



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

**ИНСТИТУТ МАТЕМАТИКИ И КОМПЬЮТЕРНЫХ
ТЕХНОЛОГИЙ**

Департамент математического и компьютерного
моделирования

Курс «Компьютерные методы анализа больших данных»

Лабораторная работа №2
Контрольное мероприятие по рейтингу

на тему «Очистка и сравнительный анализ интернет-данных»
вариант № 3

Выполнил студент Б9122-01.03.02мкт
Вершинин Д.А.

Проверил доцент Достовалов В.Н.

г. Владивосток
2025

Содержание

Введение	3
1 Подготовка данных	4
2 Формирование и анализ моделей	8
2.1 Создание моделей	8
2.2 Оценка моделей	9
2.3 Графики моделей	11
2.4 Доверительные и предиктивные интервалы	11
Заключение	15
Список использованных источников	16
А Код графика для регрессионных моделей m1 и m2	17
Б Код предсказаний модели m1	19

Введение

Анализ и обработка данных в современном мире занимают центральное место в сфере информационных технологий. Современные методы анализа позволяют выявлять закономерности и зависимости между различными характеристиками в данных, что делает их крайне востребованными в научных и прикладных задачах.

Цель данной работы — продемонстрировать применение методов линейной регрессии на реальных данных. В качестве объекта исследования выбран открытый датасет `Stat100_Fall2018_Survey02`, предоставленный Университетом Иллинойса. Этот набор данных содержит информацию, собранную в ходе учебного опроса, и хорошо подходит для целей демонстрационного анализа.

В рамках работы будут реализованы две модели линейной регрессии: базовая и расширенная. Сравнение этих моделей позволит выявить, насколько добавление дополнительных признаков влияет на качество предсказания.

Для достижения поставленной цели решаются следующие задачи:

- а) Загрузка и предварительная обработка датасета.
- б) Построение двух моделей линейной регрессии — базовой и расширенной.
- в) Визуализация результатов моделирования для наглядной интерпретации.
- г) Сравнительный анализ моделей по ключевым метрикам качества и интерпретация полученных результатов.

В качестве инструмента анализа используется язык `R`[1]. Для визуализации данных - библиотека `ggplot2`[2]. Разделение данных произведено с помощью библиотеки `caret`[3].

1 Подготовка данных

Для начала работы требуется загрузить найденный датасет. Однако заголовки, заданные по умолчанию, содержат поясняющую информацию, поэтому требуется отдельно загрузить данные и их заголовки. Заранее все заголовки были перенесены в файл `headers.txt`. Загрузка основной информации производится с помощью кода, представленного в листинге 1.1.

Листинг 1.1 — Считывание датасета

```
1 setwd('/home/ndavs/study/math/DS/r/5_lab_Stats_100')
2
3 data <- read.table('data.dat', sep=' ')
4 dim(data)
5 head(data)
```

Данные о размерности: (1458, 35). Первые пять записей представлены в таблице 1.1.

Таблица 1.1 — Сводная таблица наблюдений: первые и последние 5 переменных (остальные обозначены как ...)

Obs	V1	V2	V3	V4	V5	...	V31	V32	V33	V34	V35
Obs1	1	1	1	0	2	...	0	100	7	1	NA
Obs2	2	1	1	0	3	...	0	100	9	1	NA
Obs3	3	1	1	1	2	...	0	100	3	1	NA
Obs4	4	1	1	0	2	...	15	0	1	1	NA
Obs5	5	0	0	1	2	...	0	50	5	1	NA

Для обработки заголовков был использован следующий подход (листинг 1.2):

- считывание строки с заголовками;
- использование регулярного выражения для удаления скобок и значений в них;
- удаление лишних точек и пробелов;
- разделение значений на элементы массива.

Листинг 1.2 — Работа с заголовками

```
1 base_headers <- readLines('headers.txt', n=1)
2 cleaned_header <- gsub("\\s*\\([^(]*)\\s*", "", base_headers)
3 cleaned_header <- gsub("\\.", "", cleaned_header)
```

```

4 cleaned_header <- gsub("\\s+", " ", cleaned_header)
5 cleaned_header <- trimws(cleaned_header)
6 cleaned_header <- unlist(strsplit(cleaned_header, " "))
7 length(cleaned_header)

```

Длина полученного массива: 34 (на один меньше, чем количество столбцов датафрейма, так как последнее значение — NULL).

После проделанных шагов значения заголовков добавляются к датасету (листинг 1.3).

Листинг 1.3 — Применение заголовков

```

1 colnames(data) <- cleaned_header
2 dim(data)
3 head(data)

```

Размерность не изменилась: (1458, 35). Первые пять строк полученного датасета представлены в таблице 1.2.

Таблица 1.2 — Социально-академические характеристики студентов: только первые и последние 4 переменные

№	Gender	Gender_ID	Greek	Home_Town	...	Work_Hours	Tuition	Career	Section
1	1	1	1	0	...	0	100	7	1
2	2	1	1	0	...	0	100	9	1
3	3	1	1	1	...	0	100	3	1
4	4	1	1	0	...	15	0	1	1
5	5	0	0	1	...	0	50	5	1
6	6	1	1	0	...	0	90	5	1

Далее подготовим только нужные значения из датафрейма (листинг 1.4).

Листинг 1.4 — Сформируем датафрейм нужных признаков

```

1 dd <- data[, c("Drinks_per_week", "Party_Hours_per_week", "Gender",
2               "Home_Town")]
3 head(dd)
4 dim(dd)
5 colSums(is.na(dd))
6 summary(dd)

```

Пропущенных значение не обнаружено. Размерность полученного датафрейма: (1458, 4).

Описание данных:

а) Пол (Gender):

— 0 = Мужской

— 1 = Женский

б) **Родной город (Home_Town)**: тип населённого пункта:

— 0 = Маленький город

— 1 = Средний город

— 2 = Большой город (пригород)

— 3 = Большой город (без пригородов)

в) **Часы вечеринок в неделю (Party_Hours_per_week)**:

среднее количество часов, проведённых на вечеринках в неделю.

г) **Алкогольные напитки в неделю (Drinks_per_week)**:

среднее количество алкогольных напитков, потребляемых за неделю.

Преобразуем параметр города следующим образом: 0 = Small Town + Medium City, 1 = Big City. Также преобразуем пол в факторную переменную (листинг 1.5).

Листинг 1.5 — Преобразование нужных признаков

```
1 dd$Home_Town[dd$Home_Town < 2] <- 0
2 dd$Home_Town[dd$Home_Town >= 2] <- 1
3 dd$Gender <- factor(dd$Gender)
```

Первые пять строк полученного датасета представлены в таблице 1.3.

Таблица 1.3 — Параметры студентов, связанные с поведением и происхождением

№	Drinks_per_week	Party_Hours_per_week	Gender	Home_Town
1	1	1	1	1
2	0	0	1	1
3	10	6	1	1
4	35	20	1	1
5	25	12	0	1
6	3	8	1	1

Размерность полученного датафрейма: (1458, 4).

Сводная информация по полученному датафрейму представлена в таблице 1.4.

Таблица 1.4 — Описательная статистика по ключевым переменным

Показатель	Drinks_per_week	Party_Hs_per_wk	Gender	Home_Tn
Min	0.000	0.000	0.0000	0.000
1st Quartile	0.000	1.000	0.0000	1.000
Median	3.000	4.000	1.0000	2.000
Mean	6.475	5.496	0.6454	1.853
3rd Quartile	10.000	8.000	1.0000	3.000
Max	50.000	50.000	1.0000	3.000

Разделим полученный датафрейм на обучающую и тестовую выборки (70:30) (листинг 1.6).

Листинг 1.6 — Разделение на выборки

```
1 library(caret)
2
3 set.seed(2004)
4
5 train_index <- caret::createDataPartition(dd$Drinks_per_week, p = 0.7,
      list = FALSE)
6 train_data <- dd[train_index, ]
7 test_data <- dd[-train_index, ]
8
9 dim(train_data)
10 dim(test_data)
```

Размерность тренировочной выборки - (1022, 4).

Тестовой - (436, 4).

Суммарно строк - 1458 (ничего не потеряно).

На этом подготовка данных завершена.

2 Формирование и анализ моделей

2.1 Создание моделей

На основе обучающей выборки, были построены две модели линейной регрессии[4]: базовая (Drinks_per_week от Party_Hours_per_week + Gender) и расширенная (Drinks_per_week от Party_Hours_per_week + Gender + Home_Town). Программный код представлен в листинге 2.1.

Листинг 2.1 — Создание моделей в R

```
1 m1 <- lm("Drinks_per_week ~ Party_Hours_per_week + Gender",
  data=train_data)
2 summary(m1)
3
4 m2 <- lm("Drinks_per_week ~ Party_Hours_per_week + Gender +
  Home_Town", data=train_data)
5 summary(m2)
```

Таблица 2.1 — Сводная информация по первой модели: Drinks_per_week ~ Party_Hours_per_week + Gender

Коэффициент	Оценка	Ст. ошибка	t-стат.	p-значение
(Intercept)	1.63	0.37	4.39	1.27e-06***
Party_Hours_per_week	1.14	0.03	34.31	<2e-16***
Gender1	-2.25	0.39	-5.73	1.34e-08***

Остатки:

Минимум: -46.40 1 кв.: -2.78 Медиана: -0.511 3 кв.: 1.19

Максимум: 43.36

Качество модели:

Стандартная ошибка остатков: 5.98 при 1019 степенях свободы

Множественный R^2 : 0.54 Скорректированный R^2 : 0.55

F-статистика: 614.3 при 2 и 1019 СС p-значение: <2.2e-16

Сводная информация по построенным моделям представлена в таблицах 2.1–2.2. Как видно, включение переменной Home_Town не является статистически значимым ($p = 0.39$). Однако изменилось

Таблица 2.2 — Сводная информация по второй модели: $\text{Drinks_per_week} \sim \text{Party_Hours_per_week} + \text{Gender} + \text{Home_Town}$

Коэффициент	Оценка	Ст. ошибка	t-стат.	p-значение
(Intercept)	1.85	0.45	4.09	4.56e-05***
Party_Hours_per_week	1.15	0.03	34.28	<2e-16***
Gender1	-2.23	0.40	-5.63	2.38e-08***
Home_Town	-0.36	0.41	-0.86	0.39

Остатки:

Минимум: -46.35 1 кв.: -2.65 Медиана: -0.420 3 кв.: 1.25

Максимум: 48.50

Качество модели:

Стандартная ошибка остатков: 5.97 при 1018 степенях свободы

Множественный R^2 : 0.55 Скорректированный R^2 : 0.55

F-статистика: 409.7 при 3 и 1018 СС p-значение: <2.2e-16

значение свободного члена (Intercept), что указывает на возможное влияние города происхождения на базовый уровень потребления.

2.2 Оценка моделей

Проведём оценку обеих моделей на всей выборке (листинг 2.2).

Листинг 2.2 — Вычисление F-статистики

```

1 predictions <- predict(m1, newdata = dd)
2
3 residuals <- dd$Drinks_per_week - predictions
4 RSS <- sum(residuals^2)
5
6 TSS <- sum((dd$Drinks_per_week - mean(dd$Drinks_per_week))^2)
7
8 SSR <- TSS - RSS
9
10 k <- length(coef(m1)) - 1
11 n <- nrow(dd)
12
13 F_statistic <- (SSR / k) / (RSS / (n - k - 1))
14 F_statistic

```

Для первой модели значение F-статистики [4] равно 865.1868. Для второй - 866.3685.

Значения F-статистики намного больше 1, что указывает на то, что обе модели в целом значимы. Вторая модель имеет чуть более высокую F-статистику (866.3685 против 865.1868). Но не стоит говорить о её лучшей объясняющей способности.

На тестовой выборке оценим модели с помощью средне-квадратичного отклонения[4] (листинг 2.3).

Листинг 2.3 — Вычисление MSE

```
1 predictions_m1 <- predict(m1, newdata=test_data)
2 predictions_m2 <- predict(m2, newdata=test_data)
3 mse_m1 <- mean((test_data$Drinks_per_week - predictions_m1)^2)
4 mse_m2 <- mean((test_data$Drinks_per_week - predictions_m2)^2)
5
6 cat("MSE for Model 1 (m1):", mse_m1, "\n")
7 cat("MSE for Model 2 (m2):", mse_m2)
```

Для первой модели получено значение 41.61. Для второй - 41.58. Вторая модель чуть лучше описывает тестовые данные.

2.3 Графики моделей

Построим графики обеих моделей (рисунок 2.1, код представлен в приложении А).

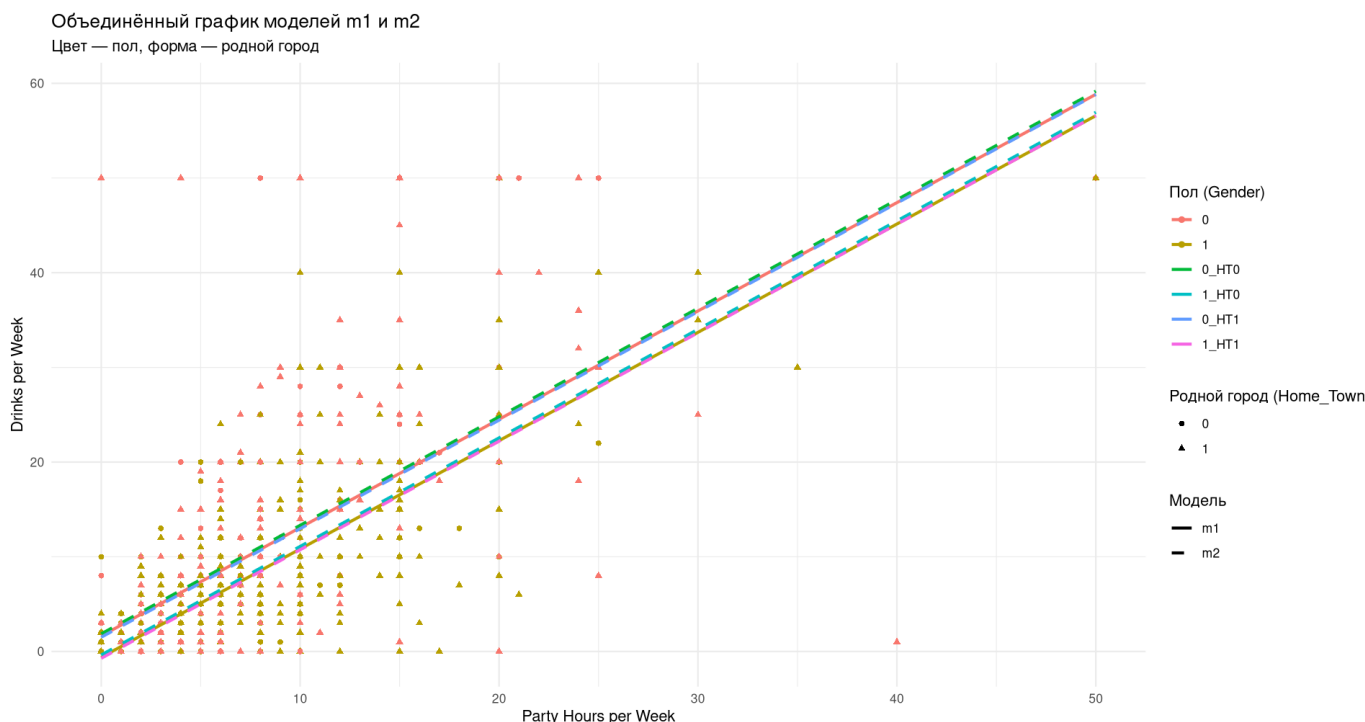


Рисунок 2.1 — Графики регрессионных моделей

На рисунке 2.1 видно, что обе модели дают параллельные прямые для каждого пола. Это объясняется одинаковым коэффициентом переменной `Party_Hours_per_week`, который отражает вклад вечернего поведения вне зависимости от пола.

2.4 Доверительные и предиктивные интервалы

Расчёт новых данных

Рассчитаем предсказания для значений переменной `Party_Hours_per_week`, увеличенных на 5%, 10% и 15% от её максимума, для обоих полов, используя первую модель (таб. 2.3, листинг 2.4).

Листинг 2.4 — Предсказания новых значений

```
1 x_max <- max(dd$Party_Hours_per_week)
2 x_ext <- x_max * c(1.05, 1.10, 1.15)
3
4 new_rows <- expand.grid(
```

```

5 Party_Hours_per_week = x_ext,
6 Gender = levels(train_data$Gender)
7 )
8
9 new_rows <- expand.grid(Party_Hours_per_week = x_ext, Gender =
  levels(train_data$Gender))
10
11 new_rows$Drinks_per_week <- predict(m1, newdata = new_rows)
12 new_rows$Home_Town <- NA
13
14 conf <- predict(m1, newdata = new_rows, interval = "confidence")
15 pred <- predict(m1, newdata = new_rows, interval = "prediction")
16
17 results <- cbind(new_rows, conf, pred_Lwr = pred[, "lwr"], pred_Upr =
  pred[, "upr"])
18 print(results)

```

Таблица 2.3 — Предсказания и интервалы для переменных

Party_Hours	Gender	Fit	Lwr	Upr	Pred_Lwr	Pred_Upr
52.5	0	61.70	58.58	64.82	49.55	73.86
55.0	0	64.57	61.28	67.85	52.37	76.76
57.5	0	67.43	63.98	70.87	55.19	79.66
52.5	1	59.45	56.32	62.57	47.29	71.60
55.0	1	62.31	59.02	65.59	50.11	74.50
57.5	1	65.17	61.72	68.61	52.93	77.41

Построим доверительные и предиктивные интервалы^[4] на всём промежутке (рисунок 2.2, код представлен в приложении Б).

Доверительные интервалы (Confidence)

Доверительные интервалы отображают неопределённость оценки среднего значения зависимой переменной при фиксированном значении независимой переменной. Видно, что интервалы более узкие в центральной части, где плотность наблюдений выше, и расширяются к краям, что типично при экстраполяции.



Рисунок 2.2 — Графики доверительных и предиктивных интервалов модели (на фоне точки всей выборки. Справа сверху шесть точек - новые предсказание +5-15% от максимального значения всего датасета)

Предиктивные интервалы (Prediction)

Предиктивные интервалы шире, так как учитывают не только ошибку оценки среднего, но и вариацию индивидуальных значений. Это позволяет:

- Адекватно отображать разброс реальных наблюдений вокруг линии регрессии.
- Учесть существенно большую неопределённость при попытке предсказать индивидуальное поведение.

Область экстраполяции

Интервалы построены с учётом расширения по оси `Party_Hours_per_week` на 15%, что наглядно демонстрирует:

- Резкое расширение обоих типов интервалов за пределами основной области наблюдений.

— Повышенную неопределённость прогнозов вне обучающего диапазона данных.

Выводы

— Модель демонстрирует адекватную уверенность в оценке среднего, но справедливо расширяет предиктивные интервалы, что указывает на реалистичную оценку неопределённости.

— Наблюдаемое различие между полами при сохранении общей тенденции подтверждает обоснованность стратификации по полу.

Заключение

Таким образом, была продемонстрирована практическая реализация методов линейной регрессии. Основная цель исследования была успешно достигнута, что подтверждается выполнением всех поставленных задач.

Был предобработан датасет: преобразован показатели городов и пола, выбраны признаки для анализа.

Построены две модели и проведён анализ с использованием F-статистики и MSE.

Также были созданы и проанализированы графики, которые иллюстрируют поведение моделей.

Список использованных источников

1. *Мэтлофф, Норман*. Искусство программирования на R. Погружение в большие данные / Норман Мэтлофф. Библиотека программиста. — СПб.: Питер, 2019. — С. 416.
2. *Мастицкий, Сергей*. Визуализация данных с помощью ggplot2 / Сергей Мастицкий. — Москва: ДМК Пресс, 2017. — С. 222.
3. *Kuhn, Max*. Caret: Classification and Regression Training. — R package version 6.0-86. — 2019. URL: <https://cran.r-project.org/package=caret/vignettes/caret.html>.
4. *Wasserman, Larry*. All of Statistics: A Concise Course in Statistical Inference / Larry Wasserman. Springer Texts in Statistics. — Springer, 2004.

Приложение А Код графика для регрессионных моделей m1 и m2

Листинг А.1 — Код для формирования графика обеих моделей

```
1 library(ggplot2)
2
3 x_seq <- seq(
4   from = min(dd$Party_Hours_per_week),
5   to = max(dd$Party_Hours_per_week),
6   length.out = 100
7 )
8
9 newdata_m1 <- expand.grid(
10   Party_Hours_per_week = x_seq,
11   Gender = levels(dd$Gender)
12 )
13 newdata_m1$Home_Town <- NA # Добавляем для совместимости
14 newdata_m1$Pred <- predict(m1, newdata_m1)
15 newdata_m1$Model <- "m1"
16 newdata_m1$Group <- newdata_m1$Gender
17
18 newdata_m2 <- expand.grid(
19   Party_Hours_per_week = x_seq,
20   Gender = levels(dd$Gender),
21   Home_Town = c(0, 1)
22 )
23 newdata_m2$Pred <- predict(m2, newdata_m2)
24 newdata_m2$Model <- "m2"
25 newdata_m2$Group <- interaction(newdata_m2$Gender,
26   newdata_m2$Home_Town, sep = "_HT")
27
28 pred_all <- rbind(newdata_m1, newdata_m2)
29 dd$Home_Town <- factor(dd$Home_Town)
30
31 ggplot(dd, aes(x = Party_Hours_per_week, y = Drinks_per_week)) +
32   geom_point(
33     aes(color = Gender, shape = Home_Town),
34     alpha = 1
35   ) +
36   geom_line(
37     data = pred_all,
38     aes(
39       x = Party_Hours_per_week,
40       y = Pred,
41       color = Group,
```

```
41     linetype = Model
42   ),
43   size = 1
44 ) +
45 scale_linetype_manual(
46   values = c(m1 = "solid", m2 = "dashed"),
47   name = "Модель"
48 ) +
49 labs(
50   title = "Объединённый график моделей m1 и m2",
51   subtitle = "Цвет - пол, форма - родной город",
52   x = "Party Hours per Week",
53   y = "Drinks per Week",
54   color = "Пол (Gender)",
55   shape = "Родной город (Home_Town)"
56 ) +
57 theme_minimal()
```

Приложение Б Код предсказаний модели m1

Листинг Б.1 — Код для формирования доверительного и предиктивного предсказания модели m1

```
1 range_x <- range(dd$Party_Hours_per_week)
2 x_min <- range_x[1]
3 x_max <- range_x[2]
4 x_range <- x_max - x_min
5
6 x_ext <- seq(
7   from = x_min,
8   to   = x_max + 0.15 * x_range,
9   length.out = 200
10 )
11
12 newdata <- expand.grid(
13   Party_Hours_per_week = x_ext,
14   Gender = levels(dd$Gender)
15 )
16
17 conf_pred <- predict(m1, newdata = newdata, interval = "confidence",
18   level = 0.95)
19 conf_df <- cbind(newdata, conf_pred)
20 conf_df$Type <- "Confidence"
21
22 pred_pred <- predict(m1, newdata = newdata, interval = "prediction",
23   level = 0.95)
24 pred_df <- cbind(newdata, pred_pred)
25 pred_df$Type <- "Prediction"
26
27 all_preds <- rbind(conf_df, pred_df)
28
29 ggplot() +
30   geom_point(data = dd, aes(x = Party_Hours_per_week, y =
31     Drinks_per_week, color = Gender), alpha = 1) +
32   geom_point(data = new_rows, aes(x = Party_Hours_per_week, y =
33     Drinks_per_week, color = Gender, ), alpha = 1) +
34
35   geom_line(data = all_preds, aes(x = Party_Hours_per_week, y = fit,
36     color = Gender), size = 1) +
37
38   geom_ribbon(
39     data = all_preds,
40     aes(x = Party_Hours_per_week, ymin = lwr, ymax = upr, fill =
41       Gender),
42     alpha = 0.2
```

```
37 ) +
38
39 facet_wrap(~Type) + # Разделение графика: confidence vs prediction
40
41 labs(
42   title = "Доверительный и предиктивный интервалы модели m1",
43   subtitle = "Прогноз с расширением области на 15% по
      Party_Hours_per_week",
44   x = "Party Hours per Week",
45   y = "Drinks per Week"
46 ) +
47 theme_minimal()
```