

Количественная оценка погрешности  
определения газонасыщенности по  
подтверждениям качественного характера

*(Численное исследование на основе теории вероятностей и решение задачи  
количественной оценки погрешности определения газонасыщенности  $K_g$   
пластов методом «ММНК- $K_g$ » по качественным результатам испытаний)*

Поляченко Юрий Анатольевич

1 августа 2020 г.

## Содержание

<b>1</b>	<b>Постановка задачи</b>	<b>3</b>
1.1	Геофизическая постановка задачи . . . . .	3
1.2	Математическая формулировка . . . . .	4
<b>2</b>	<b>Предлагаемое решение</b>	<b>6</b>
2.1	Идея и приближения . . . . .	6
2.2	Аналитика . . . . .	9
2.3	Продолжение примера аналитикой . . . . .	10
<b>3</b>	<b>Результат применения</b>	<b>11</b>
3.1	Типичные значения . . . . .	11
3.2	Теоретический анализ . . . . .	15
<b>4</b>	<b>Описание и рекомендации к программе</b>	<b>16</b>
4.1	Исследуемый интервал, ввод имеющихся данных . . . . .	17
4.2	Основной рисунок . . . . .	19
	<b>Список литературы</b>	<b>21</b>

# 1 Постановка задачи

## 1.1 Геофизическая постановка задачи

Газовая скважина исследована аппаратурой мультиметодного многозондового нейтронного каротажа (ММНК) и по методике ММНК-Кг определены  $N$  значений коэффициента газонасыщенности  $K_g$  в разных пластах разреза. После этого с целью проверки корректности всей технологии проводятся испытания скважины на приток, в которых на качественном уровне измеряется состав фактически добываемой продукции газовой-водной смеси. Это означает, что продукция классифицируется на небольшое число градаций  $M \sim 3-5$  разбиений шкалы  $K_g$  на эквидистантные широкие интервалы протяженностью по  $\Delta K_g = [K_g] / M$ , где  $[K_g]$  – максимально возможный диапазон изменения  $K_g$  в исследуемых геолого-промысловых условиях. Например, наиболее часто используемыми интервалами изменения  $\Delta K_g$  для типовых  $[K_g] = [0,1]$ ,  $M = 4$  являются по терминологии газовиков: «вода (0–0.25), вода+газ (0.25–0.5), газ+вода (0.5–0.75), газ (0.75–1)». Затем проверяются доли  $p_0$ , % правильных попаданий предсказанных ММНК-Кг численных значений  $K_g$  в соответствующие им широкие интервалы  $\Delta K_g$ . Если большинство этих долей  $p_0 > 80-90\%$ , то технология признается корректной, т.к. получила качественное подтверждение по результатам испытаний, считающихся одним из наиболее прямых и убедительных способов тестирования методик в скважинной геофизике. Представляется очевидным, что при достаточно большом числе  $N \gg M$  определений  $K_g$  и, разумеется, при условии выполнения достаточно жесткого критерия подтверждаемости по  $p_0$  любая разумно введенная оценка фактической средней погрешности определения  $K_g$  должна дать величину, существенно меньшую широких интервалов разбиения  $\Delta K_g$ . Другими словами, это означает, что технологию ММНК можно будет переквалифицировать из качественной по способу ее подтверждения в количественную по фактически достигаемому уровню погрешности определения  $K_g$ .

Поэтому задачами работы явились следующие:

- численное обоснование этого утверждения на основе теории вероятностей с разработкой алгоритма и программы расчета фактической средней погрешности определения  $K_g$  по всем имеющимся данным определений и подтверждений в исследуемой скважине;
- численное изучение поведения погрешности в зависимости от варьируемых параметров  $M$ ,  $N$ ,  $[K_g]$ ,  $p_0$  и характера вероятностного

распределения найденных значений  $K_2$  на  $[K_2]$ ,  $P(K_2)$  - от равномерного до гауссового с большой дисперсией. Диапазоны изменения варьируемых параметров:

- $M = 3, 4, 5$
- $M \leq N \leq 30$
- $[K_2] = [0.5, 1], [0.25, 1], [0, 1]$
- $p_0 = 80\%, 90\%, \approx 100\%$

- выдача практических рекомендаций по выбору единственного управляемого параметра  $N$  в зависимости от априори задаваемых геофизиками и газовиками параметров  $M$  и  $[K_2]$ , а также от фактически получившихся характеристик – распределения  $P(K_2)$  и значений  $p_0$  в результате сопоставления определений и подтверждений  $K_2$ .

## 1.2 Математическая формулировка

Цель – предсказать погрешность  $\varepsilon_0$  выдаваемых нашей программой значений  $x$  искомого параметра  $y$ .

Область изменения  $x \in [K_2]$  разбита на интервалы  $[a_j; b_j]$ ,  $j \in \{\overline{1, M}\}$ , для каждого из которых есть  $N_j$  экспериментов. Считается, что искомая погрешность  $\varepsilon_0$  может меняться от интервала к интервалу, но постоянная внутри интервала (т.е. точнее писать  $\varepsilon_{0j}$ ). Фиксируем  $j$  и работаем в выбранном интервале, поэтому далее индекс интервала  $j$  опущен.

Есть  $N$  экспериментов, про которые известно, что в каждом из них истинное значение  $y_i$  попало в интервал, т.е.  $y_i \in [a; b] \forall i \in \{\overline{1, N}\}$ . На каждый из этих экспериментов у нас есть результат работы нашей программы  $x_i$ . Предполагается, что истинное значение  $y_i$  распределено по Гауссу со средним  $x_i$  и некой дисперсией  $\varepsilon$ , т.е.

$$\frac{\mathcal{P}(y_i \in [t, t + dt])}{dt} = g(t, x_i, \varepsilon), \quad (1)$$

где  $\mathcal{P}(A)$  - вероятность того, что  $A$  верно, а

$$g(x, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right] \quad (2)$$

– Гауссово распределение.

Ищем зависимость

$$\varepsilon_0(a, b, \{x_i\}_{i=1}^N, p_0) \quad (3)$$

такую, что вероятность реализации описанной выше ситуации (т.е. что все истинные значения попали в интервал)  $= p_0$ , т.е.

$$\mathcal{P}(\forall x \in [a; b] \ |x - y| < \varepsilon_0) = p_0 \quad (4)$$

Из сторонних соображений считается известным минимально возможная погрешность  $\varepsilon_{min}$ , т.е. если метод выдает  $\varepsilon_0 < \varepsilon_{min}$ , то считаем  $\varepsilon_0 = \varepsilon_{min}$ .

## 2 Предлагаемое решение

### 2.1 Идея и приближения

Задав  $\varepsilon$ , можно посчитать вероятность реализации ситуации, описанной в постановке – попадания всех истинных значений параметра  $y_i$ , распределенных по Гауссу каждый около своего  $x_i$ , в интервал  $[a; b]$ . Далее предположение – эта вероятность равна нашей целевой вероятности  $p_0$ . Не очевидно, почему это должно выполняться точно (скорее всего это не выполняется), но для оценки предложено использовать такую модель.

Поясним разумность данного выбора. Будем брать пробные  $\varepsilon$  и смотреть как от этого зависит ожидаемое поведение истинных значений  $y_i$  относительно наших точек  $x_i$ . Для примера возьмем весь интервал  $[0.2; 0.5]$  и предположим что у нас имеются 5 точек, для которых наша программа выдала ответы 0.22, 0.3, 0.35, 0.4, 0.43. Если предположить, что погрешность наших предсказаний  $\varepsilon = 0.02$ , то плотность вероятности для каждого из 5 истинных значений будет выглядеть так:

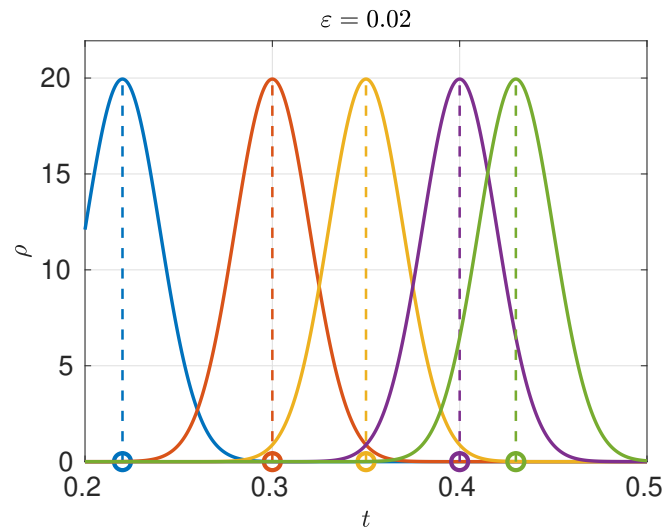


Рис. 1: Разным цветам отвечают разные эксперименты. Для каждого эксперимента: проколотый круг на оси X – наше предсказание ответа, купол – распределение плотности вероятности того, что истинное значение ответа примет значение  $t$ , в зависимости от  $t$ .

На Рис. 1 видно, что для всех точек кроме синей почти вся кривая (вероятность = площадь под кривой) находится в исследуемом интервале. Это значит, что при  $\varepsilon = 0.02$  для всех точек кроме синей вероятность

того, что истинное значение параметра попадет в интервал, равна почти 100%. Синяя же точка находится на расстоянии  $\sim 1\sigma$  (в данном случае 0.02), что значит, что вероятность того, что истинное значение параметра в синем эксперименте попадет в интервал  $[0.2; 0.5]$  будет  $\approx 16\%$ . Попадания истинных значений в интервал – события независимые, поэтому вероятность реализации картины в целом будет произведением вероятностей попадания каждого значения в интервал по отдельности. В нашем случае все вероятности кроме синей  $\approx 1$ , поэтому общая вероятность  $\approx$  синяя вероятность  $\approx 84\%$ . Это значит, что если бы погрешность нашей программы была 0.02, то вероятность случатся тому что случилось на рассматриваемых 5 экспериментах в совокупности была бы 84 %. Поняв это, можно решить обратную задачу: сказать, что мы верим эксперименту на скажем 95%, и найти такое  $\varepsilon$ , при котором вероятность его реализации как раз будет 95%. Понятно, что задача решается – если в предыдущем примере мы возьмем  $\varepsilon = 0.008$ , то распределение вероятностей для истинных значений будет

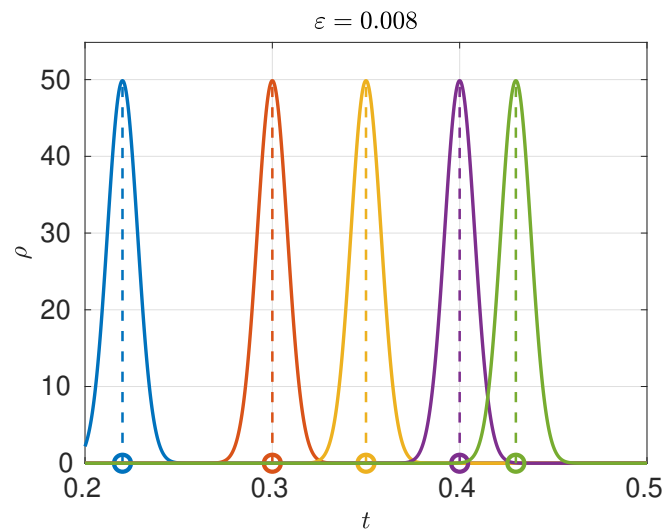


Рис. 2: Обозначения аналогичны Рис. 1. Вероятность реализации эксперимента  $>99\%$ , что эквивалентно практически полному доверию эксперименту.

Если же все наши экспериментальные точки лежат ближе к центру исследуемой области, то оценка на погрешность выходит грубой. Это понятно из следующего примера. Сдвинем точку 0.22 из предыдущего примера в точку 0.28. График для  $\varepsilon = 0.04$  будет выглядеть как

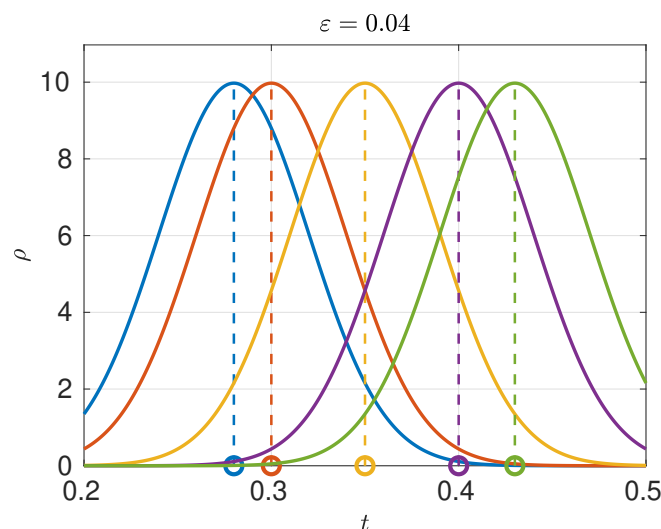


Рис. 3: Обозначения аналогичны Рис. 1. Вероятность реализации эксперимента 92.5%.

Т.е. при наличии точки 0.22, близкой к левой границе исследуемого интервала 0.2, вероятность реализации эксперимента уже при  $\varepsilon = 0.02$  была 84%, что говорило о том, что в реальности скорее всего погрешность была меньше. Здесь же даже при  $\varepsilon = 0.04$  вероятность все еще  $>90\%$ . Это на самом деле логично, т.к. если все наши точки у центра интервала, то единственный известный нам факт (на котором и строится вся оценка опгрешности) о том, что все истинные значения попали в интервал, позволяет отбросить только самые большие значения погрешностей.

Заинтересовавшийся читатель может найти весь инходный код проекта и мои контакты для вопросов и предложений [здесь](#). В частности по ссылке лежит программа для оценки ошибок по описываемому здесь методу.

Теперь приведем аналитическое выражение описанной идеи:



## 2.2 Аналитика

Вероятность попадания  $i$ -ой истинной точки в интервал

$$p_i(\varepsilon) = \int_a^b g(x, x_i, \varepsilon) dx = \int_{(a-x_i)/\varepsilon}^{(b-x_i)/\varepsilon} g(x, 0, 1) dx. \quad (5)$$

Введем функцию (известную как функция ошибок):

$$\text{erf}(x) = \int_{-\infty}^x g(t, 0, 1) dt. \quad (6)$$

Попадание каждого истинного значения в интервал - независимое событие, поэтому вероятность реализации нашей совокупности экспериментов

$$p(\varepsilon) = \prod_{i=1}^N p_i(\varepsilon) = \prod_{i=1}^N \left[ \text{erf}\left(\frac{b-x_i}{\varepsilon}\right) - \text{erf}\left(\frac{a-x_i}{\varepsilon}\right) \right] \quad (7)$$

Для нахождения желаемого  $\varepsilon_0$  решаем уравнение на  $\varepsilon_0$  при заданном  $p_0$ .

$$p(\varepsilon_0) = p_0 \quad (8)$$

Уравнения явно не решаются аналитически. Но несложно показать, что функция  $p(\varepsilon)$  монотонна, а интервал изменения  $\varepsilon$  известен и невелик, откуда следует, что уравнение легко решается численно даже самыми простейшими методами вроде деления отрезка пополам. В примерах использован алгоритм, реализованный в функции *fzero* в Matlab и описанный в [1]

## 2.3 Продолжение примера аналитикой

Продолжим использовать 5 точек из раздела 2.1. В разделе разделе 2.1 был описан алгоритм как мы задавшись определенным  $\varepsilon$  можем оценить вероятность  $p$ , с которой при этом  $\varepsilon$  реализовались бы имеющиеся у нас экспериментальные данные. Сделав так для многих различных  $\varepsilon$ , можно для каждого из них получить свое значение  $p(\varepsilon)$  (Синяя кривая на Рис. 4).

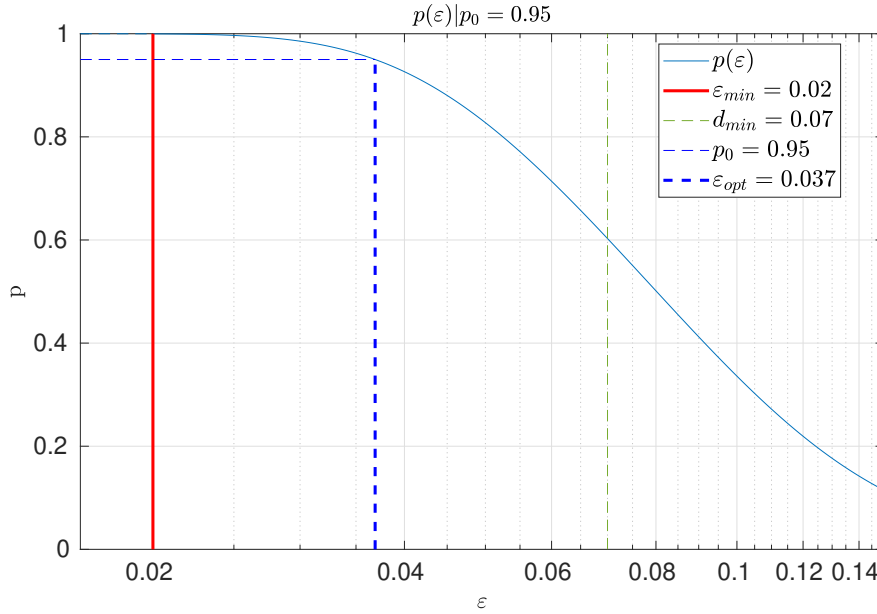


Рис. 4: Зависимость  $p(\varepsilon)$  вероятности реализации эксперимента при данной погрешности программы. Красная линия – принятая минимально возможная погрешность  $\varepsilon_{min} = 0.02$ , Синяя вертикаль – найденная оценка, синяя горизонталь – наш выбор  $p_0 = 95\%$ , зеленый пунктир - минимальное расстояние точек до границы.

Видно, что наличие множества ( $> 1$ ) точек позволяет улучшить оценку с очевидного значения минимального расстояния до границы – синяя линия левее зеленой, т.е. оценка по предлагаемому методу лучше чем наивная оценка сверху на глаз.

## 3 Результат применения

### 3.1 Типичные значения

Можно исследовать, как оценка погрешности зависит от количества имеющихся экспериментальных данных в «усредненном» случае, когда ответы нашей программы расположены в интервале на равных промежутках:

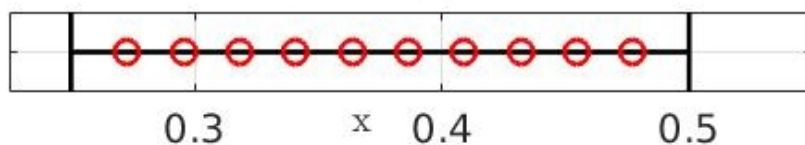


Рис. 5: Равномерное расположение 10 пробных точек в интервале  $[0.25; 0.5]$ .

Сгенерировав таким равномерным образом несколько наборов «предсказаний» нашей программы, можно получить зависимость получаемой оценки погрешности от количества имеющихся экспериментальных точек:

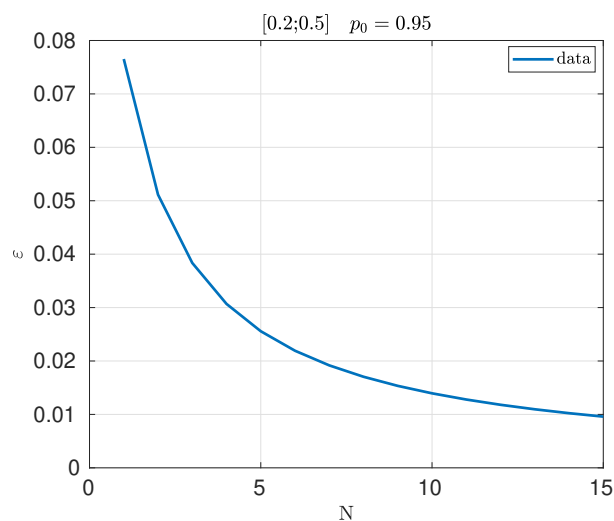


Рис. 6: Зависимость  $\varepsilon(N)$  оценки погрешности программы от количества экспериментальных точек при их равномерном распределении в интервале.

Можно делать так для разных наборов параметров, описанных в геофизической постановке задачи. Приведем некоторые полученные таким образом зависимости для типичных значений параметров:

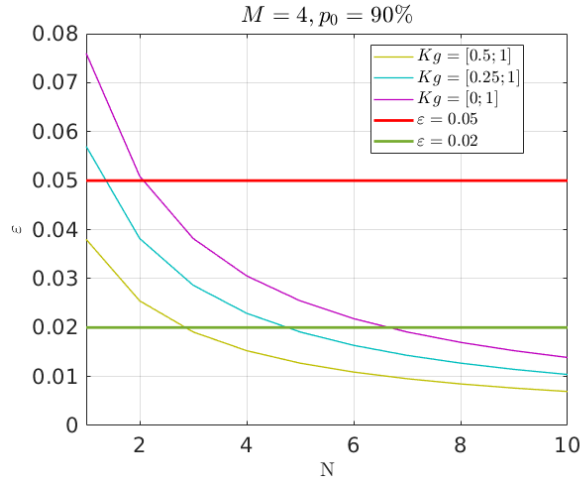


Рис. 7: Зависимость  $\varepsilon(N)$  оценки погрешности программы от количества экспериментальных точек. Различные интервалы допустимых  $K_2$ .

$N(K_2, \varepsilon)$		
$K_2 \backslash \varepsilon$	0.05	0.02
[0; 1.0]	3	7
[0.25; 1.0]	2	5
[0.5; 1.0]	1	3

Таблица 1: Минимальные значения  $N$ , необходимые для достижения данного  $\varepsilon$  при данном  $K_2$ .

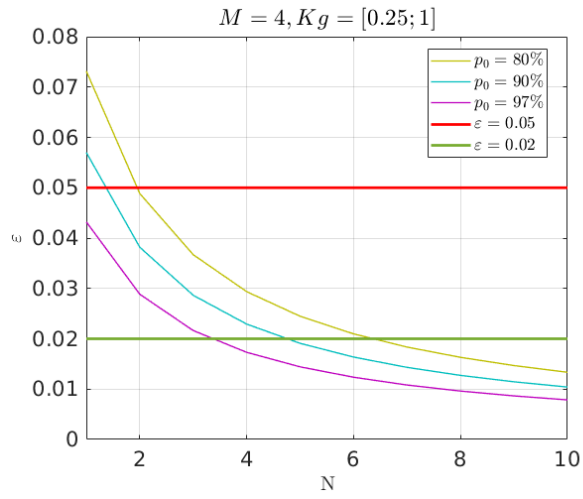
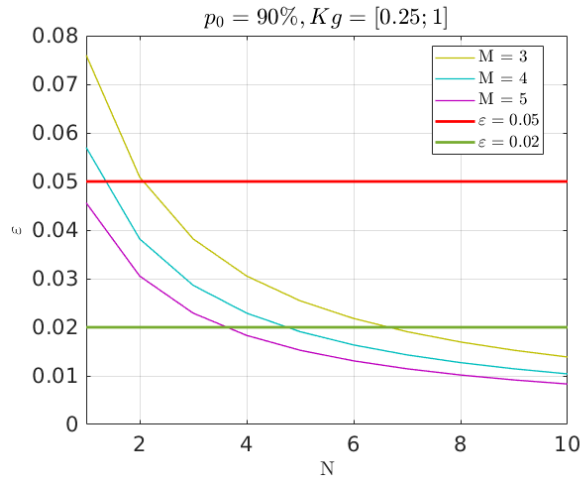


Рис. 8: Зависимость  $\varepsilon(N)$  оценки погрешности программы от количества экспериментальных точек. Различные степени доверия эксперименту.

$N(p_0, \varepsilon)$

$\varepsilon$	0.05	0.02
$p_0$		
80%	2	7
90%	2	5
97%	1	4

Таблица 2: Минимальные значения  $N$ , необходимые для достижения данного  $\varepsilon$  при данном  $p_0$ .



$N(M, \varepsilon)$		
$M \backslash \varepsilon$	0.05	0.02
3	3	7
4	2	5
5	1	4

Рис. 9: Зависимость  $\varepsilon(N)$  оценки погрешности программы от количества экспериментальных точек. Различные разбиения типичного интервала.

Таблица 3: Минимальные значения  $N$ , необходимые для достижения данного  $\varepsilon$  при данном  $M$ .

Может быть также интересен наиболее сложный случай:

$M = 3, p_0 = 80\%, Kg = [0; 1.0]$

При таких параметрах для достижения  $\varepsilon = 0.05$  нужно  $N = 5$ , а для  $\varepsilon = 0.02$  нужно  $N = 12$ . Т.е. при наличии 12 точек в интервале можно с уверенностью говорить о подтверждении хорошей точности метода в данной области  $Kg$ .

### 3.2 Теоретический анализ

На глаз зависимость на Рис. 6 близка к  $1/N$ , что ожидаемо, т.к. погрешность в основном определяется минимальным расстоянием до границы, которое при выбранной расстановке точек убывает как  $1/N$ .

Можно проверить отклонения от закона  $1/N$  – рис.(10).

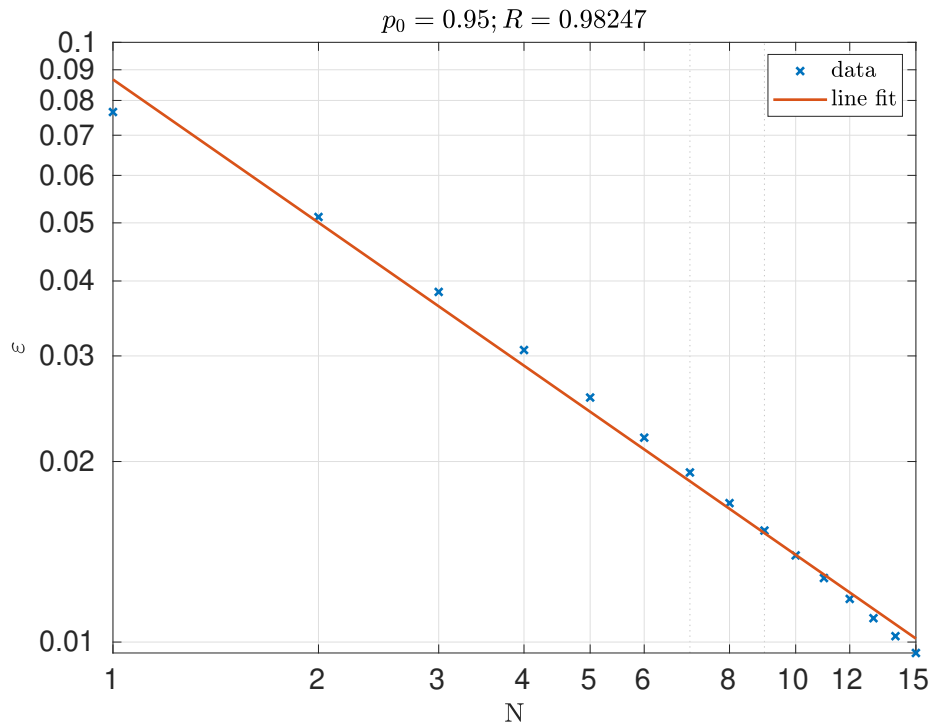


Рис. 10:  $\varepsilon(N)$ , логарифмический масштаб, попытка линеаризации

Видно, что наклон с правда близок к  $-1$ , но небольшие отклонения от линейности есть.

## 4 Описание и рекомендации к программе

Для удобства применения изложенного метода создана программа, позволяющее получить оценку погрешности для любого заданного набора экспериментов.

Окно программы выглядит так:

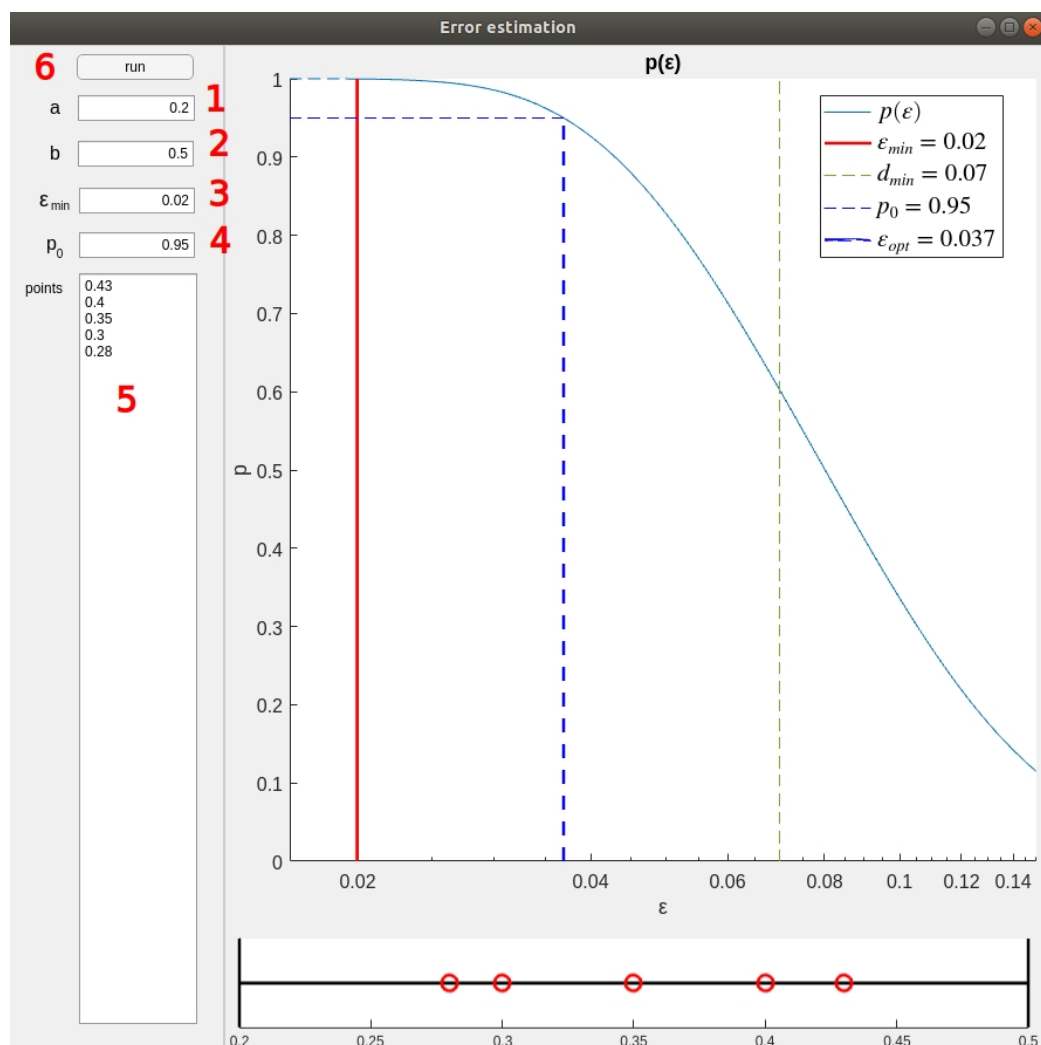


Рис. 11: Окно программы. 1,2 – границы анализируемого интервала. 3 – априорная минимально допустимая погрешность. 4 – степень доверия эксперименту. 5 – имеющиеся результаты работы программы. 6 – сделать расчет с введенными данными.



Поясним, что изображено на рисунках:

#### 4.1 Исследуемый интервал, ввод имеющихся данных

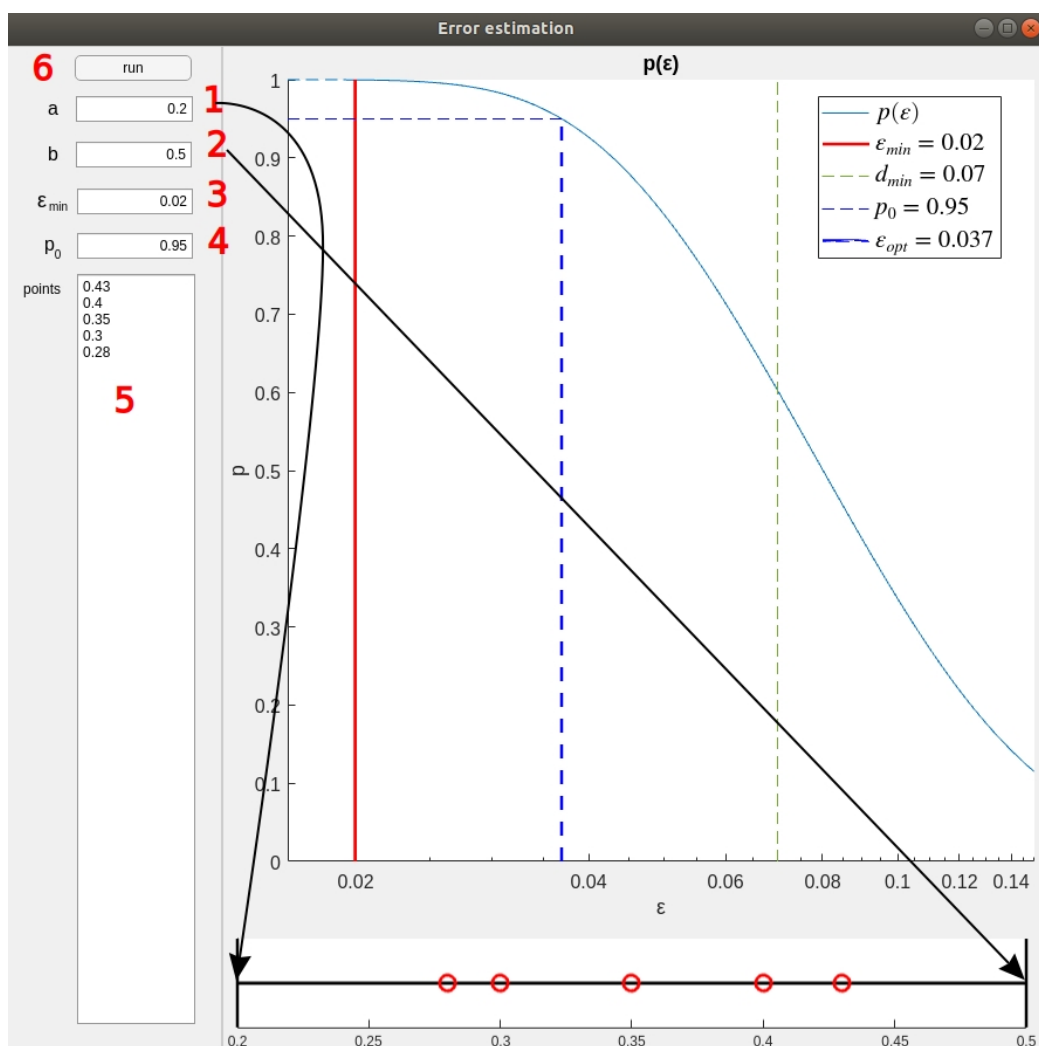


Рис. 12: Значения полей 1 и 2 – значения  $a$  и  $b$ , т.е. границ исследуемого интервала.

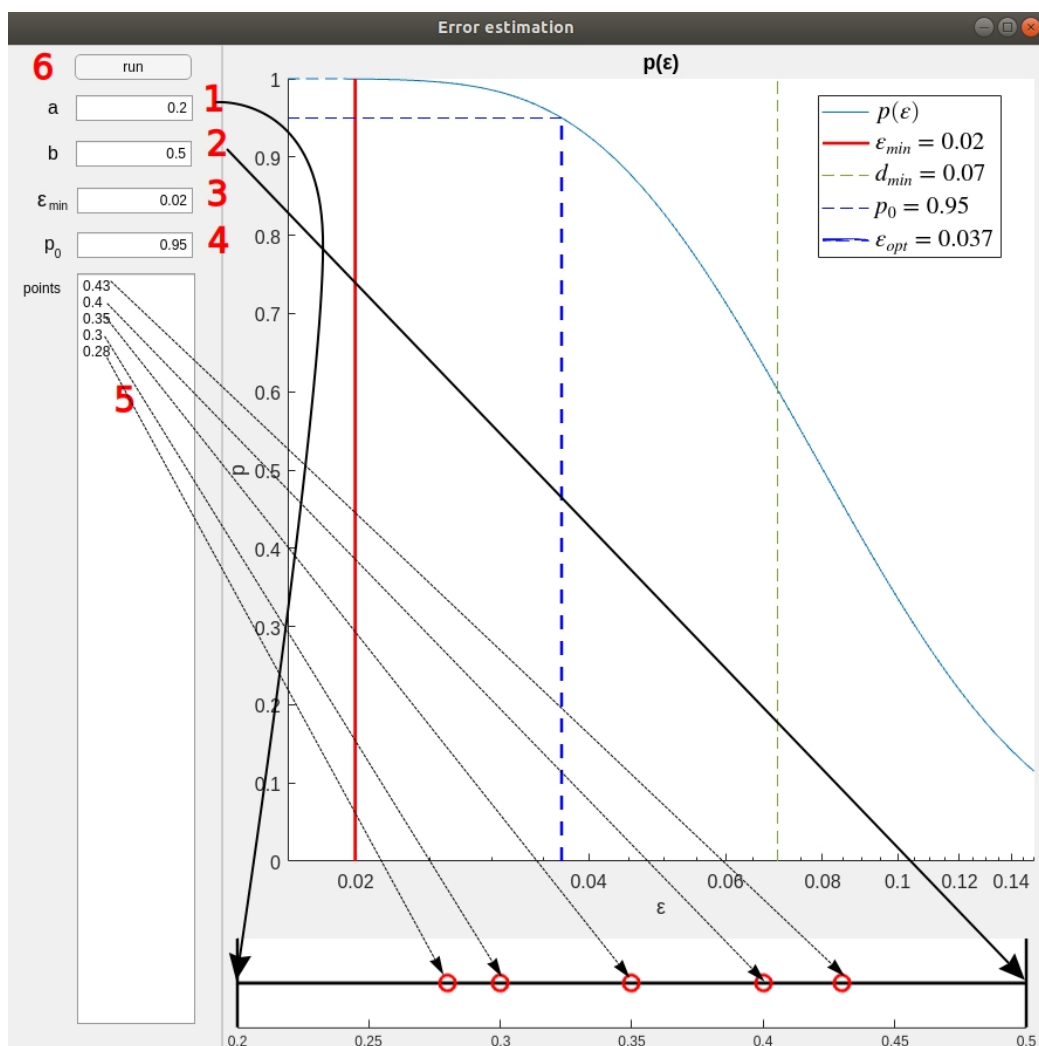


Рис. 13: Таблица 5 – имеющиеся ответы программы.

## 4.2 Основной рисунок

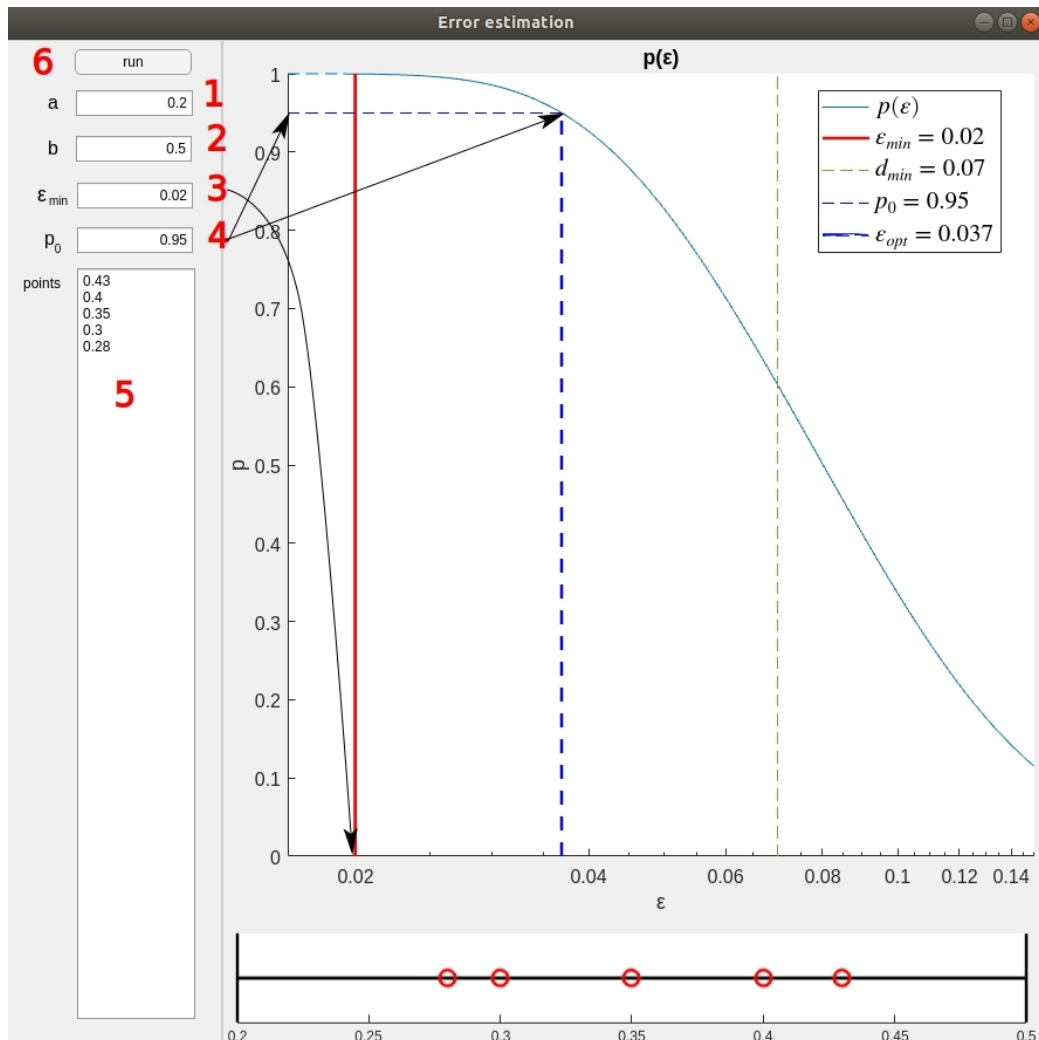


Рис. 14: Поле 3 – априорное значение минимально допустимой погрешности. Оно задается по соображениям пользователя исходя из сторонней информации. Поле 4 – степень доверия эксперименту. Если мы уверены в эксперименте на 95% – пишем 0.95. Основная задача программы состоит в использовании этой величины – по точкам 5 и интервалу 1,2 программа строит голубую кривую и находит тот  $X$  (т.е. тот  $\epsilon$ ), при котором кривая принимает значение в поле 4 (т.е. вероятность реализации эксперимента совпадает с нашими ожиданиями).

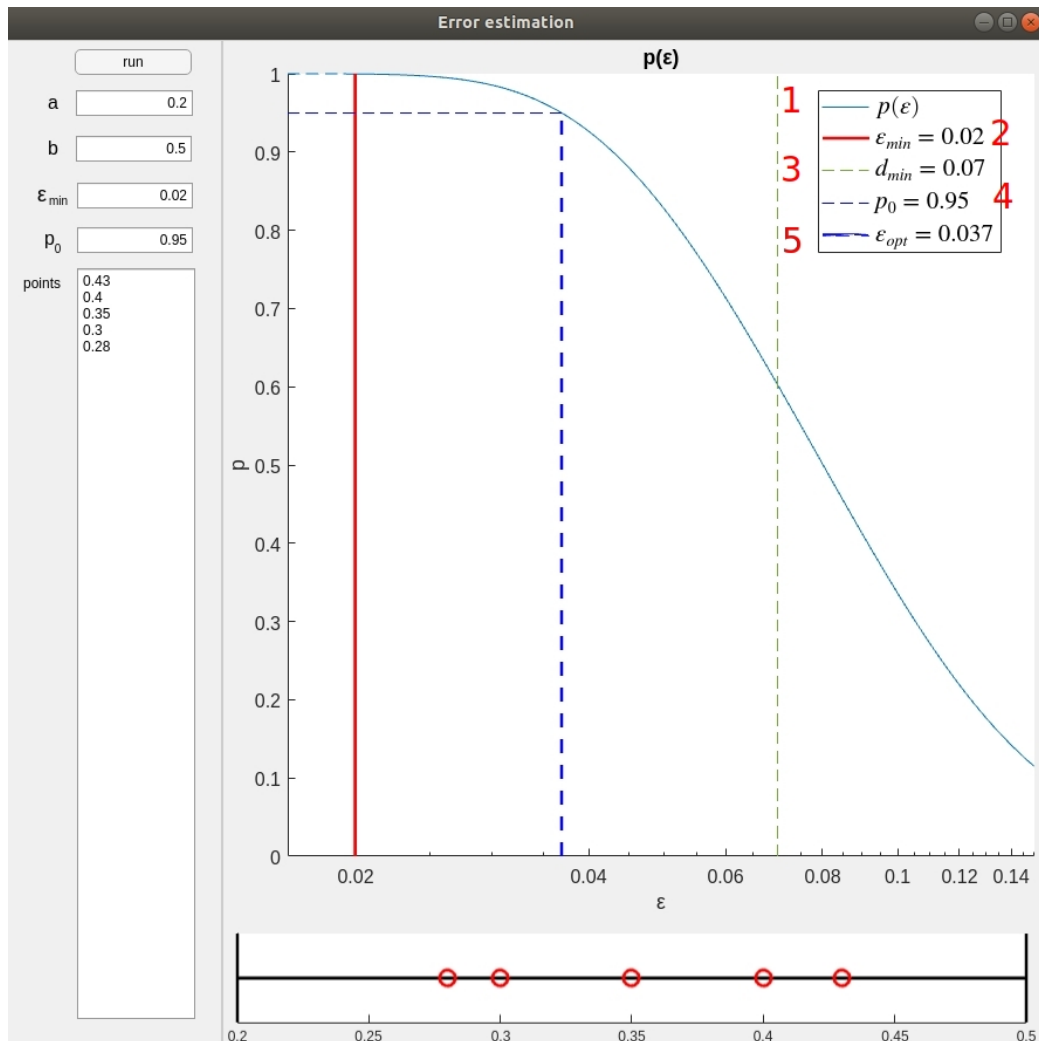


Рис. 15: Описание легенды. 1 – зависимость вероятности реализации имеющейся совокупности экспериментов от предполагаемой погрешности предсказаний. 2 – априорно заданная минимальная погрешность. 3 – минимальное расстояние от заданного множества точек до ближайшей границы интервала. 4 – степень доверия эксперименту. 5 – ответ программы, получаемая оценка погрешности предсказаний.

## Список литературы

- [1] F Grund. “Forsythe, GE/Malcolm, MA/Moler, CB, Computer Methods for Mathematical Computations. Englewood Cliffs, New Jersey 07632. Prentice Hall, Inc., 1977. XI, 259 S”. В: *ZaMM* 59.2 (1979), с. 141—142.