

Рекомендации по применению тестов эквивалентности

Роберт А. Крибби

Йоркский университет

Джейми А. Грuman

Виндзорский университет

Шанталь А. Арпен-Крибби

Йоркский университет

Исследователи в области психологии часто выбирают традиционные тесты значимости нулевой гипотезы (например, t-критерий Стьюдента), независимо от того, относятся ли исследуемые гипотезы к тому, эквивалентны ли средние значения групп или различаются ли средние значения групп. Тесты эквивалентности, которые были популярны в биофармацевтических исследованиях в течение многих лет, недавно были введены и рекомендованы исследователям в области психологии для демонстрации эквивалентности двух групп средних. Однако существует очень мало рекомендаций по применению тестов эквивалентности. Исследование методом Монте-Карло использовалось для сравнения критерия эквивалентности, предложенного Цуирманом, с традиционным t-критерием Стьюдента для определения эквивалентности двух групповых средних. Было обнаружено, что критерий эквивалентности Цуирмана более эффективен, чем t-критерий Стьюдента, для определения эквивалентности средних значений генеральной совокупности при больших размерах выборки; однако критерий эквивалентности Цуирмана плохо работает по сравнению с t-критерием Стьюдента при небольших размерах выборки и/или завышенных дисперсиях. © 2003 Wiley Periodicals, Inc. *J Clin Psychol* 60: 1–10, 2004.

Ключевые слова: проверка нулевой гипотезы; тесты эквивалентности; дисперсионная неоднородность; Студенческая t

Исследователи в области психологии, которые заинтересованы в сравнении средних значений двух групп по какому-либо показателю, часто выбирают традиционные тесты значимости нулевой гипотезы (например, t-критерий Стьюдента), в которых нулевая гипотеза относится к эквивалентности населения.

Корреспонденцию по поводу этой статьи следует направлять по адресу: Robert A. Cribbie, факультет психологии Йоркского университета, 4700 Keele Street, Toronto, ON M3J 1P3, Canada; электронная почта: cribbie@yorku.ca.

означает. Кроме того, традиционные тесты нулевой гипотезы надежно применялись независимо от относительности исходной гипотезы к тому, отличаются ли групповые средние значения или эквивалент. Например, клинический исходной гипотезы может быть интересовав оценке исходной гипотезы. гипотеза о том, что анорексика, страдающие анорексией, страдающие передачей информации, и анорексика, ограничивающие себя, имеют одинаковые уровни устойчивости к лечению. Эта исходная гипотеза явно отличается от гипотезы о том, что анорексика, страдающие анорексией от передаления информации, и анорексика, ограничивающие анорексией, имеют разные уровни устойчивости к лечению и, как будет показано ниже, могут иметь важные последствия для статистических данных. принята процедура.

Андерсон и Хук (1983), Руанет (1996), Ширманн (1987), Селвин и Холл (1984), Westlake (1976) и другие предложили статистические методы для определения того, две группы эквивалентны по конкретной зависимой мере, где исходной гипотезы априори определяет минимальное различие, приемлемое для обобщения групповых средних эквивалентными. Эти методы недавно были представлены исходной гипотезы в области психологии через влияющие статьи Роджерс, Ховарда и Весса (1993) и Симона и Серлина (1998). Обе статьи посвящены критерию эквивалентности, предложенному Ширманном (хотя, как Seaman & Serlin отмечают, Rogers et al. отдайте должное Westlake, 1976 за этот метод). Тест эквивалентности Ширмана был чрезвычайно популярен в биофармацевтике. исходной гипотезы для демонстрации биоеквивалентности, хотя до недавнего времени тесты эквивалентности редко принимались исходной гипотезы в области психологии, хотя исходной гипотезы, касающиеся эквивалентности, часто исходуются (Rogers et al., 1993).

Есть по крайней мере две основные причины рекомендовать тесты эквивалентности.

Во-первых, цель исходной гипотезы может заключаться не в том, чтобы показать, что методы лечения идентичны, а только в том, чтобы показать, что различия между методами лечения слишком малы, чтобы их можно было считать значимыми. Расмотрим, например, клинический опросник, заинтересованного в изучении двух конкурирующих методов лечения для депрессии - один для долгосрочной терапии и один для краткосрочной терапии. Исходной гипотезы может быть интересно продемонстрировать, что результаты лечения при краткосрочной терапии эквивалентны результатам долгосрочной терапии, в дополнение к тому, что она занимает меньше времени, дешевле и т. д. В этом примере исходной гипотезы может не понадобиться показывать что методы лечения «абсолютно эквивалентны» (как и в случае с традиционной нулевой гипотезой, Но: $m_1 = m_2$! но только то, что различия в результатах лечения недосаточно велики, чтобы гарантировать принятие более трудоемкой и дорогой терапии (т. е. $m_1 = m_2 \neq D$, где D представляет собой априорную критическую разницу для определения эквивалентности). Во-вторых, известно, что по мере увеличения размера выборки вероятность нахождения даже мельчайших (и потенциально бессмысленных) средних различий статистически значимыми приближается к единице с традиционными тестами нулевой гипотезы эквивалентности, такими как две независимые выборки Стьюдента тест. Это особенно важно, учитывая возросшее количество запросов, например, от учебников, редакторов журналов и статистиков, чтобы исходной гипотезы обосновали значимость своих результатов, включая такие меры величины эффекта, как d , h_2 и т. д. (см. Боуг и Томпсон, 2001; Томпсон, 2002а, 2002б). Поэтому при большой выборке Исходной гипотезы размеров могут захотеть указать критическую разницу между методами лечения, которые могли бы считать клинически значимыми.

Роджерс и др. (1993) продемонстрировали, как результаты исследований эквивалентности двух экспериментальных групп с использованием методов проверки нулевой гипотезы неэквивалентности часто противоречат результатам, полученным с помощью традиционной нулевой гипотезы эквивалентности тесты. Например, Роджерс и др. сравнивали испытуемых, пристрастившихся к алкоголю, и испытуемых пристрастие к наркотикам по шкале коррекции MMPI (на основе исследования Кэннона, Белла, Fowler, Penk, & Finkelstein, 1990), используя как критерий эквивалентности, предложенный Ширманом, так и критерий Стьюдента. Авторы обнаружили, что две группы средних были обобщены статистически разные по критерию Стьюдента, но статистически эквивалентные по критерию Ширмана

проверка эквивалентности. Авторы также обнаружили, что по шкале шизофрении MMPI две группы оказались статистически неразличимыми по t-критерию Стьюдента, но не эквивалентными по критерию эквивалентности Цуирмана.

Учитывая недавние рекомендации в поддержку проверки эквивалентности и увеличение наличие тестов эквивалентности в психологических исследованиях, важно, чтобы исследователи имели четкие руководящие принципы для применения тестов эквивалентности, а также когда тесты эквивалентности могут быть неуместными. Например, на иллюстрации, представленной ранее, нет смысла узнавать, эквивалентны ли баллы по шкале MMPI у субъектов, пристрастившихся к алкоголю, и у субъектов, зависимых от наркотиков, и, таким образом, нет никакого смысла узнавать, являются ли результаты критерия эквивалентности Цуирмана верными или если результаты традиционного теста Стьюдента тест верный. Поэтому целью настоящего исследования является сравнение рекомендуемых в настоящее время методов оценки эквивалентности двух групповых средних, предложенного Цуирманом (1987), с традиционным методом t-критерия Стьюдента. Следующее обсуждение будет (а) рассмотреть критерий эквивалентности Цуирмана, (б) обсудить применение критерия Цуирмана тест эквивалентности в примере, представленном Симаном и Серлином (1998, с. 405), и (с) использовать исследование Монте-Карло для сравнения метода Цуирмана с t-критерием Стьюдента. Метод оценки средних эквивалентности генеральной совокупности.

Критерий эквивалентности Цуирмана

Первым шагом в проведении теста эквивалентности Цуирмана является установление критического среднего значения. Разница для объединения двух популяций означает эквивалентность $\sim D!$. Любая средняя разница меньше чем D считалась бы бессмысленным в рамках эксперимента. Выбор интервала эквивалентности $\sim D!$ является важным аспектом проверки эквивалентности, который в первую очередь зависит от субъективного «уровня уверенности», с которым можно заявить о двух (или более) эквивалентности. Этот уровень достоверности может принимать множество различных форм, включая необработанное значение (например, средние результаты теста, отличные от десяти баллов), процентную разницу (например, 610%), процент разницы объединенного стандартного отклонения и так далее. Трион (2001) описал этот уровень достоверности как «сумму, которая считается несущественной» (с. 379), и Роджерс и др. (1993) заявили, что «любая разница достаточно малая, чтобы попасть в этот интервал эквивалентности, будет считаться клинически и/или практически неважной» (с. 553). Рекомендуются исследователям, обсуждающим подход к значению D, учитывать природу исследования. Например, если длительная терапия, о которой говорилось ранее, занимала в три раза больше времени, и была в три раза дороже, чем краткосрочная терапия, то разница более существенная в их отношении (например, 20%) может потребоваться, чтобы сделать вывод, что методы лечения эквивалентны, если бы длительная терапия занимала в полтора раза больше времени и была примерно эквивалентна по стоимости к краткосрочной терапии (где увеличение расходов на 5% может быть целесообразным для сделать вывод, что методы лечения эквивалентны).

Предполагается, что две выборки случайным образом и независимо выбираются из нормально распределенных совокупностей с одинаковой дисперсией. Можно провести две односторонние проверки гипотез. Используются для установления эквивалентности, где нулевая гипотеза относится к неэквивалентности населения означает и может быть выражено в виде двух отдельных составных гипотез:

$$H_01: m_1 \leq m_2 \text{ и } H_02: m_1 \geq m_2 \quad D^+$$

Отказ от H_01 означает, что $m_1 > m_2$, а отказ от H_02 означает, что $m_1 < m_2$

Далее, отказ от обеих гипотез означает, что $m_1 \approx m_2$ попадает в границы из $\sim D$, и с редства считаются эквивалентными.

Но t_1 отклоняется, если $t_1 \geq t_{\text{кр}}$:

$$t_1 = \frac{\bar{X}_1 - \bar{X}_2 - D}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}, \quad \text{где } D = \frac{n_1 n_2 (s_1^2 - s_2^2)}{n_1 n_2 - 1}$$

и t_2 отвергается, если $t_2 \geq t_{\text{кр}}$:

$$t_2 = \frac{\bar{X}_1 - \bar{X}_2 - D}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}, \quad \text{где } D = \frac{n_1 n_2 (s_1^2 - s_2^2)}{n_1 n_2 - 1}$$

\bar{X}_1 и \bar{X}_2 — групповые средние значения, n_1 и n_2 — объемы групповой выборки, s_1^2 и s_2^2 — стандартные отклонения группы X_1 и X_2 — критическое значение t верхнего уровня с $n_1 + n_2 - 2$ степенями свободы.

Морьяк и Серлин (1998) Пример

Критерий эквивалентности Цуирмана был продемонстрирован Симаном и Серлином (1998, с. 405).

В следующем примере. Критическая разница D для обаяния населения

средний эквивалент был установлен на уровне 5, номинальный уровень значимости был установлен на уровне 0,05 и:

$$X_{P1} = 65+7, n_1 = 25, n_1 = 9$$

$$X_{PC} = 65+0, n_2 = 25, n_2 = 8+$$

Подстановка выборочной статистики генерирует следующую тестовую статистику:

$$t_1 = \frac{-65+7 - 65+0 - 5}{\sqrt{2+4}} = 1+8$$

$$t_2 = \frac{-65+7 - 65+0 - 5}{\sqrt{2+4}} = 2+4+$$

Так как $t_1 \sim 1+8$ и $t_2 \sim 2+4$ и $t_{\text{кр}} = 1.68$ с редства населения обаяются эквивалентны (т. е. разница между средними не превышает критической разницы).

Симан и Серлин также объясняют, что для этого примера такое же решение было бы было достигнуто, если применялся традиционный t -критерий Стьюдента с нулевой эквивалентностью гипотеза. В частности, нулевая гипотеза $H_0: \mu_1 - \mu_2 = 0$ не была бы отвергнута в этом примере, и с редства населения были бы обаяны статистически эквивалентными.

Интересный вопрос, вытекающий из предыдущего примера, заключается в том, как часто два метода (критерий эквивалентности Цуирмана и t -критерий Стьюдента) приводят к одним и тем же выводам (или, что более важно, к тому, как часто методы приводят к разным выводам).

Этот вопрос имеет первостепенное значение, учитывая, что исследователи в области психологии регулярно применяют t -критерий Стьюдента, даже если исследователю кажется эквивалентность означает. Имитационное исследование использовалось для определения вероятности обаяния средних значений населения эквивалентными критерию эквивалентности Цуирмана (т. е. не отвергается) и Критерий Стьюдента (т. е. не отвергается), где размеры выборки и статистика из предыдущего примера использовались в качестве параметров совокупности в моделировании. Это ожидается, что критерий эквивалентности Цуирмана обаяет, что с редства населения эквивалентны более часто, чем критерий Стьюдента, учитывая, что альтернативная гипотеза критерия Цуирмана

$H_a: \mu_1 - \mu_2 \neq 0$ охватывает большую область, чем нулевая гипотеза для Стьюдента

Таблица 1

Имитационное исследование вероятности обнаружения эквивалентности с использованием параметров (m1 65, m2 65/7, s1 8, s2 9, D 5) Из Seaman and Serlin (1997, с. 405)

Статистический тест	Размер выборки							
	n1 n2 25a		n1 n2 50		n1 n2 75		n1 n2 100	
Критерий эквивалентности Цуирмана	.311		0,768		0,913		0,974	
t - критерий Стьюдента	0,939		0,934		0,914		0,912	

a Фактические размеры выборки, использованные в примере Симона и Серлина.

t test –Ho: m1 m2 0), а также потому, что различия между средними значениями (65,7 65,0 0,7) являются односторонними в пределах критической разницы D 5. Было выполнено пять тысяч симуляций. проводится с использованием номинального уровня значимости 0,05. Результаты моделирования представлены в таблице 1. Вопреки логике, лежащей в основе двух методов, результаты показывают, что средние значения, дисперсия и размеры выборки Симона и Серлина. Например, вероятность обнаружить редства эквивалентности Стьюдента (0,939) была значительно больше, чем с критерием эквивалентности Цуирмана (0,311). Далее, даже когда размеры выборки были увеличены до n 50 и n 75, критерий Стьюдента был по крайней мере так же вероятен, как критерий эквивалентности Цуирмана, чтобы обнаружить редства эквивалентными. Нет, это не так пока размеры выборки не были увеличены до n 100, критерий эквивалентности Цуирмана стал более вероятным, чем t - критерий Стьюдента, чтобы обнаружить редства эквивалентными. Однако эти результаты основаны на конкретных параметрах населения и могут не отражать общую эффективность процедур. Поэтому требуется более подробное сравнение критерия эквивалентности Цуирмана и t - критерия Стьюдента.

Исследование Монте-Карло

Имитационное исследование использовалось для сравнения критерия эквивалентности Цуирмана с традиционным t - критерием Стьюдента для обнаружения эквивалентности населения в условиях, обычно опыт исследователей в области психологии. При этом манипулировались несколькими переменными. исследование, включая (а) размер выборки, (б) среднюю конфигурацию населения и (в) население отклонения. Критическая средняя разница для установления эквивалентности популяции с помощью критерия эквивалентности Цуирмана поддерживалась на уровне 1 во всех условиях.

Одной из основных причин использования тестов эквивалентности является то, что в качестве выборки размер увеличивается, вероятность того, что даже тривиальные средние различия будут статистически значимыми, становится большой. Таким образом, размеры выборки и группления были изменены в этом исследовании. исследование. В частности, размеры групповой выборки были установлены на уровне n1 n2 10, n1 n2 25, n1 n2 50 и n1 n2 100.

Эффективность проверки эквивалентности напрямую зависит от неоднородности значительных групп. В этом исследовании оценивались пять средних конфигураций, включая эквивалентные средние значения совокупности ~ m1 m2 ! а четыре неэквивалентные совокупности означают ~ m1 m2 +4, m1 m2 +8, m1 m2 1+2 и m1 m2 1,6). Учитывая, что критическая разница если эквивалентность совокупности установлена равной 1, конфигурация эквивалентных средних и неэквивалентных конфигураций с m2 m1 1 предполагают подальтернативную гипотезу критерия эквивалентности Цуирмана (т. е. разность средних значений совокупности не превышает критическую среднюю разницу, и, таким образом, ожидается, что средние значения будут объявлены эквивалентными

критерий Цуирмана), а ненулевые конфигурации с $m_2 \neq m_1$ попадают под нулевой гипотез критерия эквивалентности Цуирмана (т.е. разность средних значений генеральной совокупности превышает разность критических средних, и, таким образом, ожидается, что средние значения будут объявлены неэквивалентными критерию Цуирмана). Для t -критерия Стьюдента разницы средних значений в любой популяции больше чем ноль подпадает под альтернативную гипотезу $\mu_1 \neq \mu_2$ и предполагается, что средства объявлен неэквивалентным.

Другой важный вопрос заключается в том, какой эффект будет иметь увеличение изменчивости населения? Имеют по критерию эквивалентности Цуирмана. Также в отношении дисперсии населения, в обзор опубликованных исследований в области образования и психологии, Keselman et al. (1998) найдено что неравные отклонения были нормой, а не исключением. В частности, исследователи часто сообщают об отношении наибольшей дисперсии к наименьшей величине, равной четырем, и наибольшей к наименьшей коэффициент дисперсии, достигающие восьми, не были редкостью. Таким образом, в дополнение к исследованию случая, когда дисперсия обеих групп была установлена равной единице, влияние популяции инфляцией дисперсии также была исследована в этом исследовании путем установки дисперсии одной из группы до четырех-восьми.

Для каждого условия было проведено пять тысяч симуляций с использованием номинального уровня значимости 0,05.

Полученные результаты

Априорная эквивалентность

В табл.

априорные средние различия в популяции были меньше, чем критическая средняя разница, выборка размер был основным фактором при сравнении критерия эквивалентности Цуирмана и критерия Стьюдента. критерий (Здесь следует отметить, что для t -критерия Стьюдента ложная популяционная разность средних

Таблица 2
Вероятность обнаружения популяционной эквивалентности для критерия эквивалентности Цуирмана (критическая разница для эквивалентности 1) и критерия Стьюдента

n	m1 m2	2 c 1 2 c 2 1		2 c 1 0 c 2 4		2 c 1 0 c 2 8	
		C	T	C	T	C	T
10	0	.3930	.9488	.0248	.9466	.0026	.9412
	0,4	.2764	.8576	.0220	.9076	.0016	.9234
	0,8	.0966	.5962	.0114	.7940	.0008	0,8582
25	0	.9384	.9534	.4426	.9506	.0698	0,9488
	0,4	.6778	.7174	.3016	.8314	.0526	0,8922
	0,8	.1698	.2036	.1056	.5832	.0308	.7528
50	0	.9996	.9510	.8660	.9494	.5230	.9514
	0,4	.9100	.4962	.5922	.7676	.3512	.8530
	0,8	.2562	.0202	.1544	.2974	.1174	.5320
100	0	1.000	.9496	.9946	.9510	.9083	0,9498
	0,4	.9956	.1922	.8444	.5688	.6324	.7330
	0,8	.3998	.0000	.2388	.0544	.1714	.2526

Критерий эквивалентности С. Цуирмана; t -критерий Стьюдента.

делает нулевую гипотезу ложной, и поэтому цель проверки при всех ненулевых условиях состоит в том, чтобы выявить неэквивалентность, а не эквивалентность средних.

Доля неотвержений нулевой гипотезы для t-Стюдента с средними различиями больше нуля [т.е. ошибки типа II] представлены для сравнения только с процедурой Цуирмана. При 10 или 25 наблюдениях в группе критерий Стюдента был более вероятным. Объединить две группы означает эквивалентность, чем критерий эквивалентности Цуирмана, независимо от размера разности средних значений генеральной совокупности. Фактически, с 10 субъектами в группе тест эквивалентности Цуирмана никогда не обнаруживал эквивалентности более чем в 50% случаев, тогда как тест Стюдента никогда не обнаруживал эквивалентности менее чем в 50% случаев.

С другой стороны, при 50 или 100 субъектах в группе критерий эквивалентности Цуирмана с большей вероятностью, чем критерий Стюдента, выявлял эквивалентность.

Когда дисперсии генеральной совокупности были завышены в одной группе, критерий Цуирмана очень плохо определял эквивалентность средних значений генеральной совокупности. В частности, когда дисперсия одной группы была установлена равной восьми, критерий эквивалентности Цуирмана никогда не был более эффективным для обнаружения эквивалентности, чем t-критерий Стюдента, независимо от размера выборки или конфигурации средних совокупности. Когда дисперсия одной группы была установлена равной четырем, критерий эквивалентности Цуирмана был более эффективным, чем t-критерий Стюдента, только когда в группе было не менее 100 субъектов. Когда средние группы эквивалентны примерно в 95% случаев, независимо от размера выборки. С другой стороны, критерий эквивалентности Цуирмана часто обнаруживал, что средние значения группы неэквивалентны, даже когда средние значения совокупности были идентичными. Например, при 10 субъектах в группе, равных средних значениях популяции и дисперсии одной группы, равной четырем, тест эквивалентности Цуирмана показал, что средние значения эквивалентны только в 2,48% случаев, тогда как критерий Стюдента обнаружил, что средние значения эквивалентны в 2,48% случаев. 94,66% случаев. Увеличение размера выборки улучшило производительность теста Цуирмана, хотя тест по-прежнему работал плохо по сравнению с t-Стюдента. Например, при 50 субъектах в группе, равных средних значениях популяции и дисперсии одной группы, равной четырем, тест эквивалентности Цуирмана показал, что средние значения эквивалентны в 86,60% случаев, тогда как t-Стюдента обнаружил, что средние значения эквивалентны в 94,84. % случаев.

Априорная неэквивалентность

Вероятности обнаружения неэквивалентности популяции с помощью критерия эквивалентности Цуирмана и t-критерия Стюдента для смоделированных условий представлены в табл. 3. (все случаи) при обнаружении различий.

Однако ясно, что превосходная способность критерия Цуирмана обнаруживать средние различия, превышающие критическую разницу, является функцией систематической ошибки теста для объединения популяций неэквивалентными, даже когда различия были меньше критической разницы. Мощность t-Стюдента для обнаружения средних различий, как и ожидалось, зависит от размера выборки и дисперсии, при этом мощность максимизируется при больших размерах выборки и небольших дисперсиях.

Обсуждение

В настоящей статье исследована альтернатива традиционному t-критерию Стюдента для определения эквивалентности средних значений двух лечебных групп. Есть много примеров (представленных здесь и в других местах) парадигм исследования клинической психологии, в которых интерес представляет вопрос, является ли однолечебное средство практически эквивалентным второму среднему, или в другом

Таблица 3

Вероятность обнаружения неэквивалентности популяции для критерия эквивалентности Цуирмана (критическая разница для эквивалентности 1) и критерия Стьюдента

n	m1 m2	2 c 1 2 c 2 1		2 c 10 c 2 2 4		2 c 10 c 2 2 8	
		C	t	C	t	C	t
10	1,2	.9832	.7254	.9928	.3598	.9982	.2290
	1,6	.9980	.9280	.9974	.5890	.9992	.3780
25	1,2	.9900	.9874	.9834	.7486	.9886	.5006
	1,6	.9998	1.000	.9988	.9414	.9980	.7446
50	1,2	.9942	.9998	.9856	.9584	.9798	0,7928
	1,6	1.000	1.000	.9998	.9990	.9990	0,9588
100	1,2	.9984	1.000	.9924	.9994	.9876	0,9764
	1,6	1.000	1.000	1.000	1.000	1.000	.9998

Критерий эквивалентности С. Цуирмана t - критерий Стьюдента.

Другими словами, разница между двумя лечебными средствами не настолько велика, чтобы считаться осмысленной. В недавних статьях исследователи в области психологии были представлены тесты эквивалентности и рекомендованы эти процедуры для ответов на вопросы, касающиеся эквивалентности двух групповых средних. Эти статьи увеличили как достоверность, так и популярность этих процедур. Однако статистических исследований мало. Практически отсутствуют вопросы процедуры и руководства по применению тестов эквивалентности.

В этой статье были проведены два исследования моделирования: (а) расширение примера, использованного Симаном и Серпином (1998) для демонстрации применения теории Цуирмана критерия эквивалентности и (б) сравнение критерия эквивалентности Цуирмана и t - критерия Стьюдента во многих условиях, обычно встречающихся в экспериментальных поведенческих науках. Цель исследований моделирования состояла в том, чтобы выделить статистические свойства тестов эквивалентности и как они соотносятся с традиционной проверкой нулевой гипотезы. Так должно быть отмечено, что многие из результатов, представленных в этой статье, можно было бы предсказать на основе лежащих в основе выборочных распределений тестовой статистики, хотя важно, чтобы эти результаты должны были количественно оценены, чтобы позволить исследователям принимать обоснованные решения при выборе соответствующую тестовую статистику. Оба исследования показали, что размер выборки является решающим фактором при выборе между критерием эквивалентности Цуирмана и традиционным t - критерием Стьюдента. Если число субъектов на одну группу велико (25 или более), критерий эквивалентности Цуирмана может быть более предпочтительным для обнаружения эквивалентности населения, чем критерий Стьюдента, особенно, когда различия средних популяций присутствуют, но меньше критической разницы. По мере увеличения размера выборки и разницы между средними t критерий Стьюдента как и ожидалось, становится более мощным в обнаружении различий и, следовательно, с меньшей вероятностью объявить различия бесмысленными. С другой стороны, по мере увеличения размера выборки критерий эквивалентности Цуирмана становится более мощным при обнаружении этих различий. меньше, чем критическая разница и поэтому рекомендуется для t Стьюдента с большими размерами выборки.

Однако требуется важная оговорка к предыдущей рекомендации относительно размера выборки; когда групповые дисперсии даже умеренно завышены, с поправкой Процедура Цуирмана для обнаружения эквивалентности существенно окрашена. Хотя это открытие частично является результатом снижения мощности (т.е. увеличения стандартной ошибки)

Стюдента для обнаружения различий средних с завышенными дисперсиями, не следует упускать из виду экстремальный эффект завышения дисперсий на критерий эквивалентности Шуирмана (особенно при небольших размерах выборки). В частности, крайне проблематично то, что при отсрочке различий между средними значениями (и вероятность вывода об эквивалентности групп с t -Стюдента составляет приблизительно 0,95) критерий эквивалентности Шуирмана так плохо работает с малыми размерами выборок и/или увеличение отклонений.

Кроме того, также важно признать, что, хотя это и не является предметом этой статьи, неравные размеры выборки в сочетании с неравными дисперсиями могут значительно повлиять на частоту ошибок типа I и типа II для t -Стюдента, помимо того, что наблюдалось в этом исследовании. Исследование. Ожидается, что влияние неравных размеров выборок и дисперсий также значительно повлияет на частоту ошибок типа I и типа II критерия эквивалентности Шуирмана, и поэтому необходимо исследование надежного теста эквивалентности для неравных размеров выборок и дисперсий.

Хотя результаты этого исследования важны с точки зрения представления рекомендаций исследователя относительно выбора подхода к тестированию статистики для оценки эквивалентности двух групповых средних, существуют два важных ограничения текущего исследования, которые следует учитывать. Во-первых, исследование Монте-Карло не могло изучить все возможные условия размера выборки, дисперсионного неравенства и т. д.; поэтому, хотя мы ожидаем, что результаты этого исследования будут обобщены на многие распространённые ряды тестирования, эти результаты специфичны для условий, изучаемых в этом исследовании. Во-вторых, хотя в этой статье мы сосредоточились на схеме проверки гипотез для оценки эквивалентности средних, были достигнуты важные успехи в применении подходов доверительного интервала к проверке эквивалентности, которые здесь не рассматривались, но могут представлять интерес читателям (см., например, Seaman & Serlin, 1998; Tryon, 2001).

Подводя итог, можно сказать, что тесты на эквивалентность чрезвычайно популярны в биофармацевтических исследованиях для демонстрации того, что эффекты двух лекарств практически эквивалентны. Ожидается, что по мере роста числа исследований, описывающих методологию тестов эквивалентности, популярность тестов эквивалентности в области психологии будет расти, учитывая, что исследователи будут более подготовлены к выявлению ситуаций, в которых тесты эквивалентности уместны. Поэтому важно, чтобы существовали четкие рекомендации по применению этих тестов. Результаты этого исследования подчеркивают необходимость признания того, что критерий эквивалентности Шуирмана и критерий Стюдента диаметрально противоположны в своем подходе к проверке гипотез, и, следовательно, одни и те же факторы, которые существенно влияют на последствия критерия Стюдента (обнаруживать различия между средними значениями (например, размер выборки, вариабельность ошибок) также существенно влияют на последствия теста Шуирмана обнаруживать эквивалентность.

использованная литература

- Андерсон, С.А., и Хук, В.В. (1983). Новая процедура проверки эквивалентности в сравнительных исследованиях бидосупности и других клинических исследованиях. Коммуникация в статистике: теория и методы, 12, 2663–2692.
- Бог, Ф., и Томпсон, Б. (2001). Использование размеров эффекта в исследованиях в области социальных наук: новые требования APA и журналов для улучшения методов методологии. Журнал исследований в области образования, 11, 120–129.
- Кэннон, Д.С., Белл, В.Е., Фаулер, Д.Р., Пенк, В.Е., и Финкельштейн, А.С. (1990). Различия MMPI между алкоголиками и наркоманами: влияние возраста и расы. Психологическая оценка: журнал клинической и консультационной психологии, 2, 51–55.
- Кесельман, Х. Дж., Хьюберти, С. Дж., Лиск, Л. М., Олейник, С., Крибби, Р. А., Донахью, Б., Ковальчук, Р. К., Лоуман, Л. Л., Петоски, М. Д., Кесельман, Дж. К., и Левин, младший (1998). Статистические практики исследователей в области образования: анализ их ANOVA, MANOVA и ANCOVA анализов. Обзор образовательных исследований, 68, 350–386.

- Роджерс, Дж. Л., Ховард, К. И. и Весс, Дж. Т. (1993). Использование критериев значимости для оценки эквивалентности между двумя экспериментальными группами. Психологический бюллетень, 113, 553-565.
- Руане, Х. (1996). Байесовские методы оценки важности эффектов. Психологический вестник, 119, 149-158.
- Цуирманн, Дж. (1987). Сравнение процедуры двух односторонних тестов и степенного подхода для оценки эквивалентности с редней биодоступности. Журнал фармакокинетики и биофармацевтики, 15, 657-680.
- Морьяк, Масанусетс, и Серлин, Р.К. (1998). Доверительные интервалы эквивалентности для двух групповых сравнений с учетом различия. Психологические методы, 3, 403-411.
- Селвин, М.Р., и Холл, Н.Р. (1984). Обобщенных методов биоэквивалентности. Биометрия, 40, 1103-1108 гг.
- Томпсон, Б. (2002а). Какими могут быть будущие количественные исследования в области социальных наук: доверительные интервалы для величины эффекта. Исследователь в области образования, 31, 24-31.
- Томпсон, Б. (2002b). «Статистические», «практические» и «клинические»: сколько видов значимости консультанты должны учитывать? Журнал консультирования и развития, 80, 64-71.
- Трион, В.В. (2001). Оценка статистической разницы, эквивалентности и неопределенности с использованием доверительных интервалов вывода интегрированный альтернативный метод выдвижения нулевой гипотезы статистические тесты. Психологические методы, 6, 371-386.
- Вестлейк, В. (1976). Симметричные доверительные интервалы для исследований биоэквивалентности. Биометрика, 37, 589-594.