



МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ
Федеральное государственное автономное образовательное учреждение
высшего образования
«Дальневосточный федеральный университет»
(ДВФУ)

**ИНСТИТУТ МАТЕМАТИКИ И КОМПЬЮТЕРНЫХ
ТЕХНОЛОГИЙ**
Департамент информационных и компьютерных систем

Курс «Компьютерные методы анализа больших данных»

Лабораторная работа №2
Контрольное мероприятие по рейтингу

на тему «Очистка и сравнительный анализ интернет-данных»
вариант №14

Выполнил студент Б9122-01.03.02мкт
Пелагеев Д.И.

Проверил доцент Достовалов В.Н.

г. Владивосток
2025

Оглавление

Введение	3
1 Подготовка данных	4
2 Построение моделей и визуализация	8
3 Прогнозирование по модели m1 (доверительный и предиктивный интервалы)	12
Заключение	14
Список использованных источников	16
A Полный код на R	17

Введение

В рамках данной работы проводится анализ данных опроса студентов курса «Stat 100» Университета Иллинойса за осенний семестр 2018 года. Основная цель анализа — изучить взаимосвязь между количеством напитков, потребляемых в неделю (`drinks_per_week`), временем, проводимым на вечеринках (`party_hours_per_week`), полом (`gender`) и типом населённого пункта (`home_bc`), в котором проживает студент.

В ходе выполнения работы были реализованы следующие этапы:

а) Загрузка и предварительная обработка исходного файла данных Stat 100 Survey 2, Fall 2018.dat, включая очистку и структурирование информации.

б) Преобразование и кодирование категориальных переменных, в том числе модификация столбца Home town для формирования базового уровня.

в) Построение двух моделей линейной регрессии:

- Модель m1: зависимость `drinks_per_week` от `party_hours_per_week` и `gender`;
- Модель m2: зависимость `drinks_per_week` от `party_hours_per_week`, `gender` и `home_bc`.

г) Визуализация результатов обеих регрессий на одном графике для сравнения.

д) Проведение прогнозирования по модели m1 с построением доверительных и предиктивных интервалов.

Анализ подобных данных позволяет выявить, какие социально-демографические и поведенческие факторы оказывают влияние на уровень потребления алкогольных напитков среди студентов, что важно для последующего построения профилактических и образовательных программ.

1 Подготовка данных

Исходный файл Stat 100 Survey 2, Fall 2018.dat содержит анкетные данные студентов. Для начала были считаны исходные заголовки, которые содержали лишние пояснения и пробелы. После первичной очистки (удаления скобок и лишних пробелов) они выглядели так:

Исходные имена признаков после очистки:

```
1 raw_header <- readLines(DATA_FILE, n = 1)
2 clean_header <- gsub("\\s*\\([~()]*\\)", "", raw_header)
3 clean_header <- gsub("\\s+", " ", clean_header)
4 clean_header <- trimws(clean_header)
5
6
7 header_tokens <- strsplit(clean_header, " ")[[1]]
8 col_names_raw <- header_tokens[-1]
9 col_names_raw <- col_names_raw[-2]
10 col_names_raw
```

Первично обработанный список содержит 32 переменные, включая как числовые, так и категориальные.

"Gender", "Greek", "Home_Town", "In_State", "Ethnicity",
"Religion", "Religious", "like_Math", "Calculus", "ACT",
"GPA", "Party_Hours_per_week", "Drinks_per_week",
"Sex_Partners", "Relationships", "Call_Parents",
"Cell_Phone", "Social_Media", "Texts", "Good_or_Well",
"Expected_Income", "President", "Illinois_Commitment",
"Family_Income", "Liberal", "Political_Party", "Grade_vs_Learning",
"Parent_Relationship", "Work_Hours", "Tuition", "Career", "Section"

Преобразование имён к стилю snake_case:

Для удобства анализа имена признаков были приведены к единому стилю, без пробелов и с нижним подчёркиванием:

```
1 col_names <- make_snake(col_names_raw)
```

Результат:

"gender", "greek", "home_town", "in_state", "ethnicity",
"religion", "religious", "like_math", "calculus", "act", "gpa",

"party_hours_per_week", "drinks_per_week", "sex_partners",
 "relationships", "call_parents", "cell_phone", "social_media",
 "texts", "good_or_well", "expected_income", "president",
 "illinois_commitment", "family_income", "liberal",
 "political_party", "grade_vs_learning", "parent_relationship",
 "work_hours", "tuition", "career", "section"

Пошаговое формирование таблицы:

Сначала данные были считаны из файла. Первая строка с названиями была пропущена, а лишние технические столбцы удалены:

```
1 df <- read.table(DATA_FILE, header = FALSE, sep = ",", stringsAsFactors =  
  FALSE, skip = 1, fill = TRUE)  
2 head(df)
```

Таблица 1 — Фрагмент таблицы после чтения файла (ещё без имён признаков)

V1	V2	V3	...	V32	V33	V34	V35
Obs1	1	1	...	10	40	5	1
Obs2	2	1	...	20	100	3	1
Obs3	3	1	...	0	100	7	1
...

Затем удалялись две первые технические колонки (номера и дублирующий столбец):

```
1 df <- df[, -c(1:2)]  
2 df <- df[, -2]  
3 head(df)
```

Таблица 2 — Фрагмент таблицы после удаления служебных столбцов

V3	V5	V6	...	V32	V33	V34	V35
1	0	3	...	10	40	5	1
1	0	2	...	20	100	3	1
1	1	2	...	0	100	7	1
...

Далее таблице были присвоены очищенные и стандартизированные имена столбцов:

```
1 names(df) <- col_names
2 head(df)
```

Таблица 3 — Фрагмент таблицы с корректными названиями признаков

gender	greek	home_town	...	family_income	career	section
1	0	3	...	85	5	1
1	0	2	...	80	3	1
1	1	2	...	50	7	1
...

Следующим шагом была обработка и перекодировка категориальных переменных. Были выделены основные факторы и преобразованы в `factor`:

```
1 factor_cols <- c("gender", "greek", "home_town", "in_state", "ethnicity"
2               ,
3               "religion", "calculus", "cell_phone", "president",
4               "political_party", "section")
5 factor_cols <- intersect(factor_cols, names(df))
df[factor_cols] <- lapply(df[factor_cols], factor)
```

Проверка уровней факторных переменных дала следующий результат:

```
$gender: [1] "0" "1"
$home_town: [1] "0" "1" "2" "3"
$cell_phone: [1] "0" "1" "2"
и т.д.
```

После этого проведена фильтрация: из всех числовых столбцов были удалены строки с некорректными (отрицательными или равными 999) значениями. Затем была проверена и подтверждена полная отсутствия пропущенных значений:

```
1 numeric_mask <- sapply(df, is.numeric)
2 problem_mat <- df[, numeric_mask, drop = FALSE] < 0 | df[, numeric_mask
, drop = FALSE] == 999
```

```

3 rows_bad      <- rowSums(problem_mat, na.rm = TRUE) > 0
4 df <- df[!rows_bad, ]
5 colSums(is.na(df))

```

В завершение был введён дополнительный фактор — бинарная переменная `home_bc`, где "big_city" выделен как отдельная группа, а остальные категории объединены:

```

1 df$home_bc <- factor(ifelse(df$home_town == "3", 1, 0),
2                       levels = c(0, 1), labels = c("other", "big_city"))

```

Для построения моделей рабочий датафрейм был сокращён до ключевых переменных:

```

1 df_work <- df[, c("gender", "party_hours_per_week", "drinks_per_week", "
2   home_bc")]
3 head(df_work)

```

Таблица 4 — Финальная таблица для моделирования (фрагмент)

gender	party_hours_per_week	drinks_per_week	home_bc
1	4	12	big_city
1	3	1	other
1	10	50	other
1	1	1	other
1	0	0	big_city
0	4	3	other

Таким образом, процесс подготовки данных включает в себя пошаговое преобразование структуры, удаление некорректных и неинформативных наблюдений, а также перекодировку категориальных признаков — всё это позволило получить чистый, структурированный датасет, полностью готовый к последующему анализу и моделированию.

2 Построение моделей и визуализация

На следующем этапе были построены две линейные регрессионные модели, в которых целевой переменной выступает количество напитков, потребляемых в неделю (`drinks_per_week`), а объясняющими — количество часов, проводимых на вечеринках (`party_hours_per_week`), пол (`gender`), а также тип населённого пункта (`home_bc`).

Для первой модели (**m1**) использовались переменные `party_hours_per_week` и `gender`:

```
1 m1 <- lm(drinks_per_week ~ party_hours_per_week + gender, data = df_work
  )
```

Вторая модель (**m2**) дополнительно учитывает переменную `home_bc` (тип населённого пункта):

```
1 m2 <- lm(drinks_per_week ~ party_hours_per_week + gender + home_bc, data
  = df_work)
```

В обеих моделях пол и тип населённого пункта выступают категориальными предикторами.

Для наглядной интерпретации были построены графики, на которых визуализированы зависимости для обеих моделей.

Визуализация модели m1

На графике 1 представлены результаты регрессии `m1`. По оси X отложено число часов на вечеринках в неделю, по оси Y — количество потребляемых напитков в неделю. Для каждой группы по полу отображены свои линии регрессии, а также области доверительных и предиктивных интервалов (полупрозрачные полосы).

Доверительный интервал показывает, где находится истинная средняя линия регрессии для каждого пола с заданной вероятностью (обычно 95%). Предиктивный интервал шире и показывает, где с большей вероятностью окажется индивидуальное новое наблюдение.

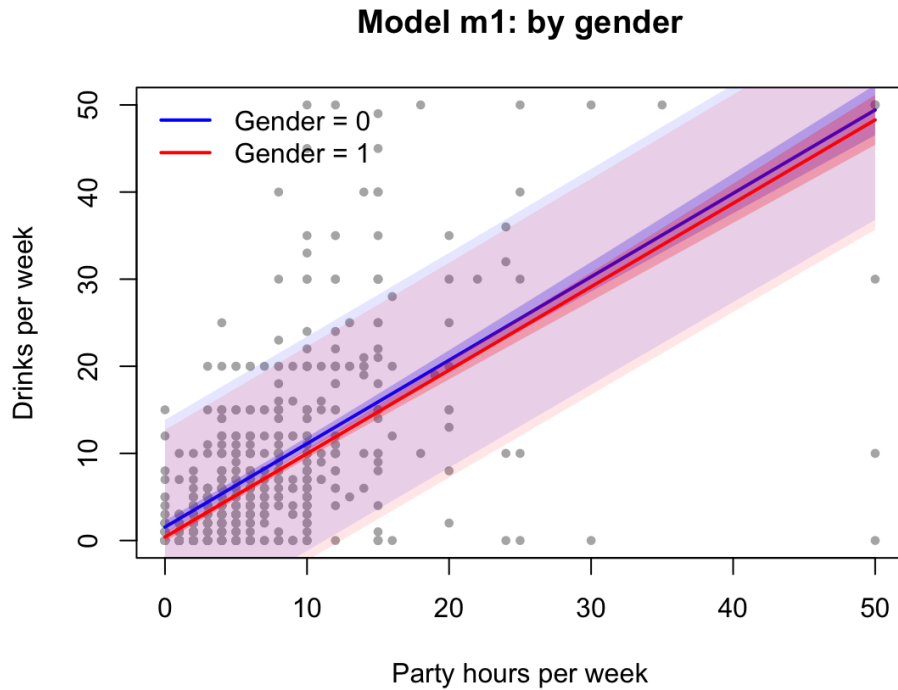


Рисунок 1 — Регрессия m1: зависимость `drinks_per_week` от `party_hours_per_week` и `gender`. Доверительные (шире) и предиктивные интервалы для каждого пола.

Визуализация модели m2

На втором графике представлены результаты второй модели, которая учитывает и пол, и тип населённого пункта. На одной оси X — часы на вечеринках, на оси Y — напитки в неделю. Для каждой из четырёх комбинаций "пол × `home_bc`" проведена отдельная линия.

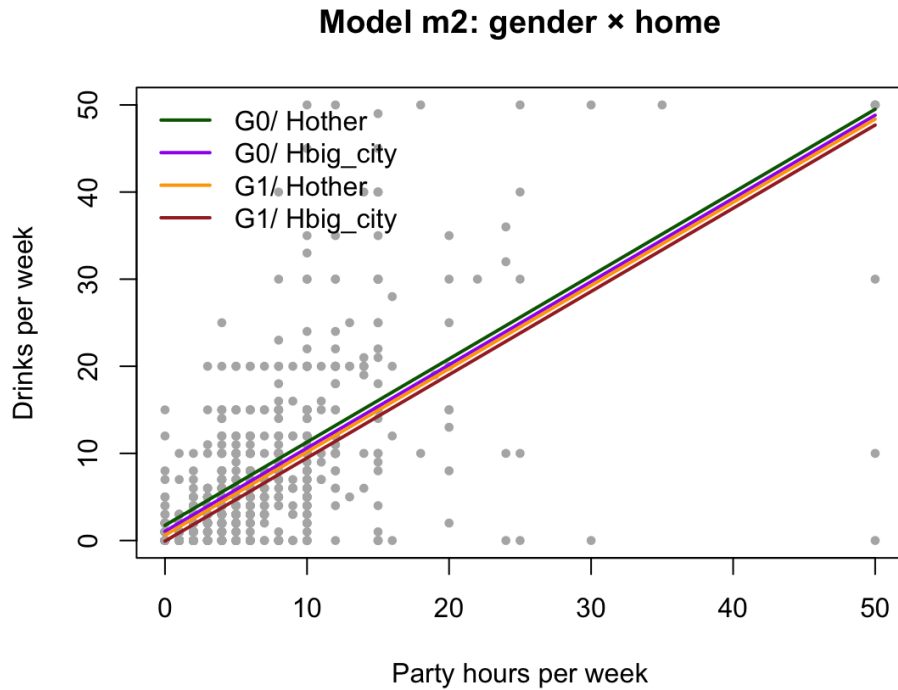


Рисунок 2 — Регрессия m2: зависимость `drinks_per_week` от `party_hours_per_week`, `gender` и `home_bc`.

Как видно из графиков, во всех моделях основной вклад вносит переменная `party_hours_per_week`, а разница между группами по полу и типу города визуально минимальна — линии расположены очень близко друг к другу.

Коэффициенты моделей

В результате оценки первой модели (`m1`), в которую входят переменные `party_hours_per_week` и `gender`, получены следующие значения коэффициентов:

Таблица 5 — Коэффициенты линейной модели `m1`

Переменная	Оценка	Std. Error	t value	p-value
(Intercept)	1.54	0.39	3.97	< 0.001
<code>party_hours_per_week</code>	0.96	0.03	29.99	< 0.001
<code>gender1</code>	−1.14	0.42	−2.74	0.006

Добавление в модель m2 переменной home_bc практически не изменяет оценки других коэффициентов. Эффект принадлежности к группе big_city оказывается статистически незначимым:

Таблица 6 — Коэффициенты линейной модели m2

Переменная	Оценка	Std. Error	t value	p-value
(Intercept)	1.74	0.41	4.27	< 0.001
party_hours_per_week	0.96	0.03	29.91	< 0.001
gender1	−1.13	0.42	−2.71	0.007
home_bcbig_city	−0.68	0.43	−1.60	0.109

Во всех моделях наибольший вклад вносит переменная party_hours_per_week: увеличение на 1 час связано практически с одним дополнительным напитком в неделю. Эффект пола также значим — при gender = 1 количество напитков в среднем меньше на 1.13–1.14, чем при gender = 0. Переменная home_bc (big city) статистически незначима: её коэффициент отрицательный, но уровень значимости превышает 0.05. Доля объяснённой дисперсии (R^2) для обеих моделей составляет около 46%, что подтверждает — основная тенденция хорошо описывается, однако остаётся значительный индивидуальный разброс в данных.

3 Прогнозирование по модели m1 (доверительный и предиктивный интервалы)

Для анализа практического применения модели m1 было выполнено построение доверительных и предиктивных интервалов для ожидаемого значения и для отдельных новых наблюдений соответственно. Это позволяет оценить, насколько устойчивы прогнозы и каков возможный разброс индивидуальных значений.

Доверительный интервал (confidence interval) показывает диапазон, в котором с заданной вероятностью находится истинное среднее значение отклика для заданных значений факторов. Предиктивный интервал (prediction interval) шире и показывает диапазон, в котором с этой же вероятностью окажется значение отклика для нового индивидуального наблюдения с заданными характеристиками.

Вычисление и визуализация интервалов выполнялись для обеих групп по полу, на всём диапазоне переменной party_hours_per_week. В качестве нового "X" бралась равномерная сетка значений, покрывающая весь диапазон наблюдений.

Пример кода для вычисления интервалов:

```
1 xseq <- seq(min(df_work$party_hours_per_week, na.rm = TRUE),
2             max(df_work$party_hours_per_week, na.rm = TRUE), length.out
3             = 120)
4 # Для группы gender = 0:
5 nd1 <- data.frame(party_hours_per_week = xseq,
6                   gender = factor(rep(levels(df_work$gender)[1], length(
7                                   xseq))), levels = levels(df_work$gender)))
8 pred_ci1 <- predict(m1, newdata = nd1, interval = "confidence")
9 pred_pi1 <- predict(m1, newdata = nd1, interval = "prediction")
10 # Для gender = 1 аналогично.
```

На графике (см. рис. 1) эти интервалы представлены полупрозрачными областями: доверительный интервал (уже), предиктивный (шире), отдельно для каждой группы по полу. Центральная линия — это прогноз модели для среднего значения отклика.

Практическое значение:

— Доверительный интервал используется для оценки того, насколько точно мы можем предсказать среднее значение потребления алкоголя среди студентов с определёнными характеристиками.

— Предиктивный интервал показывает, какой разброс можно ожидать для индивидуального студента при тех же характеристиках (например, если взять случайного студента с заданным числом часов вечеринок и полом).

В данном случае оба интервала довольно узкие при средних значениях и расширяются на краях диапазона, что типично для линейной регрессии. Разница между полами (gender) остаётся минимальной во всём диапазоне.

Заключение

В ходе данной работы была реализована последовательная обработка и анализ данных опроса студентов Университета Иллинойса с целью изучить взаимосвязь между потреблением алкогольных напитков, временем, проводимым на вечеринках, полом и типом населённого пункта проживания.

В соответствии с поставленными целями:

- Проведена очистка и структурирование исходного файла данных, в том числе перекодировка и сокращение переменных, что позволило получить качественный и пригодный для анализа датафрейм без пропусков и аномалий.

- Построены две регрессионные модели: первая учитывала только число часов на вечеринках и пол, вторая — дополнительно тип населённого пункта.

- Выполнена визуализация результатов, что позволило наглядно сравнить влияние разных факторов.

- Проведено прогнозирование с использованием доверительных и предиктивных интервалов, что дало возможность оценить точность модели и возможный разброс индивидуальных значений.

Основные результаты анализа:

- Главным фактором, определяющим уровень потребления алкоголя среди студентов, является количество времени, проведённого на вечеринках. В среднем, каждый дополнительный час на вечеринке увеличивает ожидаемое количество напитков на один.

- Пол и тип населённого пункта оказывают минимальное влияние: различия между группами незначимы как статистически, так и визуально.

- Модель даёт достаточно узкие доверительные интервалы для среднего значения, однако предиктивные интервалы (для индивидуальных прогнозов) шире, что отражает высокую вариабельность индивидуального поведения студентов.

Таким образом, поставленные во введении цели были полностью достигнуты: была проанализирована роль социально-демографических и поведенческих факторов в формировании потребления алкоголя среди студентов, и подтверждено, что именно стиль досуга (а не пол или город) является основным предиктором. Результаты работы могут быть полезны для разработки профилактических и образовательных программ, направленных на снижение уровня потребления алкоголя среди студенческой молодёжи.

Список использованных источников

1. Официальная страница среды статистического моделирования R. — 2024. — Режим доступа: [R: The R Project for Statistical Computing](#) (дата обращения: 21.01.2025).
2. Официальная страница интегрированной среды разработки RStudio. — 2024. — Режим доступа: [RStudio | Open source & professional software for data science teams - RStudio](#) (дата обращения: 21.01.2025).
3. *Wickham, Hadley*. ggplot2: Elegant Graphics for Data Analysis (Use R!) / Hadley Wickham. — New York: Springer, 2009. — Режим доступа: [ggplot2](#) (дата обращения: 21.01.2025).
4. *Max Kuhn*. The caret Package. — 2019. — Режим доступа: [The caret Package](#) (дата обращения: 21.01.2025).
5. *Alexandros Karatzoglou*. Kernel-Based Machine Learning Lab. — 2024. — Режим доступа: [Kernel-Based Machine Learning Lab](#) (дата обращения: 21.01.2025).

Приложение А Полный код на R

```
1 setwd("/Users/daniil/Desktop/Обучение/BigData")
2
3 DATA_FILE <- "lab2_data.dat"
4
5 raw_header <- readLines(DATA_FILE, n = 1)
6 clean_header <- gsub("\\s*\\([~()]*\\)", "", raw_header)
7 clean_header <- gsub("\\s+", " ", clean_header)
8 clean_header <- trimws(clean_header)
9
10
11 header_tokens <- strsplit(clean_header, " ")[[1]]
12 col_names_raw <- header_tokens[-1]
13 col_names_raw <- col_names_raw[-2]
14 col_names_raw
15 length(col_names_raw)
16
17 make_snake <- function(x) {
18   x <- tolower(gsub("[^A-Za-z0-9]+", "_", x))
19   # ensure uniqueness (append _n where needed)
20   dup <- duplicated(x)
21   if (any(dup)) x[dup] <- paste0(x[dup], "_", seq_along(which(dup)))
22   x
23 }
24 col_names <- make_snake(col_names_raw)
25 length(col_names)
26
27
28 df <- read.table(
29   DATA_FILE,
30   header      = FALSE,
31   sep         = " ",
32   stringsAsFactors = FALSE,
33   skip        = 1,
34   comment.char = " ",
35   fill        = TRUE
36 )
37
38 head(df)
39 df <- df[, -c(1:2)]
40 df <- df[, -2]
41 head(df)
42 dim(df)
43
44
45 names(df) <- col_names
```

```

46 head(df)
47
48
49 factor_cols <- c(
50   "gender", "greek", "home_town", "in_state", "ethnicity", "religion",
51   "calculus", "cell_phone", "president", "political_party", "section"
52 )
53
54
55 factor_cols <- intersect(factor_cols, names(df))
56
57 df[factor_cols] <- lapply(df[factor_cols], factor)
58
59 numeric_mask <- sapply(df, is.numeric)
60 problem_mat <- df[, numeric_mask, drop = FALSE] < 0 | df[, numeric_mask
61   , drop = FALSE] == 999
62
63 rows_bad <- rowSums(problem_mat, na.rm = TRUE) > 0
64
65 df <- df[!rows_bad, ]
66
67
68 cat("\n\n--- Проверка уровней факторов ---\n")
69 print(lapply(df[factor_cols], levels))
70
71
72 cat("\n--- Кол-во пропущенных значений ---\n")
73 print(colSums(is.na(df)))
74
75
76 df$home_bc <- factor(ifelse(df$home_town == "3", 1, 0),
77   levels = c(0, 1),
78   labels = c("other", "big_city"))
79
80 head(df)
81 df_work <- df[, -c(2:11)]
82 head(df_work)
83 df_work <- df_work[, -c(4:22)]
84 head(df_work)
85
86
87 m1 <- lm(drinks_per_week ~ party_hours_per_week + gender, data = df_work
88   )
89 summary(m1)
90 m2 <- lm(drinks_per_week ~ party_hours_per_week + gender + home_bc, data
91   = df_work)
92 summary(m2)
93
94
95 xseq <- seq(min(df_work$party_hours_per_week, na.rm = TRUE),

```

```

89         max(df_work$party_hours_per_week, na.rm = TRUE), length.out
          = 120)
90
91
92 col_gender <- c("blue", "red")
93 names(col_gender) <- levels(df_work$gender)[1:2]
94
95 col_combo <- c("darkgreen", "purple", "orange", "brown")
96 combo_labels <- NULL
97
98
99 plot(df_work$party_hours_per_week, df_work$drinks_per_week,
100      xlab = "Party hours per week", ylab = "Drinks per week",
101      main = "Model m1: by gender", pch = 19, cex = 0.6, col = gray(0.7))
102
103 g1 <- levels(df_work$gender)[1]
104 g2 <- levels(df_work$gender)[2]
105
106
107 nd1 <- data.frame(party_hours_per_week = xseq,
108                  gender = factor(rep(g1, length(xseq)), levels = levels
109                                (df_work$gender)))
109 pred_ci1 <- predict(m1, newdata = nd1, interval = "confidence")
110 pred_pi1 <- predict(m1, newdata = nd1, interval = "prediction")
111 base_col1 <- col_gender[g1]
112 ci_col1 <- adjustcolor(base_col1, alpha.f = 0.25)
113 pi_col1 <- adjustcolor(base_col1, alpha.f = 0.10)
114
115 polygon(c(xseq, rev(xseq)), c(pred_pi1[, "lwr"], rev(pred_pi1[, "upr"])),
116        , col = pi_col1, border = NA)
117 polygon(c(xseq, rev(xseq)), c(pred_ci1[, "lwr"], rev(pred_ci1[, "upr"])),
118        , col = ci_col1, border = NA)
119 lines(xseq, pred_ci1[, "fit"], col = base_col1, lwd = 2)
120
121 nd2 <- data.frame(party_hours_per_week = xseq,
122                  gender = factor(rep(g2, length(xseq)), levels = levels
123                                (df_work$gender)))
124 pred_ci2 <- predict(m1, newdata = nd2, interval = "confidence")
125 pred_pi2 <- predict(m1, newdata = nd2, interval = "prediction")
126 base_col2 <- col_gender[g2]
127 ci_col2 <- adjustcolor(base_col2, alpha.f = 0.25)
128 pi_col2 <- adjustcolor(base_col2, alpha.f = 0.10)
129
130 polygon(c(xseq, rev(xseq)), c(pred_pi2[, "lwr"], rev(pred_pi2[, "upr"])),
131        , col = pi_col2, border = NA)

```

```

128 polygon(c(xseq, rev(xseq)), c(pred_ci2[, "lwr"], rev(pred_ci2[, "upr"])))
      , col = ci_col2, border = NA)
129 lines(xseq, pred_ci2[, "fit"], col = base_col2, lwd = 2)
130
131 legend("topleft", legend = paste("Gender =", c(g1, g2)), lwd = 2, col =
      c(base_col1, base_col2), bty = "n")
132
133
134 plot(df_work$party_hours_per_week, df_work$drinks_per_week,
135       xlab = "Party hours per week", ylab = "Drinks per week",
136       main = "Model m2: gender      home", pch = 19, cex = 0.6, col = gray
      (0.7))
137
138
139 h1 <- levels(df_work$home_bc)[1]
140 h2 <- levels(df_work$home_bc)[2]
141 combo_cols <- c("darkgreen", "purple", "orange", "brown")
142
143 nd11 <- data.frame(party_hours_per_week = xseq,
144                   gender = factor(rep(g1, length(xseq)), levels =
      levels(df_work$gender)),
145                   home_bc = factor(rep(h1, length(xseq)), levels =
      levels(df_work$home_bc)))
146 lines(xseq, predict(m2, newdata = nd11), col = combo_cols[1], lwd = 2)
147
148
149 nd12 <- nd11
150 nd12$home_bc <- factor(rep(h2, length(xseq)), levels = levels(
      df_work$home_bc))
151 lines(xseq, predict(m2, newdata = nd12), col = combo_cols[2], lwd = 2)
152
153
154 nd21 <- nd11
155 nd21$gender <- factor(rep(g2, length(xseq)), levels = levels(
      df_work$gender))
156 nd21$home_bc <- factor(rep(h1, length(xseq)), levels = levels(
      df_work$home_bc))
157 lines(xseq, predict(m2, newdata = nd21), col = combo_cols[3], lwd = 2)
158
159
160 nd22 <- nd21
161 nd22$home_bc <- factor(rep(h2, length(xseq)), levels = levels(
      df_work$home_bc))
162 lines(xseq, predict(m2, newdata = nd22), col = combo_cols[4], lwd = 2)
163
164 legend("topleft",

```

```
165     legend = c(  
166         paste("G", g1, "/ H", h1, sep = ""),  
167         paste("G", g1, "/ H", h2, sep = ""),  
168         paste("G", g2, "/ H", h1, sep = ""),  
169         paste("G", g2, "/ H", h2, sep = "")  
170     ),  
171     lwd = 2, col = combo_cols, bty = "n")
```