

«Методы статистического анализа, картографирования и имитационного моделирования»

~~Компьютерные технологии и статистические методы в
экологии и природопользовании~~

(геоэкологи, 98=13+13+82 часов)

1 курс магистратуры, весенний семестр 2016 г.

Преподаватель:

- Даниил Николаевич Козлов: daniilkozlov@landscape.edu.ru
- кафедра физической географии и ландшафтоведения

Информационная поддержка:

- <http://landscape.edu.ru>
- лекционные и практические материалы, задания, статьи, ссылки

Занятия:

- четверг 4 пара, 14:55-16:30, ауд. 2017

Задания:

- реферат статьи 2014-16 года из каталога ELSEVIER
- тематические задания

Проверка знаний:

- практические (70%), вопросы экзамена (20%), активная работа (10%)



ФГМ

Астапенко Юлия

Spatial distribution of nitrate health risk associated with groundwater use as drinking water in Merida, Mexico (2015)

11.02 ЦЕЛИ, ЗАДАЧИ И СОДЕРЖАНИЕ КУРСА
Экспертные и формальные модели. ДЗ

18.02 Проблемы статистического анализа данных в
25.02 экологии и природопользовании ДЗ
03.03

10.03 Проблемы **цифрового картографического моделирования**:
17.03 геостатистика и индикационное картографирование ДЗ

24.03 Проблемы **моделирования процессов самоорганизации** в
31.03 экологии и природопользовании ДЗ

07.04 Семинар по проблемам (доклады по статьям)
14.04

21.04 Резерв

R – свободная программная среда вычислений



The R Project for Statistical Computing



- язык программирования для статистической обработки данных и работы с графикой,
- свободная программная среда вычислений
- открытый исходный код под лицензией GNU GPL
- Росс Айхэк (англ. Ross Ihaka) и Роберт Джентлмен (англ. Robert Gentleman), статистический факультет Оклендского университета, Новая Зеландия
- R Foundation, Австрия
- <http://www.r-project.org/>

This server is hosted by the [Institute for Statistics and Mathematics](#) of [WU \(Wirtschaftsuniversität Wien\)](#).

About R
[What is R?](#)
[Contributors](#)
[Screenshots](#)
[What's new?](#)

Download, Packages
[CRAN](#)

R Project
[Foundation](#)
[Members & Donors](#)
[Mailing Lists](#)
[Bug Tracking](#)
[Developer Page](#)
[Conferences](#)
[Search](#)

Documentation
[Manuals](#)
[FAQs](#)
[The R Journal](#)
[Wiki](#)
[Books](#)
[Certification](#)
[Other](#)

Misc
[Bioconductor](#)
[Related Projects](#)
[User Groups](#)
[Links](#)

- R is free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS.
- To [download](#) or [purchase](#) R, please contact your preferred provider.
- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#).

News :

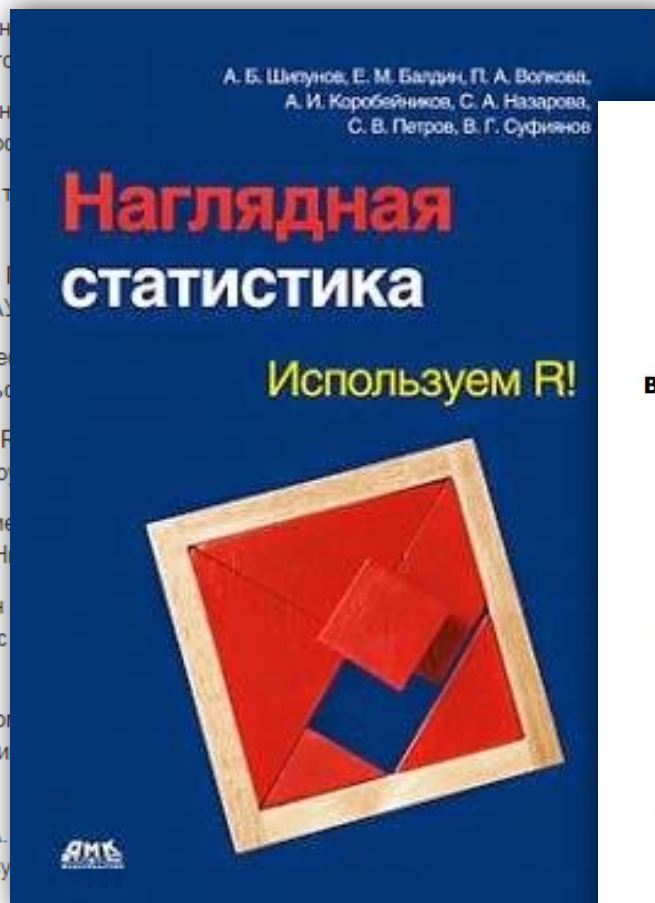
- R Foundation: The R Foundation has been established on 2014-10-31.
- [The R Journal Volume 6/1](#) is available.
- [useR! 2014](#) took place at the University of California, Los Angeles, USA June 30 - July 3, 2014.
- R version 3.0.3 (Warm Pudding) has been released on 2014-03-06.
- [R 3.0.3](#) will take place at the University of Copenhagen, Denmark, June 30 - July 3, 2015.

R – свободная программная среда вычислений

R: Анализ и визуализация данных

Книги и учебные пособия

- Мاستицкий С. Э., Шитиков В. К. (2014) Статистический анализ и визуализация данных с помощью R. - Электронная книга, 400 с. [PDF](#) | [Скрипты и данные](#)
- Савельев А. А., Мухарамова С. С., Пилюгин А. И. (2014) Методическое пособие. Казань: Казанский государственный университет, 120 с. [PDF](#)
- Савельев А. А., Мухарамова С. С., Пилюгин А. И. (2014) Методическое пособие по обработке данных. Учебно-методическое пособие. Казань: Казанский государственный университет, 120 с. [PDF](#)
- Волкова П. А., Шипунов А. Б. (2008) Статистический анализ данных в работах. М.: Экспресс, 60 с. [PDF](#)
- Буховец А. Г., Москалев П. В., Богатова В. И. (2010) Статистический анализ данных в системе R. Учебное пособие. Воронеж: ВГАУ, 120 с. [PDF](#)
- Зарядов И. С. (2010) Введение в статистический анализ информации, графика. М.: Издательство Московского государственного университета, 120 с. [PDF](#)
- Зарядов И. С. (2010) Статистический пакет R. М.: Издательство Российского университета дружбы народов, 120 с. [PDF](#)
- Зорин А. В., Федоткин М. А. (2010) Введение в статистический анализ данных. Методическое пособие. Нижний Новгород: Нижегородский государственный университет им. Г.И. Удальцова, 120 с. [PDF](#)
- Савельев А. А., Мухарамова С. С., Пилюгин А. И. (2014) Статистический анализ данных в экологии и природопользовании (с применением пакета R). Казань: Казанский государственный университет, 120 с. [PDF](#)
- Шитиков В. К., Розенберг Г. С. (2012) Рандомизированный статистический анализ данных по биологии. М.: Издательство Московского государственного университета, 120 с. [PDF](#)
- Шипунов А. Б., Балдин Е. М., Волкова П. А., Петров С. В., Суфиев Г. Г. (2012) Наглядная статистика. Используем R! М.: АМК, 120 с. [PDF](#)
- Савельев А. А., Мухарамова С. С., Чижикова Н. А., Пилюгин А. И. (2014) Теория пространственных точечных процессов в задачах экологии и природопользования (с применением пакета R). - Казань: Казанский государственный университет, 146 с. [PDF](#) | [Скрипты](#)

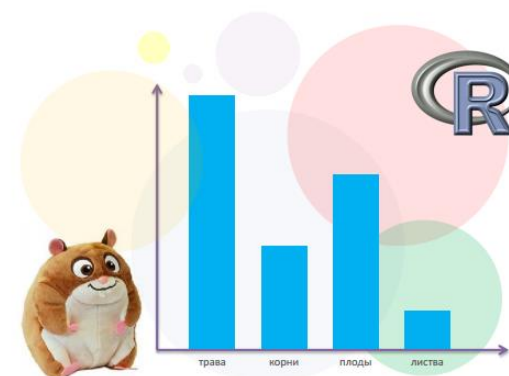


Поиск по блогу:

Поиск

С.Э. Мاستицкий, В.К. Шитиков

СТАТИСТИЧЕСКИЙ АНАЛИЗ И ВИЗУАЛИЗАЦИЯ ДАННЫХ С ПОМОЩЬЮ R



Хайделберг – Лондон – Тольятти
2014

А. Б. Шипунов, Е. М. Балдин, П. А. Волкова, А. И. Коробейников, С. А. Назарова, С. В. Петров, В. Г. Суфиянов

Наглядная статистика. Используем R!

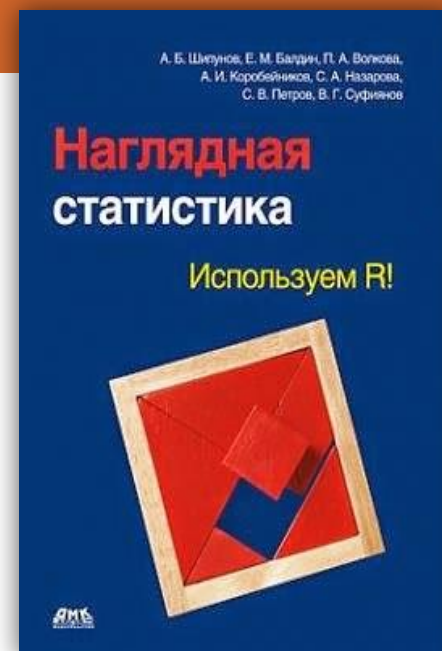
<http://ashipunov.info/shipunov/school/books/rbook.pdf>

ПРОЧИТАТЬ РАЗДЕЛЫ

Глава 2. Как обрабатывать данные

2.3. Из истории S и R	24
2.4. Применение, преимущества и недостатки R	25
2.5. Как скачать и установить R	27
2.6. Как начать работать в R	28
2.6.1. Запуск	28
2.6.2. Первые шаги	29

Приложение А. Пример работы в R, стр. 196-206



УСТАНОВИТЬ R И R-STUDIO

The screenshot displays the RStudio environment with the following components:

- Source Editor:** Contains R code for data manipulation and visualization. The code includes comments in Russian and performs operations like creating a new variable, plotting histograms, and density plots.
- Environment Pane:** Shows the current workspace with objects: data (1000 obs. of 3 variables), data.3 (100 obs. of 3 variables), datagen (1000 obs. of 2 variables), pedproph (150 obs. of 1 variable), and vars (1000 obs. of 2 variables).
- Console:** Shows the execution of the code, including several error messages: "Error: object 'h' not found".
- Plots Pane:** Displays a histogram of data\$var1 with blue bars and a density curve overlaid.

```
22 # str(data) показать структуру таблицы
23 # write.table(data, "data/zhuki_new.txt", quote=FALSE) сохранить таблицу
24 # summary(data) сводная информация по таблице
25 # summary(data$VES) или отдельной колонке
26 dev.new()
27
28 data$var2.dx <- log(data$var2+1000) #создание новой колонки
29 # гистограмма
30 hist(data$var1, breaks=50, main="", col = "red")
31 hist(data$var2, breaks=50, main="", col = "blue")
32 hist(data$var1, breaks=50, main="", col = "lightblue", border = "darkgreen")
33 hist(data$var2.dx, breaks=50, main="", col = "blue", add=TRUE)
34 plot(density(data$var1, adjust=2), main="", col = "red")
35 plot(density(data$var1, adjust=2), main="", col = "blue", lwd=3, add=TRUE)
36 rug(data$var1)
37
38 head(data)
39
40
41 # проверка на нормальность
42 qqnorm(data$var1)
43
```

Console output:

```
there were 12 warnings (use warnings() to see them)
> hist(data$var2, breaks=50, main="", col = "blue", add=TRUE)
> help(hist)
> hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
> h
Error: object 'h' not found
> h
Error: object 'h' not found
> hist(data$var1, breaks=50, main="", col = "red")
> hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
> hist(data$var2, breaks=50, main="", col = "blue", add=TRUE)
> hist(data$var1, breaks=50, main="", col = "red")
> #hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
```

ПРЕДЫДУЩЕЕ ДОМАШНЕЕ ЗАДАНИЕ

на основе темы магистерской диссертации

1. СФОРМУЛИРОВАТЬ ПРОБЛЕМУ И ЦЕЛЬ ИССЛЕДОВАНИЯ
2. СПРОЕКТИРОВАТЬ ТАБЛИЦУ ДЛЯ СТАТИСТИЧЕСКОГО АНАЛИЗА
3. ОБОСНОВАТЬ ШКАЛУ ИЗМЕРЕНИЙ КАЖДОГО СВОЙСТВА
4. ОБОСНОВАТЬ МЕТОДЫ СТАТИСТИЧЕСКОГО АНАЛИЗА
(сравнение выборок, линейная связь между переменными)
5. ЗАПОЛНИТЬ ТАБЛИЦУ в Excel
6. СОХРАНИТЬ в формате .csv (разделитель “;”)
7. ИМПОРТИРОВАТЬ ТАБЛИЦУ в R

```
11 # загрузка файла данных
12 data <- read.csv2("D:/Dropbox/R/Monte.csv", header = TRUE, sep = ";",
13                 quote = "\"", dec = ".", fill = TRUE,
14                 comment.char = "", stringsAsFactors = FALSE)
15
16 view(data) # показать таблицу данных
17 str(data) # показать структуру таблицы
18 head(data) # показать первые 6 строк таблицы
19 summary(data) # рассчитать мин, среднее, медиану, макс.
```

8. ОФОРМИТЬ ПРЕЗЕНТАЦИЮ и выслать ее преподавателю

ЛИТЕРАТУРА к заданию #1

А. Б. Шипунов, Е. М. Балдин, П. А. Волкова, А. И. Коробейников, С. А. Назарова, С. В. Петров, В. Г. Суфиянов

Наглядная статистика. Используем R!

<http://ashipunov.info/shipunov/school/books/rbook.pdf>

ПРОЧИТАТЬ РАЗДЕЛЫ

Глава 2. Как обрабатывать данные

2.4 Применение, преимущества и недостатки R, стр. 25-26

Глава 4. Великое в малом: одномерные данные

4.1. Как оценивать общую тенденцию, стр. 72-82

4.3. Одномерные статистические тесты, стр. 83-87

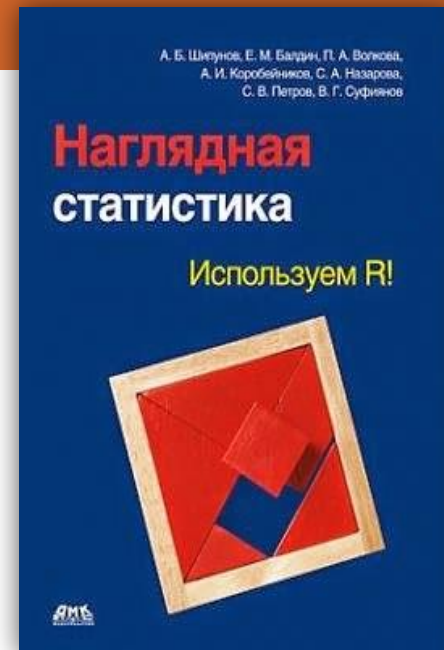
Глава 5. Анализ связей: двумерные данные

5.1. Что такое статистический тест, стр. 94-102

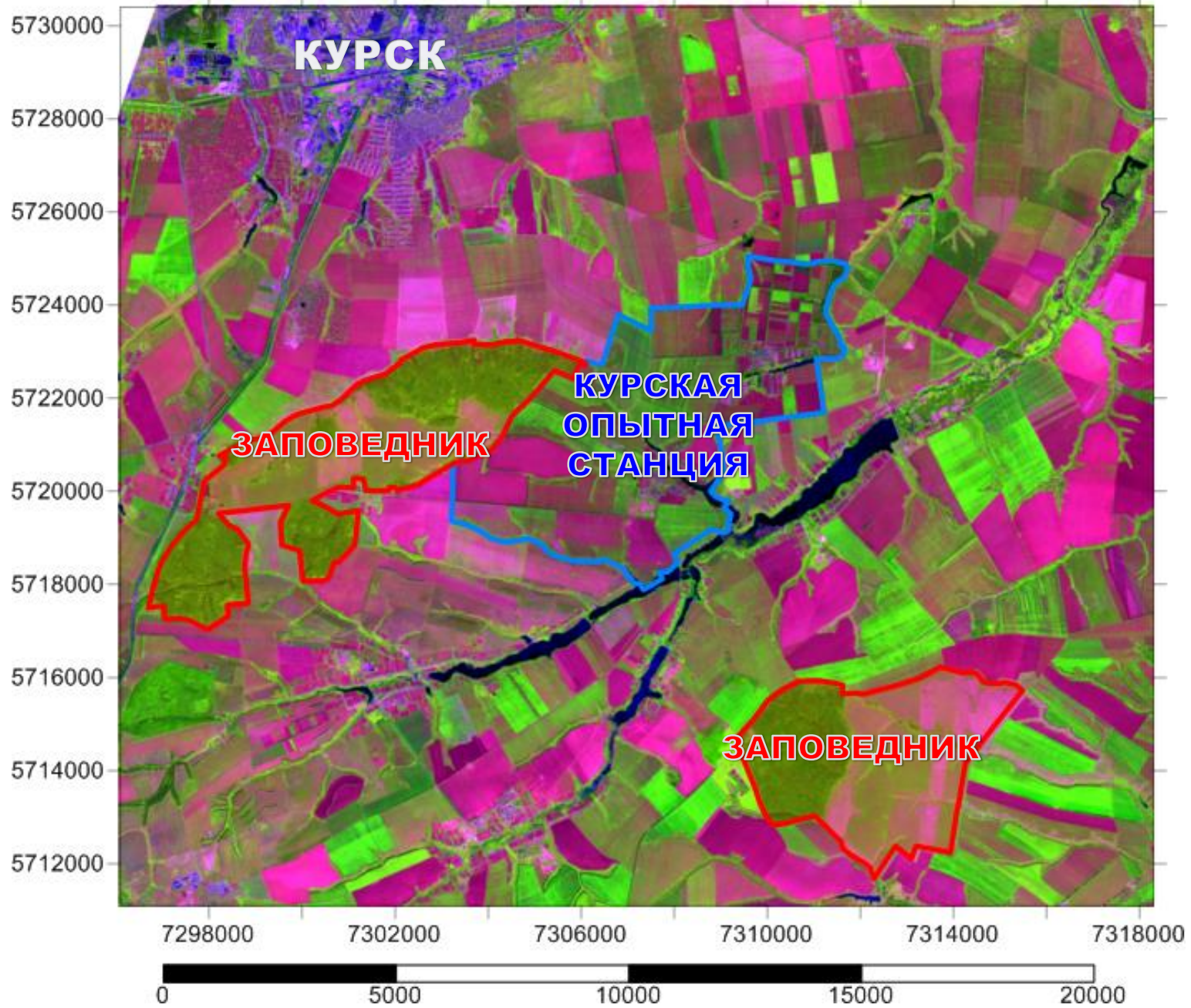
5.4. Есть ли взаимосвязь, или Анализ корреляций, стр. 109-114

5.5. Какая связь, или Регрессионный анализ, стр. 114-117

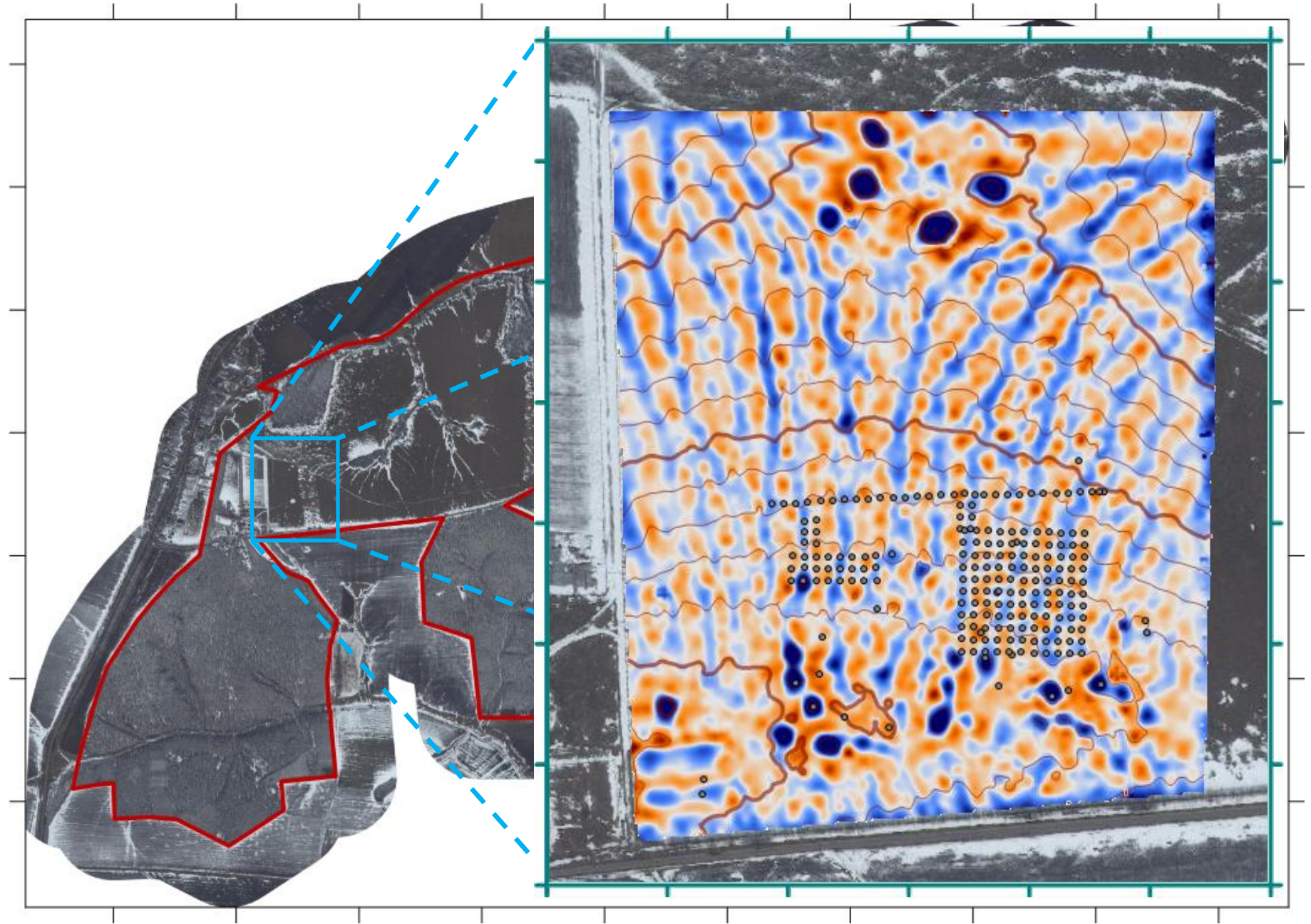
Приложение А. Пример работы в R, стр. 196-206



ТЕРРИТОРИЯ ИССЛЕДОВАНИЯ



КЛЮЧЕВОЙ УЧАСТОК «МЕТЕОСТАНЦИЯ»



ЗАДАНИЕ

- СФОРМУЛИРОВАТЬ ПРОБЛЕМУ И ЦЕЛЬ ИССЛЕДОВАНИЯ**
ИССЛЕДОВАТЬ пространственную изменчивость почв северной лесостепи Среднерусской возвышенности (Центрально-Черноземный заповедник)
- СПРОЕКТИРОВАТЬ ТАБЛИЦУ ДЛЯ СТАТИСТИЧЕСКОГО АНАЛИЗА**
строки – почвенные разрезы, колонки – свойства почвы (мощности горизонтов, таксон, запасы гумуса, кг/м²)
- ОБОСНОВАТЬ ШКАЛУ ИЗМЕРЕНИЙ КАЖДОГО СВОЙСТВА**
мощности горизонтов, запасы – интервальная, таксон - номинальная
- ЗАПОЛНИТЬ ТАБЛИЦУ в Excel**

id	A	AB	CaCO ₃	SID	HUM050	HUM100
m14-001	70	90	75	Чт	31.0	58.6
m14-002	75	90	120	Чтв	28.6	54.0
m14-003	85	95	135	Чтв	26.0	51.0
m14-004	70	85	80	Чл	38.0	62.0
...

ЗАДАНИЕ

5. ОБОСНОВАТЬ МЕТОДЫ СТАТИСТИЧЕСКОГО АНАЛИЗА (сравнение выборок, линейная связь между переменными)

- проверить H_0 равенства интервальных значений для номинальных категорий (запасов органического вещества для четырех подтипов черноземов)
- проверить H_0 отсутствия связи двух показателей (запасы органического вещества в слоях 0-50 и 0-100 см)

id	A	AB	CaCO ₃	SID	HUM050	HUM100
m14-001	70	90	75	ЧТ	31.0	58.6
m14-002	75	90	120	ЧТВ	28.6	54.0
m14-003	85	95	135	ЧТВ	26.0	51.0
m14-004	70	85	80	Чл	38.0	62.0
...

УСТАНОВИТЬ R И R-STUDIO

The screenshot displays the RStudio environment with the following components:

- Source Editor:** Contains R code for data manipulation and visualization. The code includes comments in Russian and performs operations like creating a new variable, plotting histograms, and density plots.
- Environment Pane:** Shows the current workspace with objects: data (1000 obs. of 3 variables), data.3 (100 obs. of 3 variables), datagen (1000 obs. of 2 variables), pedproph (150 obs. of 1 variable), and vars (1000 obs. of 2 variables).
- Console:** Shows the execution of the code, including several error messages: "Error: object 'h' not found".
- Plots Pane:** Displays a histogram of data\$var1 with blue bars, showing a distribution centered around 100.

```
22 # str(data) показать структуру таблицы
23 # write.table(data, "data/zhuki_new.txt", quote=FALSE) сохранить таблицу
24 # summary(data) сводная информация по таблице
25 # summary(data$VES) или отдельной колонке
26 dev.new()
27
28 data$var2.dx <- log(data$var2+1000) #создание новой колонки
29 # гистограмма
30 hist(data$var1, breaks=50, main="", col = "red")
31 hist(data$var2, breaks=50, main="", col = "blue")
32 hist(data$var1, breaks=50, main="", col = "lightblue", border = "darkgreen")
33 hist(data$var2.dx, breaks=50, main="", col = "blue", add=TRUE)
34 plot(density(data$var1, adjust=2), main="", col = "red")
35 plot(density(data$var1, adjust=2), main="", col = "blue", lwd=3, add=TRUE)
36 rug(data$var1)
37
38 head(data)
39
40
41 # проверка на нормальность
42 qqnorm(data$var1)
43
```

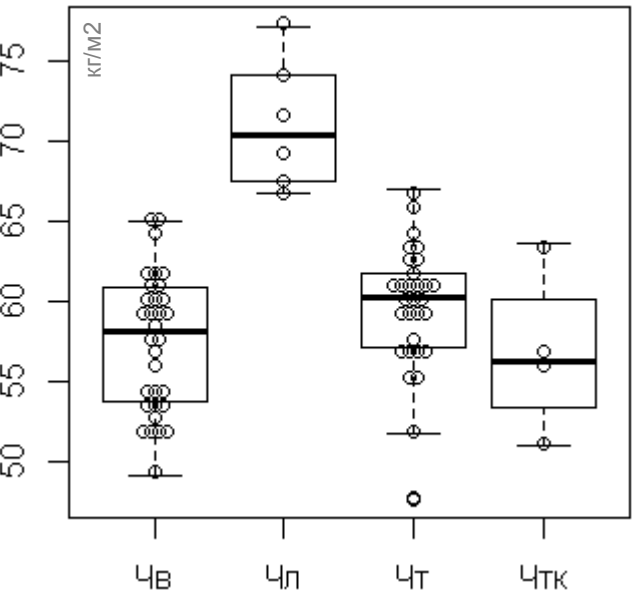
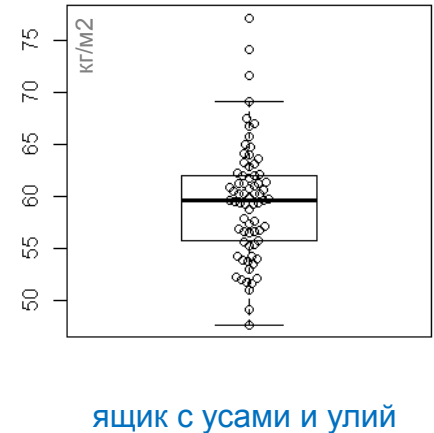
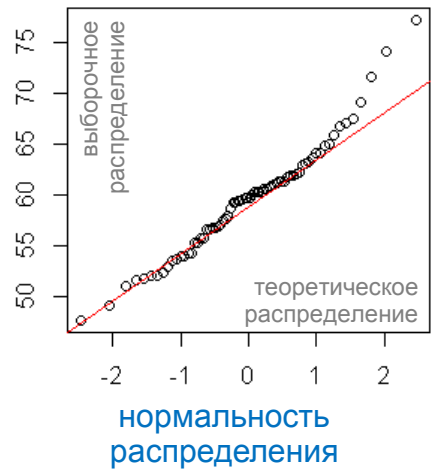
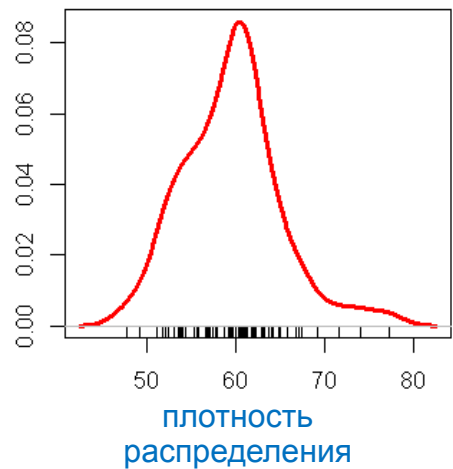
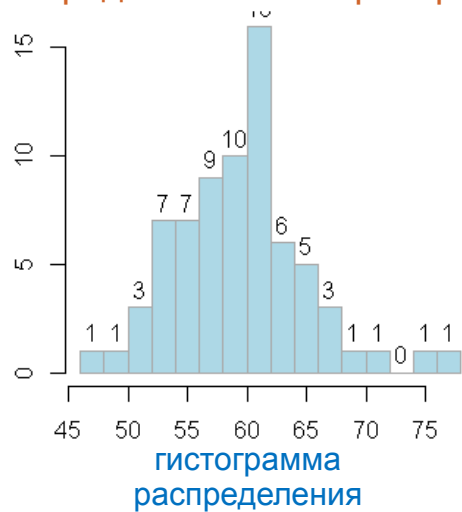
Console output:

```
there were 12 warnings (use warnings() to see them)
> hist(data$var2, breaks=50, main="", col = "blue", add=TRUE)
> help(hist)
> hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
> h
Error: object 'h' not found
> h
Error: object 'h' not found
> hist(data$var1, breaks=50, main="", col = "red")
> hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
> hist(data$var2, breaks=50, main="", col = "blue", add=TRUE)
> hist(data$var1, breaks=50, main="", col = "red")
> #hist(data$var1, breaks=50, main="", col = "lightblue", border = "pink")
```

ПРИМЕР ОФОРМЛЕНИЯ ЗАДАНИЯ

Н0: запасы органического вещества в метровом слое черноземов равны
 Чтк – типичный карбонатный, Чт – типичный, Чв – выщелоченный, Чл – лугово-черноземная почва

Порядок: 1. анализ распределений; 2. параметрический или непараметрический тест на различие средних



почва	N	сред.	тесты, p-level		
			норм. распр.*	Чт**	Чв**
Чтк	4	56.8	недостаточно данных		
Чт	30	59.7	0.009	–	–
Чв	32	57.6	0.209	0.047	–
Чл	6	71.1	0.668	0.001	0.0001

* тест на нормальность, ** - тест на достоверность различий средних

Вывод: с уровнем значимости 0.95 запасы гумуса в метровом слое различаются Чл – 71.1 кг/м2; Чт – 59.7 кг/м2; Чв – 57.6 кг/м2; для статистически обоснованного суждения по Чтк недостаточно исходных данных

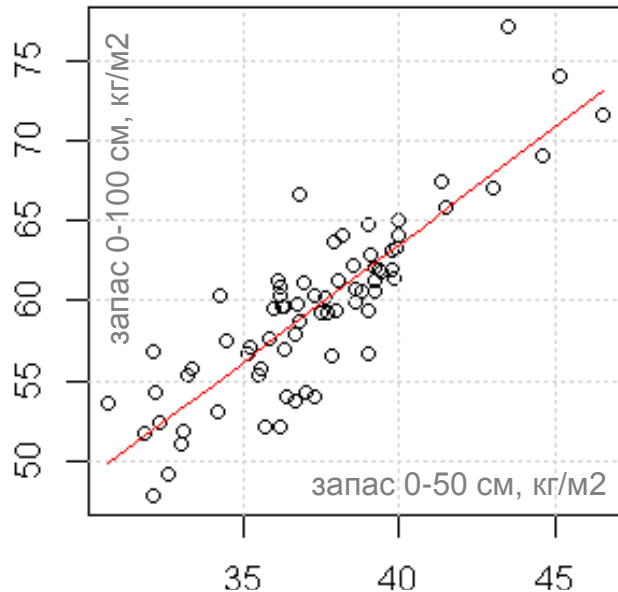
ПРИМЕР ОФОРМЛЕНИЯ ЗАДАНИЯ

H0: запасы органического вещества с глубиной меняются независимо

зависимая переменная (Y) — запасы гумуса в слое 0-100 см (h001m)

модель $Y = b_0 + b_1 * X$

независимая переменная (X) — запасы гумуса в слое 0-50 см (h0050)



```
lm(formula = data$h001m ~ data$h0050)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.7598	-1.5048	-0.0329	1.8413	8.4789

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.4578	4.0250	1.108	0.272
data\$h0050	1.4763	0.1075	13.733	<2e-16 ***

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

Residual standard error: 2.873 on 70 degrees of freedom

Multiple R-squared: 0.7293

F-statistic: 188.6 on 1 and 70 DF, p-value: < 2.2e-16

- линейная модель вида $[\text{запас } 0-100 \text{ см}] = 0 + 1.5 * [\text{запас } 0-50 \text{ см}]$ объясняет 72% ($R^2=0.72$) изменчивости запасов органического вещества в метровом слое при среднеквадратической ошибке прогноза 2.9 кг/м² (Residual standard error: 2.873);
- значимость модели подтверждает высокое значение F-статистики, равное 189, и общий уровень значимости: p-value: 2.2e-16, что много меньше 0:001;
- увеличение запасов в слое 0-50 на 10 кг соответствует увеличению запасов в слое 0-100 см примерно на 15 кг. Значение константы линейной модели (b_0) недостоверно отличается от 0;
- Наибольшее положительное отклонение истинного значения от модельного составляет 8.5 кг, наибольшее отрицательное -5.8 кг;
- Почти половина остатков находится в пределах от первой квантили (1Q = -1.5 кг/м²) до третьей (3Q = 1.8 кг/м²)

ПРАКТИЧЕСКАЯ РАБОТА №1

http://www.landscape.edu.ru/edu_help5_KiS.shtml#w1

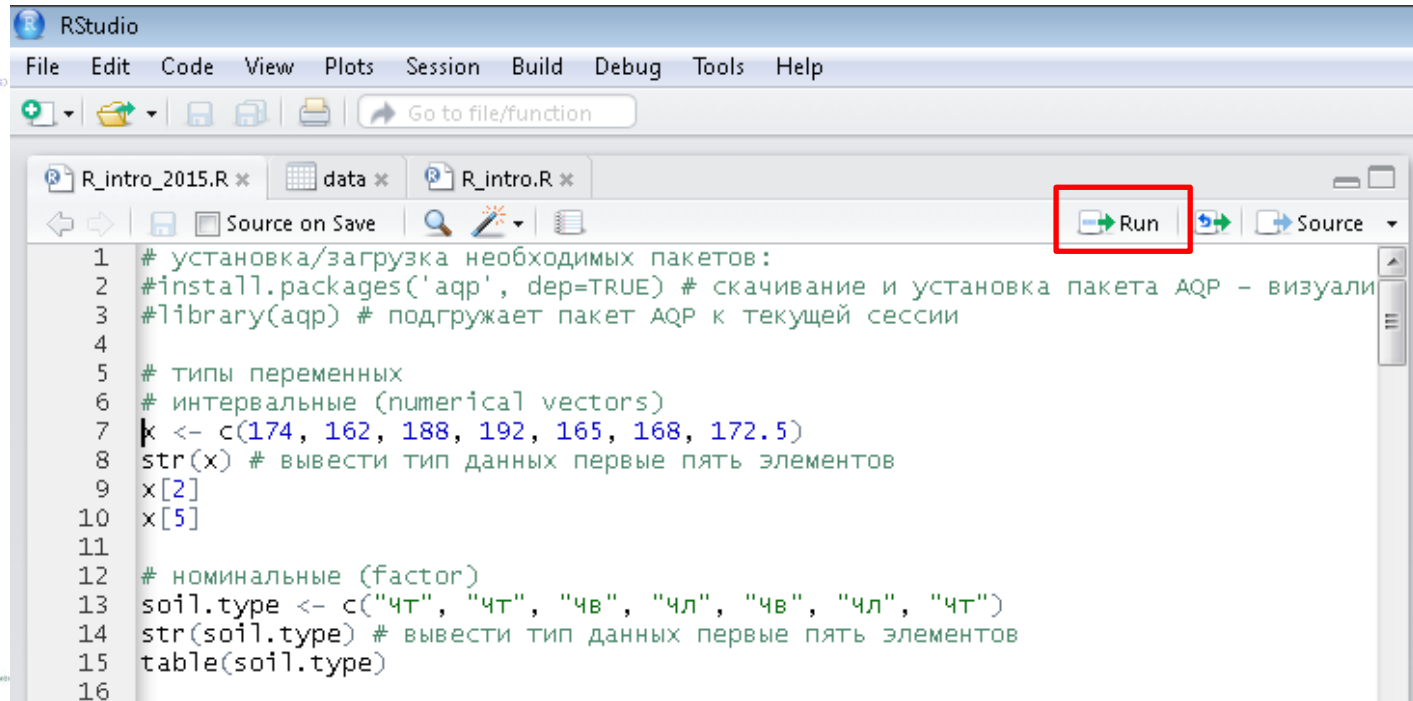
Задание № 1. «Статистический анализ в среде R»

Цель: исследовать пространственную изменчивость почв северной лесостепи Среднерусской возвышенности (Центрально-Черноземный заповедник).

 - инструкция к выполнению задания

 - zip-архив с комплектом материалов

1. скачать и распаковать архив с материалами задания
2. открыть в R-Studio код R_intro_2015.R
3. выполнить задание: курсор в строку команды + Run



The screenshot shows the RStudio environment with the R code editor open. The code in the editor is as follows:

```
1 # установка/загрузка необходимых пакетов :
2 #install.packages('aqp', dep=TRUE) # скачивание и установка пакета AQP - визуализация и анализ профилей
3 #library(aqp) # подгружает пакет AQP к текущей сессии
4
5 # типы переменных
6 # интервальные (numerical vectors)
7 k <- c(174, 162, 188, 192, 165, 168, 172.5)
8 str(x) # вывести тип данных первые пять элементов
9 x[2]
10 x[5]
11
12 # номинальные (factor)
13 soil.type <- c("чт", "чт", "чв", "чл", "чв", "чл", "чт")
14 str(soil.type) # вывести тип данных первые пять элементов
15 table(soil.type)
```

The 'Run' button in the RStudio toolbar is highlighted with a red box.

РЕЗУЛЬТАТА — КОНЕЧНОЙ ЦЕЛИ ИССЛЕДОВАНИЯ — НЕ СУЩЕСТВУЕТ

НУЛЕВАЯ ГИПОТЕЗА (H_0) – основное проверяемое предположение об отсутствии различий, отсутствие влияние фактора, отсутствие эффекта, равенство нулю значений характеристик модели и т.п.

ОШИБКА ПЕРВОГО РОДА —

отказаться от H_0 , в то время как она верна;

ОШИБКА ВТОРОГО РОДА —

принять H_0 , когда она на самом деле не верна

Уровень значимости (P -level) – вероятность ошибки первого рода при принятии решения (вероятность ошибочного отклонения нулевой гипотезы), $p = 5\%$, $p = 1\%$, $p = 0.1\%$

МАЛО = НЕ ДОСТОВЕРНО vs. МНОГО = ТРУДОЕМКО

НУЛЕВАЯ ГИПОТЕЗА

		ГЕНЕРАЛЬНАЯ СОВОКУПНОСТЬ («ИСТИНА»)	
		ВЕРНА НУЛЕВАЯ ГИПОТЕЗА (H0)	ВЕРНА АЛЬТЕРНАТИВНАЯ ГИПОТЕЗА (H1)
ВЫБОРКА («МОДЕЛЬ»)	ПРИНИМАЕМ НУЛЕВУЮ ГИПОТЕЗУ H0	<p>ПРАВИЛЬНО !</p> <p>H0 верно принята</p>	<p>H0 неверно принята</p> <p>ОШИБКА</p> <p>ВТОРОГО РОДА</p>
	ПРИНИМАЕМ АЛЬТЕРНАТИВНУЮ ГИПОТЕЗУ H1	<p>H0 неверно отвергнута</p> <p>ОШИБКА</p> <p>ПЕРВОГО РОДА</p>	<p>ПРАВИЛЬНО !</p> <p>H0 верно отвергнута</p>

Уровень значимости (*p-level, p-value*) – вероятность ошибки первого рода при принятии решения (вероятность ошибочного отклонения нулевой гипотезы),
 $p = 0.05, p = 0.01, p = 0.001$

H0: средние двух выборок (n=1000) не различаются (принадлежат одной генеральной совокупности)

