

Отчет о проверке на заимствования №1

Автор: Админ АПИ НГТУ antipl_api@corp.nstu.ru / ID: 9258
Проверяющий: Админ АПИ НГТУ (antipl_api@corp.nstu.ru / ID: 9258)
Организация: Новосибирский Государственный Технический Университет

Отчет предоставлен сервисом «Антиплагиат» - <http://nstu.antiplagiat.ru>

ИНФОРМАЦИЯ О ДОКУМЕНТЕ

№ документа: 2945
Начало загрузки: 13.06.2018 20:04:43
Длительность загрузки: 00:00:55
Имя исходного файла: Fedosov_PMM-61_432073a.pdf
Размер текста: 1128 кБ
Символов в тексте: 55368
Слов в тексте: 6091
Число предложений: 367
Способ извлечения текста: OCR

ИНФОРМАЦИЯ ОБ ОТЧЕТЕ

Последний готовый отчет (ред.)
Начало проверки: 13.06.2018 20:05:39
Длительность проверки: 00:00:08
Комментарии: не указано
Модули поиска: Сводная коллекция ЭБС, Цитирование, Модуль поиска Интернет, Модуль поиска "НГТУ", Модуль поиска перефразирований Интернет, Модуль поиска общепотребительных выражений, Кольцо вузов

ЗАИМСТВОВАНИЯ	ЦИТИРОВАНИЯ	ОРИГИНАЛЬНОСТЬ
15,07% <div><div></div></div>	0,35% <div><div></div></div>	84,58% <div><div></div></div>

Заимствования — доля всех найденных текстовых пересечений, за исключением тех, которые система отнесла к цитированиям, по отношению к общему объему документа.
Цитирования — доля текстовых пересечений, которые не являются авторскими, но система посчитала их использование корректным, по отношению к общему объему документа. Сюда относятся оформленные по ГОСТу цитаты; общепотребительные выражения; фрагменты текста, найденные в источниках из коллекций нормативно-правовой документации.
Текстовое пересечение — фрагмент текста проверяемого документа, совпадающий или почти совпадающий с фрагментом текста источника.
Источник — документ, проиндексированный в системе и содержащийся в модуле поиска, по которому проводится проверка.
Оригинальность — доля фрагментов текста проверяемого документа, не обнаруженных ни в одном источнике, по которым шла проверка, по отношению к общему объему документа.
Заимствования, цитирования и оригинальность являются отдельными показателями и в сумме дают 100%, что соответствует всему тексту проверяемого документа.
Обращаем Ваше внимание, что система находит текстовые пересечения проверяемого документа с проиндексированными в системе текстовыми источниками. При этом система является вспомогательным инструментом, определение корректности и правомерности заимствований или цитирований, а также авторства текстовых фрагментов проверяемого документа остается в компетенции проверяющего.

№	Доля в отчете	Доля в тексте	Источник	Ссылка	Актуален на	Модуль поиска	Блоков в отчете	Блоков в тексте
[01]	5,97%	5,97%	Критерий Смирнова — allRefs.net	http://allrefs.net	30 Янв 2017	Модуль поиска перефразирований Интернет	13	13
[02]	0,85%	2,96%	Критерий Смирнова — allRefs.net	http://allrefs.net	27 Фев 2016	Модуль поиска Интернет	9	22
[03]	0,92%	1,56%	скачать файл PDF	http://nstu.ru	27 Авг 2017	Модуль поиска Интернет	6	12
[04]	0,85%	1,55%	5 ПРИМЕНЕНИЕ КОМПЬЮТЕРНОГО М...	http://libed.ru	02 Сен 2017	Модуль поиска Интернет	7	12
[05]	0,45%	1,43%	3 ПРИМЕНЕНИЕ КОМПЬЮТЕРНОГО М...	http://libed.ru	19 Авг 2017	Модуль поиска Интернет	3	10
[06]	1,3%	1,42%	Расширение прикладных возможност...	http://disland.com	29 Янв 2017	Модуль поиска перефразирований Интернет	4	4
[07]	0,28%	1,28%	скачать файл PDF	http://nstu.ru	06 Сен 2017	Модуль поиска Интернет	3	10
[08]	0%	1,27%	Расширение прикладных возможност...	http://tekhnosfera.com	01 Янв 2017	Модуль поиска перефразирований Интернет	0	3
[09]	0,85%	1,26%	Лабораторная работа № 1. Проверка с...	http://pandia.ru	05 Янв 2017	Модуль поиска перефразирований Интернет	3	4
[10]	0,84%	1,11%	Применение компьютерного моделир.	http://dslib.net	30 Янв 2017	Модуль поиска перефразирований Интернет	2	3
[11]	0,21%	0,97%	скачать файл PDF	http://nstu.ru	25 Авг 2017	Модуль поиска Интернет	2	7
[12]	0%	0,92%	234537	http://biblioclub.ru	19 Апр 2016	Сводная коллекция ЭБС	0	9
[13]	0,2%	0,86%	Прикладная статистика	https://book.ru	03 Июл 2017	Сводная коллекция ЭБС	2	8
[14]	0,49%	0,61%	Моделирование в контроллинге - 2016..	http://bmstu.ru	13 Ноя 2017	Модуль поиска Интернет	5	7
[15]	0,44%	0,6%	Факультет экономики/41 ФФР 1 Шишо..	не указано	30 Мая 2012	Кольцо вузов	3	5
[16]	0%	0,51%	Реконструкция генных сетей на основ...	http://knowledge.allbest.ru	29 Янв 2017	Модуль поиска перефразирований Интернет	0	1
[17]	0,32%	0,51%	Functional approach to cluster municipa...	http://smta.net	26 Янв 2017	Модуль поиска Интернет	4	6
[18]	0%	0,49%	Вестник Томского государственного у...	http://ibooks.ru	09 Дек 2016	Сводная коллекция ЭБС	0	4
						Модуль поиска		

[19]	0%	0,37%	Проверка качественных характери...	http://allrefs.net	13 Янв 2017	перефразирований Интернет	0	1
[20]	0,19%	0,37%	Статистический анализ числовых вели..	http://reftrend.ru	04 Фев 2017	Модуль поиска Интернет	2	5
[21]	0,03%	0,35%	скачать файл PDF	http://nstu.ru	25 Авг 2017	Модуль поиска Интернет	1	2
[22]	0,28%	0,28%	volkova_e_s_model-hegselmann---kraus...	не указано	09 Мая 2018	Кольцо вузов	1	1
[23]	0,1%	0,26%	Исследование процессов формирован..	https://diss.unn.ru	15 Дек 2016	Модуль поиска Интернет	1	2
[24]	0%	0,21%	120224	http://biblioclub.ru	15 Апр 2016	Сводная коллекция ЭБС	0	1
[25]	0%	0,21%	221969	http://biblioclub.ru	19 Апр 2016	Сводная коллекция ЭБС	0	1
[26]	0,21%	0,21%	8718	http://e.lanbook.com	09 Мар 2016	Сводная коллекция ЭБС	1	1
[27]	0%	0,17%	ISBN9785922113755.txt	не указано	26 Окт 2017	Кольцо вузов	0	2
[28]	0,16%	0,16%	https://sibsutis.ru/workgroups/w/group...	https://sibsutis.ru	14 Ноя 2017	Модуль поиска Интернет	1	1
[29]	0%	0,12%	300.Шишкина Л.А.Математика учебно...	http://docme.ru	29 Июн 2017	Модуль поиска Интернет	0	2
[30]	0%	0,11%	Компьютерное моделирование и иссл..	http://journals.tsu.ru	раньше 2011	Модуль поиска Интернет	0	2
[31]	0,11%	0,11%	67302	http://biblioclub.ru	раньше 2011	Сводная коллекция ЭБС	1	1
[32]	0,02%	0,09%	РД 50-705-91: Методические указания. ...	http://standartgost.ru	01 Мая 2017	Модуль поиска Интернет	1	1
[33]	0,35%	0%	не указано	не указано	раньше 2011	Модуль поиска общеупотребительных выражений	6	8

Текст документа

Аннотация

Объем работы — 46 страниц, состоит из введения, 3 глав основного содержания, заключения, списка литературы И приложения с программным кодом, включая 8 рисунков, 28 таблиц И списка литературы из 16 источников 10 .

Объект исследования — критерии однородности Смирнова, Андерсона-Дарлинга, Лемана-Розенблатта, данные ограниченной точности, оценки мощности критериев.

Цель работы — исследование распределений статистик И мощности критериев однородности Смирнова, Лемана-Розенблатта, Андерсона-Дарлинга на данных ограниченной точности 6 .

В результате работы были проведены исследования распределения статистик. Показано влияние близости функции распределения статистики к функции предельного распределения от размерности выборок И количества уникальных значений в объединенной выборке. Исследованы оценки мощности критериев на данных ограниченной точности.

Практическая ценность работы заключается в полученных результатах исследования критериев однородности на данных ограниченной точности.

Сформулированы рекомендаЦИИ по использованию критериев Андерсона-Дарлинга, Лемана-Розенблатта, Смирнова на данных ограниченной точности.

Разработанная программа может применяться во многих сферах для решения ПРИКЛЗДНЬ1Х задач, СВЯЗЗННЬ1Х С ВЬ1ЯВЛСНИСМ ОДНОРОДНОСТИ данных.

Оглавление

Введение	5
1. Критерии проверки однородности законов распределения	9
1.1. Гипотеза об однородности распределений	9
1.2. Критершй Смирнова	9
1.3. Критершй Лемана-Розенблатта	1 1
1 .4. Критершй Андерсона-Дарлинга	12
1.5. Выводы	13
2. Исследование распределений статистик критериев однородности на данных ограниченной точности	14

2.1. Исследование распределений статистик	14
2.2. Исследование распределения статистики критерия Смирнова	14
2.3. Исследование распределения статистики критерия Лемана-Розенблатта	18
2.4. Исследование распределения статистики критерия Андерсона-Дарлинга	21
2.5. Выводы	24
3. Исследование мощностей критериев однородности на данных ограниченной точности	26
3.1. Исследование мощностей критериев	26
3.2. Исследование мощности критерия Смирнова	29
3.3. Исследование мощности критерия Лемана-Розенблатта	32
3.4. Исследование мощности критерия Андерсона-Дарлинга	35
3.5. Выводы	38
Заключение	39
Список литературы	40
Приложение А. Программные модули	42

Введение

Современное состояние и актуальность темы исследования.

В прикладных исследованиях довольно часто возникает необходимость

выяснить, имеют ли различия генеральные совокупности, из которых взяты

две независимые выборки. В [14] математической статистике данная задача

формулируется как проверка гипотезы об однородности законов

распределения вероятностей. Необходимость проверки данных гипотез

появляется в различных ситуациях, когда хотят удостовериться в

неизменности (или напротив в изменении) статистических свойств некоторого

объекта или процесса после целенаправленного изменения фактора или

факторов (методики, технологии и т.д.), неявным образом влияющих на

исследуемый объект. Иными словами, проверяется изменение или наоборот

сохранение статистических показателей объекта или процесса до некоторого

оказанного воздействия и после с течением времени. Например, надо

выяснить, влияет ли способ упаковки некоторых деталей на заводе на их

потребительские качества через год после хранения. Или [13] другой пример

применения исследований однородности: в маркетинге важно выделить

сегменты потребительского рынка.

В [14] случае если установлена однородность двух выборок, то [13] вполне

вероятно группировка сегментов, из которых они взяты, в один.

В последующем это позволит, например, воплотить в жизнь по отношению к

ним схожую рекламную политику (проводить одни и те же маркетинговые [20]

процедуры и т.п.). В случае если же установлено [20] отличие, то поведение

потребителей в двух сегментах различно, объединять эти сегменты [14]

невозможно, и могут понадобиться различные [14] рекламные компании, своя для

каждого из этих сегментов. [14]

Для решения данной задачи широко используются критерии

однородности. Критерии однородности призваны определить, взяты ли две

(или более) выборки из одного распределения вероятностей. На данный момент

5

существуют множество таких критериев. Критерий однородности Смирнова

предложен в работе [1] и рассмотрен в работах [2, 3]. В русскоязычной

литературе трудно найти упоминания о критерии Андерсона-Дарлинга. Тем не

менее, критерий однородности Андерсона-Дарлинга был подробно рассмотрен

в работах [4, 5]. На ряду с критерием Смирнова на практике частое

применение находит критерий Лемана-Розенблатта [6, 7].

На практике всегда приходится иметь дело с данными ограниченной

точности. Зачастую, это целые числа, или данные с одним, двумя знаками

после запятой. При больших объемах выборок, количество повторений в

выборках становится большим. По сути, в этом случае мы имеем дело с

некоторой дискретной случайной величиной. Становится интересно, можно ли

руководствоваться результатами исследования критериев однородности для таких выборок. Подчиняются ли статистики критериев предельным распределениям, и при каких объемах выборок можно реально пользоваться этими предельными распределениями статистик критериев. Исследования распределений статистик и мощности критериев однородности подробно рассматривались в работах [8 - 11].

Цель и задачи исследований. 3 Целью данной диссертационной работы является 3 исследование распределений статистик и мощности критериев однородности на данных ограниченной точности.

Для достижения сформулированной цели 33 были поставлены и решены следующие задачи 33 :

- разработка программы для исследования методами статистического моделирования распределений статистик критериев и вычисления мощности критериев однородности;
- исследование распределений статистик критериев однородности Андерсона-Дарлинга, Лемана-Розенблатта, Смирнова на данных ограниченной точности;
- сравнительный анализ распределений статистик критериев с предельными функциями распределения;

6

- сравнительный анализ мощности рассматриваемых критериев на данных ОГРАНИЧЕННОЙ ТОЧНОСТИ .

Методы исследования. Для решения поставленных задач 33 использовались 33 методы статистического анализа, теории вероятности, математической статистики и компьютерного моделирования 6 .

Научная новизна диссертационной работы заключается:

- В результатах исследований распределений статистик по данным ограниченной точности;
- В результатах исследования мощности критериев однородности на данных ограниченной точности и В сравнительном анализе с мощностями, полученными по выборкам без ограничений на точность.

Достоверность и обоснованность научных положений, рекомендаций и выводов подтверждается:

- корректным применением математического аппарата и методов статистического моделирования для исследования свойств и распределений статистик критериев 6 ;

11 СОВПЗДСНИСМ результатов СТАТИСТИЧЕСКОГО МОДСЛИРОВЗНИЯ С ИЗВССТНВІМІ/І теоретическими результатами.

Личный творческий вклад автора заключается:

- В 3 формулировании этапов исследования распределений статистик рассматриваемых критериев однородности на данных ограниченной точности;
- в исследовании распределений статистик критериев однородности (проверка близости к предельной функции распределения);
- В вычислении мощности критериев на данных ограниченной точности и сравнение с мощностями на данных без округления;
- В реализации рассматриваемых критериев однородности на языке разработки программного обеспечения Python.

Практическая ценность И реализация результатов работы.

Полученные в работе результаты могут быть использованы в прикладных задачах статистического анализа в задачах по выявлению однородности. Сформулированы рекомендации по использованию критериев

Андерсона-Дарлинга, Лемана-Розенблатта, Смирнова на данных ограниченной точности. Разработанная программа может применяться во многих сферах для решения прикладных задач, связанных с выявлением однородности данных.

Структура работы. 21 Диссертация состоит из введения, 3 глав основного содержания, заключения, списка литературы и приложения с программным кодом. Основная часть содержания изложена на 10-44 страницах, включая 8 рисунков, 28 таблиц и списка литературы из 16 источников 10 . 7 Краткое содержание работы. В первой главе представлены основные определения, необходимые теоретические выкладки, 11 используемые в работе, формулируются задачи исследования.

Во второй главе исследуются 7 распределения статистик критериев Смирнова, Андерсона-Дарлинга, Лемана-Розенблатта, полученные на данных ограниченной точности.

В третьей главе исследуются мощности вышеизложенных критериев.

В заключении приводится перечень основных результатов ИССЛЕДОВАНИЙ, В ПРИЛОЖЕНИИ ПРЕДСТАВЛЕНЬ1 ФРАГМЕНТЫ ИСХОДНЫ1Х ТСКСТОВ

1. Критерии проверки однородности законов распределения

1.1. Гипотеза об однородности распределений

При анализе случайных ошибок средств измерений, при статистическом управлении качеством процессов 3 часто возникают вопросы решения задачи

проверки гипотез о принадлежности двух выборок случайных величин одной и той же генеральной совокупности. 3 Такая задача, естественно, возникает при проверке средств измерений, когда пытаются убедиться в том, что закон распределения случайных ошибок измерений не претерпел существенных изменений с течением времени 6 .

Задача проверки однородности двух выборок формулируется следующим образом. Пусть имеются две выборки: X_1, X_2, \dots, X_n из распределения 9 $F(x)$ и Y_1, Y_2, \dots, Y_m из распределения $G(x)$. Обозначим упорядоченные по неубыванию элементы выборок следующим образом:

$x_1 \leq x_2 \leq \dots \leq x_n$ и $y_1 \leq y_2 \leq \dots \leq y_m$.

Для определенности обычно полагают, что $m \leq n$, но это совсем необязательно. Проверяется гипотеза о том, что обе 32 выборки извлечены из одной и той же генеральной совокупности 31 , т. е.

$H_0: F(x) = G(x)$

при любом x .

1.2. Критерий Смирнова

Критерий Смирнова — это правосторонний критерий проверки нулевой гипотезой о том, что из одного и того же непрерывного распределения извлекаются 2 независимые выборки. Критерий однородности Смирнова предложен в работе [1]. Предполагается, что функции распределения $F(x)$ и $G(x)$ являются непрерывными. Статистика критерия Смирнова измеряет расстояние МСЖДУ ЭМПИРИЧЕСКИМИ ФУНКЦИЯМИ распределения, построенными по выборкам 1 [1]

$D_n = \sup |F_n(x) - P_n(x)|$ -

9

На практике ЭНЗЧНИС СТАТИСТИКИ D_n рекомендуется ВЫЧИСЛЯТЬ В соответствии с соотношениями [8]:

$D_n = \max_{1 \leq k \leq n} |F_n(x_k) - P_n(x_k)|$

$D_n = \max_{1 \leq k \leq n} |F_n(x_k) - P_n(x_k)|$

$D_n = \max_{1 \leq k \leq n} |F_n(x_k) - P_n(x_k)|$

$D_n = \max_{1 \leq k \leq n} |F_n(x_k) - P_n(x_k)|$

D =maX(D+ Dj)

m,n

Если гипотеза H0 справедлива, то при неограниченном увеличении

объемов выборок [12] $\lim_{m,n \rightarrow \infty} P(D_{mn} \leq S) = K(S)$, т. е. статистика

$T \rightarrow_{\infty} T + n'$

$SC : T \cap D_m \cap T$

$m + n'$

В пределе подчиняется распределению Колмогорова $K(S)$ [12] с функцией

распределения

$K(S) = \tilde{E}(-1)^{\circ} e_{-2}{}^{\circ} 2^2$.

$[(= - \infty$

ОДНЗКО при ОФРАНИИНСНВІХ значениях т и п СЛУЧАЙНЫС ВСЛИЧИНЬ1 В;" И

Dm п ЯВЛЯЮТСЯ ДИСКРСТНЁЛМИ, И МНОЖССТВО ИХ ВОЗМОЖНЁЛХ значений

представляет собой решетку с шагом 1/ k, где k — наименьшее общее кратное

т и п [12]. Условное распределение G(SC | H0) статистики SC при верности

гипотезы H0 медленно сходится к K(S) и имеет существенное отличие от

него 1 при 2 малых значениях т и п .

Гладкость распределения статистики сильно зависит от величины k.

Поэтому предпочтительнее применять критерий, когда объемы выборок т и

п не равны и представляют собой взаимно простые числа. В таких случаях

наименьшее общее кратное т и п максимально и равно k : тп, а распределе-

10 1

ние статистики больше напоминает непрерывную функцию распределения 1 .

1.3. Критерий Лемана-Розенблатта

Критерий однородности Лемана—Розенблатта представляет собой

критерий типа а)2. Критерий 2 предложен В работе [13] и исследован В [14].

Статистика критерия имеет вид [12]

$C_l)$

$m \cap 2$

$T = \text{Палю} \text{—} \text{шт}] \text{ от} \dots \text{.} (\text{х} \text{ъ}$

$t + n_w$

$m \cap$

где $H_{\dots}(x) = Gm(x) + F,1(x)$ — эмпирическая функция

$t + n \cap t + n$

распределения, построенная по вариационному ряду объединения двух

выборок. Статистика T используется В форме 9 [12]

$1^{\circ} \text{ } .2 \cap .2 \text{ } 4 \text{—} 1$

$T = \text{—} [: \text{'ZZ} \text{I} (\text{' ; —} \text{I}) + m \text{Z} = 1 \text{1} < s \text{j —} \text{j}]] \text{—} 6 (n n : L + n) \text{j} \text{ (1.1)}$

где r l. — порядковыи номер (ранг) y i; S j. — порядковыи номер (ранг) x j. В объе-

ДИНСННОМ ВарI/IaIII/ИОННОМ ряде.

В [15] было показано, что статистика (1.1) В пределе распределена как

$a_1(t) :$

$\lim P\{T < 1\} = a_1(t).$

Функция распределения a1(t) имеет вид [12]:

$\text{—} {}^{\circ} 1^{\circ} (\text{' j + 1/2} \text{j/4} \text{j + 1 } _ { \text{' 4} \text{j + 1} } \text{ } 2 \text{ } \times$

ако—Ы,; $\Gamma(1/2)\Gamma(1+1)$ "p[161 }

$x1111\text{—}11111\text{—}1011$

$7 \text{ } 161 \text{ } д \text{ } 161$

где l 1(-),11(-) — модифицированные функции Бесселя вида

$7 \text{ } Я$

11

@121

...И): $20 F(k + 1)r(k + v + 1)$

$< \infty, < 7 \text{г}.$

В отличие от критерия Смирнова распределение статистики T быстро

2.1. Исследование распределений статистик

Так как цель исследования заключается в исследовании распределения статистик на данных ограниченной точности, нужно моделировать такие данные. Значения моделируемых выборок ограничивались до целого числа, до одного, двух знаков после запятой: сначала генерируется выборка заданного размера и производится округление значений.

Целью данной главы является проведение исследования, с целью выяснить, можно ли использовать критерии, если данные получены с ограниченной точностью, подчиняются ли статистики, вычисленные по таким данным, соответствующим предельным законам распределения рассматриваемых критериев однородности.

Зададимся величиной расстояния, равной 0.05, при котором будем считать, что распределение статистик все еще подчиняется предельному закону распределения.

Обозначим некоторые величины для таблиц с результатами исследований:

- количество выборок $N = 16600$,
- $\rho = \sup |F_n(x) - F(x)|$ - расстояние между эмпирическими и предельными функциями распределения статистик КРИТСРИЯ В метрике Колмогорова.

2.2. Исследование распределения статистики критерия Смирнова

В таблицах 2.1, 2.2 исследования проводились на сгенерированных данных, обе выборки, в которых, подчинялись стандартному нормальному закону распределения с плотностью

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

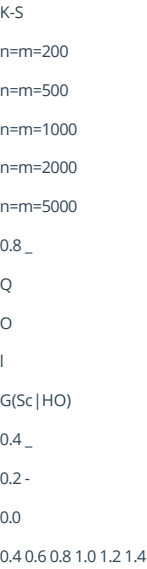
и параметрами сдвига ($\mu = 0$ и масштаба $\sigma^2 = 1$).

14 2
Таблица 2.1 — Результаты для критерия однородности Смирнова, округление до 2 знаков, $n = m$, выборки из нормального закона распределения с параметрами ($\mu = 0$, $\sigma^2 = 1$).

n, m среднее число различных значений x_i в объединенной выборке

200, 200	0.09 246.5
500, 500	0.07 368.5
1000, 1000	0.07 449.0
2000, 2000	0.07 507.0
5000, 5000	0.06 580.5
10000, 10000	0.05 626.5

На рисунке 2.1 представлена графическая иллюстрация результатов, представленных в таблице 2.1.



SC

Рисунок 2.1 — Распределения статистики критерия Смирнова при справедливости H0 В зависимости от m и 11, округление до 2 знаков

Как видно из таблицы 2.1, наблюдается уменьшение расстояния с ростом объема выборок при одинаковых размерах обеих выборок, в отличие от других критериев. В связи с этим были проведены дополнительные исследования критерия Смирнова на объемах выборок 10000. Даже при таких размерах моделируемых выборок расстояние имеет тенденцию к уменьшению. И тем не менее, заданное расстояние 0.05 между функциями распределения начинает достигаться лишь при объеме выборок 10000.

15

В таблице 2.2 обе выборки также принадлежали стандартному нормальному закону распределения, но при различных объемах выборок.

Таблица 2.2 — Результаты для критерия однородности Смирнова, округление до 2 знаков, p фт, выборки из нормального закона распределения с параметрами (90 : 0, (91 : 1.

n, m,0 среднее число различных значений в объединенной

выборке

500, 500 0.07 368.5

500, 1000 0.07 421.05

500, 2000 0.06 468.0

500, 5000 0.05 533.5

Суммируя результаты по таблице 2.2, можно заметить, что при различных объемах выборок, с увеличением объема второй выборки и при зафиксированном значении объема первой, расстояния оказываются меньшими, чем когда объемы двух выборок одинаковые (табл. 2.1).

В предыдущих исследованиях было замечено, что расстояния между эмпирической функцией распределения и предельной функцией распределения статистики критерия оказывались неприемлемо большими на данных ограниченных до целых чисел и до одного знака. Это могло быть связано с большим количеством повторений в выборке. Поэтому, для данных, ограниченных до целых чисел и одного знака, были проведены исследования на данных с большим количеством уникальных значений при тех же объемах выборок, что и в исследовании на данных ограниченных до двух знаков. С этой целью, выборки генерировались из распределения, с большей дисперсией, чем стандартное нормальное. Величина дисперсии подбиралась эмпирическим путем, чтобы ее величина была максимально приближена к единице и, чтобы расстояние не превышало 0.05.

16

Таблица 2.3 — Результаты для критерия однородности Смирнова, округление до одного знака, 11=ш, выборки из нормального закона распределения с параметрами до : 0, d] : 50 .

n, m,0 среднее число различных значений в объединенной

выборке

200, 200 0.08 360.5

500, 500 0.05 772.5

1000, 1000 0.04 1228.5

2000, 2000 0.04 1719.5

5000, 5000 0.03 2243.0

10000, 10000 0.03 2546.0

Таблица 2.4 — Результаты для критерия однородности Смирнова, округление до целых, 11=ш, выборки из нормального закона распределения с параметрами 6'0 : 0, 6'1 : 100 .

n, m,0 среднее число различных значений в объединенной

выборке

200, 200 0.1 243.5

500, 500 0.08 368.5

1000, 1000 0.08 448.5

2000, 2000 0.07 508.5

5000, 5000 0.06 578.0

10000, 10000 0.06 628.0

Для данных, ограниченных до одного знака и до целых (табл. 2.3, 2.4), при увеличении количества уникальных значений, за счет увеличения дисперсии закона распределения моделируемых выборок, в объединенной выборке расстояния становятся схожими с результатами, полученными на данных, ограниченных до Двух знаков (табл. 2.1).

В силу особенностей критерия Смирнова, упомянутых в главе 1, было необходимо провести исследования на выборках, размеры которых представляются как взаимно простые числа. Объемы выборок подбирались с максимальной схожестью объемов выборок из предыдущих исследований.

17

Таблица 2.5 — Результаты для критерия однородности Смирнова, округление до двух знаков, объемы выборок взаимно простые, выборки из нормального закона распределения с параметрами

a=цa=г

11, m ,0 среднее число различных значений в объединенной выборке

199, 201 0.05 248.5

499, 501 0.06 375.5

999, 1001 0.05 455.5

1999, 2001 0.06 507.5

4999, 50001 0.06 573.5

9999, 10001 0.06 624.0

Для взаимно простых n и m расстояния от функции распределения статистик до предельного не имеют существенных отличий в сравнении с предыдущими исследованиями из табл. 2.1.

2.3. Исследование распределения статистики критерия Лемана-Розенблатта

В таблицах ниже (2.6-2.12) представлены значения расстояний между эмпирическими и предельными функциями распределения статистик, рассчитанные по метрике Колмогорова для критерия Лемана-Розенблатта.

В таблицах 2.6, 2.7 обе выборки также принадлежали стандартному нормальному закону распределения, но при различных объемах выборок.

Таблица 2.6 — Результаты для критерия однородности Лемана-Розенблатта, округление до 2 знаков,

11=ш, выборки из нормального закона распределения с параметрами (90 = 0,191 =1 .

n, m,0 среднее число различных значений в объединенной выборке

200, 200 0.01 243.5

500, 500 0.01 369.5

1000, 1000 0.01 448.5

2000, 2000 0.01 510.5

5000, 5000 0.01 578.5

18

На рисунке 2.2 представлена графическая иллюстрация результатов, представленных в таблице 2.6.

G(Sc[H0])

0:1 0:2 0:3 0:4 035

SC

Рисунок 2.2 — Распределения статистики критерия Лемана-Розенблатта при

справедливости H0 В зависимости от m и n, округление до 2 знаков

Таблица 2.7 — Результаты для критерия однородности Лемана-Розенблатта, округление до 1 знака, 11=ш, выборки из нормального закона распределения с параметрами (90 : 0,191 : 1.

n, m,0 среднее число различных

значенш `и`1 в объединенной

выборке

200, 200 0.01 49.0

500, 500 0.01 56.5

1000, 1000 0.01 62.5

2000, 2000 0.01 67.5

5000, 5000 0.01 72.5

Судя по результатам из таблиц 2.6 и 2.7, распределение статистик для критерия Лемана-Розенблатта довольно близко располагается с предельным распределением. Для выборок, округленных до двух и одного знаков, выполняется условие не превышения расстояния в 0.05.

В таблицах 2.8-2.11 обе выборки также принадлежали стандартному нормальному закону распределения, но при различных объемах выборок.

19

Таблица 2.8 — Результаты для критерия однородности Лемана-Розенблатта, округление до 2 знаков,

nm, выборки из нормального закона распределения с параметрами (90 20,191 :1, при малых

объемах выборок.

n, m,0 среднее число различных

значеншй В объединенной

выборке

30, 30 0.02 55.5

30, 40 0.01 63.0

3 0, 50 0.01 71.5

Таблица 2.9 — Результаты для критерия однородности Лемана-Розенблатта, округление до 2 знаков,

n ft, выборки из нормального закона распределения с параметрами (90 20,191 :1, при больших

объемах выборок.

n, m,0 среднее число различных

значеншй В объединенной

выборке

500, 500 0.01 369.5

500, 1000 0.01 418.0

500, 2000 0.03 469.0

500, 5000 0.24 535

Таблица 2.10 — Результаты для критерия однородности Лемана-Розенблатта, округление до 1 знака,

nm, выборки из нормального закона распределения с параметрами (90 20,191 :1, при малых

объемах выборок.

n, m,0 среднее число различных

значеншй В объединенной

выборке

30, 30 0.02 30.5

30, 40 0.02 33.0

30, 50 0.03 34.5

Таблица 2.11 — Результаты для критерия однородности Лемана-Розенблатта, округление до 1 знака,

n ft, выборки из нормального закона распределения с параметрами (90 20,191 :1, при больших

объемах выборок.

n, m,0 среднее число различных

значеншй В объединенной

выборке

500, 500 0.01 56.5

500, 1000 0.35 60.0

500, 2000 0.94 63.0

500, 5000 0.99 70.0

Судя по результатам Данных таблиц Для критерия Лемана-Розенблатта не
20

наблюдается приближения распределения статистик к предельному закону при
различных объемах выборок в сравнении с результатами, полученными при
одинаковых объемах выборок 9 . Для таблицы 2.11 эти выводы проявляются в
наибольшей степени.

Для данных, ограниченных до целых чисел, были проведены исследования
на данных с большим количеством уникальных значений при тех же объемах
выборок, что и в исследовании на данных ограниченных до двух и одного
знаков. С этой целью, выборки генерировались из распределения, с большей
дисперсией, чем стандартное нормальное.

Таблица 2.12 — Результаты для критерия однородности Лемана-Розенблатта, округление до целых,
11=ш, выборки из нормального закона распределения с параметрами до : 0, 61 = 10.

n, m,0 среднее число различных
значенш` и'1 в объединенной
выборке

200, 200 0.01 51.0

500, 500 0.01 57.0

1000, 1000 0.01 63.5

2000, 2000 0.01 67.5

5000, 5000 0.01 69.0

Для данных, ограниченных до целых, при увеличении количества
уникальных значений, за счет увеличения дисперсии закона распределения
моделируемых выборок, в объединенной выборке расстояния становятся
схожими с результатами, полученными на данных, ограниченных до одного и
двух знаков.

2.4. Исследование распределения статистики критерия Андерсона-
Дарлинга

В таблицах ниже (2.13-2.18) представлены значения расстояний между
эмпирическими и предельными функциями распределения статистик,
рассчитанные по метрике Колмогорова для критерия Андерсона-Дарлинга.

В таблицах 2.13-2.16 исследования проводились на сгенерированных

21

Данных, обе выборки, в которых, подчинялись стандартному нормальному
закону.

Таблица 2.13 — Результаты для критерия однородности Андерсона-Дарлинга, округление до 2 знаков,
11=ш, выборки из нормального закона распределения с параметрами (90 = 0,191 =1 .

n, m,0 среднее число различных
значеншй в объединенной
выборке

200, 200 002 241.0

500, 500 002 377.0

1000, 1000 0.03 442.0

2000, 2000 0.04 510.0

5000, 5000 008 576.5

Как видно из таблицы, с увеличением объемов выборок расстояние между
эмпирической функцией распределения и предельной функцией распределения
статистики критерия увеличивалось. По результатам, представленным в
таблице 2.13, видно, что между n=2000 и n=5000 расстояние
становится большим чем 0.05 на данных, округленных до двух знаков.

На рисунке 2.3 представлена графическая иллюстрация результатов,
представленных в таблице 2.3.

0.8 -

0.6 -

G(SCINO)

0.4 -

0.2 -

0.0'..

I

0.5 1.0 1.5 2.0 2.5

SC

Рисунок 2.3 — Распределения статистики критерия Андерсона-Дарлинга при

справедливости H0 В зависимости от m и n, округление до 2 знаков

При округлении до целых и до одного знака после запятой наблюдалась

такая же ТСНДСНЦИЯ УВСЛИЧСНИЯ расстояния С УВСЛИЧСНИСМ ОБЪСМОВ выборок.

22

Но величина расстояния была около единицы и около 05 соответственно, что является показателем, что функции распределения лежат Далеко Друг от друга.

В таблице 2.14 обе выборки также принадлежали стандартному нормальному закону распределения, но при различных объемах выборок.

Таблица 2.14 — Результаты для критерия однородности Андерсона-Дарлинга, округление до 2 знаков,

n 75 т, выборки из нормального закона распределения с параметрами (90 = 0, 191 = 1.

n, m,0 среднее число различных

значений в объединенной

выборке

500, 500 0.02 377.0

500, 1000 0.01 422.0

500, 2000 0.01 465.0

500, 5000 0.01 532.5

Из результатов таблицы 2.14 наблюдается схожая картина с аналогичными

исследованиями критерия Смирнова. При различных объемах выборок, с увеличением объема второй выборки при зафиксированном значении объема

первой, расстояния оказываются меньшими, чем когда объемы двух выборок

одинаковые (табл. 2. 13).

Для данных, ограниченных до целых чисел и одного знака, были

проведены исследования на данных с большим количеством уникальных

значений при тех же объемах выборок, что и в исследовании на данных

ограниченных до двух знаков. С этой целью, выборки генерировались из

распределения, с большей дисперсией, чем стандартное нормальное.

Таблица 2.15 — Результаты для критерия однородности Андерсона-Дарлинга, округление до 1 знака,

11=ш, выборки из нормального закона распределения с параметрами до : 0, 61 = 10.

n, m,0 среднее число различных

значений в объединенной

выборке

200, 200 0.02 249.0

500, 500 0.02 374.5

1000, 1000 0.02 442.5

2000, 2000 0.04 503.5

5000, 5000 0.09 579.0

23

Таблица 2.16 — Результаты для критерия однородности Андерсона-Дарлинга, округление до целых,

11=ш, выборки из нормального закона распределения с параметрами до : 0, 61 = 80 .

n, m,0 среднее число различных

значений в объединенной

выборке

200, 200 0.01 221.0

500, 500 0.03 321.0

1000, 1000 0.04 374.0

2000, 2000 0.06 421.5

5000, 5000 0.12 475.0

Анализируя результаты, представленные в таблицах для критерия

Андерсона-Дарлинга, можно заметить тенденцию, что при уменьшении

отношения числа различных значений в объединенной выборке к общему

объему объединенной выборки, увеличивается расстояние между

распределениями эмпирической функции распределения статистик и

предельным распределением.

2.5. Выводы

Суммируя полученные результаты исследования распределения статистик

по всем критериям на данных ограниченной точности, можно сделать

следующие выводы:

— для критерия Смирнова наблюдается сходимость распределения

статистики к предельному распределению статистики при

увеличении объема выборок. Также, было замечено, что с

увеличением объема второй выборки и при зафиксированном

значении объема первой, расстояния оказываются меньшими, чем

когда объемы двух выборок одинаковые. Для данных, ограниченных

до одного знака и до целых, при увеличении количества уникальных

значений, за счет увеличения дисперсии закона распределения

моделируемых выборок, в объединенной выборке расстояния

24

становятся схожими с результатами, полученными на данных,

ограниченных до двух знаков.

Для критерия Лемана-Розенблатта была замечена следующая

особенность, что при зафиксированном распределении

генерируемых выборок и при различных n (200, 500, ...) не

меняется расстояние между эмпирическим распределением

статистики и предельным. Расстояние меняется лишь при изменении

дисперсии распределения выборок.

Для критерия Андерсона-Дарлинга можно заметить тенденцию, что

при уменьшении числа различных значений в объединенной выборке

к общему объему объединенной выборки, увеличивается расстояние

между эмпирической функцией распределения статистик и

предельным распределением. При более точном исследовании

можно попытаться получить предельное соотношение числа

уникальных элементов в объединенной выборке к общему числу

элементов при заданном расстоянии между эмпирической функцией

распределения и предельной функцией распределения статистики

критерия Андерсона-Дарлинга.

25

3. Исследование мощностей критериев однородности на данных

ограниченной точности

3.1. Исследование мощностей критериев

Очевидно, что при проверке любой статистической гипотезы

предпочтительней использовать наиболее мощный критерий. Статистическая

мощность в математической статистике является показателем вероятности

отклонения основной (или нулевой) гипотезы при проверке статистических

гипотез в случае, когда нулевая гипотеза неверна (верна альтернативная

гипотеза).

Мощность критериев проверки однородности исследовалась в случае

ряда альтернатив. Для определенности проверяемой гипотезе H_0

соответствовала принадлежность выборок одному и тому же стандартному

НОРМАЛЬНОМУ ЗАКОНУ РАСПРЕДЕЛЕНИЯ С ПЛОТНОСТЬЮ

$\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

$N(\mu, \sigma^2)$ — д1 & exp —295

и параметрами сдвига (90 = 0 и масштаба 61 = 1.

При всех конкурирующих гипотезах первая выборка всегда

соответствовала стандартному нормальному закону, а вторая — некоторому

ДРУГОМУ.

В частности, при альтернативе сдвига В случае конкурирующей гипотезы

H1 вторая выборка соответствовала нормальному закону 1 с параметром сдвига

60 = 0.1 и параметром масштаба 61 = 1, В случае 2 конкурирующей 5 гипотезы H 2

— нормальному закону с параметрами 90 = 0.5 и 61 = 1. 2

При изменении масштаба В случае конкурирующей гипотезы H3 вторая

выборка соответствовала нормальному закону с 5 параметрами (90 = 0 и 61 = 1 2 .1,

В случае конкурирующей гипотезы H 4 — нормальному закону с параметрами

(90 = 0 и 191 = 1.5 .

В случае конкурирующей гипотезы H5 вторая выборка соответствовала

26 1

ЛОГИСТИЧЕСКОМУ ЗАКОНУ С ПЛОТНОСТЬЮ 1

2

$f(x) = \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) + \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

61 1

И параметрами (90 = 0 И 61 = 1.

На рисунках ниже (3.1 — 3.5) представлены графики теоретических распределений, из которых генерировались выборки для исследования мощностей.

1.0 — F(x | H0)

— F(x | H1)

0.8 -

0.6 <

F(x | H)

0.4 '

0.2 _

0.0 —

—4 —2 0

X

Рисунок 3.1 — Графики функций распределения законов, из которых генерировалась вторая выборка в гипотезах H 0, H 1.

24

1.0 < — f(x | H0)

— F(x | H2)

0.8 <

0.6 _

F(x | H)

0.4 -

0.2 _

0.0 -

|||||

—4 —2 0 2 4

X

Рисунок 3.2 — Графики функций распределения законов, из которых генерировалась вторая выборка в гипотезах H 0, H 2 .

27

F(x | H)

F(x | H)

F(x | H)

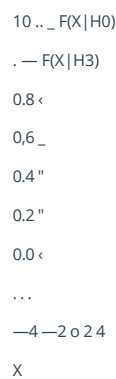


Рисунок 3.3 — Графики функций распределения законов, из которых генерировалась вторая выборка в гипотезах H_0 , H_3 .

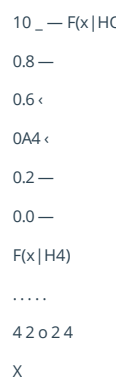


Рисунок 3.4 — Графики функций распределения законов, из которых генерировалась вторая выборка в гипотезах H_0 , H_4 .

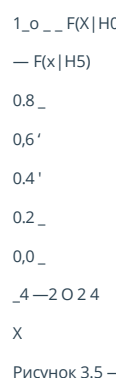


Рисунок 3.5 — Графики функций распределения законов, из которых генерировалась вторая выборка в гипотезах H_0 , H_5 .

28

3.2. Исследование мощности критерия Смирнова

В таблицах 3.1 — 3.4 представлены рассчитанные оценки мощностей

критерия однородности Смирнова ($1 - \beta$, где β - вероятность ошибки второго

вида). относительно конкурирующих

H_1 — H_5 для различных значений объемов генерируемых выборок. Значения

ОЦЕНОК МОЩНОСТИ также ПРЕДСТАВЛЕНЫ В ЗАВИСИМОСТИ ОТ различных ЗНАЧЕНИЙ

заданных уровней

$\alpha : 0.1, 0.05, 0.25$.

значимости (вероятностей ошибок первого рода): 2

Таблица 3.1 — 2 Мощность критерия Смирнова относительно гипотез H_1 — H_5 в зависимости от

объемов выборок ($n : \tau$) на данных без округления

Уровень

ЗНАЧИМОСТИ

α

$n=500$

$n=1000$

$n=2000$

Относительно альтернативы 1 H_1

0.1 0.38 0.61 0.87

0.05 0.27 0.48 0.78

0.025 0.18 0.36 0.69

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.19 0.30 0.55

0.05 0.10 0.17 0.36

0.025 0.05 0.09 0.21

Относительно альтернативы Н 4

0.1 0.99 1.0 1.0

0.05 0.99 1.0 1.0

0.025 0.98 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 0.99 1.0 1.0

29

Таблица 3.2 — Мощность критерия Смирнова относительно гипотез Н1 —Н5 в зависимости от объемов выборок (п : т) на данных, округленных до целых чисел

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы 1 Н1

0.1 0.37 0.58 0.84

0.05 0.25 0.45 0.74

0.025 0.17 0.34 0.63

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.19 0.30 0.54

0.05 0.11 0.17 0.33

0.025 0.06 0.10 0.20

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 0.99 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Таблица 3.3 — Мощность критерия Смирнова относительно гипотез Н1 —Н5 в зависимости от объемов выборок (п : т) на данных, округленных до 1 знака после запятой

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы 1 Н1

0.1 0.39 0.61 0.87

0.05 0.28 0.49 0.78

0.025 0.19 0.37 0.69

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0
30

Продолжение таблицы 3.3
Относительно альтернативы Н 3

0.1 0.20 0.31 0.57
0.05 0.11 0.18 0.37
0.025 0.05 1.0 0.22

Относительно альтернативы Н 4
0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 0.99 1.0 1.0

Относительно альтернативы Н 5
0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Таблица 3.4 — Мощность критерия Смирнова относительно гипотез Н1 —Н5 в зависимости от
объемов выборок (n : т) на данных, округленных до 2 знаков после запятой

Уровень
значимости 11 = 500 11 = 1000 11 = 2000

а
Относительно альтернативы 1 Н1

0.1 0.38 0.60 0.87
0.05 0.27 0.47 0.78
0.025 0.18 0.37 0.68

Относительно альтернативы Н 2
0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Относительно альтернативы Н 3
0.1 0.19 0.29 0.56
0.05 0.11 0.16 0.35
0.025 0.05 0.09 0.27

Относительно альтернативы Н 4
0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 0.98 1.0 1.0

Относительно альтернативы Н 5
0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

31

Суммируя результаты, полученные по таблицам 3.1 — 3.4, можно сказать,
что на округленных Данных мощность получалась выше для гипотезы Н3
ПОЧТИ ВО ВСЕХ случаях.

3.3. Исследование мощности критерия Лемана-Розенблатта
В таблицах 3.5 — 3.8 представлены рассчитанные оценки мощностей
критерия однородности Лемана-Розенблатта. Значения оценок мощности
представлены относительно конкурирующих гипотез Н1—Н5 для различных
значений объемов выборок, также в зависимости от различных значений
заданных уровней значимости: а : 0.1,0.05,0.25 .

Таблица 3.5 — Мощность критерия Лемана-Розенблатта относительно гипотез Н1—Н5 в
зависимости от объемов выборок (n : т) на данных без округления

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы **1** Н1

0.1 0.44 0.68 0.91

0.05 0.32 0.56 0.85

0.025 0.23 0.44 0.78

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.19 0.32 0.62

0.05 0.10 0.16 0.41

0.025 0.05 0.07 0.23

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

32

Таблица 3.6 — Мощность критерия Лемана-Розенблатта относительно гипотез Н1—Н5в
заВІ/ІСІ/МОСТІ/І ОТ ОБЪСМОВ ВЪ16ОРОК (п : т) на ДЗННЪ1Х, ОртІЅСННВІХ ДО ЦСЛЪ1Х ЧИССЛ

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы Н1

0.1 0.41 0.63 0.88

0.05 0.30 0.52 0.81

0.025 0.22 0.41 0.73

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.21 0.32 0.60

0.05 0.18 0.18 0.40

0.025 0.05 0.10 0.25

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Таблица 3.7 — Мощность критерия Лемана-Розенблатта относите.

1ьно гипотез Н1—Н5в

1 знака

зависимости от объемов выборки (п : т) на данных, округленных до

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы Н1

0.1 0.44 0.68 0.91
0.05 0.32 0.56 0.84
0.025 0.22 0.45 0.77

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

33

Продолжение таблицы 3 .7

Относительно альтернативы Н 3

0.1 0.20 0.32 0.62
0.05 0.10 0.17 0.40
0.025 0.05 0.08 0.24

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Таблица 3.8 — Мощность критерия Лемана-Розенблатта относите.

зависимости от объемов выборок (п : т) на данных, округленных до

1ьно гипотез Н1—Н5в

2 знаков после запятой

Уровень

ЗНАЧИМОСТИ

а

п=500

11: 1000

11 = 2000

Относительно альтернативы Н1

0.1 0.44 0.68 0.91
0.05 0.32 0.57 0.85
0.025 0.23 0.46 0.77

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.19 0.33 0.61
0.05 0.09 0.17 0.40
0.025 0.05 0.08 0.22

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0
0.05 1.0 1.0 1.0
0.025 1.0 1.0 1.0

34

Суммируя результаты, полученные по оценкам мощностей критерия

Лемана-Розенблатта, можно сказать, что на округленных Данных мощность

ПОЛУЧЗЛЗСЬ ВЫШС ДЛЯ ГИПОТСЗЫ1 НЗ ПОЧТИ ВО ВССХ СJlyanX.

3.4. Исследование мощности критерия Андерсона-Дарлинга

В таблицах 3.9 — 3.12 представлены рассчитанные оценки мощностей

критерия однородности Значения

Андерсона-Дарлинга. представлены

относительно конкурирующих гипотез Н1—Н5- Значения оценок мощности

также представлены В ЗЗВИСИМОСТИ ОТ различных ЗНЗЧСНИЙ заданных УРОВНСЙ

значимости: а = 01,005,025.

Таблица 3.9 — Мощность критерия Андерсона-Дарлинга относительно гипотез Н1—Н5в

зависимости от объемов выборки (п : т) на данных без округления

Уровень

ЗНЗЧИМОСТИ

а

п=500

11: 1000

11 = 2000

Относительно альтернативы 1 Н1

0.1 0.44 0.69 0.92

0.05 0.33 0.57 0.86

0.025 0.24 0.46 0.79

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.28 0.53 0.86

0.05 0.15 0.34 0.71

0.025 0.08 0.19 0.54

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

35

Таблица 3.10 — Мощность критерия Андерсона-Дарлинга относительно гипотез Н1 —Н5 в

ЗЗВИСИМОСТИ ОТ ОБЪСМОВ ВЫБОРОК (п : т) на ДЗННЪ1Х, ОКРУГЛСННЪ1Х ДО ЦСЛЬ1Х ЧИССЛ

Уровень

значимости 11 = 500 11 = 1000 11 = 2000

а

Относительно альтернативы Н1

0.1 0.50 0.71 0.91

0.05 0.35 0.57 0.82

0.025 0.23 0.44 0.73

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.03 0.01 0.00

0.05 0.01 0.00 0.00

0.025 0.00 0.00 0.00

Относительно альтернативы Н 4

0.1 0.0 0.0 0.0

0.05 0.0 0.0 0.0

0.025 0.0 0.0 0.0

Относительно альтернативы Н 5

0.1 0.0 0.0 0.0

0.05 0.0 0.0 0.0

0.025 0.0 0.0 0.0

Из таблицы 3.10 видно, что критерий Андерсона-Дарлинга оказывается смещенным (мощность меньше задаваемого уровня значимости) в случае альтернатив с пересечением функций распределения (рис. 3.3-3.5).

Таблица 3.11 — Мощность критерия Андерсона-Дарлинга относительно гипотез Н1 —Н5 в зависимости от объемов выборок (п : т) на данных, округленных до 1 знака после запятой

Уровень

ЗНАЧИМОСТИ

а

п=500 11: 1000 п=2000

Относительно альтернативы 1 Н1

0.1 0.60 0.81 0.96

0.05 0.46 0.71 0.93

0.025 0.34 0.60 0.88

36

Продолжение таблицы 3.11

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.13 0.17 0.22

0.05 0.07 0.09 0.12

0.025 0.03 0.045 0.07

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 0.98 1.0 1.0

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Таблица 3.12 — Мощность критерия Андерсона-Дарлинга относительно гипотез Н1—Н5 в зависимости от объемов выборок (п : т) на данных, округленных до 2 знаков после запятой

Уровень

ЗНАЧИМОСТИ

а

п=500

11: 1000

11 = 2000

Относительно альтернативы 1 Н1

0.1 0.50 0.75 0.95

0.05 0.39 0.64 0.91

0.025 0.29 0.54 0.87

Относительно альтернативы Н 2

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

Относительно альтернативы Н 3

0.1 0.29 0.50 0.84

0.05 0.16 0.31 0.70

0.025 0.08 0.18 0.52

Относительно альтернативы Н 4

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

37

Продолжение таблицы 3.12

Относительно альтернативы Н 5

0.1 1.0 1.0 1.0

0.05 1.0 1.0 1.0

0.025 1.0 1.0 1.0

По данным, полученным из таблиц исследования мощности критерия

Лемана-Розенблатта, также, можно заметить, что на данных, округленных до одного и двух знаков после запятой по альтернативе Н1 мощность оказалась выше чем на Данных без округления.

3.5. Выводы

Суммируя полученные результаты оценки мощностей по всем критериям

на данных ограниченной точности, можно сделать следующие выводы:

— На данных, округленных до целых чисел, как наиболее мощный

критерий себя показал критерий Лемана-Розенблатта по всем

предложенным альтернативам, кроме гипотезы Н1, где наибольшую

мощность продемонстрировал критерий Андерсона-Дарлинга;

— На данных, округленных до одного знака после запятой, на

предложенных альтернативах оказалось трудно явно определить

наиболее мощный критерий. По альтернативной гипотезе Н1

наибольшую мощность, как и на данных, округленных до целых,

показал критерий Андерсона-Дарлинга. По альтернативе Н3

наибольшую мощность проявил критерий Смирнова. По всем

остальным гипотезам наиболее мощным оказался критерий Лемана-

Розенблатта;

— На данных, округленных до двух знаков после запятой, наибольшую

мощность по всем представленным альтернативам показал критерий

Андерсона-Дарлинга;

— На данных, округленных до целых, критерий Андерсона-Дарлинга

оказался СМСЦСННЬ1М.

38

Заключение

В СООТВЕТСТВИИ С НСНЮ данной работы ПОЛУЧСНЫ СЛСДУЮЩИС ОСНОВНЫС
результаты:

1)

2)

3)

РаЗраб0ТаНа программа для проведения исследований с помощью

методов имитационного моделирования распределений статистик и

мощности критериев однородности в случае данных ограниченной

точности.

В результате исследования распределений статистик показано, что:

— для критерия Андерсона-Дарлинга расстояние между

эмпирической функцией распределения статистики и предельным

уменьшается с ростом отношения числа различных значений в

объединенной выборке к объему объединенной выборки;

— для критерия Лемана-Розенблатта распределения статистики

остаются близкими к предельному закону при равных объемах выборок, однако при пст расстояние между эмпирической функцией распределения статистики и предельным увеличивается с ростом объема объединенной выборки;

— для критерия Смирнова наблюдается медленная сходимость распределения статистики к предельному закону при увеличении объемов выборок.

На данных ограниченной точности наибольшую мощность среди рассмотренных критериев показали критерии Андерсона-Дарлинга и Лемана-Розенблатта. Однако В случае округления наблюдений В выборках до целых критерий Андерсона-Дарлинга оказался смещенным относительно конкурирующих гипотез с пересечением функций распределения.

Обобщая полученные результаты, можно сделать вывод 33 о предпочтительности использования критерия Лемана-Розенблатта при равных объемах выборок $n : t$.

39

Список литературы

- 1) Смирнов Н.В. Оценка расхождения между эмпирическими кривыми распределения в двух независимых выборках / Н.В. 4 Смирнов // Бюллетень МГУ, серия 4 А. — 1939. — Т.2. №2. — С.3-14.
- 2) Massey, F. J. The Kolmogorov-Smimov Test for Goodness of Fit. / F. J. Massey/ Journal of the American Statistical Association. Vol. 46, No. 253, 1951, pp. 68—78.
- 3) 15 Miller, L. H. Table of Percentage Points of Kolmogorov Statistics. / L. H. Miller/ Journal of the American Statistical Association. Vol. 51, No. 273, 1956, pp. 111—121.
- 4) 15 Anderson T. W. Asymptotic theory of certain «goodness of 23 fit» criteria based on stochastic processes / T. W. Anderson, D. A. Darling // Ann. Math. Statist 7 . — 1952. — V. 23. — P. 193—212.
- 5) Anderson T. W. A 17 test of goodness of 17 fit / T. W. Anderson, D. A. Darling // J. Amer. Statist. Assoc., 1954. — V. 29. — P. 765—769.
- 6) Lehman S. Exact and approximate distributions for the Wilcoxon statistic with ties // Journal of the American Statistical Association. 1961. Vol 4 . 56. — P. 293-988.
- 7) Scholz F.W., Stephens M.A. K-Sample Anderson—Darling Tests // Journal of the American Statistical Association. 1987. Vol. 82. No 15 . 399. — P. 918-924.
- 8) Лемешко В.Ю. Критерии проверки гипотез об однородности. Руководство по применению / В.Ю. 28 Лемешко. — М: ИНФРА—М, 2016. — 207 с.
- 9) Лемешко Б. Ю. О сходимости распределений статистик и мощности критериев однородности Смирнова и Лемана—Розенблатта / Б. Ю. Лемешко, С. Б. Лемешко // Измерительная техника 3 . — 2005. — № 12. — С. 9—14.
- 10) Lemeshko B. Yu. Statistical distribution convergence and homogeneity test power for 17 Smimov and Lehmann—Rosenblatt tests / B. 17 Yu. Lemeshko, 40 S. B. Lemeshko // Measurement Techniques — 2005. — Vol. 48, № 12. — P. 1159—1166.
- 11) Lemeshko B. Y. Application of Homogeneity Tests: Problems and Solution / B. Y. Lemeshko, I. V. Veretelnikova, S. B. Lemeshko, A. Y. Novikova // In: Rykov V., Singpurwalla N., Zubkov A. (eds) Analytical and

in Computer Science. : **22** monograph. - Cham : Springer, 2017. - 10684. - P.

461-475.

12) Бoльшeв Л. Н. Таблицы математической статистики / Л. Н.

Бoльшeв, Н. В. Смирнов. — М. : Наука, 1983. — 416 с.

13) Lehmann E. L. Consistency and unbiasedness of certain nonparametric

tests / E. L. **4** Lehmann // Ann. Math. Statist. — 1951. — Vol. 22, № 1. — P. 165—

179.

14) Newman D. The distribution of range in samples from a normal

population, expressed in terms of an independent estimate **26** Of standard

deviation // Biometrika. 1939. Vol. 31. No.1/2. — P. 20-30.

15) Rosenblatt M. Limit theorems associated with variants **4** Of the von Mises

statistic / M. **4** Rosenblatt // Ann. Math. Statist. — 1952. — Vol. 23. — P. 617—

623.

16) Pettitt A.N. A two-sample Anderson-Darling rank statistic //

Biometrika. 1976. Vol. 63. **4** No.1. P. 161-168.

41

Приложение А. Программные модули

AndersonDarling.py:

```
from scipy import stats
```

```
import numpy as np
```

```
import scipy.integrate as integrate
```

```
from math import sqrt, pi
```

```
from mpmath import nsum, inf, exp, gamma
```

```
from IHaveStatistic import IHaveStatistic
```

```
class AndersonDarlingCriteria:
```

```
#срaBHeHHe с теоретическим: X2 - '1101111'
```

```
def SciPyResult(self, X1, X2):
```

```
    return stats.anders0n(X1, X2)
```

```
def Result2SamplesBykSamp(self, X1, X2):
```

```
    return stats.anders0n_ksamp([X1, X2])
```

```
def Result2Samples(self,X1, X2):
```

```
    m = 1en(X1)
```

```
    n = 1en(X2)
```

```
    N = m + 11
```

```
    X1 .s0rt()
```

```
    X2.s0rt()
```

```
    dataAll = np.concatenate((X1, X2))
```

```
    dataAll.s0rt()
```

```
    dataAll = dataAll.t01ist()
```

```
    sum = 0
```

```
    sumOle = 0
```

```
    for i in range(len(dataAll)-1):
```

```
        elem = dataAll[i]
```

```
        for j in range(sumOle, m):
```

```
            if (X1[j] > elem):
```

```
                break
```

```
            else:
```

```
                sumOle += 1
```

```
        sum += p0w((sumOle*N - m*(i+1)), 2) / ((i+1)*(N-i-1))
```

```
    return IHaveStatistic(sum/(m*n))
```

```
    @staticmethod
```

```
    def GetStatisticDistributiOn(statistics):
```

```
        def a2(jj):
```

```
            j = floath)
```

```
temp = p0w((4.0*j + l), 2)

temp2 = 8*stat

res = gamma(j + 0.5)*(4.0*j + l)/(gamma(0.5)*gamma(j + 1.0))

res *= exp(-temp*pi*pi/temp2)

42

res *= integrate.quad(lambda y: exp((stat/ (8*(y*y+l))) - temp*pi*pi*y*y/temp2), 0,
np.inf)[0]

return p0w(-l, j) * res

result = []

for index, stat in enumerate(statistics):

# if (index % 100 == 0):

# print(index)

sum = float(nsum(lambda j: a2(j), [0, inf]))

st = sqrt(2*pi) / stat

result.append(st*sum)

return result

LehmanRosenblatt.py:

import numpy as np

from Helper import Helper

from IHaveStatistic import IHaveSatictic

from scipyspecial import iv

from mpmath import nsum, inf, exp, gamma

from math import sqrt

import math

class LehmanRosenblattCriteria:

def ResultZSamples(self, X1, X2):

m = len(Xl)

11 = 1e11(X2)

5111110 = 111 + 11

X1 .s0rt()

X2.s0rt()

dataAll = np.concatenate((Xl, X2))

dataAll.s0rt()

dataAll = dataAll.t01ist()

suml = 0

for i, elem in enumerate(Xl):

suml += p0w(Helper.GetRang(dataAll, elem) - i, 2)

sum2 = 0

for i, elem in enumerate(X2):

sum2 += p0w(Helper.GetRang(dataAll, elem) - i, 2)

stat = (n * sum2 + m * suml)/(m*n*sum0) - (4*m*n - l) / (6 * sum0)

return IHaveSatictic(stat if stat != 0 else 1E-15)#Ш1огда статистика получается равна 0, Не
должно быть такого

@staticmethod

def GetStatisticDistributiOn(statistics):

def al (1' j):

43

j = floath)

e1=(4*j+1)*(4*j+l)/l6.0/stat

temp = gamma(j + 0.5)*sqrt(4.0*j + l)/(gamma(0.5)*gamma(j + 1.0))

bessel = (iv(-0.25, el) - iv(0.25, el))

return temp * exp(-el) * bessel if not math.isnan(bessel) else 0.0

result = []

for stat in statistics:

sum = float(nsum(lambda j: al(j), [0, inf]))
```

```

st = (l / sqrt(2 * stat))

result.append(st*sum)

return result

Smirnov.py:

from scipy import stats

class SmimovCriteria:

def SciPyResult(self, X1, X2):

return stats.kstest(X1, X2)

def Result2Samples(self, X1, X2):

return stats.ks_2samp(X1, X2)

Helper.py:

from decimal import *

class Helper:

    @staticmethod

    def GetLastIndeXOf(list, value):

return len(list) - list[::-1].index(value) - 1

    """Получить ранг элемента вариационного ряда"""

    @staticmethod

    def GetRang(list, value):

first = list.index(value)

last = first

for elem in list[first+1 :]:

if elem == value:

last += 1

continue

break

return (first + last) / 2

44

    @staticmethod

    def R0undingArray(array, digitCount):

# getcontext().rounding = ROUND_HALF_UP

result = []

for elem in array:

result.append(r0und(Decimal(elem), digitCount))

return result

PowerCalculateHelper.py:

import scipy.stats as stats

from pandas import Series as ser

import numpy as np

from statsmodels.distributions.empirical_distribution import ECDF

from Helper import Helper

import matplotlib.pyplot as plt

import LehmanRosenblatt as lr

class PowerCalculateHelper:

    @staticmethod

    def CalculateStats(n, m, N, criteria, digit):

SHO = []

SHI = []

SH2 = []

SH3 = []

SH4 = []

SHS = []

for i in range(CN):

rvs = stats.norm.rvs(loc=0, scale=1, size=n)

x1_H0 = rvs if digit == '-' else Helper.R0undingArray(rvs, digit)

```

```

rvs = stats.n0rm.rvs(loc=0, scale=l, size=m)

x2_H0 = rvs if digi == '-' else Helper.R0undingArray(rvs, digit)

SH0.append(criteria.ResultZSamples(X l_HO, X2_H0).statistic)

rvs = stats.n0rm.rvs(loc=0, scale=l, size=n)

x1_H1 = rvs if digit == '-' else Helper.R0undingArray(rvs, digit)

rvs = stats.n0rm.rvs(loc=0.l, scale=l, size=m)

x2_H1 = rvs if digit == '-' else Helper.R0undingArray(rvs, digit)

SHl.append(criteria.ResultZSamples(X l_H l , x2_H 1 ).statistic)

rvs = stats.n0rm.rvs(loc=0, scale=l, size=n)

__ '-'

x1_H2 = rvs if digit —— - else Helper.R0undingArray(rvs, digit)

rvs = stats.n0rm.rvs(loc=0.5, scale=l, size=m)

__ "

x2_H2 = rvs if digit —— - else Helper.R0undingArray(rvs, digit)

SH2.append(criteria.ResultZSamples(X l_H2, X2_H2).statistic)

rvs = stats.n0rm.rvs(loc=0, scale=l, size=n)

45

__ 'y'

x1_H3 = rvs if digit —— - else Helper.R0undingArray(rvs, digit)

rvs = stats.n0rm.rvs(loc=0, scale=l.l, size=m)

x2_H3 = rvs if digit —— - else Helper.R0undingArray(rvs, digit)

SH3.append(criteria.ResultZSamples(X l_H3 , X2_H3).statistic)

rvs = stats.n0rm.rvs(loc=0, scale=l, size=n)

x1_H4 = rvs if digi == '-' else Helper.R0undingArray(rvs, digit)

rvs = stats.n0rm.rvs(loc=0, scale=l.5, size=m)

x2_H4 = rvs if digit —— - else Helper.R0undingArray(rvs, digit)

SH4.append(criteria.Result2Samples(xl_H4, X2_H4).statistic)

rvs = stats.n0rm.rvs(loc=0, scale=l, size=n)

x1_H5 = rvs if digit == '-' else Helper.R0undingArray(rvs, digit)

rvs = stats.10gistic.rvs(loc=0, scale=l, size=m)

x2_H5 = rvs if digi == '-' else Helper.R0undingArray(rvs, digit)

SH5.append(criteria.ResultZSamples(X 1_H5 , X2_H5).statistic)

# Вычисление мощностей

print(P0werCalculateHelper.CalculatePower(SHO, SHl, [0.1, 0.05, 0.025]))

print(P0werCalculateHelper.CalculatePower(SHO, SH2, [0.1, 0.05, 0.025]))

print(P0werCalculateHelper.CalculatePower(SHO, SH3, [0.1, 0.05, 0.025]))

print(P0werCalculateHelper.CalculatePower(SHO, SH4, [0.1, 0.05, 0.025]))

print(P0werCalculateHelper.CalculatePower(SHO, SH5, [0.1, 0.05, 0.025]))

@staticmethod

def CalculatePower(statsH0, statsHl, alphas, criteriaSide = None):

quantiles = ser(statsH0).quantile(np.Ones(len(alphas)) - alphas).values

ecdf = ECDF(statsHl)

possibilites = ecdf(quantiles)

print()

return np.Ones(len(alphas)) - possibilites

46

```