

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

Коллектив кафедры  
математической экономики и эконометрики,  
фольклор

14 марта 2015 г.

## Содержание

### 1 Описание

### 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 призванье 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
  - (a) РЭШ
  - (b) СМЕРШ
  - (c) МГУ
  - (d) На Луну
  - (e) .....
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 .....
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 домашних творческий подход, но мне было сказано, что я несу всякую чушь
17. Все остальные преподаватели согласны повысить мою оценку
18. Семинары и лекции начинались:
  - (a) слишком рано, я еще спал

- (b) слишком поздно, я уже спал
- (c) 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 .....

## 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

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 ТГМ. Детерминированные регрессоры

### 3.3 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:  $E(\varepsilon_i | \text{все } x_{ij}) = 0$   
В матричном виде:  $E(\varepsilon_i | X) = 0$

5. Условная гомоскедастичность:  $E(\varepsilon_i^2 | \text{все } x_{ij}) = \sigma^2$   
В матричном виде:  $E(\varepsilon_i^2 | X) = \sigma^2$
6.  $Cov(\varepsilon_i, \varepsilon_j | X) = 0$  при  $i \neq j$
7. вектора  $(x_i, y_i)$  — независимы и одинаково распределены
8. с вероятностью 1 среди регрессоров нет линейно зависимых  $rank(X) = k$   $det(X'X) \neq 0$   
 $(X'X)^{-1}$  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,  $E(\hat{\beta} | X) = \beta$ , и в частности  $E(\hat{\beta}) = \beta$
3. (тГМ) МНК оценки эффективны среди линейных несмещённых оценок. Для любой альтернативной оценки  $\hat{\beta}^{alt}$  удовлетворяющей свойствам 1 и 2:  $Var(\hat{\beta}_j^{alt} | X) \geq Var(\hat{\beta}_j | X)$   
 $Var(\hat{\beta}_j^{alt}) \geq Var(\hat{\beta}_j)$
4.  $Var(\hat{\beta} | X) = \sigma^2 (X'X)^{-1}$
5.  $Cov(\hat{\beta}, \hat{\varepsilon} | X) = 0$
6.  $E(\hat{\sigma}^2 | X) = \sigma^2$ , и  $E(\hat{\sigma}^2) = \sigma^2$  ?остается ли при условной ГК?

Если дополнительно к предпосылкам теоремы Гаусса-Маркова известно, что  $\varepsilon | X \sim N$ , 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  $F | X \sim F$

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

1.  $\hat{\beta} \rightarrow \beta$  по вероятности
2.  $t \rightarrow N(0, 1)$
3.  $rF \rightarrow \chi_r^2$ ,  $r$  — число ограничений
4.  $nR^2 \rightarrow \chi_{k-1}^2$
5.  $\frac{RSS}{n-k} \rightarrow \sigma^2$

### 3.4 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,  $\text{tr}(A) = \sum a_{ii}$ .

Свойства следа:

1.  $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$
2.  $\text{tr}(AB) = \text{tr}(BA)$ , если  $AB$  и  $BA$  существуют. При этом  $A$  и  $B$  могут не быть квадратными матрицами.
3.  $\text{tr}(A) = \sum \lambda_i$ , 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  $\text{tr}(A)$  — размерность этого подпространства. Действительно, если  $A$  — проектор, то  $A^2 = A$  и собственные числа матрицы  $A$  равны нулю или единице. Поэтому  $\text{tr}(A)$  равен количеству собственных чисел равных единице. И, следовательно,  $\text{tr}(A)$  равен  $\text{rank}(A)$ , то есть размерности пространства, на которое матрица  $A$  проецирует вектора. У следа матрицы есть и другие смыслы [mathoverflow0trace].

### 3.5 Ожидание от RSS

**Теорема 1.** След и математическое ожидание можно переставлять,  $\mathbb{E}(\text{tr}(A)) = \text{tr}(\mathbb{E}(A))$ .

**Теорема 2.** Математическое ожидание квадратичной формы

$$\mathbb{E}(x'Ax) = \text{tr}(A \text{Var}(x)) + \mathbb{E}(x')A\mathbb{E}(x) \quad (1)$$

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

Поехали,

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

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

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

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

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

□

### 3.6 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.

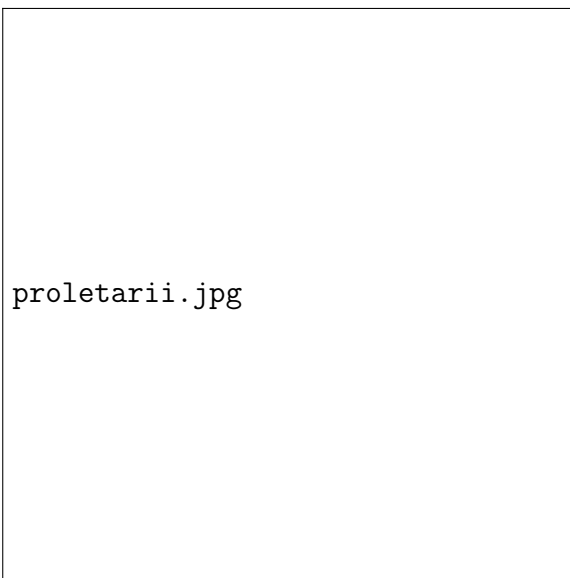
Гипотезы имеет смысл проверять о ненаблюдаемых величинах, а величина  $R^2$  является наблюдаемой. И если уж на то пошло, то проверить гипотезу о том, что  $R^2 = 0$  тривиально. Для этого не нужно знать ничего из теории вероятностей, достаточно просто сравнить посчитанное значение  $R^2$  с нулём.

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 «Н naught», а не «Н zero» или «Н null». Также корректно говорить «the null hypothesis».

## 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

$$A = \begin{pmatrix} 2 & 3 & 0 \\ 3 & 10 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$

- (a) Найдите собственные числа и собственные векторы матрицы.
- (b) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.
- (c) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.
- (d) 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.
- (a) Какова вероятность того, что Петя правильно ответил на оба вопроса?
- (b) Какова вероятность того, что Петя правильно ответил на оба вопроса, если его ответы совпали с Васиными?
- (c) Чему равно математическое ожидание числа Петиних верных ответов?
- (d) 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:
- (a)  $\mathbb{E}(U)$ ,  $\text{Var}(U)$ ,  $\mathbb{E}(V)$ ,  $\text{Var}(V)$ ,  $\text{Cov}(U, V)$
- (b) 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

1. Случайные величины  $Z_i$  независимы и нормально распределены  $N(0, 1)$ . Для их суммы  $S = \sum_{i=1}^n Z_i$  найдите  $\mathbb{E}(S)$  и  $\text{Var}(S)$ .
2. 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?»
- (a) Какова вероятность, что случайно выбранный респондент ответит «Да»?
- (b) Какова вероятность того, что он действительно доверяет, если известно, что он ответил «Да»?
3. 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

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

- (a) Укажите число наблюдений.
  - (b) Укажите число регрессоров с учетом свободного члена.
  - (c) Рассчитайте при помощи метода наименьших квадратов  $\hat{\beta}$ , оценку для вектора неизвестных коэффициентов.
  - (d) Рассчитайте  $TSS = \sum (y_i - \bar{y})^2$ ,  $RSS = \sum (y_i - \hat{y}_i)^2$  и  $ESS = \sum (\hat{y}_i - \bar{y})^2$ .
  - (e) Чему равен  $\hat{\varepsilon}_4$ , МНК-остаток регрессии, соответствующий 4-ому наблюдению?
  - (f) Чему равен  $R^2$  в модели?
  - (g) Рассчитайте несмещенную оценку для неизвестного параметра  $\sigma^2$  регрессионной модели.
  - (h) Рассчитайте  $\widehat{\text{Var}}(\hat{\beta})$ , оценку для ковариационной матрицы вектора МНК-коэффициентов  $\hat{\beta}$ .
  - (i) Найдите  $\widehat{\text{Var}}(\hat{\beta}_1)$ , несмещенную оценку дисперсии МНК-коэффициента  $\hat{\beta}_1$ .
  - (j) Найдите  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2)$ , несмещенную оценку ковариации МНК-коэффициентов  $\hat{\beta}_1$  и  $\hat{\beta}_2$ .
  - (k) Найдите  $\widehat{\text{Var}}(\hat{\beta}_1 + \hat{\beta}_2)$
  - (l) Найдите  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2)$ , оценку коэффициента корреляции МНК-коэффициентов  $\hat{\beta}_1$  и  $\hat{\beta}_2$ .
  - (m) 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 3. Дню рождения буквы «ё» посвящается...

1. Выберите верные варианты.
  - (a) Побасёнка — Побасенка
  - (b) Вёдро — Ведро
  - (c) Гренадёр — Гренадер
  - (d) Новорождённый — Новорожденный
  - (e) Бытиё — Бытие
  - (f) Опёка — Опека
  - (g) Сёрфинг — Серфинг
  - (h) Пафнутий Львович Чебышёв — Пафнутий Львович Чебышев
  - (i) 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.

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



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

- (a) Заполните пропуски в таблице.
- (b) Укажите коэффициенты, значимые на 10% уровне значимости.
- (c) 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:  $Fertility_i = \gamma_1 + \gamma_2(Education_i + Examination_i) + \gamma_3 Catholic_i + u_i$

Таблица 1:

	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  $p < 0.05$

- 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?

#### 4.4 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
  - Постройте 95%-ый доверительный интервал для параметра  $a$

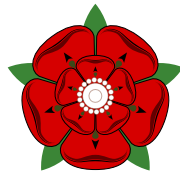
2. [R] Фактическое распределение часовой и десятиминутной скорости ветра хорошо приближается распределением Вейбулла. Случайная величина имеет распределение Вейбулла, если её функция плотности при  $x > 0$  имеет вид

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

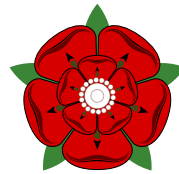
- (a) Оцените параметры  $k$  и  $\lambda$  методом максимального правдоподобия
- (b) Постройте 95%-ые доверительные интервалы для  $k$  и  $\lambda$

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

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

- (a) Найдите функцию распределения  $F(x)$
- (b) Выразите медиану распределения Вейбулла,  $m$ , через параметры  $k$  и  $\lambda$
- (c) Оцените параметры  $k$  и  $\lambda$  методом максимального правдоподобия
- (d) Постройте 95%-ые доверительные интервалы для  $k$  и  $\lambda$
- (e) Выпишите функцию плотности распределения Вейбулла через  $m$  и  $k$
- (f) 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).

## 4.5 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.
- (a) Почему это могло случиться?
  - (b) 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.

- (a) Отсортировав имеющиеся наблюдения по количеству осадков в день, Настоячивый 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?
- (b) 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.
- (c) 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.

## 4.6 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 приличные. 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

```
## Error in read.spss(file.name, to.data.frame = TRUE): unable to open file: 'No  
such file or directory'
```

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`.

```
## NULL
```

### 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.

### 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. Напишите свои пожелания и комментарии.

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

## 4.7 Домашнее задание №(n + 1) по эконометрике-1.

### Задача 1. «САРМ»

Оценим модель САРМ по реальным данным:

1. Коротко сформулируйте теоретические положения модели САРМ. 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  $R_i - R_{f,i}^{true} = \beta(R_{m,i} - R_{f,i}^{true}) + \varepsilon_i$ .
  - (a) Выпишите функцию правдоподобия для оценки данной модели
  - (b) Найдите оценки  $\hat{\beta}$ ,  $\hat{\sigma}_u^2$ ,  $\hat{\sigma}_\varepsilon^2$
  - (c) Постройте 95%-ые доверительные интервалы
  - (d) Прodelайте аналогичные действия для модели с константой
  - (e) 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.
  - (a) Выпишите функцию правдоподобия
  - (b) Оцените параметры  $\beta_1$  и  $\beta_2$
  - (c) Оцените ковариационную матрицу оценок  $\hat{\beta}_1$  и  $\hat{\beta}_2$
  - (d) Постройте 95%-ые асимптотические доверительные интервалы для параметров
  - (e) Проверьте гипотезу о том, что  $\beta_2 = 0$ . Альтернативную гипотезу сформулируйте самостоятельно.

PS. Своё смелое творчество в задачах поощряется!

## 4.8 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.



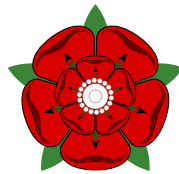
### Версия Белой Розы



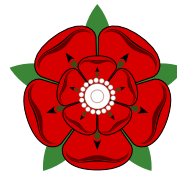
1. Домашнее задание можно делать в одиночку или группой из двух человек.
2. А можно всё-таки группой из трёх человек? Нет :)
3. Письменный отчёт должен содержать как-минимум:

- (a) Логин группы
- (b) Графический анализ имеющихся данных
- (c) Результаты оценивания logit и probit моделей
- (d) Графический анализ logit и probit моделей
- (e) «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:
- (a) Логин
  - (b) Графический анализ имеющихся данных
  - (c) Результаты оценивания logit и probit моделей
  - (d) Прогнозирование с использованием Random Forest
  - (e) Прогнозирование с использованием метода опорных векторов (SVM)
  - (f) Графический анализ оценённых моделей
  - (g) «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

$$A = \begin{pmatrix} -2 & 0 & 0 \\ 0 & 3 & 4 \\ 0 & 4 & 9 \end{pmatrix}$$

- (a) Найдите собственные числа и собственные векторы матрицы.
  - (b) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.
  - (c) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.
  - (d) 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.
    - (a) Какова вероятность того, что Петя правильно ответил на оба вопроса?
    - (b) Какова вероятность того, что Петя правильно ответил на оба вопроса, если его ответы совпали с Васиными?
    - (c) Чему равно математическое ожидание числа Петиних верных ответов?
    - (d) 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:
    - (a)  $\mathbb{E}(U)$ ,  $\text{Var}(U)$ ,  $\mathbb{E}(V)$ ,  $\text{Var}(V)$ ,  $\text{Cov}(U, V)$
    - (b) 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:

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

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 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  $D =$

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}. \text{ Рассмотрите «новую» регрессионную модель } y = Z\alpha + u, \text{ где } \alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}.$$

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.

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

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

- (a) Заполните пропуски в таблице  
 (b) Укажите коэффициенты, значимые на 10% уровне значимости.  
 (c) 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)

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

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

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

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

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

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

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

(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  $H_a : |\beta_2 - 1| + |\beta_3 - 2| \neq 0$ .

## 5.5 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  $\sigma^2$
  - С помощью LR, LM и W теста проверьте гипотезу о том, что  $\sigma^2 = 1$ .
  - 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 **Аксинья-полухлебница**.

«На Аксинью 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

## 5.6 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.

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

## 5.7 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:
  - (a) Обычную МНК-оценку параметра  $\beta$
  - (b) Самую эффективную среди несмещённых оценку параметра  $\beta$
  - (c) Во сколько раз отличается истинная дисперсия этих двух оценок?
  - (d) Во сколько раз отличаются оценки дисперсий этих оценок, если дисперсии оцениваются без поправки на гетероскедастичность в обоих случаях?

## 5.8 Большой Устный Зачёт

1. Метод Наименьших Квадратов.
  - (a) МНК-картинка
  - (b) Нахождение всего-всего, если известен вектор  $y$  и матрица  $X$
2. Теорема Гаусса-Маркова
  - (a) Формулировка с детерминистическими регрессорами
  - (b) Доказательство с детерминистическими регрессорами
  - (c) Формулировки со стохастическими регрессорами
  - (d) Что даёт дополнительное предположение о нормальности  $\varepsilon$ ?

<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.

### 3. Проверка гипотез о линейных ограничениях

- (a) Проверка гипотезы о значимости коэффициента
- (b) Проверка гипотезы о значимости регрессии в целом
- (c) Проверка гипотезы об одном линейном соотношении с помощью ковариационной матрицы
- (d) Ограниченная и неограниченная модель
- (e) Тест Чоу на стабильность коэффициентов
- (f) Тест Чоу на прогнозную силу

### 4. Метод максимального правдоподобия

- (a) Свойства оценок
- (b) Два способа получения оценки дисперсии
- (c) Три теста (LM, Wald, LR)
- (d) Выписать функцию ML для обычной регрессии
- (e) для AR(1) процесса
- (f) для MA(1) процесса
- (g) для логит модели
- (h) для пробит модели
- (i) для модели с заданным видом гетероскедастичности

### 5. Мультиколлинеарность

- (a) Определение, последствия
- (b) Величины, измеряющие силу мультиколлинеарности
- (c) Методы борьбы
- (d) Сюда же: метод главных компонент, хотя он используется и для других целей

### 6. Гетероскедастичность

- (a) Определение, последствия
- (b) Тесты, график
- (c) Стьюдентизированные остатки
- (d) НС оценки ковариации
- (e) GLS и FGLS

### 7. Временные ряды

- (a) Стационарный временной ряд
- (b) ACF, PACF
- (c) Модель ARMA
- (d) Модель GARCH (не будет, не успели)

### 8. Логит и пробит

- (a) Описание моделей

- (b) Предельные эффекты
- (c) Чувствительность, специфичность
- (d) Кривая ROC

9. Эндогенность

- (a) Три примера: одновременность, пропущенные переменные, ошибки измерения
- (b) IV, двухшаговый МНК

10. Модели панельных данных

- (a) RE, FE, сквозная регрессии
- (b) Тест Хаусмана

11. Альтернативные методы. Уметь объяснить суть метода. Уметь реализовать его в R.

- (a) Метод опорных векторов (не будет, не успели)
- (b) Классификационные деревья и случайный лес

12. R. Можно принести файл со своей заготовкой, можно пользоваться Интернетом для поиска информации, но не для общения.

- (a) Загрузить данные из .csv файла в R
- (b) Посчитать описательные статистики (среднее, мода, медиана и т.д.)
- (c) Построить подходящие описательные графики для переменных
- (d) Оценить линейную регрессию с помощью МНК. Провести диагностику на что-нибудь (гетероскедастичность, автокорреляцию, мультиколлинеарность).
- (e) Оценить logit, probit модели, посчитать предельные эффекты
- (f) Оценить ARMA модель
- (g) Выделить главные компоненты

## 5.9 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

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

- (a) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
- (b) Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
- (c) 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.

	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

- (a) Оцените вероятность выжить для женщины 20 лет
- (b) Оцените предельный эффект увеличения возраста для женщины 20 лет
- (c) С помощью какого метода оценивается логит-модель? Каким образом при этом получают оценки стандартных ошибок коэффициентов?
3. Теорема Гаусса-Маркова.
- (a) Аккуратно сформулируйте теорему Гаусса-Маркова для нестохастических регрессоров.
- (b) Поясните каждое из свойств оценок, фигурирующих в теореме.
- (c) 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  $\text{Var}(\varepsilon_t) = \sigma^2$ .

- (a) Является ли данный процесс стационарным? Что такое стационарный процесс?
- (b) Найдите автокорреляционную функцию данного процесса,  $\rho(k) = \text{Corr}(y_t, y_{t-k})$ .
- (c) 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

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

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

## 5.10 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

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

- (a) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - (b) Найдите несмещенную оценку для неизвестного параметра  $\sigma^2$ .
  - (c) 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.

	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

- (a) Оцените вероятность выжить для женщины 20 лет
- (b) Оцените предельный эффект увеличения возраста для женщины 20 лет



- (с) С помощью какого метода оценивается логит-модель? Каким образом при этом получаются 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.  
 (с) 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%.  
 (с) Какой тест можно на гетероскедастичность можно было бы использовать, если бы не было уверенности в нормальности остатков? Опишите пошагово процедуру этого теста.

## 5.11 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. Попадают довольно приличные. 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

```
library(foreign)
file.name<-"/home/boris/downloads/r20hall23c.sav"
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`.

```
var.labels<-attr(h,"variable.labels")
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://rpubs.com/boris\\_demeshev/r\\_cycle\\_12](http://rpubs.com/boris_demeshev/r_cycle_12)

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 + L<sup>A</sup>T<sub>E</sub>X. 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()`.

## 5.12 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 + L<sup>A</sup>T<sub>E</sub>X. 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 <https://github.com/bdemeshev/em301/wiki/R>.

## 6 2014-2015

### 6.1 Праздник номер 1

Вперёд, в рукопашную!

1. Сформулируйте теорему о трёх перпендикулярах и обратную к ней.

2. Для матрицы

$$A = \begin{pmatrix} 3 & 4 \\ 4 & 9 \end{pmatrix}$$

(a) Найдите собственные числа и собственные векторы матрицы.

(b) Найдите обратную матрицу,  $A^{-1}$ , ее собственные векторы и собственные числа.

(c) Представьте матрицу  $A$  в виде  $A = CDC^{-1}$ , где  $D$  — диагональная матрица.

(d) Найдите  $A^{42}$

(e) 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.

(a) Какова вероятность, что у него как минимум два туза?

- (b) Каково ожидаемое количество тузов у игрока?
  - (c) Какова вероятность, что у него как минимум два туза, если известно, что у него есть хотя бы один туз?
  - (d) 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.
- (a) Постройте 95% доверительный интервал для математического ожидания времени проводимого сотрудниками на работе
  - (b) 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.

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

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

- (a) Заполните пропуски в таблице
  - (b) Укажите коэффициенты, значимые на 10% уровне значимости.
  - (c) 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  $\text{Cov}(\hat{\varepsilon}_1, \hat{\varepsilon}_2) = 0$ ?
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...

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

- (a) Помогите Ефросинье найти количество наблюдений,  $\bar{z}$ ,  $\sum x_i z_i$ ,  $\sum (x_i - \bar{x})(z_i - \bar{z})$
- (b) (\*) 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 3

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

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

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

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

- Напишите спецификацию регрессии с ограничениями для проверки статистической гипотезы  $H_0 : \beta_5 = 2\beta_4$
  - Дайте интерпретацию проверяемой гипотезе
  - Для регрессии с ограничением был вычислен коэффициент  $R_R^2 = 0.296$ . На уровне значимости 5% проверьте нулевую гипотезу
- 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:

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

- Постройте 95%-ый доверительный интервал для коэффициента  $\beta_2$
  - Спрогнозируйте стоимость экспорта на 2010 год и постройте 90%-ый предиктивный интервал для прогноза.
- 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.

- (a) Найдите истинную дисперсию МНК оценки коэффициента  $\beta$
  - (b) Предложите более эффективную оценку  $\hat{\beta}^{alt}$
  - (c) Чему равна истинная дисперсия новой оценки?
  - (d) 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:

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

- (a) Найдите  $\mathbb{E}(\varepsilon)$ ,  $\mathbb{E}(x_1 \cdot \varepsilon_1)$
  - (b) Найдите  $\text{plim } \frac{1}{n} X' \varepsilon$
  - (c) Найдите  $\text{plim } \frac{1}{n} X' X$
  - (d) 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  $\sigma_u^2$ .
- (a) Будут ли оценки Эвридики несмещенными?
  - (b) Могут ли дисперсии оценок Эвридики быть ниже чем дисперсии МНК оценок при использовании настоящего  $y_i$ ?
  - (c) Могут ли оценки дисперсий оценок Эвридики быть ниже чем оценок дисперсий МНК оценок при использовании настоящего  $y_i$ ?

## 6.4 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.
  - (a) Укажите вид матрицы  $X$
  - (b) Выпишите формулу для МНК оценки  $\hat{\beta}$
  - (c) 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?

- Какую гипотезу можно проверить, зная отношение  $ESS/RSS$ ? Укажите  $H_0$ ,  $H_a$ , тестовую статистику и её распределение при верной  $H_0$ .
- Опишите подробно тест Чоу на прогнозную силу по двум наборам данных из  $n_1$  и  $n_2$  наблюдений соответственно. Число оцениваемых коэффициентов равно  $k$ .

## 6.5 Зачет. Базовый поток

- Сформулируйте теорему Гаусса-Маркова применительно к модели  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ . Поясните смысл каждого используемого термина.
- 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  $H_0$ .
- 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

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

- 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  $y_i$ ?
- Рассмотрим модель  $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ , где  $\varepsilon_i \sim N(0; \sigma^2)$ .

- Напишите выражения для оценок дисперсий и ковариации коэффициентов, т.е.  $\widehat{\text{Var}}(\hat{\beta}_1)$ ,  $\widehat{\text{Var}}(\hat{\beta}_2)$ ,  $\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2)$
- Найдите математическое ожидание и дисперсию каждой из выписанных оценок
- Какой закон распределения с точностью до масштабирования имеют эти оценки?

- Для модели данных  $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + \varepsilon_i$  по 100 наблюдениям получены результаты:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.9571	0.0942	20.79	0.0000
x	3.1947	0.1004	31.80	0.0000
z	0.0543	0.1040	0.52	0.6025

- (a) Выпишите полученное уравнение регрессии
- (b) Укажите, какие коэффициенты значимы при  $\alpha = 0.05$
- (c) 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  $\varepsilon_i \sim N(0; \sigma^2)$ .

- (a) Выпишите формулу для оценок коэффициентов и оценки дисперсии ошибок
- (b) Укажите математическое ожидание и дисперсию выписанных оценок
- (c) Для оценок коэффициентов укажите закон распределения
- (d) Для оценки дисперсии ошибок укажите закон распределения с точностью до масштабирования

## 6.6 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  $\mathbb{E}(\varepsilon) = 0$  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

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

- (a) Найдите вектор МНК-оценок коэффициентов  $\hat{\beta}$ .
  - (b) Найдите коэффициент детерминации  $R^2$
  - (c) 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.

- (a) Предполагая гомоскедастичность остатков на уровне значимости 5% проверьте гипотезу, что исследуемая зависимость одинакова на всех трёх частях всей выборки.
- (b) Протестируйте ошибки на гетероскедастичность на уровне значимости 5%.



- (с) 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.

	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

Оценка стандартной ошибки случайной составляющей,  $\hat{\sigma} = 33.0252513$ .

- (a) Можно ли интерпретировать коэффициент при переменной *totsp* как стоимость одного метра нежилой площади?
- (b) Проверьте гипотезу о том, что коэффициенты при регрессорах *totsp* и *livesp* равны.
- (c) 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>.
- (d) 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
- (a) для нестохастических регрессоров
- (b) для стохастических регрессоров в предположении, что наблюдения являются случайной выборкой

## 6.7 Домашняя работа 1. RLMS и гетероскедастичность

- Прочитайте про RLMS, <http://www.hse.ru/rlms/>  
Посмотрите описание проекта. Проллистайте вестник RLMS, чтобы иметь представление о том, какие исследования можно строить на основе RLMS.
- Скачайте любую волну RLMS по своему выбору. Скачайте описание переменных.  
Проллистайте описание переменных. Там их больше тысячи. Попадаются довольно прикольные. Мне нравится `pc9.6.5a`, «У Вас есть GPRS навигатор?»
- 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`.

```
library("rlms")
h <- read.rlms("/home/boris/downloads/r20hall123c.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:

```
var_meta <- attr(h, "var_meta")
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 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://rpubs.com/boris\\_demeshev/r\\_cycle\\_12](http://rpubs.com/boris_demeshev/r_cycle_12)

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 + L<sup>A</sup>T<sub>E</sub>X 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()`.