
УНИВЕРСИАДА ПО ЭКОНОМЕТРИКЕ 2024

РЕШЕНИЯ ЗАДАЧ ВТОРОГО ТУРА

Задача 1. (25 баллов) Размышления о p -значении

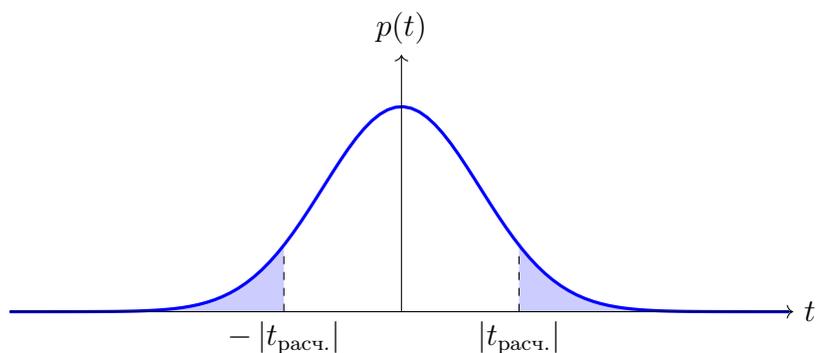
①

Если p -значение теста Стьюдента на значимость для параметра β_2 составило 0.007, то вероятность получить такое (или большее по модулю) значение t -статистики при $\beta_2 = 0$ оказывается меньше 0.7%. Это означает, что на уровне значимости 0.7% (и любом большем) гипотеза $H_0 : \beta_2 = 0$ должна быть отвергнута. Коэффициент статистически значимо отличается от нуля (2 балла).

②

Для (двустороннего) теста Стьюдента p -значение определяется как

$$p = \mathbb{P}(|t| \geq |t_{\text{расч.}}| \mid H_0 \text{ верна}) = 2 \mathbb{P}(t \geq |t_{\text{расч.}}| \mid H_0 \text{ верна}), \quad t \sim t_{n-k}$$



Фактически p -значение теста — это минимальный уровень значимости, на котором гипотеза H_0 может быть отвергнута.

Пункт оценивался в 3 балла. Принимались также объяснения, что p -значение — это площадь «хвостов» в распределении t -статистики (см. рисунок) или вероятность отвергнуть H_0 при условии, что она верна.

③

Пункт оценивался в 6 баллов, по 2 балла за каждый подпункт. Точный ответ на этот вопрос можно дать только при условии, что верна нулевая гипотеза¹ (при отсутствии указания на это снимался 1 балл).

- $\mathbb{P}(p\text{-value} < 10\% \mid H_0 \text{ верна, } \varepsilon \sim \mathcal{N}(0, \sigma^2)) = 0.1$ (ясно из вопроса пункта ④)
- $\mathbb{P}(p\text{-value} < 10\% \mid H_0 \text{ верна, } \varepsilon \not\sim \mathcal{N}(0, \sigma^2))$ не может быть однозначно определена и зависит от распределения ε (например, эта вероятность может существенно отличаться для дискретно распределенных ε)
- $\mathbb{P}(p\text{-value} < x \mid H_0 \text{ верна, } \varepsilon \sim \mathcal{N}(0, \sigma^2)) = \begin{cases} 0, & \text{при } x < 0 \\ x, & \text{при } 0 \leq x \leq 1 \\ 1, & \text{при } x > 1 \end{cases}$ – снова ясно из формулировки ④ (и вновь, если распределение ε не является нормальным, то о распределении p -значений судить нельзя)

④

Рассуждение для (двустороннего!) теста Стьюдента.

$$\mathbb{P}(p\text{-value} < x \mid H_0 \text{ верна,}) = \begin{cases} 0, & \text{при } x < 0 \\ x, & \text{при } 0 \leq x \leq 1 \\ 1, & \text{при } x > 1 \end{cases}$$
 В выкладках использованы

монотонность и непрерывность F_0 для распределения $t \sim t_{n-2}$):

$$\begin{aligned} \mathbb{P}(p\text{-value} < \alpha \mid H_0 \text{ верна}) &= \mathbb{P}(2\mathbb{P}(t > |t_{\text{расч.}}|) < \alpha) = \mathbb{P}\left(\mathbb{P}(t > |t_{\text{расч.}}|) < \frac{\alpha}{2}\right) = \\ &= \mathbb{P}\left(1 - F_0(t_{\text{расч.}}) < \frac{\alpha}{2}\right) = \mathbb{P}\left(F_0(|t_{\text{расч.}}|) > 1 - \frac{\alpha}{2}\right) = \\ &= \mathbb{P}\left(|t_{\text{расч.}}| > F_0^{-1}\left(1 - \frac{\alpha}{2}\right)\right) = 2\mathbb{P}\left(t_{\text{расч.}} > F_0^{-1}\left(1 - \frac{\alpha}{2}\right)\right) = \\ &= 2\left(1 - \mathbb{P}\left(t_{\text{расч.}} < F_0^{-1}\left(1 - \frac{\alpha}{2}\right)\right)\right) = 2\left(1 - F_0\left(F_0^{-1}\left(1 - \frac{\alpha}{2}\right)\right)\right) = \\ &= 2\left(1 - \left(1 - \frac{\alpha}{2}\right)\right) = \alpha \quad (8 \text{ баллов}) \end{aligned}$$

Пункт оценивался в 8 баллов. При рассмотрении одностороннего теста (без двойки, без модуля) баллы снижались. При рассмотрении только случая $\alpha \in [0, 1]$ баллы снижались.

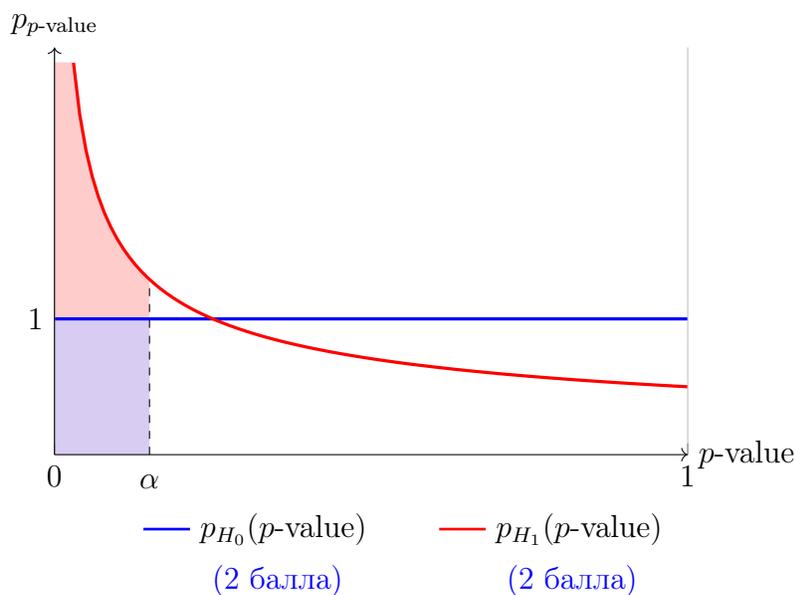
¹ Подробное обоснование здесь не требовалось (за его отсутствие балл не снижался), так как оно необходимо в пункте ④.

⑤—⑥

Сначала посмотрим, что меняется в рассуждении выше при верной альтернативной гипотезе H_1 . Обозначим плотность X при верной H_0 как $f_0(x)$, а при верной H_1 – как $f_1(x)$. Из определения p -значения мы имеем $p\text{-value} = 1 - F_0(X_{\text{расч.}})$. Тогда плотность p -значения имеет вид² (второй переход использует теорему о производной обратной функции)

$$p_{p\text{-value}}(\alpha) = f_1(F_0^{-1}(1 - \alpha)) \left| \frac{d}{d\alpha} F_0^{-1}(1 - \alpha) \right| = \frac{f_1(F_0^{-1}(1 - \alpha))}{f_0(F_0^{-1}(1 - \alpha))}$$

Видно, что распределение в общем случае уже не будет равномерным. При верной альтернативной гипотезе маленькие p -значения будут встречаться чаще, а большие – реже (то есть плотность p -значения будет убывать). Пример с некоторым заданным уровнем значимости α изображен на рисунке ниже.



Уровень значимости α при этом является отсечкой по оси абсцисс и отсекает закрашенную площадь (синюю – для случая верной H_0 и красную – для H_1) (2 балла).

² Это получается из следующего факта. Пусть X и Y – случайные величины, причем $Y = g(X)$ и X имеет плотность $p_X(x)$. Тогда из соображений теоремы о замене переменной (в общей ситуации – в интеграле Лебега) при некоторых требованиях на функцию $g(x)$ плотность Y имеет вид

$$p_Y(y) = p_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

Дополнительный комментарий, не влияющий на баллы. В контексте примера из задания можно также упомянуть, что для случая теста Стьюдента распределение тестовой статистики при верной гипотезе H_1 имеет так называемое [нецентрированное распределение Стьюдента](#), поэтому при большом желании точную форму плотности p -значения можно воспроизвести с его помощью. Соответствующий пример на сгенерированных данных можно посмотреть [здесь](#).

Также отдельно можно было бы отметить работу [\[Murdoch, Tsai & Adcock, 2008\]](#), в которой идея задачи проиллюстрирована дополнительными симуляциями и несколько развита.

Задача 2. (25 баллов)

Boosted ridge or Ridged boosting?

(a)

При оценивании линейной регрессии с помощью МНК вектор y ортогонально проектируется на линейную оболочку столбцов X , поэтому вектор ошибок $e^{(0)} = y - \hat{y}^{(0)}$ будет лежать в ортогональном дополнении к этой линейной оболочке. Поэтому предсказанный линейной регрессией $\hat{e}^{(0)}$ окажется тождественно нулевым, и все остальные $\hat{y}^{(m)}$ в точности совпадут с МНК-оценкой нулевого шага³:

$$\hat{y}^{(m)} = X (X^\top X)^{-1} X^\top y \quad (3 \text{ балла})$$

Выходит, никакого смысла в процедуре улучшения не было.

(b)

В векторно-матричном виде задачу Андрея по оцениванию ridge-модели можно записать в виде

$$\begin{aligned} (y - X\beta)^\top (y - X\beta) + \lambda\beta^\top \beta &\rightarrow \min_{\beta} \\ (y^\top - \beta^\top X^\top) (y - X\beta) + \lambda\beta^\top \beta &\rightarrow \min_{\beta} \\ y^\top y + \beta^\top X^\top X\beta - 2\beta^\top X^\top y + \lambda\beta^\top \beta &\rightarrow \min_{\beta} \end{aligned}$$

Оптимизируемая функция очевидно является выпуклой, поэтому для поиска оптимального β достаточно найти вектор, удовлетворяющий необходимому условию, которое записывается в матричном виде как

$$\begin{aligned} \nabla (y^\top y + \beta^\top X^\top X\beta - 2\beta^\top X^\top y + \lambda\beta^\top \beta) &= 0 \\ 2X^\top X\beta - 2X^\top y + 2\lambda\beta &= 0 \implies \hat{\beta}_{\text{ridge}} = (X^\top X + \lambda I)^{-1} X^\top y \end{aligned}$$

А тогда прогноз нулевого шага имеет вид

$$\hat{y}^{(0)} = X\hat{\beta}_{\text{ridge}} = X (X^\top X + \lambda I)^{-1} X^\top y \quad (4 \text{ балла})$$

³ То же самое можно получить напрямую, вычислив

$$\begin{aligned} \hat{e}^{(0)} &= X (X^\top X)^{-1} X^\top e^{(1)} = X (X^\top X)^{-1} X^\top (y - \hat{y}^{(1)}) = \\ &= X (X^\top X)^{-1} X^\top y - X (X^\top X)^{-1} X^\top X (X^\top X)^{-1} X^\top y = 0 \end{aligned}$$

(c)

Попробуем сначала вычислить $\hat{y}^{(1)}$ и $\hat{y}^{(2)}$, а затем обобщим результат. Вычисляем ошибки нулевого шага и предсказываем их тем же методом:

$$e^{(0)} = y - \hat{y}^{(0)} \implies \hat{e}^{(0)} = \underbrace{X (X^\top X + \lambda I)^{-1} X^\top}_B (y - \hat{y}^{(0)})$$
$$\hat{y}^{(1)} = \hat{y}^{(0)} + e^{(0)} = By + B(y - By) = B(I + (I - B))y$$

Теперь улучшим этот прогноз еще раз

$$e^{(1)} = y - \hat{y}^{(1)} \implies \hat{e}^{(1)} = B(y - \hat{y}^{(1)})$$
$$\hat{y}^{(2)} = \hat{y}^{(1)} + \hat{e}^{(1)} = B(I + (I - B))y + B(y - B(I + (I - B))y)$$
$$= B(I + (I - B) + (I - B)^2)y$$

Паттерн понятен и легко доказывается методом математической индукции (что, однако, не требовалось для получения полного балла за пункт). Заключаем, что

$$\hat{y}^{(m)} = B \sum_{k=0}^m (I - B)^k y \quad (3 \text{ балла})$$

Сумма выше – частичная сумма матричной геометрической прогрессии (так называемого [ряда Неймана](#)), поэтому ее можно свернуть⁴

$$\hat{y}^{(m)} = B(B^{-1}(I - (I - B)^{m+1}))y = (I - (I - B)^{m+1})y$$

Подставляя сюда B , получаем ответ в виде⁵

$$\hat{y}^{(m)} = \left(I - (I - X(X^\top X + \lambda I)^{-1} X^\top)^{m+1} \right) y \quad (1 \text{ балл})$$

(d)

Пусть мы считаем, что $y = X\beta + \varepsilon$ и ковариационная матрица ошибок имеет вид $\mathbb{V}\varepsilon = \Sigma$. Тогда ковариационная матрица y получается как

$$\mathbb{V}y = \mathbb{V}(X\beta + \varepsilon) = \mathbb{V}\varepsilon = \Sigma$$

⁴ Для таких матриц A , что $I - A$ обратима, верно равенство $\sum_{k=0}^m A^k = (I - A)^{-1}(I - A^{m+1})$. При этом ниоткуда не следует, что этот ряд сходится при $m \rightarrow \infty$.

⁵ То же можно было получить, раскрыв скобки в сумме и собрав результат по биному Ньютона.

Теперь, пользуясь свойством ковариационной матрицы⁶, получаем

$$\mathbb{V} \hat{y}^{(m)} = \mathbb{V} (I - (I - B)^{m+1}) y = (I - (I - B)^{m+1}) \Sigma (I - (I - B)^{m+1})^\top$$

Если дополнительно предположить, что $\Sigma = \sigma^2 I$, и заметить, что матрица B и, как следствие, матрица $I - (I - B)^{m+1}$ симметричны, то ответ упростится до вида (для полного балла такое упрощение в решении не требовалось)

$$\mathbb{V} \hat{y}^{(m)} = \sigma^2 (I - (I - B)^{m+1})^2 \quad (3 \text{ балла})$$

(e)

Прогноз одношаговой ridge-регрессии, как известно, смещен, поэтому естественно ожидать смещения и от комбинации таких моделей. Рассмотрим вектор ошибок итогового прогноза

$$e^{(m)} = \hat{y}^{(m)} - y = -(I - B)^{m+1} y = -(I - B)^{m+1} (X\beta + \varepsilon)$$

Тогда ожидаемое смещение⁷ прогноза имеет вид (в общем случае оно ненулевое)

$$\mathbb{E} e^{(m)} = -\mathbb{E} ((I - B)^{m+1} (X\beta + \varepsilon)) = -(I - B)^{m+1} X\beta \quad (3 \text{ балла})$$

(f)

Рассмотрим [сингулярное разложение](#) для матрицы X , то есть представим ее в виде $X = UDV^\top$, где U – ортогональная⁸ матрица $n \times n$, V – ортогональная матрица $k \times k$, а D – прямоугольная диагональная матрица $n \times k$, на главной диагонали которой стоят корни из собственных чисел матрицы $X^\top X$, а за пределами главной диагонали – нули (если различных собственных чисел окажется меньше k , то следующие числа на диагонали также будут нулевыми). Тогда получаем

$$\begin{aligned} B &= X (X^\top X + \lambda I)^{-1} X^\top = UDV^\top (VD^\top U^\top UDV^\top + \lambda I)^{-1} VD^\top U^\top = \\ &= UDV^\top (VD^\top DV^\top + \lambda VV^\top)^{-1} VD^\top U^\top = UDV^\top V (D^\top D + \lambda I)^{-1} V^\top VD^\top U^\top = \\ &= U \underbrace{D (D^\top D + \lambda I)^{-1} D^\top}_{\tilde{D}} U^\top \end{aligned}$$

⁶ А именно тем, что $\mathbb{V}(A\xi) = A\mathbb{V}(\xi)A^\top$.

⁷ В некоторых источниках смещение может определяться как $\mathbb{E}(y - \hat{y})$, а не как $\mathbb{E}(\hat{y} - y)$, однако, поскольку сути это не меняет, оба варианта оценивались полным баллом.

⁸ Ортогональными называются (вещественные) матрицы U , такие что $U^\top U = UU^\top = I$.

В записи выше матрица \tilde{D} также диагональная, имеет размер $n \times n$, причем на ее диагонали стоят числа $\frac{d_i^2}{d_i^2 + \lambda}$ вплоть до $i = k$, а затем – нули. Отсюда получаем

$$I - B = UU^\top - U\tilde{D}U^\top = U(I - \tilde{D})U^\top \implies (I - B)^{m+1} = U(I - \tilde{D})^{m+1}U^\top$$

Теперь рассмотрим вычисленное нами ранее смещение

$$-\mathbb{E} e^{(m)} = (I - B)^{m+1}X\beta = U(I - \tilde{D})^{m+1}U^\top UDV^\top\beta = U(I - \tilde{D})^{m+1}DV^\top\beta$$

На диагонали (при $i \leq k$) матрицы $(I - \tilde{D})^{m+1}$ стоят числа $\left(1 - \frac{d_i^2}{d_i^2 + \lambda}\right)^{m+1}$ (которые бесконечно малы при $m \rightarrow \infty$, поскольку $1 - \frac{d_i^2}{d_i^2 + \lambda} = \frac{\lambda}{\lambda + d_i^2} \in [0, 1)$ при всех $\lambda > 0$) и единицы при $i > k$. У матрицы D – наоборот – в этих строках нули.

$$(I - \tilde{D})^{m+1}D = \begin{pmatrix} \left(\frac{\lambda}{\lambda + d_1^2}\right)^{m+1} & & & & 0 \\ & \ddots & & & \\ & & \left(\frac{\lambda}{\lambda + d_k^2}\right)^{m+1} & & \\ \hline & & & 1 & 0 \\ & 0 & & & \ddots \\ & & & & & 1 \end{pmatrix} \begin{pmatrix} d_1 & 0 & & & \\ & \ddots & & & 0 \\ & & d_k & & \\ \hline & & & 0 & 0 \\ 0 & & & & \ddots \\ & & & & & 0 \end{pmatrix}$$

Из этого представления ясно, что

$$(I - \tilde{D})^{m+1}D \xrightarrow{m \rightarrow \infty} 0 \implies \mathbb{E} e^{(m)} \xrightarrow{m \rightarrow \infty} 0 \quad (4 \text{ балла})$$

(g)

При любом $\lambda > 0$ прогноз $\hat{y}^{(m)}$ линеен по y и притом его смещение при увеличении m стремится к нулю. В этот момент может возникнуть мысль о том, что прогноз $\hat{y}^{(m)}$ должен с ростом m стремиться к МНК-оценке \hat{y}_{OLS} . Попробуем это обосновать. Вычисленная ранее ковариационная матрица после замены через сингулярные числа примет вид

$$\begin{aligned} \mathbb{V} \hat{y}^{(m)} &= \sigma^2 (I - (I - B)^{m+1})^2 = \sigma^2 \left(I - U(I - \tilde{D})^{m+1}U^\top \right)^2 = \\ &= \sigma^2 \left(UU^\top - U(I - \tilde{D})^{m+1}U^\top \right)^2 = \sigma^2 U \left(I - (I - \tilde{D})^{m+1} \right)^2 U^\top \end{aligned}$$

Теперь выразим через U , D и V ковариационную матрицу МНК-оценок обычной линейной регрессии (предполагается, что матрица X имеет полный ранг, иначе МНК-оценка не будет существовать):

$$\begin{aligned} \mathbb{V} \hat{y}_{\text{OLS}} &= \mathbb{V} X\hat{\beta}_{\text{OLS}} = \mathbb{V} X(X^\top X)^{-1}X^\top(X\beta + \varepsilon) = \sigma^2 X(X^\top X)^{-1}X^\top = \\ &= \sigma^2 UDV^\top (VD^\top U^\top UDV^\top)^{-1}VD^\top U^\top = \sigma^2 UD(D^\top D)^{-1}D^\top U^\top \end{aligned}$$

Синяя матрица из записи выше состоит из k единиц, идущих по главной диагонали, а остальные ее элементы равны нулю. Что (не)удивительно, при $m \rightarrow \infty$ ровно такой же вид приобретает и синяя матрица из выражения для $\mathbb{V} \hat{y}^{(m)}$, что следует из последних рассуждений пункта f.

Получается, смещение прогноза $\hat{y}^{(m)}$ стремится к нулю, при этом он имеет линейный вид, а его ковариационная матрица сходится к матрице $\mathbb{V} \hat{y}_{\text{OLS}}$ при любом $\lambda > 0$. Тогда из соображений того, что МНК дает наиболее эффективную среди несмещенных оценок по теореме Гаусса–Маркова, в предпосылках этой теоремы имеет место сходимость

$$\hat{y}^{(m)} \xrightarrow{m \rightarrow \infty} \hat{y}_{\text{OLS}} = X (X^\top X)^{-1} X^\top y \quad (4 \text{ балла})$$

В задаче не предполагалось настолько формальное обоснование. Для получения полного балла достаточно было сказать, что прогноз сходится к МНК-оценке из соображений асимптотической несмещенности, постепенного возрастания дисперсии и линейности.

Дополнительный комментарий. Задача, решение которой только что было представлено, фактически описывает процедуру градиентного бустинга над линейными ridge-регрессиями. Эта тема уже освещалась в литературе, см., например, работы [Bühlmann & Yu, 2003] и [Tutz & Binder, 2007]. Алгоритм градиентного бустинга (хотя он и показал себя не с лучшей стороны в нашей задаче) является чрезвычайно мощным методом ансамблирования моделей машинного обучения, поэтому кажется достаточно интересным для рассмотрения.

Если участники не были знакомы с сингулярным разложением ранее, но заинтересовались им, автор может порекомендовать [визуализированное объяснение](#) и [более подробное и формальное обоснование](#) из открытого доступа.

Задача 3. (25 баллов)

Стать мракоборцем

①

Поскольку по условию переменные можно считать нормально распределенными, то пусть условные плотности $X(1, 2)$ имеют следующий вид

$$p_1(X) = p(X | y = 1) = \frac{1}{2\pi\sqrt{|\Sigma|}} e^{-\frac{1}{2}(X-\mu_1)^\top \Sigma^{-1}(X-\mu_1)}$$

$$p_2(X) = p(X | y = 2) = \frac{1}{2\pi\sqrt{|\Sigma|}} e^{-\frac{1}{2}(X-\mu_2)^\top \Sigma^{-1}(X-\mu_2)}$$

Теперь преобразуем $\mathbb{P}(y = 1 | X)$ – апостериорную вероятность обучения на мракоборца при условии на X – по формуле Байеса

$$\mathbb{P}(y = 1 | X) = \frac{p(X | y = 1) \overbrace{\mathbb{P}(y = 1)}^{\pi}}{p(X)} = \frac{\pi p_1(X)}{\pi p_1(X) + (1 - \pi)p_2(X)} = \frac{\frac{\pi}{1-\pi} \frac{p_1(X)}{p_2(X)}}{1 + \frac{\pi}{1-\pi} \frac{p_1(X)}{p_2(X)}}$$

$$\frac{\pi}{1-\pi} \frac{p_1(X)}{p_2(X)} = e^{\ln(\frac{\pi}{1-\pi})} \frac{2\pi\sqrt{|\Sigma|}}{2\pi\sqrt{|\Sigma|}} e^{-\frac{1}{2}((X-\mu_1)^\top \Sigma^{-1}(X-\mu_1) - (X-\mu_2)^\top \Sigma^{-1}(X-\mu_2))} = e^{D(X)}$$

Попробуем вычислить это в нашем конкретном численном примере

$$X = \begin{pmatrix} 60 \\ 4.2 \end{pmatrix}, \quad \mu_1 = \begin{pmatrix} 50 \\ 4 \end{pmatrix}, \quad \mu_2 = \begin{pmatrix} 40 \\ 3 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 100 & 3 \\ 3 & 1 \end{pmatrix} \implies \Sigma^{-1} = \frac{1}{91} \begin{pmatrix} 1 & -3 \\ -3 & 100 \end{pmatrix}$$

$$(X - \mu_1)^\top \Sigma^{-1} (X - \mu_1) = \frac{1}{91} (60 - 50)^2 - 2 \times \frac{3}{91} (60 - 50) (4.2 - 4) + \frac{100}{91} (4.2 - 4)^2 = \frac{92}{91}$$

$$(X - \mu_2)^\top \Sigma^{-1} (X - \mu_2) = \frac{1}{91} (60 - 40)^2 - 2 \times \frac{3}{91} (60 - 40) (4.2 - 3) + \frac{100}{91} (4.2 - 3)^2 = \frac{400}{91}$$

$$\implies D(X) = \ln\left(\frac{0.3}{1-0.3}\right) - \frac{1}{2} \left(\frac{92}{91} - \frac{400}{91}\right) = \ln \frac{3}{7} + \frac{22}{13} \approx 0.845 \implies e^{D(X)} \approx 2.328$$

Теперь, подставляя это в формулу условной вероятности, получаем

$$\mathbb{P}(y = 1 | X) = \frac{e^{D(X)}}{1 + e^{D(X)}} = \frac{2.328}{1 + 2.328} \approx 0.7$$

- Если участники давали (и разумно обосновывали) ответ только на основании средних значений, то за это ставилось **5 баллов**
- Если участники развивали мысль о формуле Байеса или условной вероятности, то за это ставилось **5 баллов**
- Если участники формализовали свои рассуждения в виде статистической модели, пробовали строить МНК-модель и получали ответ из $[0, 1]$, но никак не использовали ковариационные матрицы, то за это ставилось **10 баллов**
- Если студенты пробовали оценивать логит-модель, формируя и максимизируя соответствующее правдоподобие, то за это ставилось **15 баллов**

②

Вернемся к виду условной вероятности, который мы получили в пункте ①

$$\mathbb{P}(y = 1 | X) = \frac{e^{D(X)}}{1 + e^{D(X)}} = \frac{1}{1 + e^{-D(X)}}$$

Видно, что это логистически преобразованное выражение $D(X)$. Покажем, что перед нами на самом деле выражение для вероятности из логистической регрессии. Еще раз рассмотрим

$$D(X) = \ln\left(\frac{\pi}{1 - \pi}\right) - \frac{1}{2} \left((X - \mu_1)^\top \Sigma^{-1} (X - \mu_1) - (X - \mu_2)^\top \Sigma^{-1} (X - \mu_2) \right)$$

Поскольку ковариационные матрицы Σ для обоих классов используются одни и те же, то все элементы с квадратами и попарными произведениями взаимно уничтожаются, поэтому функция $D(X)$ оказывается **линейной** по X с точностью до константы, то есть имеет вид $D(X) = \beta_0 + X\beta$. А это означает, что перед нами в самом деле логит-модель (**5 баллов**)

$$\mathbb{P}(y = 1 | X) = \frac{e^{\beta_0 + X\beta}}{1 + e^{\beta_0 + X\beta}} = \frac{1}{1 + e^{-(\beta_0 + X\beta)}}$$

Задача 4. (25 баллов)

О взвешивании

(А)

Как известно, МНК-оценка параметров линейной регрессии имеет вид $\hat{\beta}_{\text{OLS}} = (X^\top X)^{-1} X^\top y$. В нашем случае матрицы X и y имеют блочный вид, поэтому

$$\begin{aligned} \hat{\beta}_{\text{OLS}} &= \left[\begin{pmatrix} X_1^\top & X_2^\top \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \right]^{-1} \begin{pmatrix} X_1^\top & X_2^\top \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = (X_1^\top X_1 + X_2^\top X_2)^{-1} (X_1^\top y_1 + X_2^\top y_2) = \\ &= (X_1^\top X_1 + X_2^\top X_2)^{-1} \left(\underbrace{X_1^\top X_1 (X_1^\top X_1)^{-1}}_I X_1^\top y_1 + \underbrace{X_2^\top X_2 (X_2^\top X_2)^{-1}}_I X_2^\top y_2 \right) = \\ &= (X_1^\top X_1 + X_2^\top X_2)^{-1} (X_1^\top X_1 \hat{\beta}_1 + X_2^\top X_2 \hat{\beta}_2) = \underbrace{(X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1}_{A_1} \hat{\beta}_1 + \\ &\quad + \underbrace{(X_1^\top X_1 + X_2^\top X_2)^{-1} X_2^\top X_2}_{A_2} \hat{\beta}_2 \end{aligned}$$

Осталось проверить, что $A_1 + A_2 = I$. Убедимся в этом непосредственно

$$\begin{aligned} A_1 + A_2 &= (X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 + (X_1^\top X_1 + X_2^\top X_2)^{-1} X_2^\top X_2 = \\ &= \underbrace{(X_1^\top X_1 + X_2^\top X_2)^{-1} (X_1^\top X_1 + X_2^\top X_2)}_I = I \quad (15 \text{ баллов}) \end{aligned}$$

(Б)

Теперь исследуем, как себя ведет GLS-оценка модели. Нам известно, что истинная ковариационная матрица ошибок имеет блочный вид

$$\Omega = \begin{pmatrix} \sigma_1^2 & 0 & & & \\ & \ddots & & & \\ 0 & & \sigma_1^2 & & \\ & & & \sigma_2^2 & 0 \\ 0 & & & & \ddots & \\ & & & & & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 I & 0 \\ 0 & \sigma_2^2 I \end{pmatrix} \implies \Omega^{-1} = \begin{pmatrix} \frac{1}{\sigma_1^2} I & 0 \\ 0 & \frac{1}{\sigma_2^2} I \end{pmatrix}$$

Получим аналитический вид GLS-оценки (для получения полного балла за задачу этого не требовалось). Обобщенный метод наименьших квадратов решает задачу

$$e^\top \Omega^{-1} e = (y - X\beta)^\top \Omega^{-1} (y - X\beta) \rightarrow \min_{\beta}$$

Вновь пользуясь векторно-матричным дифференцированием, выписываем необходимое условие экстремума в этой задаче⁹

$$\begin{aligned} \nabla (y - X\beta)^\top \Omega^{-1} (y - X\beta) = 0 &\iff \nabla (y^\top \Omega^{-1} y + \beta^\top X^\top \Omega^{-1} X \beta - 2\beta^\top X^\top \Omega^{-1} y) = 0 \\ 2X^\top \Omega^{-1} X \beta - 2X^\top \Omega^{-1} y = 0 &\implies \hat{\beta}_{\text{GLS}} = (X^\top \Omega^{-1} X)^{-1} X^\top \Omega^{-1} y \end{aligned}$$

Теперь еще раз воспользуемся блочной структурой матриц, данных в условии

$$\begin{aligned} \hat{\beta}_{\text{GLS}} &= \left[\begin{pmatrix} X_1^\top & X_2^\top \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma_1^2} I & 0 \\ 0 & \frac{1}{\sigma_2^2} I \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \right]^{-1} \begin{pmatrix} X_1^\top & X_2^\top \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma_1^2} I & 0 \\ 0 & \frac{1}{\sigma_2^2} I \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \\ &= \left(\frac{X_1^\top X_1}{\sigma_1^2} + \frac{X_2^\top X_2}{\sigma_2^2} \right)^{-1} \left(\frac{X_1^\top y_1}{\sigma_1^2} + \frac{X_2^\top y_2}{\sigma_2^2} \right) \quad (10 \text{ баллов}) \end{aligned}$$

В случае, если участник записывал обратную матрицу в знаменателе дроби или обращался с матрицами и векторами, как с числами, то за соответствующий пункт выставлялось не более 5 баллов.

⁹ И вновь оптимизируемая функция выпукла, поэтому необходимое условие оказывается одновременно и достаточным.