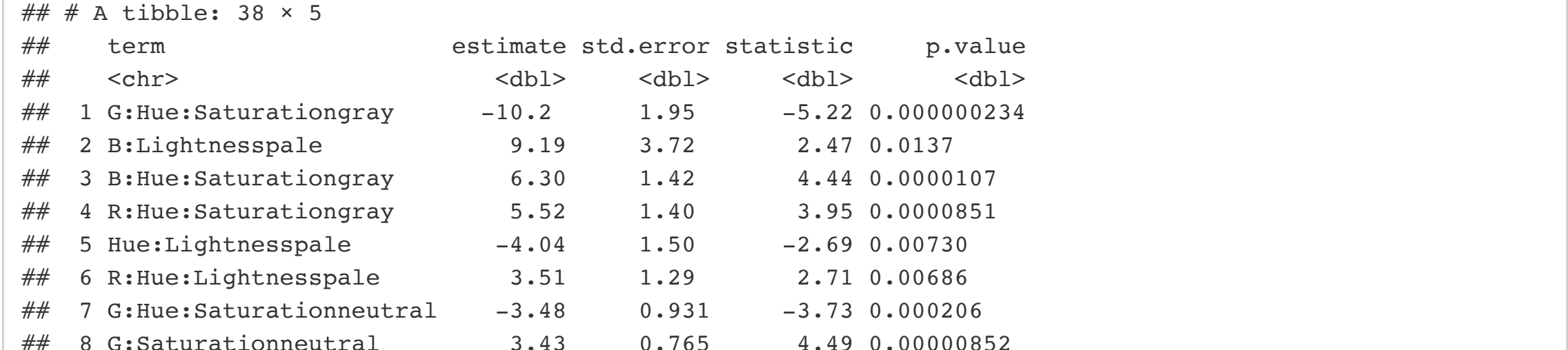
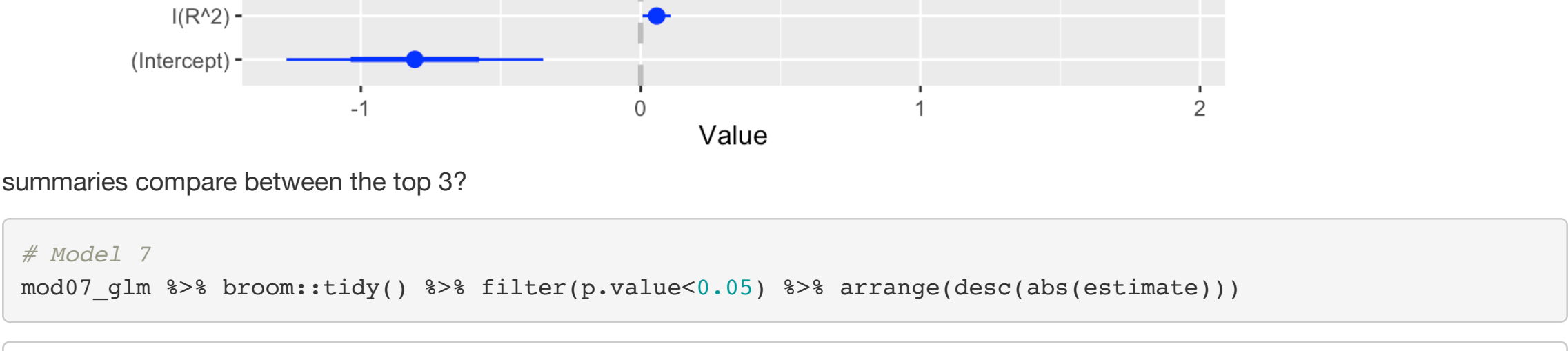
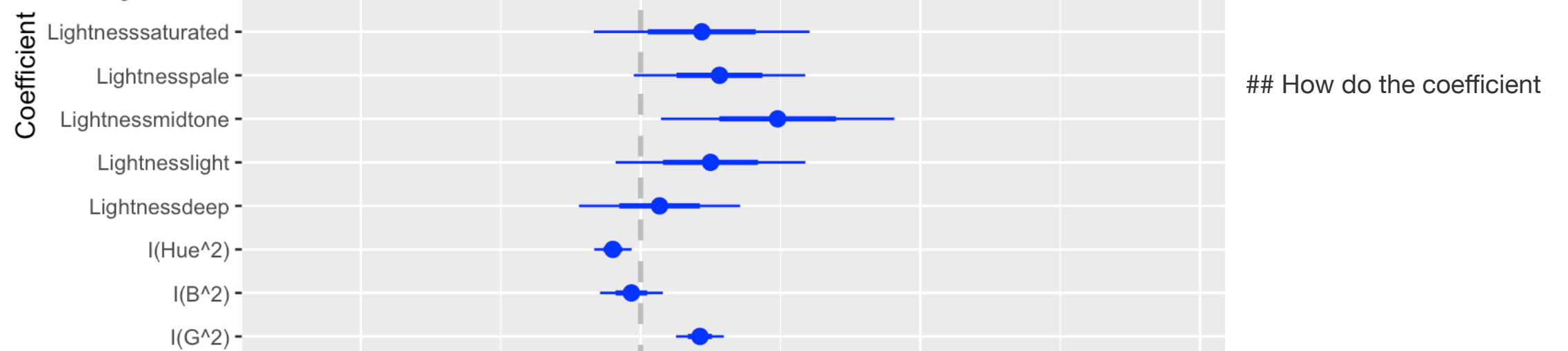
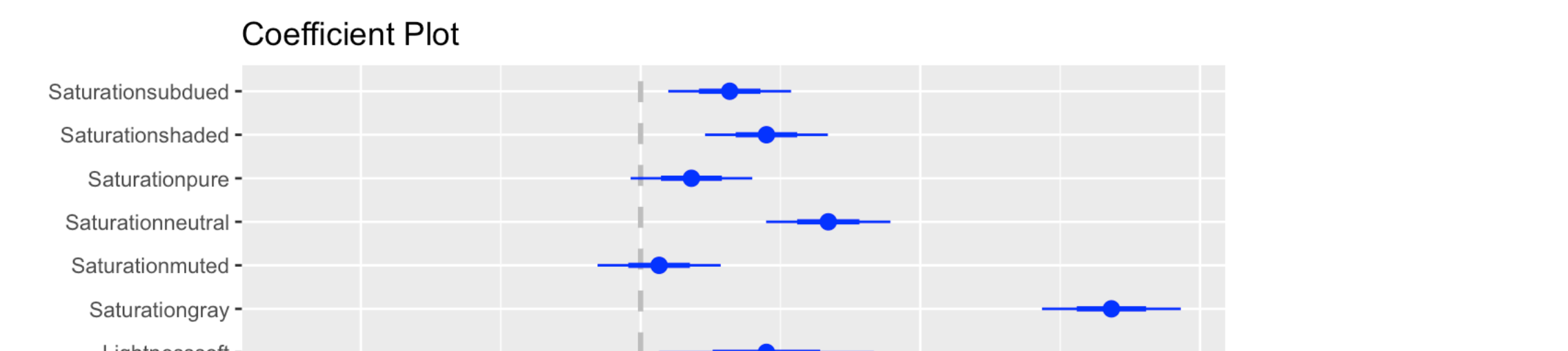
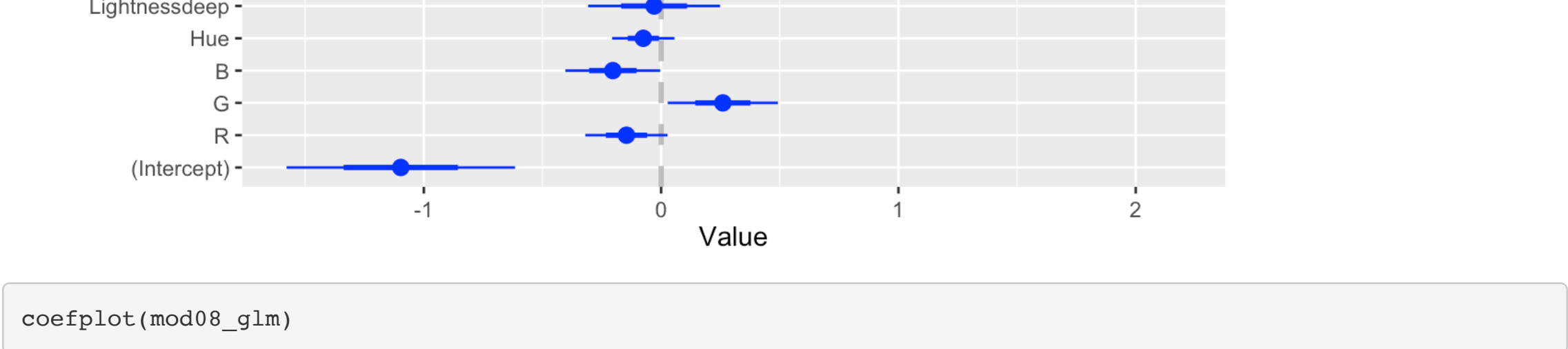
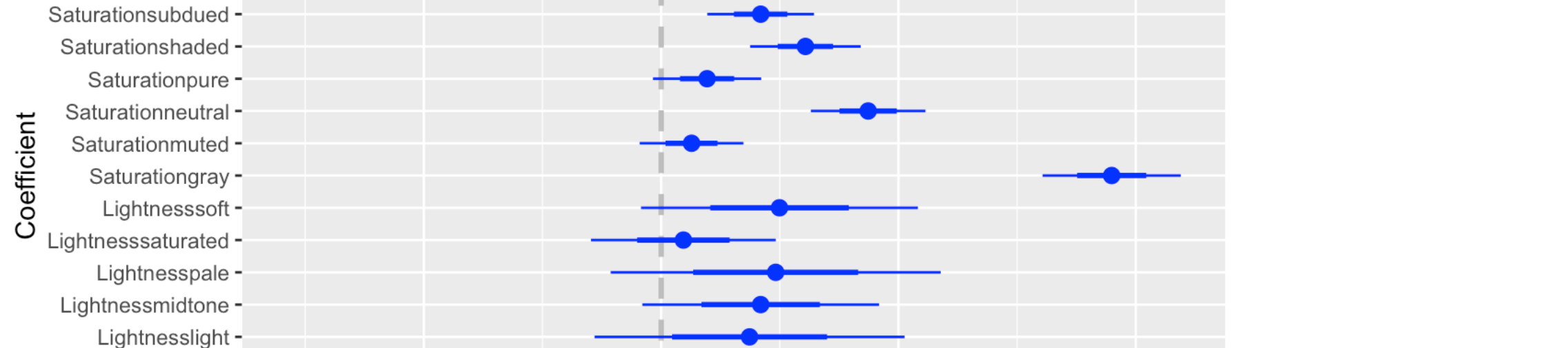
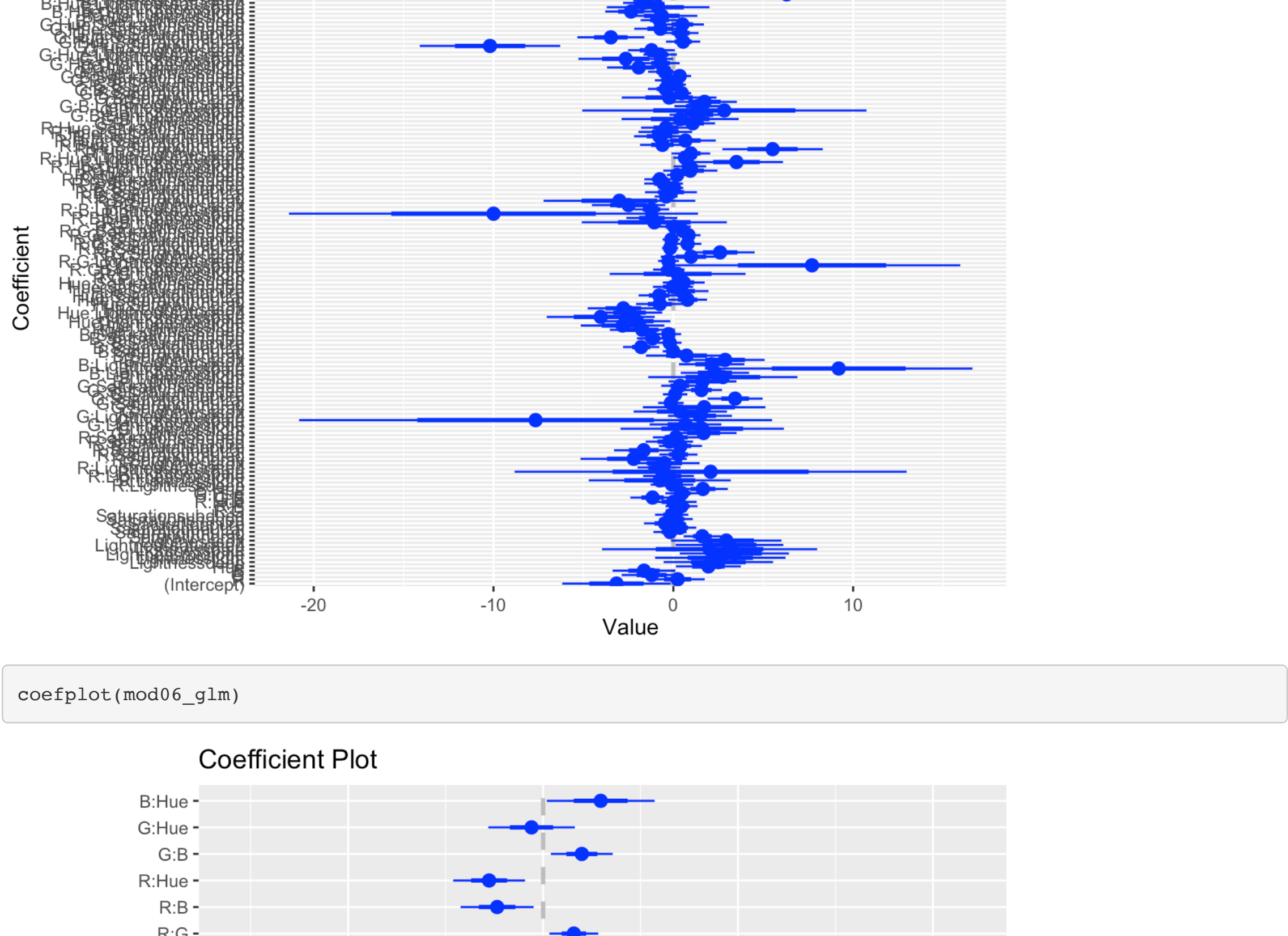
glm_mle_results %>%
  select(model_name, AIC, BIC) %>%
  pivot_longer(c(AIC, BIC)) %>%
  ggplot(mapping = aes(x = model_name, y = value)) +
  facet_wrap(size = 5) +
  geom_vline(nume, scales = 'free_y') +
  theme_bw()

Which of the 10 models is
```



the best? Ans: model 7 What performance metric did you use to make your selection? Ans: AIC

Visualize the coefficient summaries for your top 3 models.



Which inputs seem important?

Ans: So all the inputs continuous and categorical are important