

# Todo list

Проверить! Нет ли у $\beta_1$ особого положения? . . . . .	2
на картинке три $s$ : очень большое — дающее мнк решение, меньше — ненулевые $\beta$ , маленькое — одна из $\beta$ равна 0 . . . . .	2
Может ли появиться мультимодальность? В точности ли на моду или только примерно? .	3

## 1 Конвенция

$y$  — вектор столбец зависимых переменных, наблюдаемый случайный

$\beta$  — вектор столбец неизвестных параметров, ненаблюдаемый, неслучайный

$\hat{y}$  — прогноз  $y$  полученный по некоторой модели, наблюдаемый, случайный

$\hat{\beta}$  — оценки  $\beta$

$X$  — матрица всех объясняющих переменных

$\varepsilon$

$\hat{\varepsilon}$

Некоторые авторы используют обозначения:

$Y$  и  $y$  для разных вещей,  $y = Y - \bar{Y}$ .

## 2 Семинар 1

Неформальное определение. 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.  $\det(A) = \prod \lambda_i$

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

Свойства следа:

1.  $\text{tr}(AB) = \text{tr}(BA)$ , если  $AB$  и  $BA$  существуют. При этом  $A$  и  $B$  могут не быть квадратными матрицами.
2.  $\text{tr}(A) = \sum \lambda_i$

Добавить про геометрический смысл следа, <http://mathoverflow.net/questions/13526/geometric-interpretation-of-trace>.

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

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

Метод наименьших квадратов (МНК), ordinary least squares (OLS):

Есть  $n$  наблюдений,  $y_1, \dots, y_n$ . Есть модель, которая даёт прогнозы,  $\hat{y}_1, \dots, \hat{y}_n$ . Эта модель зависит от вектора неизвестных параметров,  $\beta$ . МНК предлагает в качестве оценок неизвестных параметров взять такое  $\hat{\beta}$ , чтобы минимизировать  $\sum (y_i - \hat{y}_i)^2$ .

## 3 Семинар 2

## 4 Разное

1. Гипотеза  $H_0$  по-английски читается как «H naught»

2. При проверке гипотезы об адекватности регрессии НЕЛЬЗЯ писать  $H_0 : R^2 = 0$ .

Гипотезы имеет смысл проверять о ненаблюдаемых неизвестных константах. Проверить гипотезу о том, что  $R^2 = 0$  легко. Для этого не нужно знать ничего из теории вероятностей, достаточно просто сравнить посчитанное значение  $R^2$  с нулём.

Более того, даже корректировка  $\mathbb{E}(R^2) = 0$  неверна. Случайная величина  $R^2$  всегда неотрицательна, поэтому при любых разумных предположениях на  $\varepsilon$  окажется, что  $\mathbb{P}(R^2 > 0) > 0$ . А это приведёт к тому, что  $\mathbb{E}(R^2) > 0$  даже если  $Y$  никак не зависит от  $X$ .

Единственный правильный вариант —  $H_0 : \beta_2 = \beta_3 = \dots = \beta_k = 0$  и  $H_a : \exists i \geq 2 : \beta_i \neq 0$ .

## 5 Ridge/Lasso regression

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

Вспомним обычный МНК:

$$\min_{\beta} (y - X\beta)'(y - X\beta) \quad (1)$$

LASSO вместо исходной задачи решает задачу условного экстремума:

$$\min_{\beta} (y - X\beta)'(y - X\beta) \quad (2)$$

при ограничении  $\sum_{j=1}^k |\beta_j| \leq c$ .

**Проверить! Нет ли у  $\beta_1$  особого положения?**

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.

$$L(\beta, \lambda) = (y - X\beta)'(y - X\beta) + \lambda \left( \sum_{j=1}^k |\beta_j| - c \right) \quad (3)$$

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

$$\min_{\beta} (y - X\beta)'(y - X\beta) + \lambda \sum_{j=1}^k |\beta_j| \quad (4)$$

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:

$$\min_{\beta} (y - X\beta)'(y - X\beta) + \lambda \sum_{j=1}^k \beta_j^2 \quad (5)$$

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 приходится в точности (?) на оценки LASSO.

Может ли появиться мультимодальность? В точности ли на моду или только примерно?