



Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Московский государственный технический университет  
имени Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)

---

ФАКУЛЬТЕТ \_\_\_\_\_ «Информатика и системы управления»

КАФЕДРА \_\_\_\_\_ «Теоретическая информатика и компьютерные технологии»

**Домашняя работа №2**  
**по курсу «Моделирование»**  
**«Проверка статистических гипотез»**

Студент группы ИУ9-81Б Окутин Денис Алексеевич

Преподаватель Домрачева Анна Борисовна

*Москва 2025*

## Цель работы

Проверить гипотезу о возможности описания стохастической зависимости между переменными двупараметрической степенной функцией связи. Для оценки параметров использовать модифицированный метод моментов, для проверки гипотез - критерий Колмогорова-Смирнова..

## Постановка задачи

Для исследования использован датасет, содержащий медицинские показатели пациентов, включая уровень глюкозы (Glucose) и инсулина (Insulin) в крови. Датасет включает 768 наблюдений. Glucose — уровень глюкозы в плазме крови (мг/дл) через 2 часа после приёма глюкозы. Insulin — уровень инсулина в сыворотке крови (мкЕд/мл) через 2 часа после теста. Нулевые значения в данных были скорректированы для возможности логарифмирования. Для реализации работы нужно выполнить следующие шаги:

- Загрузка и предварительная обработка данных;
- Анализ данных и удаление выбросов;
- Оценка параметров и с помощью модифицированного метода моментов;
- Проверка гипотезы о соответствии модели данным с помощью критерия Колмогорова-Смирнова;
- Визуализация результатов.

## Теоретические сведения

**Стохастическая зависимость** характеризуется вероятностным влиянием одной переменной на распределение другой без однозначной детерминации. В отличие от функциональной связи, здесь присутствует случайная компонента. Такие зависимости анализируются методами математической статистики и эконометрики.

Для моделирования нелинейных взаимосвязей используется **двухпараметрическая функция**:

$$y = \alpha \cdot x^\beta$$

Линеаризация модели через логарифмирование:

$$\ln y = \ln \alpha + \beta \ln x$$

позволяет применять методы линейной регрессии. Параметр  $\beta$  интерпретируется как эластичность, а  $\alpha$  — как масштабный коэффициент.

#### **Обоснование применимости степенной зависимости:**

Рассмотрим, что экспоненциальное распределение и S-образные кривые могут быть аппроксимированы функцией Вейбулла:

$$F(y) = 1 - \exp(-ay)$$

Предположим, что стохастическая зависимость между двумя случайными величинами  $y$  и  $x$  описывается функциональной зависимостью вида:

$$y = y(x) = bx^c$$

Тогда функция распределения примет следующий вид:

$$F(x) = F(y(x)) = 1 - \exp(-abx^c)$$

Обозначив  $d = ab$ , получим:

$$F(x) = 1 - \exp(-dx^c)$$

Таким образом, получается классическая форма функции распределения Вейбулла, где параметры  $d$  и  $c$  могут быть оценены по данным. Это демонстрирует, что степенная функция  $y = bx^c$  естественным образом приводит к Вейбулловскому распределению, что обосновывает её применимость для описания стохастической зависимости между переменными.

**Модифицированный метод моментов** — это статистический метод оценки параметров вероятностных распределений. В контексте данной работы

этот метод адаптирован для оценки параметров степенной функции с использованием логарифмического преобразования. Метод основан на сопоставлении теоретических моментов распределения с их эмпирическими аналогами. Модифицированный метод моментов для оценки параметров включает:

$$\beta = \frac{\mathbb{D}[\ln y]}{\mathbb{D}[\ln x]} \quad (1)$$

$$\alpha = \exp(\mathbb{E}[\ln y] - \beta \cdot \mathbb{E}[\ln x]) \quad (2)$$

где  $\mathbb{D}[\cdot]$  — выборочная дисперсия,  $\mathbb{E}[\cdot]$  — выборочное среднее.

**Критерий Колмогорова-Смирнова** — это непараметрический статистический тест, который используется для проверки гипотезы о соответствии эмпирического распределения теоретическому. Критерий применяется для определения, насколько хорошо предсказанные значения зависимой переменной соответствуют её фактическому распределению. Статистика критерия:

$$D_n = \sup_x |F_n(x) - F(x)|$$

где  $F_n(x)$  — эмпирическая функция распределения,  $F(x)$  — теоретическая функция распределения. Модель считается адекватной при  $D_n < D_{\text{крит}}$  для заданного уровня значимости.

## Практическая реализация

В фрагменте кода в листинге 1 подготавливаем данные для обучения модели.

Нулевые значения в колонках Glucose и Insulin могут быть артефактами измерения (например, неучтенные пределы чувствительности приборов). Их замена на  $x+1$  позволяет сохранить наблюдения и избежать математических ошибок при логарифмировании. Сдвиг на  $+1$  выполнен для корректного применения логарифмического преобразования  $\ln(x)$ , так как логарифм нуля не определён. Это стандартный приём при работе с данными, содержащими нулевые или отрицательные значения.

Разделение данных на обучающую (80%) и тестовую (20%) выборки проведено для проверки обобщающей способности модели. Это позволяет оценить, насколько модель адекватна на новых данных, не участвовавших в обучении.

### Листинг 1: Подготовка наборов данных

```
1 df = pd.read_csv('diabetes.csv')
2 df['Glucose'] = df['Glucose'] + 1
3 type1_data = df['Glucose'].values
4
5 df['Insulin'] = df['Insulin'] + 1
6 type2_data = df['Insulin'].values
7
8 type1_train, type1_test = train_test_split(type1_data, test_size=0.2,
      random_state=42)
9 type2_train, type2_test = train_test_split(type2_data, test_size=0.2,
      random_state=42)
10
11 def distribution_function(sample):
12     min_value = sample.min()
13     max_value = sample.max()
14     probs = []
15
16     values = np.arange(min_value, max_value, 0.2)
17     for value in values:
18         prob = np.sum(sample < value)
19         probs.append(prob / len(sample))
20
21     return values, probs
```

В фрагменте кода в листинге 2 получаем значения коэффициентов для двухпараметрической степенной функцией связи с помощью модифицированного метода моментов.

### Листинг 2: Модифицированный метод моментов

```
1 def moment_method(sample_1, sample_2):
2     log_x = np.log(sample_1)
3     log_y = np.log(sample_2)
4
5     A = np.vstack([log_x, np.ones(len(log_x))]).T
6     b_lin, ln_a = np.linalg.lstsq(A, log_y, rcond=None)[0]
7     a = np.exp(ln_a)
8     return a, b_lin
9
10
11 alpha, beta = moment_method(type1_train, type2_train)
```

```

12
13 to_plot_s1 = [x - 1 for x in type1_test]
14 to_plot_s2 = [x - 1 for x in type2_test]
15
16
17 to_plot_s2.sort()
18 to_plot_s1.sort()
19 approxed = alpha * to_plot_s1 ** beta

```

В фрагменте кода в листинге 3 проверяем гипотезу о том, что между переменными возможно описать стохастическую зависимость с помощью двухпараметрической степенной функции связи.

### Листинг 3: Критерий Колмогорова-Смирнова

```

1 def calculate_Kolmogorov_statistic(sample_1, sample_2, alpha, beta):
2     y_sample_1 = list(map(lambda point: calculate_distribution_func(
3         sample_1, point), sample_1))
4     y_sample_2 = list(map(lambda point: calculate_distribution_func(
5         alpha * sample_2 ** beta, point), alpha * sample_2 ** beta))
6     max_dif = -1e10
7     for y_1 in y_sample_1:
8         for y_2 in y_sample_2:
9             if np.abs(y_1 - y_2) > max_dif:
10                 max_dif = np.abs(y_1 - y_2)
11     return np.max(np.abs(y_approx_interp - y_type1_interp))
12
13 kolmogorov_stat = calculate_Kolmogorov_statistic(type1_test, type2_test,
14     alpha, beta)
15 alpha_ = 0.05
16 K_alpha = np.sqrt(-1/2 * np.log((1 - alpha_) / 2))
17
18 print(f"                                     : {kolmogorov_stat}")
19 print(f"K_alpha: {K_alpha}")

```

## Результаты работы

В ходе исследования проведен анализ стохастической зависимости уровня инсулина от концентрации глюкозы в крови с использованием двухпараметрической степенной функции связи были получены следующие результаты:

На рисунке 1 и рисунке 2 представлены отдельные эмпирические функции распределения для уровня глюкозы и инсулина соответственно. Визуальный ана-

лиз подтверждает несоответствие их форм, а также они имеют разные единицы измерения, что исключает возможность совмещения на одном графике.

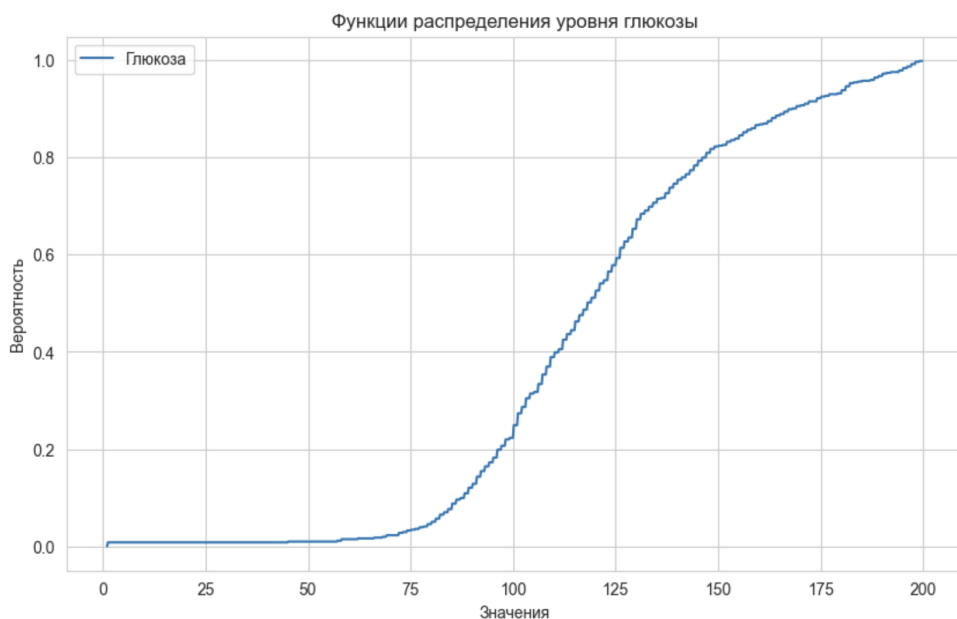


Рис. 1 — Функции распределения уровня глюкозы

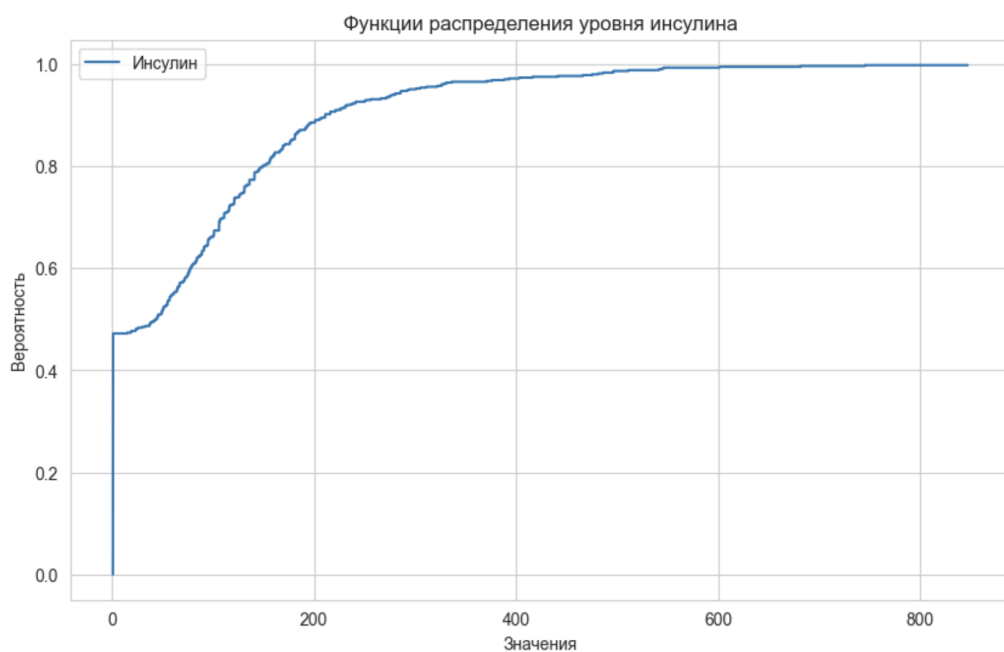


Рис. 2 — Функции распределения уровня инсулина

С помощью метода моментов были получены следующие значения для  $\alpha = 0.00008$  и  $\beta = 2.52$ :

$$y = 0.00008 \cdot x^{2.52}$$

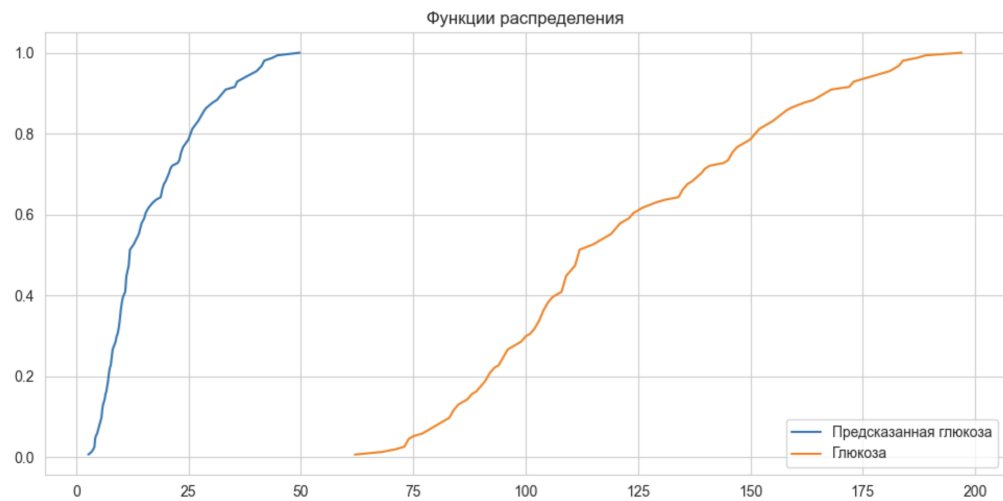


Рис. 3 — Функции распределения глюкозы



Рис. 4 — Абсолютная разница между распределениями глюкозы

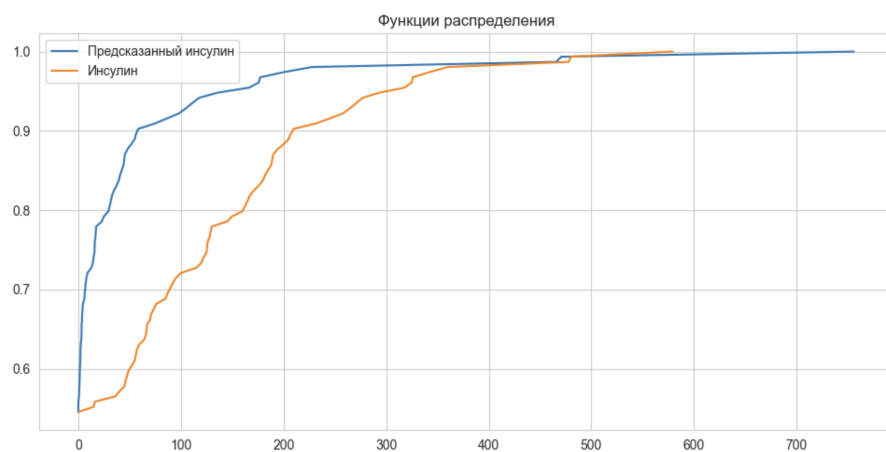


Рис. 5 — Функции распределения инсулина

Проверка соответствия модели данным на тестовой выборке с использованием критерия Колмогорова-Смирнова показала:





Рис. 6 — Абсолютная разница между распределениями инсулина

- Максимальное расхождение между эмпирическими распределениями:  
 $D_n = 0.93$
- Критическое значение при  $\alpha = 0.05$ :  $K_\alpha = 0.61$

Так как  $D_n > K_\alpha$ , гипотеза о согласии распределений *отвергается* на заданном уровне значимости.

Визуализация результатов подтвердила плохое соответствие модели данным на тестовой выборках.

## Заключение

В данной работе применена двупараметрическая степенная функция для попытки описания стохастической зависимости между концентрацией инсулина и уровнем глюкозы. Как на тестовой, так и на обучающей выборках, гипотеза о возможности описания стохастической зависимости с помощью степенной функции с двумя параметрами отвергается на уровне значимости 0.05 с помощью критерия Колмогорова-Смирнова.