

О НОРМАЛЬНОСТИ ВРЕМЕННЫХ РАСПРЕДЕЛЕНИЙ КВАЗИСТАЦИОНАРНЫХ, ЛИНЕЙНЫХ ИНДЕКСОВ В СОЦИОЛОГИИ

ON NORMALITY OF THE TIME DISTRIBUTIONS OF QUASI-STATIONARY, LINEAR INDICES IN SOCIOLOGY

Аннотация

Заметную роль в социологии играют индексы, которые исчисляются по данным опросов, образующих выборки с размерами порядка 1 000 - 10 000, и линейно выражаются через доли положительных и отрицательных ответов. Примерами служат индексы социальных настроений, потребительских настроений и ожиданий безработицы, которые публикуются «Левада-центром», различные индексы социальных оценок от ВЦИОМ. А также многие другие показатели в различных областях знания (экономика, медицина и т.д).

В статье выдвинута и теоретически обоснована гипотеза о том, что такие индексы имеют нормальные законы распределения на тех промежутках времени, где индексы являются квазистационарными. Это условие проявляет себя в том, что сезонная компонента индекса отсутствует, а его тренд является (приблизительно) стационарным. Полученный результат можно использовать для предварительных оценок размеров выборок при опросах, что могло бы снизить расходы на социальный мониторинг.

Summary

In sociology a notable role is played by the indices, which are calculated according to the surveys with samples of the sizes of the orders 1,000 - 10,000, and are linearly expressed through the shares of positive and negative answers. Examples are the indices of social moods, consumer moods and unemployment expectations, which are published by the "Levada Center", social estimates indices from the RPORC. As well as many other indicators in various areas of knowledge (economics, medicine, etc.).

The article has put forward and theoretically substantiated the hypothesis that such indices have normal distribution laws on those time intervals where the trends are quasi-stationary. This condition manifests itself in the fact that the seasonal component of an index is absent, and the trend is (approximately) stationary. The result might be used for preliminary estimates of the sample sizes for surveys, which could reduce the cost of social monitoring.

Ключевые слова: Социологический индекс, индекс социальных настроений, нормальное распределение, теорема Ляпунова, временной ряд.

Keywords: Sociological index, index of social moods, normal distribution, Gaussian random variable, Lyapunov theorem, time series.

JEL: C12, C43, C65.

§ 1. Введение

Социологические показатели (*индексы*) исчисляются на основании опросов или иных исследований членов социальных групп. Будем называть их респондентами независимо от того, какой именно метод сбора индивидуальных данных применяется. В дальнейшем этот метод называется опросом, хотя фактически он может быть другим.

Значения $x(t_1), x(t_2), \dots, x(t_m)$ некоторого социологического индекса X , исчисленные в последовательные моменты времени t_1, t_2, \dots, t_m , образуют т.н. временной ряд. В нем выделяют тренд $y(t_i)$ (линию регрессии), сезонную и циклическую компоненты $s(t_i)$ и $c(t_i)$, а также ошибку вычислений $e(t_i)$ [Крыштановский, 2000]. Пренебрегая этой ошибкой предположим, что сезонная компонента $s(t_i)$ отсутствует или незначительна, а на промежутке времени $[t_1; t_m]$ тренд $y(t)$ меняется так медленно, что его можно считать постоянным. Тогда вариация временного ряда обусловлена циклической компонентой $c(t_i)$. Возникает естественная гипотеза о том, что индекс X является приблизительно нормальной, случайной величиной, хотя бы на протяжении времени от t_1 до t_m . В таком случае можно было бы экстраполировать значения X на ближайшее будущее, используя метод доверительных интервалов, а также определить достаточные размеры предстоящих выборок. Последнее имеет важное значение для экономии ресурсов на опросы. Одновременно возникает возможность строго обосновать репрезентативность выборок.

. Многие социологические индексы линейно выражаются через частоты ответов на вопросы. Для таких индексов можно обосновать нормальный закон распределения на тех промежутках времени, где они являются *квазистационарными* (§2). В сущности, это условие сводится к тому, что за период времени $[t_1; t_m]$ не было социальных изменений, способных существенно повлиять на результаты опросов при исчислении индекса X . При этом опросы должны происходить на случайных выборках, что означает равновероятность попадания в опрос любого представителя изучаемой, социальной группы. Случайность выборок трудно

осуществима, поэтому в социологической практике активно применяются методы неслучайных выборок

[Neyman, 1934:565]

В этой статье теоретически доказано, что квазистационарный, линейный, социологический индекс X должен иметь приблизительно нормальный (т.е. гауссов) закон распределения. Нормальные случайные величины занимают центральное место в теории вероятностей и математической статистике [Гмурман, 2003]. Для них можно по эмпирическим выборкам найти доверительные интервалы и получить оценки статистических показателей, используя такие пакеты прикладных программ, как STATISTICA и SPSS.

Результат настоящей статьи в какой-то мере отвергает точку зрения о том, что нормальные распределения не находят места в социологии. Согласно [Давыдов, 1995]: «Опыт автора по анализу результатов опросов общественного мнения в Институте социологии РАН за 1985—1989 гг. показывает, что нормальное распределение встречается довольно редко». В [Ильясов, 2014] утверждается: «Не удалось найти аргументов, обосновывающих гипотезу о том, почему определенный вид социологического распределения должен (может) подчиняться нормальному распределению». И далее. «Понятно, что какие-то эмпирические распределения могут иметь вид, близкий к нормальному, но это может происходить не в силу закономерности, а вследствие вариативности распределений, т.е. по совпадению. Таким образом, можно предположить, что идея нормального распределения в социологии не обоснована».

По-видимому, сказанное выше относится к распределениям значений показателей среди членов социальных групп. Нормальные распределения, о которых идет речь в этой статье, возникают в связи с временными рядами индексов социологии. В [Ильясов, 2014] сказано следующее. «Временные ряды, как представляется, некорректно соотносить с законами Гаусса или Ципфа. ... Тем не менее проблема соотношения нормального распределения с временными рядами обсуждается». Результат настоящей статьи непосредственно касается этой проблемы.

Представленный результат получен из классической теоремы Ляпунова [Ляпунов, 1948: 245], согласно которой величину X можно считать нормальной, если она складывается из большого числа взаимно независимых, случайных величин, каждая из которых значительно меньше X . Используются методы теории вероятностей [Гмурман, 2003; Вентцель, 1962].

§ 2. Линейные социологические индексы

Пусть проводится опрос в социальной группе, при этом членам выборки размера n задают по Q вопросов, на каждый из которых можно дать один из предлагаемых ответов. Пронумеруем респондентов $i = 1, \dots, n$, вопросы $q = 1, \dots, Q$, ответы $r = 1, \dots, R_q$ (число ответов на вопрос № q в общем зависит от q). Обозначим S_{qr} долю ответов № r на вопрос № q , выраженную в % от общего числа полученных ответов, где $r \in \{1, \dots, R_q\}$. Пусть некоторый индекс X , связанный с этим опросом, выражается формулой:

$$X = A + \sum_{q=1}^Q \sum_{r=1}^{R_q} K_{qr} S_{qr} \quad (1)$$

Постоянные K_{qr} и A определяются при разработке показателя. Такие индексы X мы будем называть *линейными*.

Примером служит индекс социальных настроений (ИСН¹), который ежемесячно публикуется «Левада-Центром». Респондентам из выборки размером $n = 1600$ предлагается по $Q = 12$ вопросов, на каждый из которых нужно дать один из нескольких, заданных ответов. ИСН вычисляется, как разность между долями положительных и отрицательных ответов, выраженными в % от общего числа ответов, к которой прибавляют 100 для исключения отрицательных значений. Таким образом, если X есть ИСН, то $A = 100$, для положительных ответов $K_{qr} = 1$, для отрицательных $K_{qr} = -1$, для нейтральных $K_{qr} = 0$. Аналогично выражаются индекс потребительских настроений (ИПН), индекс ожидания безработицы (ИБ) от «Левада-Центра», индексы социальных оценок от ВЦИОМ (где $A = 0$).

Чтобы принять во внимание разницу между ответами по важности или степени уверенности респондентов, можно присвоить ответам на вопрос № q уровни значимости l_{q1}, \dots, l_{qP_q} , где P_q – число вариантов положительных ответов, которое равно числу вариантов отрицательных ответов, $q = 1, \dots, Q$. Уровни значимости определяются интуитивно – методом экспертного оценивания, обычно по шкале баллов от 1 до 5 или от 1 до 10, хотя возможны и другие шкалы. Для такого *дифференцированного*, социального индекса формула (1) приобретает следующий вид:

¹ Официальный сайт «Левада-Центра» <https://www.levada.ru/obnovlennaya-metodika-izmereniya-indeksa-sotsialnykh-nastroenii-isn/>

$$X = A + \sum_{q=1}^Q \sum_{s=1}^{P_q} w_{qs} \cdot (S_{qs}^+ - S_{qs}^-) \quad w_{qs} = \frac{l_{qs} P_q}{l_{q1} + \dots + l_{qP_q}} \quad (2)$$

где S_{qs}^+ и S_{qs}^- есть доли положительных и отрицательных ответов на вопрос № q , имеющих номер s и общий уровень значимости l_{qs} . Как положительные, так и отрицательные ответы пронумерованы от 1 до P_q , нейтральные ответы не принимаются в расчет ($R_q = 2P_q + 1$). При каждом q весовые коэффициенты w_{qs} пропорциональны l_{qs} , где $\sum_{s=1}^{P_q} w_{qs} = P_q$.

Если $l_{q1} = \dots = l_{qP_q}$ для каждого q , т.е., все ответы на один вопрос равнозначны (кроме нейтральных), то $w_{qs} = 1$ и формула (2) сводится к следующей:

$$X = A + \sum_{q=1}^Q (S_q^+ - S_q^-) = A + \sum_{q=1}^Q S_q^+ - \sum_{q=1}^Q S_q^- \quad (3)$$

где S_q^\pm - доли положительных и отрицательных ответов на вопрос № q , выраженные в % от числа всех ответов: $S_q^+ = S_{q1}^+ + \dots + S_{qP_q}^+$ и $S_q^- = S_{q1}^- + \dots + S_{qP_q}^-$.

Формула (3) отвечает ИСН, ИПН и ИБ от «Левада-Центра», индексам социальных оценок от ВЦИОМ и другим показателям.

Формула (1) может выражать не только социологические индексы. Например, в [Дуганов, 2011] описан показатель уровня преждевременной смертности от болезней (ППЖ). Здесь $Q = 1$ и параметр R_1 на 1 больше т.н. базового возраста, все смерти раньше достижения которого считаются преждевременными, S_{1r} равно доле случаев смерти в возрасте $r - 1$ лет в % от числа N умерших преждевременно, $A = 0$ и $K_{1r} = N \cdot (R_1 - r)/100$.

Другой пример связан с композитным индексом I материального благосостояния, описанным в [Балацкий, Саакянц, 2006]. Он вычисляется по формуле (1) при $Q = 1$ и $R_1 = 5$, где S_{1r} равно доле в % респондентов, относящих себя к $r - й$ группе благосостояния (определенной в таблице 1 [Балацкий, Саакянц, 2006]), $A = 0$, весовые коэффициенты $K_{11} = 0$ $K_{12} = 0,25$ $K_{13} = 0,5$ $K_{14} = 0,75$ $K_{15} = 1$.

Примером социологического показателя, не выражаемого формулой (1), является индекс Джини [Балацкий, Саакянц; 2006], характеризующий материальное неравенство.

Рассмотрим произвольный индекс X , определяемый формулой (1). Пусть он вычисляется в последовательные моменты времени t_1, t_2, \dots, t_m . Для определенности будем считать, что респонденты нумеруются в порядке занесения ответов в базу данных службы

социального мониторинга. Введем случайную величину X_{qr}^i , которая зависит от ответа респондента № i на вопрос № q . Пусть она равна 1, если дан ответ № r , и равна нулю при любом другом ответе. Тогда из (1) следует, что

$$X = A + \sum_{q=1}^Q \sum_{r=1}^{R_q} K_{qr} \cdot \frac{\sum_{i=1}^n X_{qr}^i}{nQ} \cdot 100 = A + 100 \cdot \frac{\sum_{i=1}^n X_i}{n} \quad X_i = \frac{1}{Q} \sum_{q=1}^Q \sum_{r=1}^{R_q} K_{qr} X_{qr}^i \quad (4)$$

где nQ есть общее число полученных ответов.

Умозаключение о (приблизительной) нормальности индекса X основано на том, что все случайные величины X_i/n независимы между собой и $|X_i/n| \ll |X|$, т.к. $n \gg 1$. Поэтому величина $Y = \sum_{i=1}^n X_i/n$ имеет близкое к гауссову распределение [Гмурман, 2003: 135], а линейная функция $X = 100 \cdot Y + A$ от нормальной величины Y также является нормальной [Гмурман, 2003: 141]. Мы должны еще показать, что при всех $i \neq j$ случайные величины $X(t_i)$ и $X(t_j)$ независимы между собой и имеют близкие законы распределения, которые можно считать совпадающими.

Каждая величина X_i представляет собой результат «измерения», т.е., определения в результате опроса респондента № i значения одной и той же, случайной величины

$$\chi = \frac{1}{Q} \sum_{q=1}^Q \sum_{r=1}^{R_q} K_{qr} x_{qr} \quad (5)$$

где $x_{qr} = 1$, если на вопрос № q был дан ответ № r , а при других ответах $x_{qr} = 0$.

Линейный, социологический индекс X называется *квазистационарным* на данном промежутке времени, если закон распределения случайной величины χ является неизменным или меняется за это время незначительно. В таком случае можно считать, что случайные величины $\chi(t_1), \dots, \chi(t_m)$ имеют общий закон распределения или, другими словами, случайная величина χ не зависит от времени (это справедливо лишь в некотором, практически достаточном приближении на определенном промежутке времени).

Тогда, как видно из (4) и (5), индекс X является случайной величиной, которая также не зависит от времени. Предположим, что респонденты нумеруются индексом i , который изменяется от 1 до бесконечности. Ни одна реальная, социальная группа не состоит из бесконечного числа членов. Но если она достаточно велика, то такая идеализация не приведет к существенной ошибке. Рассмотрим последовательность величин $Y_n = \sum_{i=1}^n X_i/n$, где $n = 1, 2, \dots$ и X_i есть результат измерения показателя χ для i -го респондента. Каждая X_i есть случайная величина, поскольку респонденты выбираются случайным образом. Согласно следствию из теоремы [Ляпунов 1948: 245], при $n \rightarrow \infty$ закон распределения величины Y_n неограниченно приближается к нормальному.

Условия теоремы Ляпунова также требуют, чтобы для некоторого $\delta > 0$ математическое ожидание величины $|\chi - M(\chi)|^{2+\delta}$ было конечным [Ляпунов 1948: 245]. В данном случае это верно для всех $\delta > 0$, т.к. из (5) следует, что полученное при любом измерении значение x величины χ удовлетворяет неравенству $|x| < \sum_{q=1}^Q \sum_{r=1}^{R_q} |K_{qr}|/Q$. Следовательно, величина $|\chi - M(\chi)|^{2+\delta}$ является ограниченной при любом $\delta > 0$. Поэтому ее математическое ожидание конечно. Таким образом, условия теоремы Ляпунова выполняются.

Теоретически обоснована гипотеза о том, что любой социологический индекс X , определяемый формулой (1), должен иметь приблизительно нормальные распределения на тех промежутках времени, где его можно считать квазистационарным.

Для проверки последнего условия нужна статистическая выборка значений величины χ (5) в различные моменты времени для достаточно большого числа респондентов. Размер этой выборки может быть сопоставимым с опросом при исчислении индекса X (учитывая, что необходимо получить значения χ в различных регионах страны). Поэтому о выполнении условия квазистационарности следует судить по косвенным признакам, описанным в § 1.

Зависимость индекса X от времени определяет случайный процесс $X(t)$, который для случайных выборок респондентов должен иметь нулевую автокорреляционную функцию. Последнее означает, что при всех $t_i \neq t_j$ коэффициент корреляции случайных величин $X(t_i)$ и $X(t_j)$ равен нулю. Если же в каждый момент t_i респонденты выбираются не случайно (а из одной и той же, узкой группы), то автокорреляционная функция $X(t)$ будет существенно отличной от нуля. Поэтому $x(t_1), \dots, x(t_m)$ нельзя рассматривать, как выборку значений одной и той же, случайной величины X , полученных в результате m независимых измерений.

Если при этом X выражается формулой (1), то функция $X(t) - A$ разлагается в сумму случайных процессов $100 \cdot K_{qr} X_{qr}^i(t)/(nQ)$, каждый из которых связан с *одним и тем же* респондентом № i (4). Вывод об эргодичности случайного процесса делают только при уверенности в том, что он не разлагается в сумму нескольких процессов [Вентцель, 1962: 451]. Поэтому процесс $X(t)$ нельзя считать эргодическим. Так мы снова приходим к заключению о последовательности $x(t_1), \dots, x(t_m)$, которое было сформулировано выше.

Таким образом, если при исчислении социологического индекса периодически опрашиваются одни и те же респонденты, то постановка вопроса о нормальности временного распределения этого индекса не имеет смысла.

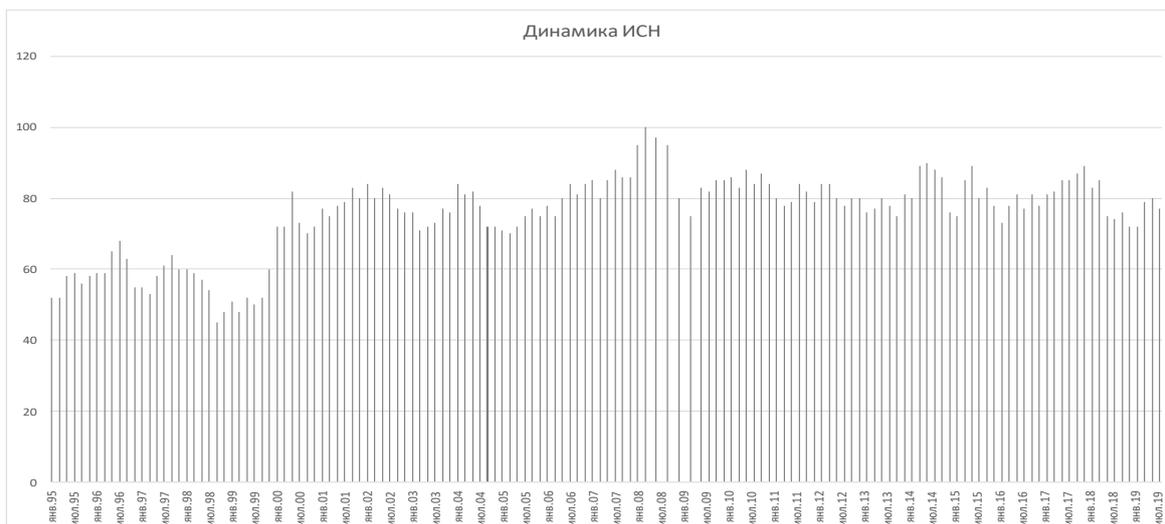


Рис. 1. Значения ИСН в % от показателя марта 2008 с января 1995 по август 2019.

§ 4. Индекс социальных настроений

Используя данные, опубликованные на сайте «Левада-Центра»², проверим гипотезу о нормальности индекса социальных настроений (ИСН) на промежутке времени с марта 2009 по август 2019. Абсолютные значения на сайте не указаны, но даны индексы семьи, России, ожиданий и власти за каждые 2 месяца. Значения ИСН вычислены, как среднеарифметические этих 4-х индексов, округленные до целых. Получен ряд чисел от 101 до 128, за исключением 103, 121 и 123. Сезонная компонента на рис. 1 не усматривается (локальные максимумы и минимумы приходятся как на летние, так и на зимние месяцы). Гистограмма производит впечатление, что с марта 2009 по август 2019 ИСН был квазистационарным (§ 1). Это соответствует официальной точке зрения о том, что за первое десятилетие XXI века Россия достигла стабильного, социально-экономического состояния. Однако, квадратичный тренд данного временного ряда, построенный в пакете STATISTICA, снижается на 22 % (Рис. 2).

² См. ссылку № 1.

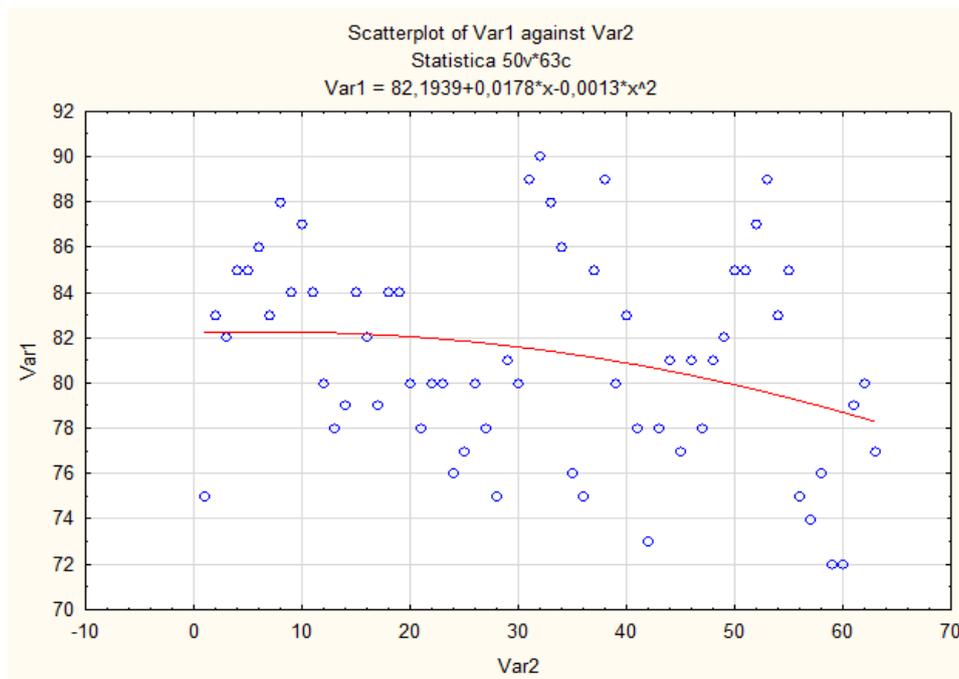


Рис. 2. Квадратичный тренд временного ряда ИСН с марта 2009 по август 2019.

Между отметками 0 и 40 снижение тренда близко к 6 %, поэтому будем считать, что ИСН - квазистационарный. Согласно представленному выше результату, индекс должен иметь (приблизительно) нормальное распределение. Проверим это с помощью теста Шапиро-Уилка (для применения теста Колмогорова-Смирнова число измерений должно быть не меньше 60, а в данном случае их 40). Получим критическую вероятность $p = 0,268$, что дает основание не отвергать гипотезу нормальности ($p > 0,05$).

§ 5. Заключение.

В статье теоретически обоснована гипотеза о том, что социологические индексы, исчисляемые по формулам вида (1), должны иметь (приблизительно) нормальные законы распределения на тех промежутках времени, где эти индексы квазистационарны. Условие квазистационарности означает следующее.

- 1) каждый опрос, по результатам которого вычисляются значения индекса, проводится среди респондентов, выбираемых случайно в соответствующей, социальной группе;
- 2) случайная величина (5), определяемая ответами случайного респондента, является приблизительно стационарной (т.е. не зависящей от времени).

Надежно проверить условие (2) непросто, поэтому о его справедливости можно судить косвенно - по значениям временного ряда, отвечающего индексу. В квазистационарном случае сезонная компонента ряда отсутствует или незначительна, а тренд меняется на данном промежутке времени так, что его можно считать постоянным.

Одно из практических приложений факта нормальности распределения состоит в том, что возникает возможность экстраполировать доступную статистику значений этого индекса или детализировать ее для более частых отметок дат. Например, «Левада-центр» публикует статистику индекса социальных настроений (ИСН) с частотой 2 месяца, начиная с 1995 года. Если на каком-то промежутке времени ИСН был квазистационарным, то знание (нормального) закона распределения соответствующей, случайной величины позволяет моделировать значения ИСН на этом промежутке времени с частотой, скажем, 1 неделя. Или экстраполировать ИСН на будущий период в предположении, что значительных социальных изменений за это время не произойдет.

После надежной, эмпирической проверки результат настоящей статьи может оказаться полезным для анализа методик исчисления социологических индексов, а также использоваться для предварительных оценок размеров выборок при опросах, что могло бы снизить расходы на социальный мониторинг.

В обзоре истории выборочных исследований в США Франкель и Франкель [Frankel and Frankel, 1987, p. 129] писали: «После 1948 года споры сторонников квотной и случайной выборки закончились, и случайная выборка стала для США предпочтительным методом».

. Тем не менее, как уже было отмечено рабочей группой Американской ассоциации исследователей общественного мнения по вопросам онлайн-панелей (AAPOR Task Force on Online Panels), такого рода опросы несомненно представляют ценность для некоторых видов исследований, но исследователям «следует избегать опт-ин-панелей с неслучайным методом отбора респондентов, когда основная задача заключается в точной оценке характеристик всей генеральной совокупности <...> при использовании таких источников для выборки следует избегать претензии на „репрезентативность“».

«Мы считаем, что методы сбора данных и расчёта показателей, не имеющие под собой теоретической основы, не подходят для получения статистических выводов. Конформная выборка является одним из таких методов. В докладе мы не будем говорить о конформной выборке именно по этой причине – из-за отсутствия теории, но для полноты картины ниже коротко её опишем.»

Список литературы

AAPOR (Американская ассоциация исследователей общественного мнения). Отчёт рабочей группы AAPOR о неслучайных выборках: июнь 2013. Перевод с англ. Рогозина Д.М.,

Ипатовой. А. Москва: Общероссийский общественный фонд «Общественное мнение». 2016.

Балацкий Е.В., Саакянц К.М. Индексы социального неравенства // Мониторинг общественного мнения. 2006. № 2. С. 122-128.

Вентцель Е.С. Теория вероятностей: учебное пособие для вузов. Москва: «ФизматГИЗ», 1962.

Гмурман В.Е. Теория вероятностей и математическая статистика: учебное пособие для вузов. Москва: Высшая школа, 2003.

Давыдов А.А. Анализ одномерных частотных распределений в социологии: эволюция подходов // Социологические исследования. 1995. № 5. С. 113-116.

Дуганов М.Д., Калашников К.Н. Методологические подходы к оценке эффективности регионального здравоохранения // Экономические и социальные перемены: факты, тенденции, прогноз. 2011. № 6. С. 93-105.

Ильясов Ф.Н. Типы шкал и распределений в социологии. // Мониторинг общественного мнения. 2014. № 4. С. 24-40.

Крыштановский А. О. Методы анализа временных рядов // Мониторинг общественного мнения : экономические и социальные перемены. 2000. № 2. С. 44-51.

Ляпунов А.М. Избранные труды. Москва: Издательство АН СССР. 1948.

Neuman, J. On the Two Different Aspects of the Representative Method: The Method of Stratified Sampling and the Method of Purposive Selection // Journal of the Royal Statistical Society. 1934. № 97. P. 558–625.

References

Baker R., Brick J.M., Bates N.A., Battaglia M., Couper M.P., Dever J.A., Gile K.J., Tourangeau R. Summary Report of the AAPOR Task Force on Non-probability Sampling // Journal of Survey Statistics and Methodology. 2013. V. 1, № 2. P. 90–143.

Balatskii E.V., Saakyants K.M. Indices of Social Inequality // Public opinion monitoring. 2006. № 2. С. 122- 128. (In Russian)

Gmurman V.E. Probability Theory and Mathematical Statistics: Textbook. manual for universities. M.: Higher school. 2003. (In Russian)

Davidov A.A. Analysis of one-dimensional frequency distributions in sociology: the evolution of approaches // Sociological Studies. 1995. № 5. P. 113-116. (In Russian)

Duganov M.D., Kalashnikov K.N. Methodological approaches to assessing the effectiveness of regional health care // Economic and social changes: facts, trends, forecast. 2011. № 6, P. 93-105. (In Russian)

Iliassov F.N. Types of scales and analysis of distributions in sociology // Public opinion monitoring. 2014. № 4. P. 24-40. (In Russian)

Kryshtanovskii A.O. Methods of analysis of time series // Public opinion monitoring. 2000. № 2. P. 44-51. (In Russian)

Lyapunov A.M. Selected Works, M: Publishing House of the Academy of Sciences of the USSR, 1948. (In Russian)

Neyman, J. On the Two Different Aspects of the Representative Method: The Method of Stratified Sampling and the Method of Purposive Selection // Journal of the Royal Statistical Society. 1934. № 97. P. 558–625.

Wentzel E.S., Probability Theory: textbook for universities. M.: "FizmatGIZ", 1962. (In Russian).