

Смешанные линейные модели (случайный интерсепт и случайный угол наклона)

Линейные модели...

Марина Варфоломеева, Вадим Хайтов
Осень 2022

Вы узнаете

- Что такое смешанные модели и когда они применяются
- Что такое фиксированные и случайные факторы

Вы сможете

- Рассказать чем фиксированные факторы отличаются от случайных
- Привести примеры факторов, которые могут быть фиксированными или случайными в зависимости от задачи исследования
- Рассказать, что оценивает коэффициент внутриклассовой корреляции и вычислить его для случая с одним случайным фактором
- Подобрать смешанную линейную модель со случайным отрезком и случайным углом наклона в R при помощи методов максимального правдоподобия

“Многоуровневые” данные

Независимость наблюдений

Обычные линейные модели предполагают, что наблюдения должны быть независимы друг от друга.

Но так происходит совсем не всегда.

Независимость наблюдений

Обычные линейные модели предполагают, что наблюдения должны быть независимы друг от друга.

Но так происходит совсем не всегда.

Многоуровневые (multilevel), сгруппированные (clustered) данные

Иногда наблюдения бывают сходны по каким-то признакам:

- измерения в разные периоды времени
 - измерения, сделанные в химической лаборатории в разные дни
- измерения в разных участках пространства
 - урожай на участках одного поля
 - детали, произведенные на одном из нескольких аналогичных станков
- повторные измерения на одних и тех же субъектах
 - измерения до и после какого-то воздействия
- измерения на разных субъектах, которые сами объединены в группы
 - ученики в классах, классы в школах, школы в районах, районы в городах и т.п.

Внутригрупповые корреляции

Детали, произведенные на одном станке будут более похожи, чем детали, сделанные на разных.

Аналогично, у учеников из одного класса будет более похожий уровень подготовки к какому-нибудь предмету, чем у учеников из разных классов.

Таким образом, можно сказать, что есть корреляции значений внутри групп.

Внутригрупповые корреляции

Детали, произведенные на одном станке будут более похожи, чем детали, сделанные на разных.

Аналогично, у учеников из одного класса будет более похожий уровень подготовки к какому-нибудь предмету, чем у учеников из разных классов.

Таким образом, можно сказать, что есть корреляции значений внутри групп.

Последствия для анализа

Игнорировать такую группирующую структуру данных нельзя – можно ошибиться с выводами.

Моделировать группирующие факторы обычными методами тоже нельзя – придется подбирать очень много параметров.

Решение – случайные факторы.

Сейчас давайте на примере убедимся в том, что без случайных факторов бывает сложно справиться с анализом.



Недосып и время реакции

Недосып и время реакции

В ночь перед нулевым днем всем испытуемым давали спать нормальное время, а в следующие 9 ночей — только по 3 часа. Каждый день измеряли время реакции в серии тестов. (Данные Belenky et al., 2003)

Как время реакции людей зависит от бессонницы?

- `Reaction` — среднее время реакции в серии тестов в день наблюдения, мс
- `Days` — число дней депривации сна
- `Subject` — номер испытуемого

```
library(lme4)
data(sleepstudy)

sl <- sleepstudy
str(sl)

'data.frame': 180 obs. of 3 variables:
 $ Reaction: num 250 259 251 321 357 ...
 $ Days     : num 0 1 2 3 4 5 6 7 8 9 ...
 $ Subject  : Factor w/ 18 levels "308","309","310",...: 1 1 1 1 1 1 1 1 1 1 ...
```

Знакомимся с данными

```
# Есть ли пропущенные значения?
```

```
colSums(is.na(sl))
```

```
Reaction      Days   Subject  
      0          0         0
```

```
# Сколько субъектов?
```

```
length(unique(sl$Subject))
```

```
[1] 18
```

```
# Сколько наблюдений для каждого субъекта?
```

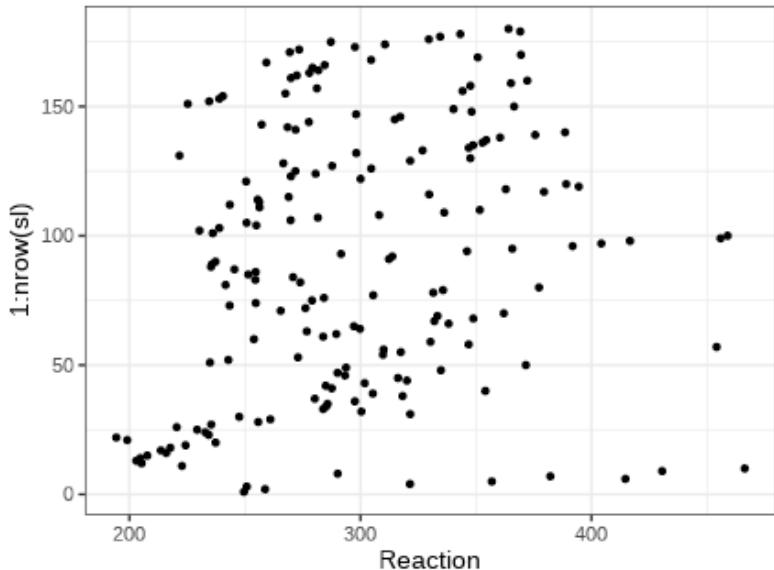
```
table(sl$Subject)
```

```
308 309 310 330 331 332 333 334 335 337 349 350 351 352 369 370 371 372  
10  10  10  10  10  10  10  10  10  10  10  10  10  10  10  10  10  10  10
```

Есть ли выбросы?

```
library(ggplot2)
theme_set(theme_bw(base_size = 14))

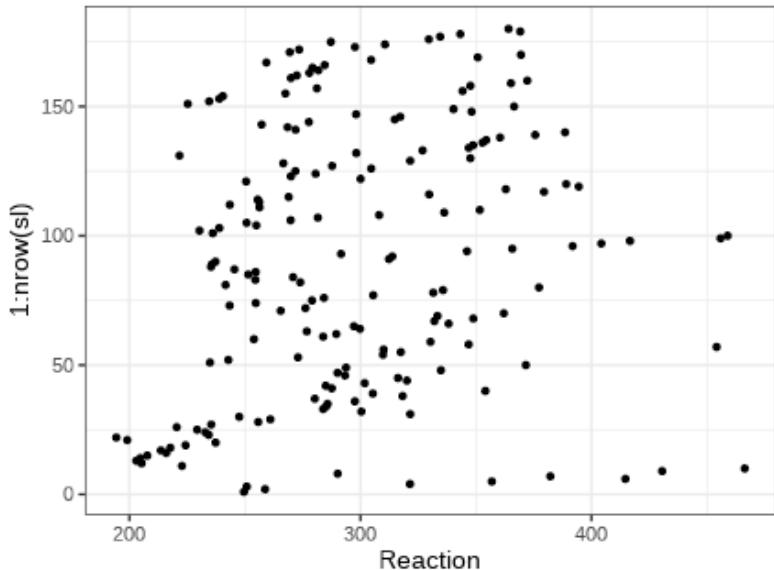
ggplot(sl, aes(x = Reaction, y = 1:nrow(sl))) +
  geom_point()
```



Есть ли выбросы?

```
library(ggplot2)
theme_set(theme_bw(base_size = 14))

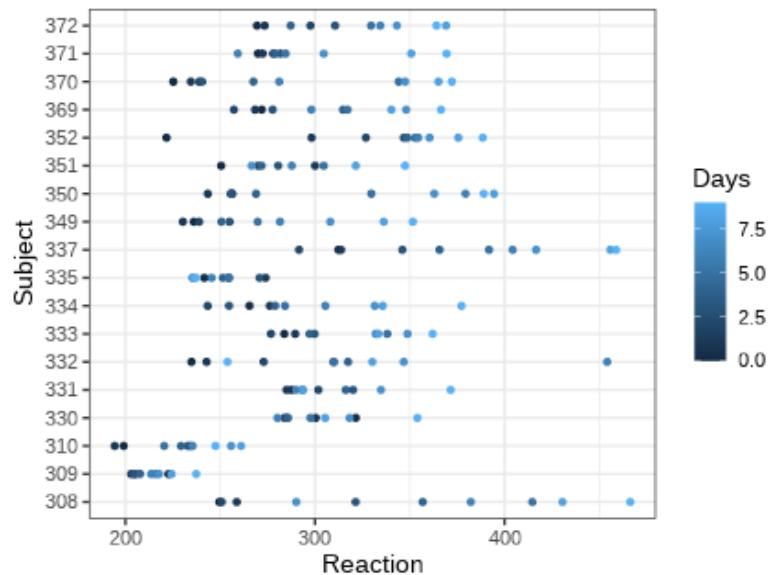
ggplot(sl, aes(x = Reaction, y = 1:nrow(sl))) +
  geom_point()
```



Мы пока еще не учли информацию о субъектах...

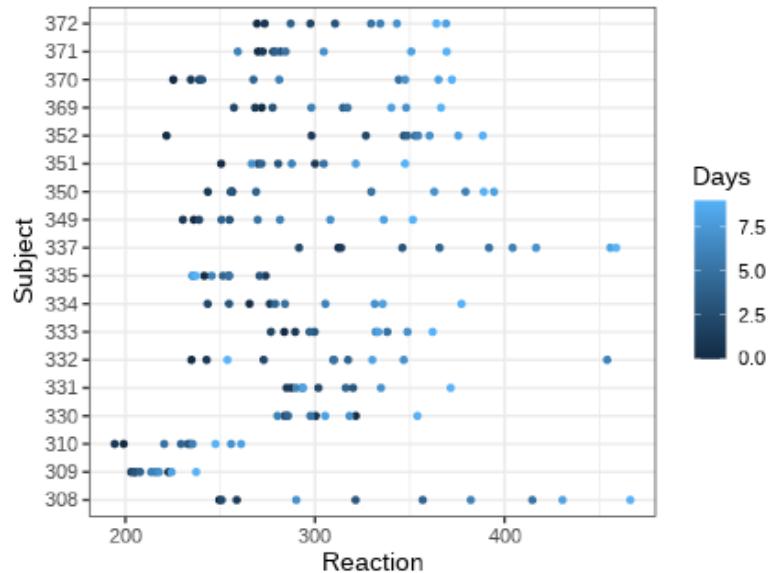
Как меняется время реакции разных людей?

```
ggplot(sl, aes(x = Reaction, y = Subject, colour = Days)) +  
  geom_point()
```



Как меняется время реакции разных людей?

```
ggplot(sl, aes(x = Reaction, y = Subject, colour = Days)) +  
  geom_point()
```



- У разных людей разное время реакции.
- Межиндивидуальную изменчивость нельзя игнорировать.

Что делать с разными субъектами?

Что делать с разными субъектами?



The Good — подбираем смешанную модель, в которой есть фиксированный фактор *Days* и случайный фактор *Subject*, который опишет межиндивидуальную изменчивость.

Что делать с разными субъектами?



The Good — подбираем смешанную модель, в которой есть фиксированный фактор `Days` и случайный фактор `Subject`, который опишет межиндивидуальную изменчивость.



The Bad — игнорируем структуру данных, подбираем модель с единственным фиксированным фактором `Days`. (Не учитываем группирующий фактор `Subject`). Неправильный вариант.

Что делать с разными субъектами?



The Good — подбираем смешанную модель, в которой есть фиксированный фактор `Days` и случайный фактор `Subject`, который опишет межиндивидуальную изменчивость.



The Bad — игнорируем структуру данных, подбираем модель с единственным фиксированным фактором `Days`. (Не учитываем группирующий фактор `Subject`). Неправильный вариант.



The Ugly — подбираем модель с двумя фиксированными факторами: `Days` и `Subject`. (Группирующий фактор `Subject` описывает межиндивидуальную изменчивость как обычный фиксированный фактор).

Плохое решение: не учитываем группирующий фактор

$$Reaction_i = \beta_0 + \beta_1 Days_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma)$$

Плохое решение: не учитываем группирующий фактор

$$Reaction_i = \beta_0 + \beta_1 Days_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma)$$

В матричном виде это можно записать так:

$$\begin{pmatrix} Reaction_1 \\ Reaction_2 \\ \vdots \\ Reaction_{180} \end{pmatrix} = \begin{pmatrix} 1 & Days_1 \\ 1 & Days_2 \\ \vdots & \\ 1 & Days_{180} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{180} \end{pmatrix}$$

что можно сокращенно записать так:

$$\mathbf{Reaction} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Плохое решение: не учитываем группирующий фактор

```
W1 <- glm(Reaction ~ Days, data = sl)
summary(W1)
```

```
Call:
glm(formula = Reaction ~ Days, data = sl)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-110.85 -27.48    1.55   26.14  139.95 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept)  251.41     6.61   38.03 < 2e-16 ***
Days         10.47     1.24    8.45  9.9e-15 ***
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 2277)

Null deviance: 567954  on 179  degrees of freedom
Residual deviance: 405252  on 178  degrees of freedom
AIC: 1906

Number of Fisher Scoring iterations: 2
```

Плохое решение: не учитываем группирующий фактор

```
W1 <- glm(Reaction ~ Days, data = sl)
summary(W1)
```

```
Call:
glm(formula = Reaction ~ Days, data = sl)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-110.85 -27.48    1.55   26.14  139.95 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 251.41      6.61   38.03 < 2e-16 ***
Days         10.47      1.24    8.45  9.9e-15 ***
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 2277)

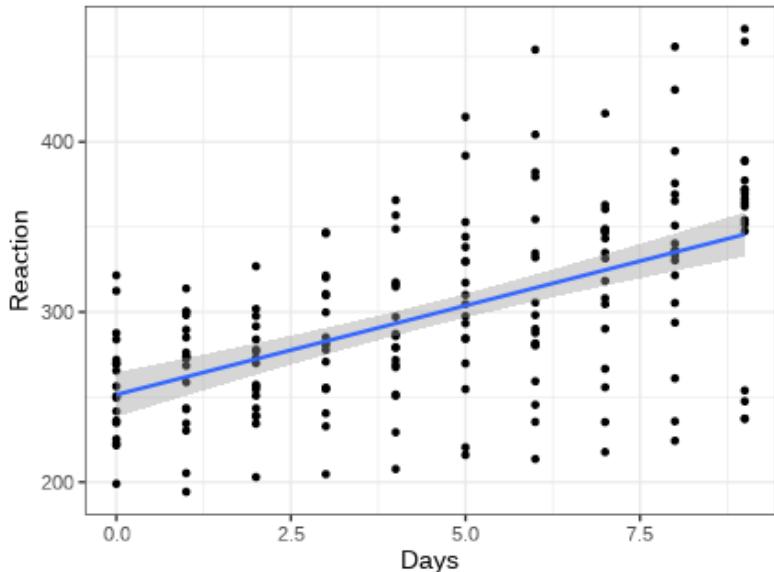
Null deviance: 567954  on 179  degrees of freedom
Residual deviance: 405252  on 178  degrees of freedom
AIC: 1906

Number of Fisher Scoring iterations: 2
```

- Объем выборки завышен (180 наблюдений, вместо 18 субъектов). Стандартные ошибки и уровни значимости занижены. Увеличивается вероятность ошибок I рода.
- Нарушено условие независимости наблюдений.

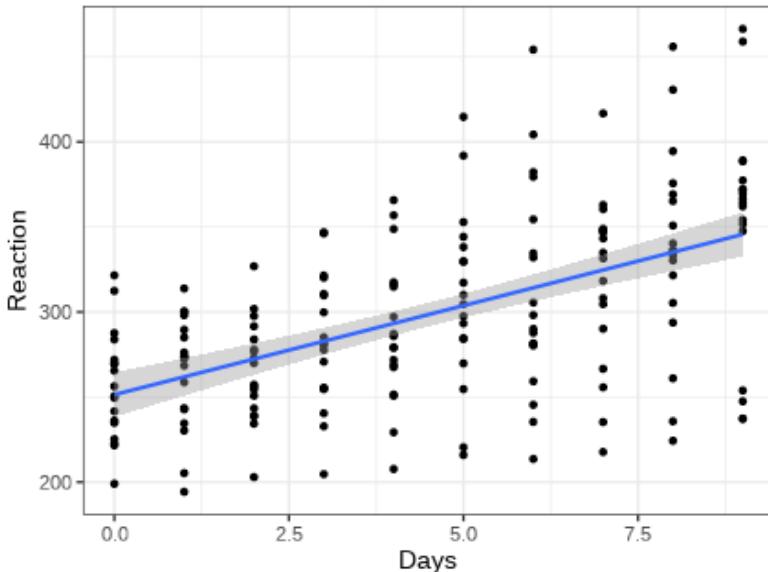
Плохое решение: не учитываем группирующий фактор

```
ggplot(sl, aes(x = Days, y = Reaction)) +  
  geom_point() +  
  geom_smooth(se = TRUE, method = "lm", size = 1)
```



Плохое решение: не учитываем группирующий фактор

```
ggplot(sl, aes(x = Days, y = Reaction)) +  
  geom_point() +  
  geom_smooth(se = TRUE, method = "lm", size = 1)
```



- Доверительная зона регрессии “заужена”.
- Большие остатки, т.к. неучтенная межиндивидуальная изменчивость “ушла” в остаточную.

Громоздкое решение: группирующий фактор как фиксированный

$$Reaction_i = \beta_0 + \beta_1 Days_i + \beta_2 Subj_{2\,i} + \dots + \beta_{18} Subj_{18\,i} + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma)$$

Громоздкое решение: группирующий фактор как фиксированный

$$Reaction_i = \beta_0 + \beta_1 Days_i + \beta_2 Subj_{2\,i} + \dots + \beta_{18} Subj_{18\,i} + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma)$$

В матричном виде это можно записать так:

$$\begin{pmatrix} Reaction_1 \\ Reaction_2 \\ \vdots \\ Reaction_{180} \end{pmatrix} = \begin{pmatrix} 1 & Days_1 & Subj_{2\,1} & \dots & Subj_{18\,1} \\ 1 & Days_2 & Subj_{2\,2} & \dots & Subj_{18\,2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & Days_{180} & Subj_{2\,180} & \dots & Subj_{18\,180} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{18} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{180} \end{pmatrix}$$

То есть: $\mathbf{Reaction} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

Громоздкое решение: группирующий фактор как фиксированный

```
W2 <- glm(Reaction ~ Days + Subject, data = sl)
coef(W2)
```

	Days	Subject309	Subject310
295.031	10.467	-126.901	-111.133
Subject330	Subject331	Subject332	Subject333
-38.912	-32.698	-34.832	-25.976
Subject334	Subject335	Subject337	Subject349
-46.832	-92.064	33.587	-66.299
Subject350	Subject351	Subject352	Subject369
-28.531	-52.036	-4.712	-36.099
Subject370	Subject371	Subject372	
-50.432	-47.150	-24.248	

Фрагмент `summary(W2)`:

```
Residual standard error: 30.99 on 161 degrees of freedom
Multiple R-squared:  0.7277,    Adjusted R-squared:  0.6973
F-statistic: 23.91 on 18 and 161 DF,  p-value: < 2.2e-16
```

Громоздкое решение: группирующий фактор как фиксированный

```
W2 <- glm(Reaction ~ Days + Subject, data = sl)
coef(W2)
```

	Days	Subject309	Subject310
295.031	10.467	-126.901	-111.133
Subject330	Subject331	Subject332	Subject333
-38.912	-32.698	-34.832	-25.976
Subject334	Subject335	Subject337	Subject349
-46.832	-92.064	33.587	-66.299
Subject350	Subject351	Subject352	Subject369
-28.531	-52.036	-4.712	-36.099
Subject370	Subject371	Subject372	
-50.432	-47.150	-24.248	

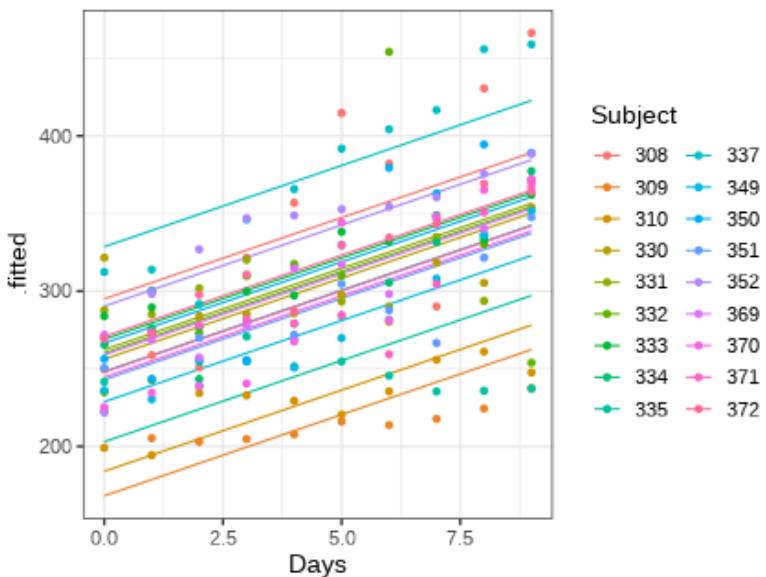
Фрагмент `summary(W2)`:

```
Residual standard error: 30.99 on 161 degrees of freedom
Multiple R-squared:  0.7277,    Adjusted R-squared:  0.6973
F-statistic: 23.91 on 18 and 161 DF,  p-value: < 2.2e-16
```

- 20 параметров (18 для `Subject`, один для `Days` и σ), а наблюдений всего 180.
- Нужно минимум 10–20 наблюдений на каждый параметр (Harrell, 2013) — у нас всего 9.

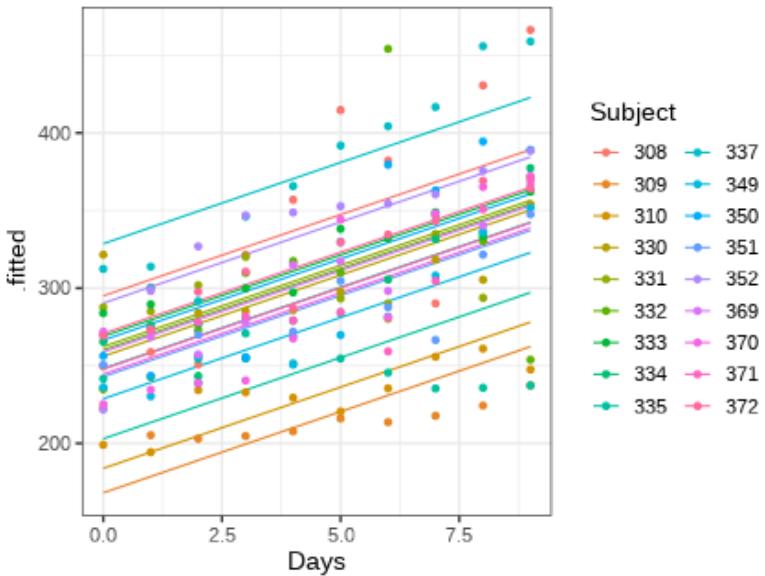
Громоздкое решение: что нам делать с этим множеством прямых?

```
ggplot(fortify(W2), aes(x = Days, colour = Subject)) +  
  geom_line(aes(y = .fitted, group = Subject)) +  
  geom_point(data = sl, aes(y = Reaction)) +  
  guides(colour = guide_legend(ncol = 2))
```



Громоздкое решение: что нам делать с этим множеством прямых?

```
ggplot(fortify(W2), aes(x = Days, colour = Subject)) +  
  geom_line(aes(y = .fitted, group = Subject)) +  
  geom_point(data = sl, aes(y = Reaction)) +  
  guides(colour = guide_legend(ncol = 2))
```



- В модели, где субъект — фиксированный фактор, для каждого субъекта есть “поправка” для значения свободного члена в уравнении регрессии.
- Универсальность модели теряется: предсказания можно сделать только с учетом субъекта.

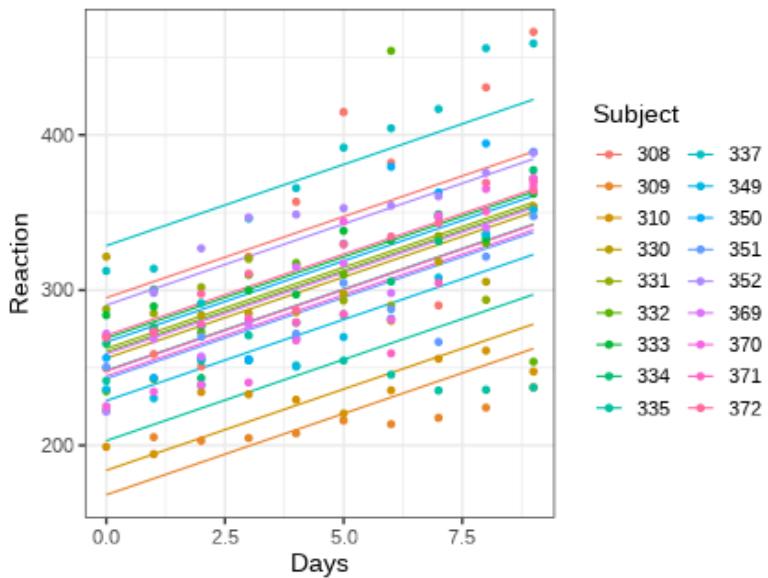
Фиксированные и случайные факторы

Фиксированные факторы

До сих пор мы имели дело только с фиксированными факторами.

Мы моделировали средние значения для уровней фиксированного фактора. Если групп было много, то приходилось моделировать много средних значений.

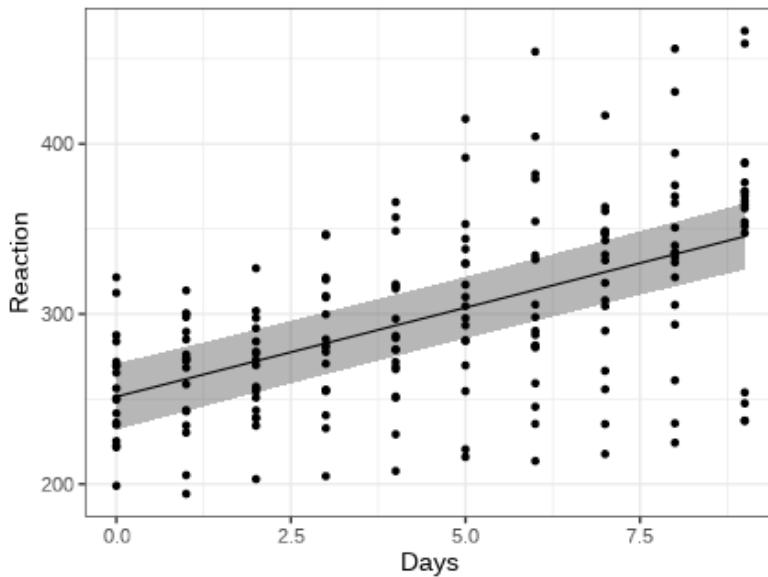
Поступая так, мы считали, что сравниваемые группы – фиксированные, и нам интересны именно сравнения между ними.



Можно посмотреть на группирующий фактор иначе!

Когда нам не важны конкретные значения интерсептов для разных уровней фактора, мы можем представить, что эффект фактора (величина “поправки”) — случайная величина, и можем оценить дисперсию между уровнями группирующего фактора.

Такие факторы называются **случайными факторами**.



Случайные факторы

- измерения в разные периоды времени
 - измерения, сделанные в химической лаборатории в разные дни
- измерения в разных участках пространства
 - урожай на участках одного поля
 - детали, произведенные на одном из нескольких аналогичных станков
- повторные измерения на одних и тех же субъектах
 - измерения до и после какого-то воздействия
- измерения на разных субъектах, которые сами объединены в группы
 - ученики в классах, классы в школах, школы в районах, районы в городах и т.п.

Случайные факторы в моделях

На один и тот же фактор можно посмотреть и как на фиксированный и как на случайный в зависимости от целей исследователя.

Поскольку моделируя случайный фактор мы оцениваем дисперсию между уровнями, то хорошо, если у случайного фактора будет минимум пять градаций.

GLMM со случайным отрезком

GLMM со случайным отрезком

$$Reaction_{ij} = \beta_0 + \beta_1 Days_{ij} + b_i + \varepsilon_{ij}$$

$b_i \sim N(0, \sigma_b)$ — случайный эффект субъекта (случайный отрезок, интерсепт)

$\varepsilon_{ij} \sim N(0, \sigma)$ — остатки модели

i — субъекты, j — дни

GLMM со случайным отрезком

$$Reaction_{ij} = \beta_0 + \beta_1 Days_{ij} + b_i + \varepsilon_{ij}$$

$b_i \sim N(0, \sigma_b)$ — случайный эффект субъекта (случайный отрезок, интерсепт)

$\varepsilon_{ij} \sim N(0, \sigma)$ — остатки модели

i — субъекты, j — дни

В матричном виде это записывается так:

$$\begin{pmatrix} Reaction_1 \\ Reaction_2 \\ \vdots \\ Reaction_{180} \end{pmatrix} = \begin{pmatrix} 1 & Days_1 \\ 1 & Days_2 \\ \vdots & \vdots \\ 1 & Days_{180} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} (b) + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{180} \end{pmatrix}$$

То есть: $\mathbf{Reaction} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$

Подберем модель со случайным отрезком

Используем `lmer` из пакета `lme4`.

```
?lmer # справка о lmer
```

`lmer` по умолчанию использует REML для подбора параметров. Это значит, что случайные эффекты будут оценены более точно, чем при использовании ML.

```
M1 <- lmer(Reaction ~ Days + (1 | Subject), data = sl)
```

Уравнение модели со случайным отрезком

```
summary(M1)
```

```
Linear mixed model fit by REML ['lmerMod']
Formula: Reaction ~ Days + (1 | Subject)
Data: sl
```

```
REML criterion at convergence: 1786
```

```
Scaled residuals:
```

Min	1Q	Median	3Q	Max
-3.226	-0.553	0.011	0.519	4.251

```
Random effects:
```

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	1378	37.1
Residual		960	31.0

```
Number of obs: 180, groups: Subject, 18
```

```
Fixed effects:
```

	Estimate	Std. Error	t value
(Intercept)	251.405	9.747	25.8
Days	10.467	0.804	13.0

```
Correlation of Fixed Effects:
```

```
(Intr)
```

Уравнение модели со случайным отрезком

```
summary(M1)
```

```
Linear mixed model fit by REML ['lmerMod']
Formula: Reaction ~ Days + (1 | Subject)
Data: sl
```

```
REML criterion at convergence: 1786
```

```
Scaled residuals:
```

Min	1Q	Median	3Q	Max
-3.226	-0.553	0.011	0.519	4.251

```
Random effects:
```

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	1378	37.1
Residual		960	31.0

```
Number of obs: 180, groups: Subject, 18
```

```
Fixed effects:
```

	Estimate	Std. Error	t value
(Intercept)	251.405	9.747	25.8
Days	10.467	0.804	13.0

```
Correlation of Fixed Effects:
```

```
(Intr)
```

$$Reaction_{ij} = 251.4 + 10.5 Days_{ij} + b_i + \varepsilon_{ij}$$

$b_i \sim N(0, 37.12)$ — случайный эффект субъекта

$\varepsilon_{ij} \sim N(0, 30.99)$ — остатки модели

i — субъекты, j — дни

Предсказания смешанных моделей бывают двух типов

- Предсказания с учетом лишь фиксированных эффектов,
- Предсказания с одновременным учетом как фиксированных, так и случайных эффектов.

Данные для графика предсказаний фиксированной части модели:

```
library(dplyr)
NewData <- sl %>% group_by(Subject) %>%
  do(data.frame(Days = seq(min(. $Days), max(. $Days), length = 10)))
head(NewData, 3)

# A tibble: 3 × 2
# Groups:   Subject [1]
  Subject  Days
  <fct>    <dbl>
1 308      0
2 308      1
3 308      2
```

Предсказания фиксированной части модели при помощи predict()

```
?predict.merMod
```

Функция `predict()` в `lme4` не считает стандартные ошибки и доверительные интервалы. Это потому, что нет способа адекватно учесть неопределенность, связанную со случайными эффектами.

```
NewData$fit <- predict(M1, NewData, type = 'response', re.form = NA)  
head(NewData, 3)
```

```
# A tibble: 3 × 3  
# Groups:   Subject [1]  
  Subject  Days    fit  
  <fct>    <dbl> <dbl>  
1 308        0 251.  
2 308        1 262.  
3 308        2 272.
```

Предсказания фиксированной части модели в матричном виде

Стандартные ошибки, рассчитанные обычным методом, позволяют получить **приблизительные** доверительные интервалы.

```
# Предсказанные значения при помощи матриц
X <- model.matrix(~ Days, data = NewData)
betas <- fixef(M1)
NewData$fit <- X %*% betas

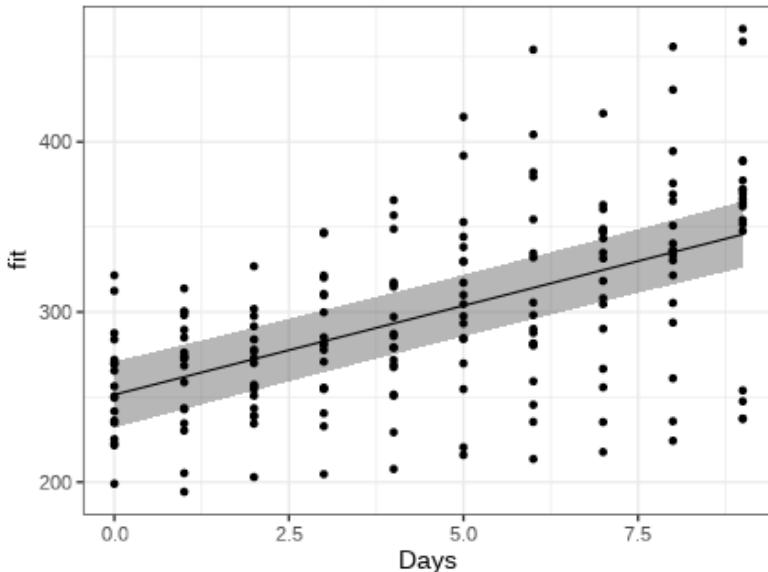
# Стандартные ошибки
NewData$SE <- sqrt( diag(X %*% vcov(M1) %*% t(X)) )

NewData$lwr <- NewData$fit - 2 * NewData$SE
NewData$upr <- NewData$fit + 2 * NewData$SE
```

Более точные доверительные интервалы можно получить при помощи бутстрепа. Мы сделаем это позже для финальной модели.

График предсказаний фиксированной части модели

```
ggplot(data = NewData, aes(x = Days, y = fit)) +  
  geom_ribbon(alpha = 0.35, aes(ymin = lwr, ymax = upr)) +  
  geom_line() + geom_point(data = sl, aes(x = Days, y = Reaction))
```

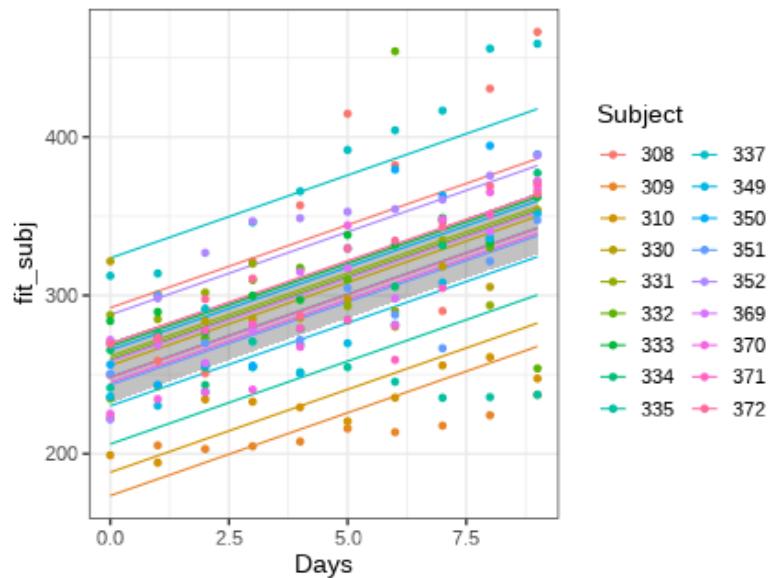


Зависимость времени реакции от продолжительности периода бессонницы без учета субъекта:

$$\widehat{Reaction}_{ij} = 251.4 + 10.5 Days_{ij}$$

Предсказания для уровней случайного фактора

```
NewData$fit_subj <- predict(M1, NewData, type = 'response')
ggplot(NewData, aes(x = Days, y = fit_subj)) +
  geom_ribbon(alpha = 0.3, aes(ymin = lwr, ymax = upr)) +
  geom_line(aes(colour = Subject)) +
  geom_point(data = sl, aes(x = Days, y = Reaction, colour = Subject)) +
  guides(colour = guide_legend(ncol = 2))
```



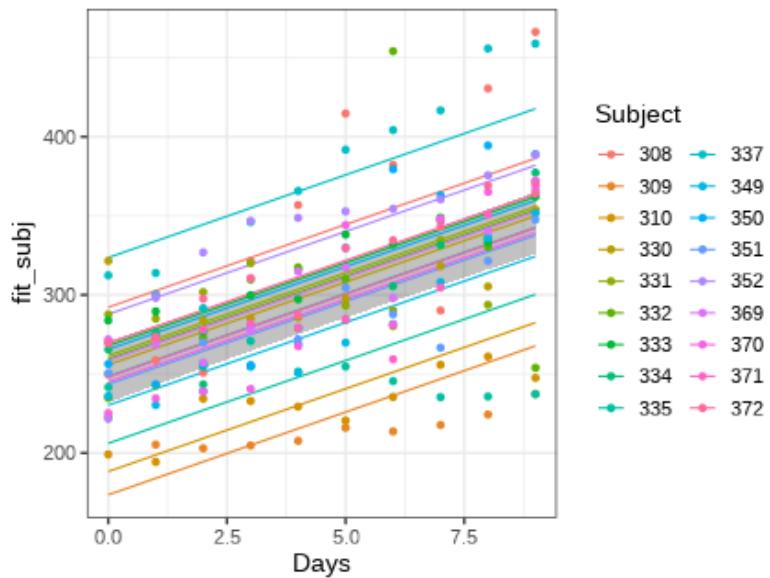
Зависимость времени реакции от продолжительности периода бессонницы для обследованных субъектов:

$$\widehat{Reaction}_{ij} = 251.4 + 10.5 Days_{ij} + b_i$$

Важные замечания

Случайный фактор помогает учесть взаимозависимость наблюдений для каждого из субъектов – “индуцированные” корреляции.

После анализа остатков модели можно будет понять, стоит ли с ней работать дальше. Одна из потенциальных проблем – время реакции разных субъектов может меняться непараллельно. Возможно, модель придется переформулировать.



Индукционная корреляция

Разберемся со случайной частью модели

$$\text{Reaction} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon}$$

$\mathbf{b} \sim N(0, \mathbf{D})$ - случайные эффекты b_i нормально распределены со средним 0 и матрицей ковариаций \mathbf{D}

$\boldsymbol{\epsilon} \sim N(0, \boldsymbol{\Sigma})$ - остатки модели нормально распределены со средним 0 и матрицей ковариаций $\boldsymbol{\Sigma}$

Разберемся со случайной частью модели

$$\text{Reaction} = \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \epsilon$$

$\mathbf{b} \sim N(0, \mathbf{D})$ - случайные эффекты b_i нормально распределены со средним 0 и матрицей ковариаций \mathbf{D}

$\epsilon \sim N(0, \Sigma)$ - остатки модели нормально распределены со средним 0 и матрицей ковариаций Σ

Матрица ковариаций остатков: $\Sigma = \sigma^2 \cdot \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} = \begin{pmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{pmatrix}$

Остатки модели должны быть независимы друг от друга

В матрице ковариаций остатков вне диагонали стоят нули, т.е. ковариация разных остатков равна нулю. Т.е. остатки независимы друг от друга.

$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

В то же время, отдельные значения переменной-отклика \mathbf{Y} уже не будут независимы друг от друга при появлении в модели случайного фактора.

Матрица ковариаций переменной-отклика

$$\text{Reaction} = \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \varepsilon$$

$$\mathbf{b} \sim N(0, \mathbf{D})$$

$$\varepsilon \sim N(0, \Sigma)$$

Можно показать, что переменная-отклик \mathbf{Y} нормально распределена:

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \mathbf{V})$$

Матрица ковариаций переменной-отклика

$$\text{Reaction} = \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \varepsilon$$

$$\mathbf{b} \sim N(0, \mathbf{D})$$

$$\varepsilon \sim N(0, \Sigma)$$

Можно показать, что переменная-отклик \mathbf{Y} нормально распределена:

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \mathbf{V})$$

Матрица ковариаций переменной-отклика:

$$\mathbf{V} = \mathbf{Z}\mathbf{D}\mathbf{Z}' + \Sigma, \text{ где } \mathbf{D} \text{ — матрица ковариаций случайных эффектов.}$$

Т.е. добавление случайных эффектов приводит
к изменению ковариационной матрицы \mathbf{V}

Добавление случайных эффектов приводит к изменению ковариационной матрицы

$$\mathbf{V} = \mathbf{ZDZ}' + \Sigma$$

Для простейшей смешанной модели со случайным отрезком:

$$\begin{aligned}\mathbf{V} &= \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \sigma_b^2 \cdot (1 \ 1 \ \cdots \ 1) + \sigma^2 \cdot \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} = \\ &= \begin{pmatrix} \sigma^2 + \sigma_b^2 & \sigma_b^2 & \cdots & \sigma_b^2 \\ \sigma_b^2 & \sigma^2 + \sigma_b^2 & \cdots & \sigma_b^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma^2 + \sigma_b^2 \end{pmatrix}\end{aligned}$$

Индуцированная корреляция – следствие включения в модель случайных эффектов

$$\mathbf{V} = \begin{pmatrix} \sigma^2 + \sigma_b^2 & \sigma_b^2 & \cdots & \sigma_b^2 \\ \sigma_b^2 & \sigma^2 + \sigma_b^2 & \cdots & \sigma_b^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma^2 + \sigma_b^2 \end{pmatrix}$$

- σ_b^2 — ковариация между наблюдениями одного субъекта.
- $\sigma^2 + \sigma_b^2$ — дисперсия.

Т.е. корреляция между наблюдениями одного субъекта $\sigma_b^2 / (\sigma^2 + \sigma_b^2)$

Коэффициент внутриклассовой корреляции (intra-class correlation, ICC)

$$ICC = \sigma_b^2 / (\sigma^2 + \sigma_b^2)$$

Способ измерить, насколько коррелируют друг с другом наблюдения из одной и той же группы, заданной случайным фактором.

Если ICC низок, то наблюдения очень разные внутри групп, заданных уровнями случайного фактора.

Если ICC высок, то наблюдения очень похожи внутри каждой из групп, заданных случайногм фактором.

ICC можно использовать как описательную статистику, но интереснее использовать его при планировании исследований.

Если в пилотном исследовании ICC маленький, то для надежной оценки эффекта этого случайного фактора нужно брать больше наблюдений в группе.

Если в пилотном исследовании ICC большой, то можно брать меньше наблюдений в группе.

Вычислим коэффициент внутриклассовой корреляции

```
summary(M1)
```

```
Linear mixed model fit by REML ['lmerMod']
Formula: Reaction ~ Days + (1 | Subject)
Data: sl
```

```
REML criterion at convergence: 1786
```

```
Scaled residuals:
```

Min	1Q	Median	3Q	Max
-3.226	-0.553	0.011	0.519	4.251

```
Random effects:
```

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	1378	37.1
	Residual	960	31.0

```
Number of obs: 180, groups: Subject, 18
```

```
Fixed effects:
```

	Estimate	Std. Error	t value
(Intercept)	251.405	9.747	25.8
Days	10.467	0.804	13.0

```
Correlation of Fixed Effects:
```

(Intr)	Days
-0.371	

```
# Случайные эффекты отдельно
VarCorr(M1)
```

Groups	Name	Std.Dev.
Subject	(Intercept)	37.1
	Residual	31.0

```
37.124^2 / (37.124^2 + 30.991^2)
```

```
[1] 0.5893
```

Диагностика модели

Условия применимости

- Случайность и независимость наблюдений.
- Линейная связь.
- Нормальное распределение остатков.
- Гомогенность дисперсий остатков.
- Отсутствие коллинеарности предикторов.

Данные для анализа остатков

```
M1_diag <- data.frame(  
  sl,  
  .fitted = predict(M1),  
  .resid = resid(M1, type = 'pearson'),  
  .scresid = resid(M1, type = 'pearson', scaled = TRUE))  
  
head(M1_diag, 4)
```

	Reaction	Days	Subject	.fitted	.resid	.scresid
1	249.6	0	308	292.2	-42.629	-1.3755
2	258.7	1	308	302.7	-43.951	-1.4182
3	250.8	2	308	313.1	-62.323	-2.0110
4	321.4	3	308	323.6	-2.151	-0.0694

- `.fitted` — предсказанные значения.
- `.resid` — Пирсоновские остатки.
- `.scresid` — стандартизованные Пирсоновские остатки.

График остатков от предсказанных значений

```
gg_resid <- ggplot(M1_diag, aes(y = .scresid)) +  
  geom_hline(yintercept = 0)  
gg_resid + geom_point(aes(x = .fitted))
```

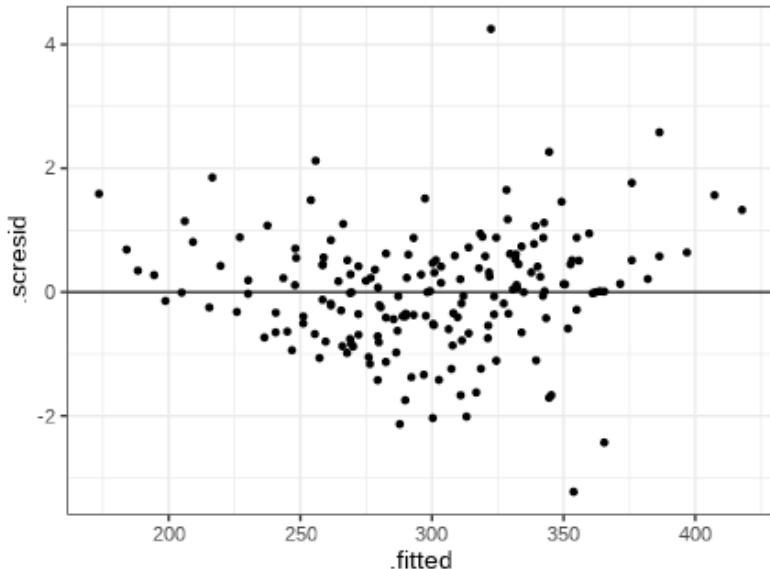
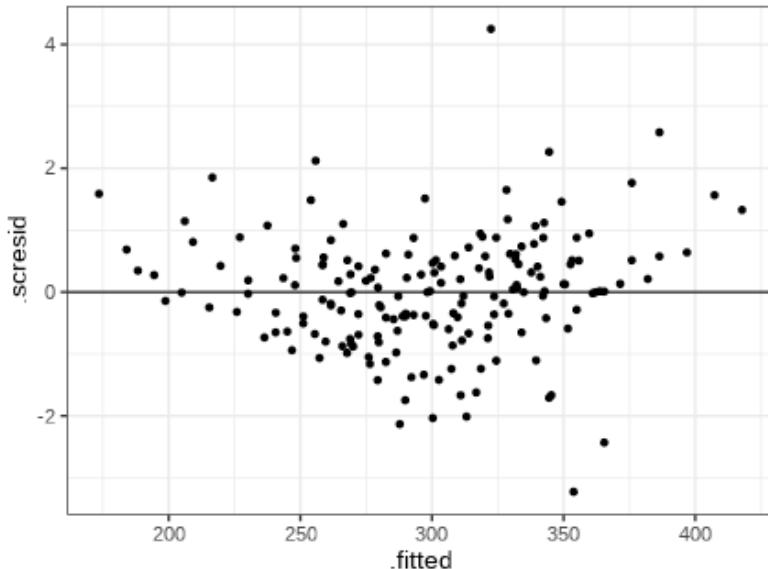


График остатков от предсказанных значений

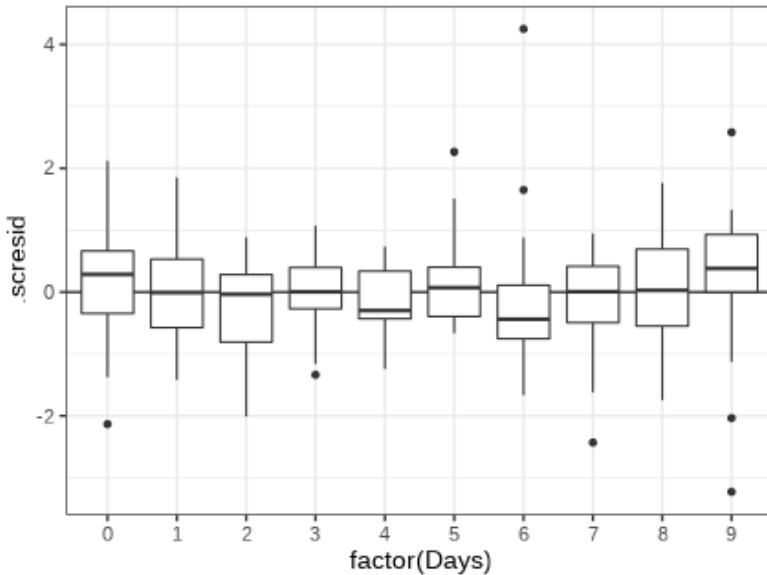
```
gg_resid <- ggplot(M1_diag, aes(y = .scresid)) +  
  geom_hline(yintercept = 0)  
gg_resid + geom_point(aes(x = .fitted))
```



- Большие остатки.
- Гетерогенность дисперсий.

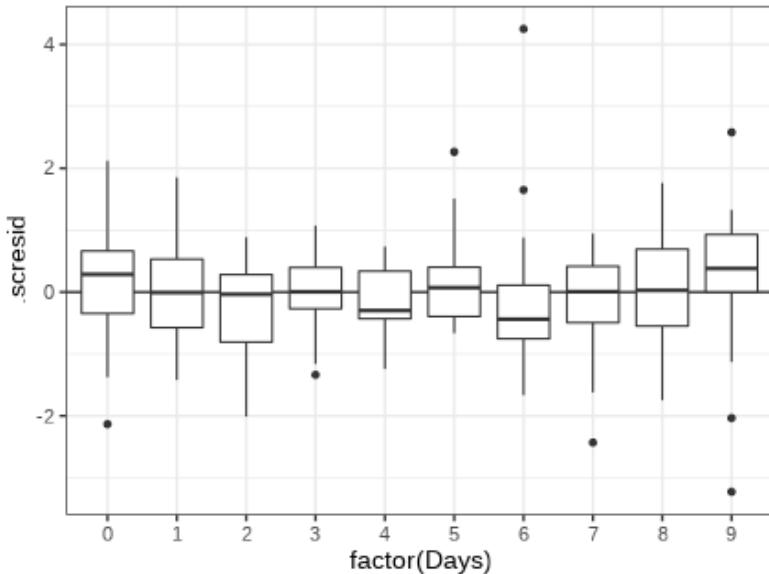
Графики остатков от ковариат в модели и не в модели

```
gg_resid + geom_boxplot(aes(x = factor(Days)))
```



Графики остатков от ковариат в модели и не в модели

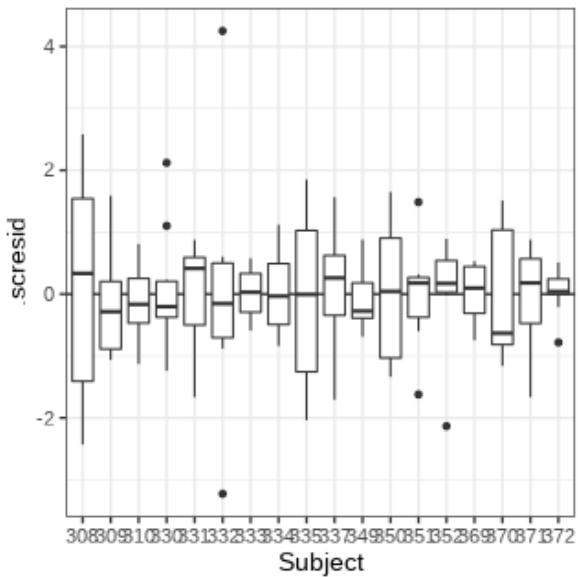
```
gg_resid + geom_boxplot(aes(x = factor(Days)))
```



- Большие остатки в некоторые дни.
- Гетерогенность дисперсий.

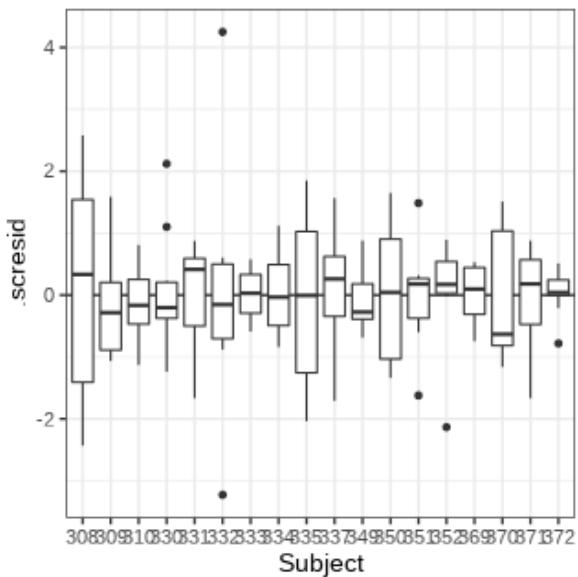
Графики остатков от ковариат в модели и не в модели

```
gg_resid + geom_boxplot(aes(x = Subject))
```



Графики остатков от ковариат в модели и не в модели

```
gg_resid + geom_boxplot(aes(x = Subject))
```

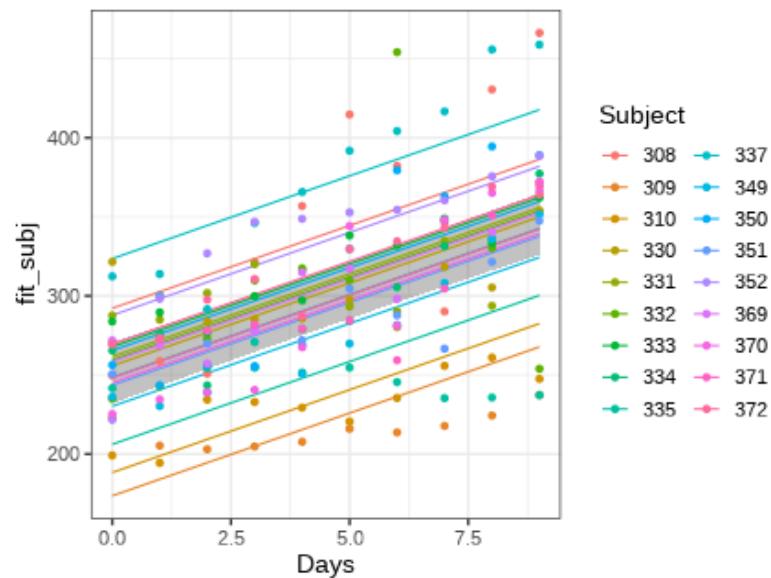


- Большие остатки у 332 субъекта.
 - Гетерогенность дисперсий.

GLMM со случайным отрезком и углом наклона

GLMM со случайным отрезком и углом наклона

На графике индивидуальных эффектов было видно, что измерения для разных субъектов, возможно, идут непараллельно. Усложним модель — добавим случайные изменения угла наклона для каждого из субъектов.



Это можно биологически объяснить. Возможно, в зависимости от продолжительности бессонницы у разных субъектов скорость реакции будет ухудшаться разной скоростью: одни способны выдержать 9 дней почти без потерь, а другим уже пары дней может быть достаточно.

Уравнение модели со случайным отрезком и углом наклона

$$Reaction_{ij} = \beta_0 + \beta_1 Days_{ij} + b_i + c_{ij} Days_{ij} + \varepsilon_{ij}$$

$b_i \sim N(0, \sigma_b)$ — случайный интерсепт для субъекта

$c_{ij} \sim N(0, \sigma_c)$ — случайный угол наклона для субъекта

$\varepsilon_{ij} \sim N(0, \sigma)$ — остатки модели

i — субъекты, j — дни

Уравнение модели со случайным отрезком и углом наклона

$$Reaction_{ij} = \beta_0 + \beta_1 Days_{ij} + b_i + c_{ij} Days_{ij} + \varepsilon_{ij}$$

$b_i \sim N(0, \sigma_b)$ — случайный интерсепт для субъекта

$c_{ij} \sim N(0, \sigma_c)$ — случайный угол наклона для субъекта

$\varepsilon_{ij} \sim N(0, \sigma)$ — остатки модели

i — субъекты, j — дни

В матричном виде это записывается так:

$$\begin{pmatrix} Reaction_1 \\ Reaction_2 \\ \vdots \\ Reaction_{180} \end{pmatrix} = \begin{pmatrix} 1 & Days_1 \\ 1 & Days_2 \\ \vdots & \\ 1 & Days_{180} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} 1 & Days_1 \\ 1 & Days_2 \\ \vdots & \\ 1 & Days_{180} \end{pmatrix} \begin{pmatrix} b \\ c \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{180} \end{pmatrix}$$

То есть: $\mathbf{Reaction} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$

Подберем модель со случайным отрезком и углом наклона

Формат записи формулы для случайных эффектов в `lme4`

```
(1 + угловой_коэффициент | отрезок)
```

```
MS1 <- lmer(Reaction ~ Days + ( 1 + Days|Subject), data = sl)
```

Уравнение модели со случайным отрезком и углом наклона

```
summary(MS1)
```

```
Linear mixed model fit by REML ['lmerMod']
Formula: Reaction ~ Days + (1 + Days | Subject)
Data: sl
```

```
REML criterion at convergence: 1744
```

```
Scaled residuals:
```

Min	1Q	Median	3Q	Max
-3.954	-0.463	0.023	0.463	5.179

```
Random effects:
```

Groups	Name	Variance	Std.Dev.	Corr
Subject	(Intercept)	612.1	24.74	
	Days	35.1	5.92	0.07
Residual		654.9	25.59	

```
Number of obs: 180, groups: Subject, 18
```

```
Fixed effects:
```

Estimate	Std. Error	t value
----------	------------	---------

$$Reaction_{ij} = 251.4 + 10.5Days_{ij} + b_i + c_{ij}Days_{ij} + \varepsilon_{ij}$$

$b_i \sim N(0, 24.74)$ — случайный интерсепт для субъекта

$c_{ij} \sim N(0, 5.92)$ — случайный угол наклона для субъекта

$\varepsilon_{ij} \sim N(0, 25.59)$ — остатки модели

i — субъекты, j — дни

Данные для графика предсказаний фиксированной части модели

```
library(dplyr)
NewData <- sl %>% group_by(Subject) %>%
  do(data.frame(Days = seq(min(. $Days), max(. $Days), length = 10)))
NewData$fit <- predict(MS1, NewData, type = 'response', re.form = NA)
head(NewData, 3)

# A tibble: 3 × 3
# Groups:   Subject [1]
  Subject  Days    fit
  <fct>    <dbl> <dbl>
1 308        0 251.
2 308        1 262.
3 308        2 272.
```

Предсказания фиксированной части модели в матричном виде

Вычислим **приблизительные** доверительные интервалы.

```
# Предсказанные значения при помощи матриц
X <- model.matrix(~ Days, data = NewData)
betas <- fixef(MS1)
NewData$fit <- X %*% betas

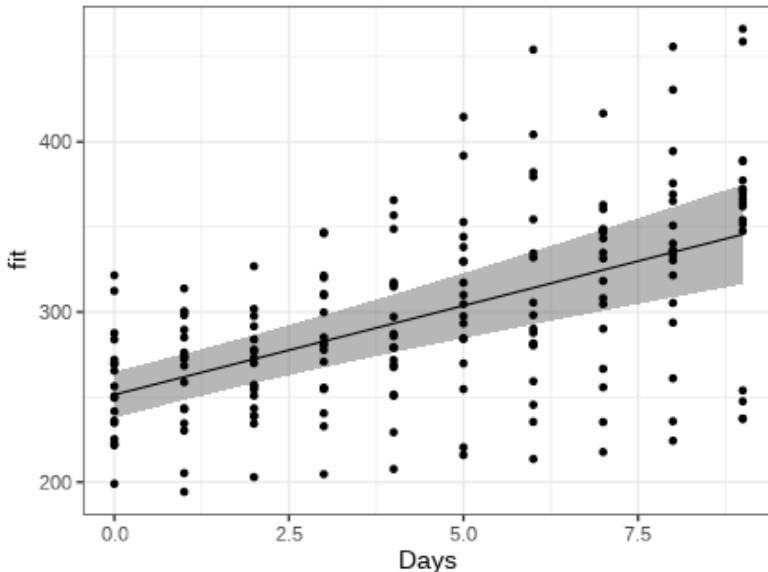
# Стандартные ошибки
NewData$SE <- sqrt( diag(X %*% vcov(MS1) %*% t(X)) )

NewData$lwr <- NewData$fit - 2 * NewData$SE
NewData$upr <- NewData$fit + 2 * NewData$SE
```

Более точные доверительные интервалы можно получить при помощи бутстрепа.

График предсказаний фиксированной части модели

```
ggplot(data = NewData, aes(x = Days, y = fit)) +  
  geom_ribbon(alpha = 0.35, aes(ymin = lwr, ymax = upr)) +  
  geom_line() + geom_point(data = sl, aes(x = Days, y = Reaction))
```

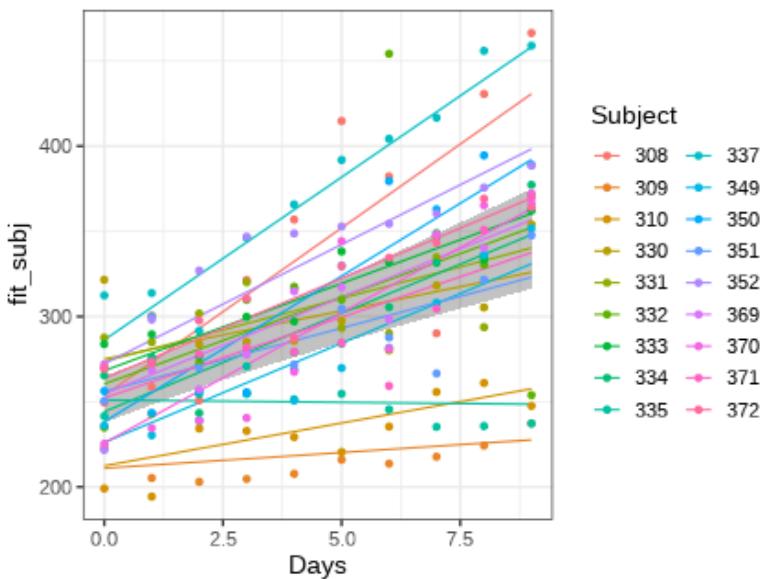


Зависимость времени реакции от продолжительности периода бессонницы без учета субъекта:

$$\widehat{Reaction}_{ij} = 251.4 + 10.5 Days_{ij}$$

Предсказания для уровней случайного фактора

```
NewData$fit_subj <- predict(MS1, NewData, type = 'response')
ggplot(NewData, aes(x = Days, y = fit_subj)) +
  geom_ribbon(alpha = 0.3, aes(ymin = lwr, ymax = upr)) +
  geom_line(aes(colour = Subject)) +
  geom_point(data = sl, aes(x = Days, y = Reaction, colour = Subject)) +
  guides(colour = guide_legend(ncol = 2))
```



Зависимость времени реакции от продолжительности периода бессонницы для обследованных субъектов:

$$\widehat{Reaction}_{ij} = 251.4 + 10.5 Days_{ij} + b_i + c_{ij} Days_{ij}$$

Диагностика модели со случайным отрезком и углом наклона

Данные для анализа остатков

```
MS1_diag <- data.frame(  
  sl,  
  .fitted = predict(MS1),  
  .resid = resid(MS1, type = 'pearson'),  
  .scresid = resid(MS1, type = 'pearson', scaled = TRUE))  
  
head(MS1_diag, 4)
```

	Reaction	Days	Subject	.fitted	.resid	.scresid
1	249.6	0	308	253.7	-4.104	-0.1604
2	258.7	1	308	273.3	-14.625	-0.5715
3	250.8	2	308	293.0	-42.196	-1.6488
4	321.4	3	308	312.7	8.777	0.3430

График остатков от предсказанных значений

```
gg_resid <- ggplot(MS1_diag, aes(y = .scresid)) +  
  geom_hline(yintercept = 0)  
gg_resid + geom_point(aes(x = .fitted))
```

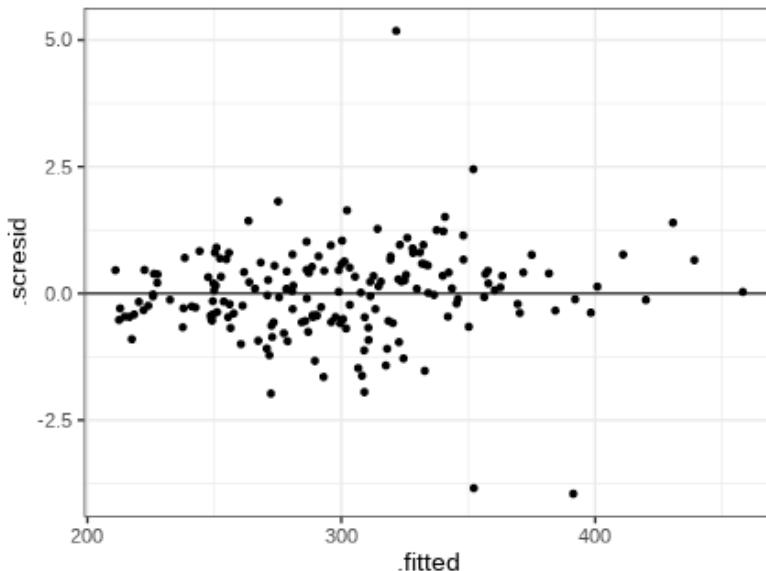
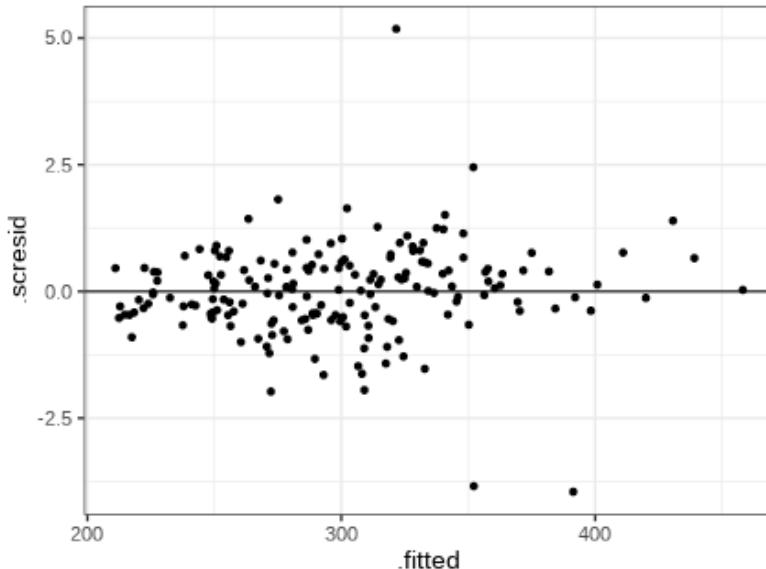


График остатков от предсказанных значений

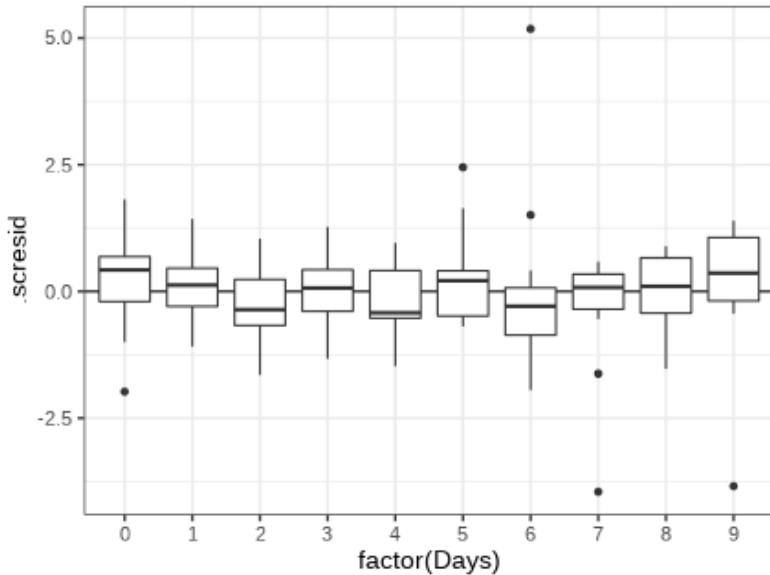
```
gg_resid <- ggplot(MS1_diag, aes(y = .scresid)) +  
  geom_hline(yintercept = 0)  
gg_resid + geom_point(aes(x = .fitted))
```



- Несколько больших остатков.
- Гетерогенность дисперсий не выражена.

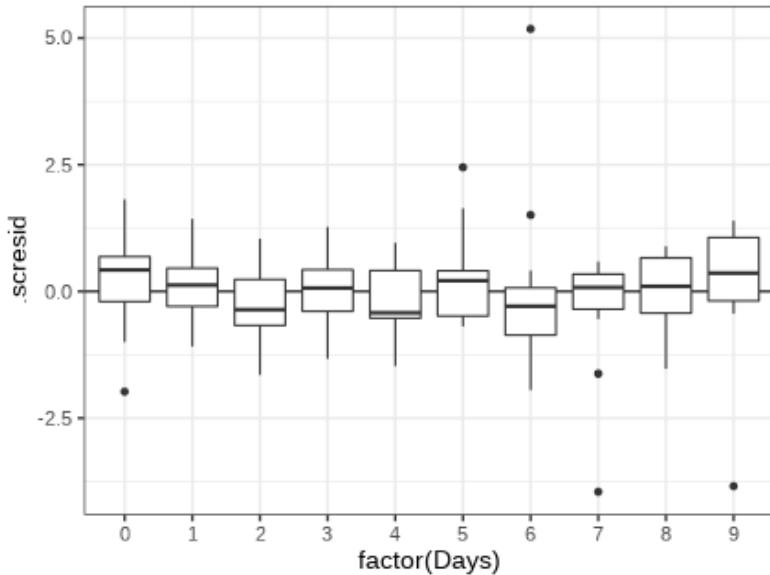
Графики остатков от ковариат в модели и не в модели

```
gg_resid + geom_boxplot(aes(x = factor(Days)))
```



Графики остатков от ковариат в модели и не в модели

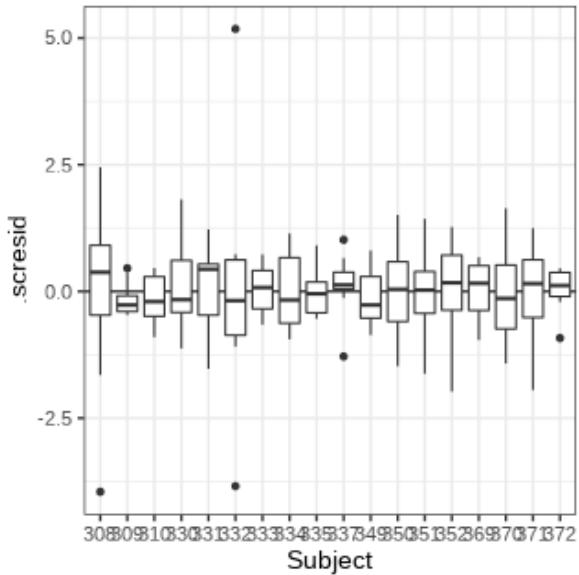
```
gg_resid + geom_boxplot(aes(x = factor(Days)))
```



- Большие остатки в некоторые дни.
- Нет гетерогенности дисперсий остатков.

Графики остатков от ковариат в модели и не в модели

```
gg_resid + geom_boxplot(aes(x = Subject))
```



- Большие остатки у 332 субъекта.
- Гетерогенность дисперсий не выражена.

Смешанные линейные модели

Смешанные модели (Mixed Models)

Смешанными называются модели, включающие случайные факторы.

- Общие смешанные модели (General Linear Mixed Models) — только нормальное распределение зависимой переменной.
- Обобщенные смешанные модели (Generalized Linear Mixed Models) — распределения зависимой переменной могут быть другими (из семейства экспоненциальных распределений).

Смешанная линейная модель в общем виде

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

$\mathbf{b} \sim N(0, \mathbf{D})$ — случайные эффекты нормально распределены со средним 0 и матрицей ковариаций \mathbf{D} (ее диагональные элементы — стандартное отклонение σ_b).

$\boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Sigma})$ — остатки модели нормально распределены со средним 0 и матрицей ковариаций $\boldsymbol{\Sigma}$ (ее диагональные элементы — стандартное отклонение σ).

$\mathbf{X}\boldsymbol{\beta}$ — фиксированная часть модели.

$\mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$ — случайная часть модели.

В зависимости от устройства модельной матрицы для случайных эффектов \mathbf{Z} смешанные модели делят на модели со случайным отрезком и случайным углом наклона.

Методы подбора параметров в смешанных моделях

Метод максимального правдоподобия (Maximum Likelihood, ML)

Метод ограниченного максимального правдоподобия (Restricted Maximum Likelihood, REML)

Метод максимального правдоподобия, ML

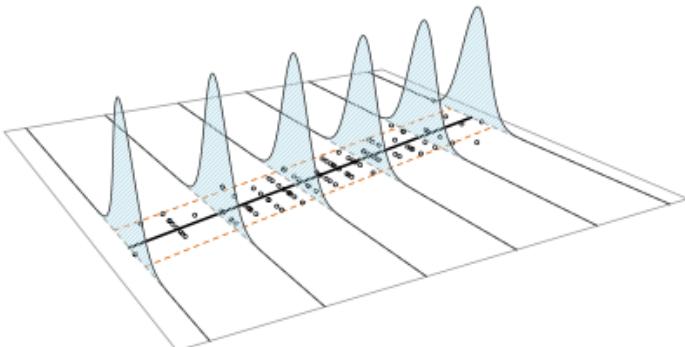
Правдоподобие (likelihood) — измеряет соответствие реально наблюдаемых данных тем, что можно получить из модели при определенных значениях параметров.

Это произведение вероятностей данных:

$$L(\theta|y_i) = \prod_{i=1}^n f(y_i|\theta)$$

Параметры модели должны максимизировать значение логарифма правдоподобия (loglikelihood):

$$\ln L(\theta|y_i) \rightarrow \max$$



ML-оценки для дисперсий – смещенные

Например, ML оценка обычной выборочной дисперсии будет смещенной:

$$\hat{\sigma}^2 = \frac{\sum(x_i - \bar{x})^2}{n},$$

т.к. в знаменателе не $n - 1$, а n .

Аналогичные проблемы возникают при ML оценках для случайных эффектов в линейных моделях.

Это происходит потому, что в смешанной линейной модели

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}), \text{ где } \mathbf{V} = \mathbf{Z}\mathbf{D}\mathbf{Z}' + \boldsymbol{\Sigma}$$

т.е. одновременно приходится оценивать $\boldsymbol{\beta}$ и \mathbf{V} .

Метод ограниченного максимального правдоподобия, REML

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

$$Y \sim N(\mathbf{X}\beta, \mathbf{V}), \text{ где } \mathbf{ZDZ}' + \Sigma$$

Если найти матрицу \mathbf{A} , ортогональную к \mathbf{X}' (т.е. $\mathbf{A}'\mathbf{X} = 0$), то умножив ее на \mathbf{Y} можно избавиться от β :

$$\mathbf{A}'\mathbf{Y} = \mathbf{A}'\mathbf{X}\beta + \mathbf{A}'\mathbf{V} = \mathbf{0} + \mathbf{A}'\mathbf{V} = \mathbf{A}'\mathbf{V}$$

$$A'Y \sim N(\mathbf{0}, \mathbf{A}'\mathbf{V}\mathbf{A})$$

Тогда можно воспользоваться ML, чтобы найти \mathbf{V} .

В результате получатся несмешенные оценки дисперсий.

REML-оценки β будут стремится к ML-оценкам при увеличении объема выборки.

ML или REML?

Если нужны точные оценки фиксированных эффектов – ML.

Если нужны точные оценки случайных эффектов – REML.

Если нужно работать с правдоподобиями – следите, чтобы в моделях, подобранных REML была одинаковая фиксированная часть.

Для обобщенных (негауссовых) смешанных линейных моделей REML не определен – там используется ML.

Тестирование гипотез в смешанных моделях

Использование смешанных моделей для получения выводов

Тесты, которые традиционно применяются для GLM, дадут лишь **приблизительные результаты** для GLMM:

- t-(или z-) тесты Вальда для коэффициентов,
- тесты отношения правдоподобий (Likelihood ratio tests, LRT).

Поэтому для отбора моделей применяют подход, не связанный с тестами:

- информационные критерии (AIC, BIC и т.п.).

Наиболее точные результаты тестов можно получить, используя **“золотой стандарт”**:

- параметрический бутстреп.

t-(или -z) тесты Вальда

$$H_0 : \beta_k = 0, \quad H_A : \beta_k \neq 0$$

$\frac{b_k}{SE_{b_k}} \sim N(0, 1)$ или $\frac{b_k}{SE_{b_k}} \sim t_{(df=n-p)}$, если нужно оценивать σ

b_k — оценка коэффициента, n — объем выборки, p — число параметров модели.

t-(или -z) тесты Вальда дают лишь приблизительный результат, поэтому в пакете `lme4` даже не приводят уровни значимости в `summary()`. Не рекомендуется ими пользоваться.

```
coef(summary(MS1))
```

	Estimate	Std. Error	t value
(Intercept)	251.41	6.825	36.838
Days	10.47	1.546	6.771

Тесты отношения правдоподобий (LRT)

$$LRT = 2 \ln \left(\frac{L_{M_1}}{L_{M_2}} \right) = 2(\ln L_{M_1} - \ln L_{M_2})$$

- M_1 и M_2 — вложенные модели (M_1 — более полная, M_2 — уменьшенная),
- L_{M_1}, L_{M_2} — правдоподобия моделей и $\ln L_{M_1}, \ln L_{M_2}$ — логарифмы правдоподобий.

Распределение LRT **аппроксимируют** распределением χ^2 с $df = df_{M_2} - df_{M_1}$.

- LRT консервативен для случайных эффектов, т.к. тест гипотезы вида $H_0 : \hat{\sigma}_k^2 = 0$ происходит на границе области возможных значений параметра.
- LRT либерален для фиксированных эффектов, дает заниженные уровни значимости.

LRT для случайных эффектов

Модели с одинаковой фиксированной частью, подобранные REML, вложенные. Уровни значимости будут завышены.

```
MS1 <- lmer(Reaction ~ Days + (1 + Days | Subject), data = sl, REML = TRUE)
MS0 <- lmer(Reaction ~ Days + (1 | Subject), data = sl, REML = TRUE)
anova(MS1, MS0, refit = FALSE)
```

```
Data: sl
Models:
MS0: Reaction ~ Days + (1 | Subject)
MS1: Reaction ~ Days + (1 + Days | Subject)
      npar   AIC   BIC logLik deviance Chisq Df Pr(>Chisq)
MS0     4 1794 1807    -893      1786
MS1     6 1756 1775    -872      1744  42.8  2      5e-10 ***
---
Signif. codes:  0 '****' 0.001 '***' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Время реакции у разных людей по-разному зависит от продолжительности бессонницы.
Обычно тесты не делают, набор случайных эффектов определяется устройством данных.

LRT для фиксированных эффектов

Модели с одинаковой случайной частью, подобранные ML, вложенные. Уровни значимости будут занижены.

```
MS1.ml <- lmer(Reaction ~ Days + (1 + Days | Subject), data = sl, REML = FALSE)
MS0.ml <- lmer(Reaction ~ 1 + (1 + Days | Subject), data = sl, REML = FALSE)
anova(MS1.ml, MS0.ml)
```

```
Data: sl
Models:
MS0.ml: Reaction ~ 1 + (1 + Days | Subject)
MS1.ml: Reaction ~ Days + (1 + Days | Subject)
      npar  AIC  BIC logLik deviance Chisq Df Pr(>Chisq)
MS0.ml     5 1785 1801    -888      1775
MS1.ml     6 1764 1783    -876      1752  23.5  1  0.0000012 ***
---
Signif. codes:  0 '****' 0.001 '***' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Время реакции зависит от продолжительности бессонницы.

Сравнение моделей по AIC

Модели с одинаковой случайной частью, подобранные ML, вложенные или невложенные.

```
AIC(MS1.ml, MS0.ml)
```

	df	AIC
MS1.ml	6	1764
MS0.ml	5	1785

Время реакции зависит от продолжительности бессонницы (AIC).

Бутстреп для тестирования значимости и для предсказаний

Бутстреп

Способ тестирования гипотез, при котором из данных многократно получают выборки (с повторениями).

Если по таким многократным выборкам построить распределение статистики, то оно будет стремиться к истинному распределению этой статистики.

Сгенерированное бутстрепом распределение статистики можно использовать для тестирования гипотез.

Бутстреп особенно удобен, когда невозможно получить распределение статистики аналитическим путем.

Параметрический бутстреп для LRT

Чтобы при помощи бутстрапа получить оценку уровня значимости для LRT при сравнении двух моделей M_{full} и $M_{reduced}$, нужно

1. Многоократно повторить:

- сгенерировать новые данные из уменьшенной модели,
- по сгенерированным данным подобрать полную и уменьшенную модели и рассчитать LRT.

2. Построить распределение LRT по всем итерациям бутстрапа.

Уровень значимости – это доля итераций, в которых получено LRT больше, чем данное.

Параметрический бутстреп для LRT фиксированных эффектов

В строке PBtest – значение LRT и его уровень значимости, полученный бутстрепом.

```
library(pbkrtest)
pmod <- PBmodcomp(MS1.ml, MS0.ml, nsim = 100) # 1000 и больше для реальных данных
summary(pmod)
```

```
Bootstrap test; time: 1.82 sec; samples: 100; extremes: 0;
large : Reaction ~ Days + (1 + Days | Subject)
Reaction ~ 1 + (1 + Days | Subject)
      stat   df   ddf   p.value
LRT     23.5  1.0    0.0000012 ***
PBtest  23.5          0.00990 **
Gamma   23.5          0.0000048 ***
Bartlett 20.1  1.0    0.0000072 ***
F       23.5  1.0 13.9   0.00026 ***
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Бутстреп-оценка доверительной зоны регрессии

Чтобы при помощи бутстрапа оценить положение зоны, где с 95% вероятностью будут лежать предсказанные значения, нужно:

1. Многократно повторить:

- сгенерировать новые данные из модели
- по сгенерированным данным подобрать модель и получить ее предсказания

2. Построить распределение предсказанных значений по всем итерациям бутстрапа.

95% доверительная зона регрессии — это область, куда попали предсказанные значения в 95% итераций (т.е. между 0.025 и 0.975 персентилями).

Бутстреп-оценка доверительной зоны регрессии

```
NewData <- sl %>% group_by(Subject) %>%
  do(data.frame(Days = seq(min(.Days), max(.Days), length = 10)))
NewData$fit <- predict(MS1, NewData, type = 'response', re.form = NA)

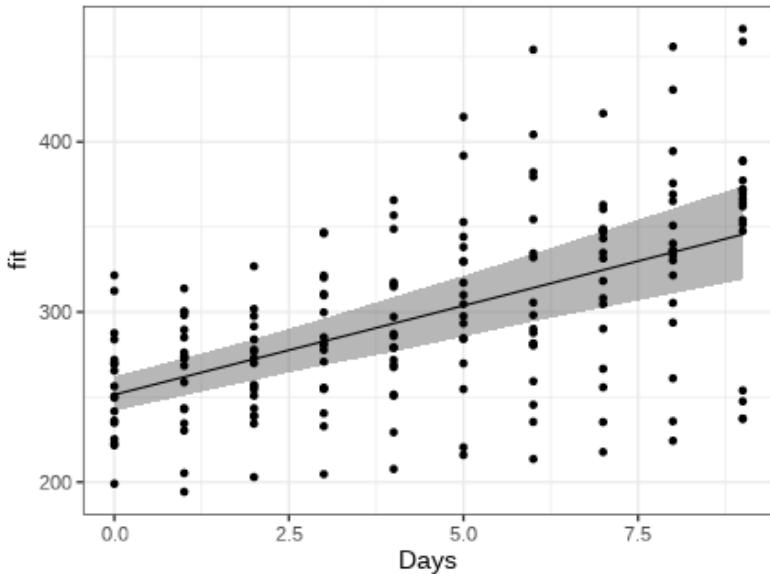
# Многократно симулируем данные из модели и получаем для них предсказанные значения
bMS1 <- bootMer(x = MS1,
                  FUN = function(x) predict(x, new_data = NewData, re.form = NA),
                  nsim = 100)

# Рассчитываем квантили предсказанных значений для всех итераций бутстрепа
b_se <- apply(X = bMS1$t,
               MARGIN = 2,
               FUN = function(x) quantile(x, probs = c(0.025, 0.975), na.rm = TRUE))

# Доверительная зона для предсказанных значений
NewData$lwr <- b_se[1, ]
NewData$upr <- b_se[2, ]
```

График предсказаний фиксированной части модели

```
ggplot(data = NewData, aes(x = Days, y = fit)) +  
  geom_ribbon(alpha = 0.35, aes(ymin = lwr, ymax = upr)) +  
  geom_line() + geom_point(data = sl, aes(x = Days, y = Reaction))
```



Take-home messages

Свойства	Фиксированные факторы	Случайные факторы
Уровни фактора	фиксированные, заранее определенные и потенциально воспроизводимые уровни	случайная выборка из всех возможных уровней
Используются для тестирования гипотез	о средних значениях отклика между уровнями фактора $H_0 : \mu_1 = \mu_2 = \dots = \mu_i = \mu$	о дисперсии отклика между уровнями фактора $H_0 : \sigma_{rand.\ fact.}^2 = 0$
Выводы можно экстраполировать	только на уровни из анализа	на все возможные уровни
Число уровней фактора	Осторожно! Если уровней фактора слишком много, то нужно подбирать слишком много коэффициентов — должно быть много данных	Важно! Для точной оценки σ нужно много уровней фактора — не менее 5

Take-home messages

- Смешанные модели могут включать случайные и фиксированные факторы.
 - Градации фиксированных факторов заранее определены, а выводы можно экстраполировать только на такие уровни, которые были задействованы в анализе. Тестируется гипотеза о равенстве средних в группах.
 - Градации случайных факторов — выборка из возможных уровней, а выводы можно экстраполировать на другие уровни. Тестируется гипотеза о дисперсии между группами.
- Есть два способа подбора коэффициентов в смешанных моделях: ML и REML. Для разных этапов анализа важно, каким именно способом подобрана модель.
- Коэффициент внутриклассовой корреляции оценивает, насколько коррелируют друг с другом наблюдения из одной и той же группы случайного фактора.
- Случайные факторы могут описывать вариацию как интерсептов, так и коэффициентов угла наклона.
- Модели со смешанными эффектами позволяют получить предсказания как общем уровне, так и на уровне отдельных субъектов.

Дополнительные ресурсы

- Crawley, M.J. (2007). The R Book (Wiley).
- Faraway, J. J. (2017). Extending the linear model with R: generalized linear, mixed effects and nonparametric regression models (Vol. 124). CRC press.
- Zuur, A. F., Hilbe, J., & Ieno, E. N. (2013). A Beginner's Guide to GLM and GLMM with R: A Frequentist and Bayesian Perspective for Ecologists. Highland Statistics.
- Zuur, A.F., Ieno, E.N., Walker, N., Saveliev, A.A., and Smith, G.M. (2009). Mixed Effects Models and Extensions in Ecology With R (Springer)
- Pinheiro, J., Bates, D. (2000). Mixed-Effects Models in S and S-PLUS. Springer