

# Подборка экзаменов по эконометрике. Факультет экономики, НИУ-ВШЭ

Коллектив кафедры  
математической экономики и эконометрики,  
фольклор,  
очень умные студенты

21 мая 2018 г.

## Содержание

<b>1</b>	<b>Описание</b>	<b>3</b>
<b>2</b>	<b>Вечное</b>	<b>3</b>
2.1	Гимн-памятка для эконометриста . . . . .	3
2.2	Прошение о повышении оценки . . . . .	4
2.3	Цитаты . . . . .	5
<b>3</b>	<b>Немного теории</b>	<b>5</b>
3.1	Конвенция об обозначениях . . . . .	5
3.2	Свойства ковариационных матриц . . . . .	6
3.3	Картинка . . . . .	6
3.4	ТГМ. Детерминированные регрессоры . . . . .	6
3.5	ТГМ. Стохастические регрессоры . . . . .	6
3.6	Ликбез по линейной алгебре . . . . .	8
3.7	Ожидание от RSS . . . . .	8
3.8	Устоявшиеся слова . . . . .	9
3.9	Ridge/Lasso regression . . . . .	9
3.10	Заповеди научного программирования . . . . .	10
<b>4</b>	<b>2012-2013</b>	<b>11</b>
4.1	Праздник 1. Пролетарий на коня! . . . . .	11
4.2	Праздник 2. Базовая задача . . . . .	12
4.3	Праздник 2. Базовая задача, ответы . . . . .	13
4.4	Праздник 3. Дню рождения буквы «ё» посвящается... . . . .	14
4.5	Праздник 4, ML . . . . .	15
4.6	Праздник 5, 01.04.2013, Гетероскедастичность . . . . .	16
4.7	Домашнее задание 3. Знакомство с RLMS . . . . .	17
4.8	Домашнее задание 1 ( $n + 1$ ) по эконометрике-1. . . . .	18
4.9	Домашнее задание. Титаник. . . . .	19
<b>5</b>	<b>2013-2014</b>	<b>21</b>
5.1	Праздник 1. Вперед в рукопашную! . . . . .	21
5.2	Праздник 2. Мегаматрица . . . . .	21

5.3	Праздник 3. Базовая задача . . . . .	22
5.4	Праздник 3. Ответы . . . . .	22
5.5	Праздник 4 . . . . .	23
5.6	Праздник 5. Максимальное правдоподобие . . . . .	24
5.7	Переписывание кр 5. Максимальное правдоподобие . . . . .	25
5.8	Праздник 6. Гетероскедастичность . . . . .	26
5.9	Большой Устный Зачёт . . . . .	26
5.10	Экзамен. . . . .	28
5.11	Пересдача экзамена . . . . .	29
5.12	Домашняя работа 1. RLMS и гетероскедастичность . . . . .	31
5.13	Домашняя работа 2. Титаник . . . . .	32
<b>6</b>	<b>2014-2015</b>	<b>32</b>
6.1	Праздник номер 1 . . . . .	32
6.2	Праздник номер 2 . . . . .	33
6.3	Праздник номер 2, ответы . . . . .	35
6.4	Праздник номер 3 . . . . .	35
6.5	Миникр . . . . .	36
6.6	Зачет. Базовый поток . . . . .	37
6.7	Зачет, 26.12.2014. Ликвидация безграмотности . . . . .	38
6.8	Домашняя работа 1. RLMS и гетероскедастичность . . . . .	39
6.9	Экзамен. Демо-1 . . . . .	41
6.10	Экзамен. Демо-2 . . . . .	42
6.11	Экзамен. 15.06.15 . . . . .	45
<b>7</b>	<b>2015-2016</b>	<b>46</b>
7.1	Праздник номер 1. Вспомнить всё! 15.09.2015 . . . . .	46
7.2	Праздник номер 1. Вспомнить всё! 15.09.2015, решение . . . . .	47
7.3	Праздник номер 2, 10 ноября 2015 . . . . .	49
7.4	Праздник номер 2, 10 ноября 2015, некоторые ответы . . . . .	51
7.5	Блокбастер, 28-12-2015 . . . . .	52
7.6	Блокбастер, 28-12-2015, ответы . . . . .	53
7.7	Максимально правдоподобно, 25-02-2016 . . . . .	53
7.8	Решение задач КР по Эконометрике, 2015-2016 . . . . .	54
7.9	Большой Устный Зачёт 2016 . . . . .	58
7.10	Экзамен 20.06.2016. Вариант 1 . . . . .	60
7.11	Экзамен 20.06.2016. Вариант 2 . . . . .	62
<b>8</b>	<b>2016-2017</b>	<b>64</b>
8.1	ИП. Подготовка к празднику . . . . .	64
8.2	ИП. Праздник «Вспомнить всё!» 12.09.2016 . . . . .	65
8.3	ИП. Праздник «Вспомнить всё!» 12.09.2016, ответы . . . . .	66
8.4	Кр 1, демо . . . . .	66
8.5	Кр 1, демо, решения . . . . .	67
8.6	Кр 1, 24.10.2016 . . . . .	69
8.7	ИП. Подготовка к битве . . . . .	69
8.8	ИП. Битва под Малоярославцем, 24.10.2016/24.10.1812 . . . . .	70
8.9	ИП. Битва — решения . . . . .	72
8.10	Кр 2, экзамен за I семестр, демо . . . . .	75
8.11	Кр 2, экзамен за I семестр, демо, решения . . . . .	78
8.12	Кр 2, экзамен за I семестр, 24.12.2016 . . . . .	80

8.13	Кр 2, экзамен за I семестр, 24.12.2016, решения	84
8.14	Кр 3, задачи для подготовки	86
8.15	Кр 3, 20.03.2017	86
8.16	Кр 3, 20.03.2017, решения	87
8.17	ИП. Комоедица, 24.03.2016	87
8.18	Кр 4, финальный экзамен, демо	90
8.19	Кр 4, финальный экзамен	91
8.20	ИП. БУЗА	93
<b>9</b>	<b>2017-2018</b>	<b>94</b>
9.1	ИП, вспомнить всё!	94
9.2	ИП, вспомнить всё!, ответы	95
9.3	Контрольная 1, 26.10.2017	96
9.4	Контрольная 1, 26.10.2017, решения	98
9.5	Контрольная 2, 26.10.2017	100
9.6	Контрольная 2, 26.10.2017, решения	104
9.7	Кр 3, 2018-03-28, бп часть	106
9.7.1	Тест	106
9.7.2	Задачи	106
9.8	Кр 3, 2018-03-28, бп часть, решения	106
9.9	КР 3, 2018-03-28, ип часть	109

## 1. Описание

Свежая версия на [https://github.com/bdemeshev/em301/tree/master/metrics\\_exams](https://github.com/bdemeshev/em301/tree/master/metrics_exams). Часть задач и многие решения придуманы студентами. Огромное спасибо всем тем, кто в мучениях, своей кровью дописывал этот документ :)

## 2. Вечное

### 2.1. Гимн-памятка для эконометриста

Эмилю Борисовичу Ершову посвящается

Ничего на свете лучше нету,  
Чем оценивать параметр «бета»!  
Лучшее оружие демократа —  
Метод наименьшего квадрата!

Если вдруг подавит вас депрессия,  
Виновата, значит, здесь дисперсия.  
Убери гетероскедастичность,  
Это придаёт оптимистичность.

Если в данных автокорреляция,  
Всё, что посчитал ты, — профанация.  
Применяй, не глядя исподлобья,  
Максимальное правдоподобие.

Если ощутил ты свою бренность,  
Не иначе это эндогенность.

Соглашайся выдать алименты  
Тем, кто знает, где взять инструменты.

Где б ты ни был, в саклях и ярангах  
Применяй везде условия ранга.  
Помни также: лучшая зарядка —  
Выполнить условие порядка.

Мы своё призвание не забудем!  
BLUE-оценки мы предъявим людям!  
Нам законов априорных своды  
Не понизят степеней свободы!

## 2.2. Прошение о повышение оценки

От .....

Группа .....

Я считаю, что моя итоговая оценка по курсу ..... должна быть исправлена с .... на ... по следующим причинам (обведите нужные).

1. Это единственная плохая оценка в моей зачетке
2. Тот, кто полностью списал мою работу, получил более высокую оценку
3. Тот, у кого я полностью списал работу, получил более высокую оценку
4. Из-за низкого рейтинга меня могут не взять в
  - а) РЭШ
  - б) СМЕРШ
  - в) МГУ
  - г) На Луну
  - д) .....
5. Мне нужно получить 10, чтобы компенсировать 4 по .....
6. Меня лишат стипендии
7. Я не успел договориться с тетечками из копировального отдела и раздобыть варианты контрольной, потому что .....
8. Я не посещал лекции, а тот, чьими конспектами я пользовался, не записал материал, необходимый для сдачи контрольных и домашних
9. Я отлично понимаю теорию, просто не умею решать задачи
10. Я умею решать все задачи, а на контрольной требовалось знание теории
11. У лектора/семинариста были предрассудки против негров/евреев/лесбиянок/.....
12. Все вопросы на экзамене допускали двойную трактовку. Я считаю, что не должен нести наказание за то, что мое мнение — особенное

13. Если я получу плохую оценку, отец отберет у меня ключи от машины
14. Я не мог/могла заниматься из-за необходимости разгружать вагоны по ночам
15. Нам сказали использовать творческий подход, но не объяснили, что это означает
16. Я использовал в домашках творческий подход, но мне было сказано, что я несу всякую чушь
17. Все остальные преподаватели согласны повысить мою оценку
18. Семинары и лекции начинались:
  - а) слишком рано, я еще спал
  - б) слишком поздно, я уже спал
  - в) в обеденное время, я был голодный
19. Причина по которой я получил низкую оценку проста — я очень честный. Не хочу ничего говорить о моих одноклассниках
20. У меня нет особой причины, я просто хочу оценку повыше

Дата .....

Подпись .....

## 2.3. Цитаты

«В выборке из ста муравьёв и одного кита средняя масса муравья может превышать тонну.» (?)

«Methodology, like sex, is better demonstrated than discussed, though often better anticipated than experienced.»  
(Leamer, 1983)

— Иван, ты знаешь, у нас в Чили есть нестандартные обозначения. Мы используем десятичную запятую вместо точки, умножение пишем точкой, а не крестиком, деление — двумя точками, а при измерении температуры пользуемся градусами Цельсия. Я уверен, что ты сможешь поправить все по-своему, как привыкли дети в России.

— Разумеется, Раймундо, сделаем. Это не составит труда. (со стены Ивана Высоцкого в контакте)

«Если на лекции всё понятно, это хорошо, если на докладе всё понятно, то уважать не будут!»  
(Шведов)

## 3. Немного теории

### 3.1. Конвенция об обозначениях

- $y$  — вектор-столбец зависимых переменных размера  $(n \times 1)$ , наблюдаемый случайный
- $\beta$  — вектор-столбец неизвестных коэффициентов размера  $(k \times 1)$ , ненаблюдаемый, случайный
- $\hat{y}$  — вектор столбец прогнозов для  $y$ , полученных по некоторой модели, размера  $(n \times 1)$ , наблюдаемый, случайный

- $\hat{\beta}$  — вектор-столбец оценок  $\beta$  размера  $(k \times 1)$ , наблюдаемый, случайный
- $X$  — матрица всех объясняющих переменных, размера  $(n \times k)$ . Известная, стохастическая или детерминированная в зависимости от парадигмы.
- $\varepsilon$  — вектор-столбец случайных ошибок размера  $(n \times 1)$ , ненаблюдаемый случайный
- $\hat{\varepsilon}$  — вектор-столбец остатков модели размера  $(n \times 1)$ , наблюдаемый случайный
- $c$  — вектор из единиц

В некоторых учебниках используется обозначение  $Y$  для исходного вектора зависимых переменных, а  $y$  — для центрированного, т.е.  $y = Y - \bar{Y}$ . В этом документе  $y$  обозначает исходный вектор  $y$ .

### 3.2. Свойства ковариационных матриц

Здесь  $y$  — вектор-столбец  $n \times 1$ ,  $z$  — вектор-столбец  $k \times 1$ ,  $A$  — матрица констант подходящего размера,  $b$  — вектор констант подходящего размера.

1.  $\mathbb{E}(Ay + b) = A\mathbb{E}(y) + b$ ,  $\mathbb{E}(yA + b) = \mathbb{E}(y)A + b$
2.  $\text{Cov}(y, z) = \mathbb{E}(yz') - \mathbb{E}(y)\mathbb{E}(z')$
3.  $\text{Var}(y) = \mathbb{E}(yy') - \mathbb{E}(y)\mathbb{E}(y')$
4.  $\text{Cov}(Ay + b, z) = A \text{Cov}(y, z)$ ,  $\text{Cov}(y, Az + b) = \text{Cov}(y, z)A'$
5.  $\text{Var}(Ay + b) = A \text{Var}(y)A'$
6.  $\text{Cov}(y, z) = \text{Cov}(z, y)'$

### 3.3. Картинка

Утверждение.  $\text{sCorr}^2(y, \hat{y}) = R^2$

Доказательство. По определению,  $\text{sCorr}(y, \hat{y}) = \frac{(y - \bar{y})'(\hat{y} - \bar{\hat{y}})}{|y - \bar{y}||\hat{y} - \bar{\hat{y}}|}$ . Поскольку в регрессии присутствует свободный член,  $\bar{\hat{y}} = \bar{y}$ . Значит,

$$\text{sCorr}(y, \hat{y}) = \frac{(y - \bar{y})(\hat{y} - \bar{y})}{|y - \bar{y}||\hat{y} - \bar{y}|} = \cos(y - \bar{y}, \hat{y} - \bar{y}) \quad (1)$$

По определению,  $R^2 = \frac{|\hat{y} - \bar{y}|^2}{|y - \bar{y}|^2} = \cos^2(y - \bar{y}, \hat{y} - \bar{y})$

### 3.4. ТГМ. Детерминированные регрессоры

### 3.5. ТГМ. Стохастические регрессоры

Если:

1. Истинная зависимость имеет вид  $y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i$   
В матричном виде:  $y = X\beta + \varepsilon$
2. С помощью МНК оценивается регрессия  $y$  на константу,  $x_{.2}, x_{.3}, \dots, x_{.k}$   
В матричном виде:  $\hat{\beta} = (X'X)^{-1}X'y$

3. Наблюдений больше, чем оцениваемых коэффициентов  $\beta$ :  $n > k$
4. Строгая экзогенность:  $\mathbb{E}(\varepsilon_i | \text{все } x_{ij}) = 0$   
В матричном виде:  $\mathbb{E}(\varepsilon_i | X) = 0$
5. Условная гомоскедастичность:  $E(\varepsilon_i^2 | \text{все } x_{ij}) = \sigma^2$   
В матричном виде:  $\mathbb{E}(\varepsilon_i^2 | X) = \sigma^2$
6.  $\text{Cov}(\varepsilon_i, \varepsilon_j | X) = 0$  при  $i \neq j$
7. вектора  $(x_i, y_i)$  — независимы и одинаково распределены
8. с вероятностью 1 среди регрессоров нет линейно зависимых  $\text{rank}(X) = k$   $\det(X'X) \neq 0$   
 $(X'X)^{-1}$  существует

То:

1. (тГМ) МНК оценки  $\hat{\beta}$  линейны по  $y$ :  $\hat{\beta}_j = c_1 y_1 + \dots + c_n y_n$
2. (тГМ) МНК оценки несмещенные. А именно,  $\mathbb{E}(\hat{\beta} | X) = \beta$ , и в частности  $\mathbb{E}(\hat{\beta}) = \beta$
3. (тГМ) МНК оценки эффективны среди линейных несмещённых оценок. Для любой альтернативной оценки  $\hat{\beta}^{alt}$  удовлетворяющей свойствам 1 и 2:  $\text{Var}(\hat{\beta}_j^{alt} | X) \geq \text{Var}(\hat{\beta}_j | X)$   $\text{Var}(\hat{\beta}_j^{alt}) \geq \text{Var}(\hat{\beta}_j)$
4.  $\text{Var}(\hat{\beta} | X) = \sigma^2 (X'X)^{-1}$
5.  $\text{Cov}(\hat{\beta}, \hat{\varepsilon} | X) = 0$
6.  $\mathbb{E}(\hat{s}^2 | X) = \sigma^2$ , и  $\mathbb{E}(\hat{s}^2) = \sigma^2$  ?остается ли при условной ГК?

Если дополнительно к предпосылкам теоремы Гаусса-Маркова известно, что  $\varepsilon | X \sim \mathcal{N}$ , то:

1. МНК оценки эффективны среди всех несмещённых оценок.
2.  $t | X \sim t_{n-k}$ ,  $t \sim t_{n-k}$
3.  $RSS/\sigma^2 | X \sim \chi_{n-k}^2$ ,  $RSS/\sigma^2 \sim \chi_{n-k}^2$
4.  $F$  тест  $F | X \sim F_{r, n-k_{UR}}$  при выполнении  $r$  ограничений
5.  $R^2 \sim \mathcal{B}(\dots, \dots)$  при  $\beta_2 = \dots = \beta_k = 0$

Если дополнительно к предпосылкам теоремы Гаусса-Маркова известно, что  $n \rightarrow \infty$ , то:

1.  $\hat{\beta} \rightarrow \beta$  по вероятности (состоятельность)

*Доказательство.* Разложим  $\hat{\beta}$  в виде  $\hat{\beta} = (X'X)^{-1}X'y = (X'X)^{-1}X'(X\beta + \varepsilon) = \beta + (X'X)^{-1}X'\varepsilon$

Заметим, что  $(X'X)^{-1}X'\varepsilon = \left(\frac{1}{n}X'X\right)^{-1} \frac{1}{n}X'\varepsilon$ .

$\text{plim} \left(\frac{1}{n}X'X\right) = \text{Var}(X_i.)$

$\text{plim} \frac{1}{n}X'\varepsilon = 0$

□

2.  $t \rightarrow \mathcal{N}(0, 1)$

3.  $rF \rightarrow \chi_r^2$ ,  $r$  — число ограничений при выполнении  $r$  ограничений
4.  $nR^2 \rightarrow \chi_{k-1}^2$  при  $\beta_2 = \dots = \beta_k = 0$
5.  $\frac{RSS}{n-k} \rightarrow \sigma^2$

#### Сравнение двух парадигм

	детерминированные $X$	случайные $X$
$\mathbb{E}(y_i)$	разные, $X_i \beta$	одинаковые
$sVar(y)$ — несмещенная оценка для $Var(y_i)$	Нет	Да

### 3.6. Ликбез по линейной алгебре

**Определение.** Неформально. Если матрица  $A$  квадратная, то её определителем называется площадь/объём параллелограмма/параллелепипеда образованного векторами-столбцами матрицы. Знак определителя задаётся порядком следования векторов.

Свойства определителя:

1.  $\det(AB) = \det(A) \det(B) = \det(BA)$ , если  $A$  и  $B$  квадратные
2.  $\det(A) = \prod \lambda_i$ , где  $\lambda_i$  — собственное число матрицы  $A$ , возможно комплексное.

**Определение.** Ненулевой вектор  $x$  называется собственным вектором матрицы  $A$ , если при умножении на матрицу  $A$  он остаётся на той же прямой, т.е.  $Ax = \lambda x$ .

**Определение.** Число  $\lambda$  называется собственным числом матрицы  $A$ , если существует вектор  $x$ , который при умножении на матрицу  $A$  изменяется в  $\lambda$  раз, т.е.  $Ax = \lambda x$ .

**Определение.** Если матрица  $A$  квадратная, то её следом называется сумма диагональных элементов,  $\text{trace}(A) = \sum a_{ii}$ .

Свойства следа:

1.  $\text{trace}(A + B) = \text{trace}(A) + \text{trace}(B)$
2.  $\text{trace}(AB) = \text{trace}(BA)$ , если  $AB$  и  $BA$  существуют. При этом  $A$  и  $B$  могут не быть квадратными матрицами.
3.  $\text{trace}(A) = \sum \lambda_i$ , где  $\lambda_i$  — собственное число матрицы  $A$ , возможно комплексное.

Смысл следа. Если умножение на матрицу  $A$  — это проецирование, то есть  $Ax$  — есть проекция вектора  $x$  на некоторое подпространство, то  $\text{trace}(A)$  — размерность этого подпространства. Действительно, если  $A$  — проектор, то  $A^2 = A$  и собственные числа матрицы  $A$  равны нулю или единице. Поэтому  $\text{trace}(A)$  равен количеству собственных чисел равных единице. И, следовательно,  $\text{trace}(A)$  равен  $\text{rank}(A)$ , то есть размерности пространства, на которое матрица  $A$  проецирует вектора. У следа матрицы есть и другие смыслы [1].

### 3.7. Ожидание от RSS

**Теорема.** След и математическое ожидание можно переставлять,  $\mathbb{E}(\text{tr}(A)) = \text{tr}(\mathbb{E}(A))$ .

**Теорема.** Математическое ожидание квадратичной формы

$$\mathbb{E}(x'Ax) = \text{tr}(A \text{Var}(x)) + \mathbb{E}(x')A\mathbb{E}(x) \quad (2)$$



*Доказательство.* Мы будем пользоваться простым приёмом. Если  $u$  — это скаляр, вектор размера 1 на 1, то  $\text{tr}(u) = u$ .

Поехали,

$$\mathbb{E}(x'Ax) = \mathbb{E}(\text{tr}(x'Ax)) = \mathbb{E}(\text{tr}(Axx')) = \text{tr}(\mathbb{E}(Axx')) = \text{tr}(A\mathbb{E}(xx')) \quad (3)$$

По определению дисперсии,  $\text{Var}(x) = \mathbb{E}(xx') - \mathbb{E}(x)\mathbb{E}(x')$ . Поэтому:

$$\text{tr}(A\mathbb{E}(xx')) = \text{tr}(A(\text{Var}(x) + \mathbb{E}(x)\mathbb{E}(x'))) = \text{tr}(A\text{Var}(x)) + \text{tr}(A\mathbb{E}(x)\mathbb{E}(x')) \quad (4)$$

И готовимся снова использовать приём  $\text{tr}(u) = u$ :

$$\text{tr}(A\text{Var}(x)) + \text{tr}(A\mathbb{E}(x)\mathbb{E}(x')) = \text{tr}(A\text{Var}(x)) + \text{tr}(\mathbb{E}(x')A\mathbb{E}(x)) = \text{tr}(A\text{Var}(x)) + \mathbb{E}(x')A\mathbb{E}(x) \quad (5)$$

□

### 3.8. Устоявшиеся слова

Выражение «гипотеза о значимости отдельного коэффициента» на самом деле означает «гипотеза о незначимости отдельного коэффициента», т.к. де-факто проверяется гипотеза  $H_0: \beta_j = 0$ .

Выражение «гипотеза о значимости регрессии в целом» или «гипотеза об адекватности регрессии» на самом деле означает «гипотеза о незначимости регрессии в целом», т.к. проверяется  $H_0: \beta_2 = \dots = \beta_k = 0$ .

В некоторых источниках гипотезу об адекватности регрессии ошибочно обозначают  $H_0: R^2 = 0$ . Эту ошибку не нужно повторять.

Гипотезы имеет смысл проверять о ненаблюдаемых величинах, а величина  $R^2$  является наблюдаемой. И если уж на то пошло, то проверить гипотезу о том, что  $R^2 = 0$  тривиально. Для этого не нужно знать ничего из теории вероятностей, достаточно просто сравнить посчитанное значение  $R^2$  с нулём.

Более того, даже корректировка  $H_0: \mathbb{E}(R^2) = 0$  неверна. В модели, где в регрессоры включена только константа, величина  $R^2$  тождественно равна нулю, поэтому  $\mathbb{E}(R^2) = 0$  и проверять такую гипотезу бессмысленно. В модели, где в регрессоры включено что-то помимо константы,  $R^2$  является неотрицательной случайной величиной с  $\mathbb{P}(R^2 > 0) > 0$ . Поэтому а-приори  $\mathbb{E}(R^2) > 0$  и проверка гипотезы  $H_0: \mathbb{E}(R^2) = 0$  снова бессмысленна.

Кстати, обозначение  $H_0$  по-английски читается как «H naught», а не «H zero» или «H null». Также корректно говорить «the null hypothesis».

### 3.9. Ridge/Lasso regression

LASSO — Least Absolute Shrinkage and Selection Operator. Метод построения регрессии, предложенный Robert Tibshirani в 1995 году.

Вспомним обычный МНК:

$$\min_{\beta} (y - X\beta)'(y - X\beta) \quad (6)$$

LASSO вместо исходной задачи решает задачу условного экстремума:

$$\min_{\beta} (y - X\beta)'(y - X\beta) \quad (7)$$

при ограничении  $\sum_{j=2}^k |\beta_j| \leq c$ .

Естественно, при больших значениях  $c$  результат LASSO совпадает с МНК. Что происходит при малых  $c$ ?

Для наглядности рассмотрим задачу с двумя коэффициентами  $\beta: \beta_1$  и  $\beta_2$ . Линии уровня целевой функции — эллипсы. Допустимое множество имеет форму ромба с центром в начала координат.

на картинке три  $c$ : очень большое — дающее мнк решение, меньше — ненулевые  $\beta$ , маленькое — одна из  $\beta$  равна 0

То есть при малых  $c$  LASSO обратит ровно в ноль некоторые коэффициенты  $\beta$ .

Применим метод множителей Лагранжа для случая, когда ограничение  $\sum_{j=1}^k |\beta_j| \leq c$  активно, то есть выполнено как равенство.

$$L(\beta, \lambda) = (y - X\beta)'(y - X\beta) + \lambda \left( \sum_{j=1}^k |\beta_j| - c \right) \quad (8)$$

Необходимым условием первого порядка является  $\partial L / \partial \beta = 0$ . Это условие первого порядка не изменится, если мы зачеркнём  $c$  в выражении. Таким образом мы получили альтернативную формулировку метода LASSO:

$$\min_{\beta} (y - X\beta)'(y - X\beta) + \lambda \sum_{j=1}^k |\beta_j| \quad (9)$$

LASSO пытается минимизировать взвешенную сумму  $RSS = (y - X\beta)'(y - X\beta)$  и «размера» коэффициентов  $\sum_{j=1}^k |\beta_j|$ .

Мы не будем вдаваться в численные алгоритмы, которые используются при решении этой задачи. Ridge regression отличается от LASSO ограничением  $\sum \beta_j^2 \leq c$ . Также как и LASSO Ridge regression допускает альтернативную формулировку:

$$\min_{\beta} (y - X\beta)'(y - X\beta) + \lambda \sum_{j=1}^k \beta_j^2 \quad (10)$$

Также как и LASSO Ridge regression тоже приближает значения коэффициентов  $\beta_j$  к нулю. Принципиальное отличие LASSO и RR. В LASSO крайнее решение с несколькими коэффициентами равными нулю является типичной ситуацией. В RR коэффициент  $\beta_j$  может оказаться точно равным нулю только по чистой случайности.

LASSO допускает байесовскую интерпретацию...

Предположим, что априорное распределение параметров следующее:

...

Тогда мода апостериорного распределения будут приходиться в точности на оценки LASSO.

### 3.10. Заповеди научного программирования

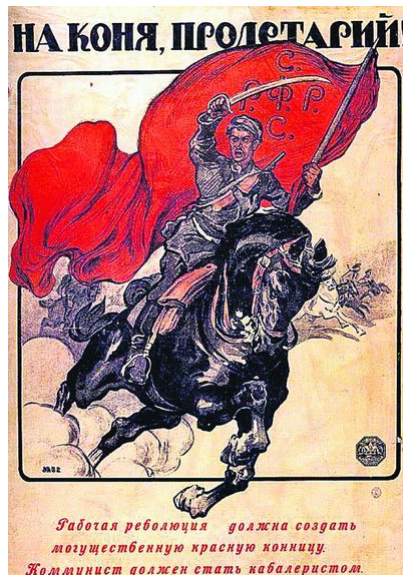
1. Не помяни русские буквы или пробелы в имени файла твоего
2. Не помяни запятую в качестве десятичного разделителя числа твоего
3. Не используй никаких форматов хранения данных кроме csv пока не превышен будет объём диска твоего
4. Почитай кодировку UTF-8 для русскоязычных текстов, чтобы продлились дни их на земле
5. Проверяй целостность данных после загрузки и преобразований
6. Комментируй свой код щедро и обильно
7. Руководствуйся стилевым гидом при оформлении кода твоего
8. Сохраняй seed в случайных экспериментах, дабы были они воспроизводимы
9. Используй систему контроля версий, дабы не быть в горе и печали

## Список литературы

[1] Смысл следа матрицы. URL: <http://mathoverflow.net/questions/13526/>.

### 4. 2012-2013

#### 4.1. Праздник 1. Пролетарий на коня!



1. Найдите длины векторов  $a = (1, 2, 3)$  и  $b = (1, 0, -1)$  и косинус угла между ними.
2. Сформулируйте теорему о трёх перпендикулярах.
3. Сформулируйте и докажите теорему Пифагора.
4. Для матрицы
$$A = \begin{pmatrix} 2 & 3 & 0 \\ 3 & 10 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$
  - а) Найдите собственные числа и собственные векторы матрицы.
  - б) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.
  - в) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.
  - г) Представьте  $A^{2012}$  в виде произведения трёх матриц.
5. Вася и Петя независимо друг от друга решают тест по теории вероятностей. В тесте всего два вопроса. На каждый вопрос два варианта ответа. Петя знает решение каждого вопроса с вероятностью 0,7. Если Петя не знает решения, то он отвечает равновероятно наугад. Вася знает решение каждого вопроса с вероятностью 0,5. Если Вася не знает решения, то он отвечает равновероятно наугад.
  - а) Какова вероятность того, что Петя правильно ответил на оба вопроса?
  - б) Какова вероятность того, что Петя правильно ответил на оба вопроса, если его ответы совпали с Васиными?
  - в) Чему равно математическое ожидание числа Петиних верных ответов?

- г) Чему равно математическое ожидание числа Петиних верных ответов, если его ответы совпали с Васиными?
6. Для случайных величин  $X$  и  $Y$  заданы следующие значения:  $\mathbb{E}(X) = 1$ ,  $\mathbb{E}(Y) = 4$ ,  $\mathbb{E}(XY) = 8$ ,  $\text{Var}(X) = \text{Var}(Y) = 9$ . Для случайных величин  $U = X + Y$  и  $V = X - Y$  вычислите:
- $\mathbb{E}(U)$ ,  $\text{Var}(U)$ ,  $\mathbb{E}(V)$ ,  $\text{Var}(V)$ ,  $\text{Cov}(U, V)$
  - Можно ли утверждать, что случайные величины  $U$  и  $V$  независимы?
7. Вася ведёт блог. Обозначим  $X_i$  — количество слов в  $i$ -ой записи. После первого года он по своим записям обнаружил, что  $\bar{X}_{200} = 95$  и выборочное стандартное отклонение равно 282 слова. На уровне значимости  $\alpha = 0.10$  проверьте гипотезу о том, что  $\mu = 100$  против альтернативной гипотезы  $\mu \neq 100$ . Найдите также точное Р-значение.

## 4.2. Праздник 2. Базовая задача

Плывут облака

Отдыхать после знойного дня,

Стремительных птиц

Улетела последняя стая.

Гляжу я на горы,

И горы глядят на меня,

И долго глядим мы,

Друг другу не надоедая.

Ли Бо, Одиноко сижу в горах Цзинтиншань

- Случайные величины  $Z_i$  независимы и нормально распределены  $\mathcal{N}(0, 1)$ . Для их суммы  $S = \sum_{i=1}^n Z_i$  найдите  $\mathbb{E}(S)$  и  $\text{Var}(S)$ .
- Социологическим опросам доверяют 70% жителей. Те, кто доверяют опросам, на все вопросы отвечают искренне; те, кто не доверяют, отвечают равновероятно наугад. Социолог Петя в анкету очередного опроса включил вопрос «Доверяете ли Вы социологическим опросам?»
  - Какова вероятность, что случайно выбранный респондент ответит «Да»?
  - Какова вероятность того, что он действительно доверяет, если известно, что он ответил «Да»?
- Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 2 & 1 \\ 2 & 2 & 0 \\ 1 & 0 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & -1 \\ -1 & 2 & 1 \\ -1 & 1 & 3 \end{pmatrix}.$$

- а) Укажите число наблюдений.
  - б) Укажите число регрессоров с учетом свободного члена.
  - в) Рассчитайте при помощи метода наименьших квадратов  $\hat{\beta}$ , оценку для вектора неизвестных коэффициентов.
  - г) Рассчитайте  $TSS = \sum (y_i - \bar{y})^2$ ,  $RSS = \sum (y_i - \hat{y}_i)^2$  и  $ESS = \sum (\hat{y}_i - \bar{y})^2$ .
  - д) Чему равен  $\hat{\varepsilon}_4$ , МНК-остаток регрессии, соответствующий 4-ому наблюдению?
  - е) Чему равен  $R^2$  в модели?
  - ж) Рассчитайте несмещенную оценку для неизвестного параметра  $\sigma^2$  регрессионной модели.
  - з) Рассчитайте  $\widehat{\text{Var}}(\hat{\beta})$ , оценку для ковариационной матрицы вектора МНК-коэффициентов  $\hat{\beta}$ .
  - и) Найдите  $\widehat{\text{Var}}(\hat{\beta}_1)$ , несмещенную оценку дисперсии МНК-коэффициента  $\hat{\beta}_1$ .
  - к) Найдите  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2)$ , несмещенную оценку ковариации МНК-коэффициентов  $\hat{\beta}_1$  и  $\hat{\beta}_2$ .
  - л) Найдите  $\widehat{\text{Var}}(\hat{\beta}_1 + \hat{\beta}_2)$
  - м) Найдите  $\widehat{\text{Corr}}(\hat{\beta}_1, \hat{\beta}_2)$ , оценку коэффициента корреляции МНК-коэффициентов  $\hat{\beta}_1$  и  $\hat{\beta}_2$ .
  - н) Найдите  $se(\hat{\beta}_1)$ , стандартную ошибку МНК-коэффициента  $\hat{\beta}_1$ .
4. В классической линейной модели предполагается, что  $\mathbb{E}(\varepsilon) = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2 I$ . Найдите  $\text{Cov}(y, \hat{\varepsilon})$ ,  $\text{Cov}(\hat{y}, \hat{\varepsilon})$ .

### 4.3. Праздник 2. Базовая задача, ответы

1.  $\mathbb{E}(S) = 0$ ,  $\text{Var}(S) = n$ .
2. а) 0.85  
б) 0.7/0.85
3. а)  $n = 5$   
б)  $k = 3$   
в)  $\hat{\beta}_1 = 1.5$ ,  $\hat{\beta}_2 = 3$ ,  $\hat{\beta}_3 = 1.5$   
г)  $TSS = 10$ ,  $RSS = 1$ ,  $ESS = 9$   
д)  $\hat{\varepsilon}_4 = -0.5$   
е)  $R^2 = 0.9$   
ж)  $\hat{\sigma}^2 = \frac{RSS}{n-k} = 0.5$
- з)  $\widehat{\text{Var}}(\hat{\beta}) = 0.25(X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & -1 \\ -1 & 2 & 1 \\ -1 & 1 & 3 \end{pmatrix}$
- и)  $\widehat{\text{Var}}(\hat{\beta}_1) = 0.25$
- к)  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) = -0.25$
- л)  $\widehat{\text{Var}}(\hat{\beta}_1 + \hat{\beta}_2) = 0.5$
- м)  $\widehat{\text{Corr}}(\hat{\beta}_1, \hat{\beta}_2) = -\frac{1}{\sqrt{2}}$
- н)  $se(\hat{\beta}_1) = 0.5$
4.  $\text{Cov}(y, \hat{\varepsilon}) = \sigma^2(I - H)$ ,  $\text{Cov}(\hat{y}, \hat{\varepsilon}) = 0$

#### 4.4. Праздник 3. Дню рождения буквы «ё» посвящается...

1. Выберите верные варианты.

- а) Побасёнка — Побасенка
- б) Вёдро — Ведро
- в) Гренадёр — Гренадер
- г) Новорождённый — Новорожденный
- д) Бытиё — Бытие
- е) Опёка — Опека
- ж) Сёрфинг — Серфинг
- з) Пафнутий Львович Чебышёв — Пафнутий Львович Чебышев
- и) Лёв Николаевич Толстой — Лев Николаевич Толстой

2. По 47 наблюдениям оценивается зависимость доли мужчин занятых в сельском хозяйстве от уровня образованности и доли католического населения по Швейцарским кантонам в 1888 году.

$$Agriculture_i = \beta_1 + \beta_2 Examination_i + \beta_3 Catholic_i + \varepsilon_i$$

---

```

1 library("lmtest")
2 library("apsrtable")
3 library("xtable")
4 h <- swiss
5 model1 <- glm(Agriculture ~ Examination + Catholic, data = h)
6 coef.t <- coeftest(model1)
7 dimnames(coef.t)[[2]] <- c("Оценка", "Ст. ошибка", "t-статистика", "Р-значение")
8 coef.t <- coef.t[, -4]
9 coef.t[1, 1] <- NA
10 coef.t[2, 2] <- NA
11 coef.t[3, 3] <- NA
12 xtable(coef.t)

```

---

	Оценка	Ст. ошибка	t-статистика
(Intercept)		8.72	9.44
Examination	-1.94		-5.08
Catholic	0.01	0.07	

- а) Заполните пропуски в таблице.
  - б) Укажите коэффициенты, значимые на 10% уровне значимости.
  - в) Постройте 95%-ый доверительный интервал для коэффициента при переменной Catholic
3. Оценивается зависимость уровня фертильности всё тех же швейцарских кантонов в 1888 году от ряда показателей. В таблице представлены результаты оценивания двух моделей.

Модель 1:  $Fertility_i = \beta_1 + \beta_2 Agriculture_i + \beta_3 Education_i + \beta_4 Examination_i + \beta_5 Catholic_i + \varepsilon_i$

Модель 2:  $Fertility_i = \gamma_1 + \gamma_2 (Education_i + Examination_i) + \gamma_3 Catholic_i + u_i$

---

```

1 m1 <- lm(Fertility ~ Agriculture + Education + Examination + Catholic, data = h)
2 m2 <- lm(Fertility ~ I(Education + Examination) + Catholic, data = h)
3 apsrtable(m1, m2)

```

---

	Model 1	Model 2
(Intercept)	91.06*	80.52*
	(6.95)	(3.31)
Agriculture	-0.22*	
	(0.07)	
Education	-0.96*	
	(0.19)	
Examination	-0.26	
	(0.27)	
Catholic	0.12*	0.07*
	(0.04)	(0.03)
I(Education + Examination)		-0.48*
		(0.08)
$N$	47	47
$R^2$	0.65	0.55
adj. $R^2$	0.62	0.53
Resid. sd	7.74	8.56

Standard errors in parentheses

\* indicates significance at  $p < 0.05$

Таблица 1:

- Посчитайте  $RSS$  для каждой модели.
- Какая модель является ограниченной (короткой), какая — неограниченной (длинной)?
- Какие ограничения нужно добавить к неограниченной модели, чтобы получить ограниченную?
- Найдите наблюдаемое значение  $F$  статистики.
- Отвергается или не отвергается гипотеза об ограничениях?

#### 4.5. Праздник 4, ML



Версия Белой Розы



- Наблюдения  $X_1, X_2, \dots, X_n$  независимы и одинаково распределены с функцией плотности  $f(x) = \frac{a(\ln(x))^{a-1}}{x}$  при  $x \in [1; e]$ . По 100 наблюдениям известно, что  $\sum_{i=1}^{100} \ln(\ln(X_i)) = -20$ 
  - Оцените параметр  $a$  методом максимального правдоподобия
  - Проверьте гипотезу о том, что  $a = 5$  против альтернативной  $a \neq 5$  с помощью теста отношения правдоподобия, теста Вальда, теста множителей Лагранжа

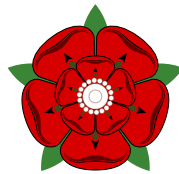
- в) Постройте 95%-ый доверительный интервал для параметра  $a$
2. [R] Фактическое распределение часовой и десятиминутной скорости ветра хорошо приближается распределением Вейбулла. Случайная величина имеет распределение Вейбулла, если её функция плотности при  $x > 0$  имеет вид

$$f(x) = \frac{1}{\lambda^k} k x^{k-1} \exp(-x^k / \lambda^k)$$

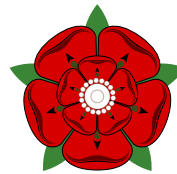
- а) Оцените параметры  $k$  и  $\lambda$  методом максимального правдоподобия
- б) Постройте 95%-ые доверительные интервалы для  $k$  и  $\lambda$

Часовые данные я не нашёл, нашёл дневные. Данные по среднедневной скорости ветра содержатся в weather\_nov\_2012\_moscow.csv в столбике wind. Данные взяты с сайта [http://www.atlas-yakutia.ru/weather/climate\\_russia-I.html](http://www.atlas-yakutia.ru/weather/climate_russia-I.html).

Hint: read.csv("filename.csv")



Версия Алой Розы



1. Купив пачку мэндэмс я насчитал в ней 1 жёлтую, 7 зелёных, 4 оранжевых, 3 коричневых, 2 синих и 1 красную мэндэмсину. С помощью теста отношения правдоподобия проверьте гипотезу, что мэндэмсины всех цветов встречаются равновероятно.
2. [R] Фактическое распределение часовой и десятиминутной скорости ветра хорошо приближается распределением Вейбулла. Случайная величина имеет распределение Вейбулла, если её функция плотности при  $x > 0$  имеет вид

$$f(x) = \frac{1}{\lambda^k} k x^{k-1} \exp(-x^k / \lambda^k)$$

- а) Найдите функцию распределения  $F(x)$
- б) Выразите медиану распределения Вейбулла,  $m$ , через параметры  $k$  и  $\lambda$
- в) Оцените параметры  $k$  и  $\lambda$  методом максимального правдоподобия
- г) Постройте 95%-ые доверительные интервалы для  $k$  и  $\lambda$
- д) Выпишите функцию плотности распределения Вейбулла через  $m$  и  $k$
- е) Проверьте гипотезу о том, что медиана равна 1 м/сек с помощью трёх тестов

Часовые данные я не нашёл, нашёл дневные. Данные по среднедневной скорости ветра содержатся в weather\_nov\_2012\_moscow.csv в столбике wind. Данные взяты с сайта [http://www.atlas-yakutia.ru/weather/climate\\_russia-I.html](http://www.atlas-yakutia.ru/weather/climate_russia-I.html).

#### 4.6. Праздник 5, 01.04.2013, Гетероскедастичность

С 1-м апреля!!!

1. Рождается старичком, умирает младенцем, сегодня празднует день рождения, но не Гоголь. Кто это? Опишите внешний вид, характер, или нарисуйте его :)
2. Для борьбы с гетероскедастичностью в модели  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$  исследователь перешёл к модели  $\tilde{y}_i = \beta_1 \frac{1}{z_i} + \beta_2 \tilde{x}_i + \tilde{\varepsilon}_i$ , где  $\tilde{x}_i = x_i / z_i$ ,  $\tilde{y}_i = y_i / z_i$ ,  $\tilde{\varepsilon}_i = \varepsilon_i / z_i$ . Какой вид гетероскедастичности предполагался?



3. Василий Аспушкин провёл два разных теста на гетероскедастичность на одном уровне значимости. Оказалось, что в одном из них  $H_0$  отвергается, а в другом — нет.
  - а) Почему это могло случиться?
  - б) Какой же вывод о гетероскедастичности следует сделать Василию? Что можно сказать об уровне значимости предложенного Вами способа сделать вывод?
4. Писатель Василий Аспушкин пишет Большой Роман. Количество страниц, которое он пишет ежедневно, зависит от количества съеденных пирожков, выпитого лимонада и числа посещений Музы.

$$Stranitsi_i = \beta_1 + \beta_2 Pirojki_i + \beta_3 Limonad_i + \beta_4 Musa_i + \varepsilon_i$$

Когда идёт дождь, Василий Аспушкин очень волнуется: он ошибочно считает, что музы плохо летают в дождь. Поэтому в дождливые дни дисперсия  $\varepsilon_i$  может быть выше.

- а) Отсортировав имеющиеся наблюдения по количеству осадков в день, Настойчивый издатель построил регрессию по 40 самым дождливым дням и получил  $RSS = \sum_i (y_i - \hat{y}_i)^2 = 360$ . В регрессии по 40 самым сухим дням  $RSS = 252$ . Всего имеется 100 наблюдений. Проверьте гипотезу о гомоскедастичности. Как называется соответствующий тест?
- б) Василий Аспушкин оценил по 100 наблюдениям исходную модель с помощью МНК. А затем построил регрессию квадратов студентизированных остатков на количество осадков и константу. Во второй регрессии  $R^2 = 0.3$ . Проверьте гипотезу о гомоскедастичности.
- в) Предположим, что дисперсия ошибок линейно зависит от количества осадков.
  - i. Как будет выглядеть функция максимального правдоподобия для оценивания коэффициентов исходной модели?
  - ii. Опишите процедуру доступного обобщенного метода наименьших квадратов (FGLS, feasible generalized least squares) применительно к данной ситуации

Hint: Функция плотности одномерного нормального распределения имеет вид

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

5. В курсе теории вероятностей изучался тест о равенстве математических ожиданий по двум нормальным выборкам при предположении о равенстве дисперсий. Предложите состоятельный способ тестировать гипотезу о равенстве математических ожиданий без предположения равенства дисперсий.

#### 4.7. Домашнее задание 3. Знакомство с RLMS

1. Прочитайте про RLMS, <http://www.hse.ru/rlms/>  
Посмотрите описание проекта. Прочитайте вестник RLMS, чтобы иметь представление о том, какие исследования можно строить на основе RLMS.
2. Скачайте любую волну RLMS по своему выбору. Скачайте описание переменных. Прочитайте описание переменных. Там их больше тысячи. Попадают довольно прикольные. Мне нравится rs9.6.5a, «У Вас есть GPRS навигатор?»

### 3. Загрузите данные в R.

Данные RLMS выложены на сайте в формате SPSS. SPSS это потихоньку погибающий статистический пакет для домохозяек. Для чтения формата .sav в таблицу данных R можно сделать так

---

```
1 library(foreign)
2 file.name <- "/home/boris/downloads/r20hall23c.sav"
3 h <- read.spss(file.name, to.data.frame = TRUE)
```

---

Первая команда, library(foreign), подгружает библиотеку R, в которой содержатся команды для чтения вражеских форматов, spss, stata, etc

Описания переменных при этом также загружаются в таблицу данных. Можно их выделить в отдельный вектор и прочитать, например, про переменную pc9.631a.

---

```
1 var.labels <- attr(h, "variable.labels")
2 var.labels["pc9.631a"]
```

---

### 4. Выберите любую количественную переменную в качестве зависимой и несколько переменных в качестве объясняющей.

Цель этой домашки скорее ознакомится с наличием мониторинга RLMS, поэтому можно не сильно заморачиваться с этим этапом. Хотя в реальности тут-то всё самое интересное и начинается. За оригинальные гипотезы будут плюшки.

### 5. Опишите выбранные переменные.

Постройте симпатичные графики. Посчитайте описательные статистики. Много ли пропущенных наблюдений? Есть ли что-нибудь интересенькое?

### 6. Постройте регрессию зависимой переменной на объясняющие.

Проверьте гипотезу о значимости каждого полученного коэффициента. Проверьте гипотезу о значимости регрессии в целом. Для нескольких коэффициентов (двух достаточно) постройте 95%-ый доверительный интервал.

### 7. Напишите свои пожелания и комментарии.

Какие домашки хочется сделать? Что не ясно в курсе эконометрики? Содержательные комментарии позволяют получить бонус. Искусная лесть оценивается :)

## 4.8. Домашнее задание 1 ( $n + 1$ ) по эконометрике-1.

### Задача 1. «CAPM»

Оценим модель CAPM по реальным данным:

1. Коротко сформулируйте теоретические положения модели CAPM. За корректное отделение выводов от предпосылок — дополнительный бонус.
2. Соберите реальные данные по трём показателям:  $R_i$  — доходность некоей акции за  $i$ -ый период,  $R_{m,i}$  — рыночная доходность за  $i$ -ый период,  $R_{f,i}$  — безрисковая доходность за  $i$ -ый период. Статья [quantile.ru/06/06-AT.pdf](http://quantile.ru/06/06-AT.pdf) в помощь.
3. Представьте информацию графически
4. С помощью МНК оцените модель без константы,  $R_i - R_{f,i} = \beta(R_{m,i} - R_{f,i}) + \varepsilon_i$ . Предположим, что  $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ .

5. Прокомментируйте результаты оценивания. В частности, проверьте гипотезы о значимости коэффициента и регрессии в целом.
6. С помощью МНК оцените модель с константой,  $R_i - R_{f,i} = \beta_1 + \beta_2(R_{m,i} - R_{f,i}) + \varepsilon_i$ . Предположим, что  $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ .
7. Прокомментируйте результаты оценивания. В частности, проверьте гипотезы о значимости коэффициентов и регрессии в целом.
8. Труднее всего измерить безрисковую ставку процента. Поэтому предположим, что имеющиеся у нас наблюдения — это безрисковая ставка, измеренная с ошибкой. Т.е. имеющиеся у нас наблюдения  $R_{f,i}$  представимы в виде  $R_{f,i} = R_{f,i}^{true} + u_i$ , где  $u_i \sim N(0, \sigma_u^2)$ . Величина  $R_{f,i}^{true}$  ненаблюдаема, но именно она входит в модель CAPM. Получается, что оцениваемая модель имеет вид  $R_i - R_{f,i}^{true} = \beta(R_{m,i} - R_{f,i}^{true}) + \varepsilon_i$ .
  - а) Выпишите функцию правдоподобия для оценки данной модели
  - б) Найдите оценки  $\hat{\beta}$ ,  $\hat{\sigma}_u^2$ ,  $\hat{\sigma}_\varepsilon^2$
  - в) Постройте 95%-ые доверительные интервалы
  - г) Прodelайте аналогичные действия для модели с константой
  - д) Сделайте выводы

#### Задача 2. «Цифёрки на мониторе»

При входе на каждую станцию метро есть турникеты. Рядом с турникетами в будке сидит бабушка божий одуванчик. В будке у бабушки висит монитор. На этом мониторе — прямоугольники с цифёрками.

1. Понаблюдав за изменением цифёрок, догадайтесь, что они означают.
2. Вечером какого-нибудь буднего дня запишите все цифёрки с монитора на своей родной станции метро.
3. Представьте информацию графически
4. Будем моделировать величину  $i$ -ой цифёрки пуассоновским распределением с математическим ожиданием  $\lambda_i$ . Предположим также, что  $\lambda_i = \beta_1 + \beta_2 \cdot i$ , где  $i$  — номер турникета считая от будки с бабушкой.
  - а) Выпишите функцию правдоподобия
  - б) Оцените параметры  $\beta_1$  и  $\beta_2$
  - в) Оцените ковариационную матрицу оценок  $\hat{\beta}_1$  и  $\hat{\beta}_2$
  - г) Постройте 95%-ые асимптотические доверительные интервалы для параметров
  - д) Проверьте гипотезу о том, что  $\beta_2 = 0$ . Альтернативную гипотезу сформулируйте самостоятельно.

PS. Своё смелое творчество в задачах поощряется!

#### 4.9. Домашнее задание. Титаник.

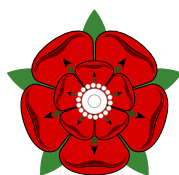
Нужно зарегистрироваться на сайте [www.kaggle.com](http://www.kaggle.com) и принять участие в конкурсе «Titanic: Machine Learning from Disaster». Крайний срок сдачи отчёта: в ночь с 14 на 15 апреля 2013 года.



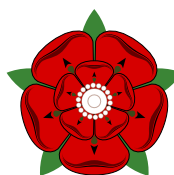
### Версия Белой Розы



1. Домашнее задание можно делать в одиночку или группой из двух человек.
2. А можно всё-таки группой из трёх человек? Нет :)
3. Письменный отчёт должен содержать как-минимум:
  - а) Логин группы
  - б) Графический анализ имеющихся данных
  - в) Результаты оценивания logit и probit моделей
  - г) Графический анализ logit и probit моделей
  - д) «Если бы я был пассажиром Титаника, то я спасся бы с вероятностью...».  
С помощью logit и probit моделей необходимо построить 95%-ый доверительный интервал для вероятности спасения каждого из участников группы, сдающей домашку. Пол и возраст взять фактические, а остальные объясняющие переменные — по своему желанию.



### Версия Алой Розы



1. Домашнее задание можно делать только в одиночку :)
2. Нет, нельзя :)
3. Письменный отчёт должен содержать как-минимум:
  - а) Логин
  - б) Графический анализ имеющихся данных
  - в) Результаты оценивания logit и probit моделей
  - г) Прогнозирование с использованием Random Forest
  - д) Прогнозирование с использованием метода опорных векторов (SVM)
  - е) Графический анализ оценённых моделей
  - ж) «Если бы я был пассажиром Титаника, то я спасся бы с вероятностью...».  
С помощью логит и пробит моделей необходимо построить 95%-ый доверительный интервал для своей вероятности спасения. Для Random Forest требуется только точечная оценка вероятности спасения. Пол и возраст взять фактические, а остальные объясняющие переменные — по своему желанию.

## 5. 2013-2014

### 5.1. Праздник 1. Вперед в рукопашную!

1. Найдите длины векторов  $a = (2, 1, 1)$  и  $b = (-2, 0, 1)$  и косинус угла между ними.
2. Сформулируйте теорему о трёх перпендикулярах
3. Для матрицы
$$A = \begin{pmatrix} -2 & 0 & 0 \\ 0 & 3 & 4 \\ 0 & 4 & 9 \end{pmatrix}$$
  - а) Найдите собственные числа и собственные векторы матрицы.
  - б) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.
  - в) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.
  - г) Представьте  $A^{2013}$  в виде произведения трёх матриц.
4. Матрицы  $A$  и  $B$  таковы, что  $\det(AB)$ ,  $\det(BA)$ ,  $\text{tr}(AB)$  и  $\text{tr}(BA)$  определены. Возможно ли что  $\det(AB) \neq \det(BA)$ ? Возможно ли, что  $\text{tr}(AB) \neq \text{tr}(BA)$ ? Если неравенство возможно, то приведите пример.
5. Вася и Петя независимо друг от друга решают тест по теории вероятностей. В тесте всего два вопроса. На каждый вопрос два варианта ответа. Петя знает решение каждого вопроса с вероятностью 0,4. Если Петя не знает решения, то он отвечает равновероятно наугад. Вася знает решение каждого вопроса с вероятностью 0,7. Если Вася не знает решения, то он отвечает равновероятно наугад.
  - а) Какова вероятность того, что Петя правильно ответил на оба вопроса?
  - б) Какова вероятность того, что Петя правильно ответил на оба вопроса, если его ответы совпали с Васиними?
  - в) Чему равно математическое ожидание числа Петиних верных ответов?
  - г) Чему равно математическое ожидание числа Петиних верных ответов, если его ответы совпали с Васиними?
6. Для случайных величин  $X$  и  $Y$  заданы следующие значения:  $\mathbb{E}(X) = 1$ ,  $\mathbb{E}(Y) = 4$ ,  $\mathbb{E}(XY) = 8$ ,  $\text{Var}(X) = \text{Var}(Y) = 9$ . Для случайных величин  $U = X + Y$  и  $V = X - Y$  вычислите:
  - а)  $\mathbb{E}(U)$ ,  $\text{Var}(U)$ ,  $\mathbb{E}(V)$ ,  $\text{Var}(V)$ ,  $\text{Cov}(U, V)$
  - б) Можно ли утверждать, что случайные величины  $U$  и  $V$  независимы?
7. Вася ведёт блог. Обозначим  $X_i$  — количество слов в  $i$ -ой записи. После первого года он по своим записям обнаружил, что  $\bar{X}_{200} = 95$  и выборочное стандартное отклонение равно 282 слова. На уровне значимости  $\alpha = 0.10$  проверьте гипотезу о том, что  $\mu = 100$  против альтернативной гипотезы  $\mu \neq 100$ . Найдите также точное Р-значение.

### 5.2. Праздник 2. Мегаматрица

В рамках классической линейной модели с детерминистическими регрессорами найдите  $\text{Var}(\hat{\beta})$ ,  $\text{Cov}(\hat{\varepsilon}, \hat{\beta})$ ,  $\text{Cov}(\hat{\varepsilon}, \hat{y})$ .

### 5.3. Праздник 3. Базовая задача

Пусть регрессионная модель  $y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$ ,  $i = 1, \dots, n$ , задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1 \ \beta_2 \ \beta_3)^T$ . Известно, что ошибки  $\varepsilon$  нормально распределены с  $\mathbb{E}\varepsilon = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что:

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

Для удобства расчётов ниже приведены матрицы:

$$X^T X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X^T X)^{-1} = \begin{pmatrix} 0.5 & -0.5 & 0 \\ -0.5 & 1 & -0.5 \\ 0 & -0.5 & 1.5 \end{pmatrix}.$$

1. Оценки  $\hat{\beta}$
2. Спрогнозируйте  $y$ , если  $x_2 = 1$  и  $x_3 = -2$
3.  $TSS, ESS, RSS, R^2$
4.  $\mathbb{E}(\hat{\sigma}^2)$
5.  $\hat{\sigma}^2$
6.  $\text{Var}(\varepsilon_1)$
7.  $\text{Var}(\beta_1)$
8.  $\text{Var}(\hat{\beta}_1)$
9.  $\widehat{\text{Var}}(\hat{\beta}_1)$
10.  $\text{Cov}(\hat{\beta}_2, \hat{\beta}_3)$
11.  $\widehat{\text{Cov}}(\hat{\beta}_2, \hat{\beta}_3)$
12.  $\text{Var}(\hat{\beta}_2 - \hat{\beta}_3)$
13.  $\widehat{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3)$
14.  $\text{Var}(\beta_2 - \beta_3)$
15. Проверьте гипотезу  $H_0: \beta_1 = 1$  против гипотезы  $H_a: \beta_1 \neq 1$  на уровне значимости 5%
16. Проверьте гипотезу  $H_0: \beta_2 = 0$  против гипотезы  $H_a: \beta_2 \neq 0$  на уровне значимости 10%
17. Проверьте гипотезу  $H_0: \beta_2 = \beta_3$  против гипотезы  $H_a: \beta_2 \neq \beta_3$  на уровне значимости 5%

### 5.4. Праздник 3. Ответы

1.  $\hat{\beta}_1 = 1.5, \hat{\beta}_2 = 3, \hat{\beta}_3 = 1.5$
2.  $\hat{y} = 1.5$
3.  $TSS = 10, RSS = 1, ESS = 9, R^2 = 0.9$
4.  $\mathbb{E}(\hat{\sigma}^2) = \sigma^2$
5.  $\hat{\sigma}^2 = \frac{RSS}{n-k} = 0.5$

6.  $\text{Var}(\varepsilon_1) = 0.5$
7.  $\text{Var}(\beta_1) = 0$
8.  $\text{Var}(\hat{\beta}_1) = \sigma^2 \cdot 0.5$
9.  $\widehat{\text{Var}}(\hat{\beta}_1) = 0.25$
10.  $\text{Cov}(\hat{\beta}_2, \hat{\beta}_3) = \sigma^2 \cdot (-0.5)$
11.  $\widehat{\text{Cov}}(\hat{\beta}_2, \hat{\beta}_3) = -0.25$
12.  $\text{Var}(\hat{\beta}_2 - \hat{\beta}_3) = \sigma^2 \cdot 3.5$
13.  $\widehat{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3) = 1.75$
14.  $\text{Var}(\beta_2 - \beta_3) = 0$
15.  $\frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)} \sim t_2, t_{\text{obs}} = \frac{1.5-1}{\sqrt{0.5}} \approx 0.7, t_{\text{crit}} = 4.3$ , нет оснований отвергать  $H_0$
16.  $\frac{\hat{\beta}_2 - \beta_2}{\text{se}(\hat{\beta}_2)} \sim t_2, t_{\text{obs}} = \frac{3}{1} = 3, t_{\text{crit}} = 2.9$ , основная гипотеза отвергается
17.  $\frac{\hat{\beta}_2 - \hat{\beta}_3 - (\beta_2 - \beta_3)}{\text{se}(\hat{\beta}_2 - \hat{\beta}_3)} \sim t_2, t_{\text{obs}} = \frac{3-1.5}{\sqrt{1.75}} \approx 1.1, t_{\text{crit}} = 4.3$ , нет оснований отвергать  $H_0$

## 5.5. Праздник 4

1. Пусть  $y = X\beta + \varepsilon$  — регрессионная модель, где  $\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$ . Пусть  $Z = XD$ , где

$$D = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}. \text{ Рассмотрите «новую» регрессионную модель } y = Z\alpha + u, \text{ где}$$

$$\alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}. \text{ Определите, как выражаются «новые» МНК-коэффициенты через «старые»}.$$

2. Рассмотрим модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 w_i + \beta_4 z_i + \varepsilon_i$ . При оценке модели по 24 наблюдениям оказалось, что  $RSS = 15$ ,  $\sum (y_i - \bar{y} - w_i + \bar{w})^2 = 20$ . На уровне значимости 1% протестируйте гипотезу

$$H_0 : \begin{cases} \beta_2 + \beta_3 + \beta_4 = 1 \\ \beta_2 = 0 \\ \beta_3 = 1 \\ \beta_4 = 0 \end{cases}$$

3. По 47 наблюдениям оценивается зависимость доли мужчин занятых в сельском хозяйстве от уровня образованности и доли католического населения по Швейцарским кантонам в 1888 году.

$$\text{Agriculture}_i = \beta_1 + \beta_2 \text{Examination}_i + \beta_3 \text{Catholic}_i + \varepsilon_i$$

---

```

1 h <- swiss
2 model1 <- glm(Agriculture ~ Examination + Catholic, data = h)
3 coef.t <- coeftest(model1)
4 dimnames(coef.t)[[2]] <-

```

```

5 c("Оценка", "Ст. ошибка", "t-статистика", "P-значение")
6 coef.t <- coef.t[, -4]
7 coef.t[1, 1] <- NA
8 coef.t[2, 2] <- NA
9 coef.t[3, 3] <- NA
10 xtable(coef.t)

```

	Оценка	Ст. ошибка	t-статистика
(Intercept)		8.72	9.44
Examination	-1.94		-5.08
Catholic	0.01	0.07	

- а) Заполните пропуски в таблице
- б) Укажите коэффициенты, значимые на 10% уровне значимости.
- в) Постройте 99%-ый доверительный интервал для коэффициента при переменной Catholic
4. Рассмотрим модель:  $y_i = \beta_1 + \beta_2 x_{1i} + \beta_3 x_{2i} + \beta_4 x_{3i} + \beta_5 x_{4i} + \varepsilon_i$ . По 20 наблюдениям оценены следующие регрессии:

$$\hat{y}_i = 10.01 + 1.05x_1 + 2.06x_2 + 0.49x_3 - 1.31x_4, RSS = 6.85$$

(s.e.)      (0.15)      (0.06)      (0.04)      (0.06)      (0.06)

$$y_i - \widehat{x_1 - 2x_2} = 10.00 + 0.50x_3 - 1.32x_4, RSS = 8.31$$

(s.e.)      (0.15)      (0.07)      (0.06)

$$y_i + \widehat{x_1} + 2x_2 = 9.93 + 0.56x_3 - 1.50x_4, RSS = 4310.62$$

(s.e.)      (3.62)      (1.48)      (1.42)

$$y_i - \widehat{x_1} + 2x_2 = 10.71 + 0.09x_3 - 1.28x_4, RSS = 3496.85$$

(s.e.)      (3.26)      (1.33)      (1.28)

$$y_i + \widehat{x_1 - 2x_2} = 9.22 + 0.97x_3 - 1.54x_4, RSS = 516.23$$

(s.e.)      (1.25)      (0.51)      (0.49)

На уровне значимости 5% проверьте гипотезу  $H_0 : \begin{cases} \beta_2 = 1 \\ \beta_3 = 2 \end{cases}$  против альтернативной гипотезы  $H_a : |\beta_2 - 1| + |\beta_3 - 2| \neq 0$ .

## 5.6. Праздник 5. Максимальное правдоподобие

- Случайные величины  $X_1, \dots, X_n$  — независимы и одинаково распределены с функцией плотности  $f(t) = \frac{\theta \cdot (\ln t)^{\theta-1}}{t}$  при  $t \in [1; e]$ . По выборке из 100 наблюдений оказалось, что  $\sum \ln(\ln(X_i)) = -30$ 
  - Найдите ML оценку параметра  $\theta$
  - Постройте 95% доверительный интервал для  $\theta$
  - С помощью LR, LM и W теста проверьте гипотезу о том, что  $\theta = 1$ .
- Величины  $X_1, \dots, X_n$  — независимы и нормально распределены,  $N(\mu, \sigma^2)$ . По 100 наблюдениям  $\sum X_i = 100$  и  $\sum X_i^2 = 900$ .
  - Найдите ML оценки неизвестных параметров  $\mu$  и  $\sigma^2$ .
  - Постройте 95%-ые доверительные интервалы для  $\mu$  и  $\sigma^2$



- в) С помощью LR, LM и W теста проверьте гипотезу о том, что  $\sigma^2 = 1$ .
- г) С помощью LR, LM и W теста проверьте гипотезу о том, что  $\sigma^2 = 1$  и одновременно  $\mu = 2$ .

Всех участников правдоподобной контрольной с древнерусским эконометрическим праздником!

Сегодня Акси́нья-полухлебница.

«На Акси́нью гадали о ценах на хлеб в ближайшее время и на будущий урожай: брали печёный хлеб и взвешивали его сначала вечером, а потом утром. Коли вес оставался неизменным — цена на хлеб не изменится. Если за ночь вес уменьшался — значит, хлеб подешевеет, а если увеличивался, то подорожает»

Wikipedia

## 5.7. Переписывание кр 5. Максимальное правдоподобие

- По совету Лисы Волк опустил в прорубь хвост и поймал 100 чудо-рыб. Веса рыбин независимы и имеют распределение Вейбулла,  $f(x) = 2 \exp(-x^2/a^2) \cdot x/a^2$  при  $x \geq 0$ . Известно, что  $\sum x_i^2 = 120$ .
  - Найдите ML оценку параметра  $a$
  - Постройте 95% доверительный интервал для  $a$
  - С помощью LR, LM и W теста проверьте гипотезу о том, что  $a = 1$ .
- Как известно, Фрекен-Бок пьет коньяк по утрам и иногда видит привидения. За 110 дней имеются следующие статистические данные

Рюмок	1	2	3
Дней с привидениями	10	25	20
Дней без привидений	20	25	10

Вероятность увидеть привидение зависит от того, сколько рюмок коньяка было выпито утром, а именно,  $p = \exp(a + bx)/(1 + \exp(a + bx))$ , где  $x$  — количество рюмок, а  $a$  и  $b$  — неизвестные параметры.

- Найдите<sup>1</sup> ML оценки неизвестных параметров  $a$  и  $b$ .
- Постройте 95%-ые доверительные интервалы для  $a$  и  $b$
- С помощью LR, LM и W теста проверьте гипотезу о том, что  $b = 0$ .
- С помощью LR, LM и W теста проверьте гипотезу о том, что  $a = 0$  и одновременно  $b = 0$ .

Всем участникам переписывания правдоподобной контрольной счастья! Много!

Сегодня, 20 марта, **Международный День счастья**.

<sup>1</sup>Здесь потребуется максимизировать функцию в R. Если этот пункт не получился, то в последующих пунктах можно считать, что  $\hat{a} = -1.5$ , а  $\hat{b} = 0.5$ . Это сильно округленные значения коэффициентов.

## 5.8. Праздник 6. Гетероскедастичность

1. Желая протестировать наличие гетероскедастичности в модели  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \beta_4 w_i + \varepsilon_i$ , эконометресса Глафира решила провести тест Уайта и получила во вспомогательной регрессии  $R^2 = 0.50$ . Глафира строит модель удоя по 200 коровам. Помогите ей провести тест на уровне значимости 5%.
2. На всякий случай эконометресса Глафира решила подстраховаться и провести тест Голдфельда-Квандта. Но она совсем забыла, как его делать. Напомните Глафире, как провести тест Голдфельда-Квандта, если она подозревает, что дисперсия  $Var(\varepsilon_i)$  возрастает с ростом  $z_i$ . Чётко напишите гипотезы  $H_0$ ,  $H_a$ , методику проведения теста, правило согласно которому отвергается или не отвергается  $H_0$ .
3. Имеются три наблюдения,  $x = (1, 2, 2)'$ ,  $y = (2, 1, 0)'$ . Предполагая, что в модели  $y_i = \beta x_i + \varepsilon_i$  имеется гетероскедастичность вида  $Var(\varepsilon_i) = \sigma^2 x_i^4$  найдите:
  - а) Обычную МНК-оценку параметра  $\beta$
  - б) Самую эффективную среди несмещённых оценок параметра  $\beta$
  - в) Во сколько раз отличается истинная дисперсия этих двух оценок?
  - г) Во сколько раз отличаются оценки дисперсий этих оценок, если дисперсии оцениваются без поправки на гетероскедастичность в обоих случаях?

## 5.9. Большой Устный ЗАчёт

1. Метод Наименьших Квадратов.
  - а) МНК-картинка
  - б) Нахождение всего-всего, если известен вектор  $y$  и матрица  $X$
2. Теорема Гаусса-Маркова
  - а) Формулировка с детерминистическими регрессорами
  - б) Доказательство с детерминистическими регрессорами
  - в) Формулировки со стохастическими регрессорами
  - г) Что даёт дополнительное предположение о нормальности  $\varepsilon$ ?
3. Проверка гипотез о линейных ограничениях
  - а) Проверка гипотезы о значимости коэффициента
  - б) Проверка гипотезы о значимости регрессии в целом
  - в) Проверка гипотезы об одном линейном соотношении с помощью ковариационной матрицы
  - г) Ограниченная и неограниченная модель
  - д) Тест Чоу на стабильность коэффициентов
  - е) Тест Чоу на прогнозную силу
4. Метод максимального правдоподобия
  - а) Свойства оценок
  - б) Два способа получения оценки дисперсии
  - в) Три теста (LM, Wald, LR)
  - г) Выписать функцию ML для обычной регрессии
  - д) для AR(1) процесса
  - е) для MA(1) процесса

- ж) для логит модели
  - з) для пробит модели
  - и) для модели с заданным видом гетероскедастичности
5. Мультиколлинеарность
- а) Определение, последствия
  - б) Величины, измеряющие силу мультиколлинеарности
  - в) Методы борьбы
  - г) Сюда же: метод главных компонент, хотя он используется и для других целей
6. Гетероскедастичность
- а) Определение, последствия
  - б) Тесты, график
  - в) Стюдентизированные остатки
  - г) НС оценки ковариации
  - д) GLS и FGLS
7. Временные ряды
- а) Стационарный временной ряд
  - б) ACF, PACF
  - в) Модель ARMA
  - г) Модель GARCH (не будет, не успели)
8. Логит и пробит
- а) Описание моделей
  - б) Предельные эффекты
  - в) Чувствительность, специфичность
  - г) Кривая ROC
9. Эндогенность
- а) Три примера: одновременность, пропущенные переменные, ошибки измерения
  - б) IV, двухшаговый МНК
10. Модели панельных данных
- а) RE, FE, сквозная регрессии
  - б) Тест Хаусмана
11. Альтернативные методы. Уметь объяснить суть метода. Уметь реализовать его в R.
- а) Метод опорных векторов (не будет, не успели)
  - б) Классификационные деревья и случайный лес
12. R. Можно принести файл со своей заготовкой, можно пользоваться Интернетом для поиска информации, но не для общения.
- а) Загрузить данные из .csv файла в R
  - б) Посчитать описательные статистики (среднее, мода, медиана и т.д.)
  - в) Построить подходящие описательные графики для переменных
  - г) Оценить линейную регрессию с помощью МНК. Провести диагностику на что-нибудь (гетероскедастичность, автокорреляцию, мультиколлинеарность).
  - д) Оценить logit, probit модели, посчитать предельные эффекты
  - е) Оценить ARMA модель
  - ж) Выделить главные компоненты

## 5.10. Экзамен.

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{pmatrix}.$$

- Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
  - Проверьте гипотезу  $\beta_2 = 0$  против альтернативной о неравенстве на уровне значимости 5%
2. По данным о пассажирах Титаника оценивается логит-модель. Зависимая переменная *survived* равна 1, если пассажир выжил. Объясняющая переменная *sexmale* равна 1 для мужчин.

---

```

1 library(texreg)
2 # here I need data source
3 mod.tit <- glm(data = titanic, survived ~ age + sex, family = "binomial")
4 texreg(mod.tit, float.pos = "h!", label = "table:titanic")

```

---

	Model 1
(Intercept)	1.92*** (0.28)
age	-0.01 (0.01)
sexmale	-2.84*** (0.21)
AIC	633.45
BIC	646.80
Log Likelihood	-313.72
Deviance	627.45
Num. obs.	633

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

Таблица 2: Statistical models

- Оцените вероятность выжить для женщины 20 лет
  - Оцените предельный эффект увеличения возраста для женщины 20 лет
  - С помощью какого метода оценивается логит-модель? Каким образом при этом получают оценки стандартных ошибок коэффициентов?
3. Теорема Гаусса-Маркова.

- а) Аккуратно сформулируйте теорему Гаусса-Маркова для нестохастических регрессоров.
  - б) Поясните каждое из свойств оценок, фигурирующих в теореме.
  - в) Как меняются свойства оценок МНК при нарушении предпосылки теоремы о том, что дисперсия  $\varepsilon_i$  постоянна?
4. Рассмотрим временной ряд, описываемый МА(2) моделью,

$$y_t = \gamma + \varepsilon_t + \alpha_1 \varepsilon_{t-1} + \alpha_2 \varepsilon_{t-2},$$

где  $\varepsilon_t$  — белый шум с  $\text{Var}(\varepsilon_t) = \sigma^2$ .

- а) Является ли данный процесс стационарным? Что такое стационарный процесс?
  - б) Найдите автокорреляционную функцию данного процесса,  $\rho(k) = \text{Corr}(y_t, y_{t-k})$ .
  - в) Выпишите функцию правдоподобия для данной модели в предположении нормальности  $\varepsilon_t$ .
5. Рассмотрите модель  $y_i = \beta x_i + \varepsilon_i$ . Предположим, что все предпосылки классической линейной регрессионной модели выполнены. Модель оценивается с помощью МНК и получается оценка  $\hat{\beta}_{OLS}$ . В условиях мультиколлинеарности для снижения дисперсии оценки  $\hat{\beta}$  можно применять ряд методов, например, алгоритм «ridge regression». Он состоит в том, что при некотором фиксированном  $\lambda \geq 0$  минимизируется по  $\hat{\beta}$  величина

$$Q(\hat{\beta}) = \sum_i (y_i - \hat{\beta} x_i)^2 + \lambda \hat{\beta}^2$$

- а) Как выглядит МНК оценка  $\hat{\beta}_{OLS}$ ?
- б) Как выглядит оценка методом «ridge regression»,  $\hat{\beta}_{RR}$ ?
- в) Верно ли, что оценка  $\hat{\beta}_{RR}$  является несмещенной только при  $\lambda = 0$ ?
- г) (\*) Верно ли, что всегда найдется такое  $\lambda$ , что среднеквадратичная ошибка оценки  $\hat{\beta}_{RR}$  будет меньше, т.е.  $\mathbb{E}((\hat{\beta}_{RR} - \beta)^2) < \mathbb{E}((\hat{\beta}_{OLS} - \beta)^2)$ ?

## 5.11. Пересдача экзамена

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{pmatrix}.$$

- а) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - б) Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
  - в) Проверьте гипотезу  $\beta_2 = 0$  против альтернативной о неравенстве на уровне значимости 5%
2. По данным о пассажирах Титаника оценивается логит-модель. Зависимая переменная *survived* равна 1, если пассажир выжил. Объясняющая переменная *sexmale* равна 1 для мужчин.

---

```

1 library(texreg)
2 # data source
3 mod.tit <- glm(data = titanic, survived ~ age + sex, family = "binomial")
4 texreg(mod.tit, float.pos = "h!", label = "table:titanic-2")

```

---

	Model 1
(Intercept)	1.92*** (0.28)
age	-0.01 (0.01)
sexmale	-2.84*** (0.21)
AIC	633.45
BIC	646.80
Log Likelihood	-313.72
Deviance	627.45
Num. obs.	633

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

Таблица 3: Statistical models

- а) Оцените вероятность выжить для женщины 20 лет
  - б) Оцените предельный эффект увеличения возраста для женщины 20 лет
  - в) С помощью какого метода оценивается логит-модель? Каким образом при этом получают оценки стандартных ошибок коэффициентов?
3. Теорема Гаусса-Маркова.
- а) Аккуратно сформулируйте теорему Гаусса-Маркова для нестохастических регрессоров.
  - б) Поясните каждое из свойств оценок, фигурирующих в теореме.
  - в) Как меняются свойства оценок МНК при нарушении предпосылки теоремы о том, что дисперсия  $\varepsilon_i$  постоянна?
4. Для линейной регрессии  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  была выполнена сортировка наблюдений по возрастанию переменной  $x$ . Исходная модель оценивалась по разным частям выборки:

Выборка	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$RSS$
$i = 1, \dots, 50$	0.93	2.02	3.38	145.85
$i = 1, \dots, 21$	1.12	2.01	3.32	19.88
$i = 22, \dots, 29$	0.29	2.07	2.24	1.94
$i = 30, \dots, 50$	0.87	1.84	3.66	117.46

Известно, что ошибки в модели являются независимыми нормальными случайными величинами с нулевым математическим ожиданием.

- а) Предполагая гомоскедастичность остатков на уровне значимости 5% проверьте гипотезу, что исследуемая зависимость одинакова на всех трёх частях всей выборки.
- б) Протестируйте ошибки на гетероскедастичность на уровне значимости 5%.
- в) Какой тест можно на гетероскедастичность можно было бы использовать, если бы не было уверенности в нормальности остатков? Опишите пошагово процедуру этого теста.

## 5.12. Домашняя работа 1. RLMS и гетероскедастичность

1. Прочитайте про RLMS, <http://www.hse.ru/rlms/>  
Посмотрите описание проекта. Пролистайте вестник RLMS, чтобы иметь представление о том, какие исследования можно строить на основе RLMS.
2. Скачайте любую волну RLMS по своему выбору. Скачайте описание переменных. Пролистайте описание переменных. Там их больше тысячи. Попадаются довольно прикольные. Мне нравится pc9.6.5a, «У Вас есть GPRS навигатор?»
3. Загрузите данные в R.  
Данные RLMS выложены на сайте в формате SPSS. SPSS это потихоньку погибающий статистический пакет для домохозяек. Для чтения формата .sav в таблицу данных R можно сделать так

---

```
1 library(foreign)
2 file.name <- "/home/boris/downloads/r20hall23c.sav"
3 h <- read.spss(file.name, to.data.frame = TRUE)
```

---

Первая команда, library(foreign), подгружает библиотеку R, в которой содержатся команды для чтения вражеских форматов, spss, stata, etc  
Описания переменных при этом также загружаются в таблицу данных. Можно их выделить в отдельный вектор и прочитать, например, про переменную pc9.631a.

---

```
1 var.labels <- attr(h, "variable.labels")
2 var.labels["pc9.631a"]
```

---

4. Выберите любую количественную переменную в качестве зависимой и несколько переменных в качестве объясняющей.  
Цель этой домашки скорее ознакомится с наличием мониторинга RLMS, поэтому можно не сильно заморачиваться с этим этапом. Хотя в реальности тут-то всё самое интересное и начинается. За оригинальные гипотезы будут плюшки. Кстати, неплохо бы дать выбранным переменным понятные названия.
5. Опишите выбранные переменные.  
Постройте симпатичные графики. Посчитайте описательные статистики. Много ли пропущенных наблюдений? Есть ли что-нибудь интересенькое?
6. Постройте регрессию зависимой переменной на объясняющие.  
Проверьте гипотезу о значимости каждого полученного коэффициента. Проверьте гипотезу о значимости регрессии в целом. Для нескольких коэффициентов (двух достаточно) постройте 95%-ый доверительный интервал.
7. Разберитесь с возможным наличием гетероскедастичности в данных.  
С какой переменной может быть связана дисперсия  $\text{Var}(\varepsilon_i)$ ? Проведите визуальный анализ на гетероскедастичность. Проведите формальные тесты на гетероскедастичность. Примените оценки дисперсии  $\hat{\beta}$  устойчивые к гетероскедастичности. Прокомментируйте. Может помочь [http://bdemeshev.github.io/r\\_cycle/cycle\\_files/12\\_hetero.html](http://bdemeshev.github.io/r_cycle/cycle_files/12_hetero.html)
8. Покажите буйство своей фантазии и аккуратность!  
Не стоит думать, что побуквенное выполнение этих инструкций гарантирует оценку в десять баллов. Эконометрика — это не ремесло, а искусство! Фантазируйте! Убедите меня в работе, что вы были на лекциях, даже если это так :) Аккуратность в виде подписанных осей на графиках, указанных единицах измерения также не повредит.

9. Срок сдачи — 27 февраля 2014 года.

Работа принимается исключительно в печатном виде с применением грамотного программирования R +  $\text{\LaTeX}$ . Каждый день более поздней сдачи умножает оценку за работу на 0.8. Работа должна представлять слитный текст, код скрывать не нужно. В конце должна быть команда `sessionInfo()`.

### 5.13. Домашняя работа 2. Титаник

1. Зарегистрируйтесь на сайте [www.kaggle.com](http://www.kaggle.com) в конкурсе «Titanic: Machine Learning from Disaster». В работе укажите login, использованный при регистрации.
2. Проанализируйте данные графически и с помощью описательных статистик (среднее, мода, медиана и т.д.)  
Прокомментируйте графики, обратите внимание на количество пропущенных значений.
3. Оцените logit и probit модели.  
Приведите оценки моделей. Какие коэффициенты значимы? Прокомментируйте знак коэффициентов. Посчитайте и сравните предельные эффекты.
4. Оцените random forest и SVM модели.  
Параметры методов подберите с помощью кросс-валидации. Можно применять любые другие подходы, не только random forest и SVM. Другой подход следует описать в тексте.
5. «Если бы я был пассажиром Титаника, то я спасся бы с вероятностью...».  
С помощью логит и пробит моделей постройте 95%-ый доверительный интервал для вероятности своего спасения. Для random forest — только точечный прогноз вероятности, для svm — только прогноз типа «да»/«нет».
6. Подумайте, чем можно заполнить пропущенные значения. Заполните пропущенные значения и заново оцените logit, random forest и svm. Насколько сильно меняется качество оцененных моделей?
7. Сравните все использованные подходы по прогнозной силе на тестовой выборке с сайта. Какой оказался наилучшим?
8. При прогнозировании и расчете предельных эффектов используйте свои фактические пол и возраст, а остальные объясняющие переменные — выбирайте согласно своей фантазии :)
9. Срок сдачи — 30 апреля 2014 года.  
Работа принимается исключительно в печатном виде с применением грамотного программирования R +  $\text{\LaTeX}$ . Каждый день более поздней сдачи умножает оценку за работу на 0.8. Работа должна представлять слитный текст, код скрывать не нужно. В конце должна быть команда `sessionInfo()`.
10. Популярные ошибки прошлой домашки будут караться со всей строгостью военного времени!  
Цикл заметок про R в помощь [http://bdemeshev.github.io/r\\_cycle/](http://bdemeshev.github.io/r_cycle/).

## 6. 2014-2015

### 6.1. Праздник номер 1

Вперёд, в рукопашную!



1. Сформулируйте теорему о трёх перпендикулярах и обратную к ней.
2. Для матрицы
 
$$A = \begin{pmatrix} 3 & 4 \\ 4 & 9 \end{pmatrix}$$
  - а) Найдите собственные числа и собственные векторы матрицы.
  - б) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.
  - в) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.
  - г) Найдите  $A^{42}$
  - д) Не находя  $A^{100}$  найдите  $\text{tr}(A^{100})$  и  $\det(A^{100})$
3. Игрок получает случайным образом 13 карт из колоды в 52 карты.
  - а) Какова вероятность, что у него как минимум два туза?
  - б) Каково ожидаемое количество тузов у игрока?
  - в) Какова вероятность, что у него как минимум два туза, если известно, что у него есть хотя бы один туз?
  - г) Каково ожидаемое количество тузов у игрока, если известно, что у него на руках хотя бы один туз?
4. В ходе анкетирования 100 сотрудников банка «Омега» ответили на вопрос о том, сколько времени они проводят на работе ежедневно. Среднее выборочное оказалось равно 9.5 часам при выборочном стандартном отклонении 0.5 часа.
  - а) Постройте 95% доверительный интервал для математического ожидания времени проводимого сотрудниками на работе
  - б) Проверьте гипотезу о том, что в среднем люди проводят на работе 10 часов, против альтернативной гипотезы о том, что в среднем люди проводят на работе меньше 10 часов, укажите точное Р-значение.

## 6.2. Праздник номер 2

### Паниковать на контрольной строго воспрещается! :)

1. По 47 наблюдениям оценивается зависимость доли мужчин занятых в сельском хозяйстве от уровня образованности и доли католического населения по Швейцарским кантонам в 1888 году.

$$Agriculture_i = \beta_1 + \beta_2 Examination_i + \beta_3 Catholic_i + \varepsilon_i$$

---

```

1 h <- swiss
2 model1 <- glm(Agriculture ~ Examination + Catholic, data = h)
3 coef.t <- coeftest(model1)
4 dimnames(coef.t)[[2]] <-
5   c("Оценка", "Ст. ошибка", "t-статистика", "Р-значение")
6 coef.t <- coef.t[, -4]
7 coef.t[1, 1] <- NA
8 coef.t[2, 2] <- NA
9 coef.t[3, 3] <- NA
10 xtable(coef.t)

```

---

	Оценка	Ст. ошибка	t-статистика
(Intercept)		8.72	9.44
Examination	-1.94		-5.08
Catholic	0.01	0.07	

- Заполните пропуски в таблице
  - Укажите коэффициенты, значимые на 10% уровне значимости.
  - Постройте 99%-ый доверительный интервал для коэффициента при переменной Catholic
2. В рамках классической линейной модели с неслучайными регрессорами найдите  $\text{Var}(\hat{\varepsilon})$ ,  $\text{Cov}(\hat{\beta}, \hat{\varepsilon})$ . Верно ли, что  $\text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_2) = 0$ ?

---

```

1 X <- model.matrix(model1)
2 B <- t(X) %*% X
3 colnames(B) <- NULL
4 rownames(B) <- NULL
5 XXm <- solve(B)
6 x <- xtable(B, align = rep("", ncol(B) + 1), digits = 0)
7 xm <- xtable(XXm, align = rep("", ncol(B) + 1), digits = 5)
8 print(xm, floating = FALSE, tabular.environment = "bmatrix",
9       hline.after = NULL, include.rownames = FALSE, include.colnames = FALSE)

```

---

3. Эконометресса Ефросинья оценивала модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$ . Найдя матрицы  $X'X$  и  $(X'X)^{-1}$ , она призадумалась...

$$X'X = \begin{bmatrix} 47 & 775 & 1934 \\ 775 & 15707 & 23121 \\ 1934 & 23121 & 159570 \end{bmatrix}, (X'X)^{-1} = \begin{bmatrix} 0.26653 & -0.01067 & -0.00168 \\ -0.01067 & 0.00051 & 0.00006 \\ -0.00168 & 0.00006 & 0.00002 \end{bmatrix}$$

- Помогите Ефросинье найти количество наблюдений,  $\bar{z}$ ,  $\sum x_i z_i$ ,  $\sum (x_i - \bar{x})(z_i - \bar{z})$
  - (\*) Ефросинья решила зачем-то также оценить модель  $x_i = \gamma_1 + \gamma_2 z_i + u_i$ . Как она может найти RSS в новой модели в одно арифметическое действие?
4. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 2 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{pmatrix}.$$

- Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
- Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
- Проверьте гипотезу  $\beta_2 = 0$  против альтернативной о неравенстве на уровне значимости 5%

### 6.3. Праздник номер 2, ответы

1. а)  $\hat{\beta}_1 = 82.3, se(\hat{\beta}_2) = 0.38, t_{obs} = 1/7$   
б)  $\hat{\beta}_1, \hat{\beta}_2$   
в)  $[0.01 - 0.07 \cdot 2.704; 0.01 + 0.07 \cdot 2.704]$
2.  $\text{Var}(\hat{\varepsilon}) = \sigma^2(I - H), \text{Cov}(\hat{\beta}, \hat{\varepsilon}) = 0$
3. а)  $n = 47, \bar{z} = 1934/47, \sum x_i z_i = 23121, \sum (x_i - \bar{x})(z_i - \bar{z}) = 23121 - 1934 \cdot 775/47$   
б)  $RSS = \sigma^2 / \text{Var}(\hat{\beta}_2)$
4. а)  $\hat{\beta}_1 = -0.5, \hat{\beta}_2 = 4, \hat{\beta}_3 = -1.5$   
б)  $\hat{\sigma}^2 = 4.5$   
в) Нет оснований отвергать основную гипотезу.

### 6.4. Праздник номер 3

Примечание: во всех задачах, если явно не сказано обратное, предполагается, что выполнены стандартные предпосылки классической линейной регрессионной модели.

1. Рассмотрим следующую модель зависимости почасовой оплаты труда  $W$  от уровня образования  $Educ$ , возраста  $Age$ , уровня образования родителей  $Fathedu$  и  $Mothedu$ :

$$\widehat{\ln W} = \hat{\beta}_1 + \hat{\beta}_2 Educ + \hat{\beta}_3 Age + \hat{\beta}_4 Age^2 + \hat{\beta}_5 Fathedu + \hat{\beta}_6 Mothedu$$

$$R^2 = 0.341, n = 27$$

- а) Напишите спецификацию регрессии с ограничениями для проверки статистической гипотезы  $H_0 : \beta_5 = 2\beta_4$
  - б) Дайте интерпретацию проверяемой гипотезе
  - в) Для регрессии с ограничением был вычислен коэффициент  $R_R^2 = 0.296$ . На уровне значимости 5% проверьте нулевую гипотезу
2. По ежегодным данным с 2002 по 2009 год оценивался тренд в динамике общей стоимости экспорта из РФ:  $Exp_t = \beta_1 + \beta_2 t + \varepsilon_t$ , где  $t$  — год ( $t = 0$  для 2002 г.,  $t = 1$  для 2003 г., ...,  $t = 7$  для 2009 г.),  $Exp_t$  — стоимость экспорта из РФ во все страны в млрд. долл. Оценённое уравнение выглядит так:  $\widehat{Exp}_t = 111.9 + 43.2t$ . Получены также оценки дисперсии случайной ошибки  $\hat{\sigma}^2 = 4009$  и ковариационной матрицы оценок коэффициентов:

$$\widehat{Var}(\hat{\beta}) = \begin{pmatrix} 1671 & -334 \\ -334 & 95 \end{pmatrix}$$

- а) Постройте 95%-ый доверительный интервал для коэффициента  $\beta_2$
  - б) Спрогнозируйте стоимость экспорта на 2010 год и постройте 90%-ый предиктивный интервал для прогноза.
3. Имеется 100 наблюдений. Исследователь Вениамин предполагает, что дисперсия случайной ошибки в последних 50-ти наблюдениях в 4 раза выше, чем в первых 50-ти, в частности  $\text{Var}(\varepsilon_1) = \sigma^2$ , а  $\text{Var}(\varepsilon_{100}) = 4\sigma^2$ . Вениамин оценивает модель  $y_i = \beta x_i + \varepsilon_i$  с помощью МНК.
  - а) Найдите истинную дисперсию МНК оценки коэффициента  $\beta$
  - б) Предложите более эффективную оценку  $\hat{\beta}^{alt}$

- в) Чему равна истинная дисперсия новой оценки?
  - г) Подробно опишите любой способ, который позволяет протестировать гипотезу о гомоскедастичности против предположения Вениамина о дисперсии.
4. Закон больших чисел гласит, что если  $z_i$  независимы и одинаково распределены, то  $\text{plim } \bar{z}_n = \mathbb{E}(z_1)$ . Предположим, что регрессоры — стохастические, а именно, наблюдения являются случайной выборкой (то есть отдельные наблюдения независимы и одинаково распределены), и  $\mathbb{E}(\varepsilon|X) = 0$ . Модель имеет вид:

$$y_i = \beta_1 + \beta_2 x_i + \beta_3 w_i + \varepsilon_i$$

- а) Найдите  $\mathbb{E}(\varepsilon)$ ,  $\mathbb{E}(x_1 \cdot \varepsilon_1)$
  - б) Найдите  $\text{plim } \frac{1}{n} X' \varepsilon$
  - в) Найдите  $\text{plim } \frac{1}{n} X' X$
  - г) Докажите, что вектор МНК оценок  $\hat{\beta}$  является состоятельным
5. Эконометресса Эвридика хочет оценить модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$ . К сожалению, она измеряет зависимую переменную с ошибкой. Т.е. вместо  $y_i$  она знает значение  $y_i^* = y_i + u_i$  и использует его в качестве зависимой переменной при оценке регрессии. Ошибки измерения  $u_i$  некоррелированы между собой и с  $\varepsilon_i$ , имеют нулевое математическое ожидание и постоянную дисперсию  $\sigma_u^2$ .
- а) Будут ли оценки Эвридики несмещенными?
  - б) Могут ли дисперсии оценок Эвридики быть ниже чем дисперсии МНК оценок при использовании настоящего  $y_i$ ?
  - в) Могут ли оценки дисперсий оценок Эвридики быть ниже чем оценок дисперсий МНК оценок при использовании настоящего  $y_i$ ?

## 6.5. Миникр

### Миникр 5

1. Как проверить гипотезу об одновременной незначимости всех коэффициентов регрессии кроме свободного члена? Укажите  $H_0$ ,  $H_a$ , тестовую статистику и её распределение при верной  $H_0$ .
2. Рассмотрим модель  $y = X\beta + \varepsilon$ , где  $n$  — количество наблюдений,  $k$  — количество коэффициентов и  $\varepsilon_i$  — одинаково распределены и независимы.
  - а) Укажите вид матрицы  $X$
  - б) Выпишите формулу для МНК оценки  $\hat{\beta}$
  - в) Выпишите формулу для ковариационной матрицы оценок  $\hat{\beta}$
3. Опишите подробно тест Чоу на стабильность коэффициентов по двум наборам данных из  $n_1$  и  $n_2$  наблюдений соответственно. Число оцениваемых коэффициентов равно  $k$ .
4. Рассмотрим модель со свободным членом. Как вычисляются  $R^2$  и скорректированный  $R_{adj}^2$ ? Что может произойти с этими величинами при увеличении количества регрессоров? При уменьшении?
5. Какую гипотезу можно проверить, зная отношение  $ESS/RSS$ ? Укажите  $H_0$ ,  $H_a$ , тестовую статистику и её распределение при верной  $H_0$ .
6. Опишите подробно тест Чоу на прогнозную силу по двум наборам данных из  $n_1$  и  $n_2$  наблюдений соответственно. Число оцениваемых коэффициентов равно  $k$ .

## 6.6. Зачет. Базовый поток

1. Сформулируйте теорему Гаусса-Маркова применительно к модели  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ . Поясните смысл каждого используемого термина.
2. Как проверить гипотезу о том, что эксцесс случайной выборки совпадает с эксцессом нормально распределенной случайной величины? Аккуратно укажите проверяемые  $H_0$ ,  $H_a$ , используемую статистику и её асимптотический закон распределения при верной  $H_0$ .
3. Рассмотрим модель спроса на продукцию трёх фирм  $y_i = \beta_1 + \beta_2 x_i + \beta_3 d_{i1} + \beta_4 d_{i2} + \varepsilon_i$ . Здесь  $x_i$  — цена, а  $y_i$  — величина спроса.

Дамми переменные определены следующим образом:

	$d_{i1}$	$d_{i2}$
Фирма 1	0	1
Фирма 2	0	0
Фирма 3	1	0

- а) Как проверить гипотезу, что спрос на продукцию трёх фирм совпадает? Укажите  $H_0$ ,  $H_a$ , тестовую статистику, закон распределения статистики при верной  $H_0$ .

- б) Дамми-переменные  $d_{i1}$  и  $d_{i2}$  заменяют на  $d_{i3}$  и  $d_{i4}$ :

	$d_{i3}$	$d_{i4}$
Фирма 1	0	0
Фирма 2	1	0
Фирма 3	1	1

В новых переменных модель имеет вид  $y_i = \beta'_1 + \beta'_2 x_i + \beta'_3 d_{i1} + \beta'_4 d_{i2} + \varepsilon_i$ . Как новые коэффициенты  $\beta'$  выражаются через старые коэффициенты  $\beta$ ?

4. Что можно сказать об оценке МНК  $\hat{\beta}_2$  в модели  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$  при наличии ошибок измерений  $x_i$ ? А при наличии ошибок измерений  $y_i$ ?
5. Рассмотрим модель  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ , где  $\varepsilon_i \sim N(0; \sigma^2)$ .
  - а) Напишите выражения для оценок дисперсий и ковариации коэффициентов, т.е.  $\widehat{\text{Var}}(\hat{\beta}_1)$ ,  $\widehat{\text{Var}}(\hat{\beta}_2)$ ,  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2)$
  - б) Найдите математическое ожидание и дисперсию каждой из выписанных оценок
  - в) Какой закон распределения с точностью до масштабирования имеют эти оценки?
6. Для модели данных  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  по 100 наблюдениям получены результаты:

---

```

1 set.seed(777)
2 x <- rnorm(100)
3 z <- rnorm(100)
4 x <- x - mean(x)
5 pre_model <- lm(z ~ x)
6 z <- resid(pre_model) # X'X is diagonal :)
7
8 y <- 2 + 3*x - 0.05*z + rnorm(100)
9 model <- lm(y ~ x + z)
10
11 xtable(model)

```

---

- а) Выпишите полученное уравнение регрессии

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	2.1304	0.1174	18.14	0.0000
x	2.9263	0.1228	23.83	0.0000
z	-0.0562	0.1148	-0.49	0.6252

- б) Укажите, какие коэффициенты значимы при  $\alpha = 0.05$
- в) Проверьте  $H_0: \beta_2 - \beta_3 = 3$  предполагая, что оценки коэффициентов  $\beta_2$  и  $\beta_3$  независимы
7. Рассмотрим модель данных  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ , где  $\varepsilon_i \sim N(0; \sigma^2)$ .
- а) Выпишите формулу для оценок коэффициентов и оценки дисперсии ошибок
- б) Укажите математическое ожидание и дисперсию выписанных оценок
- в) Для оценок коэффициентов укажите закон распределения
- г) Для оценки дисперсии ошибок укажите закон распределения с точностью до масштабирования

## 6.7. Зачет, 26.12.2014. Ликвидация безграмотности

В этот день, 26 декабря 1919 года, совнарком РСФСР принял декрет «О ликвидации безграмотности в РСФСР». Всем желаю отметить этот день написанием грамотного зачета по эконометрике! Удачи!

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $E(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 2 & 1 \\ 2 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{3} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 4 & -3 \\ 0 & -3 & 6 \end{pmatrix}.$$

- а) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
- б) Найдите коэффициент детерминации  $R^2$
- в) Предполагая нормальное распределение вектора  $\varepsilon$ , проверьте гипотезу  $H_0: \beta_2 = 0$  против альтернативной  $H_a: \beta_2 \neq 0$
2. Для линейной регрессии  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  была выполнена сортировка наблюдений по возрастанию переменной  $x$ . Исходная модель оценивалась по разным частям выборки:

Выборка	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	RSS
$i = 1, \dots, 50$	0.93	2.02	3.38	145.85
$i = 1, \dots, 21$	1.12	2.01	3.32	19.88
$i = 22, \dots, 29$	0.29	2.07	2.24	1.94
$i = 30, \dots, 50$	0.87	1.84	3.66	117.46

Известно, что ошибки в модели являются независимыми нормальными случайными величинами с нулевым математическим ожиданием.

- а) Предполагая гомоскедастичность остатков на уровне значимости 5% проверьте гипотезу, что исследуемая зависимость одинакова на всех трёх частях всей выборки.
  - б) Протестируйте ошибки на гетероскедастичность на уровне значимости 5%.
  - в) Какой тест можно на гетероскедастичность можно было бы использовать, если бы не было уверенности в нормальности остатков? Опишите пошагово процедуру этого теста.
3. По 2040 наблюдениям оценена модель зависимости стоимости квартиры в Москве (в 1000\$) от общего метража и метража жилой площади.

---

```

1 model1 <- lm(price ~ totsp + livesp, data = flats)
2 report <- summary(model1)
3 coef.table <- report$coefficients
4 rownames(coef.table) <- c("Константа", "Общая площадь", "Жилая площадь")
5 xtable(coef.table)
6 xtable(vcov(model1))

```

---

	Estimate	Std. Error	t value	Pr(> t )
Константа	-88.81	4.37	-20.34	0.00
Общая площадь	1.70	0.10	17.78	0.00
Жилая площадь	1.99	0.18	10.89	0.00

Оценка ковариационной матрицы  $\widehat{\text{Var}}(\hat{\beta})$  имеет вид

	(Intercept)	totsp	livesp
(Intercept)	19.07	0.03	-0.45
totsp	0.03	0.01	-0.02
livesp	-0.45	-0.02	0.03

Оценка стандартной ошибки случайной составляющей,  $\hat{\sigma} = 33.03$ .

- а) Можно ли интерпретировать коэффициент при переменной *totsp* как стоимость одного метра нежилой площади?
  - б) Проверьте гипотезу о том, что коэффициенты при регрессорах *totsp* и *livesp* равны.
  - в) Постройте 95%-ый доверительный интервал для ожидаемой стоимости квартиры с жилой площадью 30 м<sup>2</sup> и общей площадью 60 м<sup>2</sup>.
  - г) Постройте 95%-ый прогнозный интервал для фактической стоимости квартиры с жилой площадью 30 м<sup>2</sup> и общей площадью 60 м<sup>2</sup>.
4. Аккуратно сформулируйте теорему Гаусса-Маркова
- а) для нестохастических регрессоров
  - б) для стохастических регрессоров в предположении, что наблюдения являются случайной выборкой

## 6.8. Домашняя работа 1. RLMS и гетероскедастичность

1. Прочитайте про RLMS, <http://www.hse.ru/rlms/>  
Посмотрите описание проекта. Пролистайте вестник RLMS, чтобы иметь представление о том, какие исследования можно строить на основе RLMS.

2. Скачайте любую волну RLMS по своему выбору. Скачайте описание переменных. Проллистайте описание переменных. Там их больше тысячи. Попадаютсь довольно прикольные. Мне нравится rs9.6.5a, «У Вас есть GPRS навигатор?»
3. Загрузите данные в R.  
Данные RLMS выложены на сайте в формате SPSS. SPSS это потихоньку погибающий статистический пакет для домохозяек. Для удобства можно воспользоваться готовой функцией для чтения данных RLMS в пакете rlms.

---

```
1 library("rlms")  
2 h <- read.rlms("/home/boris/downloads/r20hall23c.sav")
```

---

Про установку пакета rlms можно прочитать на страничке <https://github.com/bdemeshev/rlms>

Описания переменных при этом также загружаются в таблицу данных. Можно их посмотреть:

---

```
1 var_meta <- attr(h,"var_meta")  
2 var_meta
```

---

4. Выберите любую количественную переменную в качестве зависимой и несколько переменных в качестве объясняющих.  
Цель этой домашки скорее ознакомится с наличием мониторинга RLMS, поэтому можно не сильно заморачиваться с этим этапом. Хотя в реальности тут-то всё самое интересное и начинается. За оригинальные гипотезы будут плюшки. Кстати, неплохо бы дать выбранным переменным понятные названия.
5. Опишите выбранные переменные.  
Постройте симпатичные графики. Посчитайте описательные статистики. Много ли пропущенных наблюдений? Есть ли что-нибудь интересненькое?
6. Постройте регрессию зависимой переменной на объясняющие.  
Проверьте гипотезу о значимости каждого полученного коэффициента. Проверьте гипотезу о значимости регрессии в целом. Для нескольких коэффициентов (двух достаточно) постройте 95%-ый доверительный интервал.
7. Разберитесь с возможным наличием гетероскедастичности в данных.  
С какой переменной может быть связана дисперсия  $\text{Var}(\varepsilon_i)$ ? Проведите визуальный анализ на гетероскедастичность. Проведите формальные тесты на гетероскедастичность. Примените оценки дисперсии  $\hat{\beta}$  устойчивые к гетероскедастичности. Прокомментируйте. Может помочь [http://bdemeshev.github.io/r\\_cycle/cycle\\_files/12\\_hetero.html](http://bdemeshev.github.io/r_cycle/cycle_files/12_hetero.html)
8. Покажите буйство своей фантазии и аккуратность!  
Не стоит думать, что побуквенное выполнение этих инструкций гарантирует оценку в десять баллов. Эконометрика — это не ремесло, а искусство! Фантазируйте! Убедите меня в работе, что вы были на лекциях, даже если это не так :) Аккуратность в виде подписанных осей на графиках, указанных единицах измерения также не повредит.
9. Срок сдачи — 12 января 2015 года.  
Работа принимается исключительно в печатном виде с применением грамотного программирования R +  $\text{\LaTeX}$  или markdown. Каждый день более поздней сдачи умножает оценку за работу на 0.8. Работа должна представлять слитный текст, код скрывать не нужно. В конце должна быть команда sessionInfo().



## 6.9. Экзамен. Демо-1

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{pmatrix}.$$

- Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
  - Проверьте гипотезу  $\beta_2 = 0$  против альтернативной о неравенстве на уровне значимости 5%
2. По данным о пассажирах Титаника оценивается логит-модель. Зависимая переменная `survived` равна 1, если пассажир выжил. Объясняющая переменная `sexmale` равна 1 для мужчин.

---

```

1 library(texreg)
2 mod_tit <- glm(data = titanic, survived ~ age + sex, family = "binomial")
3 texreg(mod_tit, float.pos = "h!")

```

---

	Model 1
(Intercept)	1.92*** (0.28)
age	-0.01 (0.01)
sexmale	-2.84*** (0.21)
AIC	633.45
BIC	646.80
Log Likelihood	-313.72
Deviance	627.45
Num. obs.	633

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

- Оцените вероятность выжить для женщины 20 лет
- Оцените предельный эффект увеличения возраста для женщины 20 лет
- С помощью какого метода оценивается логит-модель? Каким образом при этом получают оценки стандартных ошибок коэффициентов?

3. Теорема Гаусса-Маркова.

- а) Аккуратно сформулируйте теорему Гаусса-Маркова для нестохастических регрессоров.
- б) Поясните каждое из свойств оценок, фигурирующих в теореме.
- в) Как меняются свойства оценок МНК при нарушении предпосылки теоремы о том, что дисперсия  $\varepsilon_i$  постоянна?

4. Рассмотрим временной ряд, описываемый МА(2) моделью,

$$y_t = \gamma + \varepsilon_t + \alpha_1 \varepsilon_{t-1} + \alpha_2 \varepsilon_{t-2},$$

где  $\varepsilon_t$  — белый шум с  $\text{Var}(\varepsilon_t) = \sigma^2$ .

- а) Является ли данный процесс стационарным? Что такое стационарный процесс?
  - б) Найдите автокорреляционную функцию данного процесса,  $\rho(k) = \text{Corr}(y_t, y_{t-k})$ .
  - в) Выпишите функцию правдоподобия для данной модели в предположении нормальности  $\varepsilon_t$ .
5. Рассмотрите модель  $y_i = \beta x_i + \varepsilon_i$ . Предположим, что все предпосылки классической линейной регрессионной модели выполнены. Модель оценивается с помощью МНК и получается оценка  $\hat{\beta}_{OLS}$ . В условиях мультиколлинеарности для снижения дисперсии оценки  $\hat{\beta}$  можно применять ряд методов, например, алгоритм гребневой регрессии. Он состоит в том, что при некотором фиксированном  $\lambda \geq 0$  минимизируется по  $\hat{\beta}$  величина

$$Q(\hat{\beta}) = \sum_i (y_i - \hat{\beta} x_i)^2 + \lambda \hat{\beta}^2$$

- а) Как выглядит МНК оценка  $\hat{\beta}_{OLS}$ ?
- б) Как выглядит оценка методом «ridge regression»,  $\hat{\beta}_{RR}$ ?
- в) Верно ли, что оценка  $\hat{\beta}_{RR}$  является несмещенной только при  $\lambda = 0$ ?
- г) (\*) Верно ли, что всегда найдется такое  $\lambda$ , что среднеквадратичная ошибка оценки  $\hat{\beta}_{RR}$  будет меньше, то есть  $\mathbb{E}((\hat{\beta}_{RR} - \beta)^2) < \mathbb{E}((\hat{\beta}_{OLS} - \beta)^2)$ ?

## 6.10. Экзамен. Демо-2

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 2 & 1 \\ 2 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{3} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 4 & -3 \\ 0 & -3 & 6 \end{pmatrix}.$$

- а) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
- б) Найдите коэффициент детерминации  $R^2$ .
- в) Предполагая нормальное распределение вектора  $\varepsilon$ , проверьте гипотезу  $H_0: \beta_2 = 0$  против альтернативной  $H_a: \beta_2 \neq 0$ .

2. Для линейной регрессии  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  была выполнена сортировка наблюдений по возрастанию переменной  $x$ . Исходная модель оценивалась по разным частям выборки:

Выборка	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$RSS$
$i = 1, \dots, 50$	0.93	2.02	3.38	145.85
$i = 1, \dots, 21$	1.12	2.01	3.32	19.88
$i = 22, \dots, 29$	0.29	2.07	2.24	1.94
$i = 30, \dots, 50$	0.87	1.84	3.66	117.46

Известно, что ошибки в модели являются независимыми нормальными случайными величинами с нулевым математическим ожиданием.

- Предполагая гомоскедастичность остатков на уровне значимости 5% проверьте гипотезу, что исследуемая зависимость одинакова на всех трёх частях всей выборки.
  - Протестируйте ошибки на гетероскедастичность на уровне значимости 5%.
  - Какой тест можно на гетероскедастичность можно было бы использовать, если бы не было уверенности в нормальности остатков? Опишите пошагово процедуру этого теста.
3. По 2040 наблюдениям оценена модель зависимости стоимости квартиры в Москве (в 1000\$) от общего метража и метража жилой площади.

```

1 model1 <- lm(price ~ totsp + livesp, data = flats)
2 report <- summary(model1)
3 coef.table <- report$coefficients
4 rownames(coef.table) <- c("Константа", "Общая площадь", "Жилая площадь")
5 xtable(coef.table)
6 xtable(vcov(model1))

```

	Estimate	Std. Error	t value	Pr(> t )
Константа	-88.81	4.37	-20.34	0.00
Общая площадь	1.70	0.10	17.78	0.00
Жилая площадь	1.99	0.18	10.89	0.00

Оценка ковариационной матрицы  $\widehat{Var}(\hat{\beta})$  имеет вид

	(Intercept)	totsp	livesp
(Intercept)	19.07	0.03	-0.45
totsp	0.03	0.01	-0.02
livesp	-0.45	-0.02	0.03

- Можно ли интерпретировать коэффициент при переменной  $totsp$  как стоимость одного метра нежилой площади?
  - Проверьте гипотезу о том, что коэффициенты при регрессорах  $totsp$  и  $livesp$  равны.
  - Постройте 95%-ый доверительный интервал для ожидаемой стоимости квартиры с жилой площадью 30 м<sup>2</sup> и общей площадью 60 м<sup>2</sup>.
  - Постройте 95%-ый прогнозный интервал для фактической стоимости квартиры с жилой площадью 30 м<sup>2</sup> и общей площадью 60 м<sup>2</sup>.
4. Предположим, что в классической линейной модели ошибки имеют нормальное распределение, т.е.

$$y_i = \beta_1 + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i} + \varepsilon_i$$

где  $\varepsilon_i$  нормальны  $N(0, \sigma^2)$  и независимы

- а) Найдите оценки для  $\beta$  и  $\sigma^2$  методом максимального правдоподобия.
  - б) Являются ли полученные оценки  $\hat{\beta}_{ML}$  и  $\hat{s}_{ML}^2$  несмещенными?
  - в) Выведите формулу  $LR$ -статистики у теста отношения правдоподобия для тестирования гипотезы об адекватности регрессии  $H_0: \beta_2 = \beta_3 = \dots = \beta_k = 0$ .
5. В модели есть три регрессора,  $x_1$ ,  $x_2$  и  $x_3$ . Для удобства будем считать, что они центрированы и нормированы, т.е. выборочное среднее каждого регрессора равно нулю, а выборочная дисперсия — единице. Эти три регрессора являются столбцами матрицы  $X$ . Известно, что

$$X'X = \begin{pmatrix} 100 & 0 & 0 \\ 0 & 100 & 90 \\ 0 & 90 & 100 \end{pmatrix}$$

- а) Найдите число обусловленности матрицы  $X'X$ .
- б) Выразите первые две главные компоненты через  $x_1$ ,  $x_2$  и  $x_3$

6. По данным о пассажирах Титаника оценивается логит-модель. Зависимая переменная *survived* равна 1, если пассажир выжил.

```

1 library(texreg)
2 mod_tit <- glm(data = titanic, survived ~ age + sex,
3               family = "binomial")
4 texreg(mod_tit, float.pos = "h!")

```

	Model 1
(Intercept)	1.92*** (0.28)
age	-0.01 (0.01)
sexmale	-2.84*** (0.21)
AIC	633.45
BIC	646.80
Log Likelihood	-313.72
Deviance	627.45
Num. obs.	633

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

- Оцените вероятность выжить для мужчины 30 лет
- Оцените предельный эффект увеличения возраста для мужчины 30 лет
- С помощью какого метода оценивается логит-модель? Каким образом получаются оценки стандартных ошибок коэффициентов?

## 6.11. Экзамен. 15.06.15

Больше двух столетий назад, 15 июня 1763 года Екатерина II издала манифест, запрещающий произнесение необдуманных речей, опасных для общественного спокойствия. Давайте применим этот манифест к экзамену по эконометрике!

1. Регрессионная модель задана в матричном виде при помощи уравнения  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $E(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix} \text{ и } (X'X)^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{pmatrix}.$$

- Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
- Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
- Предполагая нормальность  $\varepsilon_i$  проверьте гипотезу  $\beta_2 = 0$  против альтернативной о неравенстве на уровне значимости 5%

2. Аккуратно сформулируйте теорему Гаусса-Маркова в парадигме стохастических регрессоров для ситуации случайной выборки
3. Немного вопросов про гетероскедастичность:
  - а) Что такое гетероскедастичность?
  - б) К каким последствиям она приводит?
  - в) Что можно предпринять в условиях гетероскедастичности и что эта мера даёт?
4. Для МА(1) модели  $y_t = 3 + \varepsilon_t + 0.5\varepsilon_{t-1}$  посчитайте автокорреляционную и частную автокорреляционную функцию.
5. По 100 наблюдениям была оценена логит модель,  $\hat{y}_i^* = 2.3 - 2x_i$ .
  - а) Оцените вероятность  $P(y_i = 1)$  для  $x_i = 1$
  - б) Оцените предельный эффект  $dP(y_i = 1)/dx$  для  $x_i = 1$
6. Рассмотрим классическую модель парной регрессии  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$  с нестохастическими регрессорами и нормальными остатками. Модель оценивается по 100 наблюдениям. Известно, что  $F$ -статистика, проверяющая гипотезу о незначимости регрессии в целом, равна 50.  
Рассчитайте значения статистики множителей Лагранжа,  $LM$ , статистики Вальда,  $W$ , для проверки гипотезы о незначимости регрессии в целом.
7. Иван Андреевич Крылов хочет оценить, как зависит количество написанных им за день строчек басни,  $y_i$ , от количества съеденных булочек,  $x_i$ . То есть его интересует коэффициент  $\beta_2$  в уравнении

$$y_i = \beta_1 + \beta_2 x_i + \varepsilon_i.$$

Однако когда его окрыляет вдохновение, он очень невнимательно считает, поэтому величина  $x_i$  ненаблюдаема. Кухарка и жена всегда готовы сказать, сколько булочек Крылов якобы съел за день. Эти величины содержат ошибку измерения, то есть  $x_i^A = x_i + u_i^A$  и  $x_i^B = x_i + u_i^B$ .

Ошибки измерения  $u_i^A, u_i^B$ , случайная составляющая  $\varepsilon_i$  независимы,  $u_i^A \sim N(0, \sigma_A^2)$ ,  $u_i^B \sim N(0, \sigma_B^2)$ ,  $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ . Иван Андреевич располагает информацией о 100 случайно выбранных днях.

- а) Рассмотрим оценку  $\hat{\beta}_2$ , получаемую с помощью обычного МНК в регрессии  $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i^A$ . Будет ли она состоятельной?
- б) Рассмотрим двухшаговый МНК. На первом шаге строим регрессию  $x_i^A$  на  $x_i^B$  и получаем прогнозные значения  $\hat{x}_i^A$ . На втором шаге оцениваем регрессию  $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 \hat{x}_i^A$ . Будет ли оценка  $\hat{\beta}_2$  состоятельной?

## 7. 2015-2016

### 7.1. Праздник номер 1. Вспомнить всё! 15.09.2015

1. Найдите длины векторов  $a = (1, 1, 1, 1)$  и  $b = (1, 2, 3, 4)$  и косинус угла между ними. Найдите один любой вектор, перпендикулярный вектору  $a$ .
2. Сформулируйте теорему о трёх перпендикулярах

3. Для матрицы

$$A = \begin{pmatrix} 10 & 15 \\ 15 & 26 \end{pmatrix}$$

- а) Найдите собственные числа и собственные векторы матрицы
  - б) Найдите  $\det(A)$ ,  $\text{tr}(A)$
  - в) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа
4. Известно, что  $X$  — матрица размера  $n \times k$  и  $n > k$ , известно, что  $X'X$  обратима. Рассмотрим матрицу  $H = X(X'X)^{-1}X'$ . Укажите размер матрицы  $H$ , найдите  $H^{2015}$ ,  $\text{tr}(H)$ ,  $\det(H)$ , собственные числа матрицы  $H$ . Штрих означает транспонирование.
5. Для случайных величин  $X$  и  $Y$  заданы следующие значения:  $\mathbb{E}(X) = 1$ ,  $\mathbb{E}(Y) = 4$ ,  $\mathbb{E}(XY) = 8$ ,  $\text{Var}(X) = \text{Var}(Y) = 16$ . Для случайных величин  $U = X + Y$  и  $V = X - Y$  вычислите:
- а)  $\mathbb{E}(U)$ ,  $\text{Var}(U)$ ,  $\mathbb{E}(V)$ ,  $\text{Var}(V)$ ,  $\text{Cov}(U, V)$
  - б) Можно ли утверждать, что случайные величины  $U$  и  $V$  независимы?
6. Вася ведёт блог. Обозначим  $X_i$  — количество слов в  $i$ -ой записи. После первого года он по 200 своим записям обнаружил, что  $\bar{X}_{200} = 95$  и выборочное стандартное отклонение равно 300 слов. На уровне значимости  $\alpha = 0.15$  проверьте гипотезу о том, что  $\mu = 100$  против альтернативной гипотезы  $\mu \neq 100$ . Постройте 85-ти процентный доверительный интервал для  $\mu$ .
7. Саша и Маша решали одну и ту же задачу. Саша правильно решает задачу с вероятностью 0.8, Маша, независимо от Саши (!), с вероятностью 0.7. Какова вероятность того, что Маша верно решила задачу, если задачу верно решил только кто-то один из них?

## 7.2. Праздник номер 1. Вспомнить всё! 15.09.2015, решение

Автор решения: Кирилл Пономарёв

1. Ну тут все понятно
2. Правильная формулировка: «Если прямая, проведенная на плоскости через основание наклонной, перпендикулярна её проекции, то она перпендикулярна и самой наклонной»
3. а) Собственные значения  $\lambda_i$ :

$$\begin{vmatrix} 10 - \lambda & 15 \\ 15 & 26 - \lambda \end{vmatrix} = \lambda^2 - 36\lambda + 35 = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 35$$

Собственный вектор  $h$  соответствующий собственному значению  $\lambda$  по определению:

$$Ah = \lambda h \Rightarrow (A - \lambda \cdot I)h = 0$$

То есть для  $\lambda_1$ :

$$\begin{pmatrix} 9 & 15 \\ 15 & 25 \end{pmatrix} \cdot \begin{pmatrix} h_{11} \\ h_{21} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Откуда:

$$h_{11} = -\frac{5}{3} \cdot h_{21} \quad \text{то есть, например, вектор: } \begin{pmatrix} 5 \\ -3 \end{pmatrix}$$

Аналогично находится собственный вектор, соответствующий  $\lambda_2 = 35$

- б) Подставив  $\lambda = 0$  в первый определитель, получим  $\det(A) = 35$ , а  $\text{tr}(A) = \lambda_1 + \lambda_2 = 36$
- в) По определению  $Ah = \lambda h$ . Домножив на  $A^{-1}$  (если она существует) обе части, получим:

$$h = \lambda A^{-1}h \Rightarrow A^{-1}h = \frac{1}{\lambda}h$$

То есть собственные значения для  $A^{-1}$  это  $1/\lambda_i$  а собственные векторы такие же, как и у матрицы  $A$ .

4. Размер матрицы  $H$ :

$$\begin{matrix} X & (X'X)^{-1} & X' & = & H \\ n \times k & k \times k & k \times n & n \times n \end{matrix}$$

Заметим что  $H^k = H$  для  $k = 1, 2, \dots$

$$H^2 = X \underbrace{(X'X)^{-1}X'X}_{=I} (X'X)^{-1}X' = X(X'X)^{-1}X' = H$$

Такие матрицы называются идемпотентными, а для них:  $\text{tr}(A) = \text{rank}(A)$ . Найдем след, пользуясь свойством что  $\text{tr}(AB) = \text{tr}(BA)$ .

$$\text{tr}(H) = \text{tr}(X(X'X)^{-1}X') = \text{tr}((X'X)^{-1}XX') = \text{tr}(I_k) = k$$

То есть в данном случае  $\text{rank}(H) = k$ . Так как  $\text{rank}(H) < n$ , определитель этой матрицы равен 0.

Для того чтобы найти собственные числа ( $\lambda$ ), снова вспомним определение:

$$Hv = \lambda v$$

Домножив обе части этого уравнения на  $H$ , получим:

$$H^2v = \lambda \cdot Hv \Leftrightarrow Hv = \lambda^2 v$$

Можем делать так сколько угодно раз, то есть если  $\lambda$  — собственное число матрицы  $H$ , то и  $\lambda^k$  тоже. Значит это могут быть только 0 или 1. Тот факт, что 0 является собственным значением сразу вытекает из неполного ранга.

5. а)

$$\mathbb{E}(U) = \mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y) = 5$$

$$\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) = 8 - 4 = 4$$

$$\text{Var}(U) = \text{Var}(X) + \text{Var}(Y) + 2 \cdot \text{Cov}(X, Y) = 16 + 16 + 2 \cdot 4 = 40$$

$$\mathbb{E}(V) = \mathbb{E}(X - Y) = \mathbb{E}(X) - \mathbb{E}(Y) = -3$$

$$\text{Var}(V) = \text{Var}(X) + \text{Var}(Y) - 2 \cdot \text{Cov}(X, Y) = 16 + 16 - 2 \cdot 4 = 24$$

$$\begin{aligned} \text{Cov}(U, V) &= \text{Cov}(X+Y, X-Y) = \text{Cov}(X, X) - \text{Cov}(X, Y) + \text{Cov}(Y, X) - \text{Cov}(Y, Y) = \\ &= \text{Var}(X) - \text{Var}(Y) = 0 \quad (11) \end{aligned}$$

- б) Нельзя, это не следует из равенства ковариации нулю. Можно было бы утверждать про независимость, если бы величины имели совместное нормальное распределение.



6. Ликбез: если известно, что сами  $X_i$  нормальные, а истинной дисперсии нет, то используется распределение Стюдента, и только в этом случае! Здесь же про распределение  $X_i$  ничего не известно, но асимптотически (по ЦПТ) можно использовать нормальное.

$$\sqrt{n} \cdot \frac{\bar{X} - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

Проверяем нулевую гипотезу против двухсторонней альтернативной:

$$\begin{cases} H_0 : & \mu = 100 \\ H_{al} : & \mu \neq 100 \end{cases}$$

Это значит что критические 15 процентов должны быть распределены по двум хвостам, поэтому нас интересуют  $z_{0.075}$  и  $z_{0.925}$ . Расчетная статистика:

$$T = \sqrt{200} \cdot \frac{95 - 100}{300} = -0.24 > -1.44 = z_{0.075}$$

Значит нулевая гипотеза не отвергается. Строим доверительный интервал:

$$\mathbb{P} \left( z_{0.925} < \sqrt{n} \cdot \frac{\bar{X} - \mu}{\sigma} < z_{0.925} \right) = 0.85$$

Так как распределение Стюдента симметрично:

$$\mathbb{P} \left( \bar{X} - z_{0.925} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{0.925} \cdot \frac{\sigma}{\sqrt{n}} \right) = 0.85$$

Таким образом:

$$\mathbb{P}(64.46 < \mu < 125.43) = 0.85$$

7. Пусть событие  $A$ : «Маша верно решила задачу», а событие  $B$ : «Задачу решил только кто-то один». По формуле Байеса:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A) \cdot \mathbb{P}(A)}{\mathbb{P}(B)}$$

Здесь:

$$\begin{aligned} \mathbb{P}(B|A) &= \mathbb{P}(\text{Саша не решил}) = 0.2 \\ \mathbb{P}(A) &= 0.7 \\ \mathbb{P}(B) &= 0.8 \cdot 0.3 + 0.7 \cdot 0.2 = 0.38 \end{aligned}$$

Поэтому:

$$\mathbb{P}(A|B) = \frac{0.2 \cdot 0.7}{0.38} = \frac{7}{19}$$

### 7.3. Праздник номер 2, 10 ноября 2015

- В рамках классической линейной регрессионной модели  $y = X\beta + \varepsilon$ ,  $\mathbb{E}(\varepsilon) = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ , найдите:  $\mathbb{E}(\hat{\varepsilon})$ ,  $\text{Var}(\hat{\varepsilon})$ ,  $\text{Cov}(\hat{\varepsilon}, y)$ .
- Имеются данные:

$y_i$	$x_i$	$z_i$
1	2	1
2	-1	2
3	-3	-3
4	2	0

Предположим, что ошибки нормальны  $\mathcal{N}(0; \sigma^2)$  и независимы.

- а) Оцените с помощью МНК модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$
  - б) Найдите  $RSS$ ,  $TSS$ ,  $ESS$  и  $R^2$
  - в) Проверьте гипотезу о незначимости коэффициента  $\hat{\beta}_2$  на уровне значимости 5%.
  - г) Найдите оценку ковариационной матрицы коэффициентов  $\widehat{\text{Var}}(\hat{\beta})$
  - д) Проверьте гипотезу  $H_0 : \beta_2 = \beta_3$  на уровне значимости 5%.
  - е) Для четвертого наблюдения постройте прогноз и найдите ошибку прогноза.
3. Как могут измениться (могут ли увеличиться? уменьшиться?)  $RSS$ ,  $TSS$  и  $ESS$  при добавлении дополнительного наблюдения? При добавлении дополнительного регрессора?
  4. Эконометресса Агнеса оценила множественную регрессию  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$ . Потом она добавила два наблюдения к своей выборке:  $y_{n+1} = 1$ ,  $x_{n+1} = 2$ ,  $z_{n+1} = 3$  и  $y_{n+2} = -1$ ,  $x_{n+2} = 1$ ,  $z_{n+2} = 2$ . Как при этом изменились матрицы  $X'X$  и  $X'y$ ?
  5. По 47 наблюдениям оценивается зависимость фертильности женщин от доли мужчин занятых в сельском хозяйстве и доли католического населения по Швейцарским кантонам в 1888 году.

$$Fertility_i = \beta_1 + \beta_2 Examination_i + \beta_3 Catholic_i + \varepsilon_i$$

---

```

1 library(lmtest)
2 library(apsrtable)
3 library(xtable)
4 h <- swiss
5 model1 <- lm(Fertility ~ Examination + Catholic, data = h)
6 coef.t <- coeftest(model1)
7 dimnames(coef.t)[[2]] <- c("Оценка", "Ст. ошибка", "t-статистика", "P-значение")
8 coef.t <- coef.t[, -4]
9 coef.t[1, 1] <- NA
10 coef.t[2, 2] <- NA
11 coef.t[3, 3] <- NA
12 xtable(coef.t)

```

---

	Оценка	Ст. ошибка	t-статистика
(Intercept)		4.98	16.68
Examination	-0.89		-4.08
Catholic	0.04	0.04	

- а) Заполните пропуски в таблице.
  - б) Укажите коэффициенты, значимые на 10% уровне значимости.
  - в) Постройте 95%-ый доверительный интервал для коэффициента при Examination
6. Аккуратно сформулируйте (с «Если» и «то») теорему Гаусса-Маркова для случая нестохастических регрессоров.
  7. Нарисуйте Самую Главную Картинку, иллюстрирующую метод наименьших квадратов для множественной регрессии. Отметьте на картинке  $RSS$ ,  $ESS$ ,  $TSS$  и  $R^2$
  8. Эконометресса Ефросинья оценивала модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$ . Найдя матрицы  $X'X$  и  $(X'X)^{-1}$ , она призадумалась...

$$X'X = \begin{bmatrix} 47 & 775 & 1934 \\ 775 & 15707 & 23121 \\ 1934 & 23121 & 159570 \end{bmatrix}, (X'X)^{-1} = \begin{bmatrix} 0.26653 & -0.01067 & -0.00168 \\ -0.01067 & 0.00051 & 0.00006 \\ -0.00168 & 0.00006 & 0.00002 \end{bmatrix}$$

- а) Помогите Ефросинье найти количество наблюдений,  $\bar{z}$ ,  $\sum x_i z_i$ ,  $\sum (x_i - \bar{x})(z_i - \bar{z})$
  - б) Ефросинья решила зачем-то также оценить модель  $x_i = \gamma_1 + \gamma_2 z_i + u_i$ . Как выглядят матрицы  $X'X$  и  $X'y$  для новой модели?
  - в) (\*) Как Ефросинья может найти RSS в новой модели в одно арифметическое действие?
9. Как известно,  $\hat{y} = Hy$ , где матрица-шляпница  $H$  задаётся формулой  $H = X(X'X)^{-1}X'$ .
- а) Является ли вектор остатков  $\hat{\varepsilon}$  собственным вектором матрицы  $H$ ? Если да, то какое собственное число ему соответствует?
  - б) Является ли вектор прогнозов  $\hat{y}$  собственным вектором матрицы  $H$ ? Если да, то какое собственное число ему соответствует?
  - в) Является ли регрессор  $z$  (скажем, второй столбец матрицы  $X$ ) собственным вектором матрицы  $H$ ? Если да, то какое собственное число ему соответствует?
10. Задача на компе с R.
- Подключите библиотеку ggplot2 командой `library("ggplot2")`. Рассмотрим набор данных по цене бриллиантов `diamonds`. Оцените линейную модель зависимости цены бриллианта `price` от массы `carat` и глубины `depth`. После оценки модели:
- а) Поместите  $R^2$  в переменную `r_sq`
  - б) Поместите  $RSS$  в переменную `rss`
  - в) Поместите оценку коэффициента при `carat` в переменную `hb_carat`
  - г) Поместите прогнозы в вектор `y_hat`

#### 7.4. Праздник номер 2, 10 ноября 2015, некоторые ответы

1.  $\mathbb{E}(\hat{\varepsilon}) = 0$ ,  $\text{Var}(\hat{\varepsilon}) = \sigma^2(I - H)$ ,  $\text{Cov}(\hat{\varepsilon}, y) = \sigma^2(I - H)$ .
- 2.
3. При добавлении нового наблюдения TSS и RSS не уменьшаются, а ESS может меняться в обе стороны. При добавлении дополнительного регрессора RSS не увеличивается, TSS не меняется, ESS не уменьшается.
5. а)  $\hat{\beta}_1 = 83$ ,  $se(\hat{\beta}_2) = 0.22$ ,  $t_{obs} = 1$   
 б)  $\hat{\beta}_1, \hat{\beta}_2$   
 в)  $[-0.89 - 2.021 \cdot 0.22; -0.89 + 2.021 \cdot 0.22]$
8. а)  $n = 47$ ,  $\sum x_i z_i = 23121$ ,  $\sum (x_i - \bar{x})(z_i - \bar{z}) = 23121 - \frac{23121}{47}$   
 б)  
 в)  $\text{Var}(\hat{\beta}_j) = \sigma^2 / RSS_j$ , где  $RSS_j$  — сумма квадратов остатков в регрессии j-ой объясняющей переменной на остальные объясняющие переменные, включая константу.
9. а) Да,  $\lambda = 0$   
 б) Да,  $\lambda = 1$   
 в) Да,  $\lambda = 1$

## 7.5. Блокбастер, 28-12-2015

В этот день, 28 декабря 1895 года, в индийском салоне «Гран-кафе» на бульваре Капуцинок в Париже состоялся публичный показ «Синематографа братьев Люмьер» :)

1. Регрессионная модель  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  задана в матричном виде  $y = X\beta + \varepsilon$ , где  $\beta = (\beta_1, \beta_2, \beta_3)'$ . Известно, что  $\mathbb{E}(\varepsilon) = 0$  и  $\text{Var}(\varepsilon) = \sigma^2 \cdot I$ . Известно также, что

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, X = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \end{pmatrix}.$$

Для удобства расчетов приведены матрицы

$$X'X = \begin{pmatrix} 5 & 2 & 2 \\ 2 & 2 & 0 \\ 2 & 0 & 2 \end{pmatrix} \text{ и } (X'X)^{-1} = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1.5 & 1 \\ -1 & 1 & 1.5 \end{pmatrix}.$$

- Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - Найдите коэффициент детерминации  $R^2$ .
  - Предполагая нормальное распределение вектора  $\varepsilon$ , проверьте гипотезу  $H_0: \beta_3 = 0$  против альтернативной  $H_a: \beta_3 \neq 0$  на уровне значимости 5%.
  - Постройте точечный прогноз и 95%-ый предиктивный интервал для  $x_6 = 2$  и  $z_6 = 0$ .
2. Рассмотрим модель со стохастическими регрессорами  $y = X\beta + \varepsilon$ . При этом  $\mathbb{E}(\varepsilon|X) = 0$ , как и положено, однако ошибки  $\varepsilon$  хитро зависят друг от друга, и поэтому  $\text{Var}(\varepsilon|X)$  есть некоторая известная недиагональная матрица  $V$ . Несмотря на нарушение предпосылок теоремы Гаусса-Маркова Чак Норрис использует обычный МНК для получения оценок  $\hat{\beta}$ .  
Найдите  $\mathbb{E}(\hat{\beta}|X)$ ,  $\text{Var}(\hat{\beta}|X)$  и  $\text{Cov}(\hat{y}, \hat{\varepsilon}|X)$ .
3. Рассмотрим классическую линейную регрессионную модель:

$$y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$$

- Огюст Люмьер утверждает, что при нестохастических регрессорах математические ожидания  $\mathbb{E}(y_i)$  различны. Луи Люмьер утверждает, что при стохастических регрессорах и предпосылке о том, что наблюдения являются случайной выборкой, все  $\mathbb{E}(y_i)$  равны между собой. Кто из них прав?
  - Помогите Луи Люмьеру найти  $\text{plim } \hat{\varepsilon}_1$  и  $\text{plim } \hat{y}_1$ .
4. Рассмотрим классическую линейную регрессионную модель:

$$y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$$

Наблюдения являются случайной выборкой. Истинная ковариация  $\text{Cov}(x_i, z_i)$  равна нулю. Мы оцениваем с помощью МНК две регрессии.

Регрессия 1:

$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i + \hat{\beta}_3 z_i$$

Регрессия 2:

$$\hat{y}_i = \hat{\gamma}_1 + \hat{\gamma}_2 x_i$$

- Верно ли, что  $\hat{\beta}_2$  совпадает с  $\hat{\gamma}_2$ ?

- б) Верно ли, что  $\text{plim } \hat{\beta}_2 = \text{plim } \hat{\gamma}_2$ ?
5. Аккуратно опишите процедуру сравнения с помощью  $F$ -теста двух вложенных (ограниченной и неограниченной) линейных моделей:
- Сформулируйте  $H_0$  и  $H_a$
  - Сформулируйте все предпосылки теста
  - Укажите способ подсчёта тестовой статистики
  - Укажите закон распределения тестовой статистики при верной  $H_0$
  - Сформулируйте правило, по которому делается вывод об  $H_0$
6. Чтобы не выдать себя, Джеймс Бонд оценивает с помощью МНК только однопараметрические регрессии вида  $y_i = \beta x_i + \varepsilon_i$ . Однако он знаком с теоремой Фриша-Бау.
- Сколько подобных однопараметрических регрессий ему придется оценить, чтобы получить оценку коэффициента  $\beta_3$  в множественной регрессии  $y_i = \beta_1 w_i + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$ ?
  - Укажите, какие именно регрессии нужно построить для данной цели



## 7.6. Блокбастер, 28-12-2015, ответы

- $(\hat{\beta}_1 \quad \hat{\beta}_2 \quad \hat{\beta}_3)' = (3 \quad 1.5 \quad -1.5)'$
  - $R^2 = 0.9$
  - $t_{obs} = \frac{-1.5}{\sqrt{0.5 \cdot 1.5}} \approx -1.73, t_{crit} = 4.3$ , нет оснований отвергать  $H_0$
  - $\hat{y}_6 = 6$
- 

## 7.7. Максимально правдоподобно, 25-02-2016

- Как известно, Фрекен Бок любит пить коньяк по утрам. За прошедшие пять дней она записала, сколько рюмочек коньяка выпила утром,  $x_i$ , и видела ли она в этот день привидение,  $y_i$ ,

$y_i$	1	0	1	0	0
$x_i$	2	1	3	1	0

Зависимость между  $y_i$  и  $x_i$  описывается пробит-моделью,  $\mathbb{P}(y_i = 1) = F(\beta_1 + \beta_2 x_i)$ .

- Выпишите логарифмическую функцию правдоподобия

- б) Выпишите условия первого порядка для оценки  $\beta_1$  и  $\beta_2$
2. Приведите пример небольшого набора данных для которого оценки логит модели  $\mathbb{P}(y_i = 1) = F(\beta_1 + \beta_2 x_i)$  не существуют. В наборе данных должны присутствовать хотя бы одно наблюдение  $y_i = 0$  и хотя бы одно наблюдение  $y_i = 1$ .
3. Почему в пробит-модели предполагается, что  $\varepsilon_i \sim \mathcal{N}(0; 1)$ , а не  $\varepsilon_i \sim \mathcal{N}(0; \sigma^2)$  как в линейной регрессии?
4. Исследователь Вениамин пытается понять, как логарифм количества решённых им по эконометрике задач зависит от количества съеденных им пирожков. Для этого он собрал 100 наблюдений. Первые 50 наблюдений — относятся к пирожкам с мясом, а последние 50 наблюдений — к пирожкам с повидлом. Вениамин считает, что ожидаемое количество решённых задач не зависит от начинки пирожков, а только от их количества, т.е.  $y_i = \beta x_i + u_i$ . Однако он полагает, что для пирожков с мясом —  $u_i \sim \mathcal{N}(0; \sigma_M^2)$ , а для пирожков с повидлом —  $u_i \sim \mathcal{N}(0; \sigma_J^2)$ .
- а) Выпишите логарифмическую функцию правдоподобия
- б) Выпишите условия первого порядка для оценки  $\beta$ ,  $\sigma_M^2$ ,  $\sigma_J^2$
5. При оценке логит модели  $\mathbb{P}(y_i = 1) = \Lambda(\beta_1 + \beta_2 x_i)$  по 500 наблюдениям оказалось, что  $\hat{\beta}_1 = 0.7$  и  $\hat{\beta}_2 = 3$ . Оценка ковариационной матрицы коэффициентов имеет вид

$$\begin{pmatrix} 0.04 & 0.01 \\ 0.01 & 0.09 \end{pmatrix}$$

- а) Проверьте гипотезу о незначимости коэффициента  $\hat{\beta}_2$
- б) Найдите предельный эффект роста  $x_i$  на вероятность  $\mathbb{P}(y_i = 1)$  при  $x_i = -0.5$
- в) Найдите максимальный предельный эффект роста  $x_i$  на вероятность  $\mathbb{P}(y_i = 1)$
- г) Постройте точечный прогноз вероятности  $\mathbb{P}(y_i = 1)$  если  $x_i = -0.5$
- д) Найдите стандартную ошибку построенного прогноза
6. После долгих изысканий Вениамин пришёл к выводу, что  $\beta = 0$ , т.е. что логарифм количества решенных им по эконометрике за вечер задач имеет нормальное распределение  $y_i$  с математическим ожиданием ноль. Однако он по прежнему уверен, что дисперсия  $y_i$  зависит от того, какие пирожки он ел в этом вечер. Для пирожков с повидлом  $y_i \sim \mathcal{N}(0; \sigma_J^2)$ , а для пирожков с мясом —  $y_i \sim \mathcal{N}(0; \sigma_M^2)$ . Всего 100 наблюдений. Первые 50 вечеров относятся к пирожкам с мясом, последние 50 вечеров — к пирожкам с повидлом:

$$\sum_{i=1}^{50} y_i = 10, \sum_{i=1}^{50} y_i^2 = 100, \sum_{i=51}^{100} y_i = -10, \sum_{i=51}^{100} y_i^2 = 300$$

- а) Найдите оценки  $\sigma_M^2$ ,  $\sigma_J^2$ , которые получит Вениамин.
- б) Помогите Вениамину проверить гипотезу  $\sigma_M^2 = \sigma_J^2$  с помощью тестов отношения правдоподобия, множителей Лагранжа и Вальда.

## 7.8. Решение задач КР по Эконометрике, 2015-2016

1. Функция максимального правдоподобия

$$L = \prod_{i=1}^5 \mathbb{P}(Y_i = j), \text{ где } j \in \{0, 1\} \implies \quad (12)$$

$$L = \mathbb{P}(Y_1 = 1) \times \mathbb{P}(Y_2 = 0) \times \mathbb{P}(Y_3 = 1) \times \mathbb{P}(Y_4 = 0) \times \mathbb{P}(Y_5 = 0) \implies \quad (13)$$

$$l = \log L = \sum_{i=1,3} \log(F(\beta_1 + x_i\beta_2)) + \sum_{i=2,4,5} \log(1 - F(\beta_1 + x_i\beta_2)) \quad (14)$$

Найдем теперь условия первого порядка

$$\frac{\partial L}{\partial \beta_1} = \sum_{i=1,3} \frac{f(\beta_1 + x_i\beta_2)}{F(\beta_1 + x_i\beta_2)} - \sum_{i=2,4,5} \frac{f(\beta_1 + x_i\beta_2)}{1 - F(\beta_1 + x_i\beta_2)} = 0, \quad (15)$$

$$\frac{\partial L}{\partial \beta_2} = \sum_{i=1,3} \frac{x_i f(\beta_1 + x_i\beta_2)}{F(\beta_1 + x_i\beta_2)} - \sum_{i=2,4,5} \frac{x_i f(\beta_1 + x_i\beta_2)}{1 - F(\beta_1 + x_i\beta_2)} = 0. \quad (16)$$

2. Пример.

Пусть  $x_1 = 12, x_2 = 23, x_3 = 1223$  и  $y_1 = y_2 = 0, y_3 = 1$ .

3. Краткий ответ. Потому что  $y^*$  сравнивается с 0.

Немного более развернутый ответ. Пробит модель в общем случае выглядит такие образом:

$$y_i^* = \beta_1 + \beta_2 x_i + \epsilon_i, \quad (17)$$

где  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , а  $y = 1$ , когда  $y^* \geq c$  и  $y = 0$  иначе. Заметим, что замена  $\tilde{y}^* = \sigma^{-1}(y^* - c)$  приводит к случаю, когда  $\tilde{y}^* \sim \mathcal{N}(0, 1)$ , следовательно всегда можно перейти к  $\epsilon_i \sim \mathcal{N}(0, 1)$  и сравнивать  $y^*$  с нулем.

4. Функция максимального правдоподобия

$$L = \prod_{i=1}^{50} \frac{1}{\sqrt{2\pi\sigma_M^2}} \exp\left\{-\frac{(y_i - \beta x_i)^2}{2\sigma_M^2}\right\} \times \prod_{i=51}^{100} \frac{1}{\sqrt{2\pi\sigma_J^2}} \exp\left\{-\frac{(y_i - \beta x_i)^2}{2\sigma_J^2}\right\} \implies \quad (18)$$

$$L = \frac{1}{(2\pi \cdot \sigma_M \cdot \sigma_J)^{50}} \exp\left\{-\sum_{i=1}^{50} \frac{(y_i - \beta x_i)^2}{2\sigma_M^2} - \sum_{i=51}^{100} \frac{(y_i - \beta x_i)^2}{2\sigma_J^2}\right\} \quad (19)$$

$$l = \log L = -50 \log(2\pi \cdot \sigma_M \cdot \sigma_J) - \sum_{i=1}^{50} \frac{(y_i - \beta x_i)^2}{2\sigma_M^2} - \sum_{i=51}^{100} \frac{(y_i - \beta x_i)^2}{2\sigma_J^2} \quad (20)$$

Запишем теперь условия первого порядка

$$\frac{\partial L}{\partial \sigma_M^2} = -\frac{25}{\sigma_M^2} + \frac{\sum_{i=1}^{50} (y_i - \beta x_i)^2}{2(\sigma_M^2)^2} = 0, \quad (21)$$

$$\frac{\partial L}{\partial \sigma_J^2} = -\frac{25}{\sigma_J^2} + \frac{\sum_{i=51}^{100} (y_i - \beta x_i)^2}{2(\sigma_J^2)^2} = 0, \quad (22)$$

$$\frac{\partial L}{\partial \beta} = \frac{\sum_{i=1}^{50} 2(y_i - \beta x_i)x_i}{2\sigma_M^2} + \frac{\sum_{i=51}^{100} 2(y_i - \beta x_i)x_i}{2\sigma_J^2} = 0. \quad (23)$$

5. (a)

$$H_0 : \hat{\beta}_2 = 0 \quad (24)$$

$$H_a : \hat{\beta}_2 \neq 0 \quad (25)$$

$t_{\text{obs}} = 3/\sqrt{0.09} = 10 \implies$ , гипотеза о незначимости коэффициента отвергается, т.е. коэффициент  $\beta_2$  значим.

(b)

$$\frac{\partial \mathbb{P}(y_i = 1)}{\partial x_i} \Big|_{x_i = -0.5} = \frac{e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}}{(1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i})^2} \hat{\beta}_2 = \frac{e^{0.8}}{(1 + e^{0.8})^2} \cdot 3 = 0.64 \quad (26)$$

(c)

$$\frac{\partial \mathbb{P}(y_i = 1)}{\partial x_i} = \frac{e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}}{(1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i})^2} \hat{\beta}_2 \rightarrow \max_{x_i} \implies \quad (27)$$

$$\frac{\hat{\beta}_2 e^{\hat{\beta}_1 + \hat{\beta}_2 x_i} (1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i})^2 - 2 \hat{\beta}_2 (1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}) (e^{\hat{\beta}_1 + \hat{\beta}_2 x_i})^2}{(1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i})^4} \hat{\beta}_2 = 0, \quad (28)$$

$$1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i} - 2e^{\hat{\beta}_1 + \hat{\beta}_2 x_i} = 0 \implies x_i = -\frac{0.7}{3}. \quad (29)$$

(d)

$$\mathbb{P}(y_i = 1) = \Lambda(\hat{\beta}_1 + \hat{\beta}_2 x_i) = \frac{1}{1 + \exp\{-0.7 + 1.5\}} = 0.3100255 \quad (30)$$

(e) Здесь нужно использовать Дельта-метод. Линеаризуем функцию  $F(\hat{\beta}) = \hat{\mathbb{P}}(y_i = 1)$  в окрестности точки  $\beta = (\beta_1, \beta_2)'$ :

$$F(\hat{\beta}) \approx F(\beta) + \nabla_F(\beta)' \cdot (\hat{\beta} - \beta)$$

Где все векторы — столбцы. Тогда оценка дисперсии прогноза:

$$\widehat{Var}(\hat{P}(y_i = 1)) = \nabla_F(\hat{\beta})' \cdot \widehat{Var}(\hat{\beta}) \cdot \nabla_F(\hat{\beta})$$

Градиент функции  $F$ :

$$\nabla_F(\hat{\beta}) = \left( \frac{e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}}{1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}} \quad x_i \cdot \frac{e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}}{1 + e^{\hat{\beta}_1 + \hat{\beta}_2 x_i}} \right)' = (0.69 \quad -0.34)'$$

Тогда стандартная ошибка прогноза:

$$s.e(\hat{P}(y_i = 1)) = (0.69^2 \cdot 0.04 - 2 \cdot 0.01 \cdot 0.69 \cdot 0.34 + 0.34^2 \cdot 0.09)^{1/2} = 0.18$$

6. Используя формулы (21) и (22) получаем

$$25 = \frac{\sum_{i=1}^{50} y_i^2}{2\sigma_M^2} \implies \hat{\sigma}_M^2 = 2, \quad (31)$$

$$25 = \frac{\sum_{i=51}^{100} y_i^2}{2\sigma_J^2} \implies \hat{\sigma}_J^2 = 6. \quad (32)$$



Проверим следующую гипотезу

$$H_0 : \sigma_M^2 = \sigma_J^2 \quad (33)$$

$$H_a : \sigma_M^2 \neq \sigma_J^2 \quad (34)$$

**Тест отношения правдоподобия**

$$LR = 2(\hat{l}_{UR} - \hat{l}_R) \sim \chi_1^2, \quad (35)$$

где  $\hat{l}_{UR}$  и  $\hat{l}_R$  значение логарифмических функций правдоподобия в точках максимума для неограниченной и ограниченной моделей соответственно. Решим ограниченную модель. Функция правдоподобия:

$$L(\sigma^2) = \prod_{i=1}^{100} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{y_i^2}{2\sigma^2} \right\} \rightarrow \max_{\sigma^2}$$

Логарифмическая функция правдоподобия:

$$l(\sigma^2) = -50 \ln(2\pi) - 50 \ln(\sigma^2) - \sum_{i=1}^{100} \frac{y_i^2}{2\sigma^2} \rightarrow \max_{\sigma^2}$$

Откуда:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{100} y_i^2}{100} = 4$$

Тогда:

$$LR = 2(-25 \ln(12) - 50 + 50 \ln(4) + 50) = 14.38 > \chi_{crit}^2 = 3.84$$

То есть нулевая гипотеза отвергается!

**Тест Вальда**

$$W = g(\hat{\theta}_{UR})' (\nabla_g \hat{I}(\hat{\theta}_{UR})^{-1} \nabla_g')^{-1} g(\hat{\theta}_{UR})$$

где градиент  $\nabla_g = (1 \ -1)$  – вектор строки. Нам пригодится информационная матрица Фишера в задачи без ограничений. Найти ее можно взяв вторые производные от функции правдоподобия в четвертой задаче по  $\sigma_M^2$  и  $\sigma_J^2$ , положив  $\beta = 0$  и домножив на -1.

$$\hat{I}(\hat{\theta}_{UR}) = \begin{pmatrix} -\frac{25\hat{\sigma}_M^2 - 100}{(\hat{\sigma}_M^2)^3} & 0 \\ 0 & -\frac{25\hat{\sigma}_J^2 - 300}{(\hat{\sigma}_J^2)^3} \end{pmatrix} = \begin{pmatrix} \frac{50}{8} & 0 \\ 0 & \frac{150}{216} \end{pmatrix}$$

Тогда:

$$W = (6 - 2)^2 \left( (1 \ -1) \begin{pmatrix} \frac{150}{216} & 0 \\ 0 & \frac{50}{8} \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right)^{-1} = 10$$

Как и ранее, нулевая гипотеза отвергается!

**Тест множителей Лагранжа**

Нужно посчитать статистику:

$$LM = \left( \frac{\partial l(\hat{\theta}_R)}{\partial \theta} \right)' \hat{I}(\hat{\theta}_R)^{-1} \left( \frac{\partial l(\hat{\theta}_R)}{\partial \theta} \right) \sim \chi_r^2,$$

При верной  $H_0$  логарифмическая функция правдоподобия будет выглядеть следующим образом:

$$l(\sigma^2) = -50\ln(2\pi) - 50\ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{100} y_i^2$$

Тогда:

$$\frac{\partial l(\hat{\sigma}^2)}{\partial \sigma^2} = -\frac{50}{\hat{\sigma}^2} + \frac{1}{2} \sum_{i=1}^{100} y_i^2 \frac{1}{(\hat{\sigma}^2)^2} = 0$$

Откуда:  $\hat{\sigma}_R^2 = 4$ . Заметим, что если подставить эту оценку в матрицу  $\hat{I}(\theta)$ , она не будет иметь обратной, что эквивалентно бесконечной  $LM$  статистике. Гипотеза  $H_0$ , конечно, отвергается!

## 7.9. Большой Устный ЗАчёт 2016

1. Метод Наименьших Квадратов.
  - а) МНК-картинка
  - б) Нахождение всего-всего, если известен вектор  $y$  и матрица  $X$
2. Теорема Гаусса-Маркова
  - а) Формулировка с детерминистическими регрессорами
  - б) Доказательство с детерминистическими регрессорами
  - в) Формулировки со стохастическими регрессорами
  - г) Что даёт дополнительное предположение о нормальности  $\varepsilon$ ?
  - д) Теорема Фриша-Вау
3. Проверка гипотез о линейных ограничениях
  - а) Проверка гипотезы о значимости коэффициента
  - б) Проверка гипотезы о значимости регрессии в целом
  - в) Проверка гипотезы об одном линейном соотношении с помощью ковариационной матрицы
  - г) Ограниченная и неограниченная модель
  - д) Тест Чоу на стабильность коэффициентов
  - е) Тест Чоу на прогнозную силу
4. Метод максимального правдоподобия
  - а) Свойства оценок
  - б) Два способа получения оценки дисперсии
  - в) Три теста (LM, Wald, LR)
  - г) Выписать функцию ML для обычной регрессии
  - д) для AR(1) процесса
  - е) для MA(1) процесса
  - ж) для логит модели
  - з) для пробит модели
  - и) для модели с заданным видом гетероскедастичности
5. Мультиколлинеарность
  - а) Определение, последствия

- б) Величины, измеряющие силу мультиколлинеарности
  - в) Методы борьбы
  - г) Сюда же: метод главных компонент, хотя он используется и для других целей
6. Гетероскедастичность
- а) Определение, последствия
  - б) Тесты, график
  - в) Стьюдентизированные остатки
  - г) НС оценки ковариации
  - д) GLS и FGLS
7. Временные ряды
- а) Стационарный временной ряд
  - б) ACF, PACF
  - в) Модель ARMA
  - г) ARIMA-SARIMA
  - д) Модель GARCH (не будет, не успели)
8. Логит и пробит
- а) Описание моделей
  - б) Предельные эффекты
  - в) Чувствительность, специфичность (не будет, не успели)
  - г) Кривая ROC (не будет, не успели)
9. Эндогенность
- а) Три примера: одновременность, пропущенные переменные, ошибки измерения
  - б) IV, двухшаговый МНК
10. Модели панельных данных (не будет, не успели)
- а) RE, FE, сквозная регрессии
  - б) Тест Хаусмана
11. И ещё алгоритмы. Уметь объяснить суть метода. Уметь реализовать его в R.
- а) Метод опорных векторов
  - б) Классификационные деревья и случайный лес
  - в) Ridge regression
  - г) LASSO
  - д) Квантильная регрессия
  - е) Байесовский подход к регрессии (не будет, не успели)
12. R. Можно принести файл со своей заготовкой, можно пользоваться Интернетом для поиска информации, но не для общения.
- а) Загрузить данные из .csv файла в R
  - б) Посчитать описательные статистики (среднее, мода, медиана и т.д.)
  - в) Построить подходящие описательные графики для переменных
  - г) Оценить линейную регрессию с помощью МНК. Провести диагностику на что-нибудь (гетероскедастичность, автокорреляцию, мультиколлинеарность).
  - д) Оценить logit, probit модели, посчитать предельные эффекты
  - е) Оценить ARMA модель
  - ж) Выделить главные компоненты

## 7.10. Экзамен 20.06.2016. Вариант 1

1. На основании опроса была оценена следующая модель:

$$\ln(wage_i) = \beta_1 + \beta_2 exper_i + \beta_3 exper_i^2 + \beta_4 married_i + \beta_5 educ_i + \beta_6 black_i + \varepsilon_i$$

где:

- $wage_i$  — величина заработной платы в долларах
- $exper_i$  — опыт работы в годах
- $educ_i$  — количество лет обучения
- $married_i$  — наличие супруга/супруги (1 — есть, 0 — нет)
- $black_i$  — принадлежность к негроидной расе (1 — да, 0 — нет)

Показатель	Значение
$R^2$	<b>B7</b>
Скорректированный $R^2$	0.219
Стандартная ошибка регрессии	<b>B6</b>
Количество наблюдений	<b>B2</b>

Результаты дисперсионного анализа:

	df	SS	MS	F	P-значение
Регрессия	<b>B1</b>	5.993	1.199	<b>B5</b>	0.000
Остаток	134	18.240	0.136		
Итого	<b>B3</b>	<b>B4</b>			

Коэффициент	Оценка	$se(\hat{\beta})$	t-статистика	P-Значение	Нижние 95%	Верхние 95%
Константа	4.529	0.331	13.688	0.000	3.874	5.183
$exper$	0.090	0.037	2.419	0.017	0.016	0.164
$exper^2$	-0.003	0.002	-1.790	0.076	-0.006	0.000
$married$	0.240	0.079	3.045	0.003	<b>B8</b>	<b>B9</b>
$educ$	0.078	0.017	<b>B10</b>	0.000	0.045	0.111
$black$	0.073	0.171	0.424	0.672	-0.266	0.411

Найдите пропущенные числа **B1--B10**.

Ответ округляйте до 3-х знаков после запятой. Кратко поясняйте, например, формулой, как были получены результаты.

2. На основании данных по ценам на квартиры в Москве были построена модель

$$\ln(price_i) = \beta_1 + \beta_2 totsp_i + \beta_3 metrdist_i + \beta_4 dist_i + \beta_5 floor_i + \varepsilon_i,$$

где:

- $\ln(price_i)$  — логарифм цены квартиры в тысячах долларов
- $totsp_i$  — общая площадь квартиры в кв. м.
- $metrdist_i$  — расстояние до метро в минутах
- $dist_i$  — расстояние до центра города в км
- $floor_i$  — дамми-переменная (1 — если квартира не на первом и последнем этажах, 0 — иначе)

Модели были оценены на пяти разных выборках, результаты представлены в таблице:

Коэффициент	Выборка А	Выборка В	Выборка С	Выборка D	Выборка Е
Константа	3.980***	3.926***	3.929***	3.719***	4.224***
<i>totsp</i>	0.0155***	0.0148***	0.0163***	0.0179***	0.0139***
<i>metrdist</i>	-0.00858***	-0.0169***	-0.00566**	-0.0108***	-0.0077
<i>dist</i>	-0.0267***	-0.0186***	-0.0253***	-0.0150***	-0.0350***
<i>floor</i>	0.0419**	0.0633*	0.0224	0.0225	0.0228
Наблюдений	460	145	315	150	150
$R^2$	0.693	0.684	0.723	0.328	0.520
$RSS$	15.120	4.503	9.408	2.163	8.545

\* — значимость на 10%, \*\* — значимость на 5%, \*\*\* — значимость на 1%.

- Для всей выборки (выборка А) проинтерпретируйте коэффициент при переменной  $dist_i$ .
- Определите на 5%-ом уровне значимости, можно ли использовать одну модель для квартир, находящихся в пешей доступности от метро (выборка С), и квартир, находящихся в транспортной доступности (выборка В).
- Исследователь предположил, что дисперсия ошибок модели возрастает с увеличением площади квартиры. Проверьте, есть ли в модели гетероскедастичность на 10% уровне значимости на основании соответствующего теста. В выборку D включены 150 квартир с наименьшей общей площадью, в выборку Е — 150 квартир с наибольшей общей площадью.

При проверке гипотез: выпишите  $H_0$ ,  $H_a$ , найдите значение тестовой статистики, укажите её распределение, найдите критическое значение, сделайте выводы

- По ежемесячным данным, 146 наблюдений, была оценена зависимость

$$\widehat{credit}_t = 362.21 - 7.50r\_credit_t - 13.09ipc_t, R^2 = 0.44$$

где:

- $credit_t$  — объём потребительских кредитов, выданных домашним хозяйствам РФ
- $r\_credit_t$  — ставка процента по кредитам
- $ipc_t$  — индекс потребительских цен

Известно, что  $\sum_{t=2}^{146} (\hat{\varepsilon}_t - \hat{\varepsilon}_{t-1})^2 = 266491$ ,  $\sum_{t=1}^{146} \hat{\varepsilon}_t^2 = 438952$ ,  $\sum_{t=2}^{146} |\hat{\varepsilon}_t - \hat{\varepsilon}_{t-1}| = 3617$ ,  $\sum_{t=1}^{146} |\hat{\varepsilon}_t| = 6382$ .

Кроме того, была оценена вспомогательная модель для остатков исходной модели:

$$\hat{\varepsilon}_t = 15.67 - 0.75r\_credit_t + 0.02ipc_t + 0.39\hat{\varepsilon}_{t-1} + 0.21\hat{\varepsilon}_{t-2} + 0.24\hat{\varepsilon}_{t-3}, R^2 = 0.56$$

- На 1%-ом уровне значимости проверьте гипотезу об адекватности исходной регрессии
- Проведите тест Дарбина-Уотсона на 5%-ом уровне значимости
- Проведите тест Бройша-Годфри на 5%-ом уровне значимости

При проверке гипотез: выпишите  $H_0$ ,  $H_a$ , найдите значение тестовой статистики, укажите её распределение, найдите критическое значение, сделайте выводы

- Домохозяйка Глаша очень любит читать романы Л.Н. Толстого и смотреть сериалы. Её сын Петя учится на третьем курсе ВШЭ. Последние 30 дней он записывал, сколько Глаша прочитала страниц «Анны Карениной»,  $pages_t$ , и посмотрела серий «Доктора Хауса»,  $series_t$ . На основании этих наблюдений при помощи МНК Петя оценил следующую модель:

$$\widehat{pages}_t = 100 - 3series_t$$

Оценка ковариационной матрицы коэффициентов,  $\widehat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} 11 & 0.5 \\ 0.5 & 1 \end{pmatrix}$

Оценка дисперсии ошибок равна  $\hat{s}^2 = 323$ .

Завтра Глаша собирается посмотреть 10 серий «Доктора Хауса».

- а) Постройте точечный прогноз количества прочитанных Глашей страниц романа
  - б) Постройте 95%-ый доверительный интервал для  $\mathbb{E}(pages_t | series_t = 10)$ , ожидаемого количества прочитанных страниц
  - в) Постройте 95%-ый предиктивный интервал для фактического количества прочитанных страниц
5. Опишите МНК для парной регрессии: выпишите целевую функцию, систему нормальных уравнений, оценки коэффициентов, оценки дисперсий коэффициентов.
  6. Сформулируйте теорему Гаусса-Маркова для детерминированных регрессоров.
  7. Опишите тест Уайта: сформулируйте нулевую и альтернативную гипотезы, способ получения тестовой статистики, её распределение при верной нулевой гипотезе, вид критической области.

## 7.11. Экзамен 20.06.2016. Вариант 2

1. На основании опроса была оценена следующая модель:

$$\ln(wage_i) = \beta_1 + \beta_2 exper_i + \beta_3 exper_i^2 + \beta_4 married_i + \beta_5 educ_i + \beta_6 black_i + \varepsilon_i$$

где:

- $wage_i$  — величина заработной платы в долларах
- $exper_i$  — опыт работы в годах
- $educ_i$  — количество лет обучения
- $married_i$  — наличие супруга/супруги (1 — есть, 0 — нет)
- $black_i$  — принадлежность к негроидной расе (1 — да, 0 — нет)

Показатель	Значение				
$R^2$	<b>B6</b>				
Скорректированный $R^2$	0.164				
Стандартная ошибка регрессии	<b>B7</b>				
Количество наблюдений	<b>B1</b>				
Результаты дисперсионного анализа:					
	df	SS	MS	F	P-значение
Регрессия	<b>B2</b>	<b>B4</b>	7.425	<b>B5</b>	0.000
Остаток	<b>B3</b>	184.954	0.145		
Итого	1279	222.079			

Коэффициент	Оценка	$se(\hat{\beta})$	t-статистика	P-Значение	Нижние 95%	Верхние 95%
Константа	4.906	0.106	46.129	0.000	4.698	5.115
<i>exper</i>	0.095	0.011	8.956	0.000	0.074	0.115
<i>exper</i> <sup>2</sup>	-0.003	0.001	-5.437	0.000	-0.004	-0.002
<i>married</i>	<b>B8</b>	<b>B9</b>	<b>B10</b>	0.234	-0.018	0.074
<i>educ</i>	0.064	0.006	11.582	0.000	0.053	0.075
<i>black</i>	-0.183	0.028	-6.490	0.000	-0.238	-0.127

Найдите пропущенные числа **B1--B10**.

Ответ округляйте до 3-х знаков после запятой. Кратко поясняйте, например, формулой, как были получены результаты.

2. Винни-Пух и Пятачок попробовали очень странный мёд. После его употребления, к ним пришли слоники в количестве 100 штук и начали водить вокруг них хороводы. Винни-Пух смог на глазок оценить вес и рост каждого слоника, а Пятачок — его возраст. Эти данные позволили им оценить следующую модель:

$$weight_i = \beta_1 + \beta_2 \ln(height_i) + \beta_3 \ln(age_i) + \varepsilon_i$$

где:

- $weight_i$  — вес слоника
- $\ln(height_i)$  — логарифм роста слоника
- $\ln(age_i)$  — логарифм возраста слоника

Выборка	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	ESS	RSS	N
1. Самые молодые слоники	43.1	1.4	3.7	243	345	40
2. Самые старые слоники	48.4	3.6	1.1	489	194	40
3. Зелёные слоники	39.6	2.1	2.4	311	268	45
4. Розовые слоники	53.1	2.9	3.1	369	307	55
5. Все слоники	45.7	2.6	2.9	615	741	100

- Для выборки розовых слоников проинтерпретируйте коэффициент  $\hat{\beta}_2$
- Определите на 5%-ом уровне значимости, можно ли использовать одну модель для розовых и зелёных слоников
- Пятачок уверен, что дисперсия ошибок модели падает с увеличением возраста слоника. Проверьте, так ли это, на 1% уровне значимости на основании соответствующего теста

При проверке гипотез: выпишите  $H_0$ ,  $H_a$ , найдите значение тестовой статистики, укажите её распределение, найдите критическое значение, сделайте выводы

3. Царевна Несмеяна по 146 дням наблюдений построила регрессию:

$$\widehat{tear}_t = 200 - 5sun_t - 10prince_t - 15chocolate_t, R^2 = 0.8$$

где:

- $tear_t$  — количество пролитых слёз в мл
- $sun_t$  — дамми-переменная для погоды (1 — солнечная, 0 — пасмурная)
- $prince_t$  — дамми-переменная для посещений Прекрасным Принцем (1 — Принц пришёл, 0 — нет)
- $chocolate_t$  — количество съеденного шоколада в плитках

Известно, что  $\sum_{t=2}^{146} (\hat{\varepsilon}_t - \hat{\varepsilon}_{t-1})^2 = 876$ ,  $\sum_{t=1}^{146} \hat{\varepsilon}_t^2 = 538$ ,  $\sum_{t=2}^{146} |\hat{\varepsilon}_t - \hat{\varepsilon}_{t-1}| = 100$ ,  $\sum_{t=1}^{146} |\hat{\varepsilon}_t| = 150$ .

- На 1%-ом уровне значимости проверьте гипотезу об адекватности регрессии
- Проведите тест Дарбина-Уотсона на 5% уровне значимости
- Кроме того, была оценена следующая модель:

$$\hat{\varepsilon}_t = 5 - 0.05sun_t + 0.02prince_t + 0.09\hat{\varepsilon}_{t-1} + 0.01\hat{\varepsilon}_{t-2} + 0.004\hat{\varepsilon}_{t-3}, R^2 = 0.06$$

Проведите тест Бройша-Годфри на 1% уровне значимости

При проверке гипотез: выпишите  $H_0$ ,  $H_a$ , найдите значение тестовой статистики, укажите её распределение, найдите критическое значение, сделайте выводы

- Ослик Иа-Иа горюет и считает количество мёда в горшочках. Сейчас у него 50 горшочков. Горшочки отличаются друг от друга цветом: есть более розовые и менее розовые. Иа-Иа считает, что степень розовости влияет на количество мёда. Он смог оценить следующую регрессию:

$$\widehat{honey}_i = 15 + 3pinkness_i$$

Оценка ковариационной матрицы коэффициентов,  $\widehat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} 3 & 1.5 \\ 1.5 & 9 \end{pmatrix}$

Оценка дисперсии ошибок равна  $\hat{\sigma}^2 = 137$ .

Ослик нашёл новый горшочек с розовостью, равной 5.

- Постройте точечный прогноз для количества мёда
  - Постройте 95%-ый доверительный интервал для  $\mathbb{E}(honey_i | pinkness_i = 5)$ , ожидаемого количества мёда в горшочке
  - Постройте 95%-ый предиктивный интервал для фактического количества мёда в горшочке
- Опишите  $F$ -тест для гипотезы о нескольких линейных ограничениях: сформулируйте нулевую и альтернативную гипотезы, способ получения тестовой статистики, её распределение при верной нулевой гипотезе, вид критической области.
  - В парной регрессии  $y_t = \beta_1 + \beta_2 x_t + \varepsilon_t$  известно, что  $\text{Var}(\varepsilon_t) = \sigma^2 x_t^4$ . Опишите процедуру получения эффективных оценок коэффициентов.
  - Опишите двухшаговый МНК: сформулируйте условия для его применения и опишите процедуру построения оценок.

## 8. 2016-2017

### 8.1. ИП. Подготовка к празднику

Линейная алгебра:

- Умножать, складывать, вычитать матрицы
- Считать определитель и след матрицы
- Находить обратную матрицу



4. Транспонировать
5. Находить собственные числа и собственные векторы
6. Находить  $A^{2012}$ , если  $A$  - диагонализируема
7. Знать свойства:
  - а)  $(A')^{-1} = (A^{-1})'$
  - б) Если  $A$  и  $B$  квадратные, то  $(AB)^{-1} = B^{-1}A^{-1}$
  - в)  $\det(A) = 0$  равносильно тому, что у матрицы есть линейно зависимые столбцы (или строки)
  - г) Если  $A$  и  $B$  квадратные, то  $\det(AB) = \det(A) \det(B)$
  - д) Если  $AB$  и  $BA$  существуют, то  $\text{trace}(AB) = \text{trace}(BA)$

Школьная программа:

1. Теорема Пифагора, формулировка и доказательство
2. Теорема о трёх перпендикулярах, формулировка
3. Скалярное произведение двух векторов, формула
4. Косинус угла между векторами через скалярное произведение

Теория вероятностей:

1. Уметь находить вероятность и условную вероятность в простейших случаях
2. Знать свойства  $\mathbb{E}(X)$ ,  $\text{Cov}(X, Y)$ ,  $\text{Var}(Y)$

Статистика:

1. Уметь проверять гипотезу о математическом ожидании (асимптотический нормальный случай, t-тест) и строить доверительный интервал
2. Знать смысл слов: ошибка I, II рода,  $P$ -значение

## 8.2. ИП. Праздник «Вспомнить всё!» 12.09.2016

Сегодня 256-ой день года, всех с днём программиста! :) А ещё 12 сентября в 490 году до нашей эры Фидиппид добежал из Марафона в Афины с криком «Νενικήκαμεν»<sup>2</sup>!

1. Найдите длины векторов  $a = (1, 1, 1)$  и  $b = (1, 2, 3)$  и косинус угла между ними. Найдите один любой вектор, перпендикулярный вектору  $b$ .
2. Сформулируйте теорему о трёх перпендикулярах и обратную к ней
3. На плоскости  $\alpha$  лежит прямая  $\ell$ . Вне плоскости  $\alpha$  лежит точка  $C$ . Ромео проецирует точку  $C$  на прямую  $\ell$  и получает точку  $R$ . Джульетта проецирует точку  $C$  сначала на плоскость  $\alpha$ , а затем проецирует полученную точку  $A$  на прямую  $\ell$ . После двух действий Джульетта получает точку  $D$ . Обязательно ли  $R$  и  $D$  совпадают?
4. Для матрицы

$$A = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}$$

- а) Найдите собственные числа и собственные векторы матрицы

---

<sup>2</sup>Ликуйте! Мы победили!

- б) Найдите  $\det(A)$ ,  $\text{tr}(A)$
- в) Найдите собственные числа матрицы  $A^{2016}$ ,  $\det(A^{2016})$  и  $\text{tr}(A^{2016})$
5. Известно, что  $X$  — матрица размера  $n \times k$  и  $n > k$ , известно, что  $X'X$  обратима. Рассмотрим матрицу  $H = X(X'X)^{-1}X'$ . Укажите размер матрицы  $H$ , найдите  $H^{2016}$ ,  $\text{tr}(H)$ ,  $\det(H)$ , собственные числа матрицы  $H$ . Штрих означает транспонирование.
6. Занудная халява: известно, что  $\text{Cov}(X, Y) = 5$ ,  $\text{Var}(X) = 10$ ,  $\text{Var}(Y) = 20$ ,  $\mathbb{E}(X) = 10$ ,  $\mathbb{E}(Y) = -10$ . Найдите  $\text{Cov}(X + 2Y, Y - X)$ ,  $\text{Var}(X + 2Y)$ ,  $\mathbb{E}(X + 2Y)$ .
7. За 100 дней Ромео посчитал все глубокие вздохи Джульетты. Настроение Джульетты столь спонтанно, что глубокие вздохи за разные дни можно считать независимыми. В сумме оказалось 890 вздохов. Сумма квадратов оказалась равна 8000. Постройте 95%-ый доверительный интервал для математического ожидания ежедневного количества глубоких вздохов Джульетты. На уровне значимости 5%-ов проверьте гипотезу, что математическое ожидание равно 9.
8. Ромео подкидывает монетку два раза. Если монетка выпадает орлом, то Ромео кладет в мешок черный шар, если решкой — белый. Джульетта не знает, как выпадала монетка, и достает шары из мешка наугад по очереди. Первый шар оказался черного цвета. Какова вероятность того, что второй шар Джульетты будет белым?

### 8.3. ИП. Праздник «Вспомнить всё!» 12.09.2016, ответы

1.  $|a| = \sqrt{3}$ ,  $|b| = \sqrt{14}$ ,  $\cos(a, b) = \frac{6}{\sqrt{42}}$ ,  $b^\perp = (-3 \ 0 \ 1)'$
2. Прямая: если прямая проходит через основание плоскости и перпендикулярна её проекции, то она перпендикулярна и самой наклонной.  
Обратная: если прямая, проведённая на плоскости через основание наклонной, перпендикулярна самой наклонной, то она перпендикулярна и самой прямой.
- 3.
4. а)  $\lambda_1 = 1$ ,  $\lambda_2 = 9$ ,  $h_1 = (1 \ -1)'$ ,  $h_2 = (1 \ 1)'$   
б)  $\det(A) = 9$ ,  $\text{tr}(A) = 10$   
в)  $\lambda_1 = 1$ ,  $\lambda_2 = 9^{2016}$ ,  $\det(A^{2016}) = 9^{2016}$ ,  $\text{tr}(A^{2016}) = 1 + 9^{2016}$
5. Размер:  $n \times n$ ,  $H^{2016} = H$ ,  $\text{tr}(H) = k$ ,  $\det(H) = 0$ ,  $\lambda_1 = \dots = \lambda_k = 1$ ,  $\lambda_{k+1} = \dots = \lambda_n = 0$ .
6.  $\text{Cov}(X + 2Y, Y - X) = 25$ ,  $\text{Var}(X + 2Y) = 110$ ,  $\mathbb{E}(X + 2Y) = -10$ .
7.  $7.15 < \mu < 10.65$ , число 9 входит в доверительный интервал, значит, оснований отвергать основную гипотезу нет.

### 8.4. Кр 1, демо

1. Эконометресса Ефросинья исследует, как зависит надой молока,  $\text{milk}_i$ , (в литрах) от возраста коровы,  $\text{age}_i$ , (в годах):

$$\text{milk}_i = \beta_1 + \beta_2 \text{age}_i + u_i$$

Показатель	Значение
$RSS$	<b>B1</b>
$ESS$	<b>B2</b>
$TSS$	<b>1240</b>
$R^2$	<b>B3</b>
Стандартная ошибка регрессии	<b>1.45</b>
Количество наблюдений	<b>340</b>

Коэффициент	Оценка	$se(\hat{\beta})$	t-статистика	P-значение	Левая (95%)	Правая (95%)
Константа	4.565	0.207	<b>B4</b>	<b>B9</b>	<b>B5</b>	<b>B6</b>
$age$	<b>B7</b>	<b>B8</b>	3.670	0.000	0.036	0.119

Найдите пропущенные числа **B1--B9**.

Ответ округляйте до 2-х знаков после запятой. Кратко поясняйте формулой, как были получены результаты.

2. Гарри Поттер и Рон Уизли активно готовятся к чемпионату мира по квиддичу. В течение 30 дней они сначала посещают Хогсмид и выпивают некоторое количество сливочного пива в пинтах,  $beer_t$ , после забивают определённое количество квоффлов в штуках,  $quaffle_t$ . Гермиона Грейнджер оценила следующую регрессию:

$$\widehat{quaffle}_t = \underset{(2.83)}{80} - \underset{(1)}{3} beer_t$$

В скобках приведены стандартные ошибки. Оценка дисперсии ошибок равна  $\hat{s}^2 = 238$ .

Сегодня Гарри и Рон выпили 4 пинты сливочного пива.

- Проверьте гипотезы о значимости каждого коэффициента на уровне значимости 5%.
  - Постройте точечный прогноз количества квоффлов, забитых сегодня Гарри Поттером и Роном Уизли
  - Постройте 90%-ый доверительный интервал для коэффициента наклона регрессии
3. Для модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$  выполнены все предпосылки теоремы Гаусса-Маркова.
- Докажите, что МНК-оценка коэффициента  $\beta_2$  является случайной величиной
  - Докажите, что эта оценка является несмещённой
  - Найдите дисперсию этой оценки
4. Для модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$  выполнены все предпосылки теоремы Гаусса-Маркова. Для МНК-оценок коэффициентов найдите  $\widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2)$ .
5. Дайте определения следующих понятий
- Несмещённая оценка
  - Эффективная оценка
  - Состоятельная последовательность оценок

## 8.5. Кр 1, демо, решения

1.

$$B1 = 1.45^2 \cdot (340 - 2)$$

$$B2 = ESS = TSS - RSS$$

$$B3 = R^2 = ESS/TSS$$

$$B4 = t_c = 4.565/0.207 = 22$$

По таблицам ( $t$ -распределение с 338 степенями свободы или примерно нормальное)  
 $t_{crit} = 1.96$

$$B5 = CI_{left} = 4.565 - 1.96 \cdot 0.207$$

$$B6 = CI_{right} = 4.565 + 1.96 \cdot 0.207$$

$$B7 = \hat{\beta}_{milk} = (0.036 + 0.119)/2 = 0.0775$$

$$B8 = se(\hat{\beta}_{milk}) = \hat{\beta}_{milk}/t_{milk} = 0.0775/3.670$$

$$B9 = P - value(22) \approx 0.000$$

2. а) Находим  $t$ -статистики:  $t_c = 80/2.83 = 28.3$ ,  $t_{beer} = -3/1 = -3$ . Если предположить нормальность ошибок, то  $t_{crit} = 2.05$ . Следовательно, в обоих случаях  $H_0$ :  $\beta = 0$  отвергается и оба коэффициента значимо отличны от нуля.

б)

$$\hat{Y}_i = 80 - 3 \cdot 4 = 80 - 12 = 68$$

- в) Для уровня доверия 90% получаем критическое значение  $t_{crit} = 1.7$ . Отсюда доверительный интервал равен

$$[-3 - 1 \cdot 1.7; -3 + 1 \cdot 1.7]$$

3. Решение изложено в лекциях

4. Решение изложено в лекциях

5. а) Несмещённая оценка

Оценка  $\hat{\theta}$  называется несмещённой, если  $\mathbb{E}(\hat{\theta}) = \theta$

- б) Эффективная оценка

Оценка  $\hat{\theta}$  называется эффективной среди множества оценок  $\Theta$ , если для любой оценки  $\hat{\theta}'$  из множества  $\Theta$  выполнено неравенство  $\text{Var}(\hat{\theta}) \leq \text{Var}(\hat{\theta}')$

- в) Состоятельная последовательность оценок

Последовательность оценок  $\hat{\theta}_1, \hat{\theta}_2, \dots$ , называется состоятельной, если

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| < \varepsilon) = 1$$

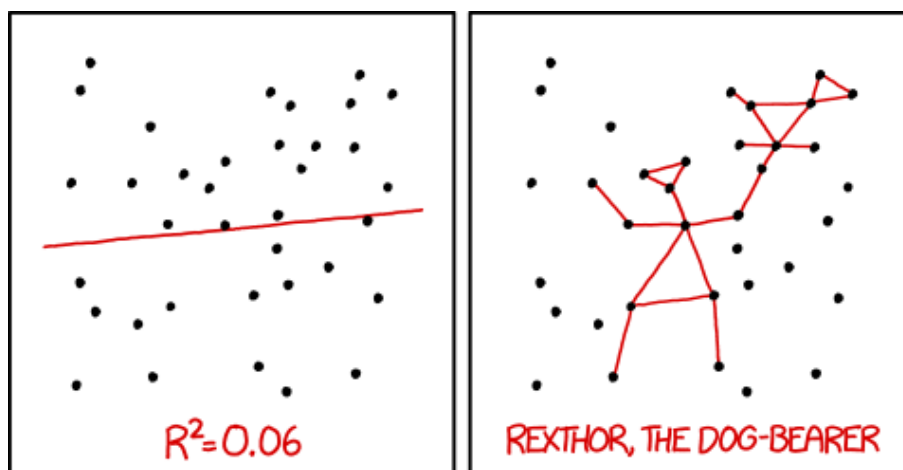
для любого числа  $\varepsilon > 0$ .

## 8.6. Кр 1, 24.10.2016

- В течение 10 дней Василий записывал количество пойманных им покемонов,  $Y_i$ , и количество решённых задач по эконометрике,  $X_i$ . Оказалось, что  $\sum X_i^2 = 44$ ,  $\sum Y_i^2 = 197$ ,  $\sum X_i = 15$ ,  $\sum Y_i = 15$  и  $\sum X_i Y_i = 44$ . Василий предполагает корректность линейной модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$ .
  - Найдите МНК-оценки коэффициентов регрессии
  - Найдите  $RSS$ ,  $ESS$ ,  $TSS$  и  $R^2$
- Для модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$  выполнены все предпосылки теоремы Гаусса-Маркова, а случайные ошибки нормально распределены. Известны все значения  $Y_i$ , все значения  $\hat{Y}_i$  и часть значений  $X_i$ :
 

$X_i$	4	5	.
$Y_i$	7	5	6
$\hat{Y}_i$	6.0	5.5	6.5

  - Найдите МНК-оценки коэффициентов регрессии
  - Найдите стандартную ошибку коэффициента  $\hat{\beta}_2$
  - Постройте 95%-ый доверительный интервал для коэффициента  $\hat{\beta}_2$
  - Проверьте гипотезу о незначимости коэффициента  $\beta_2$  на уровне значимости 5%
- Для модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$  выполнены все предпосылки теоремы Гаусса-Маркова. Докажите несмещённость МНК-оценки коэффициента  $\beta_1$ .
- Для модели  $Y_i = \beta_1 + \beta_2 X_i + u_i$  выполнены все предпосылки теоремы Гаусса-Маркова. Выведите формулу для дисперсии МНК-оценки,  $\text{Var}(\hat{\beta}_1)$ .
- Рассмотрим модель  $Y_i = \beta_1 + \beta_2 X_i + u_i$  с неслучайным регрессором. Аккуратно сформулируйте теорему Гаусса-Маркова, пояснив смысл используемых понятий.



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Randall Munroe, xkcd

Рис. 1: \*

## 8.7. ИП. Подготовка к битве

- Убедитесь, что в рамках классической модели  $y = X\beta + u$  и предпосылок  $\mathbb{E}(u) = 0$ ,  $\text{Var}(u) = \sigma^2 I$  умеете находить любые  $\mathbb{E}()$ ,  $\text{Var}()$  и  $\text{Cov}()$  для векторов  $\beta$ ,  $\hat{\beta}$ ,  $y$ ,  $\hat{y}$ ,  $u$ ,  $\hat{u}$ .

2. Вектор  $u$  размера  $4 \times 1$  имеет стандартное нормальное распределение,  $u \sim \mathcal{N}(0; I)$ . Известен вектор  $d$ ,  $d = (1, -1, 2, -2)$ .
  - а) Найдите такую матрицу  $H$ , что её умножение на произвольный вектор  $y$  означает проецирование вектора  $y$  на прямую порождённую вектором  $d$
  - б) Как распределена случайная величина  $u'Hu$ ? Чему равно её математическое ожидание и дисперсия?
3. Вектор  $u$  имеет стандартное нормальное распределение,  $u \sim \mathcal{N}(0; I)$ . Матрица  $A$  такова, что  $Au$  также имеет стандартное нормальное распределение,  $Au \sim \mathcal{N}(0; I)$ .
  - а) Выпишите уравнение, которому подчиняется матрица  $A$
  - б) Чему может равняться  $\det A$ ?
  - в) Рассмотрим  $c_1$  и  $c_2$  — первый и второй столбец матрицы  $A$ . Найдите  $c_1'c_1$  и  $c_1'c_2$
4. Предположим, что функция  $f(x) = x'Ax + Bx + c$  имеет минимум. Найдите его, используя технику дифференцирования по вектору
5. Рассмотрим систему уравнений  $X\beta = y$ . Здесь  $y$  — известный вектор размера  $n \times 1$ ,  $\beta$  — неизвестный вектор размера  $k \times 1$ , и  $X$  — известная матрица размера  $n \times k$  полного ранга. Мы хотим решить эту систему относительно  $\beta$ . Если  $n = k$ , то решать эту систему скучно, и, конечно,  $\beta = X^{-1}y$ . Гораздо интереснее решать систему, когда решений нет или когда их бесконечно много :)
  - а) Если решений нет, то найдите наилучшее приближение к решению, то есть такое  $\beta$  при котором длина  $(y - X\beta)$  минимальна.
  - б) Если решений,  $\beta$ , бесконечно много, то найдите решение с наименьшей длиной.

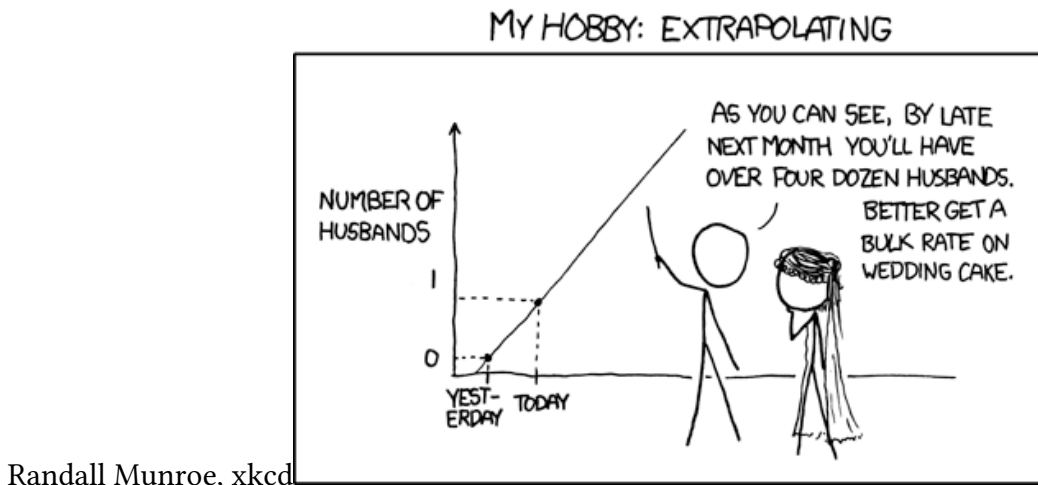


Рис. 2: \*

## 8.8. ИП. Битва под Малоярославцем, 24.10.2016/24.10.1812

1. Цель этой кампании — доказать теорему Гаусса-Маркова для множественной регрессии. Итак,  $y = X\beta + u$ , про случайные ошибки  $u$  известно, что  $\mathbb{E}(u) = 0$ ,  $\text{Var}(u) = \sigma^2 \cdot I$ . Поехали! :)

Тульский дворянин Дмитрий Сергеевич Дохтуров использует классическую МНК-оценку,  $\hat{\beta}_{OLS}$ :

- а) Вспомните формулу для МНК-оценки  $\hat{\beta}_{OLS}$

- б) Является ли оценка  $\hat{\beta}_{OLS}$  линейной по  $y$ ?
- в) Докажите, что оценка  $\hat{\beta}_{OLS}$  является несмещённой

Корсиканец Наполеон Бонапарт предлагает альтернативную несмещённую оценку  $\hat{\beta}_{alt} = A \cdot y$ :

- г) Является ли оценка  $\hat{\beta}_{alt}$  линейной по  $y$ ?
- д) Чему равняется матрица  $AX$ ? Да поможет здесь условие несмещённости!
- е) Найдите  $\text{Cov}(\hat{\beta}_{alt}, \hat{\beta}_{OLS})$  и  $\text{Cov}(\hat{\beta}_{OLS}, \hat{\beta}_{alt})$ .
- ж) Полученное в предыдущем пункте выражение должно вызывать ностальгию, так как очень похоже на... На что?

Теперь пора посмотреть на разницу  $\hat{\beta}_{alt} - \hat{\beta}_{OLS}$ :

- з) Вспомните или выведите формулу для  $\text{Var}(r + s)$ , где  $r$  и  $s$  — случайные векторы одинакового размера
  - и) Докажите, что  $\text{Var}(\hat{\beta}_{alt} - \hat{\beta}_{OLS}) = \text{Var}(\hat{\beta}_{alt}) - \text{Var}(\hat{\beta}_{OLS})$
  - к) Рассмотрим матрицу  $C = \text{Var}(\hat{\beta}_{alt}) - \text{Var}(\hat{\beta}_{OLS})$ . Что находится на её диагонали?
  - л) Является ли матрица  $C$  симметричной?
  - м) Докажите, что матрица  $C$  является положительно полуопределённой. Если кто забыл, то это означает, что для любого вектора  $a$  выполнено неравенство  $a'Ca \geq 0$
  - н) Докажите, что диагональные элементы матрицы  $C$  не меньше нуля
2. Вектор  $u$  размера  $3 \times 1$  имеет стандартное нормальное распределение,  $u \sim \mathcal{N}(0; I)$ . Дана матрица  $D$ :

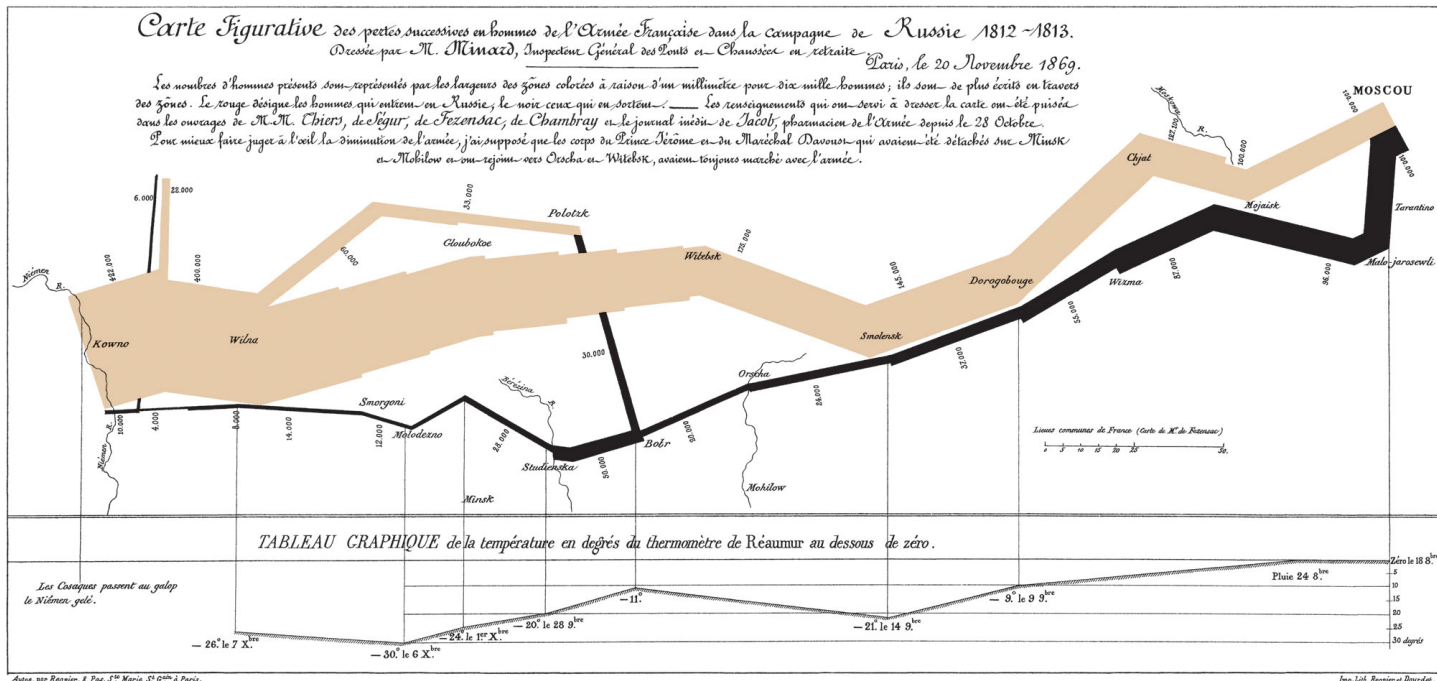
$$D = \frac{1}{14} \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{pmatrix}$$

- а) Является ли матрица  $D$  симметричной? Идемпотентной?
  - б) Найдите собственные числа матрицы  $D$  с учётом кратности
  - в) Как распределена случайная величина  $u'Du$ ?
  - г) Аккуратно объясните, в чём состоит геометрический смысл умножения произвольного вектора  $y$  на матрицу  $D$ ?
3. Найдите величины  $ESS$ ,  $RSS$ ,  $TSS$  и  $R^2$  для регрессии  $y_i = \mu + u_i$
4. В прошлом году, в курсе теории вероятностей и математической статистики, использовалось без доказательства следующее утверждение:
- Если случайные величины  $y_i$  независимы и нормально распределены  $y_i \sim \mathcal{N}(\mu; \sigma^2)$ , то  $q = \sum_{i=1}^n (y_i - \bar{y})^2 / \sigma^2$  имеет хи-квадрат распределение с  $(n - 1)$  степенью свободы.

Настала пора вернуть долг чести и доказать это утверждение :)

- а) Рассмотрим вектор центрированных  $y$ , то есть такой вектор  $\tilde{y}$ , что  $\tilde{y}_i = y_i - \bar{y}$ . Представьте вектор  $\tilde{y}$  в виде  $\tilde{y} = Ay$ . Как выглядит матрица  $A$ ?
- б) Является ли матрица  $A$  симметричной? Идемпотентной? Найдите все её собственные числа с учётом кратности.
- в) Представьте скаляр  $q$  в виде  $q = \frac{1}{\sigma^2} \tilde{y}' B \tilde{y}$ . Как выглядит матрица  $B$ ?
- г) Представьте скаляр  $q$  в виде  $q = \frac{1}{\sigma^2} y' C y$ . Как выглядит матрица  $C$ ?
- д) Представьте скаляр  $q$  в виде  $q = u' D u$ , где вектор  $u \sim \mathcal{N}(0; I)$ . Как выглядит матрица  $D$ ?

- е) Является ли матрица  $D$  симметричной? Идемпотентной? Найдите все её собственные числа с учётом кратности.
- ж) Сформулируйте теорему о законе распределение квадратичной формы нормальных случайных величин и верните долг чести.



Charles Joseph Minard, Схема потерь наполеоновской армии в компании 1812-1813 годов

Рис. 3: \*

## 8.9. ИП. Битва — решения

1. (25 points)

- а) (1)
- б) (1)
- в) (2)
- г) (2) Является, каждый элемент вектора оценок, это линейная комбинация значений  $y$ .
- д) (2)

$$\mathbb{E}(\hat{\beta}_{alt}) = \mathbb{E}(Ay) = A\mathbb{E}(y) = AX\beta = \beta \Rightarrow AX = I$$

е) (2)

$$\text{Cov}(Ay, (X'X)^{-1}X'y) = A\text{Cov}(y, y)X(X'X)^{-1} = \sigma^2AX(X'X)^{-1} = \sigma^2(X'X)^{-1}$$

$$\text{Cov}((X'X)^{-1}X'y, Ay) = (X'X)^{-1}X'\text{Cov}(y, y)A' = \sigma^2(X'X)^{-1}(AX)' = \sigma^2(X'X)^{-1}$$

ж) (2) На ковариационную матрицу  $\hat{\beta}_{OLS}$ !

з) (2)

$$\text{Var}(r + s)$$

Понятно, что на диагонали такой матрицы будут находиться элементы следующего вида

$$\text{Var}(r_i + s_i) = \text{Var}(r_i) + \text{Var}(s_i) + 2\text{Cov}(r_i, s_i)$$



Все остальные элементы

$$\text{Cov}(r_i + s_i, r_j + s_j) = \text{Cov}(r_i, r_j) + \text{Cov}(r_i, s_j) + \text{Cov}(s_i, r_j) + \text{Cov}(s_i, s_j)$$

Первое слагаемое принадлежит матрице  $\text{Var}(r)$ , третье —  $\text{Var}(s)$ , второе —  $\text{Cov}(r, s)$ , четвертое —  $\text{Cov}(s, r)$ . Получаем

$$\text{Var}(r + s) = \text{Var}(r) + \text{Var}(s) + \text{Cov}(r, s) + \text{Cov}(s, r)$$

и) (1)

$$\begin{aligned} \text{Var}(\hat{\beta}_{alt} - \hat{\beta}_{OLS}) &= \text{Var}(\hat{\beta}_{alt}) + \text{Var}(\hat{\beta}_{OLS}) - 2\text{Cov}(\hat{\beta}_{alt}, \hat{\beta}_{OLS}) = \\ &= \text{Var}(\hat{\beta}_{alt}) + \text{Var}(\hat{\beta}_{OLS}) - 2\text{Var}(\hat{\beta}_{OLS}) = \text{Var}(\hat{\beta}_{alt}) - \text{Var}(\hat{\beta}_{OLS}) \end{aligned}$$

к) (2) На диагонали находится разница в дисперсиях коэффициентов.

л) (2)

$$\text{Var}(Ay) = A \text{Var}(y) A' = \sigma^2 A A'$$

Поэтому  $C$  симметрична, как разница симметричных матриц.

м) (4) Константа не влияет на определенность матрицы, поэтому докажем этот факт без неё.

$$\begin{aligned} a' C a &= a' (A A' - (X' X)^{-1}) a = a' (A A' - A X (X' X)^{-1} X' A') a = a' A (I - H) A' a = \\ &= a' A (I - H) (I - H)' A' a = \|(I - H)' A' a\|_2^2 \geq 0 \end{aligned}$$

н) (2) В матрице  $C$  на диагонали находятся дисперсии разности оценок, а дисперсия всегда неотрицательна.

2. (10)

а) (2) Является симметричной и идемпотентной

б) (3) Можем посчитать по известной формуле, но лучше не надо :) Заметим, что ранг равен 1. Значит у нас 2 нулевых собственных значения. Матрица идемпотентна, значит третье собственное значение равно 1.

в) (3) Матрица  $D$  представима в следующем виде:

$$\frac{1}{14} (1, 2, 3)' (1, 2, 3) \Rightarrow u' D u = \frac{1}{14} u' (1, 2, 3)' (1, 2, 3) u$$

$u' (1, 2, 3)'$  — это сумма нормально распределенных случайных величин с некоторыми коэффициентами, такая сумма тоже имеет нормальное распределение с матожиданием 0 и дисперсией равной 14. Нормируя на корень из 14, мы получаем стандартную нормальную случайную величину. Значит

$$\frac{1}{14} u' (1, 2, 3)' (1, 2, 3) u \sim \chi_1^2$$

г) (2) Это проекция на пространство строк/столбцов матрицы  $D$ . В частности, на прямую, порожденную вектором (1,2,3).

3. (7)

Оценка МНК даст  $\mu = \bar{y}$ .

$$ESS = \sum (\hat{y}_i - \bar{y})^2 = \sum (\bar{y} - \bar{y})^2 = 0$$

$$\begin{aligned} RSS &= \sum (\hat{y}_i - y_i)^2 = \sum (\bar{y} - y_i)^2 = TSS \\ R^2 &= 0 \end{aligned}$$

4. (16)

- а) (2) Пусть  $G$  — матрица из одних единиц. Её смысл в том, что она будет суммировать элементы вектора, который был умножен на неё. Тогда

$$A = I - \frac{1}{n}G$$

б) (3)

$$A' = A$$

$$(I - \frac{1}{n}G)(I - \frac{1}{n}G) = I - \frac{1}{n}G - \frac{1}{n}G + \frac{1}{n^2}GG$$

Легко видеть, что  $GG$  — это матрица, где каждый элемент равен  $n$ . Получаем

$$A^2 = A$$

Заметим, что матрица  $\frac{1}{n}G$  также идемпотентна. Значит её собственные значения равны единицам и нулям. Легко также видеть, что ранг равен 1. Значит у этой матрицы одно собственное значение равно 1, и все остальные равны 0. Умножая матрицу на  $-1$ , мы меняем знак собственных значений. А прибавляя единичную матрицу, мы увеличиваем все собственные значения на 1. Значит мы имеем  $n-1$  единичное собственное значение, и одно собственное значение равно 0.

в) (1)

$$B = I$$

- г) (3) Чтобы перейти к равенству предыдущего пункта, нужно применить преобразование  $A$  на вектора  $y$ . Тогда

$$q = \frac{1}{\sigma^2} \tilde{y}' \tilde{y} = \frac{1}{\sigma^2} y' A' A y \Rightarrow C = A' A$$

- д) (3) Заметим, что деление  $y$  на  $\sigma$  нормирует вектор  $y$ .

$$\frac{1}{\sigma^2} y' A' A y = \frac{1}{\sigma^2} y' A' A y = \left( \frac{\sigma u + \mu}{\sigma} \right)' A' A \left( \frac{\sigma u + \mu}{\sigma} \right) = (u + \frac{\mu}{\sigma})' A' A (u + \frac{\mu}{\sigma}) =$$

Так как  $A$  центрирует переменные, то  $A \frac{\mu}{\sigma} = 0$ , следовательно

$$(u + \mu)' A' A (u + \mu) = u' A' A u$$

$$D = A' A = A A = A$$

- е) (2) Аналогично пункту б)

- ж) (3) Осталось понять, какое распределение имеет  $u' A u = u' (I - \frac{1}{n}G) u$ . Как известно, матрицу  $G$  можно представить в следующем виде

$$\frac{1}{n}G = S \Lambda S'$$

Где у матрицы лямбда все элементы 0 кроме элемента  $S_{1,1}$ . Пусть  $s$  — это собственный вектор, который соответствует собственному значению 1, тогда

$$u' \frac{1}{n}G u = u' s s' u$$

Понятно, что собственный вектор с собственным значением 1, это вектор констант, но так как в данном разложении матрица  $S$  ортогональна, то длина собственного вектора  $s$  равна 1. Тогда все элементы вектора  $s$  равны  $\frac{1}{\sqrt{n}}$ .

$$u'ss'u = \left( \sum_{i=1}^n \frac{u_i}{\sqrt{n}} \right)^2$$

$$\sum_{i=1}^n \frac{u_i}{\sqrt{n}} \sim N(0, 1) \Rightarrow u'ss'u \sim \chi_1^2$$

$$u'u \sim \chi_n^2 \Rightarrow q = u'Au \sim \chi_{n-1}^2$$

## 8.10. Кр 2, экзамен за I семестр, демо

1. На основании опроса была оценена следующая модель:

$$\ln(wage_i) = \beta_1 + \beta_2 exper_i + \beta_3 age_i + \beta_4 sex_i + \varepsilon_i$$

где:

- $wage_i$  — величина заработной платы в долларах
- $exper_i$  — опыт работы в годах
- $age_i$  — возраст в годах
- $sex_i$  — пол (1 — мужской, 0 — женский)

Показатель	Значение
$R^2$	0.903
Скорректированный $R^2$	<b>B7</b>
Стандартная ошибка регрессии	<b>B6</b>
Количество наблюдений	<b>B2</b>

Результаты дисперсионного анализа:

	df	SS	MS	F	P-значение
Регрессия	3	2638.3	879.4	<b>B5</b>	0.000
Остаток	<b>B1</b>	282.1	16.6		
Итого	<b>B3</b>	<b>B4</b>			

Коэффициент	Оценка	$se(\hat{\beta})$	t-статистика	P-Значение	Нижние 95%	Верхние 95%
Константа	-6.96	12.38	-0.56	0.58	-33.08	19.16
$exper$	2.65	0.32	8.38	0.000	1.98	3.32
$age$	<b>B8</b>	<b>B9</b>	8.06	0.000	4.13	7.06
$sex$	-10.73	1.79	<b>B10</b>	0.000	-14.49	-6.95

а) Найдите пропущенные числа **B1--B10**.

б) Как изменятся результаты оценки регрессии, если переменную  $sex_i$  переопределить так, чтобы 0 соответствовал мужчинам, 1 — женщинам?

Ответ округляйте до 2-х знаков после запятой. Кратко поясняйте, например, формулой, как были получены результаты.

2. Юный эконометрист Вениамин очень любит опаздывать на первую пару. Он считает, что время (в минутах), на которое он опаздывает,  $time_t$ , зависит от количества снега (в миллиметрах), выпавшего за последние сутки,  $snow_t$ . После месяца упорного сбора данных, он смог оценить следующую модель:

$$\widehat{time}_t = 12 + 0.2snow_t$$

Оценка ковариационной матрицы коэффициентов,  $\widehat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} 5 & -0.03 \\ -0.03 & 0.01 \end{pmatrix}$

Оценка дисперсии ошибок равна  $\hat{s}^2 = 1.45$ .

За последние сутки выпало 15 миллиметров снега.

- Постройте точечный прогноз времени опоздания Вениамина
  - Постройте 95%-ый доверительный интервал для  $\mathbb{E}(time_t | snow_t = 15)$  для ожидаемого опоздания Вениамина
  - Постройте 95%-ый предиктивный интервал (доверительный интервал) для фактического опоздания Вениамина
3. По данным 27 фирм, упорядоченных по капиталу,  $K_1 < K_2 < \dots < K_n$ , была оценена зависимость выпуска  $Y_i$  от труда  $L_i$  и капитала  $K_i$  с помощью моделей
- $\ln Y_i = \beta_1 + \beta_2 \ln L_i + \beta_3 \ln K_i + u_i$
  - $\ln Y_i = \beta_1 + \beta_2 \ln \frac{L_i}{K_i} + u_i$

	Модель (1)	Модель (2)
константа	1.1706 (0.326)	1.0692 (0.1317)
$\ln L$	0.6029 (0.125)	
$\ln K$	0.375 (0.085)	
$\ln \frac{L}{K}$		0.6369 (0.0754)
R-squared	0.943	0.74
F	200.24	71.351
P-value	0.0	0.0
RSS	0.851	
N	27	27

- На уровне значимости  $\alpha = 0.05$  проверьте для модели (1) гипотезы  $H_0: \beta_2 = \beta_3 = 0$  и  $H_0: \beta_3 = 0.5$ .
- Объясните, почему модель (2) является ограниченной версией модели (1). Явно выпишите ограничения. На уровне значимости  $\alpha = 0.05$  проверьте гипотезу об этих ограничениях.
- Фирмы разделили на маленькие,  $i \leq 14$ , и большие,  $i \geq 15$ . Для этих двух групп оценили отдельные регрессии. Можно ли считать, что производственные функции для больших и маленьких фирм не различаются? Ответ подтвердите подходящим тестом.

	Модель для $i \leq 14$	Модель для $i \geq 15$
константа	0.6998 (0.649)	1.4082 (0.678)
$\ln L$	0.9000 (0.133)	0.0081 (0.226)
$\ln K$	0.2100 (0.056)	0.805 (0.179)
R-squared	0.896	0.908
F	47.84	49.81
P-value	0.0	0.0
RSS	0.119	0.362
N	14	13

4. С помощью теста Бокса-Кокса оценили зависимость веса индивида (в килограммах) от его роста (в сантиметрах):

Log likelihood = -2659.5656

Number of obs = 540  
LR chi2(2) = 230.68  
Prob > chi2 = 0.000

w	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
/lambda	1.055498	1.892654	0.56	0.577	-2.654035	4.76503
/theta	-.0263371	.1471576	-0.18	0.858	-.3147607	.2620865

Estimates of scale-variant parameters

	Coef.
Notrans _cons	2.936809
Trans h	.0237224
/sigma	.1660251

Test H0:	Restricted log likelihood	chi2	Prob > chi2
theta=lambda = -1	-2680.8693	42.61	0.000
theta=lambda = 0	-2659.7618	0.39	0.531
theta=lambda = 1	-2685.5201	51.91	0.000

Какую спецификацию модели (линейную, линейную в логарифмах, полулогарифмическую) следует предпочесть и почему?

5. По данным для 23 демократических стран оценили зависимость индекса Джини<sup>3</sup> от ВВП на душу населения с учетом ППС (паритета покупательной способности). Затем провели тест Рамсея.

<sup>3</sup>Индекс Джини — это мера неравенства доходов. Чем выше индекс Джини, тем выше неравенство.

```
. reg gini gdp if democ==1
```

Source	SS	df	MS	Number of obs =	23
Model	506.853501	1	506.853501	F( 1, 21) =	13.05
Residual	815.572523	21	38.8367868	Prob > F =	0.0016
				R-squared =	0.3833
				Adj R-squared =	0.3539
				Root MSE =	6.2319
Total	1322.42602	22	60.1102738		

gini	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
gdp	-.0006307	.0001746	-3.61	0.002	-.0009937 -.0002676
_cons	44.30983	3.572733	12.40	0.000	36.87993 51.73974

```
. ovtest
```

```
Ramsey RESET test using powers of the fitted values of gini
Ho: model has no omitted variables
F(3, 18) = 5.16
Prob > F = 0.0095
```

- Сформулируйте нулевую и альтернативную гипотезу теста Рамсея
- Опишите пошагово, как проводится тест Рамсея
- Прокомментируйте результаты теста Рамсея

## 8.11. Кр 2, экзамен за I семестр, демо, решения

```
1 y_hat <- 12 + 0.2 * 15
2 se_y_hat <- 5 + 225 * 0.01 + 2 * 15 * 0.03
3 se_forecast_error <- se_y_hat + 1.45
4 t_crit <- qt(0.975, df = 30 - 2)
```

- Или руками:

**B1:** Воспользуемся следующим соотношением:

$$\frac{\text{Residual SS}}{\text{Residual df}} = \text{Residual MS} \Rightarrow \mathbf{B1} = \text{Residual df} = 17$$

$$\mathbf{B2: Residual df} = n - k \Rightarrow n = 21$$

$$\mathbf{B3: Total df} = n - 1 = 20$$

$$\mathbf{B4: TSS} = RSS + ESS \Rightarrow TSS = 2920.4$$

$$\mathbf{B5: F} = \frac{\text{Regression MS}}{\text{Residual MS}} = \frac{879.4}{16.6} = 52.98$$

$$\mathbf{B6: \hat{\sigma}^2} = \frac{RSS}{n-k} = \frac{282.1}{21-4} = 16.59$$

$$\mathbf{B7: R_{adj}^2} = 1 - \frac{RSS/(n-k)}{TSS/(n-1)} = 1 - \frac{282.1/17}{2920.4/20} = 0.87$$

**B8 и B9** получаются из следующей системы, где 2.11 берётся из таблицы распределения Стьюдента,  $t_{21-4;0.025}$ :

$$\begin{cases} \hat{\beta}_{age} - 2.11 \cdot se(\hat{\beta}_{age}) = 4.13 \\ \hat{\beta}_{age} + 2.11 \cdot se(\hat{\beta}_{age}) = 7.06 \end{cases} \Rightarrow \begin{cases} \hat{\beta}_{age} = 5.63 \\ se(\hat{\beta}_{age}) = 0.71 \end{cases}$$

$$\mathbf{B10: t_{obs}} = \frac{\hat{\beta}_{sex}}{se(\hat{\beta}_{sex})} \Rightarrow t_{obs} = -5.99$$

- $\hat{\beta}_{4,new} = -\hat{\beta}_{4,old}$ ,  $\hat{\beta}_{1,new} = \hat{\beta}_{1,old} - \hat{\beta}_{4,new}$ , остальные не изменятся.

---

```

2_1 t_obs <- (0.385 - 0.5) / 0.085
2_2 t_crit <- qt(0.975, df = 27 - 3)
3_1 F_obs <- (0.943 - 0.74) / (1 - 0.943) * (27 - 3)
4_1 F_crit <- qf(0.95, df1 = 1, df2 = 27 - 3)

```

---

Или руками:

а)  $\widehat{time}_f = 12 + 0.2 \cdot 15 = 15$

б) Доверительный интервал для среднего значения имеет вид:

$$[\widehat{time}_f - z_{cr} \cdot se(\widehat{time}_f); \widehat{time}_f + z_{cr} \cdot se(\widehat{time}_f)]$$

$$\begin{aligned} \text{Var}(\widehat{time}_f - \mathbb{E}(time_f|X)|X) &= \text{Var}(\widehat{time}_f|X) = \text{Var}(\beta_1 + 15\beta_2|X) = \\ &= \text{Var}(\beta_1|X) + 225 \text{Var}(\beta_2|X) + 2 \cdot 15 \text{Cov}(\beta_1, \beta_2|X) = 6.35 \end{aligned}$$

$$[15 - 2\sqrt{6.35}; 15 + 2\sqrt{6.35}]$$

в) Доверительный (предиктивный) интервал для конкретного значения имеет вид:

$$[\widehat{time}_f - z_{cr} \cdot se(\widehat{time}_f - \varepsilon_f); \widehat{time}_f + z_{cr} \cdot se(\widehat{time}_f - \varepsilon_f)]$$

$$\begin{aligned} \text{Var}(\widehat{time}_f - time_f|X) &= \text{Var}(\widehat{time}_f - \mathbb{E}(time_f|X) - \varepsilon_f|X) = \\ &= \text{Var}(\widehat{time}_f - \varepsilon_f|X) = \text{Var}(\widehat{time}_f|X) + \text{Var}(\varepsilon_f|X) = 6.35 + 1.45 = 7.8 \end{aligned}$$

$$[15 - 2\sqrt{7.8}; 15 + 2\sqrt{7.8}]$$

3. а) Гипотеза об адекватности множественной регрессии проверяется, с помощью следующей статистики:

$$F = \frac{ESS/(k-1)}{RSS_{UR}/(n-k)} = \frac{R^2/(k-1)}{(1-R^2)/(n-k)}$$

где  $F \sim F_{k-1, n-k}$  при верной  $H_0$ . Найдём наблюдаемое и критическое значения статистики:

$$F_{obs} = \frac{0.943/2}{(1-0.943)/(27-3)} = 198.5 \quad F_{crit} = 3.4$$

Поскольку  $F_{obs} > F_{crit}$ , основная гипотеза отвергается.

Также можно было заметить, что в таблице уже дано  $F = 200.24$  и сделать тот же вывод.

Проверим  $H_0 : \beta_3 = 0.5$ :

$$t_{obs} = \frac{\hat{\beta}_3 - \beta_3}{se(\hat{\beta}_3)} = \frac{0.375 - 0.5}{0.085} = -1.47 \quad t_{crit} = 2.06$$

Так как  $|t_{obs}| < t_{crit}$ , оснований отвергать  $H_0$  нет.

б) Модель (2) получается из модели (1) при  $\beta_3 = -\beta_2$ . Будем проверять гипотезу  $H_0 : \beta_3 + \beta_2 = 0$ .

$$F_{obs} = \frac{RSS_R - RSS_{UR}/q}{RSS_{UR}/(n - k_{UR})} = \frac{(R_{UR}^2 - R_R^2)/q}{(1 - R_{UR}^2)/(n - k_{UR})} = \frac{(0.943 - 0.74)/1}{(1 - 0.943)/(27 - 3)} = 85.47$$

$F_{obs} > F_{crit} = 4.26 \Rightarrow$  основная гипотеза отвергается.

в) Введём дополнительную переменную:

$$d_i = \begin{cases} 1, & i \leq 14 \\ 0, & i > 14 \end{cases}$$

Тогда можем записать неограниченную и ограниченную модели:

$$UR : \ln Y_i = \beta_1 + \Delta_1 d_i + (\beta_2 + \Delta_2 d_i) \ln L_i + (\beta_3 + \Delta_3 d_i) \ln K_i + u_i$$

$$R : \ln Y_i = \beta_1 + \beta_2 \ln L_i + \beta_3 \ln K_i + u_i$$

Заметим, что  $RSS_{UR} = 0.119 + 0.362 = 0.481$ .

Будем проверять следующую гипотезу:

$$H_0 : \begin{cases} \Delta_1 = 0 \\ \Delta_2 = 0 \\ \Delta_3 = 0 \end{cases} \quad H_a : \Delta_1^2 + \Delta_2^2 + \Delta_3^2 > 0$$

Для этого считаем F-статистику,  $F \sim F_{3,21}$  при верной  $H_0$ :

$$F_{obs} = \frac{(RSS_R - RSS_{UR})/q}{RSS_{UR}/(n - k_{UR})} = \frac{(0.851 - 0.481)/3}{0.481/(27 - 6)} = 5.38 \quad F_{crit} = 3.07$$

Поскольку  $F_{obs} > F_{crit}$ , основная гипотеза отвергается.

4. Гипотезы  $H_0 : \lambda = \theta = -1$  и  $H_0 : \lambda = \theta = 1$  (последняя соответствует линейной спецификации) отвергаются. Гипотеза  $H_0 : \lambda = \theta = 0$  (соответствует логарифмической спецификации) не отвергается. Значит, предпочтительной является логарифмическая модель.
5. а)  $H_0 : gini_i = \beta_1 + \beta_2 \cdot gdp_i + u_i$ , то есть в модели нет пропущенных переменных;  
 $H_A$  : есть неизвестные нам пропущенные регрессоры.
  - б) • Оценить модель  $gini_i = \beta_1 + \beta_2 \cdot gdp_i + u_i$ , получить прогнозы  $\widehat{gini}_i$ .
  - Оценить вспомогательную регрессию  $gini_i = \beta_1 + \beta_2 \cdot gdp_i + \gamma_1 \widehat{gini}_i^2 + \gamma_2 \widehat{gini}_i^3 + \dots + \gamma_p \widehat{gini}_i^{p+1} + u_i$
  - Посчитать F-статистику, проверяющую гипотезу о равенстве всех  $\gamma_i$  нулю.  
 При верной  $H_0$  и нормальности остатков  $F \sim F_{p, n-k-p}$ .
- в) Основная гипотеза отвергается на любом разумном уровне значимости. Значит, в модели есть пропущенные факторы.

## 8.12. Кр 2, экзамен за I семестр, 24.12.2016

---

```

1 set.seed(var_no)
2 n_obs <- 200
3 opros <- data_frame(exper = rnorm(n_obs, mean = 7, sd = 2),
4                       exper2 = exper^2,
5                       sex = sample(0:1, n_obs, rep = TRUE),
6                       eps = rnorm(n_obs, sd = 2),
7                       wage = 3 + 6 * exper - 0.2 * exper2 + 1.5 * sex + eps)
8 model <- lm(data = opros, wage ~ exper + exper2 + sex)
9 report <- summary(model)
10 coefs <- coef(model)

```

---



1. На основании опроса 200 человек была оценена следующая модель:

$$\ln(wage_i) = \beta_1 + \beta_2 exper_i + \beta_3 exper_i^2 + \beta_4 sex_i + \varepsilon_i$$

где:

- $wage_i$  — величина заработной платы в долларах
- $exper_i$  — опыт работы в годах
- $exper_i^2$  — опыт работы в годах
- $sex_i$  — пол (1 — мужской, 0 — женский)

Показатель	Значение
$R^2$	0.911
Скорректированный $R^2$	<b>B7</b>
Стандартная ошибка регрессии	<b>B6</b>
Количество наблюдений	<b>B2</b>

Результаты дисперсионного анализа:

	df	сумма квадратов	F	P-значение
Регрессия	3	<b>B9</b>	<b>B5</b>	0.000
Остаток	<b>B1</b>	830.1		
Итого	<b>B3</b>	<b>B4</b>		

	Оценка	Ст. ошибка	t-статистика	P-Значение
Константа	3.6869	1.1960	3.08	0.0023
$exper$	<b>B8</b>	0.3525	16.45	0.0000
$exper^2$	-0.1916	0.0254	-7.54	0.0000
$sex$	1.5745	0.2937	<b>B10</b>	0.0000

а) Найдите пропущенные числа **B1--B10**.

б) Как изменятся результаты оценки регрессии, если переменную  $sex_i$  переопределить так, чтобы 0 соответствовал мужчинам, 1 — женщинам?

Ответ округляйте до 2-х знаков после запятой. Кратко поясняйте, например, формулой, как были получены результаты.

2. Исследовательница Глафира изучает зависимость спроса на молоко от цены молока и дохода семьи. В её распоряжении есть следующие переменные:

- $price$  — цена молока в рублях за литр
- $income$  — ежемесячный доход семьи в тысячах рублей
- $milk$  — расходы семьи на молоко за последние семь дней в рублях

В данных указано, проживает ли семья в сельской или городской местности. Поэтому Глафира оценила три регрессии: (All) — по всем данным, (Urban) — по городским семьям, (Rural) — по сельским семьям.

```

1 var_no <- 1
2 set.seed(var_no)
3 n_obs <- 100
4
5 milk_demand <- data_frame(
6   eps = rnorm(n_obs, sd = 5),

```

```

7 price = rnorm(n_obs, mean = 20, sd = 3),
8 city = sample(0:1, n_obs, rep = TRUE),
9 income = rnorm(n_obs, mean = 70, sd = 10))
10
11 milk_demand <- mutate(milk_demand, milk = ifelse(city == 1,
12                                     1 + 0.2 * income - 0.1 * price + eps,
13                                     0 + 0.3 * income - 0.5 * price + eps))
14
15 model_all <- lm(data = milk_demand, milk ~ income + price)
16 model_urban <- lm(data = filter(milk_demand, city == 1), milk ~ income + price)
17 model_rural <- lm(data = filter(milk_demand, city == 0), milk ~ income + price)
18
19 model_table <- mtable("(All)" = model_all, "(Urban)" = model_urban, "(Rural)" = model_rural,
20                       summary.stats = c("R-squared", "adj. R-squared", "sigma", "F", "p", "Deviance", "N"))
21
22 model_table <- relabel(model_table, Deviance = "RSS", p = "P-value", N = "n observations")
23
24 toLatex(model_table)

```

---

	(All)	(Urban)	(Rural)
(Intercept)	1.479 (4.480)	-0.797 (7.808)	4.598 (5.121)
income	0.252*** (0.049)	0.204* (0.092)	0.262*** (0.053)
price	-0.335* (0.165)	0.001 (0.272)	-0.567** (0.194)
R-squared	0.230	0.104	0.380
adj. R-squared	0.214	0.063	0.355
sigma	4.678	5.036	4.160
F	14.464	2.541	15.326
P-value	0.000	0.090	0.000
RSS	2123.000	1115.693	865.080
n observations	100	47	53

---

- а) Проверьте значимость в целом регрессии (All) на 5%-ом уровне значимости.
- б) На 5%-ом уровне значимости проверьте гипотезу, что зависимость спроса на молоко является единой для городской и сельской местности.
3. Исследовательница Глафира продолжает изучать спрос на молоко. В её распоряжении по-прежнему данные по трём переменным:
- *price* — цена молока в рублях за литр
  - *income* — ежемесячный доход семьи в тысячах рублей
  - *milk* — расходы семьи на молоко за последние семь дней в рублях

Имеются результаты оценивания модели  $milk_i = \beta_1 + \beta_2 income_i + \beta_3 price_i + u_i$  по 100 наблюдениям:

---

```

1 xtable(model_all)

```

---

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.4791	4.4796	0.33	0.7420
income	0.2524	0.0486	5.19	0.0000
price	-0.3354	0.1649	-2.03	0.0447

Коэффициент детерминации  $R^2$  оказался равен 0.23.

Глафира рассчитала оценку ковариационной матрицы исходных переменных:

---

```
1 xtable(var(dplyr::select(milk_demand, price, income, milk)))
```

---

	price	income	milk
price	8.26	3.48	-1.89
income	3.48	95.09	22.83
milk	-1.89	22.83	27.84

---

- Постройте точечный прогноз расходов на молоко семьи с доходом 100 тысяч рублей при цене на молоко 30 рублей за литр.
  - Найдите выборочную корреляцию между фактическими расходами на молоко и их прогнозами.
  - Разложите коэффициент детерминации  $R^2$  в модели в сумму эффектов переменных *income* и *price*.
4. По квартальным данным 1958-1976 годов была оценена модель с тремя объясняющими факторами:

$$\hat{Y}_i = 2.2 + 0.104X_i - 3.48Z_i + 0.34W_i, \quad ESS = 100, \quad RSS = 120$$

- Какую модель необходимо оценить исследователю, если он считает, что в различные сезоны среднее значение зависимой переменной помимо зависимости от трёх регрессоров может отличаться на константу?
  - При оценивании модели, допускающей сезонные эффекты, оказалось, что значение  $ESS$  увеличилось до 160. На уровне значимости 5% проверьте гипотезу о наличии сезонности.
- 

```
1 set.seed(var_no)
2 hb_w <- sample(c(3, 4, 5), 1)
3 hb_r <- sample(c(6, 7, 8), 1)
4
5 X <- mvtnorm::rmvnorm(n = 24, mean = c(1, 1, 1),
6   sigma = matrix(c(0, 0, 0, 0, 0.4, -0.1, 0, -0.1, 0.9), nrow = 3))
7 XXm <- solve(crossprod(X))
8 xmatrix <- function(a, environment = "pmatrix", output = TRUE, digits = 3) {
9
10   # override default alignment for xtable
11   xa <- xtable::xtable(a, align = rep("", ncol(a) + 1), digits = digits)
12
13   res <- print(xa,
14     floating = FALSE,
15     tabular.environment = environment,
16     hline.after = NULL,
17     include.rownames = FALSE,
18     include.colnames = FALSE,
```

```

19     file = "junk.txt")
20
21   res <- paste0("\\ensuremath{" , res, "}")
22
23   if (output) {
24     cat(res)
25   }
26   return(invisible(res))
27 }
28 xmatrix(XXm)

```

---

5. По 24 наблюдениям была оценена модель:

$$\hat{Y}_i = 15 - 4Z_i + 3W_i$$

Известно, что случайные ошибки нормально распределены,  $RSS = 180$ , и

$$(X'X)^{-1} = \begin{pmatrix} 0.365 & -0.218 & -0.084 \\ -0.218 & 0.184 & 0.027 \\ -0.084 & 0.027 & 0.046 \end{pmatrix}$$

- а) Проверьте гипотезу  $H_0 : \beta_Z = 0$  против  $H_a : \beta_Z \neq 0$  на уровне значимости 5%.
- б) Проверьте гипотезу  $H_0 : \beta_Z + \beta_W = 0$  против  $H_a : \beta_Z + \beta_W \neq 0$  на уровне значимости 5%.
- в) Выпишите использованные при проверке гипотез предположки о случайных ошибках модели.

### 8.13. Кр 2, экзамен за I семестр, 24.12.2016, решения

1. а) **B1** =  $n - k = 200 - 4 = 196$

Из условия, **B2** =  $n = 200$

**B3** =  $n - 1 = 199$

Для нахождения **B4** воспользуемся знанием значения  $R^2$ :

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{830.1}{\mathbf{B4}} = 0.911 \Rightarrow \mathbf{B4} = 9326.97$$

Для того, чтобы найти **B5**, понадобится **B9** =  $TSS - RSS = 8496.87$ . Тогда:

$$\mathbf{B5} = F = \frac{\text{Regression MS}}{\text{Residual MS}} = \frac{8496.87/3}{830.1/196} = 668.75$$

$$\mathbf{B6} = \hat{\sigma}^2 = \frac{RSS}{n - k} = 4.24$$

$$\mathbf{B7} = R_{adj}^2 = 1 - \frac{RSS/(n - k)}{TSS/(n - 1)} = 0.91$$

Значение **B8** находится из соотношения:

$$\frac{\hat{\beta}}{se(\hat{\beta})} = t_{obs} \Rightarrow \mathbf{B8} = 5.8$$

Аналогично, для **B10**, **B10** = 5.36.

- б)  $\hat{\beta}_{4,new} = -\hat{\beta}_{4,old}$ ,  $\hat{\beta}_{1,new} = \hat{\beta}_{1,old} - \hat{\beta}_{4,new}$ , остальные не изменятся.
2. а) Необходимо проверить следующую гипотезу:

$$H_0 : \begin{cases} \beta_1 = 0 \\ \beta_2 = 0 \\ \beta_3 = 0 \end{cases} \quad H_a : \beta_1^2 + \beta_2^2 + \beta_3^2 > 0$$

Для этого воспользуемся F-статистикой,  $F \sim F_{2,97}$ , она дана в таблице:

$$F_{obs} = \frac{ESS/(k-1)}{RSS_{UR}/(n-k_{UR})} = 14.464$$

Находим  $F_{crit} = 3.09$ . Поскольку  $F_{obs} > F_{crit}$ , гипотеза отвергается.

- б) Для того, чтобы записать спецификацию неограниченной модели, введём дополнительную переменную  $d_i$ :

$$d_i = \begin{cases} 1 & \text{городская местность} \\ 0 & \text{сельская местность} \end{cases}$$

Пусть коэффициенты для городской местности отличаются на некоторое  $\Delta_i$ , тогда неограниченная модель имеет вид:

$$milk_i = \beta_1 + \Delta_1 d_i + (\beta_2 + \Delta_2 d_i) price_i + (\beta_3 + \Delta_3 d_i) income_i + \varepsilon_i$$

$$RSS_{UR} = RSS_{Urban} + RSS_{Rural} = 1115.693 + 865.08 = 1980.773$$

Проверяем следующую гипотезу:

$$H_0 : \begin{cases} \Delta_1 = 0 \\ \Delta_2 = 0 \\ \Delta_3 = 0 \end{cases} \quad H_a : \Delta_1^2 + \Delta_2^2 + \Delta_3^2 > 0$$

Для этого считаем F-статистику,  $F \sim F_{3,94}$  при верной  $H_0$ :

$$F_{obs} = \frac{(2123 - 1980.773)/3}{1980.773/94} \approx 2.25$$

Сравниваем с  $F_{crit} = 2.7$ ,  $F_{obs} < F_{crit} \Rightarrow$  нет оснований отвергать гипотезу.

3. а)  $\widehat{milk} = 1.4791 + 0.2524 \cdot 100 - 0.3354 \cdot 30 = 16.6571$
- б)  $sCorr(milk, \widehat{milk}) = \sqrt{R^2} = \sqrt{0.23}$
- в)  $R^2 = \sum_{j=2}^k \hat{\beta}_j \frac{sCov(x_j, y)}{sVar(y)} = 0.2524 \cdot \frac{22.83}{27.84} - 0.3354 \cdot \frac{(-1.89)}{27.84}$
4. а) В модель необходимо добавить фиктивные переменные  $Q_1, Q_2$  и  $Q_3$ , соответствующие первым трём кварталам года:

$$Y_i = \beta_1 + \beta_2 X_i + \beta_3 Z_i + \beta_4 W_i + \beta_5 Q_{1i} + \beta_6 Q_{2i} + \beta_7 Q_{3i} + \varepsilon_i$$

- б) Гипотеза о наличии (а на самом деле отсутствии) сезонности записывается следующим образом:

$$H_0 : \begin{cases} \beta_5 = 0 \\ \beta_6 = 0 \\ \beta_7 = 0 \end{cases} \quad H_a : \beta_5^2 + \beta_6^2 + \beta_7^2 > 0$$

Поскольку гипотезы оцениваются по одним и тем же данным,  $TSS$  в них совпадает и равен  $ESS_R + RSS_R = 100 + 120 = 220$

Тогда в неограниченной модели  $RSS_{UR} = 220 - 160 = 60$ .

Количество наблюдений:  $(1976 - 1958 + 1) \cdot 4 = 76$ . Считаем, что данные взяты с первого квартала 1958 года по четвёртый квартал 1976 года.

Теперь можно проверять гипотезу,  $F \sim F(3, 69)$ :

$$F_{obs} = \frac{(RSS_R - RSS_{UR})/q}{RSS_{UR}/(n - k_{UR})} = \frac{(120 - 60)/3}{40/69} = 34.5$$

Так как  $F_{obs} > F_{crit} = 2.73$ , гипотеза отвергается.

5. а) Чтобы найти  $se(\beta_Z)$ , сначала найдём стандартную ошибку регрессии  $\sigma^2$ :

$$\sigma^2 = \frac{RSS}{n - k} = \frac{180}{24 - 3} \approx 8.57$$

Тогда  $se(\beta_Z) = \sqrt{8.57 \cdot 0.184} = \sqrt{1.57688}$ . Теперь можно проверять гипотезу:

$$t_{obs} = \frac{-4}{\sqrt{1.57688}} \approx -3.19 \quad t_{crit} = -2.08$$

Основная гипотеза отвергается.

- б)  $t$ -статистика имеет вид:

$$t_{obs} = \frac{\hat{\beta}_Z + \hat{\beta}_W - 0}{se(\hat{\beta}_Z + \hat{\beta}_W)} = \frac{-4 + 3}{\sqrt{8.57(0.184 + 0.046 + 2 \cdot 0.027)}} \approx -0.64 \quad t_{crit} \approx 2.08$$

Нет оснований отвергать  $H_0$ .

- в)  $u_i \sim \mathcal{N}(0, \sigma_u^2)$

## 8.14. Кр 3, задачи для подготовки

<https://www.hse.ru/mirror/pubs/share/203792575>

## 8.15. Кр 3, 20.03.2017

Часть 1. Тест.

Часть 2. Задачи.

1. По данным для 39 районов Балтимора в 1970 г. были оценены уравнения

$$\ln \hat{Y}_i = \underset{t=54.7}{10.093} - \underset{t=-12.28}{0.239} X_i, \quad R^2 = 0.803$$

и

$$\frac{\ln \hat{Y}_i}{\sqrt{X_i}} = \underset{t=47.87}{9.093} \frac{1}{\sqrt{X_i}} - \underset{t=-15.10}{0.2258} \sqrt{X_i},$$

где  $Y_i$  — плотность населения района,  $X_i$  — расстояние до центрального делового квартала.

- а) С какой целью оценили второе уравнение? Какое при этом было сделано предположение о дисперсии ошибок?

б) Дайте интерпретацию полученным результатам.

2. Были обследованы 36 предприятий по трём показателям:  $K_i$  — основным фондам (млн. руб.),  $W_i$  — фонду оплаты труда (млн. руб.),  $R_i$  — расходам на НИОКР (млн. руб.). Получены оценки вектора средних  $\hat{\mu} = (3, 4, 1)'$  и ковариационной матрицы

$$\hat{\Sigma} = \begin{pmatrix} 2 & 3 & 0 \\ 3 & 10 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

Найдите первую главную компоненту и определите долю суммарной дисперсии, которую она объясняет.

3. Мартовский Заяц и Безумный Шляпник почти все время пьют чай. Известно, что количество выпитого за день чая (в чашках) зависит от количества пирожных (в штуках) и печенья (в штуках). Алиса, гостившая у героев в течение 25 дней, заметила, что если оценить зависимость выпитого чая от закуски для Мартовского Зайца и Шляпника

$$Tea_i = \beta_1 + \beta_2 Biscuit_i + \beta_3 Cake_i + u_i,$$

то получится регрессия с  $RSS = 17$ .

Чтобы понять, удачную ли модель она построила, Алиса оценила еще одну регрессию

$$Tea_i = \beta_1 + \beta_2 Biscuit_i + \beta_3 Cake_i + \gamma_2 \widehat{Tea_i^2} + \gamma_3 \widehat{Tea_i^3} + \gamma_4 \widehat{Tea_i^4} + \nu_i,$$

с  $RSS = 10$ .

Помогите Алисе понять, верную ли спецификацию модели она выбрала: сформулируйте основную и альтернативную гипотезы и проведите подходящий тест.

## 8.16. Кр 3, 20.03.2017, решения

1. Второе уравнение оценили для корректного построения доверительных интервалов и проверки гипотез о коэффициентах в условиях гетероскедастичности, для получения более эффективных оценок,  $\text{Var}(u_i) = \sigma^2 X_i$ .

После применения взвешенного МНК оба коэффициента значимы.

2. Собственные значения ковариационной матрицы:  $\lambda_1 = 11$ ,  $\lambda_2 = 3$ ,  $\lambda_3 = 1$ .

Первая главная компонента:  $(1/\sqrt{10} \quad 3/\sqrt{10} \quad 0)'$ .

Доля дисперсии:  $\frac{11}{15}$

- 3.

$$F = \frac{(RSS_R - RSS_{UR})/3}{RSS_{UR}/(n - k_{UR})} = \frac{(17 - 10)/3}{10/(25 - 6)} = 4.4(3)$$

## 8.17. ИП. Комоедица, 24.03.2016

Комоедица — белорусский народный праздник, посвящённый пробуждению медведя :)

1. Докажем свойства оценок максимального правдоподобия!

Пусть  $L(y|\theta)$  — функция правдоподобия,  $y$  — вектор-столбец из  $n$  наблюдений,  $\theta$  — вектор-столбец из  $m$  параметров. Кроме того, введём дополнительные обозначения,  $\ell(\theta) = \ln L(y|\theta)$ ,  $s(\theta) = \partial \ell / \partial \theta$ . Буква  $s$  сокращает слово «score».

Для наших целей мы определим информацию Фишера как  $I = \text{Var}(s(\theta))$ . То есть информация Фишера — это ковариационная матрица первых производных лог-функции правдоподобия. По определению.

- а) Чтобы взбодриться, укажите размеры векторов и матриц  $s(\theta)$ ,  $\mathbb{E}(s(\theta))$ ,  $\text{Var}(s(\theta))$ .
- б) Собрав всю силу воли в кулак, найдите  $\mathbb{E}(1)$ .
- в) Запишите  $\mathbb{E}(1)$  с помощью интеграла по  $dy$  и функции правдоподобия  $L()$ .
- г) Продифференцировав обе части найденного тождества по  $\theta_j$ , найдите  $\int \frac{\partial L}{\partial \theta_j} dy$ .
- д) Найдите  $\mathbb{E}\left(\frac{\partial \ell}{\partial \theta_j}\right)$ .
- е) Найдите  $\mathbb{E}(s(\theta))$ .
- ж) Докажите, что  $I = \mathbb{E}(s(\theta)s(\theta)')$ .
- з) Вспомнив магию дифференцирования ещё раз, найдите  $\int \frac{\partial^2 L}{\partial \theta_j \partial \theta_i} dy$ .
- и) Найдите  $\mathbb{E}\left(\frac{\partial^2 L}{\partial \theta_j \partial \theta_i} \frac{1}{L}\right)$ .
- к) Выразите  $\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_i}$  через  $\frac{\partial \ell}{\partial \theta_j}$ ,  $\frac{\partial \ell}{\partial \theta_i}$  и  $\frac{\partial^2 L}{\partial \theta_j \partial \theta_i}$ .
- л) Докажите, что  $\mathbb{E}\left(\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_i}\right) = -\mathbb{E}\left(\frac{\partial \ell}{\partial \theta_j} \frac{\partial \ell}{\partial \theta_i}\right)$ .
- м) Докажите, что  $I = -\mathbb{E}\left(\frac{\partial^2 \ell}{\partial \theta \partial \theta'}\right)$ .

## 2. ML в линейных моделях:

Можно смело считать первое упражнение сделанным, то есть использовать тот факт, что  $I = -E\left(\frac{\partial^2 \ell}{\partial \theta \partial \theta'}\right)$ .

Рассмотрим задачу линейной регрессии,  $y = X\beta + u$ , где  $u \sim \mathcal{N}(0; \sigma^2)$ . Для удобства определим  $x'_i$  —  $i$ -ую строку матрицы  $X$  и будем считать регрессоры неслучайными.

В нашем случае  $\theta = (\beta', \sigma^2)'$ .

Хинты для забывших матричное дифференцирование:  $\frac{\partial Ar}{\partial r} = A$ ,  $\frac{\partial^2 r' Ar}{\partial r \partial r'} = A + A'$ ,  $\frac{\partial r' Ar}{\partial r'} = Ar + A'r$ .

- а) Выпишите  $\ell(\theta)$  в виде суммы.
- б) Выпишите вектор  $s(\theta)$  в виде  $s(\theta) = \begin{pmatrix} ? \\ ? \end{pmatrix}$ , где первый элемент — это сразу вектор производных по всем  $\beta$  одним махом.
- в) Найдите ML оценки  $\hat{\theta}$ .
- г) Докажите, что  $L(\hat{\theta}) = a \cdot RSS^b$ , где  $a$  и  $b$  — некоторые константы. Забейте на  $a$  и найдите  $b$ .
- д) Найдите  $\frac{\partial^2 \ell}{\partial \theta \partial \theta'}$  в виде четырёх блоков:

$$\frac{\partial^2 \ell}{\partial \theta \partial \theta'} = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \beta \partial \beta'} & \frac{\partial^2 \ell}{\partial \beta \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial \sigma^2 \partial \beta'} & \frac{\partial^2 \ell}{\partial \sigma^2 \partial \sigma^2} \end{pmatrix}.$$

- е) В предыдущем пункте все блоки должны были получиться ненулевые. Однако найдите  $I$  и пара блоков занулятся :)
- ж) Найдите  $I^{-1}$ .

## 3. Выведите формулу для $R^2$ в регрессии вектора $y = X\beta + u$ .

## 4. LM-тест в линейных моделях.

Обозначим:  $\hat{\theta}_R$  и  $\hat{\theta}_{UR}$  — ограниченные и неограниченные экстремумы правдоподобия, а  $\hat{u}_R$  и  $\hat{u}_{UR}$  — соответствующие остатки.

Определим LM статистику как  $LM = s(\hat{\theta}_R)' \widehat{\text{Var}}^{-1}(s) s(\hat{\theta}_R)$ .

Будем считать первое упражнение сделанным, поэтому

$$LM = s(\hat{\theta}_R)' I^{-1}(\hat{\theta}_R) s(\hat{\theta}_R).$$



Также можно считать сделанным второе упражнение, поэтому:

$$s = \begin{pmatrix} \frac{X'u}{\sigma^2} \\ \frac{u'u - n\sigma^2}{2\sigma^4} \end{pmatrix}, \quad I^{-1} = \begin{pmatrix} \sigma^2(X'X)^{-1} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix}.$$

- а) Кстати, а чему равно  $s(\hat{\theta}_{UR})$ ?
  - б) Найдите  $s(\hat{\theta}_R)$  и  $I^{-1}(\hat{\theta}_R)$ .
  - в) Выведите формулу для  $LM$  статистики, содержащую только  $X$ ,  $\hat{u}_R$  и  $n$ .
  - г) Какую регрессию надо построить, чтобы  $R^2$  в ней оказался таким, что  $LM = nR^2$ ?
5. Исследовательница Елизавета оценила модель множественной регрессии,  $y = X\beta + u$ . Затем Елизавета проверяет гипотезу о незначимости отдельного коэффициента  $\beta_j$  двумя способами: через  $t$ -статистику и через  $F$ -статистику с ограниченной и неограниченной регрессией.

Докажите, что  $t^2 = F$ .

6. Василий обнаружил странную монетку и решил произвести над ней эксперименты. Выпадение орла он кодирует  $y_i = 1$ , решки —  $y_i = 0$ . При известном параметре  $p$ , наблюдаемые  $y_1, \dots, y_n$  независимы и имеют распределение Бернулли,  $y_i|p \sim \text{Bernoulli}(p)$ . Априорно по мнению Василия параметр  $p$  имеет бета-распределение,  $p \sim \text{beta}(\alpha, \beta)$ , где  $\alpha$  и  $\beta$  — некоторые неслучайные константы, описывающие мнения Василия.

Функция плотности бета-распределения имеет вид:

$$f(p) \propto p^{\alpha-1}(1-p)^{\beta-1}$$

Василий подкинул неизвестную монетку 100 раз и оказалось, что орёл выпал 70 раз и решка — 30.

- а) При каких  $\alpha$  и  $\beta$  априорное распределение совпадает с равномерным?
  - б) Найдите апостериорное распределение параметра  $p$ .
  - в) Найдите апостериорный прогнозный закон распределения для  $y_{n+1}$ .
  - г) Проинтерпретируйте смысл чисел  $\alpha$  и  $\beta$ .
7. В  $i$ -ый день Тимофей встречает  $y_i$  покемонов. Тимофей предполагает, что при известном параметре  $\lambda$ , наблюдаемые  $y_1, \dots, y_n$  независимы и имеют пуассоновское распределение,  $y_i|\lambda \sim \text{Pois}(\lambda)$ . Априорно по мнению Тимофея параметр  $\lambda$  имеет гамма-распределение,  $p \sim \text{Gamma}(\text{shape} = \alpha, \text{rate} = \beta)$ , где  $\alpha$  и  $\beta$  — некоторые константы, определяющие мнение Тимофея о встречаемости покемонов.

Функция плотности гамма-распределения имеет вид:

$$f(\lambda) \propto \lambda^{\alpha-1} \exp(-\lambda\beta)$$

За прошедшие 100 дней Тимофей встретил 70 покемонов.

- а) Найдите апостериорное распределение параметра  $\lambda$ .
  - б) Найдите апостериорный прогнозный закон распределения  $y_{n+1}$ .
  - в) Проинтерпретируйте смысл констант  $\alpha$  и  $\beta$ .
8. Андрей генерирует случайные величины  $X_i$  и  $Y_i$  по следующим принципам. Начинает Андрей с  $X_0 = 0$ . При  $i \geq 1$  Василий генерирует  $Y_i$  из нормального распределения  $Y_i|X_{i-1} \sim \mathcal{N}(0.5X_{i-1} + 2, 1)$ . Затем Андрей генерирует  $X_i$  из нормального распределения  $X_i|Y_i \sim \mathcal{N}(0.5Y_i + 4, 1)$ .

- а) Как в пределе распределена величина  $X_i$ ?
- б) Как в пределе распределена величина  $Y_i$ ?
9. Василий генерирует случайные величины  $X_i$  и  $Y_i$  по следующим принципам. Начинает Василий с  $X_0 = 0$ . При  $i \geq 1$  Василий генерирует  $Y_i$  из нормального распределения  $Y_i | X_{i-1} \sim \mathcal{N}(X_{i-1}, 1)$ . Затем Василий генерирует независимую от  $Y_i$  величину  $Z_i \sim \mathcal{N}(0; 1)$ . Если оказалось, что  $Z_i > Y_i$ , то Василий берёт  $X_i = 1$ , и  $X_i = 0$  иначе.
- а) Как в пределе распределена величина  $X_i$ ?
- б) Как в пределе распределена величина  $Y_i$ ?
10. Рассмотрим линейную модель  $y = X\beta + u$ , причем  $u \sim \mathcal{N}(0; \sigma^2 I)$ . Пусть априорно считается, что  $\beta \sim \mathcal{N}(0; \tau I)$ . Константы  $X$ ,  $\tau$  и  $\sigma^2$  известны. Для удобства все  $X$  и  $y$  центрированы.
- а) Найдите апостериорное распределение  $\beta$  с учётом наблюдаемых  $y$ .
- б) Найдите, при каком  $\lambda$  апостериорное среднее совпадёт с результатом гребневой регрессии
- $$\min_{\beta} \sum (y_i - \hat{y}_i)^2 + \lambda \sum \beta_j^2.$$
- в) Как надо изменить целевую функцию гребневой регрессии, чтобы результат совпал с апостериорным средним  $\beta$  при априорном распределении  $\beta \sim \mathcal{N}(0; \Sigma)$ ?

## 8.18. Кр 4, финальный экзамен, демо

### Часть 1. Тест.

### Часть 2. Задачи.

1. Величины  $X_i$  равномерны на отрезке  $[-a; 3a]$  и независимы. Есть несколько наблюдений,  $X_1 = 0.5$ ,  $X_2 = 0.7$ ,  $X_3 = -0.1$ .
- а) Найдите  $\mathbb{E}(X_i)$  и  $\mathbb{E}(|X_i|)$ ;
- б) Постройте оценку параметра  $a$  методом моментов, используя  $\mathbb{E}(|X_i|)$ ;
- в) Постройте оценку параметра  $a$  обобщённым методом моментов, используя моменты  $\mathbb{E}(X_i)$ ,  $\mathbb{E}(|X_i|)$  и взвешивающую матрицу

$$W = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}.$$

2. Рассмотрим логит-модель, задаваемую системой

$$\begin{cases} Y_i = \begin{cases} 1, & \text{если } Y_i^* \geq 0; \\ 0, & \text{иначе;} \end{cases} \\ Y_i^* = \beta_1 + \beta_2 X_i + u_i \end{cases}.$$

```
1 df_logit <- tibble(x = c(0, 2, 3, 4), y = c(0, 1, 0, 1))
2 logit <- glm(data = df_logit, family = binomial(link = "logit"), y ~ x)
3 xtable(logit)
```

---

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.9559	2.5571	-0.76	0.4443
x	0.8424	0.9329	0.90	0.3665

---

- а) Выпишите функцию правдоподобия для набора из четырёх наблюдений:  $(X_1, Y_1) = (4, 1)$ ,  $(X_2, Y_2) = (0, 0)$ ,  $(X_3, Y_3) = (2, 1)$ ,  $(X_4, Y_4) = (3, 0)$ .
- б) Оценки коэффициентов равны  $\hat{\beta}_1 = -1.95$  и  $\hat{\beta}_2 = 0.85$ . Оцените вероятность того, что  $Y_5 = 1$  при  $X_5 = 1$ .
3. Фирмы определяют необходимый запас товаров  $Y_i$  в зависимости от ожидаемых годовых продаж  $X_i^e$ , используя линейную форму зависимости  $Y_i = \beta_1 + \beta_2 X_i^e + \varepsilon_i$ . Исследователю доступны только данные о реальных продажах  $X_i = X_i^e + u_i$ , где ошибки  $u_i$  распределены независимо от  $X_i$  и удовлетворяют условию теоремы Гаусса–Маркова.
- а) Какие проблемы возникнут при оценке исходной модели с помощью МНК, если вместо данных по  $X_i^e$  будут использованы данные по  $X_i$ ?
- б) Каков возможный способ решения этих проблем?
4. Рассмотрим стационарный случайный процесс  $y_t$ , удовлетворяющий уравнению
- $$y_t = 3 + 0.7y_{t-1} - 0.1y_{t-2} + u_t,$$
- где  $u_t$  — белый шум с дисперсией 5.
- Найдите  $\mathbb{E}(y_t)$ ,  $\text{Var}(y_t)$ ,  $\text{Cov}(y_t, y_{t-1})$ ,  $\text{Cov}(y_t, y_{t-2})$ .
5. Что такое коинтегрированные временные ряды? Как проверить, являются ли два временных ряда коинтегрированными?
6. Модели панельных данных с фиксированными эффектами: определение, способы оценивания.

## 8.19. Кр 4, финальный экзамен

### Часть 1. Тест.

### Часть 2. Задачи.

1. Рассмотрим AR(2) процесс  $Y_t = 4 + Y_{t-1} - 0.4Y_{t-2} + u_t$ , где  $u_t$  — белый шум с единичной дисперсией.
- а) Является ли данный процесс стационарным?
- б) Найдите  $\text{Cov}(Y_t, Y_{t-1})$ ,  $\text{Cov}(Y_t, Y_{t-2})$ .
2. Начинающий исследователь Елисей исследует зависимость успехов в учёбе своих однокурсников,  $G_i$ , от времени, которое они тратят на учёбу,  $T_i$ . По выборке из 100 человек он смог оценить следующую регрессию:

$$\hat{G}_i = 30 + 6T_i$$

Елисей был бы рад полученному результату, но тут на лекции по эконометрике ему рассказали про эндогенность и пропущенные переменные, и он решил, что в его модели эти проблемы точно есть. Изучив литературу, он узнал, что на успехи в учёбе кроме времени влияют ещё и способности студента,  $A_i$ , при этом способности коррелированы со временем, которое студент тратит на учёбу.

- а) Проверьте, является ли найденная Елисеем оценка коэффициента при времени состоятельной;
- б) Если оценка не состоятельна, то предложите способ получения состоятельной оценки;

- в) Найдите асимптотическую величину смещения оценки, если  $\text{Cov}(G_i, A_i) = 6$ ,  $\text{Cov}(T_i, A_i) = 4$ ,  $\text{Var}(G_i) = 16$ ,  $\text{Var}(A_i) = 100$ ,  $\text{Var}(T_i) = 49$ .
3. Для определения, сколько земли следует фермеру отвести под клубнику, если ее будущие цены неизвестны, используется модель адаптивных ожиданий:

$$\begin{cases} A_t = \beta_1 + \beta_2 P_{t+1}^e + u_t \\ P_{t+1}^e - P_t^e = \lambda(P_t - P_t^e) \end{cases},$$

где  $A_t$  — количество акров, отведенное под клубнику в году  $t$ ,  $P_t$  — фактическая цена клубники, а  $P_t^e$  — ожидаемая цена клубники. Константа  $\lambda$  — коэффициент адаптации. Ошибки  $u_t$  удовлетворяют условию теоремы Гаусса-Маркова.

- а) Объясните, как исследователь перешёл от исходной модели к преобразованной модели  $A_t = \alpha_1 + \alpha_2 P_t + \alpha_3 A_{t-1} + v_t$ .
- б) Какие проблемы возникнут при оценивании коэффициентов преобразованной модели с помощью МНК? Как с ними справиться?
4. Рассмотрим систему одновременных уравнений

$$\begin{cases} c_t = \alpha_1 + \alpha_2 y_t + \alpha_3 c_{t-1} + u_{1t} \\ i_t = \beta_1 + \beta_2 r_t + \beta_3 y_t + u_{2t} \\ y_t = c_t + g_t + i_t \end{cases},$$

где  $c_t$  — потребление,  $i_t$  — инвестиции,  $y_t$  — ВНР,  $r_t$  — процентная ставка,  $g_t$  — правительственные расходы. Первые три переменные являются эндогенными.

- а) Возможно ли оценить коэффициенты данной системы уравнений и почему?
- б) Если возможно, то опишите последовательность Ваших действий.
5. Исследователь, используя данные по 870 индивидуумам, оценил вероятность получения степени бакалавра после четырехлетнего обучения в колледже в зависимости от обобщённых результатов тестов ASVABC. Переменная BACH равна 1, если индивидуум получил степень бакалавра, и равна 0 иначе. Исследователь оценил линейную модель с помощью МНК:

$$\widehat{BACH}_i = -0.8 + 0.02 ASVABC.$$

(0.04)      (0.001)

А также логит-модель:

$$\widehat{BACH^*}_i = -11.1 + 0.2 ASVABC,$$

(0.5)      (0.01)

где  $BACH_i = 1$  если  $BACH^*_i > 0$ .

- а) Как оценивается логит-модель?
- б) Каковы недостатки линейной модели в данном случае?
- в) Оцените предельный эффект объясняющего фактора для среднего значения ASVABC, равного 50.
6. Модели панельных данных со случайными эффектами: определение, способы оценивания.

## 8.20. ИП. БУЗА

1. Метод Наименьших Квадратов.
  - а) МНК-картинка
  - б) Нахождение всего-всего, если известен вектор  $y$  и матрица  $X$
2. Теорема Гаусса-Маркова
  - а) Формулировка с детерминистическими регрессорами
  - б) Доказательство с детерминистическими регрессорами
  - в) Формулировки со стохастическими регрессорами
  - г) Что даёт дополнительное предположение о нормальности  $\varepsilon$ ?
  - д) Теорема Фриша-Вау
  - е) Матрица-Мать всех регрессий
3. Проверка гипотез о линейных ограничениях
  - а) Проверка гипотезы о значимости коэффициента
  - б) Проверка гипотезы о значимости регрессии в целом
  - в) Проверка гипотезы об одном линейном соотношении с помощью ковариационной матрицы
  - г) Ограниченная и неограниченная модель
  - д) Тест Чоу на стабильность коэффициентов
  - е) Тест Чоу на прогнозную силу
4. Метод максимального правдоподобия
  - а) Свойства оценок
  - б) Два способа получения оценки дисперсии
  - в) Три теста (LM, Wald, LR)
  - г) Выписать функцию ML для обычной регрессии
  - д) для AR(1) процесса
  - е) для MA(1) процесса
  - ж) для логит модели
  - з) для пробит модели
  - и) для модели с заданным видом гетероскедастичности
5. Мультиколлинеарность
  - а) Определение, последствия
  - б) Величины, измеряющие силу мультиколлинеарности
  - в) Методы борьбы
  - г) Сюда же: метод главных компонент, хотя он используется и для других целей
6. Гетероскедастичность
  - а) Определение, последствия
  - б) Тесты, график
  - в) Стьюдентизированные остатки
  - г) НС оценки ковариации
  - д) GLS и FGLS
7. Временные ряды

- а) Стационарный временной ряд
  - б) ACF, PACF
  - в) Модель ARMA
  - г) ARIMA-SARIMA
8. Логит и пробит
- а) Описание моделей
  - б) Предельные эффекты
  - в) Чувствительность, специфичность
  - г) Кривая ROC — смотрим лекции :)
9. Эндогенность
- а) Три примера: одновременность, пропущенные переменные, ошибки измерения
  - б) IV, двухшаговый МНК
10. Модели панельных данных — смотрим лекции :)
- а) RE, FE, сквозная регрессии
  - б) Тест Хаусмана
11. Больше алгоритмов. Уметь объяснить суть метода. Уметь реализовать его.
- а) Классификационные деревья, случайный лес, xgboost
  - б) Гребневая регрессия (ridge regression)
  - в) LASSO
  - г) Квантильная регрессия
12. Байесовский подход
- а) Описать суть байесовского подхода.
  - б) Описать простенькую модель на языке STAN.
  - в) Быть готовым реализовать готовую модель с помощью пакета rstanarm в духе: байесовской линейной регрессии, байесовской логит модели.
13. R. Можно принести файл со своей заготовкой, можно пользоваться Интернетом для поиска информации, но не для общения. Примеры заданий:
- а) Загрузить данные из .csv файла в R
  - б) Посчитать описательные статистики: среднее, мода, медиана и т.д.
  - в) Построить подходящие описательные графики для переменных
  - г) Оценить линейную регрессию с помощью МНК. Провести диагностику на что-нибудь.
  - д) Оценить logit, probit модели, посчитать предельные эффекты
  - е) Оценить ARMA модель
  - ж) Выделить главные компоненты

## 9. 2017-2018

### 9.1. ИП, вспомнить всё!

1. Найдите длины векторов  $a = (1, 1, 1)$  и  $b = (1, 4, 6)$  и косинус угла между ними. Найдите длину проекции вектора  $b$  на вектор  $a$ .

2. Сформулируйте теорему Фалеса. Сформулируйте и докажите теорему Пифагора.

3. Для матрицы

$$A = \begin{pmatrix} 6 & 5 \\ 5 & 6 \end{pmatrix}$$

а) Найдите собственные числа и собственные векторы матрицы

б) Найдите  $\det(A)$ ,  $\text{tr}(A)$

в) Найдите собственные числа матрицы  $A^{2017}$ ,  $\det(A^{2017})$  и  $\text{tr}(A^{2017})$

4. Занудная халява: известно, что  $\text{Cov}(X, Y) = 5$ ,  $\text{Var}(X) = 16$ ,  $\text{Var}(Y) = 25$ ,  $\mathbb{E}(X) = 10$ ,  $\mathbb{E}(Y) = -5$ . Найдите  $\text{Cov}(X + 2Y, Y - X)$ ,  $\text{Var}(X + 2Y)$ ,  $\mathbb{E}(X + 2Y)$ .

5. Блондинка Маша 100 раз выходила на улицу и при этом 40 раз встретила динозавра. Постройте 95% доверительный интервал для вероятности встретить динозавра. На уровне 5% проверьте гипотезу о том, что данная вероятность равна 0.5 против альтернативной гипотезы об отличии данной вероятности от 0.5.

6. В кошельке 5 монеток, три золотых и две серебряных. Маша берёт наугад две монетки по очереди. Маше достались одинаковые монетки. Какова условная вероятность того, что обе золотые?

## 9.2. ИП, вспомнить всё!, ответы

1.  $|a| = \sqrt{3}$ ,  $|b| = \sqrt{53}$ ,  $\cos(a, b) = 11/\sqrt{159}$ , длина проекции равна  $11/\sqrt{3}$ .

2. Теорема Фалеса. Если параллельные прямые, пересекающие стороны угла, отсекают на одной его стороне равные отрезки, то они отсекают равные отрезки и на другой его стороне.

Теорема Пифагора. В прямоугольном треугольнике квадрат длины гипотенузы равен сумме квадратов длин катетов.

3. а)  $\lambda_1 = 11$ ,  $\lambda_2 = 1$ ,  $v_1 = (1 \ 1)'$ ,  $v_2 = (-1 \ 1)'$

б)  $\det(A) = 11$ ,  $\text{tr}(A) = 12$

в)  $\lambda_1 = 11^{2017}$ ,  $\lambda_2 = 1$ ,  $\det(A^{2017}) = 11^{2017}$ ,  $\text{tr}(A^{2017}) = 1 + 11^{2017}$

4.  $\text{Cov}(X + 2Y, Y - X) = 29$ ,  $\text{Var}(X + 2Y) = 136$ ,  $\mathbb{E}(X + 2Y) = 0$

5.  $\left[0.4 - 2\sqrt{\frac{0.4 \cdot 0.6}{100}}; 0.4 + 2\sqrt{\frac{0.4 \cdot 0.6}{100}}\right]$ , значение 0.5 не входит в доверительный интервал, значит, основная гипотеза отвергается.

6.  $3/4$

### 9.3. Контрольная 1, 26.10.2017

**Вопрос 1 ♣** Совместное распределение случайных величин  $X$  и  $Y$  задано с помощью таблицы

	$X = 3$	$X = 4$	$X = 5$
$Y = 3$	0.1	0.3	0.1
$Y = 4$	0.15	0.05	0.05
$Y = 6$	0.05	0.15	0.05

Математическое ожидание случайной величины  $Y$  при условии, что  $X = 3$ , равно

- ☐ A 2                                      ☐ D 3.4  
☐ B 2.4                                      ☒ 4  
☐ C 6                                        ☐ F Нет верного ответа.

**Вопрос 2 ♣** Оценка МНК коэффициента регрессии без свободного члена  $Y_i = \beta X_i + \varepsilon_i, i = 1, \dots, n$ , где  $x_i = X_i - \bar{X}, y_i = Y_i - \bar{Y}$ , находится по формуле

- ☐ A  $\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$                                       ☐ D  $\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$   
☒  $\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}$                                       ☐ E  $\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n y_i^2}$   
☐ C  $\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n Y_i^2}$                                       ☐ F Нет верного ответа.

**Вопрос 3 ♣** Оценка МНК неизвестного параметра  $\theta$  для модели  $Y_i = \theta X_{1i} + (1 + \theta)X_{2i} + \varepsilon_i, i = 1, \dots, n$  равна

- ☒  $\frac{\sum_{i=1}^n (X_{1i} + X_{2i})(Y_i - X_{2i})}{\sum_{i=1}^n (X_{1i} + X_{2i})^2}$                                       ☐ D  $\frac{\sum_{i=1}^n (X_{1i} + X_{2i})(Y_i - X_{2i})}{\sum_{i=1}^n (X_{1i} - Y_i)^2}$   
☐ B  $\frac{\sum_{i=1}^n (X_{1i} + X_{2i})(Y_i - X_{1i})}{\sum_{i=1}^n (X_{1i} + X_{2i})^2}$                                       ☐ E  $\frac{\sum_{i=1}^n (X_{1i} - Y_i)(Y_i - X_{2i})}{\sum_{i=1}^n (X_{1i} + X_{2i})^2}$   
☐ C  $\frac{\sum_{i=1}^n (X_{1i} + X_{2i})(Y_i - X_{2i})}{\sum_{i=1}^n (Y_i - X_{2i})^2}$                                       ☐ F Нет верного ответа.

**Вопрос 4 ♣** Для оцениваемой по 30 наблюдениям регрессии  $Y_i = \alpha + \beta X_i + \varepsilon_i, i = 1, \dots, n$  известны суммы  $\sum_{i=1}^{30} X_i = -15, \sum_{i=1}^{30} X_i^2 = 60, \sum_{i=1}^{30} X_i Y_i = 15, \sum_{i=1}^{30} Y_i = 75$ . Система нормальных уравнений для оценок коэффициентов регрессии  $\alpha, \beta$  методом наименьших квадратов равносильна системе

- ☐ A  $2\alpha - \beta = -1; \alpha - 4\beta = 5$                                       ☐ D  $4\alpha - 6\beta = 1; 6\alpha + 60\beta = 75$   
☒  $2\alpha - \beta = 5; \alpha - 4\beta = -1$                                       ☐ E  $30\alpha - 15\beta = 15; -15\alpha - 12\beta = 1$   
☐ C  $30\alpha + 15\beta = 75; 15\alpha + 60\beta = 15$                                       ☐ F Нет верного ответа.

**Вопрос 5 ♣** Для модели парной регрессии  $Y = \beta_1 I_n + \beta_2 X + \varepsilon$ , где  $Y = (Y_1, \dots, Y_n), X = (X_1, \dots, X_n), I_n = (1, \dots, 1), \varepsilon = (\varepsilon_1, \dots, \varepsilon_n), \hat{Y} = \hat{\beta}_1 I_n + \hat{\beta}_2 X, e = Y - \hat{Y}$  в пространстве  $\mathbb{R}^n$  ортогональны вектора

- ☐ A  $\varepsilon$  и  $\hat{Y}$                                       ☐ D  $X$  и  $\hat{Y}$   
☒  $e$  и  $I_n$                                       ☐ E  $Y$  и  $I_n$   
☐ C  $Y$  и  $\hat{Y}$                                       ☐ F Нет верного ответа.



**Вопрос 6 ♣** Эмманюэль и Владимир оценили зависимость стоимости подержанных Пежо (одной серии)  $Y$  от пробега  $X$  (измеряемого в км) с помощью модели парной регрессии  $Y = \alpha + \beta X + \varepsilon$  по одной и той же выборке, однако Эмманюэль измерял стоимость машин в евро, а Владимир – в рублях, 1 евро = 65 рублей. Оценки МНК коэффициента наклона регрессии, полученные Эмманюэлем  $\beta_E$  и Владимиром  $\beta_B$  связаны следующим образом:

☐ A  $\hat{\beta}_E = 4225\hat{\beta}_B$

☒  $\hat{\beta}_B = 65\hat{\beta}_E$

☐ B  $\hat{\beta}_E = 65\hat{\beta}_B$

☐ E  $\hat{\beta}_B = 4225\hat{\beta}_E$

☐ C  $\hat{\beta}_B = \hat{\beta}_E$

☐ F Нет верного ответа.

**Вопрос 7 ♣** При проверке гипотезы о значимости коэффициента линейной регрессии  $p$ -значение, соответствующее тестовой статистике, оказалось равным 0.07. Отсюда следует, что

☐ A длина 95% доверительного интервала для этого коэффициента меньше 0.07

☐ D соответствующий коэффициент значим при уровне значимости 1%

☒ соответствующий коэффициент не значим при уровне значимости 5%

☐ E длина 95% доверительного интервала для этого коэффициента равна 0.07

☐ C длина 95% доверительного интервала

☐ F Нет верного ответа.

**Вопрос 8 ♣** Для модели  $Y_i = \beta_1 + \beta_2 X_i + \varepsilon_i, i = 1, \dots, n, \varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$  тестовая статистика  $\frac{\hat{\beta}_2 - \beta_2^0}{se(\hat{\beta}_2)}$  имеет распределение

☐ A  $\chi_1^2$

☐ C  $\mathcal{N}(0, 1)$

☒  $t_{n-2}$

☐ B  $\chi_{n-2}^2$

☐ D  $\mathcal{N}(0, \sigma_\varepsilon)$

☐ F Нет верного ответа.

**Вопрос 9 ♣** С помощью  $t$ -теста проверяется гипотеза о том, что

☐ A надо ходить на семинары по эконометрике

равна единице

☒ коэффициент регрессии равен единице

☐ B оценка стандартной ошибки коэффициента регрессии равна единице

☐ E стандартная ошибка коэффициента регрессии равна единице

☐ C оценка коэффициента регрессии

☐ F Нет верного ответа.

**Вопрос 10 ♣** Для регрессии  $Y_i = \beta_1 + \beta_2 X_i + \varepsilon_i$ , оценённой по 30 наблюдениям с суммой квадратов остатков, равной 15, несмещенная оценка дисперсии случайной составляющей равна

☐ A 15/32

☐ D 13/30

☐ G Нет верного ответа.

☒ 15/28

☐ E 13/28

☐ C 0.5

☐ F 2

## Часть 2. Задачи.

- Докажите, что для модели парной регрессии  $Y_i = \beta_1 + \beta_2 X_i + \varepsilon_i$ , оцененной с помощью МНК, сумма остатков регрессии  $e_i = Y_i - \hat{Y}_i$  равна 0.
- Покажите, что для регрессий  $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$ ,  $\hat{X}_i = \hat{\alpha}_1 + \hat{\alpha}_2 Y_i$ , оценённых по одной и той же выборке  $(X_1, Y_1), \dots, (X_n, Y_n)$  коэффициенты детерминации  $R^2$  совпадают, а оценки коэффициентов наклона связаны соотношением  $\hat{\beta}_2 \hat{\alpha}_2 = R^2$ .
- Для классической регрессионной модели  $Y_i = \beta_1 + \beta_2 X_i + \varepsilon_i, i = 1, \dots, 20$  известно, что  $\sum_{i=1}^{20} X_i = 12, \sum_{i=1}^{20} X_i^2 = 12, \sum_{i=1}^{20} Y_i = 60, \sum_{i=1}^{20} Y_i^2 = 220, \sum_{i=1}^{20} X_i Y_i = 48$ . Найдите
  - $\hat{\beta}_1, \hat{\beta}_2$ ,
  - $TSS$ ,
  - $ESS$ ,
  - $\hat{\sigma}_\varepsilon^2$
- Для предыдущей задачи постройте точечный и 95% интервальный индивидуальный прогноз в точке  $X = 2$ .
- Заполните клетки с точками в приведенной ниже таблице. Клетки с XXX заполнять не надо.

Показатель	Значение
Multiple R	XXX
$R^2$	...
Standart error	XXX
Observations	60

ANOVA:

	df	SS	MS	F	Significance F
Regression	1	0.1			
Residual	59	0.4			
Total	60	...			

	Coef.	St. error	t-stat	Lower 95%	Upper 95%
Intercept	0.0045	0.015	XXX	XXX	XXX
MARKET	0.56	0.14	...	...	...

## 9.4. Контрольная 1, 26.10.2017, решения

- Воспользуемся свойством о том, что в парной регрессии точка  $(\bar{X}, \bar{Y})$  лежит на линии выборочной регрессии. Тогда

$$\begin{aligned} \sum_i e_i &= \sum_i (Y_i - \hat{\beta}_1 - \hat{\beta}_2 X_i) \\ \frac{1}{n} \sum_i e_i &= \bar{Y} - \hat{\beta}_1 - \hat{\beta}_2 \bar{X} = 0 \\ \sum_i e_i &= 0 \end{aligned}$$

2. МНК-оценки в регрессиях будут равны:

$$\hat{\beta}_2 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\hat{\alpha}_2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

Тогда их произведение равно

$$\begin{aligned}\hat{\beta}_2 \hat{\alpha}_2 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \cdot \sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2 \cdot \sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \cdot \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \cdot \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= \frac{s\text{Cov}^2(X, Y)}{s\text{Var}(X) s\text{Var}(Y)} = s\text{Corr}^2(X, Y) = R^2\end{aligned}$$

Поскольку  $s\text{Corr}^2(X, Y) = s\text{Corr}^2(Y, X)$ , коэффициенты детерминации в обеих регрессиях равны.

3. а) По формулам МНК-оценок получаем:

$$\hat{\beta}_2 = \frac{\sum_{i=1}^{20} X_i Y_i - \bar{X} \sum_{i=1}^{20} Y_i}{\sum_{i=1}^{20} (X_i - \bar{X})^2} = 2.5$$

$$\hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X} = 1.5$$

б)  $TSS = \sum_{i=1}^{20} (Y_i - \bar{Y})^2 = 40$

в)  $ESS = \sum_{i=1}^{20} (\hat{Y}_i - \bar{Y})^2 = 30$

г)  $\hat{\sigma}_\varepsilon^2 = \frac{RSS}{n-k} = \frac{5}{9}$

4. Точечный прогноз:

$$\hat{Y}_f = 1.5 + 2.5 \cdot 2 = 6.5$$

Для индивидуального прогноза понадобится:

$$\begin{aligned}\widehat{\text{Var}}(Y_f - \hat{Y}_f) &= \widehat{\text{Var}}(\beta_1 + \beta_2 X_f + \varepsilon_f - \hat{Y}_f) = \widehat{\text{Var}}(\varepsilon_f - \hat{Y}_f) = \widehat{\text{Var}}(\varepsilon_f) + \widehat{\text{Var}}(\hat{Y}_f) \\ &= \hat{\sigma}_{\varepsilon_f}^2 + \widehat{\text{Var}}(\hat{\beta}_1 + \hat{\beta}_2 \cdot 2) = \hat{\sigma}_{\varepsilon_f}^2 + \widehat{\text{Var}}(\hat{\beta}_1) + 4\widehat{\text{Var}}(\hat{\beta}_2) + 4\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) \\ &= \hat{\sigma}_{\varepsilon_f}^2 + \frac{\hat{\sigma}_{\varepsilon_f}^2 \sum_{i=1}^{20} X_i^2}{20 \sum_{i=1}^{20} (X_i - \bar{X})^2} + 4 \cdot \frac{\hat{\sigma}_{\varepsilon_f}^2}{\sum_{i=1}^{20} (X_i - \bar{X})^2} - 4 \cdot \frac{\bar{X} \hat{\sigma}_{\varepsilon_f}^2}{\sum_{i=1}^{20} (X_i - \bar{X})^2} \\ &\approx 0.81\end{aligned}$$

Из таблицы для  $t_{20-2}$  получаем  $t = 2.1$ . И интервальный прогноз примет вид:

$$[6.5 - 2.1\sqrt{0.81}; 6.5 + 2.1\sqrt{0.81}]$$

$$[4.61; 8.39]$$

5.  $R^2 = ESS/TSS = 0.2$

$$TSS = RSS + ESS = 0.5$$

$$t = (0.56 - 0)/0.14 = 4$$

Доверительный интервал:  $[0.56 - 2 \cdot 0.14; 0.56 + 2 \cdot 0.14] = [0.28; 0.84]$ , где  $t_{60-2; 0.975} = 2$ .

## 9.5. Контрольная 2, 26.10.2017

### Часть 1. Тест.

#### Вопрос 1 ♣

Рассмотрим модель  $Y = X\beta + \varepsilon$ . Условия теоремы Гаусса-Маркова выполнены, причём  $\text{Var}(\varepsilon_i) = \sigma_\varepsilon^2$ ,  $\hat{Y} = PY$ ,  $P = X(X'X)^{-1}X'$  и  $I$  - единичная матрица. Ковариационная матрица случайного вектора  $e = Y - \hat{Y}$  равна

- ☐ A  $\sigma_\varepsilon^2 I$ 
☐ D  $\sigma_\varepsilon^2 (P - I)$   
☒ B  $\sigma_\varepsilon^2 (I - P)$ 
☐ E  $\sigma_\varepsilon^2 P$   
☐ C  $\sigma_\varepsilon^2 (I + P)$ 
☐ F Нет верного ответа.

**Вопрос 2 ♣** Для регрессии  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$ , оцененной по 36 наблюдениям с  $R^2 = 0.9$ , значение тестовой статистики для проверки гипотезы об адекватности регрессии равно

- ☐ A 11/9
 ☒ B имеющимся данным
 ☐ C 96  
☐ D невозможно вычислить
 ☐ E по
 ☐ F Нет верного ответа.

#### Вопрос 3 ♣

Если в уравнение регрессии не включена константа, то

- ☒ A К этой модели применима теорема Гаусса-Маркова  
☐ B  $R_{adj}^2$  в этой модели всегда неотрицательный  
☐ C  $R^2$  является показателем качества подгонки регрессии  
☐ D Значимость коэффициентов регрессии нельзя проверять при помощи t-статистики  
☐ E Сумма остатков регрессии равна 0  
☐ F Нет верного ответа.

#### Вопрос 4 ♣

По 546 наблюдениям за 1987 г. оценили зависимость стоимости частных домов в канаде price (измеряемой в долларах США) от общей площади square (измеряемой в кв. м.) наличия подъездного пути driveway (1 — если есть, 0 — если нет):

переменная	коэффициент	ст. ошибка	t-статистика	P-значение
square	2.724	2.15	1.27	0.206
driveway	-922.563	8602.312	-0.11	0.915
square*driveway	3.479	1.43	2.42	0.016
const	38731.07	8156.39	4.75	0.000

Согласно полученным результатам, при уровне значимости 5%, наличие подъездного пути увеличивает стоимость каждого квадратного метра жилья на

- ☒ A 3.479 \$
 ☐ B 2.724 \$
 ☐ C -922.563 \$  
☐ D 0
 ☐ E 6.203 \$
 ☐ F Нет верного ответа.

### Вопрос 5 ♣

Оценена зависимость расходов потребителей на газ и электричество  $Y$  в США в 1977-1999 г. в постоянных ценах I квартала 1977 г. от времени ( $t = 1$  для 1977,  $t = 2$  для 1978 и т.д.) с учётом сезонных факторов ( $D_i = 1$ , если наблюдение относится к  $i$ -ому кварталу и 0 иначе,  $i = 1, \dots, 4$ ):

$$\hat{Y} = 8 + 0.1t - 3D_2 - 2.6D_3 - 2D_4$$

Если в качестве базовой категории будет принят не первый квартал, а третий, уравнение регрессии примет вид

☐  $\hat{Y} = 8 + 0.1t + 3D_1 + 2.6D_2 + 2D_4$

☐  $\hat{Y} = 8 + 0.1t - 3D_1 - 2.6D_2 - 2D_4$

☐  $\hat{Y} = 5.4 + 0.1t - 3D_1 - 2.6D_2 - 2D_4$

☒  $\hat{Y} = 5.4 + 0.1t + 2.6D_1 - 0.4D_2 + 0.6D_4$

☐  $\hat{Y} = 5.4 + 0.1t - 3D_1 - 0.4D_2 - D_4$

☐ Нет верного ответа.

### Вопрос 6 ♣

Для выбора между линейной и полулогарифмической моделями (где EARNINGS — почасовая заработная плата в \$, S — длительность обучения, ASVABC - результаты тестов, характеризующие успеваемость) был проведен тест Дэвидсона, Уайта и МакКиннона и получены следующие результаты:

	Зависимая: $Y$	Зависимая: $\ln Y$
(Intercept)	−26.148 (4.17)	−1.941 (3.2499)
S	2.008 (0.276)	0.087 (0.035)
ASVABC	0.393 (0.079)	0.017 (0.007)
lin_add	−15.373 (5.984)	
semilog_add		−0.029 (0.065)
$R^2$	0.2071	0.2212
F	46.59	50.74
Adj. $R^2$	0.2027	0.2168
Num. obs.	540	540
RSS	90975.57	148.1
$\hat{\sigma}$	13.04	0.5256

Где  $\text{lin\_add} = \ln(\hat{Y}) - \hat{\ln Y}$ ,  $\text{semilog\_add} = \hat{Y} - \exp(\hat{\ln Y})$  и в скобках указаны стандартные ошибки.

На уровне значимости 5% можно сделать вывод, что

☐ Между линейной и полулогарифмической моделями нет статистической разницы

☐ Лучше линейная модель

☒ Лучше полулогарифмическая модель

☐ Лучше линейная в логарифмах модель

☐ Невозможно выбрать лучшую

☐ Нет верного ответа.

### Вопрос 7 ♣

По данным для 27 фирм была оценена зависимость выпуска  $Y$  от труда  $L$  и капитала  $K$  с помощью моделей:

$$\ln Y_i = b_1 + b_2 \ln L_i + b_3 \ln K_i + \varepsilon_i \quad (1)$$

$$\ln Y_i = b_1 + b_2(\ln L_i + \ln K_i) + \varepsilon_i \quad (2)$$

Суммы квадратов остатков в этих моделях известны,  $RSS_1 = 8$  и  $RSS_2 = 10$ .  $F$ -статистика для проверки гипотезы о равенстве эластичностей по труду и по капиталу равна

☐ A 12

☐ C 4

☒ 6

☐ B 2

☐ D 8

☐ F Нет верного ответа.

### Вопрос 8 ♣

По одним и тем же наблюдениям оценили две регрессии:  $\hat{Y} = 1 + 3X_1$  и  $\hat{Y} = 2 + 5X_2$ . Известно, что  $\widehat{\text{Cov}}(X_1, X_2) > 0$ . Оценки МНК коэффициентов регрессии  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$ :

☐ A  $\hat{\beta}_0 = 3, \hat{\beta}_1 = 3, \hat{\beta}_2 = 5$

☐ D  $\hat{\beta}_0 = 1.5, \hat{\beta}_1 = 3, \hat{\beta}_2 = 5$

☐ B  $\hat{\beta}_1, \hat{\beta}_2$  найти невозможно,  $\hat{\beta}_0 = 3$

☒ оценки коэффициентов невозможно найти по имеющимся данным

☐ C  $\hat{\beta}_0$  найти невозможно,  $\hat{\beta}_1 = 3, \hat{\beta}_2 = 5$

☐ F Нет верного ответа.

### Вопрос 9 ♣

По данным для 27 фирм исследована зависимость прибыли  $Y$  от числа работников  $X$  вида

$$Y = \beta_0 + \beta_1 X + \varepsilon \text{ и получено } \hat{\beta}_0 = 8, \hat{\beta}_1 = 2, \hat{s}^2 = 25 \text{ и матрица } (X'X)^{-1} = \begin{pmatrix} 0.36 & -0.03 \\ -0.03 & 0.09 \end{pmatrix}.$$

95% доверительный интервал для  $\beta_1$ :

☐ A [-0.94; 4.94]

☐ D невозможно

☐ E [0.04; 3.96]

☒ [-1.09; 5.09]

вычислить по имеющимся данным

☐ C [1.82; 14.18]

☐ F Нет верного ответа.

### Вопрос 10 ♣ Исследователь Борис оценил параметры нескольких моделей:

Модель	Уравнение
1	$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + u$
2	$\ln Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + u$
3	$Y = \beta_1 + \beta_4 X_4 + \beta_5 X_5 + u$
4	$Y/X_2 = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + u$

С помощью  $R_{adj}^2$  можно выбрать лучшую из пар моделей

☒ 1 и 3

☐ C 2 и 4

☐ E 2 и 3

☐ B 1 и 2

☐ D 1 и 4

☐ F Нет верного ответа.

## Часть 2. Задачи.

1. По 29 наблюдениям оценили функцию спроса на яблоки

$$\widehat{\ln Q} = 14 - 6 \ln P_{apple} + 3 \ln P_{orange} + 2 \ln P_{banana},$$

где  $Q$  — спрос на яблоки, а  $P_{apple}$  — цена яблок,  $P_{orange}$  — цена апельсинов и  $P_{banana}$  — цена бананов.

Известна оценка ковариационной матрицы коэффициентов регрессии:

$$\widehat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} 1 & 0.1 & -0.2 & 0.3 \\ 0.1 & 2 & 0.5 & 0.7 \\ -0.2 & 0.5 & 3 & 0.6 \\ 0.3 & 0.7 & 0.6 & 4 \end{pmatrix}$$

На уровне значимости 5% проверьте гипотезу о том, что  $\beta_{orange} = \beta_{banana}$ .

2. По ежегодным данным за 25 лет была оценена зависимость расходов на жилье  $Y$  от доходов индивидуумов  $I$  и относительного индекса цен  $P$  с помощью трёх моделей:

$$\hat{Y} = -39.9 + 0.179P + 0.113I, R^2 = 0.987$$

(46.9)                      (0.009)                      (0.409)

$$\hat{Y} = -27.03 + 0.177I, R^2 = 0.987$$

(3.34)                      (0.004)

$$\hat{Y} = 813.3 + -7.08P, R^2 = 0.76$$

(24.2)                      (0.019)

Известно, что  $\widehat{\text{Corr}}(P, I) = -0.88$ , в скобках указаны стандартные отклонения коэффициентов.

Какую модель Вы предпочтёте и почему?

3. Для регрессии в отклонениях  $y = \beta_1 x + \beta_2 z + \varepsilon$ , оцененной по 100 наблюдениям, известны следующие суммы:

$$\sum y_i^2 = 400, \sum x_i^2 = 30, \sum z_i^2 = 3, \sum x_i y_i = 30, \sum z_i y_i = 24, \sum x_i z_i = 0$$

Найдите оценки МНК коэффициентов  $\beta_1, \beta_2$  и коэффициент детерминации  $R^2$ .

4. Приведены результаты оценки зависимости логарифма арендной платы жилья в России,  $\ln \text{PRICE}$ , от общей площади,  $\text{GENSQUARE}$  (кв. м.), наличия газа,  $\text{GAS}$ , (1 — есть, 0 — нет), наличия телефона,  $\text{PHONE}$  (1 — есть, 0 — нет).

Модели 1 и 2 оценены для городов с численностью населения более миллиона человек, модель 3 — для городов с численностью от полумиллиона до миллиона, модель 4 — для обеих выборок. В скобках в таблице приведены стандартные ошибки.

- Проверьте значимость коэффициентов в модели 2 при уровне значимости 5% и дайте интерпретацию полученным результатам.
- Проверьте гипотезу о равенстве 0.01 коэффициента при переменной  $\text{GENSQUARE}$  при уровне значимости 0.05 в модели 2.
- Можно ли утверждать, что наличие газа и телефона в крупном городе не влияет на стоимость аренды? Ответ обоснуйте формулировкой и проверкой подходящей гипотезы.
- Опираясь на результаты оценки моделей 2, 3 и 4, можно ли утверждать, что зависимость стоимости жилья от рассмотренных выше факторов одинакова для городов с численностью населения более миллиона человек и с численностью от полумиллиона до миллиона? Ответ обоснуйте подходящим тестом.

	Модель 1	Модель 2	Модель 3	Модель 4
(Intercept)	7.168 (0.1334)	6.884 (0.2398)	7.459 (0.059)	7.448 (0.0569)
GENSQUARE	0.012 (0.002687)	0.012 (0.0027)	0.0064 (0.000913)	0.00683 (0.00086)
GAS		0.119 (0.166)	-0.262 (0.038)	-0.272 (0.0354)
PHONE		0.197 (0.123)	0.372 (0.042)	0.3505 (0.0399)
$R^2$	0.1330	0.1541	0.177	0.1757
F	19.95	7.78	67.37	75.75
Adj. $R^2$	0.1264	0.1343	0.175	0.1734
Num. obs.	132	132	938	1070
RSS	23.05	22.494	309.444	334.846
$\hat{\sigma}$	0.42113	0.41921	0.575	0.956

## 9.6. Контрольная 2, 26.10.2017, решения

1. Для проверки гипотезы используем  $t$ -статистику:

$$t_{obs} = \frac{\hat{\beta}_{orange} - \hat{\beta}_{banana} - \beta_{orange} - \beta_{banana}}{se(\hat{\beta}_{orange} - \hat{\beta}_{banana})} \sim t_{n-k} = t_{29-4} = t_{25}$$

Сначала найдём стандартное отклонение:

$$\begin{aligned} se(\hat{\beta}_{orange} - \hat{\beta}_{banana}) &= \sqrt{\widehat{\text{Var}}(\hat{\beta}_{orange}) + \widehat{\text{Var}}(\hat{\beta}_{banana}) - 2\widehat{\text{Cov}}(\hat{\beta}_{orange}, \hat{\beta}_{banana})} \\ &= \sqrt{3 + 4 - 2 \cdot 0.6} \approx 2.41 \end{aligned}$$

Тогда наблюдаемое значение статистики равно  $t_{obs} = \frac{3-4}{2.41} \approx -0.42$ , а  $t_{crit} \approx 2.06$ . Значит, оснований отвергать нулевую гипотезу нет.

2. В модели (1)  $\widehat{\text{Corr}}(P, I) = -0.8$  свидетельствует о проблеме мультиколлинеарности в данных. В модели (3) значение  $R^2$  самое низкое. Предпочтительной является модель (2).
- 3.

$$\begin{aligned} \hat{\beta}_1 &= \frac{s\text{Cov}(X, Y)}{s\text{Var}(X)} = \frac{30}{30} = 1 \\ \hat{\beta}_2 &= \frac{s\text{Cov}(Z, Y)}{s\text{Var}(Z)} = \frac{20}{3} \end{aligned}$$

Для нахождения коэффициента детерминации воспользуемся факторным разложением:

$$R^2 = \hat{\beta}_1 \frac{s\text{Cov}(X, Y)}{s\text{Var}(Y)} + \hat{\beta}_2 \frac{s\text{Cov}(Z, Y)}{s\text{Var}(Y)} = 1 \cdot \frac{30}{\frac{493}{3}} + \frac{20}{3} \cdot \frac{20}{\frac{493}{3}} = \frac{490}{493}$$

4. а)  $n = 132, k = 3 \Rightarrow n - k = 128$ , в таблице находим, что при  $\alpha = 0.05$   $t_{crit} = 1.98$ .

$$t_{obs} = \frac{6.884-0}{0.2398} = 28.7 > t_{crit} \Rightarrow \hat{\beta}_{Intercept} \text{ значим.}$$

$$t_{obs} = \frac{0.012-0}{0.0027} = 4.4 > t_{crit} \Rightarrow \hat{\beta}_{GENSQUARE} \text{ значим.}$$

$$t_{obs} = \frac{0.119-0}{0.166} = 0.72 < t_{crit} \Rightarrow \hat{\beta}_{GAS} \text{ незначим.}$$

$$t_{obs} = \frac{0.197-0}{0.123} = 1.6 < t_{crit} \Rightarrow \hat{\beta}_{PHONE} \text{ незначим.}$$

При прочих равных общая площадь оказывает влияние на логарифм арендной платы, а наличие газа и наличие телефона — нет.



б) При верной  $H_0$  и  $\alpha = 0.1$   $t_{crit} = 1.658$ .

$$t_{obs} = \frac{0.199 - 0.02}{0.166} = -0.596$$

Поскольку  $|t_{obs}| < t_{crit}$ , оснований отвергать нулевую гипотезу нет.

в)

$$H_0 : \begin{cases} \beta_{GAS} = 0 \\ \beta_{PHONE} = 0 \end{cases} \quad H_a : \beta_{GAS}^2 + \beta_{PHONE}^2 > 0$$

Выпишем ограниченную и неограниченную модели:

$$R : \ln PRICE_i = \beta_{Intercept} + \beta_{GENSQUARE} GENSQUARE_i + \varepsilon_i$$

$$UR : \ln PRICE_i = \beta_{Intercept} + \beta_{GENSQUARE} GENSQUARE_i + \beta_{GAS} GAS_i + \beta_{PHONE} PHONE_i + \varepsilon_i$$

Из условия находим:  $RSS_R = 23.05$ ,  $RSS_{UR} = 22.494$ .

При верной  $H_0$   $F \sim F_{2,128}$  и  $F_{crit} \approx 3.8$  ( $\alpha = 0.05$ ).

$$F_{obs} = \frac{(RSS_R - RSS_{UR})/q}{RSS_{UR}/(n - k_{UR})} = \frac{(23.05 - 22.494)/2}{22.494/(132 - 4)} = 1.58$$

Так как  $F_{obs} < F_{crit}$ , нет оснований отвергать  $H_0$ .

г) Чтобы записать неограниченную модель, введём вспомогательную переменную

$$d_i = \begin{cases} 1, & \text{если город крупный} \\ 0, & \text{иначе} \end{cases}$$

Тогда неограниченная модель примет вид:

$$\ln PRICE_i = (\beta_{Intercept} + \Delta_1 d_i) + (\beta_{GENSQUARE} + \Delta_2 d_i) GENSQUARE_i + (\beta_{GAS} + \Delta_3 d_i) GAS_i + (\beta_{PHONE} + \Delta_4 d_i) PHONE_i + \varepsilon_i$$

$$RSS_{UR} = 22.494 + 616.709 = 639.203$$

Выпишем также ограниченную модель:

$$\ln PRICE_i = \beta_{Intercept} + \beta_{GENSQUARE} GENSQUARE_i + \beta_{GAS} GAS_i + \beta_{PHONE} PHONE_i + \varepsilon_i$$

$$RSS_R = 710.21$$

Проверятся следующая гипотеза:

$$H_0 : \begin{cases} \Delta_1 = 0 \\ \Delta_2 = 0 \\ \Delta_3 = 0 \\ \Delta_4 = 0 \end{cases} \quad H_a : \Delta_1^2 + \Delta_2^2 + \Delta_3^2 + \Delta_4^2 > 0$$

При верной  $H_0$   $F \sim F_{4,1173}$  и  $F_{crit} \approx 2.8$  ( $\alpha = 0.05$ ).

$$F_{obs} = \frac{(RSS_R - RSS_{UR})/q}{RSS_{UR}/(n - k_{UR})} = \frac{(710.21 - 639.203)/4}{639.203/(1181 - 8)} \approx 32.58$$

Поскольку  $F_{obs} > F_{crit}$ , основная гипотеза отвергается, а значит, зависимость нельзя считать единой.

## 9.7. Кр 3, 2018-03-28, бп часть

### 9.7.1. Тест

### 9.7.2. Задачи

1. На основании наблюдений получена МНК оценка уравнения регрессии  $\hat{Y}_i = 0.2Z_i + 0.3W_i$  и оценка дисперсии ошибок  $\hat{\sigma}^2 = 0.04$ . Матрица наблюдений регрессоров имеет вид

$$X^T = \begin{pmatrix} 1 & 2 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 4 & 5 & 6 \end{pmatrix}.$$

Ошибки имеют нормальное распределение.

Постройте 95% предиктивный интервал (доверительный интервал для индивидуального прогноза) в точке  $Z = -2$ ,  $W = 5$ .

2. В модели множественной регрессии  $Y = X\beta + \varepsilon$  выполнены все предпосылки классической линейной модели кроме предпосылки о гомоскедастичности. Вектор ошибок имеет нормальное распределение, а возможная гетероскедастичность имеет вид

$$\text{Var}(\varepsilon_i) = \begin{cases} \sigma_1^2, & \text{при } i \leq m; \\ \sigma_2^2, & \text{при } i > m. \end{cases}$$

Матрица  $X$  имеет размер  $n$  на  $k + 1$ .

Выведите формулу статистики LR-теста для проверки гипотезы о гомоскедастичности.

3. Рассмотрим модель  $Y_i = \beta X_i + \varepsilon_i$ , где  $\varepsilon_i$  — независимые случайные величины с  $\mathbb{E}(\varepsilon_i) = 0$  и  $\text{Var}(\varepsilon_i) = 2018i$ .

Найдите наиболее эффективную оценку для параметра  $\beta$  в классе всех линейных по  $Y$  несмещённых оценок.

4. По 1000 наблюдений Винни-Пух оценил логистическую модель  $\mathbb{P}(Y_i = 1) = F(\beta_0 + \beta_1 X_i)$ , где  $X_i$  — количество времени в часах, проведённое в гостях, а  $Y_i$  — факт застревания при выходе.

Оценки параметров равны  $\hat{\beta}_0 = 2$ ,  $\hat{\beta}_1 = 3$ , с оценкой ковариационной матрицы

$$\begin{pmatrix} 0.25 & 0.1 \\ 0.1 & 0.16 \end{pmatrix}.$$

- а) Проверьте значимость отдельных коэффициентов при уровне значимости 5%;
- б) Найдите предельный эффект времени, проведённого в гостях, на вероятность застрять при выходе для получасового визита;
- в) Найдите максимально возможный предельный эффект.

## 9.8. Кр 3, 2018-03-28, бп часть, решения

1. Найдём оценку ковариационную матрицу оценок коэффициентов:

$$\widehat{\text{Var}}(\hat{\beta}) = \hat{\sigma}^2 (X'X)^{-1} = 0.04 \begin{pmatrix} 1/14 & 0 \\ 0 & 1/77 \end{pmatrix},$$

дисперсию ошибки прогноза:

$$\begin{aligned}
\widehat{\text{Var}}(y_i - \hat{y}_f | X) &= \widehat{\text{Var}}(\beta_z z_i + \beta_w w_i + \varepsilon_i - \hat{y}_f | X) \\
&= \widehat{\text{Var}}(\varepsilon_i | X) + \widehat{\text{Var}}(-2\hat{\beta}_z + 5\hat{\beta}_w | X) - 2\widehat{\text{Cov}}(\varepsilon_i, \hat{y}_f) \\
&= 0.04 + 4\widehat{\text{Var}}(\hat{\beta}_z | X) + 25\widehat{\text{Var}}(\hat{\beta}_w | X) - 2 \cdot (-2) \cdot 5\widehat{\text{Cov}}(\hat{\beta}_z, \hat{\beta}_w) + 0 \\
&= 0.04 + 4 \cdot 0.04 \cdot \frac{1}{14} + 25 \cdot 0.04 \cdot \frac{1}{77} + 0 \\
&\approx 0.0644,
\end{aligned}$$

и сам прогноз:

$$\hat{y}_f = 0.2 \cdot (-2) + 0.3 \cdot 5 = 1.1.$$

Теперь можно выписать доверительный интервал,  $t_{0.975,5} = 2.57$ :

$$\left[ 1.1 - 2.57\sqrt{0.0644}; 1.1 + 2.57\sqrt{0.0644} \right]$$

2. Функция правдоподобия в неограниченной модели ( $\sigma_1^2 \neq \sigma_2^2$ ) имеет вид:

$$L(\sigma_1^2, \sigma_2^2) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2} \frac{(y_i - x'_i \beta)^2}{\sigma_1^2}} \cdot \prod_{i=m+1}^n \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{1}{2} \frac{(y_i - x'_i \beta)^2}{\sigma_2^2}}$$

Выпишем также логарифмическую функцию правдоподобия для неограниченной модели и найдём оценки  $\hat{\sigma}_1^2, \hat{\sigma}_2^2$ :

$$\begin{aligned}
\ell(\sigma_1^2, \sigma_2^2) &= -\frac{m}{2} \ln 2\pi - \frac{m}{2} \ln \sigma_1^2 - \frac{1}{2\sigma_1^2} \sum_{i=1}^m (y_i - x'_i \beta)^2 \\
&\quad - \frac{n-m}{2} \ln 2\pi - \frac{n-m}{2} \ln \sigma_2^2 - \frac{1}{2\sigma_2^2} \sum_{i=m+1}^n (y_i - x'_i \beta)^2 \rightarrow \max_{\sigma_1^2, \sigma_2^2} \\
\frac{\partial \ell}{\partial \sigma_1^2} &= -\frac{m}{2\sigma_1^2} + \frac{1}{2(\sigma_1^2)^2} \sum_{i=1}^m (y_i - x'_i \beta)^2 \Big|_{\sigma_1^2 = \hat{\sigma}_1^2} = 0 \\
\frac{\partial \ell}{\partial \sigma_2^2} &= -\frac{n-m}{2\sigma_2^2} + \frac{1}{2(\sigma_2^2)^2} \sum_{i=m+1}^n (y_i - x'_i \beta)^2 \Big|_{\sigma_2^2 = \hat{\sigma}_2^2} = 0 \\
\hat{\sigma}_1^2 &= \frac{\sum_{i=1}^m (y_i - x'_i \beta)^2}{m} \\
\hat{\sigma}_2^2 &= \frac{\sum_{i=m+1}^n (y_i - x'_i \beta)^2}{n-m}
\end{aligned}$$

Тогда логарифмическая функция правдоподобия примет вид:

$$\ell_{UR}(\hat{\sigma}_1^2, \hat{\sigma}_2^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} - \frac{m}{2} \ln \frac{\sum_{i=1}^m (y_i - x'_i \beta)^2}{m} - \frac{n-m}{2} \ln \frac{\sum_{i=m+1}^n (y_i - x'_i \beta)^2}{n-m}$$

Прделаем то же самое для ограниченной модели ( $\sigma_1^2 = \sigma_2^2 = \sigma_0^2$ ):

$$L(\sigma_0^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2} \frac{(y_i - x'_i\beta)^2}{\sigma_0^2}}$$

$$\ell(\sigma_0^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma_0^2 - \frac{1}{2\sigma_0^2} \sum_{i=1}^n (y_i - x'_i\beta)^2 \rightarrow \max_{\sigma_0^2}$$

$$\left. \frac{\partial \ell}{\partial \sigma_0^2} = -\frac{n}{2\sigma_0^2} + \frac{1}{2(\sigma_0^2)^2} \sum_{i=1}^n (y_i - x'_i\beta)^2 \right|_{\sigma_0^2 = \hat{\sigma}_0^2} = 0$$

$$\hat{\sigma}_0^2 = \frac{\sum_{i=1}^n (y_i - x'_i\beta)^2}{n}$$

$$\ell_R(\hat{\sigma}_0^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \frac{\sum_{i=1}^n (y_i - x'_i\beta)^2}{n} - \frac{n}{2}$$

Осталось выписать формулу статистики LR-теста:

$$LR = 2(\ell_{UR} - \ell_R)$$

$$= 2 \left( -\frac{m}{2} \ln \frac{\sum_{i=1}^m (y_i - x'_i\beta)^2}{m} - \frac{n-m}{2} \ln \frac{\sum_{i=m+1}^n (y_i - x'_i\beta)^2}{n-m} + \frac{n}{2} \ln \frac{\sum_{i=1}^n (y_i - x'_i\beta)^2}{n} \right)$$

3. Наиболее эффективная оценка коэффициента  $\beta$  может быть получена с помощью взвешенного МНК. Для этого необходимо оценить регрессию

$$\frac{y_i}{\sqrt{i}} = \beta \frac{x_i}{\sqrt{i}} + \frac{\varepsilon_i}{\sqrt{i}}.$$

Переобозначив  $\frac{y_i}{\sqrt{i}}$  за  $y_i^*$ ,  $\frac{x_i}{\sqrt{i}}$  за  $x_i^*$  и  $\frac{\varepsilon_i}{\sqrt{i}}$  за  $\varepsilon_i^*$ , получим регрессию

$$y_i^* = \beta x_i^* + \varepsilon_i^*,$$

применив к которой обычный МНК, получим эффективную оценку  $\hat{\beta}_{WLS}$  вида

$$\hat{\beta}_{WLS} = \frac{\sum_{i=1}^n x_i^* y_i^*}{\sum_{i=1}^n (x_i^*)^2} = \frac{\sum_{i=1}^n \frac{x_i}{\sqrt{i}} \cdot \frac{y_i}{\sqrt{i}}}{\sum_{i=1}^n \frac{x_i^2}{i}}.$$

4. а)  $\beta_0 : \frac{2-0}{\sqrt{0.25}} = 4 > 2 \Rightarrow$  гипотеза о незначимости коэффициента отвергается  
 $\beta_1 : \frac{3-0}{\sqrt{0.16}} = 7.5 > 2 \Rightarrow$  гипотеза о незначимости коэффициента отвергается  
 б)

$$\begin{aligned} \frac{\partial \hat{\mathbb{P}}(Y_i = 1)}{\partial X_i} &= F'(\hat{\beta}_0 + \hat{\beta}_1 \cdot X_i) \cdot \hat{\beta}_1 \\ &= F(\hat{\beta}_0 + \hat{\beta}_1 \cdot X_i) (1 - F(\hat{\beta}_0 + \hat{\beta}_1 \cdot X_i)) \cdot \hat{\beta}_1 \Big|_{X_i=0.5} \\ &= \frac{e^{3.5}}{1 + e^{3.5}} \left( 1 - \frac{e^{3.5}}{1 + e^{3.5}} \right) \cdot 3 \\ &\approx 0.085 \end{aligned}$$

- в) Предельный эффект максимален в точке, где наклон касательной к логистической функции самый крутой, то есть в нуле:

$$2 + 3x = 0 \Rightarrow x = -\frac{2}{3}$$

## 9.9. КР 3, 2018-03-28, ип часть

Ровно 42 года назад польская путешественница Кристина Хойновская-Лискевич начала первое женское одиночное кругосветное плавание на парусной яхте. Плавание продлилось примерно два года.

1. Рассмотрим задачу линейной регрессии для случая идеально точно известной дисперсии:  $y = X\beta + u$ ,  $u \sim \mathcal{N}(0; I_{n \times n})$ , где  $I_{n \times n}$  — единичная матрица. Обозначим  $s(\beta)$  — вектор-столбец, градиент логарифмической функции правдоподобия, а  $\hat{\beta}$  — оценку параметров  $\beta$  с помощью максимального правдоподобия.
  - а) Найдите  $\text{Var}(s(\beta)|X)$  и вспомните  $\text{Var}(\hat{\beta}|X)$ ;
  - б) Упростите выражение  $\text{Var}(s(\beta)|X) \cdot \text{Var}(\hat{\beta}|X)$ ;
2. В созвездии Малой Медведицы водится  $k$  видов медведепришельцев. Исследователь Миша отлавливает  $n$  медведепришельцев и классифицирует их по видам:  $y_1$  — количество медведепришельцев первого вида,  $y_2$  — второго, ...,  $y_k$  —  $k$ -го. Миша хочет оценить вектор вероятностей  $p = (p_1, \dots, p_{k-1})$ . Вероятностей на одну меньше, чем видов, чтобы избежать жёсткой линейной зависимости между ними.
  - а) Как распределена в теории величина  $y_1$ ? Чему равна её дисперсия?
  - б) Как распределена в теории величина  $y_{12} = (y_1 + y_2)$ ? Чему равна её дисперсия?
  - в) Чему равна ковариация  $y_1$  и  $y_2$ ?
  - г) Выпишите функцию правдоподобия с точностью до домножения на константу;
  - д) Найдите  $\hat{p}_{ML}$ ;
  - е) Найдите  $\text{Var}(s(\theta))$  и  $\text{Var}(\hat{p})$ ;
  - ж) Найдите предел  $\lim \text{Var}(s(\theta)) \cdot \text{Var}(\hat{p})$ ?
3. Исследовательница Несмеяна вывела хитрую формулу для  $\hat{a}$  — несмещённой оценки неизвестного векторного параметра  $a$ . Обозначим  $s(a)$  — вектор-столбец градиент логарифмической функции правдоподобия. Докажите, что для оценки Несмеяны выполнено неравенство Крамера-Рао, а именно, матрица  $M = \text{Var}(s(a)) \cdot \text{Var}(\hat{a}) - I_{k \times k}$  положительно определена.

Подсказки:

- а) Вспомните, чему равно  $\mathbb{E}(s(a))$ . Достаточно просто вспомнить, доказывать не требуется!
- б) Найдите скаляры  $\text{Cov}(\hat{a}_1, \frac{\partial \ell}{\partial a_1})$ ,  $\text{Cov}(\hat{a}_1, \frac{\partial \ell}{\partial a_2})$  и матрицу  $\text{Cov}(\hat{a}, s(a))$ .
- в) Рассмотрим два произвольных случайных вектора  $R$  и  $S$  и два вектора констант подходящей длины  $\alpha$  и  $\beta$ . Найдите минимум функции  $f(\alpha, \beta) = \text{Var}(\alpha^T R + \beta^T S)$  по  $\beta$ . Выпишите явно  $\beta^*(\alpha)$  и  $f^*(\alpha)$ .
- г) Докажите, что для произвольных случайных векторов положительно определена матрица

$$\text{Var}(R) - \text{Cov}(R, S) \text{Var}^{-1}(S) \text{Cov}(S, R)$$

- д) Завершите доказательство векторного неравенства Крамера-Рао.

Без угрызений совести можно храбро переставлять интегралы и производные :)

3-лайт! Утешительная версия задачи про Несмеяну. Если не получилось доказать векторную версию неравенства Крамера-Рао, то докажите скалярную :)

Докажите, что для несмещённой скалярной оценки  $\text{Var}(s(a)) \cdot \text{Var}(\hat{a}) \geq 1$ .

Подсказки:

- а) Вспомните, чему равно  $\mathbb{E}(\ell'(a))$ . Достаточно просто вспомнить, доказывать не требуется!
- б) Найдите  $\text{Cov}(\hat{a}, \ell'(a))$ ;
- в) Сколько корней может быть у параболы  $f(t) = \text{Var}(R + tL)$ ? Каким может быть дискриминант параболы  $f(t)$ ?
- г) Докажите для произвольных случайных величин  $R$  и  $L$  неравенство Коши-Шварца,

$$\text{Var}(R) \cdot \text{Var}(L) \geq \text{Cov}^2(R, L).$$

- д) Завершите доказательство скалярного неравенства Крамера-Рао :)
4. Идея доказательства состоятельности ML оценки :)
- Пусть наблюдения  $y_1, \dots, y_n$  независимы и одинаково распределены с функцией плотности, зависящей от параметра  $a$ . Истинное значение параметра обозначим буквой  $a_0$ . Оценку максимального правдоподобия обозначим  $\hat{a}$ .
- Рассмотрим отмасштабированную логарифмическую функцию правдоподобия  $\ell_n(a) = \ell(a)/n$ , и ожидаемую логарифмическую функцию правдоподобия<sup>4</sup>,  $\tilde{\ell}(a) = \mathbb{E}(\ell(a))$ .

- а) Что больше,  $\ln x$  или  $x - 1$ ? Докажите!
  - б) В какой точке находится максимум функции  $\ell_n(a)$ ?
  - в) В какой точке находится максимум функции  $\tilde{\ell}(a)$ ?
- Подсказка: рассмотрите выражение  $\tilde{\ell}(a) - \tilde{\ell}(a_0)$  и примените доказанное неравенство :)
- г) К чему сходится  $\ell_n(a)$  по вероятности?

5. Известна структура обратимой матрицы  $M$ ,

$$M = \begin{pmatrix} A & B \\ 0 & I_{k \times k} \end{pmatrix}.$$

- а) Найдите  $M^{-1}$ .
- б) Какие условия должны выполняться на блоки  $A$  и  $B$ , чтобы  $M$  была обратимой?

---

<sup>4</sup>Внимание: ожидание считается с помощью истинного  $a_0$  от функции, в которую входит константа  $a$ .