

## Прогнозирование успеваемости в вузе по результатам ЕГЭ

*В работе исследуются эконометрические зависимости среднего балла студентов первого курса от результатов их вступительных испытаний на примере факультета экономики ВШЭ. На основе выборки с учетом выбывших студентов, оценена модель бинарного выбора для вероятности выбытия как функции результатов ЕГЭ. Прогнозирование по сумме четырех ЕГЭ — математике, обществознанию, русскому языку и иностранному языку — в большинстве случаев дает несколько худшее качество прогноза, чем в модели без иностранного языка. Модели с использованием в качестве регрессоров результатов всех четырех ЕГЭ по отдельным дисциплинам имеют наилучшее качество прогноза. Среди дисциплин ЕГЭ наибольшее влияние оказывает математика.*

**Ключевые слова:** прием в вузы, единый государственный экзамен, прогноз успеваемости.

### 1. Введение

В Российской Федерации прием в государственные высшие учебные заведения осуществляется на конкурсной основе. Конкурсный характер отбора реализуется, в частности, через систему вступительных испытаний. С 2009 года в качестве вступительных испытаний повсеместно используются результаты ЕГЭ. Поступившие в 2009 году студенты закончили свой первый год обучения в 2010 году, и появилась возможность использовать статистические данные для исследования связи результатов ЕГЭ и последующей успеваемости в вузе. Такие исследования важны для анализа эффективности процедуры отбора студентов в вузы, которая не может считаться эффективной, если она не в состоянии обеспечить выбор наиболее предрасположенных к будущей профессии претендентов.

Значительный опыт в области изучения связи результатов довузовского тестирования и последующей успеваемости в университете накоплен в США, где многие вузы учитывают при приеме результаты стандартизированных тестов, таких как SAT или ACT, проводимых частными организациями в течение многих десятилетий. Способность SAT прогнозировать средний балл за первый год обучения замечена давно. В работе (Fishman, Pasanella, 1960) был сделан обзор 147 исследований, в них множественный коэффициент корреляции  $R$  (квадратный корень из коэффициента детерминации) находился в диапазоне от 0.34 до 0.82, средний показатель был равен 0.61. В статье (Burton, Ramist, 2001) можно найти обзор работ по изучению влияния результатов SAT на академическую успеваемость, выполненных в последние два десятилетия двадцатого столетия. В одном из последних исследований (Kobrin et al., 2008) использовались данные более чем по 150 тыс. учащихся из 110 колледжей и университетов: при этом множественный коэффициент корреляции в среднем составил 0.35. Исследования в этой области ведутся и отдельными университетами по собственным данным.

Продолжаются дискуссии о том, в какой степени предсказательную способность результатов стандартизированных тестов можно объяснить их корреляцией с другими переменными, такими как социально-экономический статус учащегося (см., например, (Rothstein, 2004; Sackett et al., 2009)). Высокий доход и хорошее образование родителей студента могут способствовать его успешной учебе в вузе, однако немногие готовы открыто поддержать идею осуществлять отбор по этим критериям. В то же время, мало сторонников найдется у предложений понижать баллы вступительных испытаний тем абитуриентам, успеху которых способствовал высокий социально-экономический статус, с целью обеспечения равных возможностей.

В России в настоящее время в качестве критерия, по которому осуществляется ранжирование абитуриентов, используется простое (с равными весами) суммирование баллов, полученных на ЕГЭ по нескольким предметам. Перечень предметов из трех или четырех вступительных испытаний по каждому образовательному направлению определяется Министерством образования и науки Российской Федерации. В их число обязательно входит экзамен по русскому языку и профильному общеобразовательному предмету.

Очевидно, что выбор дисциплин вступительных испытаний и весовых коэффициентов для разных предметов при суммировании баллов играет важнейшую роль в процедуре отбора абитуриентов. Выбор предметов, выносимых на экзамен, определяет набор тестируемых специфических способностей, знаний и навыков абитуриентов. Весовые коэффициенты определяют вклад отдельных дисциплин в итоговый критерий.

Выбор оптимального набора дисциплин и их весов может быть основан на регрессионном анализе модели, в которой в качестве объясняемой переменной выступает интегральный показатель вузовской успеваемости (средний балл или рейтинг студентов), а в роли объясняющих факторов — набор результатов вступительных испытаний. Коэффициенты при регрессорах (оценках за тот или иной предмет) могут быть основой для нахождения весов конкретной дисциплины в суммарном показателе.

При неоптимальном выборе дисциплин вступительных испытаний и их весов происходит неблагоприятный отбор абитуриентов. В вуз могут попасть такие учащиеся, средняя успеваемость которых ниже, чем ожидаемая успеваемость некоторых абитуриентов, отсеянных во время приемной кампании<sup>1</sup>.

## 2. Статистические характеристики результатов ЕГЭ абитуриентов

По решению Минобрнауки, для направления «Экономика» в перечень предметов вступительных испытаний входят четыре дисциплины — математика, обществознание, русский язык, иностранный язык. Число испытаний должно быть не менее трех, причем русский язык и математика относятся к числу обязательных. Таким образом, вуз самостоятельно выбирает, использовать ли для отбора абитуриентов ЕГЭ по обществознанию или иностранному языку, или оба экзамена.

<sup>1</sup> Ожидаемая успеваемость абитуриентов, которые не поступили в вуз на бюджетное место и отказались от платного обучения, может оцениваться по фактической успеваемости студентов с аналогичными баллами ЕГЭ, обучающихся на платной основе. Если студентов-платников с такими баллами нет (в достаточном количестве), то ожидаемая успеваемость может быть оценена экстраполяцией регрессионного уравнения.

В 2009 г. для поступления по результатам вступительных испытаний на факультет экономики в Высшую школу экономики (ВШЭ, г. Москва) необходимо было предоставить результаты ЕГЭ по всем четырем предметам. Отбор абитуриентов без льгот на оставшиеся бюджетные места проводился по критерию наибольшей суммы баллов ЕГЭ. Без вступительных испытаний зачислялись победители и призеры некоторых олимпиад.

На рисунке 1 приведены эмпирические функции распределения результатов ЕГЭ по отдельным дисциплинам для трех различных групп людей:

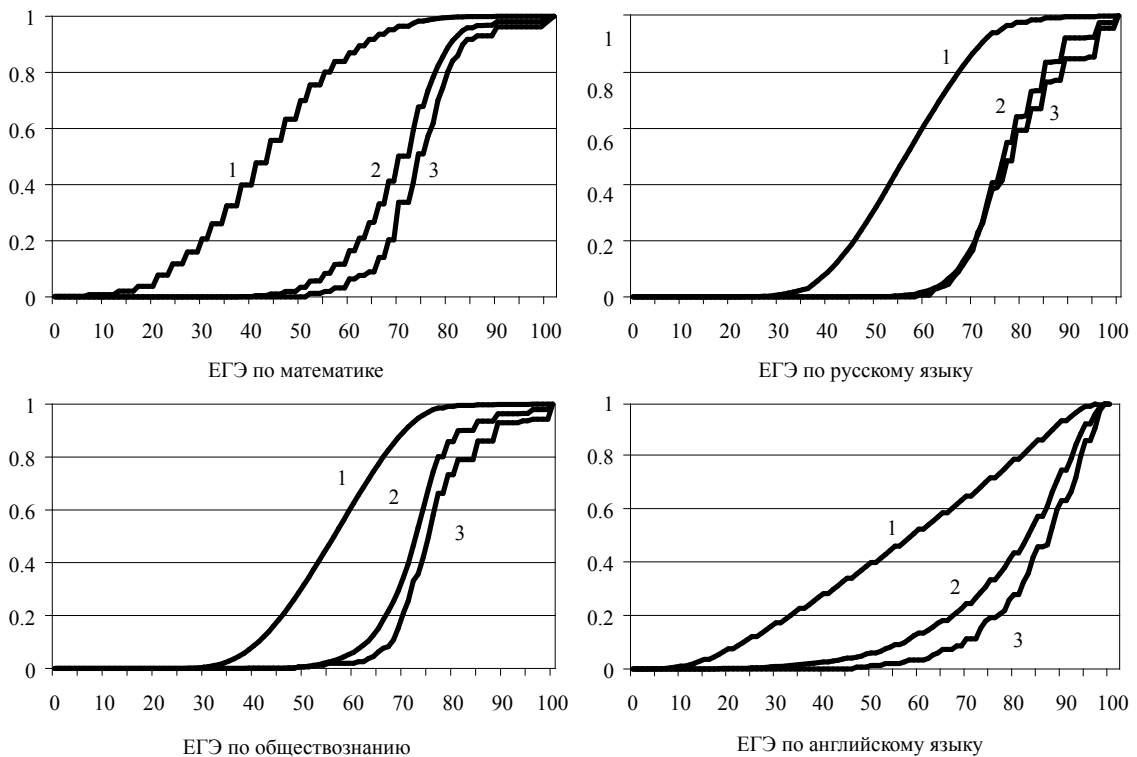
1) всех учащихся, сдававших экзамен в 2009 г. (около 1 млн человек сдавали экзамены по русскому языку и математике, 450 тыс. — по обществознанию, 77 тыс. — по английскому языку);

2) абитуриентов ВШЭ, подавших заявления о поступлении на факультет экономики (2856 человек);

3) студентов факультета экономики, зачисленных в 2009 г. на факультет экономики по результатам ЕГЭ и закончивших первый курс в 2010 г. (157 человек)<sup>2</sup>.

Данные о результатах ЕГЭ абитуриентов ВШЭ были взяты с официального сайта ВШЭ (<http://www.hse.ru/>) и находились в открытом доступе. Данные о распределении тестовых баллов ЕГЭ по всей стране в 2009 г. находились на официальном информационном портале ЕГЭ (<http://www.ege.edu.ru/>).

Прогнозирование успеваемости в вузе по результатам ЕГЭ



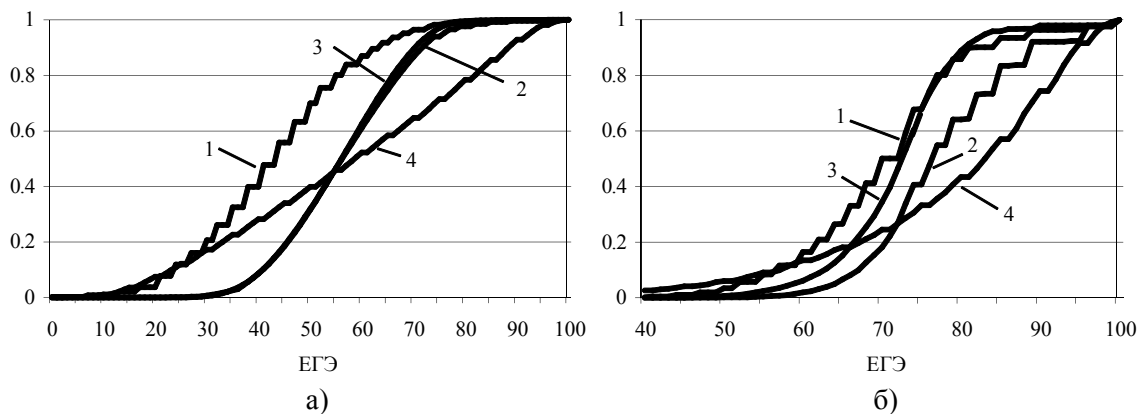
**Рис. 1.** Эмпирические функции распределения баллов по дисциплинам ЕГЭ для различных групп: (1) — все выпускники, (2) — абитуриенты, (3) — студенты

<sup>2</sup> Все данные относятся к ВШЭ в Москве, без учета филиалов.

Максимально возможное количество баллов по каждому предмету — 100. Сравнение распределений групп 1 и 2 позволяет судить о том, выпускники какого качества претендовали на поступление. Медиана распределения баллов по математике, русскому языку и обществознанию для абитуриентов соответствует примерно 95%-ому перцентилю общего распределения, медиана по английскому языку для абитуриентов — 82%-ому перцентилю общего распределения. Из тех, кто сдавал ЕГЭ по иностранному языку, подавляющее большинство выбрали английский язык, поэтому сравнение проведено только для английского языка.

Сравнение распределений групп 2 и 3 характеризует селективность того или иного направления. Из общего числа абитуриентов студентами стали те абитуриенты, которые либо прошли по конкурсу на бюджетное место, либо имели льготы быть зачисленными на бюджет при числе баллов меньше проходного, либо обучаются на платной основе. Распределение баллов для студентов не обязательно находится правее распределения для абитуриентов. Многие абитуриенты подавали заявления параллельно на разные программы и в различные вузы. Абитуриенты с высокими баллами могли выбирать между программами, на которые они прошли по конкурсу, а обучающиеся на платной основе могли принимать во внимание стоимость обучения. Поэтому возможна ситуация, когда распределение студентов находится левее распределения абитуриентов. Как видно из рис. 1, распределения ЕГЭ студентов доминируют над распределениями абитуриентов.

На рисунке 2 сведены вместе эмпирические функции распределения баллов по различным дисциплинам ЕГЭ для всех выпускников страны (рис. 2, а) и абитуриентов (рис. 2, б). Явно заметно отличие распределений баллов по математике и английскому языку от распределений по обществознанию и русскому языку. Сравнивая правые части распределений, легко заметить, что высокие баллы имеет большая доля сдавших английский язык, чем математику. Так, 80 и выше баллов получили по математике 0.5% сдававших этот экзамен, по обществознанию — менее 1%, по русскому языку — около 2%, по английскому языку — более 20%. Среди результатов с 75 и более баллами у абитуриентов также преобладает английский язык, менее всего представлена математика.



**Рис. 2.** Эмпирические функции распределения баллов по различным дисциплинам ЕГЭ для всех выпускников (а) и абитуриентов (б):  
 (1) — математика, (2) — русский язык,  
 (3) — обществознание, (4) — иностранный язык

Описательные статистики баллов ЕГЭ по предметам для 157 студентов, поступивших по ЕГЭ, представлены в таблице 1.

Одномерные распределения и характеристики не дают представления о том, как взаимосвязаны результаты ЕГЭ по различным дисциплинам. Для многомерного анализа оказалось возможным использовать только данные приемной комиссии ВШЭ. На рисунке 3 в виде нижней треугольной матрицы показаны двумерные гистограммы баллов ЕГЭ абитуриентов для всевозможных пар из четырех предметов. В таблицах 2 и 3 даны корреляционные матрицы результатов ЕГЭ для абитуриентов и студентов, соответственно.

В таблице 4 представлены результаты анализа главных компонент для распределения ЕГЭ абитуриентов. Как известно, смысл главных компонент часто трудно интерпретировать (Айвазян, Мхитарян, 1998). В рассматриваемом случае первая главная компонента с приблизительно равными весами предметов может быть интерпретирована как фактор, характеризующий общие, универсальные способности абитуриентов и их усердие в учебе. На нее приходится 56% суммарной дисперсии. Вторая компонента с положительными долями русского и иностранного языков и отрицательными для математики и обществознания скорее характеризует лингвистические способности. Третью и четвертую компоненты интерпретировать труднее.

**Таблица 1.** Описательные статистики баллов ЕГЭ по предметам, студенты

|                        | Математика | Русский язык | Обществознание | Иностранный язык |
|------------------------|------------|--------------|----------------|------------------|
| Среднее значение       | 75.03      | 79.68        | 76.88          | 84.93            |
| Максимальное значение  | 100        | 100          | 100            | 100              |
| Минимальное значение   | 52         | 59           | 49             | 47               |
| Стандартное отклонение | 8.82       | 10.00        | 9.25           | 11.30            |

**Таблица 2.** Корреляционная матрица результатов ЕГЭ, абитуриенты

| 2856 наблюдений  | Математика | Русский язык | Обществознание | Иностранный язык |
|------------------|------------|--------------|----------------|------------------|
| Математика       | 1.000      |              |                |                  |
| Русский язык     | 0.409      | 1.000        |                |                  |
| Обществознание   | 0.477      | 0.450        | 1.000          |                  |
| Иностранный язык | 0.349      | 0.428        | 0.352          | 1.000            |

**Таблица 3.** Корреляционная матрица результатов ЕГЭ, студенты

| 157 наблюдений   | Математика | Русский язык | Обществознание | Иностранный язык |
|------------------|------------|--------------|----------------|------------------|
| Математика       | 1.000      |              |                |                  |
| Русский язык     | 0.297      | 1.000        |                |                  |
| Обществознание   | 0.410      | 0.449        | 1.000          |                  |
| Иностранный язык | 0.418      | 0.503        | 0.489          | 1.000            |

Прогнозирование успеваемости в вузе по результатам ЕГЭ

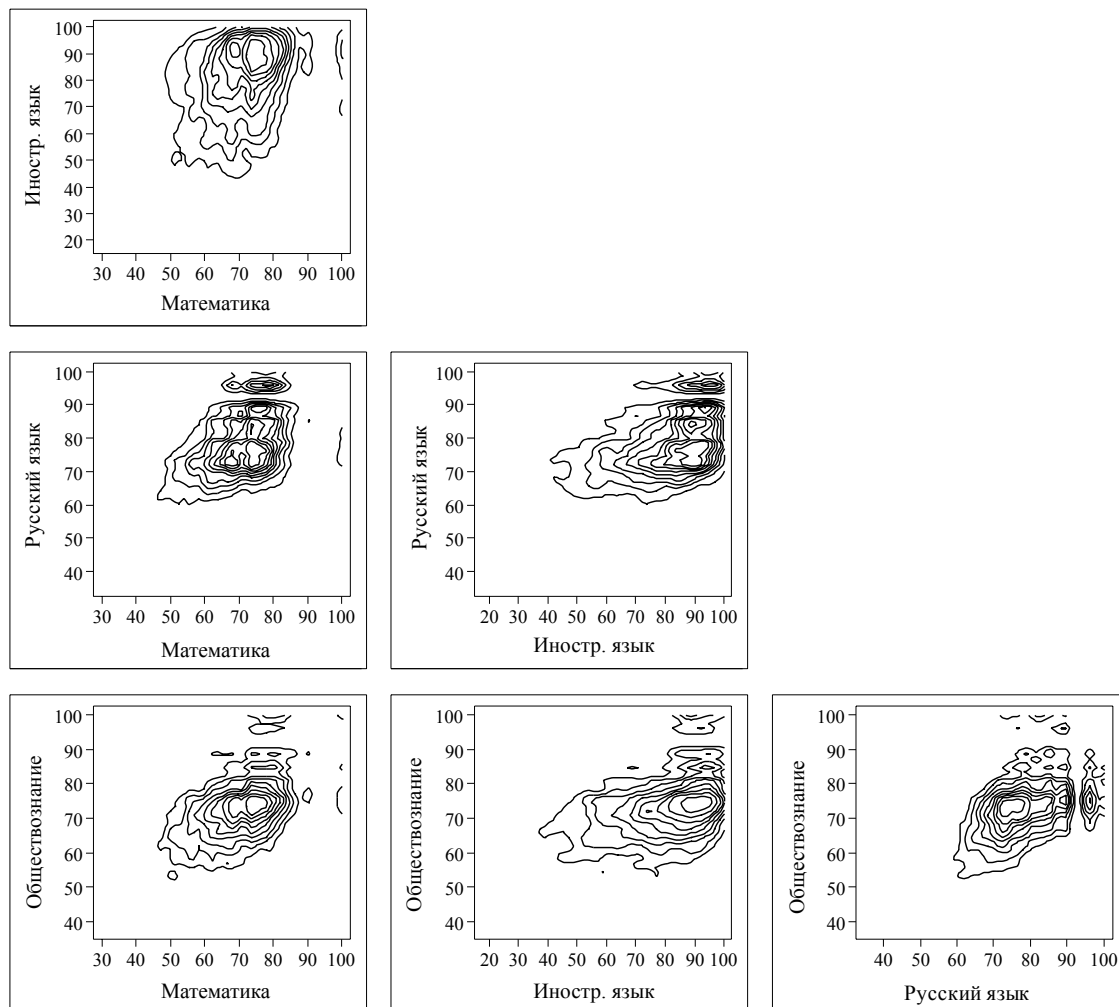


Рис. 3. Матрица двумерных гистограмм ЕГЭ, децильные контуры

Таблица 4. Анализ главных компонент для многомерного распределения результатов ЕГЭ, абитуриенты

| Собственные числа |          |  |  |
|-------------------|----------|--|--|
|                   | Значение | Относительная доля суммарной дисперсии | Накопленная относительная доля суммарной дисперсии |
| 1                 | 2.236    | 0.559                                  | 0.559  |
| 2                 | 0.696    | 0.174                                  | 0.733  |
| 3                 | 0.560    | 0.140                                  | 0.873  |
| 4                 | 0.508    | 0.127                                  | 1.000  |

| Переменная       | Собственные векторы |        |        |        |
|------------------|---------------------|--------|--------|--------|
|                  | PC 1                | PC 2   | PC 3   | PC 4   |
| Математика       | 0.501               | -0.474 | 0.541  | -0.481 |
| Русский язык     | 0.515               | 0.196  | -0.704 | -0.447 |
| Обществознание   | 0.514               | -0.416 | -0.199 | 0.723  |
| Иностранный язык | 0.467               | 0.751  | 0.415  | 0.213  |

### 3. Связь результатов ЕГЭ и вузовской успеваемости на первом году обучения

Показателем, характеризующим общую успеваемость студента ВШЭ, является рейтинг. В общих чертах, рейтинг формируется как взвешенная сумма оценок, где весами выступают величины кредитов учебной нагрузки по данному предмету. В Высшей школе экономики используется 10-балльная шкала оценок — чем выше балл, тем лучше успеваемость. Максимально возможное значение рейтинга равно 600 (10 баллов, умноженные на 60 кредитов годовой нагрузки). Сведения о текущем рейтинге студентов ВШЭ доступны на Интернет-портале ВШЭ (<http://www.hse.ru/>) на странице факультетов, данные обновляются один раз в полугодие. По итогам всего учебного года формируется кумулятивный рейтинг.

Для удобства далее используются значения рейтинга, поделенные на 6, таким образом, шкала рейтинга приводится в соответствие 100-балльной шкале ЕГЭ. Рейтинг, являющийся взвешенным средним баллом, и собственно средний балл, рассчитываемый как среднее арифметическое, имеют корреляцию 0.96. Из-за того, что термин «рейтинг» в вышеизложенном смысле используется не повсеместно, далее вместо него будет использоваться «средний балл».

Рассматривались три спецификации уравнений регрессии результатов ЕГЭ на средний балл. В первой модели в качестве регрессора выступает средний балл ЕГЭ, усредненный по четырем дисциплинам вступительных испытаний. Именно этот показатель использовался для конкурсного отбора в ВШЭ по направлению «Экономика». Вторая спецификация отличается тем, что средний балл ЕГЭ рассчитывался по трем предметам — математике, русскому языку и обществознанию. В моделях 1 и 2 ЕГЭ различные дисциплины имеют одинаковый вес при суммировании, как предусмотрено правилами приема. Однако несомненный интерес представляет вопрос о том, какая линейная комбинация ЕГЭ обладает наибольшей прогнозной способностью. Для ответа на этот вопрос рассматривалась модель 3, в которой регрессорами выступали баллы по отдельным предметам.

Регрессионные связи между рассматриваемыми переменными характеризуются негассовостью возмущений, присутствием относительно редких, но значительных, отклонений — выбросов. Метод наименьших квадратов чувствителен к выбросам, что обусловлено квадратичным характером оптимизационного критерия. Поэтому, наряду с обычным методом наименьших квадратов, применялись методы оценивания, менее чувствительные к большим возмущениям — медианная регрессия и робастная регрессия (Rousseeuw, Leroy, 1987).

В квантильной регрессии для квантиля 0.5 (медианы) минимизируется сумма абсолютных отклонений. Это объясняет устойчивость медианной регрессии к выбросам зависимой переменной. Чувствительность к выбросам регрессоров медианная регрессия не устраняет, но в данном случае экстремальных отклонений зависимых переменных в выборке нет.

В робастной регрессии используется многоэтапная процедура нахождения оценок (Rousseeuw, Leroy, 1987). После оценивания методом наименьших квадратов, для каждого наблюдения вычисляется показатель (Cook's  $D$ ), который измеряет среднеквадратичную разницу между прогнозными значениями в полной выборке и при исключении данного наблюдения. Наблюдения с  $D > 1$  исключаются. Затем совершается последовательность итераций: оценивается регрессия, для каждого наблюдения рассчитываются весовые коэффициенты, снова оценивается регрессия. Весовые коэффициенты рассчитываются так, что, чем больше ошибка прогноза, тем ниже вес наблюдения. Процедура нахождения оценок заканчивается после того, как изменение весов после очередной итерации становится ниже порогового значения.

Также рассматривались оценки МНК после удаления пяти выбросов, для которых ошибка прогноза превышала медианную ошибку более чем в четыре раза.

В таблицах 5 и 6 представлены оценки трех моделей методом наименьших квадратов, медианной регрессией, робастной регрессией (по 157 наблюдениям) и методом наименьших квадратов с удаленными наблюдениями (по 152 наблюдениям).

Коэффициент детерминации  $R^2$  является наиболее распространенной характеристикой объясняющей силы МНК-регрессии в эконометрической практике (в практике измерений в образовании также распространено применение множественного коэффициента корреляции  $R$ ). Для медианной регрессии аналогичным показателем является псевдо- $R^2$ . Сравнение по этому критерию моделей 1 и 2 показывает, что наилучшее качество подгонки дает модель 2 при оценке методом наименьших квадратов, медианной и робастной регрессией, модель 1 оказывается лучше при использовании МНК после удаления выбросов.

Таким образом, исключение иностранного языка в трех из четырех случаев улучшило качество прогноза. С учетом того, что иностранный язык сдают намного меньше выпускников школ, чем три остальные предмета, вполне допустимо, что отбор по трем дисциплинам позволил бы отобрать более подготовленных студентов. В этой гипотетической ситуации состав студентов, прошедших по конкурсу на бюджетные места, оказался бы иным. Насколько велики потенциальные искажения, можно оценить, сравнивая списки 100 лучших абитуриентов по различным критериям отбора. Из базы данных абитуриентов были выбраны 250 лучших результатов по сумме четырех ЕГЭ. Для них также была рассчитана сумма баллов по трем ЕГЭ. На рисунке 4 по горизонтали отложен ранг абитуриента по сумме трех ЕГЭ, по вертикали — ранг по сумме четырех ЕГЭ. Внутри квадрата со сторонами 0–100 находятся 84 абитуриента, которые по обоим критериям оказались в сотне лучших. Соответственно, 16 человек (16%) выбыли из списка. При отборе 50 лучших расхождение составило 10 человек (20%).

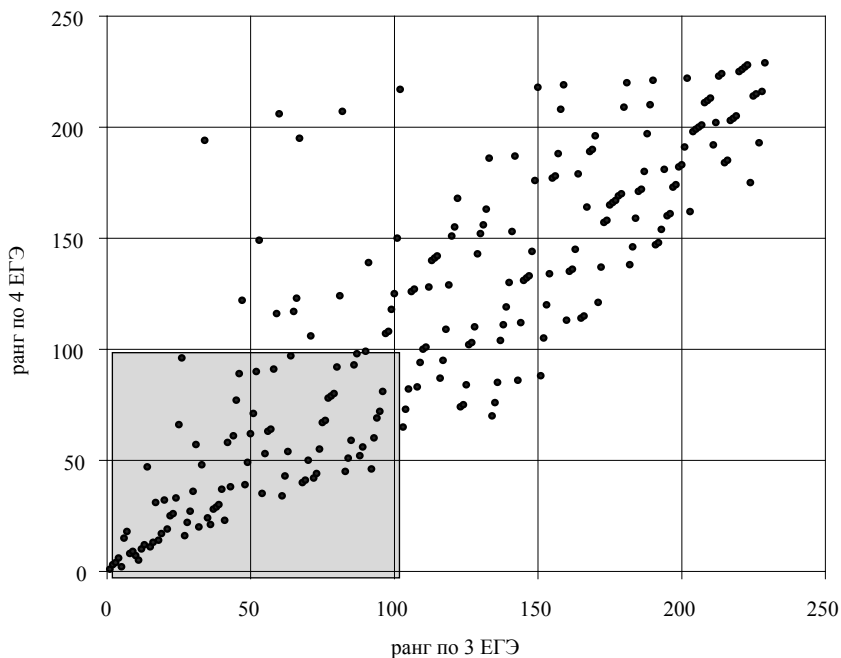
Модель 3, в которой в число регрессоров включены все четыре ЕГЭ по отдельности, по понятным причинам превосходит модели 1 и 2 по качеству прогноза. Коэффициенты модели имеют смысл приращений ожидаемого значения среднего балла в вузе при увеличении балла ЕГЭ по данному предмету на единицу. Чем выше коэффициент, тем ценнее баллы по соответствующему предмету. Наибольшим коэффициентом при всех методах оценивания обладает математика, далее идут обществознание, русский язык и иностранный язык. Коэффициент при иностранном языке незначим в трех из четырех регрессий.



**Таблица 5.** Связь среднего балла в вузе и баллов ЕГЭ, метод наименьших квадратов. Зависимая переменная — средний балл после 1 курса

| Регрессоры                       | МНК регрессии       |                     |                     | МНК регрессии (без выбросов) |                     |                     |
|----------------------------------|---------------------|---------------------|---------------------|------------------------------|---------------------|---------------------|
|                                  | 1                   | 2                   | 3                   | 1                            | 2                   | 3                   |
| Константа                        | -37.747<br>(9.271)  | -39.096<br>(9.818)  |                     | -36.546<br>(8.842)           | -36.954<br>(9.500)  | -41.580<br>(9.165)  |
| Средний балл ЕГЭ<br>(4 экзамена) | 1.222***<br>(0.117) |                     |                     | 1.208***<br>(0.112)          |                     |                     |
| Средний балл ЕГЭ<br>(3 экзамена) |                     | 1.270***<br>(0.127) |                     |                              | 1.244***<br>(0.123) |                     |
| Математика                       |                     |                     | 0.568***<br>(0.149) |                              |                     | 0.599***<br>(0.135) |
| Русский язык                     |                     |                     | 0.271**<br>(0.106)  |                              |                     | 0.215*<br>(0.102)   |
| Обществознание                   |                     |                     | 0.308**<br>(0.127)  |                              |                     | 0.296*<br>(0.119)   |
| Иностр. язык                     |                     |                     | 0.163*<br>(0.092)   |                              |                     | 0.188<br>(0.088)    |
| $R^2$                            | 0.362               | 0.363               | 0.384               | 0.402                        | 0.396               | 0.430               |

В скобках указаны робастные стандартные погрешности оценок. \*, \*\*, \*\*\* обозначают значимость на 10%, 5% и 1%-ном уровнях, соответственно.



**Рис. 4.** Ранги по сумме трех и четырех ЕГЭ

Прогнозирование успеваемости в вузе по результатам ЕГЭ

**Таблица 6.** Связь среднего балла в вузе и баллов ЕГЭ, робастная регрессия и медианная регрессия. Зависимая переменная — средний балл после 1 курса

| Регрессоры                       | Робастная регрессия |                     |                     | Медианная регрессия |                     |                     |
|----------------------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
|                                  | 1                   | 2                   | 3                   | 1                   | 2                   | 3                   |
| Константа                        | -37.098<br>(10.812) | -39.626<br>(10.800) |                     | -25.743<br>(14.440) | -39.905<br>(14.042) | -32.578<br>(14.379) |
| Средний балл ЕГЭ<br>(4 экзамена) | 1.215***<br>(0.136) |                     |                     | 1.074***<br>(0.182) |                     |                     |
| Средний балл ЕГЭ<br>(3 экзамена) |                     | 1.278***<br>(0.139) |                     |                     | 1.290***<br>(0.181) |                     |
| Математика                       |                     |                     | 0.624***<br>(0.129) |                     |                     | 0.660***<br>(0.169) |
| Русский язык                     |                     |                     | 0.250*<br>(0.120)   |                     |                     | 0.106<br>(0.165)    |
| Обществознание                   |                     |                     | 0.314*<br>(0.133)   |                     |                     | 0.306*<br>(0.181)   |
| Иностр. язык                     |                     |                     | 0.144<br>(0.112)    |                     |                     | 0.127<br>(0.153)    |
| $R^2$                            | 0.340               | 0.352               | 0.376               |                     |                     |                     |
| Псевдо- $R^2$                    |                     |                     |                     | 0.208               | 0.218               | 0.237               |

В скобках указаны робастные стандартные погрешности оценок, \*, \*\*, \*\*\* обозначают значимость на 10%, 5% и 1%-ном уровнях, соответственно.

#### 4. Оценивание моделей с учетом студентов, выбывших из учебного процесса

Регрессии, которые были рассмотрены выше, анализировались на основе выборки, включающей только окончивших первый курс студентов. Однако значительное число студентов выбывает из учебного процесса прежде всего из-за неуспеваемости во время первого года обучения. Поэтому определенный интерес представляет прогнозирование вероятности отчисления, а также оценивание связи ЕГЭ и среднего балла по полной выборке студентов, а не только в подвыборке студентов, окончивших первый курс.

Для оценивания влияния результатов ЕГЭ на вероятность отчисления использовалась логит-модель бинарного выбора (см., например, (Магнус и др., 2007)). Исходная выборка из 157 студентов, поступивших по ЕГЭ и окончивших первый курс, была дополнена 58 студентами, которым был присвоен рейтинг после первого семестра, но которые не получили его по итогам всего года. Предполагалось, что выбывшие студенты были отчислены по неуспеваемости. В таблице 7 представлены результаты оценивания модели. Для удобства интерпретации, вместо собственно коэффициентов модели, приведены так называемые маргинальные эффекты при средних значениях регрессоров. Коэффициенты имеют смысл приращения вероятности быть отчисленным при добавлении к среднему значению

регрессора единицы. В модели 3 значимым оказался эффект при математике. Это вполне объяснимо традиционными трудностями освоения дисциплин математического профиля в вузе.

**Таблица 7.** Маржинальные эффекты в модели бинарного выбора для вероятности отчисления после 1 курса

| Регрессоры                            | Логит-регрессия        |                         |                        |
|---------------------------------------|------------------------|-------------------------|------------------------|
|                                       | 1                      | 2                       | 3                      |
| Средний балл ЕГЭ (4 экзамена)         | -0.0288***<br>(0.0045) |                         |                        |
| Средний балл ЕГЭ (3 экзамена)         |                        | -0.0311***<br>(0.00466) |                        |
| Математика                            |                        |                         | -0.0156***<br>(0.0036) |
| Русский язык                          |                        |                         | -0.0039<br>(0.0035)    |
| Обществознание                        |                        |                         | -0.0071<br>(0.0056)    |
| Иностранный язык                      |                        |                         | -0.0032<br>(0.0024)    |
| Псевдо- $R^2$                         | 0.301                  | 0.309                   | 0.330                  |
| Логарифм функции псевдо-правдоподобия | -87.569                | -86.635                 | -84.002                |

В скобках указаны робастные стандартные ошибки (стандартные ошибки производных рассчитывались дельта-методом). \*, \*\*, \*\*\* обозначают значимость на 10%, 5% и 1%-ном уровнях, соответственно.

Представление о том, как могла бы выглядеть связь между переменными в полной выборке из 215 студентов, было получено в результате оценивания тобит-модели (см., например, (Магнус и др., 2007)). Выборка считалась цензурированной ниже уровня в 32 балла, всем выбывшим студентам был присвоен этот граничный балл. Выбор данного значения обусловлен тем, что именно здесь происходит резкий спад эмпирической плотности. Ниже порога находятся еще три человека со средними баллами около 20 и один студент с 26 баллами, которых также можно отнести к не очень успешным студентам.

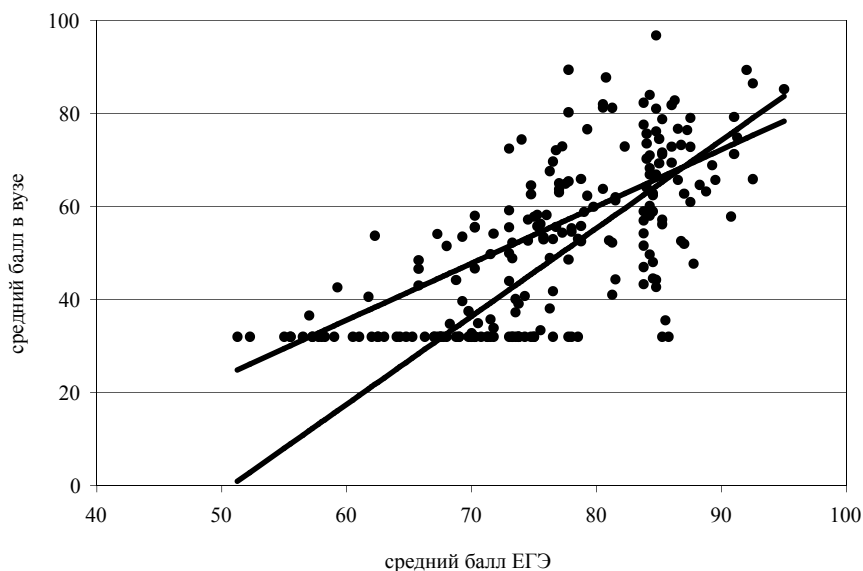
Оценки тобит-модели представлены в таблице 8. Коэффициенты моделей заметно отличаются от значений, полученных ранее методами, не учитывающими цензурирование. Тем самым проявляется явление смещения оценок, вызванное селективностью выборки. На рисунке 5 приведены линии регрессии МНК и тобит-модели для спецификации 1. Однако к данным оценкам следует относиться с осторожностью. Из-за негауссовости возмущений оценки тобит-модели не являются состоятельными. Другим ограничивающим надежность оценок фактором является неполнота информации о причинах выбытия. Для отчисления достаточно иметь неудовлетворительную оценку по одному предмету, а средний балл характеризует среднюю успеваемость. Выбытие студентов также могло происходить по причинам, отличным от академической задолженности.

Прогнозирование успеваемости в вузе по результатам ЕГЭ

**Таблица 8.** Связь среднего балла в вузе и баллов ЕГЭ, тобит-регрессия.  
Зависимая переменная — средний балл после 1 курса

| Регрессоры                            | Тобит-регрессия     |                     |                      |
|---------------------------------------|---------------------|---------------------|----------------------|
|                                       | 1                   | 2                   | 3                    |
| Константа                             | -96.060<br>(11.016) | -99.260<br>(11.160) | -98.960<br>(10.788)  |
| Средний балл ЕГЭ (4 экзамена)         | 1.892***<br>(0.138) |                     |                      |
| Средний балл ЕГЭ (3 экзамена)         |                     | 1.981***<br>(0.144) |                      |
| Математика                            |                     |                     | 0.896***<br>(0.144)  |
| Русский язык                          |                     |                     | 0.370***<br>(0.117)  |
| Обществознание                        |                     |                     | 0.438 ***<br>(0.151) |
| Иностранный язык                      |                     |                     | 0.266**<br>(0.109)   |
| Псевдо- $R^2$                         | 0.100               | 0.099               | 0.107                |
| Логарифм функции псевдо-правдоподобия | -681.577            | -681.809            | -675.834             |

В скобках указаны робастные стандартные погрешности оценок. \*, \*\*, \*\*\* обозначают значимость на 10%, 5% и 1%-ном уровнях, соответственно.



**Рис. 5.** Линии регрессии тобит-модели (сплошная линия) и МНК-модели (пунктир)

## 5. Заключение

Определение перечня предметов вступительных испытаний и способа формирования критерия из результатов отдельных экзаменов может опираться на статистический анализ связей между результатами ЕГЭ и показателями успеваемости в вузе. Эти связи специфичны для направления, по которому ведется обучение, а также зависят от особенностей образовательной программы того или иного вуза. Унифицированное определение дисциплин вступительных испытаний для крупных направлений профессиональной подготовки и негибкость при определении весовых коэффициентов способны привести к снижению эффективности отбора.

В работе исследовались статистические связи среднего балла студентов первого курса с результатами вступительных испытаний по данным приема на факультет экономики ВШЭ (г. Москва) в 2009 г. Из-за того, что регрессионные связи характеризуются присутствием небольшого числа значительных отклонений, наряду с обычным методом наименьших квадратов применялись методы оценивания, менее чувствительные к большим возмущениям — медианная и робастная регрессии.

Сумма баллов ЕГЭ по математике, русскому и иностранному языкам и обществознанию является неплохим инструментом прогноза среднего балла в вузе по итогам первого года. Однако в большинстве регрессий несколько лучший прогноз демонстрирует сумма трех ЕГЭ: по математике, обществознанию и русскому языку.

В регрессиях с использованием результатов отдельных экзаменов удастся достичь наилучшего качества прогноза. Вклады ЕГЭ по различным дисциплинам в объяснение вузовского среднего балла различаются. Наибольшее влияние имеет математика, за ней следуют обществознание и русский язык. Вклад иностранного языка оказался наименьшим, а в некоторых моделях — незначимым, что может свидетельствовать об избыточности этой переменной при условии, что имеются другие, более качественные предсказатели.

С использованием данных о вступительных испытаниях выбывших студентов, с помощью модели бинарного выбора была оценена вероятность отчисления как функция результатов ЕГЭ. Установлено, что наибольшее влияние на вероятность отчисления оказывает ЕГЭ по математике.

Полученные количественные результаты основаны на данных одного года и отражают специфику программы подготовки бакалавров на факультете экономики ВШЭ. Для выводов об устойчивости качественных результатов необходимо опираться на данные о нескольких группах студентов, поступивших на идентичную программу в разные годы. Однако есть серьезные основания предположить, что получение вузами свободы в выборе весов различных дисциплин ЕГЭ позволит улучшить отбор абитуриентов, наиболее подходящих для конкретной образовательной программы.

## Список литературы

Айвазян С. А., Мхитарян В. С. (1998). *Прикладная статистика и основы эконометрики*. М.: Юнити.

Магнус Я. Р., Катышев П. К., Пересецкий А. А. (2007). *Эконометрика. Начальный курс*. М.: Дело.

Burton N. W., Ramist L. (2001). Predicting success in college: SAT studies of classes graduating since 1980. *College Board Research Report*, 2001–2. New York, The College Board.

Fishman J. A., Pasanella A. K. (1960). College admission selection studies. *Review of Educational Research*, 30 (4), 298–310.

Kobrin J. L., Patterson B. F., Shaw E. J., Mattern K. D., Barbuti S. M. (2008). Validity of the SAT for predicting first-year college grade point average. *College Board Research Report*, 2008–5. New York, The College Board.

Rothstein J. (2004). College performance predictions and the SAT. *Journal of Econometrics*, 121, 297–317.

Rousseeuw P. J., Leroy A. M. (1987). *Robust regression and outlier detection*. New York, Wiley.

Sackett P. R., Kuncel N. R., Arneson J. J., Cooper S. R., Waters S. D. (2009). Socioeconomic status and the relationship between the SAT and Freshman GPA: An analysis of data from 41 colleges and universities. *College Board Research Report*, 2009–1. New York, The College Board.