

Statistical Modeling and Advanced Regression Analyses

R Tutorials

Holger Sennhenn-Reulen^{id}

Northwest German Forest Research Institute (NW-FVA), Germany.

November 20, 2024

Contents

1	Software	2
1.1	Organize R Session	2
2	Linear Regression Model	2
2.1	Data Simulation	2
2.1.1	Visualisations	3
2.2	Modeling	7
2.2.1	Visualisations	8
2.3	Add-Ons	9
2.3.1	Add-On Linear Model: A) Stancode	9
2.3.2	Add-On Linear Model: B) Posterior predictive check: an introduction ‘by hand’	11
3	Binary Regression Model	17
3.1	Data Simulation	17
3.1.1	Visualisations	17
3.2	Modeling	18
3.2.1	Visualisations	19
3.2.2	Estimated Expected Value	22
4	Poisson Regression Model	25
4.1	Data Simulation	25
4.1.1	Visualisations	25
4.2	Modeling	26
4.2.1	Estimated Expected Value	27
	References	30

1 Software

We use the statistical software environment *R* (R Core Team, 2024), and R add-on packages *ggplot2* (Wickham, 2016).

This document is produced using *Quarto* (Allaire et al., 2024).

1.1 Organize R Session

```
rm(list = ls())  
library("ggplot2")
```

2 Linear Regression Model

2.1 Data Simulation

Data are simulated according to the equations given in the lecture slides¹:

```
set.seed(123)  
N <- 500  
df <- data.frame(x_1 = runif(n = N),  
                 x_2 = runif(n = N))  
(beta_0 <- rnorm(n = 1, mean = 1, sd = .1))  
  
[1] 0.9398107  
  
(beta_x_1 <- rnorm(n = 1, mean = 1, sd = .1))  
  
[1] 0.9006301  
  
(beta_x_2 <- rnorm(n = 1, mean = -.5, sd = .1))  
  
[1] -0.3973215  
  
(sigma <- rgamma(n = 1, shape = 1, rate = 4))  
  
[1] 0.293026
```

¹For two covariates x_1 and x_2 .

```
df$mu <- beta_0 + beta_x_1 * df$x_1 + beta_x_2 * df$x_2
df$y <- df$mu + rnorm(n = N, mean = 0, sd = sigma)
```

2.1.1 Visualisations

```
ggplot(data = df, aes(x = x_1, y = x_2)) +
  geom_point()
```

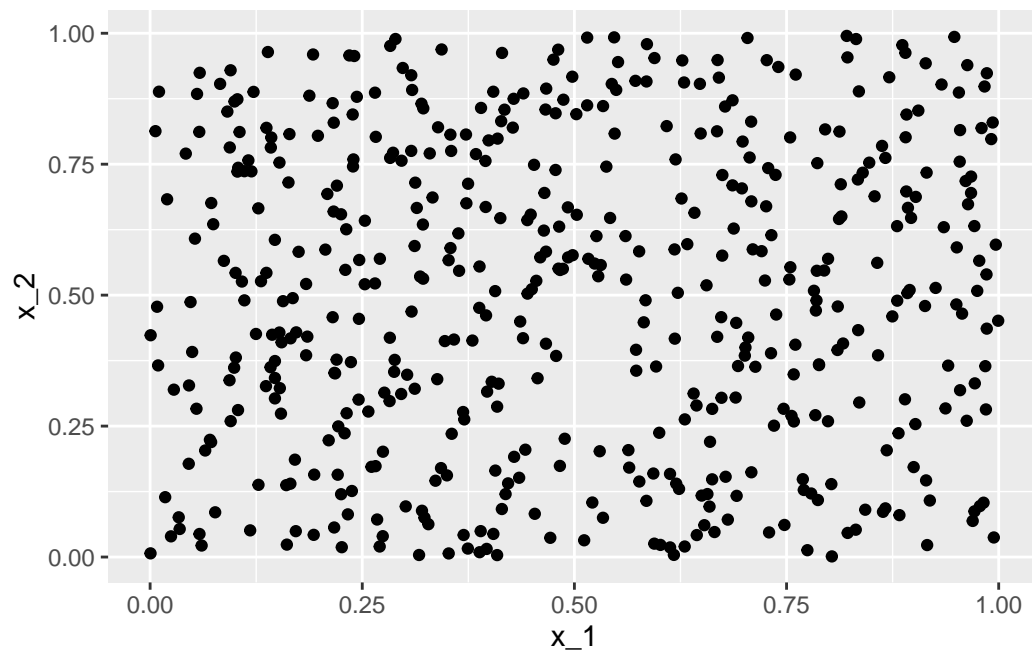


Figure 1: Scatterplot of the two simulated covariates x_1 and x_2 - each from the uniform distribution between 0 and 1.

```
ggplot(data = df, aes(x = x_1, y = mu, color = x_2)) +
  geom_point()
```

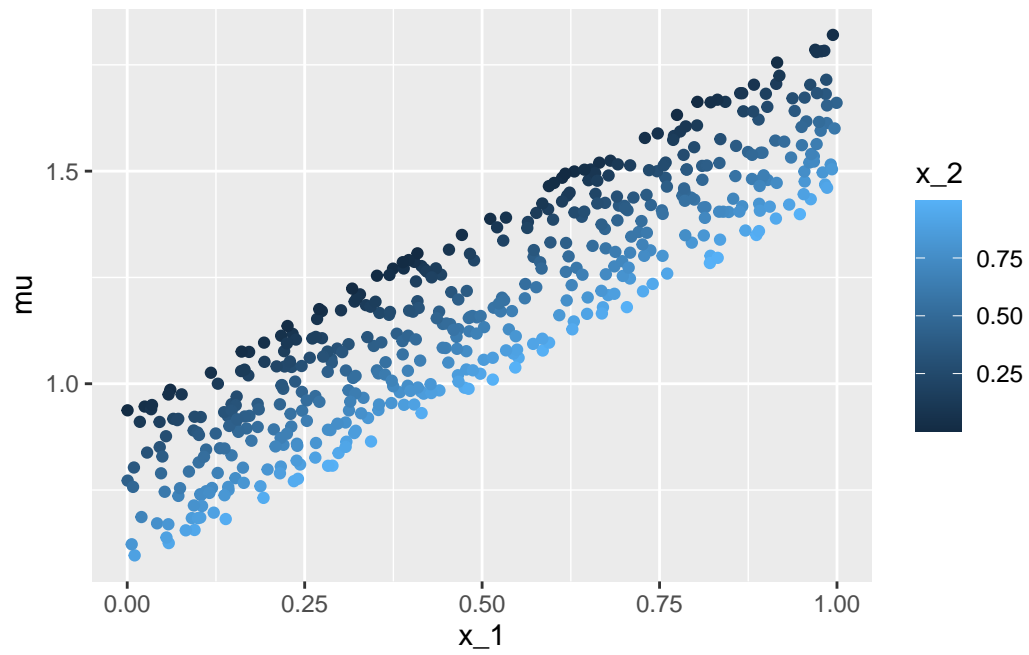


Figure 2: Scatterplot of covariate x_1 with response y - each individual observation is coloured according to the second covariate x_2 .

```
ggplot(data = df, aes(x = x_2, y = mu, color = x_1)) +  
  geom_point()
```

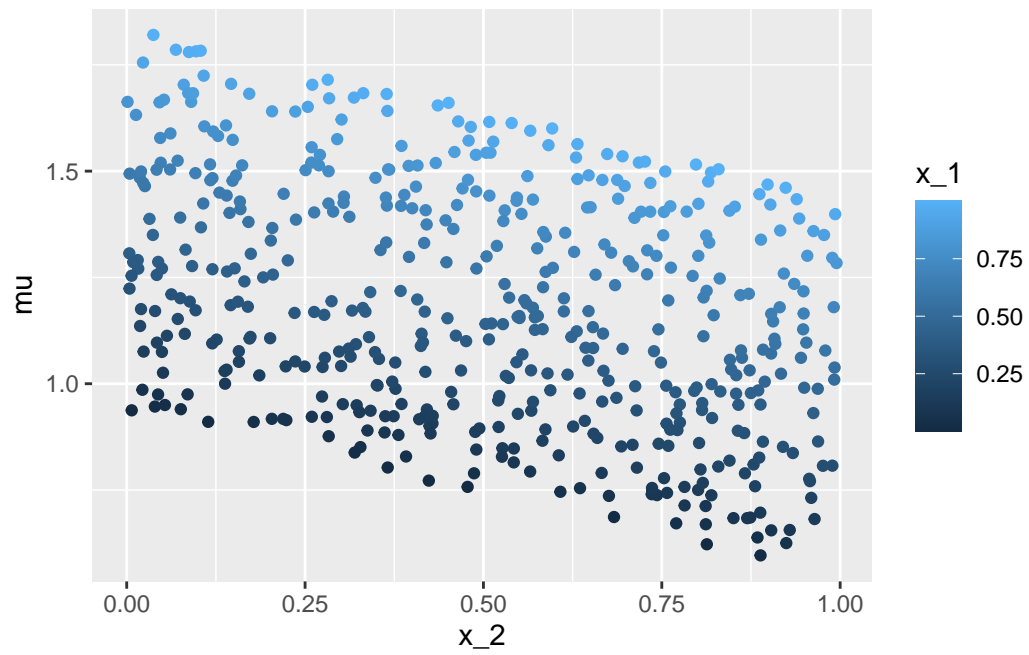


Figure 3: Scatterplot of covariate x_2 with response y - each individual observation is coloured according to the first covariate x_1 .

```
ggplot(data = df, aes(x = x_1, y = x_2, color = mu)) +  
  geom_point()
```

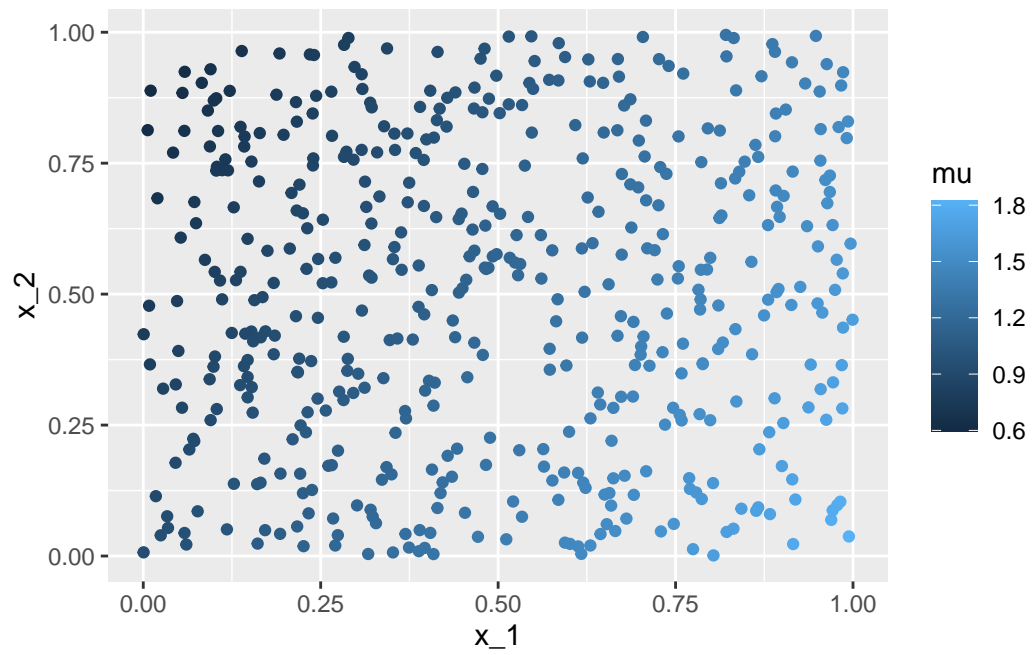


Figure 4: Scatterplot of the two simulated covariates x_1 and x_2 - each individual observation is coloured according to the underlying true conditional expectation μ .

```
ggplot(data = df, aes(x = x_1, y = x_2, color = y)) +  
  geom_point()
```

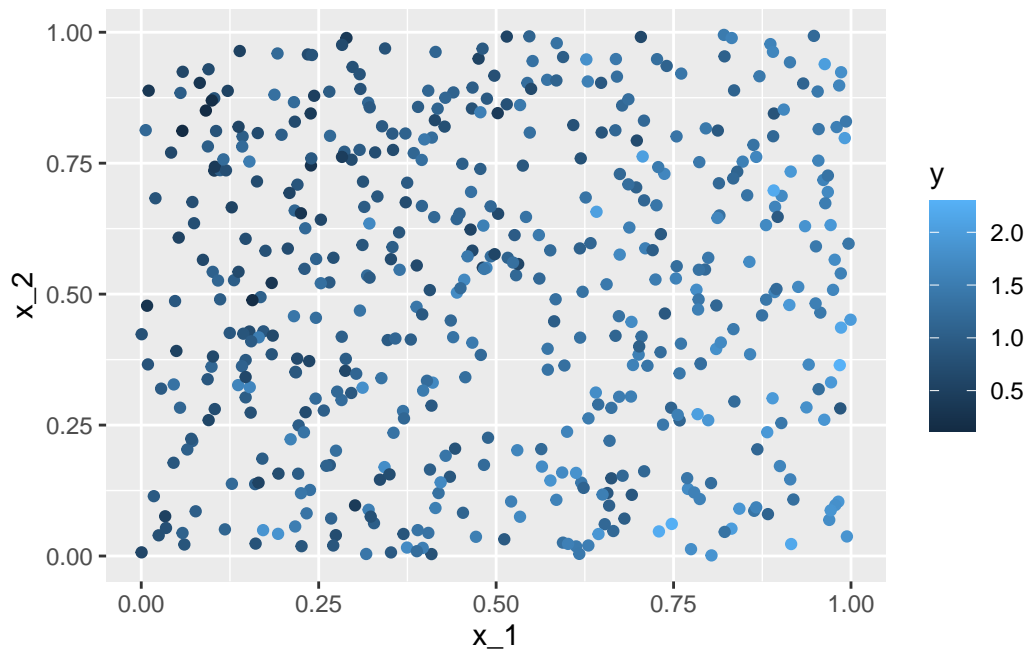


Figure 5: Scatterplot of the two simulated covariates x_1 and x_2 - each individual observation is coloured according to the response y .

2.2 Modeling

The basic R command for (frequentist) estimation of the parameters of a linear regression model is a call to the function `lm`:

Call:

```
lm(formula = y ~ x_1 + x_2, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.82082	-0.19805	0.00329	0.19051	0.81138

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.91291	0.03448	26.476	< 2e-16 ***
x_1	0.91533	0.04668	19.610	< 2e-16 ***
x_2	-0.36218	0.04566	-7.933	1.43e-14 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2963 on 497 degrees of freedom
 Multiple R-squared: 0.4674, Adjusted R-squared: 0.4652
 F-statistic: 218 on 2 and 497 DF, p-value: < 2.2e-16

2.2.1 Visualisations

```
nd <- data.frame(x_1 = seq(0, 1, by = .1),
                 x_2 = .5)
nd$mu <- predict(m, newdata = nd)
ggplot(data = df, aes(x = x_1, y = mu, color = x_2)) +
  geom_point() +
  geom_line(data = nd, aes(x = x_1, y = mu, color = x_2))
```

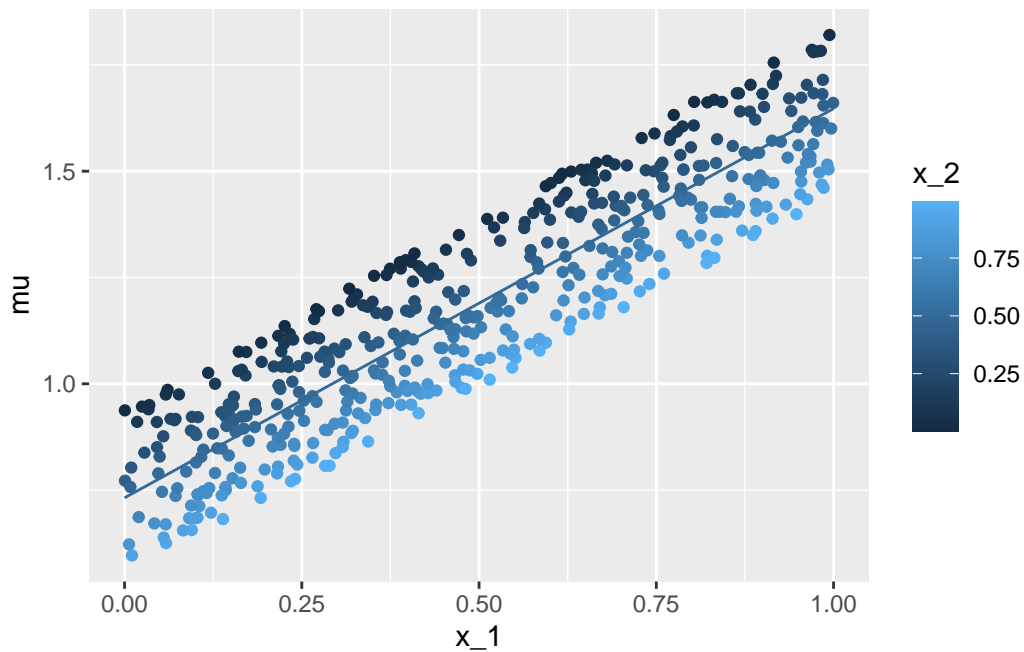


Figure 6: Scatterplot of covariate x_1 with the true conditional expectation μ - each individual observation is coloured according to the second covariate x_2 . The line gives the point estimation for the conditional expectation with the second covariate x_2 fixed to 0.5.

```
nd <- data.frame(expand.grid('x_1' = seq(0, 1, by = .1),
                             'x_2' = seq(0, 1, by = .1)))
nd$mu <- predict(m, newdata = nd)
ggplot(data = df, aes(x = x_1, y = mu, color = x_2)) +
```



```
geom_point() +
geom_line(data = nd, aes(x = x_1, y = mu, color = x_2, group = x_2))
```

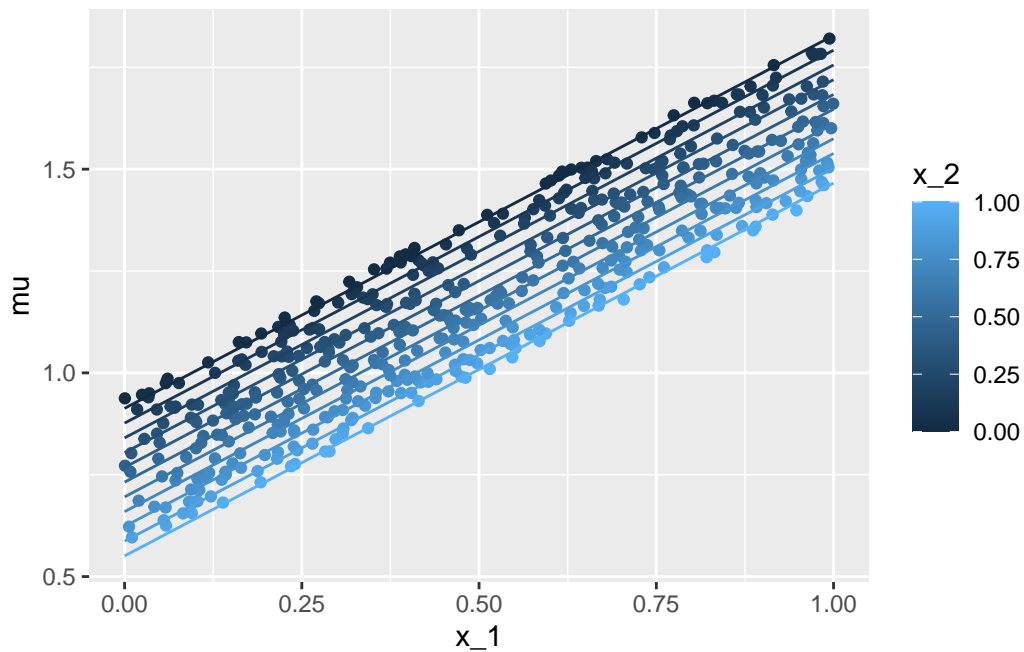


Figure 7: Scatterplot of covariate x_1 with the true conditional expectation μ - each individual observation is coloured according to the second covariate x_2 . The lines give the point estimation for the conditional expectation with the second covariate x_2 taking on values between 0 and 1 (at steps of 0.1).

2.3 Add-Ons

2.3.1 Add-On Linear Model: A) Stancode

2.3.1.1 Stan Users Guide

Probabilistic Programming Languages such as *Stan* (Carpenter et al., 2017) allow to plug together the single parts of a statistical regression model²:

The following Stan-code is published here in the Stan users guide:

```
data {
  int<lower=0> N;
  vector[N] x;
```

²Which is actually pretty 'readable' if you get used to the structure for a simple model such the linear regression model.

```

    vector[N] y;
  }
  parameters {
    real alpha;
    real beta;
    real<lower=0> sigma;
  }
  model {
    y ~ normal(alpha + beta * x, sigma);
  }

```

2.3.1.2 Stancode generated by calling `brms::brm`

The R add-on package *brms* (Bürkner, 2017, 2018) allows to implement advanced regression models without being an expert in ‘Stan-programming’.

Here is the Stan-code that is implemented by ‘brms’ for our linear regression model example:

```
brms::make_stancode(brms::bf(y ~ x_1 + x_2, center = F), data = df)
```

```

// generated with brms 2.21.0
functions {
}
data {
  int<lower=1> N; // total number of observations
  vector[N] Y; // response variable
  int<lower=1> K; // number of population-level effects
  matrix[N, K] X; // population-level design matrix
  int prior_only; // should the likelihood be ignored?
}
transformed data {
}
parameters {
  vector[K] b; // regression coefficients
  real<lower=0> sigma; // dispersion parameter
}
transformed parameters {
  real lprior = 0; // prior contributions to the log posterior
  lprior += student_t_lpdf(sigma | 3, 0, 2.5)
    - 1 * student_t_lccdf(0 | 3, 0, 2.5);
}
model {
  // likelihood including constants

```

```

    if (!prior_only) {
      target += normal_id_glm_lpdf(Y | X, 0, b, sigma);
    }
    // priors including constants
    target += lprior;
  }
generated quantities {
}

```

2.3.2 Add-On Linear Model: B) Posterior predictive check: an introduction 'by hand'

Having an `lm` object already, it is rather straightforward to get posterior samples by using function `sim` from the *arm* (Gelman & Su, 2024) package:

```
library("arm")
```

```

S <- sim(m)
str(S)

```

```

Formal class 'sim' [package "arm"] with 2 slots
..@ coef : num [1:100, 1:3] 0.882 1.014 0.904 0.978 0.958 ...
.. ..- attr(*, "dimnames")=List of 2
.. .. ..$ : NULL
.. .. ..$ : chr [1:3] "(Intercept)" "x_1" "x_2"
..@ sigma: num [1:100] 0.323 0.303 0.292 0.309 0.29 ...

```

```

S <- cbind(S@coef, 'sigma' = S@sigma)
head(S)

```

	(Intercept)	x_1	x_2	sigma
[1,]	0.8816414	0.9245094	-0.3362733	0.3227662
[2,]	1.0139849	0.7317948	-0.3398411	0.3033703
[3,]	0.9037042	0.9155575	-0.3506924	0.2922883
[4,]	0.9776909	0.8392790	-0.3845609	0.3090220
[5,]	0.9579213	0.8977625	-0.4284596	0.2900632
[6,]	0.9549211	0.8478278	-0.3937226	0.3094227

Predict the response for the covariate data as provided by the original data-frame `df` - here only by using the first posterior sample:

```
s <- 1
S[s, ]
```

```
(Intercept)      x_1      x_2      sigma
  0.8816414  0.9245094 -0.3362733  0.3227662
```

```
mu_s <- S[s, '(Intercept)'] + S[s, 'x_1'] * df$x_1 + S[s, 'x_2'] * df$x_2
y_s <- rnorm(n = nrow(df), mean = mu_s, sd = S[s, 'sigma'])
pp <- rbind(data.frame(y = df$y, source = "original", s = s),
            data.frame(y = y_s, source = "predicted", s = s))
ggplot(data = pp, aes(x = y, fill = source)) +
  geom_histogram(alpha = .5, position = "identity")
```

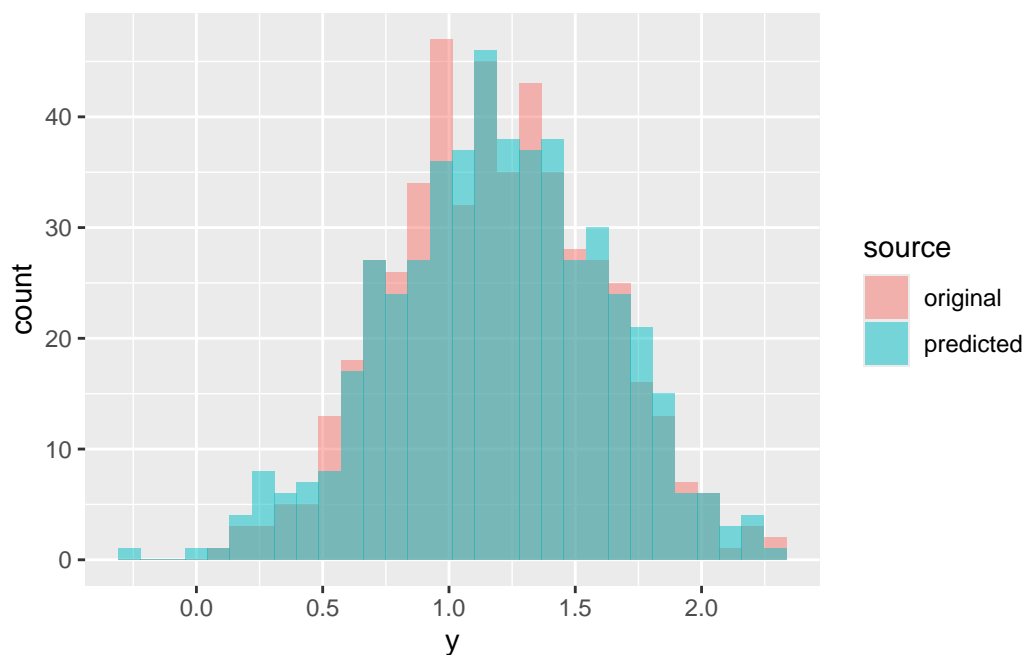


Figure 8: Histogram of the original and the posterior predicted response sample.

Now let's repeat the same for 9 different posterior samples:

```
pp <- NULL
for (s in 1:9) {
  mu_s <- S[s, '(Intercept)'] + S[s, 'x_1'] * df$x_1 + S[s, 'x_2'] * df$x_2
  y_s <- rnorm(n = nrow(df), mean = mu_s, sd = S[s, 'sigma'])
  pp <- rbind(pp,
              data.frame(y = df$y, source = "original", s = s),
```

```

    data.frame(y = y_s, source = "predicted", s = s))
}
ggplot(data = pp, aes(x = y, fill = source)) +
  geom_histogram(alpha = .5, position = "identity") +
  facet_wrap(~ s)

```

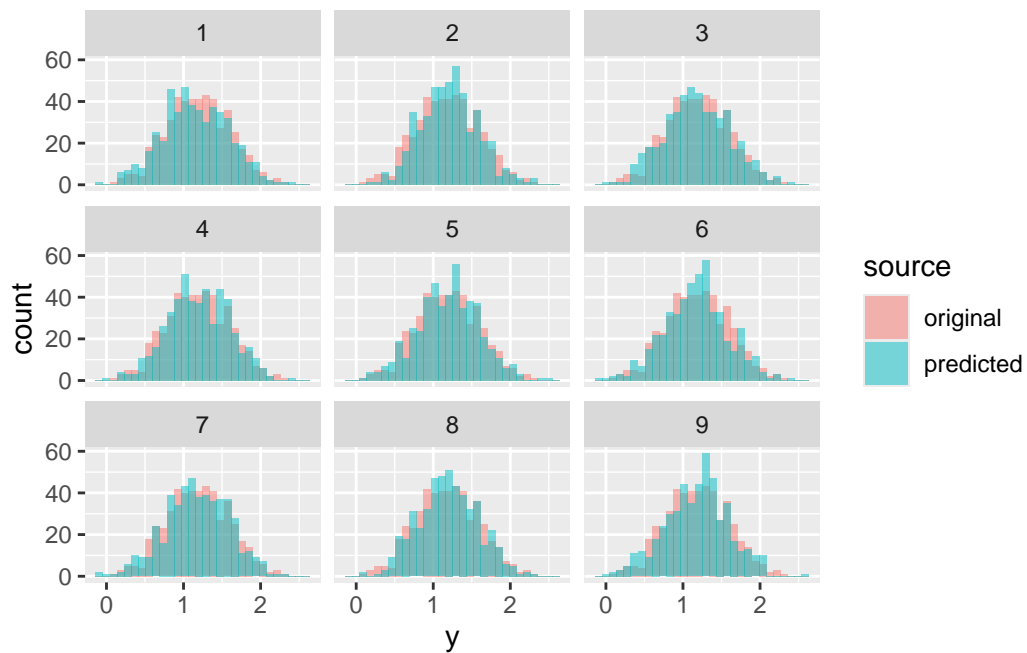


Figure 9: Histogram of the original and the posterior predicted response sample.

```

ggplot(data = pp, aes(x = y, fill = source)) +
  geom_density(alpha = .5, position = "identity") +
  facet_wrap(~ s)

```

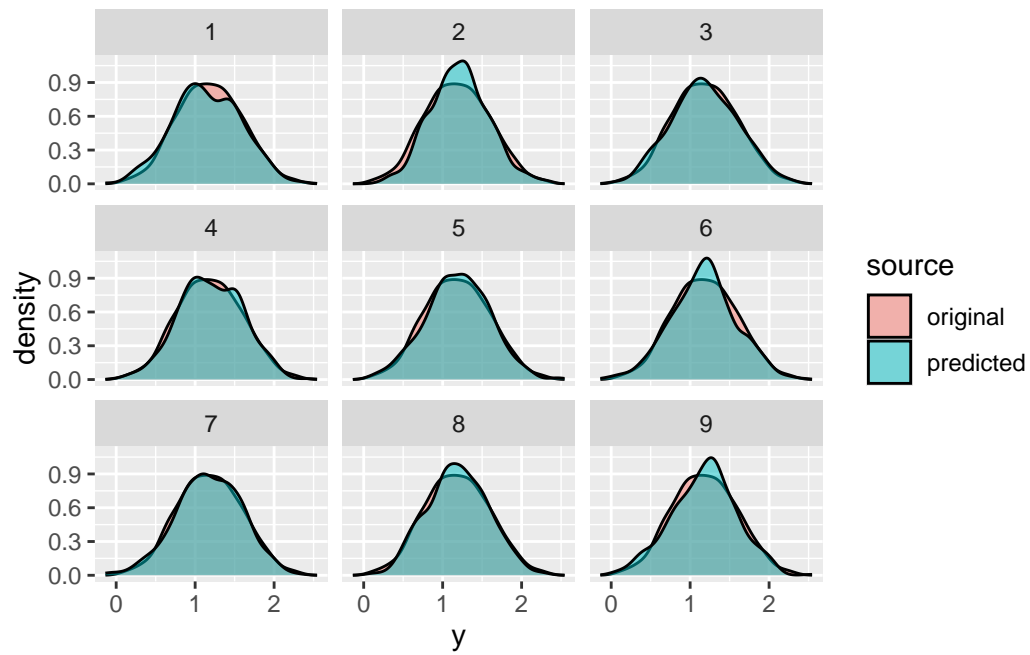


Figure 10: The same as in Figure 9, but now using kernel density visualisations.

```
ggplot(data = pp, aes(x = y, colour = source)) +
  stat_ecdf() +
  facet_wrap(~ s)
```

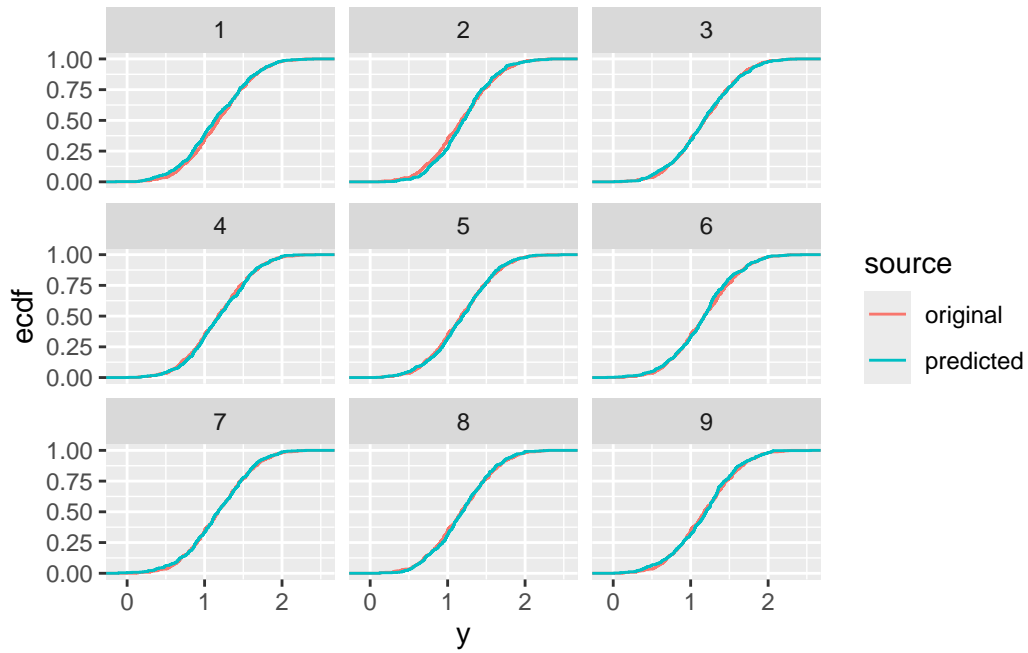


Figure 11: The same as in Figure 9 or Figure 10, but now using empirical cumulative density function visualisations.

```
ggplot(data = subset(pp, source == "predicted"),
  aes(x = y, group = s)) +
  geom_density(position = "identity", fill = NA, colour = "grey") +
  geom_density(data = subset(pp, source == "original" & s == 1),
    aes(x = y), linewidth = 1)
```

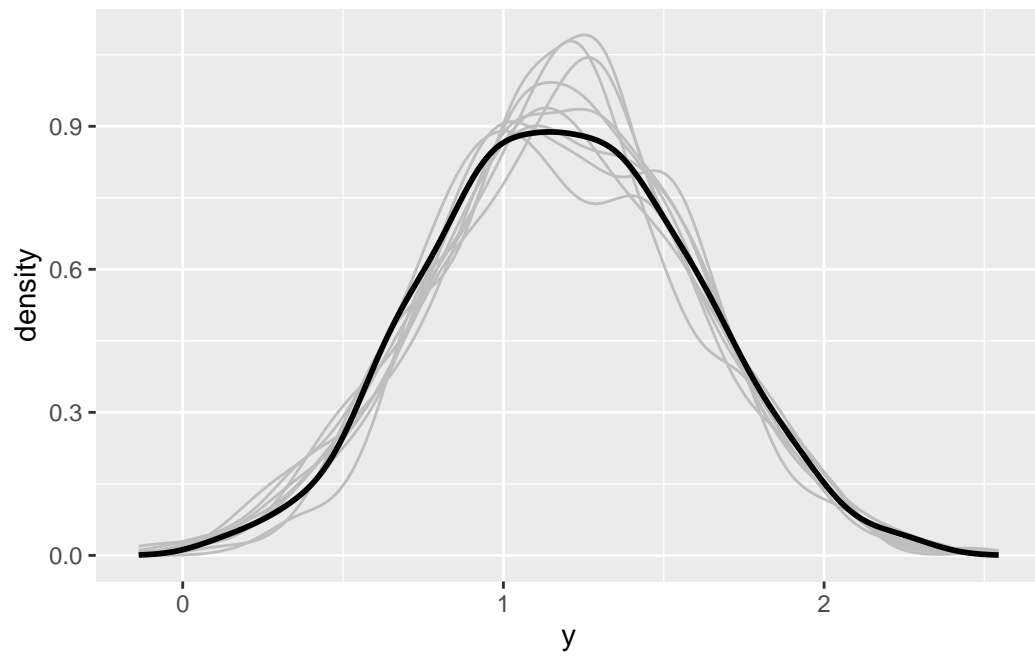


Figure 12: The same as in Figure 12, but now within one plotting window: This visualisation is what `brms : :pp_check` will produce if applied on a `brm` object.

3 Binary Regression Model

```
rm(list = ls())
library("ggplot2")
library("plyr")
```

3.1 Data Simulation

Data are simulated similarly as for the linear model:

```
set.seed(123)
N <- 500
df <- data.frame(x_1 = runif(n = N),
                 x_2 = runif(n = N))
(beta_0 <- rnorm(n = 1, mean = 0, sd = .1))

[1] -0.06018928

(beta_x_1 <- rnorm(n = 1, mean = 1, sd = .1))

[1] 0.9006301

(beta_x_2 <- rnorm(n = 1, mean = -.5, sd = .1))

[1] -0.3973215

df$eta <- beta_0 + beta_x_1 * df$x_1 + beta_x_2 * df$x_2
df$y <- rbinom(n = N, size = 1, prob = plogis(q = df$eta))
```

3.1.1 Visualisations

```
df$x_1_c <- cut(df$x_1, breaks = seq(0, 1, by = .1),
               include.lowest = T,
               labels = seq(.05, .95, by = .1))
df$x_1_c <- as.numeric(as.character(df$x_1_c))
df_p_A <- ddply(df, c("x_1_c"), summarise,
               p = mean(y > .5))
```

```

df_p_A <- data.frame('p' = rep(df_p_A$p, each = 2),
                     'x_1' = sort(c(df_p_A$x_1_c - .05,
                                     df_p_A$x_1_c + .05)))
df_p_B <- data.frame('x_1' = seq(0, 1, by = .01),
                     'p' = plogis(beta_0 +
                                   beta_x_1 * seq(0, 1, by = .01) +
                                   beta_x_2 * .5))

set.seed(0)
ggplot(data = df, aes(x = x_1, y = y)) +
  geom_jitter(aes(color = x_2), width = 0, height = .1) +
  geom_line(data = df_p_A, aes(y = p, group = p)) +
  geom_line(data = df_p_B, aes(y = p), linetype = 2)
## ... 'not as linear as it seems':
# plot(df_p_B$x_1[-1], diff(df_p_B$p))

```

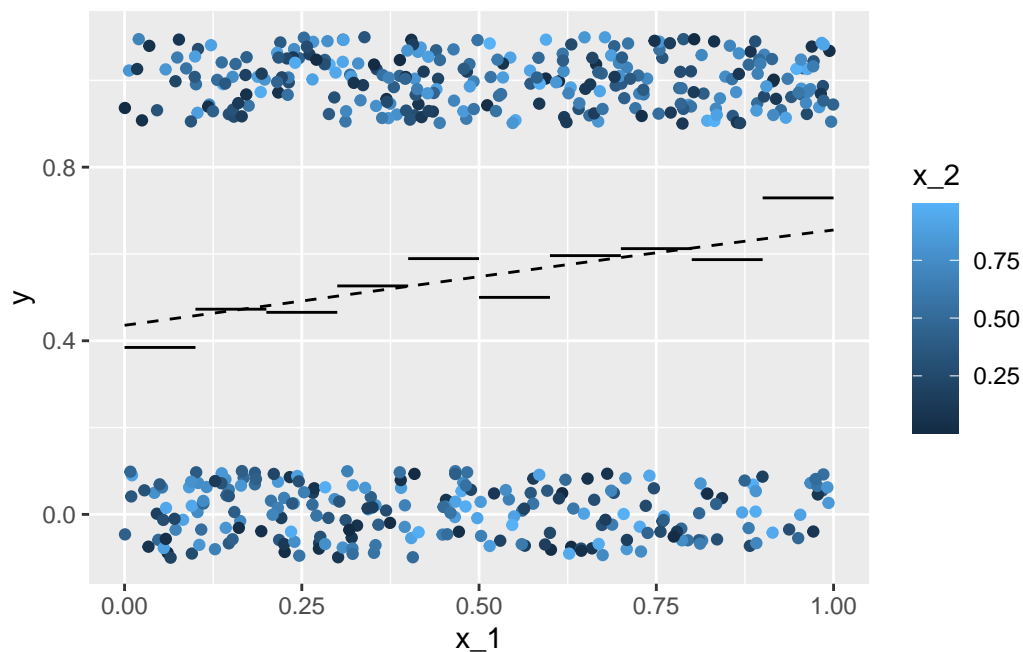


Figure 13: Scatterplot of covariate x_1 with response y - each individual observation is coloured according to the second covariate x_2 , and additionally ‘jittered’ in vertical direction.

3.2 Modeling

The basic R command for (frequentist) estimation of the parameters of a binary regression model is a call to the function `glm` with family argument `binomial`:

```
m <- glm(y ~ x_1 + x_2, data = df,
         family = binomial(link = 'logit'))
summary(m)
```

```
Call:
glm(formula = y ~ x_1 + x_2, family = binomial(link = "logit"),
    data = df)
```

```
Coefficients:
                Estimate Std. Error z value Pr(>|z|)
(Intercept)   -0.2908     0.2358  -1.233 0.217531
x_1             1.1598     0.3248   3.570 0.000356 ***
x_2            -0.1713     0.3138  -0.546 0.585034
---

```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 688.53  on 499  degrees of freedom
Residual deviance: 675.30  on 497  degrees of freedom
AIC: 681.3
```

```
Number of Fisher Scoring iterations: 4
```

3.2.1 Visualisations

```
nd <- data.frame('x_1' = 0:1, 'x_2' = .5)
(nd$eta <- predict(m, newdata = nd, type = 'link'))
```

```
      1      2
-0.3764387  0.7833343
```

```
coef(m)[1] + c(0, 1) * coef(m)[2] + .5 * coef(m)[3]
```

```
[1] -0.3764387  0.7833343
```

```
pA <- ggplot(data = df, aes(x = x_1, y = eta, color = x_2)) +
  geom_point() +
  geom_line(data = nd, aes(x = x_1, y = eta, color = x_2))
nd <- data.frame('x_1' = .5,
  'x_2' = 0:1)
(nd$eta <- predict(m, newdata = nd, type = 'link'))
```

```
      1      2
0.2891090 0.1177866
```

```
coef(m)[1] + .5 * coef(m)[2] + c(0, 1) * coef(m)[3]
```

```
[1] 0.2891090 0.1177866
```

```
pB <- ggplot(data = df, aes(x = x_2, y = eta, color = x_1)) +
  geom_point() +
  geom_line(data = nd, aes(x = x_2, y = eta, color = x_1))
cowplot::plot_grid(pA, pB, ncol = 2)
```

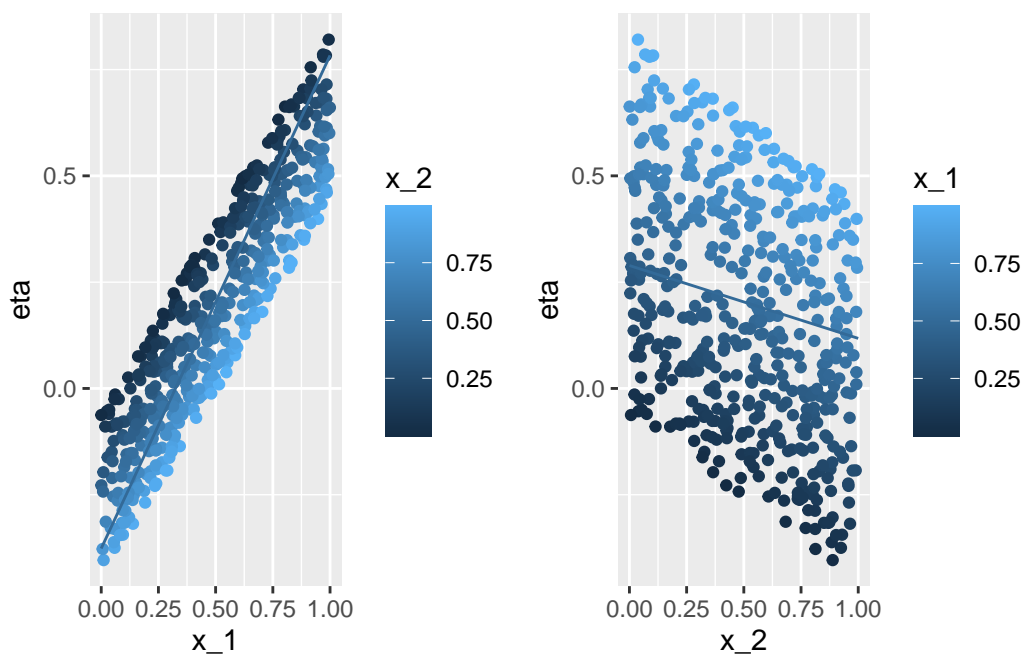


Figure 14: Scatterplot of covariate x_1 with the true linear predictor η - each individual observation is coloured according to the second covariate x_2 . The line gives the point estimation for the conditional expectation with the second covariate x_2 fixed to 0.5.

```

nd <- data.frame(x_1 = seq(0, 1, by = .1),
                 x_2 = .5)
(nd$p <- predict(m, newdata = nd, type = 'response'))

[1] 0.4069861 0.4352503 0.4639417 0.4928738 0.5218537 0.5506872 0.5791841
[8] 0.6071630 0.6344556 0.6609111 0.6863983

plogis(coef(m)[1] + seq(0, 1, by = .1) * coef(m)[2] +
       .5 * coef(m)[3])

[1] 0.4069861 0.4352503 0.4639417 0.4928738 0.5218537 0.5506872 0.5791841
[8] 0.6071630 0.6344556 0.6609111 0.6863983

ggplot(data = df, aes(x = x_1, y = y)) +
  geom_jitter(aes(color = x_2), width = 0, height = .1) +
  geom_line(data = nd, aes(x = x_1, y = p, color = x_2)) +
  geom_line(data = df_p_B, aes(y = p), linetype = 2)

```

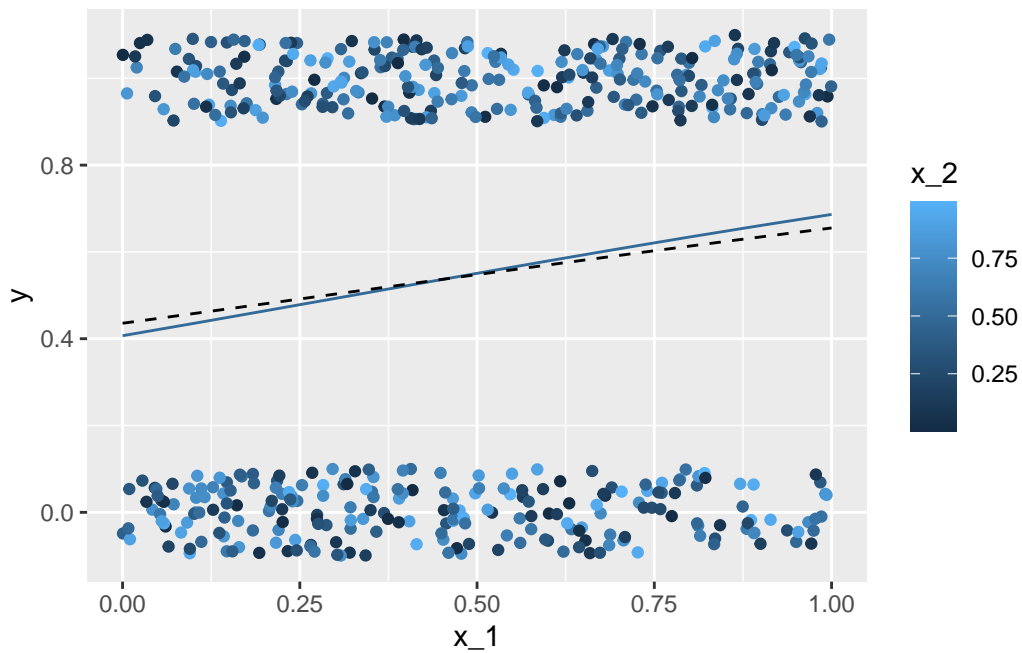


Figure 15: Scatterplot of covariate x_1 with the true conditional expectation p - each individual observation is coloured according to the second covariate x_2 . The line gives the point estimation for the conditional expectation with the second covariate x_2 fixed to 0.5.

3.2.2 Estimated Expected Value

We can apply the Bernstein-von Mises theorem to estimate the *expected value*:

- **Fit the model:** Obtain the maximum likelihood estimate for the model's coefficients (`coef`) along with their variance-covariance matrix (`vcov`).
- **Simulate coefficients:** Perform an 'informal' Bayesian posterior simulation using the multivariate normal distribution, based on the *Bernstein-von Mises theorem*.
- **Convert simulated coefficients:** Apply an appropriate transformation to the simulated coefficients to compute the *simulated quantity of interest*. This quantity typically depends on the values of all explanatory variables, and researchers may:
 - Focus on a specific observation (usually an 'average'), or
 - Average across all sample observations.

In both cases, the applied transformation incorporates the researcher's specific choice.

```
library("MASS")
coef(m)

(Intercept)      x_1      x_2
-0.2907775    1.1597730 -0.1713224

vcov(m)

              (Intercept)      x_1      x_2
(Intercept)  0.05560471 -0.048970067 -0.047028038
x_1          -0.04897007  0.105509175 -0.004560743
x_2          -0.04702804 -0.004560743  0.098439583

set.seed(0)
B <- mvrnorm(n = 100, mu = coef(m), Sigma = vcov(m))
head(B)

      (Intercept)      x_1      x_2
[1,] -0.08125910  0.6544775 -0.2581602
[2,] -0.40299145  1.3779659 -0.3263178
[3,]  0.09915843  1.0089580 -0.5398310
[4,]  0.03289839  0.8600445 -0.3880109
[5,] -0.12814786  1.3256621 -0.5036957
[6,] -0.55953065  1.4562644  0.3176658
```

```

nd <- expand.grid('x_1' = nd$x_1,
                 'x_2' = nd$x_2,
                 's' = 1:nrow(B))

head(nd)

  x_1 x_2 s
1 0.0 0.5 1
2 0.1 0.5 1
3 0.2 0.5 1
4 0.3 0.5 1
5 0.4 0.5 1
6 0.5 0.5 1

nd$p <- plogis(B[nd$s, 1] + B[nd$s, 2] * nd$x_1 +
              B[nd$s, 3] * nd$x_2)
dd <- ddply(nd, c('x_1'), summarise,
            p_mean = mean(p),
            p_lwr_95 = quantile(p, prob = .025),
            p_upr_95 = quantile(p, prob = .975),
            p_lwr_9 = quantile(p, prob = .05),
            p_upr_9 = quantile(p, prob = .95),
            p_lwr_75 = quantile(p, prob = .125),
            p_upr_75 = quantile(p, prob = .875))
set.seed(0)
ggplot(data = df, aes(x = x_1)) +
  geom_jitter(aes(y = y, color = x_2), width = 0, height = .1) +
  geom_ribbon(data = dd, aes(x = x_1, ymin = p_lwr_95,
                           ymax = p_upr_95), alpha = .4) +
  geom_ribbon(data = dd, aes(x = x_1, ymin = p_lwr_9,
                           ymax = p_upr_9), alpha = .4) +
  geom_ribbon(data = dd, aes(x = x_1, ymin = p_lwr_75,
                           ymax = p_upr_75), alpha = .4) +
  geom_line(data = dd, aes(y = p_mean))

```

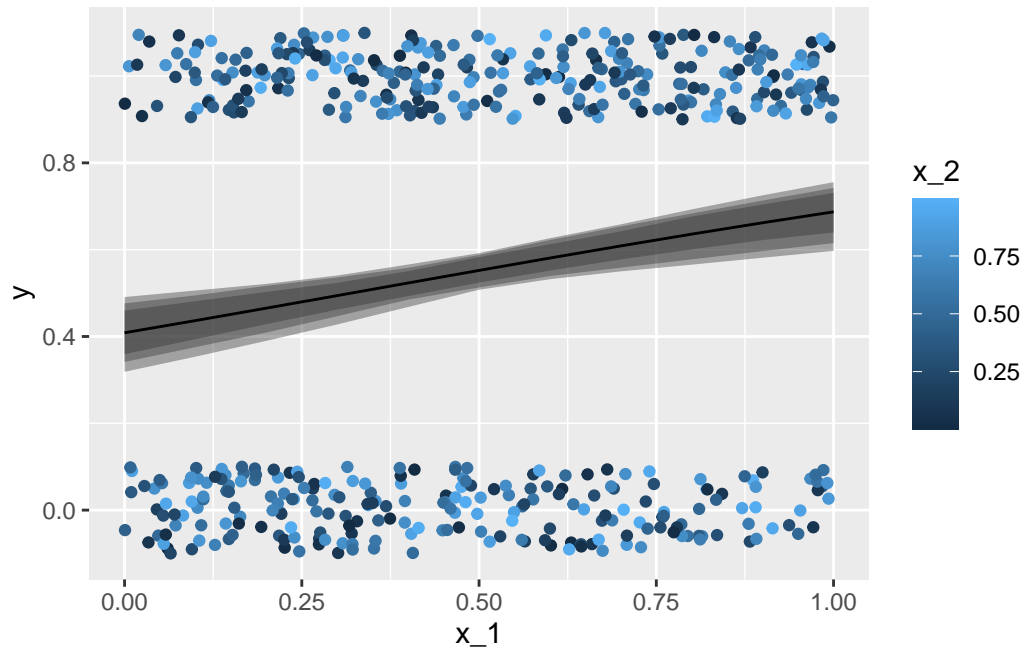


Figure 16: Scatterplot of covariate x_1 with the true conditional expectation μ - each individual observation is coloured according to the second covariate x_2 . The line gives the point estimation for the conditional expectation with the second covariate x_2 fixed to 0.5.

4 Poisson Regression Model

```
rm(list = ls())
library("ggplot2")
```

4.1 Data Simulation

Data are simulated similarly as for the linear model:

```
set.seed(123)
N <- 500
df <- data.frame(x_1 = runif(n = N),
                 x_2 = runif(n = N))
(beta_0 <- rnorm(n = 1, mean = 0, sd = .1))

[1] -0.06018928

(beta_x_1 <- rnorm(n = 1, mean = 1, sd = .1))

[1] 0.9006301

(beta_x_2 <- rnorm(n = 1, mean = -.5, sd = .1))

[1] -0.3973215

df$eta <- beta_0 + beta_x_1 * df$x_1 + beta_x_2 * df$x_2
df$y <- rpois(n = N, lambda = exp(df$eta))
```

4.1.1 Visualisations

```
df$x_1_c <- cut(df$x_1, breaks = seq(0, 1, by = .1),
               include.lowest = T,
               labels = seq(.05, .95, by = .1))
df$x_1_c <- as.numeric(as.character(df$x_1_c))
df_p_A <- ddply(df, c("x_1_c"), summarise,
               mu = mean(y))
df_p_A <- data.frame('mu' = rep(df_p_A$mu, each = 2),
```

```

      'x_1' = sort(c(df_p_A$x_1_c - .05,
                    df_p_A$x_1_c + .05)))
df_p_B <- data.frame('x_1' = seq(0, 1, by = .01),
                     'mu' = exp(beta_0 +
                                beta_x_1 * seq(0, 1, by = .01) +
                                beta_x_2 * .5))

set.seed(0)
ggplot(data = df, aes(x = x_1, y = y)) +
  geom_jitter(aes(color = x_2), width = 0, height = .1) +
  geom_line(data = df_p_A, aes(y = mu, group = mu)) +
  geom_line(data = df_p_B, aes(y = mu), linetype = 2)

```

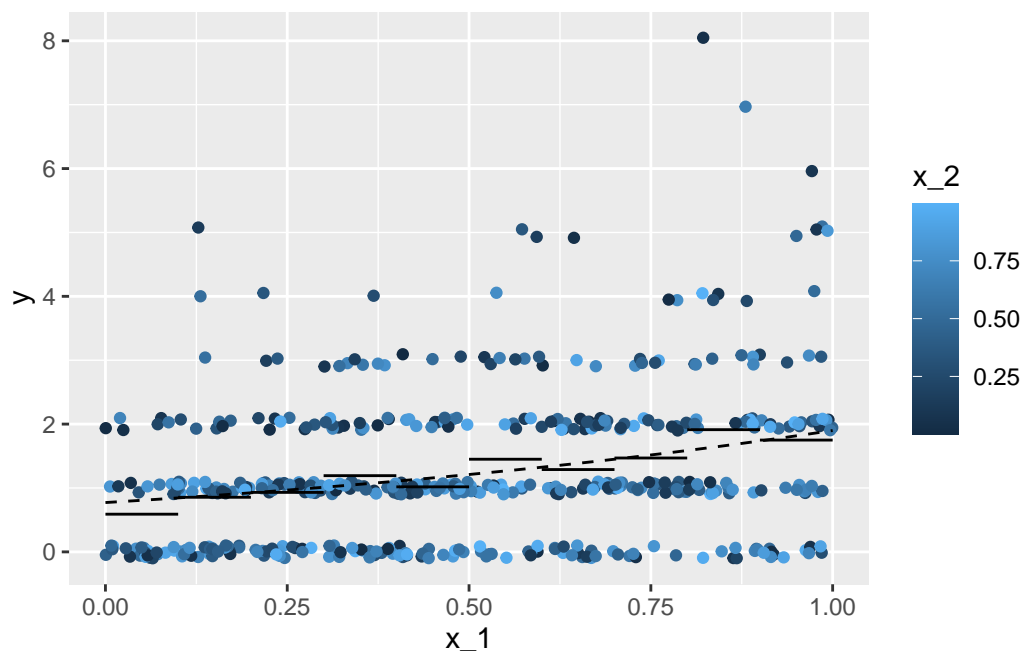


Figure 17: Scatterplot of covariate x_1 with response y - each individual observation is coloured according to the second covariate x_2 .

4.2 Modeling

The basic R command for (frequentist) estimation of the parameters of a binary regression model is a call to the function `glm` with family argument `poisson(link = 'log')`:

```

m <- glm(y ~ x_1 + x_2, data = df, family = poisson(link = 'log'))
summary(m)

```

```

Call:
glm(formula = y ~ x_1 + x_2, family = poisson(link = "log"),
    data = df)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -0.09637    0.11000  -0.876   0.381
x_1          1.05534    0.14351   7.354 1.93e-13 ***
x_2         -0.54067    0.13875  -3.897 9.74e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 1)

    Null deviance: 619.76  on 499  degrees of freedom
Residual deviance: 551.67  on 497  degrees of freedom
AIC: 1395.1

Number of Fisher Scoring iterations: 5

```

4.2.1 Estimated Expected Value

Let's again apply the Bernstein-von Mises theorem

```

library("MASS")
coef(m)

(Intercept)      x_1      x_2
-0.09636825  1.05534471 -0.54067416

vcov(m)

              (Intercept)      x_1      x_2
(Intercept)  0.012100215 -0.0115419704 -0.0083283575
x_1          -0.011541970  0.0205956476 -0.0008112633
x_2          -0.008328358 -0.0008112633  0.0192505213

set.seed(0)
B <- mvrnorm(n = 100, mu = coef(m), Sigma = vcov(m))
head(B)

```

```

      (Intercept)      x_1      x_2
[1,]  0.05743986  0.8596548 -0.5240625
[2,] -0.10120645  1.1692825 -0.5989887
[3,]  0.01910641  0.9263818 -0.7232511
[4,]  0.02386701  0.8912981 -0.6321912
[5,] -0.06117581  1.0833727 -0.7137618
[6,] -0.30539283  1.1687677 -0.3825595

nd <- expand.grid('x_1' = seq(0, 1, by = .1),
                  'x_2' = .5,
                  's' = 1:nrow(B))

head(nd)

  x_1 x_2 s
1 0.0 0.5 1
2 0.1 0.5 1
3 0.2 0.5 1
4 0.3 0.5 1
5 0.4 0.5 1
6 0.5 0.5 1

nd$mu <- exp(B[nd$s, 1] +
             B[nd$s, 2] * nd$x_1 +
             B[nd$s, 3] * nd$x_2)
dd <- ddply(nd, c('x_1'), summarise,
            mu_mean = mean(mu),
            mu_lwr_95 = quantile(mu, prob = .025),
            mu_upr_95 = quantile(mu, prob = .975),
            mu_lwr_9 = quantile(mu, prob = .05),
            mu_upr_9 = quantile(mu, prob = .95),
            mu_lwr_75 = quantile(mu, prob = .125),
            mu_upr_75 = quantile(mu, prob = .875))
df_p_B <- data.frame('x_1' = seq(0, 1, by = .01),
                     'mu' = exp(coef(m)[1] +
                                coef(m)[2] * seq(0, 1, by = .01) +
                                coef(m)[3] * .5))

set.seed(0)
ggplot(data = df, aes(x = x_1)) +
  geom_jitter(aes(y = y, color = x_2), width = 0, height = .1) +
  geom_ribbon(data = dd, aes(x = x_1, ymin = mu_lwr_95,
                           ymax = mu_upr_95), alpha = .4) +
  geom_ribbon(data = dd, aes(x = x_1, ymin = mu_lwr_9,
                           ymax = mu_upr_9), alpha = .4) +

```

```
geom_ribbon(data = dd, aes(x = x_1, ymin = mu_lwr_75,
                          ymax = mu_upr_75), alpha = .4) +
geom_line(data = dd, aes(y = mu_mean)) +
geom_line(data = df_p_B, aes(y = mu), linetype = 2)
```

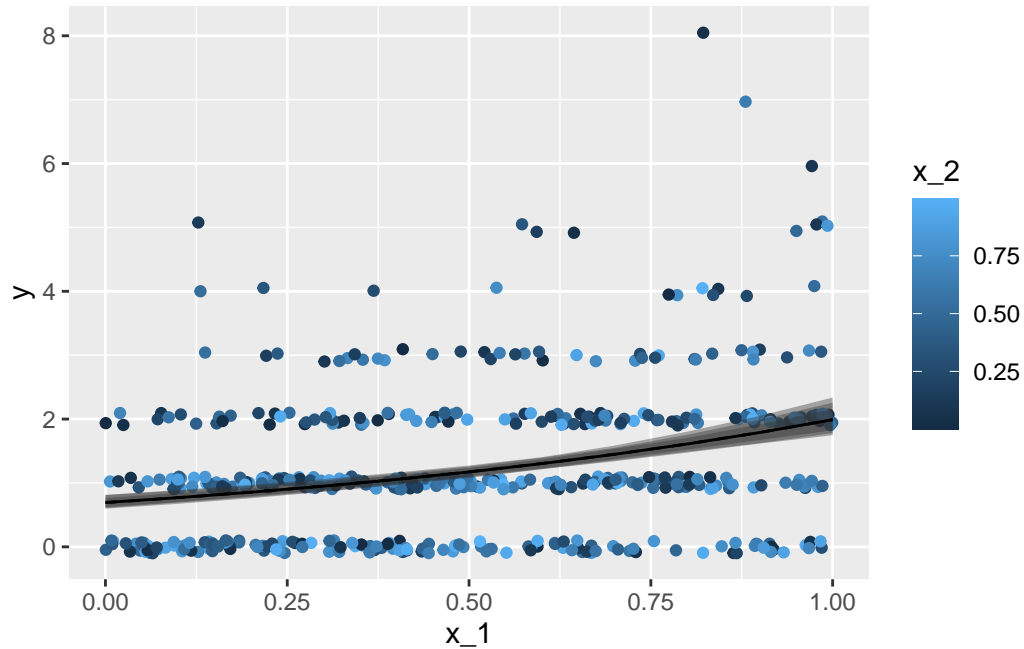


Figure 18: Scatterplot of covariate x_1 with the response observations y - each individual observation is coloured according to the second covariate x_2 . The line gives the point estimation for the conditional expectation with the second covariate x_2 fixed to 0.5, the coloured intervals give point-wise central 75%, 90%, and 95% credible intervals for the conditional expectation.

References

- Allaire, J. J., Teague, C., Scheidegger, C., Xie, Y., & Dervieux, C. (2024). *Quarto (Version 1.4.553)*. <https://doi.org/10.5281/zenodo.5960048>
- Bürkner, P.-C. (2017). Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software*, 80, 1–28. <https://doi.org/10.18637/jss.v080.i01>
- Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The R Journal*, 10(1), 395–411.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P., & Riddell, A. (2017). Stan: A Probabilistic Programming Language. *Journal of Statistical Software*, 76, 1–32. <https://doi.org/10.18637/jss.v076.i01>
- Gelman, A., & Su, Y.-S. (2024). *Arm: Data analysis using regression and multilevel/hierarchical models*. <https://CRAN.R-project.org/package=arm>
- R Core Team. (2024). *R: A Language and Environment for Statistical Computing (Version 4.4.1)*. R Foundation for Statistical Computing.
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>