

# **Регрессионный анализ, часть 1**

**Математические методы в зоологии с использованием R**

Марина Варфоломеева

- 1 **Описание зависимости между переменными**
- 2 **Линейная регрессия**
- 3 **Неопределенность оценок коэффициентов**
- 4 **Тестирование значимости модели и ее коэффициентов**
- 5 **Оценка качества подгонки модели**

## Вы сможете

- посчитать и протестировать различные коэффициенты корреляции между переменными
- подобрать модель линейной регрессии и записать ее в виде уравнения
- интерпретировать коэффициенты простой линейной регрессии
- протестировать значимость модели и ее коэффициентов при помощи t- или F-теста
- оценить долю изменчивости, которую объясняет модель, при помощи  $R^2$

# Описание зависимости между переменными

## Пример: потеря влаги личинками мучных хрущаков

Как зависит потеря влаги личинками  
малого мучного хрущака *Tribolium confusum* от влажности воздуха?

- 9 экспериментов, продолжительность 6 дней
- разная относительная влажность воздуха, %
- измерена потеря влаги, мг



Малый мучной хрущак *Tribolium confusum*, photo by Sarefo, CC BY-SA

Nelson, 1964; данные из Sokal, Rohlf, 1997, табл. 14.1 по Logan, 2010. глава 8, пример 8с; Данные в файлах nelson.xlsx и nelson.csv

## Скачиваем данные с сайта

Не забудьте войти в вашу директорию для матметодов при помощи `setwd()`

```
library(downloader)
```

```
# в рабочем каталоге создаем суб-директорию для данных
```

```
if(!dir.exists("data")) dir.create("data")
```

```
# скачиваем файл в xlsx, либо в текстовом формате
```

```
if (!file.exists("data/nelson.xlsx")) {
```

```
  download(
```

```
    url = "https://varmara.github.io/mathmethr/data/nelson.xlsx",
```

```
    destfile = "data/nelson.xlsx")
```

```
}
```

```
if (!file.exists("data/nelson.csv")) {
```

```
  download(
```

```
    url = "https://varmara.github.io/mathmethr/data/nelson.xls",
```

```
    destfile = "data/nelson.csv")
```

```
}
```

# Читаем данные из файла одним из способов

## Чтение из xlsx

```
library(readxl)
nelson <- read_excel(path = "data/nelson.xlsx", sheet = 1)
```

## Чтение из csv

```
nelson <- read.table("data/nelson.csv", header = TRUE, sep = "\t")
```

## Все ли правильно открылось?

```
str(nelson) # Структура данных
```

```
# 'data.frame': 9 obs. of 2 variables:  
# $ humidity : num 0 12 29.5 43 53 62.5 75.5 85 93  
# $ weightloss: num 8.98 8.14 6.67 6.08 5.9 5.83 4.68 4.2 3.72
```

```
head(nelson) # Первые несколько строк файла
```

```
# humidity weightloss  
# 1      0.0      8.98  
# 2     12.0      8.14  
# 3     29.5      6.67  
# 4     43.0      6.08  
# 5     53.0      5.90  
# 6     62.5      5.83
```



## Знакомимся с данными

Есть ли пропущенные значения?

```
sapply(nelson, function(x)sum(is.na(x)))
```

```
# humidity weightloss  
#           0         0
```

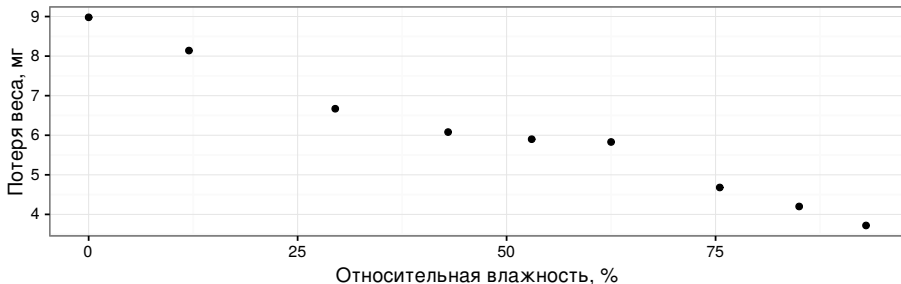
Каков объем выборки?

```
nrow(nelson)
```

```
# [1] 9
```

## Как зависит потеря веса от влажности?

```
library(ggplot2)
theme_set(theme_bw())
gg_nelson <- ggplot(data=nelson, aes(x = humidity, y = weightloss)) +
  geom_point() +
  labs(x = "Относительная влажность, %",
       y = "Потеря веса, мг")
gg_nelson
```



## **Коэффициент корреляции — способ оценки силы связи между двумя переменными**

## Коэффициент корреляции — способ оценки силы связи между двумя переменными

### Коэффициент корреляции Пирсона

- Оценивает только линейную составляющую связи
- Параметрические тесты (t-критерий) значимости применимы если переменные распределены нормально

# Коэффициент корреляции — способ оценки силы связи между двумя переменными

## Коэффициент корреляции Пирсона

- Оценивает только линейную составляющую связи
- Параметрические тесты (t-критерий) значимости применимы если переменные распределены нормально

## Ранговые коэффициенты корреляции (кор. Кендалла и кор. Спирмена)

- Не зависят от формы распределения переменных
- Тест на значимость непараметрический

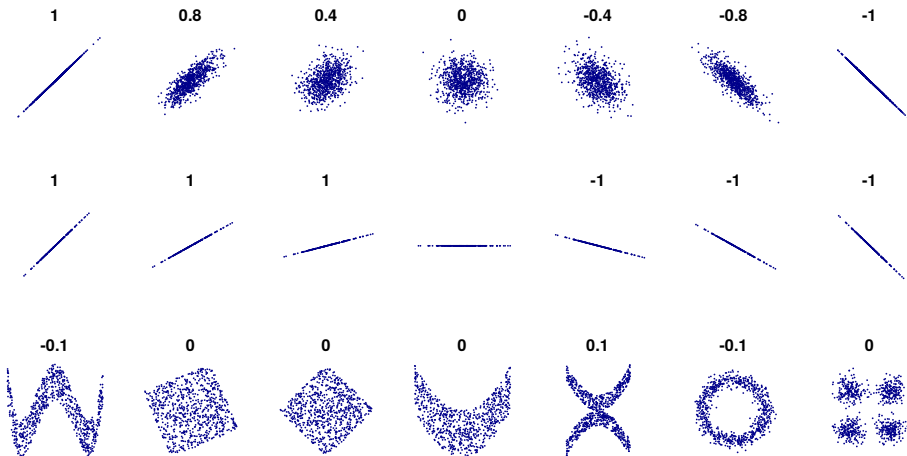
# Интерпретация коэффициента корреляции

$$-1 < \rho < 1$$

$|\rho| = 1$  — сильная связь

$\rho = 0$  — нет связи

- В тестах для проверки значимости тестируется гипотеза  $H_0 : \rho = 0$



By DenisBoigelot, original uploader was Imagecreator [CC0], via Wikimedia Commons

## Можно рассчитать значение коэффициента корреляции между потерей веса и влажностью

```
p_cor <- cor.test(nelson$humidity, nelson$weightloss,  
                 alternative = "two.sided", method = "pearson")  
p_cor
```

```
#  
# Pearson's product-moment correlation  
#  
# data: nelson$humidity and nelson$weightloss  
# t = -16.346, df = 7, p-value = 0.0000007816  
# alternative hypothesis: true correlation is not equal to 0  
# 95 percent confidence interval:  
# -0.9973935 -0.9379224  
# sample estimates:  
# cor  
# -0.9871523
```

## Можно рассчитать значение коэффициента корреляции между потерей веса и влажностью

```
p_cor <- cor.test(nelson$humidity, nelson$weightloss,
                 alternative = "two.sided", method = "pearson")
p_cor
```

```
#
# Pearson's product-moment correlation
#
# data: nelson$humidity and nelson$weightloss
# t = -16.346, df = 7, p-value = 0.0000007816
# alternative hypothesis: true correlation is not equal to 0
# 95 percent confidence interval:
# -0.9973935 -0.9379224
# sample estimates:
# cor
# -0.9871523
```

Можно описать результаты несколькими способами:

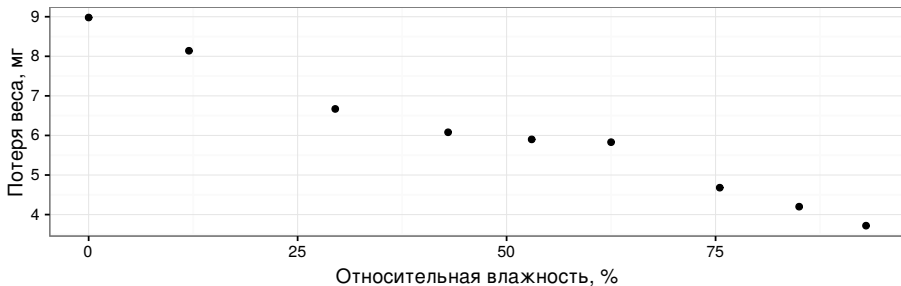
- Величина потери веса мучных хрущаков отрицательно коррелирует с относительной влажностью воздуха ( $r = -0.99, p < 0.01$ )
- Мучные хрущаки теряют вес при уменьшении относительной влажности воздуха ( $r = -0.99, p < 0.01$ )



## Коэффициент корреляции не позволяет предсказать значение одной переменной, зная значение другой

Нам бы хотелось описать функциональную зависимость

$$weightloss_i = b_0 + b_1 humidity_i$$



# Линейная регрессия

# Линейная регрессия

- простая

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- множественная

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \varepsilon_i$$

# Как провести линию регрессии?

Линейная модель:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Оценка модели:

$$\hat{y}_i = b_0 + b_1 x_i$$

Нужно оценить параметры линейной модели:

- $\beta_0$
- $\beta_1$

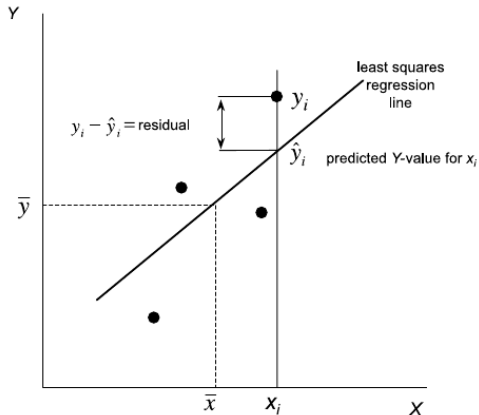
Методы оценки параметров:

- Метод наименьших квадратов (Ordinary Least Squares)
- Методы максимального правдоподобия (Maximum Likelihood, REstricted Maximum Likelihood)

# Метод наименьших квадратов

$$\hat{y}_i = b_0 + b_1 x_i$$

Оценки параметров линейной регрессии подбирают так, чтобы минимизировать остатки  $\sum (y_i - \hat{y}_i)^2$



Линия регрессии по методу наименьших квадратов

из кн. Quinn, Keough, 2002, стр. 85, рис. 5.6 а

# Оценки параметров линейной регрессии

Параметры	Оценки параметров	Стандартные ошибки оценок
$\beta_1$	$b_1 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n (x_i - \bar{x})^2}$	$SE_{b_1} = \sqrt{\frac{MS_e}{\sum_{i=1}^n (x_i - \bar{x})^2}}$
$\beta_0$	$b_0 = \bar{y} - b_1 \bar{x}$	$SE_{b_0} = \sqrt{MS_e \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]}$

Таблица из кн. Quinn, Keough, 2002, стр. 86, табл. 5.2

Стандартные ошибки коэффициентов - используются для построения доверительных интервалов - нужны для статистических тестов

# Интерпретация коэффициентов регрессии

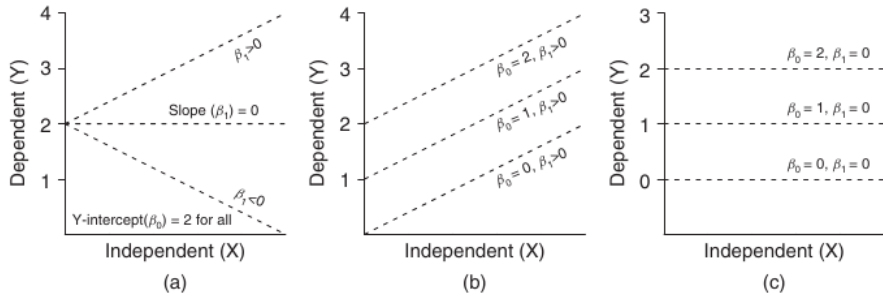


Рисунок из кн. Logan, 2010, стр. 170, рис. 8.2

# Интерпретация коэффициентов регрессии

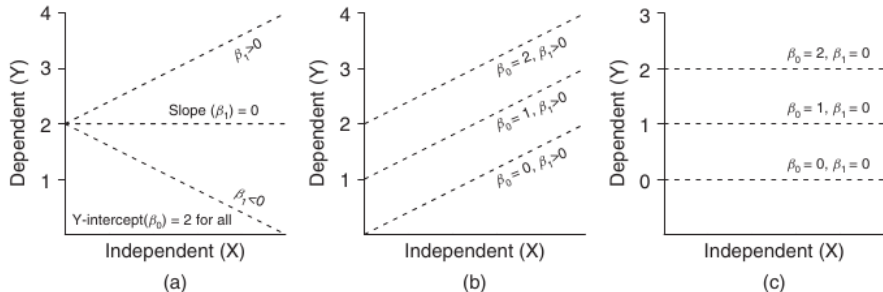


Рисунок из кн. Logan, 2010, стр. 170, рис. 8.2

- $b_0$  — Отрезок (Intercept), отсекаемый регрессионной прямой на оси  $y$ . Значение зависимой переменной  $y$ , если предиктор  $x = 0$ .
- $b_1$  — Коэффициент угла наклона регрессионной прямой. Показывает на сколько единиц изменяется отклик ( $y$ ), при увеличении значения предиктора ( $x$ ) на единицу.

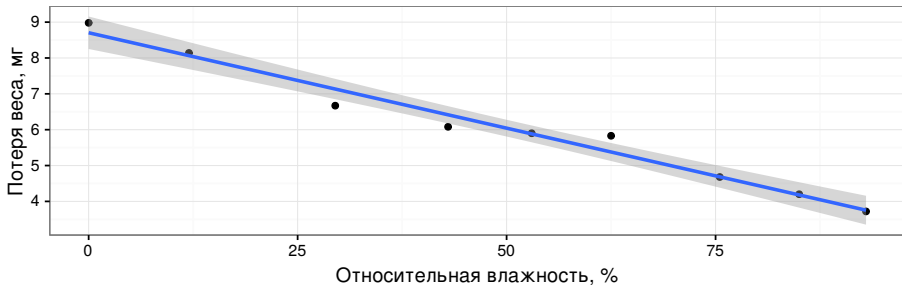


## Для сравнения разных моделей - стандартизованные коэффициенты

- Не зависят от масштаба измерений  $x$  и  $y$
- Можно вычислить, зная обычные коэффициенты и их стандартные отклонения  $b_1^* = b_1 \frac{\sigma_x}{\sigma_y}$
- Можно вычислить, посчитав регрессию по стандартизованным данным

## Добавим линию регрессии на график

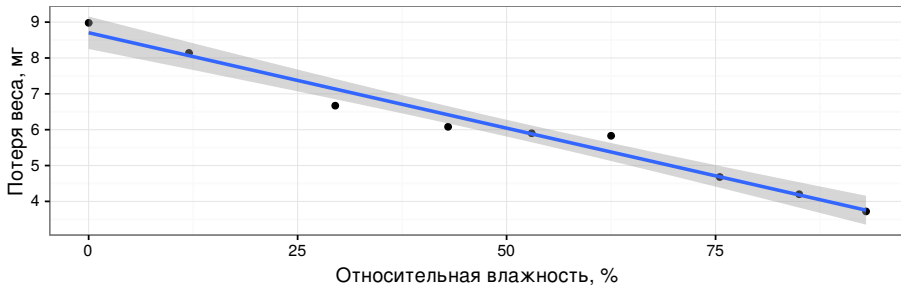
```
gg_nelson + geom_smooth(method = "lm")
```



Что это за серая область вокруг линии регрессии?

## Добавим линию регрессии на график

```
gg_nelson + geom_smooth(method = "lm")
```



Что это за серая область вокруг линии регрессии?

### Доверительная зона регрессии

- 95% доверительная зона регрессии
- В ней с 95% вероятностью лежит регрессионная прямая
- Возникает из-за неопределенности оценок коэффициентов регрессии

## Как в R задать формулу линейной регрессии

`lm(формула, данные)` - функция для подбора регрессионных моделей

Формат формулы: зависимая\_переменная ~ модель

- $\hat{y}_i = b_0 + b_1 x_i$  (простая линейная регрессия с  $b_0$  (intercept))
  - $Y \sim X$
  - $Y \sim 1 + X$
  - $Y \sim X + 1$
- $\hat{y}_i = b_1 x_i$  (простая линейная регрессия без  $b_0$ )
  - $Y \sim X - 1$
  - $Y \sim -1 + X$
- $\hat{y}_i = b_0$  (уменьшенная модель, линейная регрессия  $Y$  от  $b_0$ )
  - $Y \sim 1$
  - $Y \sim 1 - X$

## Задача

Запишите в нотации R эти модели линейных регрессий

- $\hat{y}_i = b_0 + b_1x_{1i} + b_2x_{2i} + b_3x_{3i}$

(множественная линейная регрессия с  $b_0$ )

- $\hat{y}_i = b_0 + b_1x_{1i} + b_3x_{3i}$

(уменьшенная модель множественной линейной регрессии, без  $x_2$ )

## Решение

- $\hat{y}_i = b_0 + b_1x_{1i} + b_2x_{2i} + b_3x_{3i}$

(множественная линейная регрессия с  $b_0$ )

$$Y \sim X1 + X2 + X3$$

$$Y \sim 1 + X1 + X2 + X3$$

- $\hat{y}_i = b_0 + b_1x_{1i} + b_3x_{3i}$

(уменьшенная модель множественной линейной регрессии, без  $x_2$ )

$$Y \sim X1 + X3$$

$$Y \sim 1 + X1 + X3$$

## Подбираем параметры линейной модели

```
nelson_lm <- lm(weightloss ~ humidity, nelson)
summary(nelson_lm)
```

```
#
# Call:
# lm(formula = weightloss ~ humidity, data = nelson)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.46397 -0.03437  0.01675  0.07464  0.45236
#
# Coefficients:
#              Estimate Std. Error t value      Pr(>|t|)
# (Intercept)  8.704027   0.191565   45.44 0.0000000000654 ***
# humidity    -0.053222   0.003256  -16.35 0.000000781615 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# Residual standard error: 0.2967 on 7 degrees of freedom
# Multiple R-squared:  0.9745, Adjusted R-squared:  0.9708
# F-statistic: 267.2 on 1 and 7 DF, p-value: 0.0000007816
```

## Подбираем параметры линейной модели

```
nelson_lm <- lm(weightloss ~ humidity, nelson)
summary(nelson_lm)
```

```
#
# Call:
# lm(formula = weightloss ~ humidity, data = nelson)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.46397 -0.03437  0.01675  0.07464  0.45236
#
# Coefficients:
#              Estimate Std. Error t value    Pr(>|t|)
# (Intercept)  8.704027   0.191565   45.44 0.0000000000654 ***
# humidity    -0.053222   0.003256  -16.35 0.000000781615 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# Residual standard error: 0.2967 on 7 degrees of freedom
# Multiple R-squared:  0.9745, Adjusted R-squared:  0.9708
# F-statistic: 267.2 on 1 and 7 DF, p-value: 0.0000007816
```

Коэффициенты линейной регрессии:

- $b_0 = 8.7 \pm 0.2$
- $b_1 = -0.053 \pm 0.003$



## Записываем уравнение линейной регрессии

Коэффициенты модели:

```
coef(nelson_lm)
```

```
# (Intercept)    humidity  
#  8.70402730 -0.05322215
```

Уравнение регрессии:

weightloss = 8.70 - 0.05 humidity

Более формальная запись:

$$Y = 8.70 - 0.05 X_1$$

## Неопределенность оценок коэффициентов

# Неопределенность оценок коэффициентов

## Доверительный интервал коэффициента

- зона, в которой с  $(1 - \alpha) \cdot 100\%$  вероятностью содержится среднее значение коэффициента
- $b_1 \pm t_{\alpha, df=n-2} \cdot SE_{b_1}$
- $\alpha = 0.05 \Rightarrow (1 - 0.05) \cdot 100\% = 95\%$  интервал

## Доверительная зона регрессии

- зона, в которой с  $(1 - \alpha) \cdot 100\%$  вероятностью лежит регрессионная прямая

## Находим доверительные интервалы коэффициентов

```
# оценки коэффициентов отдельно
```

```
coef(nelson_lm)
```

```
# (Intercept)    humidity
```

```
#  8.70402730 -0.05322215
```

```
# доверительные интервалы коэффициентов
```

```
confint(nelson_lm)
```

```
#                2.5 %        97.5 %
```

```
# (Intercept)  8.25104923  9.15700538
```

```
# humidity    -0.06092143 -0.04552287
```

## Предсказываем $Y$ при заданном $X$

Какова средняя потеря веса при заданной влажности?

```
newdata <- data.frame(humidity = c(50, 100)) # значения, для которых предсказываем
(pr1 <- predict(nelson_lm, newdata, interval = "confidence", se = TRUE))
```

```
# $fit
#           fit           lwr           upr
# 1 6.042920 5.809068 6.276771
# 2 3.381812 2.933952 3.829672
#
# $se.fit
#           1           2
# 0.09889579 0.18940006
#
# $df
# [1] 7
#
# $residual.scale
# [1] 0.2966631
```

## Предсказываем $Y$ при заданном $X$

Какова средняя потеря веса при заданной влажности?

```
newdata <- data.frame(humidity = c(50, 100)) # значения, для которых предсказываем
(pr1 <- predict(nelson_lm, newdata, interval = "confidence", se = TRUE))
```

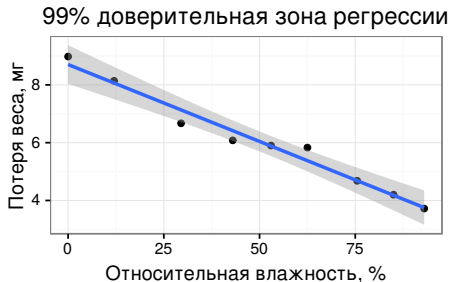
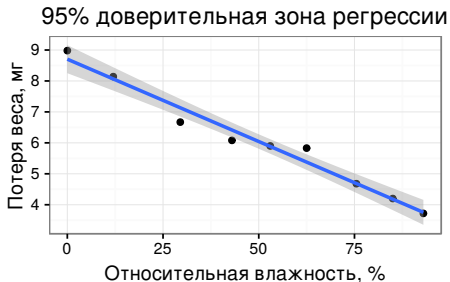
```
# $fit
#           fit           lwr           upr
# 1 6.042920 5.809068 6.276771
# 2 3.381812 2.933952 3.829672
#
# $se.fit
#           1           2
# 0.09889579 0.18940006
#
# $df
# [1] 7
#
# $residual.scale
# [1] 0.2966631
```

- При 50 и 100% относительной влажности ожидаемая средняя потеря веса жуков будет  $6 \pm 0.2$  и  $3.4 \pm 0.4$ , соответственно.

## Строим доверительную зону регрессии

```
gg_nelson + geom_smooth(method = "lm") +  
  labs (title = "95% доверительная зона регрессии")
```

```
gg_nelson + geom_smooth(method = "lm", level = 0.99) +  
  labs (title = "99% доверительная зона регрессии")
```



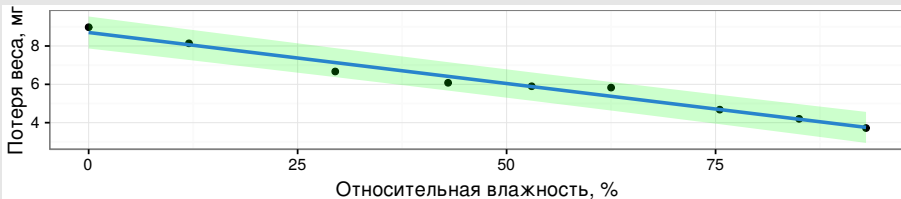
# Неопределенность оценок предсказанных значений

## Доверительный интервал к предсказанному значению

- зона в которую попадают  $(1 - \alpha) \cdot 100\%$  значений  $\hat{y}_i$  при данном  $x_i$
- $\hat{y}_i \pm t_{\alpha, n-2} \cdot SE_{\hat{y}_i}$
- $SE_{\hat{y}} = \sqrt{MS_e \left[ 1 + \frac{1}{n} + \frac{(x_{prediction} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]}$

## Доверительная область значений регрессии

- зона, в которую попадает  $(1 - \alpha) \cdot 100\%$  всех предсказанных значений





## Предсказываем изменение $Y$ для 95% наблюдений при заданном $X$

В каких пределах находится потеря веса у 95% жуков при заданной влажности?

```
newdata <- data.frame(humidity = c(50, 100)) # новые данные для предсказания значений
(pr2 <- predict(nelson_lm, newdata, interval = "prediction", se = TRUE))
```

```
# $fit
#           fit          lwr          upr
# 1 6.042920 5.303471 6.782368
# 2 3.381812 2.549540 4.214084
#
# $se.fit
#           1          2
# 0.09889579 0.18940006
#
# $df
# [1] 7
#
# $residual.scale
# [1] 0.2966631
```

## Предсказываем изменение $Y$ для 95% наблюдений при заданном $X$

В каких пределах находится потеря веса у 95% жуков при заданной влажности?

```
newdata <- data.frame(humidity = c(50, 100)) # новые данные для предсказания значений
(pr2 <- predict(nelson_lm, newdata, interval = "prediction", se = TRUE))
```

```
# $fit
#           fit          lwr          upr
# 1 6.042920 5.303471 6.782368
# 2 3.381812 2.549540 4.214084
#
# $se.fit
#           1          2
# 0.09889579 0.18940006
#
# $df
# [1] 7
#
# $residual.scale
# [1] 0.2966631
```

- У 95% жуков при 50 и 100% относительной влажности будет потеря веса будет в пределах  $6 \pm 0.7$  и  $3.4 \pm 0.8$ , соответственно.

## Данные для доверительной области значений

Предсказанные значения для исходных данных объединим с исходными данными в новом датафрейме - для графиков

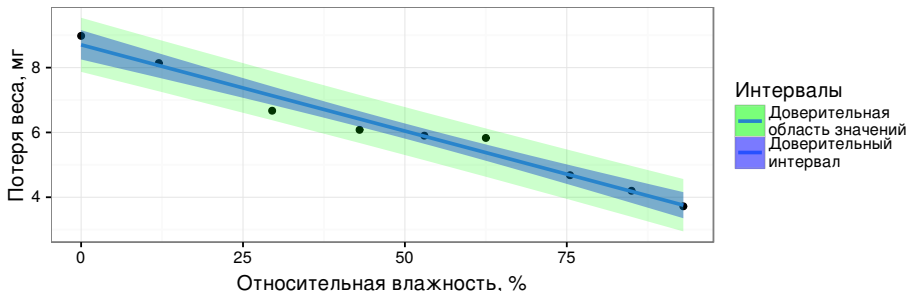
```
(pr_all <- predict(nelson_lm, interval = "prediction"))
```

```
#      fit      lwr      upr
# 1 8.704027 7.868990 9.539064
# 2 8.065361 7.269036 8.861687
# 3 7.133974 6.377243 7.890704
# 4 6.415475 5.673847 7.157102
# 5 5.883253 5.143538 6.622969
# 6 5.377643 4.632344 6.122941
# 7 4.685755 3.921455 5.450055
# 8 4.180144 3.394150 4.966139
# 9 3.754367 2.945412 4.563322
```

```
nelson_with_pred <- data.frame(nelson, pr_all)
```

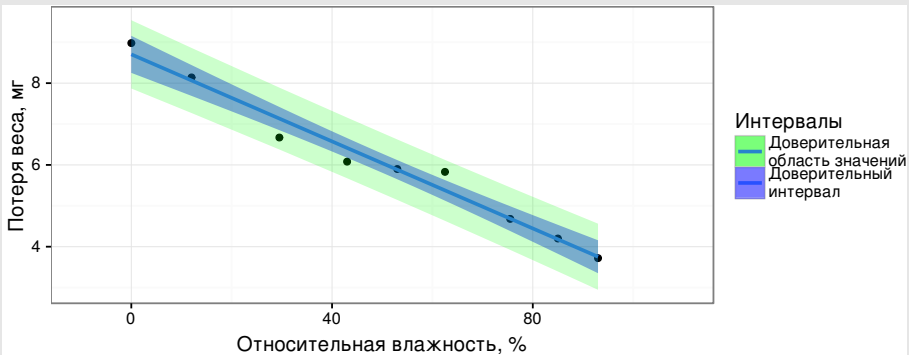
# Строим доверительную область значений и доверительный интервал одновременно

```
gg_nelson +
  geom_smooth(method = "lm",
    aes(fill = "Доверительный \n интервал"),
    alpha = 0.4) +
  geom_ribbon(data = nelson_with_pred,
    aes(y = fit, ymin = lwr, ymax = upr,
      fill = "Доверительная \n область значений"),
    alpha = 0.2) +
  scale_fill_manual('Интервалы', values = c('green', 'blue'))
```



# Осторожно!

Вне интервала значений  $X$  ничего предсказать нельзя!



# Тестирование значимости модели и ее коэффициентов

# Тестируем коэффициенты t-критерием

## t-критерий

$$t = \frac{b_1 - \theta}{SE_{b_1}}$$

$H_0 : b_1 = \theta$ , для  $\theta = 0$

Число степеней свободы  $df = n - 2$

# Тестируем значимость коэффициентов с помощью t-критерия

```
summary(nelson_lm)
```

```
#
# Call:
# lm(formula = weightloss ~ humidity, data = nelson)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.46397 -0.03437  0.01675  0.07464  0.45236
#
# Coefficients:
#              Estimate Std. Error t value    Pr(>|t|)
# (Intercept)  8.704027    0.191565   45.44 0.000000000654 ***
# humidity    -0.053222    0.003256  -16.35 0.000000781615 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# Residual standard error: 0.2967 on 7 degrees of freedom
# Multiple R-squared:  0.9745, Adjusted R-squared:  0.9708
# F-statistic: 267.2 on 1 and 7 DF, p-value: 0.0000007816
```



# Тестируем значимость коэффициентов с помощью t-критерия

```
summary(nelson_lm)
```

```
#
# Call:
# lm(formula = weightloss ~ humidity, data = nelson)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.46397 -0.03437  0.01675  0.07464  0.45236
#
# Coefficients:
#              Estimate Std. Error t value    Pr(>|t|)
# (Intercept)  8.704027    0.191565   45.44 0.000000000654 ***
# humidity     -0.053222    0.003256  -16.35 0.000000781615 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# Residual standard error: 0.2967 on 7 degrees of freedom
# Multiple R-squared:  0.9745, Adjusted R-squared:  0.9708
# F-statistic: 267.2 on 1 and 7 DF, p-value: 0.0000007816
```

Результаты можно описать в тексте так:

- Увеличение относительной влажности привело к достоверному замедлению потери веса жуками ( $b_1 = -0.053$ ,  $t = -16.35$ ,  $p < 0.01$ )

## Проверка при помощи F-критерия

### F-критерий

$$F = \frac{MS_{\text{regression}}}{MS_{\text{error}}}$$

$$H_0 : \beta_1 = 0$$

Число степеней свободы  $df_{\text{regression}}$ ,  $df_{\text{error}}$

## Общая изменчивость

Общая изменчивость -  $SS_{total}$ , отклонения от общего среднего значения

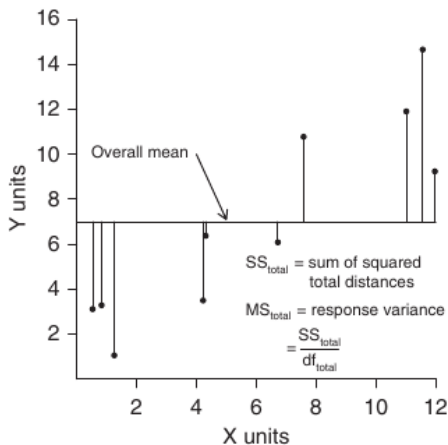
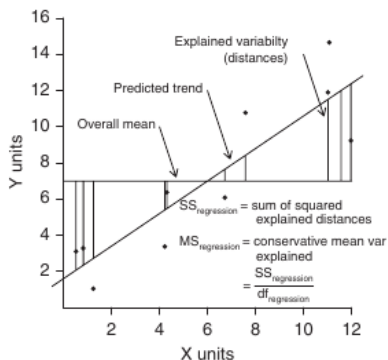


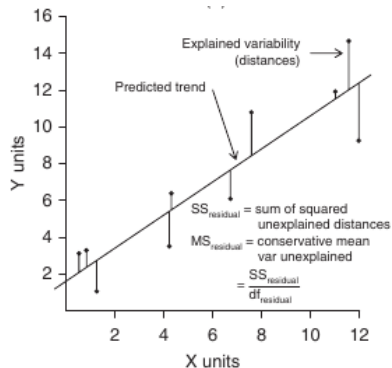
Рис. из кн. Logan, 2010, стр. 172, рис. 8.3

# Общая изменчивость

$$SS_{total} = SS_{regression} + SS_{error}$$



Объясненная изменчивость

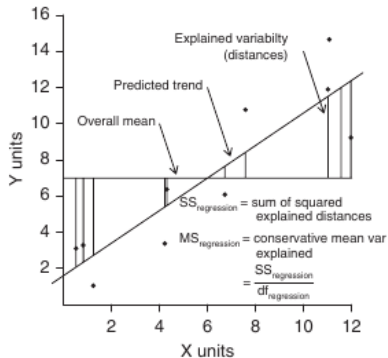


Остаточная изменчивость

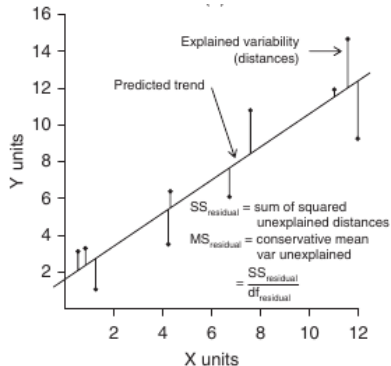
Рис. из кн. Logan, 2010, стр. 172, рис. 8.3

## Если зависимости нет, $b_1 = 0$

Тогда  $\hat{y}_i = \bar{y}_i$  и  $MS_{\text{regression}} \approx MS_{\text{error}}$



Объясненная изменчивость



Остаточная изменчивость

Рис. из кн. Logan, 2010, стр. 172, рис. 8.3

## Что оценивают средние квадраты отклонений?

Источник изменчивости	Число степеней свободы df	Суммы квадратов отклонений SS	Средний квадрат отклонений MS	Ожидаемый средний квадрат
Регрессия	1	$\sum (\bar{y} - \hat{y}_i)^2$	$\frac{\sum_{i=1}^n (\bar{y} - \hat{y}_i)^2}{1}$	$\sigma_\varepsilon^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$
Остаточная	$n - 2$	$\sum (y_i - \hat{y}_i)^2$	$\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$	$\sigma_\varepsilon^2$
Общая	$n - 1$	$\sum (\bar{y} - y_i)^2$		

Если  $b_1 = 0$ , тогда  $\hat{y}_i = \bar{y}_i$  и  $MS_r \approx MS_e$

Тестируем:

$$F = \frac{MS_{\text{regression}}}{MS_{\text{error}}}$$

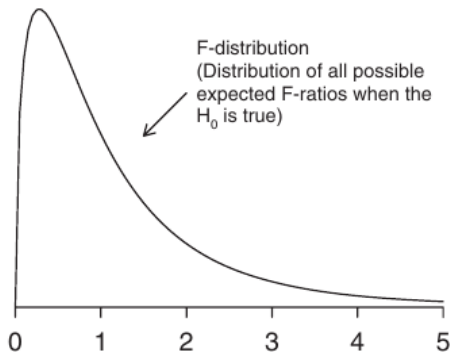
# F-критерий и распределение F-статистики

F - соотношение объясненной  
и не объясненной  
изменчивости

$$F = \frac{MS_r}{MS_e}$$

Зависит от

- $\alpha$
- $df_r$
- $df_e$



Распределение F-статистики при  
справедливой  $H_0$

Рис. из кн. Logan, 2010, стр. 172, рис. 8.3

## Таблица результатов дисперсионного анализа

Источник изменчивости	df	SS	MS	F
Регрессия	$df_r = 1$	$SS_r = \sum (\bar{y} - \hat{y}_i)^2$	$MS_r = \frac{SS_r}{df_r}$	$F_{df_r, df_e} = \frac{MS_r}{MS_e}$
Остаточная	$df_e = n - 2$	$SS_e = \sum (y_i - \hat{y}_i)^2$	$MS_e = \frac{SS_e}{df_e}$	
Общая	$df_t = n - 1$	$SS_t = \sum (\bar{y} - y_i)^2$		

Минимальное упоминание результатов в тексте должно содержать  $F_{df_r, df_e}$  и  $p$ .



## Проверяем значимость модели при помощи F-критерия

```
nelson_aov <- aov(nelson_lm)
summary(nelson_aov)
```

```
#               Df Sum Sq Mean Sq F value    Pr(>F)
# humidity      1 23.514   23.514   267.2 0.000000782 ***
# Residuals     7  0.616    0.088
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Результаты дисперсионного анализа можно описать в тексте:

- Количество влаги, потерянной жуками в период эксперимента, достоверно зависело от уровня относительной влажности ( $F_{1,7} = 267, p < 0.01$ ).

## Результаты дисперсионного анализа можно представить в виде таблицы

- Количество влаги, потерянной жуками в период эксперимента, достоверно зависело от уровня относительной влажности (Табл. 1).

**Таблица 1:** Результаты дисперсионного анализа зависимости потери веса мучных хрущаков от относительной влажности воздуха. *df* — число степеней свободы, *SS* — суммы квадратов отклонений, *MS* — средние квадраты отклонений, *F* — значение *F*-критерия, *P* — доверительная вероятность.

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>P</i>
Влажность	1	23.51	23.51	267.18	< 0.01
Остаточная	7	0.62	0.09		

## Оценка качества подгонки модели

## Коэффициент детерминации

### Коэффициент детерминации $R^2$

доля общей изменчивости, объясненная линейной связью  $x$  и  $y$

$$R^2 = \frac{SS_r}{SS_t} = 1 - \frac{SS_e}{SS_t}$$

$$0 \leq R^2 \leq 1$$

Иначе рассчитывается как квадрат коэффициента корреляции  $R^2 = r^2$   
**Не используйте обычный  $R^2$  для множественной регрессии!**

## Коэффициент детерминации можно найти в сводке модели

```
summary(nelson_lm)
```

```
#
# Call:
# lm(formula = weightloss ~ humidity, data = nelson)
#
# Residuals:
#      Min       1Q   Median       3Q      Max
# -0.46397 -0.03437  0.01675  0.07464  0.45236
#
# Coefficients:
#              Estimate Std. Error t value      Pr(>|t|)
# (Intercept)  8.704027    0.191565   45.44 0.000000000654 ***
# humidity    -0.053222    0.003256  -16.35 0.000000781615 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# Residual standard error: 0.2967 on 7 degrees of freedom
# Multiple R-squared:  0.9745, Adjusted R-squared:  0.9708
# F-statistic: 267.2 on 1 and 7 DF, p-value: 0.0000007816
```

## Сравнение качества подгонки моделей

$R_{adj}^2$  — скорректированный  $R^2$

$$R_{adj}^2 = 1 - \frac{SS_e/df_e}{SS_t/df_t}$$

где  $df_e = n - p - 1$ ,  $df_t = n - 1$

$R_{adj}^2$  учитывает число переменных в модели, вводится штраф за каждый новый параметр.

Используйте  $R_{adj}^2$  для сравнения моделей с разным числом параметров.

## Take home messages

- Модель простой линейной регрессии  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- В оценке коэффициентов регрессии и предсказанных значений существует неопределенность. Доверительные интервалы можно рассчитать, зная стандартные ошибки.
- Значимость всей регрессии и ее параметров можно проверить при помощи t- или F-теста.  $H_0 : \beta_1 = 0$
- Качество подгонки модели можно оценить при помощи коэффициента детерминации  $R^2$

## Дополнительные ресурсы

- Учебники

- Гланц, 1999, стр. 221-244
- [Open Intro to Statistics: Chapter 7. Introduction to linear regression](#), pp. 315-353.
- Quinn, Keough, 2002, pp. 78-110
- Logan, 2010, pp. 170-207
- Sokal, Rohlf, 1995, pp. 451-491
- Zar, 1999, pp. 328-355

- Упражнения для тренировки

- OpenIntro Labs, Lab 7: Introduction to linear regression (Осторожно, они используют базовую графику а не ggplot)
  - Обычный вариант, упражнения 1—4
  - Интерактивный вариант на [Data Camp](#), до вопроса 4