



УДК 159.9.075

Новые возможности статистической системы IBM SPSS Statistics для обработки данных психологических исследований

А. Н. Гусев

Московский государственный университет им. М. В. Ломоносова, г. Москва

Аннотация. Представлен обзор некоторых новых процедур обработки данных, включенных в последние годы в статистическую систему IBM SPSS Statistics (начиная с 19-й версии), которые могут быть полезны психологам при выполнении академических и прикладных исследований. Дана характеристика нового типа статистического анализа данных – многоуровневого, в котором факторы более высокого порядка рассматриваются как медиаторы и модераторы. Показано, что многоуровневый анализ позволяет количественно оценивать вклад этих факторов в характер влияния изучаемых независимых переменных. Описано назначение новых статистических процедур, основанных на логике обобщенной линейной модели. Выделены особенности обобщенной линейной модели как нового, универсального подхода к обработке эмпирических данных. Представлен краткий обзор современных регрессионных моделей, позволяющих обрабатывать разные типы данных.

Ключевые слова: статистическая система IBM SPSS Statistics, метрические и неметрические данные, обобщенная линейная модель, многоуровневый анализ данных, регрессионные модели.

Часто при проведении научных исследований психолог работает с данными низкого уровня измерения, полученными по шкале наименований или шкале порядка [2]. Это так называемые неметрические данные. Их источником является использование разнообразных процедур классификации объектов измерения по полу, национальности, профилю образования, социально-экономическому статусу респондента и другим изучаемым характеристикам, а их результатом – бинарные (например, пол, успешность/неуспешность решения задачи) или мультиномиальные (национальность, профессия, способ решения тестового задания) данные. Также весьма распространено применение процедур ранжирования с использованием числовой или графической шкал, результатом которых являются порядковые данные (академическая успеваемость, разного рода самооценки). Результат проведенных психологических измерений, как правило, выражается в числах, но крайне важно иметь в виду, что эти числа являются результатом установления отношений невысокого уровня – номинальных или порядковых [Там же]. Более того, очень часто после получения от участников исследования словесных категориальных оценок психологи преобразуют их для удобства в числовые значения,

не полностью отдавая себе отчет в том, что эти преобразованные данные остаются номинальными или порядковыми, т. е. по своей природе *неметрическими*.

Данные, полученные по шкалам отношений или интервалов, так называемые количественные, или *метрические* данные, относительно более редки в психологических исследованиях. В первую очередь это результаты тестовых оценок, полученных по стандартной шкале, антропометрические, физиологические данные, результаты приборных измерений, соотносимых с единицами измерений и даже имеющих конвенциональные нулевые значения на используемой шкале.

Невысокий уровень получаемых «сырых» данных несколько лет назад являлся причиной невозможности для исследователей использовать достаточно мощные и широко распространенные в психологии статистические процедуры дисперсионного (группа процедур «Общая линейная модель» в статистической системе IBM SPSS Statistics) и регрессионного анализа, поскольку они работают исключительно с количественными данными (шкалы интервалов и шкалы отношений), требуют нормального распределения результатов измерений и ограничены рядом других допущений [1]. Указанная причина не позволяла проводить анализ экспериментальных эффектов межфакторного взаимодействия, обрабатывать данные экспериментов, включающих одновременно как *межгрупповые*, так и *внутригрупповые* факторы (т. е. реализовывать так называемый смешанный дизайн исследования), обрабатывать эксперименты со *случайными* факторами, оценивать эффекты *ковариат*. На наш взгляд, такое ограничение весьма серьезно, поскольку в психологии весьма часто проводят повторные измерения (например, оценивая динамику академической успеваемости учащихся по четвертям или полугодиям), одновременно контролируя межгрупповые независимые переменные (например, пол, возраст респондента или вид экспериментального воздействия). Более того, психологи не имели возможности оценивать влияние на измеряемые (зависимые) переменные так называемых *разноуровневых* независимых переменных, таких как школа, вуз, район, регион и др., и таким образом обрабатывать *иерархически* организованные (кластеризованные) данные, проверяя более сложные и зачастую даже более интересные гипотезы.

Обобщенные линейные модели

В последние годы в SPSS (начиная с версии 19)¹ получили распространение новые статистические процедуры групп «Обобщенные линейные модели» и «Смешанные модели», предоставляющие психологам и педагогам более широкие возможности анализа данных низкого уровня измерения. Это группа весьма современных процедур, основанных на методологии так называемой *обобщенной* линейной модели. Ниже мы дадим общую характеристику этих процедур, описав их особенности и подчеркивая специфику решаемых ими задач по статистическому анализу данных, получаемых в исследованиях психологов.

¹ В весьма распространенной среди психологов и педагогов статистической системе Statistica эта группа процедур также имеется.

Главная *специфика* и одновременно преимущество процедур, основанных на обобщенной линейной модели, – это многоуровневый и многофакторный анализ неметрических данных. Естественно, что с помощью них мы можем обрабатывать и метрические данные тоже, но это возможно делать также с помощью уже хорошо известных психологам процедур из группы «Общая линейная модель» – одномерного дисперсионного анализа, многомерного дисперсионного анализа и дисперсионного анализа с повторными измерениями [1; 3]. Уникальная возможность, предоставляемая процедурами «Обобщенной линейной модели», состоит в том, что привычные психологам задачи оценки факторных эффектов и эффектов межфакторного взаимодействия решаются для *номинальных* (категориальные оценки типа «да – нет», «выполнил – не выполнил»), *мультиномиальных* (несколько оценочных категорий) и *порядковых* («отлично», «хорошо», «удовлетворительно», «плохо») данных. Таким образом, эти более современные модели позволяют психологу преодолеть известную дефицитарность процедур традиционного дисперсионного анализа в рамках ставшей уже классической общей линейной модели за счет *расширения* типов данных, которые подвергаются многофакторному статистическому анализу.

Еще одним принципиальным преимуществом одной из процедур, основанной на обобщенной линейной модели, является возможность обработки так называемых *многоуровневых* данных – это достаточно новый подход к статистическому анализу данных, имеющих вложенную или иерархическую структуру. Иначе говоря, это построение моделей, в которых наши эмпирические наблюдения представлены в разных контекстах. В современной литературе по статистике этот подход описывается разными терминами: случайные коэффициенты, смешанные эффекты, иерархические линейные или многоуровневые регрессионные модели (см., напр., [4]).

Поясним на примере, что означает психологическое исследование с многоуровневыми данными. Представим себе исследовательский проект, в рамках которого мы планируем изучение психофизиологических ресурсов (различные оценки функционального состояния – ФС, по индексам, по показателям периферической или центральной нервной систем) студентов большого вуза, включающего несколько факультетов и множество кафедр. Обозначим три уровня анализа данных:

1-й уровень: показатели ФС зависят от ряда индивидуальных особенностей студентов – пол, темперамент, мотивация, тревожность и др.;

2-й уровень: показатели ФС зависят от набора и сложности учебных дисциплин, требований преподавателей, возможности профессионального развития на конкретной кафедре;

3-й уровень: показатели ФС зависят от особенности обучения, организации учебного процесса, требований учебного процесса на конкретном факультете.

В данном случае с точки зрения структуры получаемых данных изучение вариации некоторой зависимой переменной (например, стабильности частоты сердечных сокращений) на одном уровне – это сравнение регрессионных пря-

мых при разных условиях этого уровня (рис. 1). Наклон регрессионной прямой мы можем рассматривать как показатель влияния некоторого условия (контролируемого исследователем фактора) нашего вымышленного эксперимента: например, уровня мотивации достижения или личностной тревожности. При многоуровневом подходе мы рассматриваем отношение зависимых и независимых переменных на разных уровнях, которые, по сути, являются *модераторами* или *медиаторами* исследуемых соотношений на нижнем уровне. Это так называемы cross-level interaction модели.

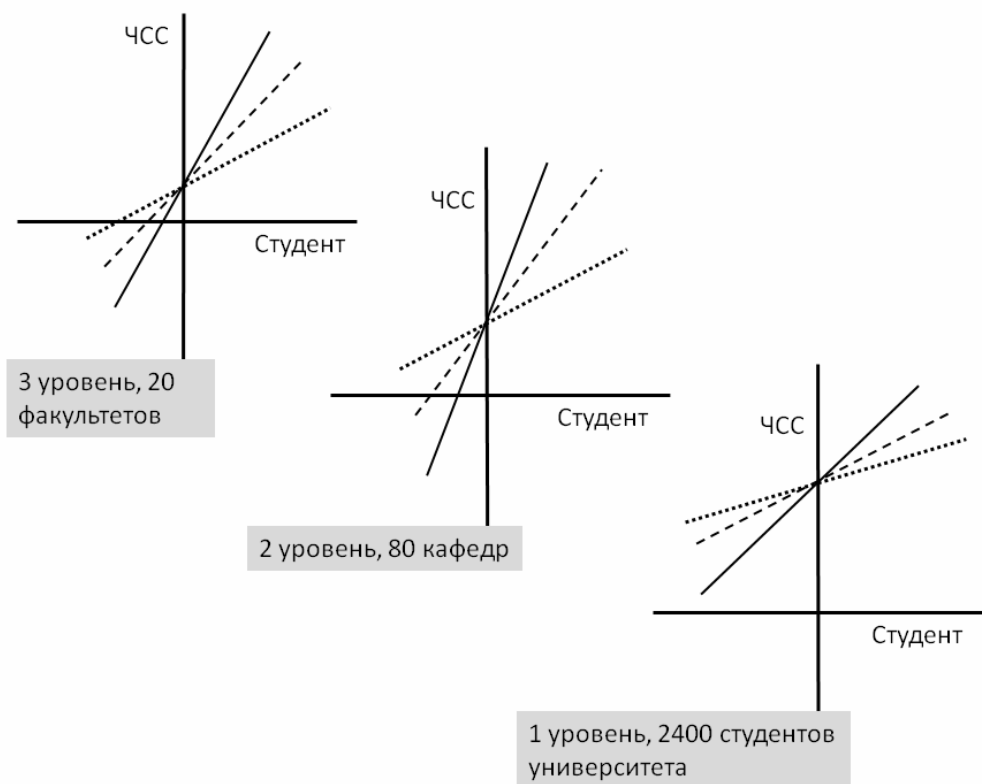


Рис. 1. Обработка данных гипотетического эксперимента по оценке влияния трех индивидуально-психологических факторов – экстраверсии (точечная линия), личностной тревожности (пунктир), мотивации достижения (сплошная линия) – на частоту сердечных сокращений (ЧСС) студентов вуза как показатель функционального состояния сердечно-сосудистой системы. 1 уровень анализа – учет индивидуально-психологических особенностей студентов, 2 уровень – учет специфики требований кафедры, 3 уровень – учет специфики требований факультета. Ось абсцисс – студенты, ось ординат – оценка ЧСС в функциональной пробе

Укажем на основные особенности многоуровневого анализа:

1) изучаемые нами факторы более высокого уровня – это не случайные (вложенные или *nesting*) факторы, а *модераторы* и *медиаторы*, т. е. факторы, опосредующие влияние индивидуальных особенностей студентов на проявление психофизиологических ресурсов обучающегося;

2) проводя сравнительный анализ вариации зависимой переменной при переходе от уровня к уровню, мы выделяем источники вариации для *каждого* уровня;

3) мы не делаем усреднение данных по группам испытуемых на каждом уровне, т. е. мы их не анализируем только на одном уровне, *исключая из анализа вариацию* зависимой переменной по испытуемым в каждой группе на каждом уровне;

4) мы рассматриваем вариацию зависимой переменной на каком-либо уровне не как признак индивидуальности испытуемого, а как характеристику уровня его анализа;

5) мы стараемся оценить *эффекты межуровневого взаимодействия* как предикторы различий между группами испытуемых на разных уровнях анализа, полагая, что вертикальные связи высокоуровневых факторов с переменными на более низком уровне *модерируют* их связи между собой.

Основа рассматриваемого нами многоуровневого подхода – это работы С. Роденбаша и А. Брика (так называемые *slopes-as-outcome model*), в которых предлагается, что в ходе статистического анализа данных нужно оценивать не только вариабельность групповых средних – интерсептов на каком-либо уровне, но и вариабельность регрессионных коэффициентов, т. е. наклонов прямых [5]. В рамках рассматриваемой модели мы предполагаем, что межуровневые взаимодействия объясняют вариацию наклонов регрессионных прямых, соответствующих разным группам испытуемых (уровням независимой переменной) на нижележащем уровне. На рисунке 1 этот аспект специально подчеркнут: на каждом уровне анализа его вклад выражен в вариации наклонов регрессионных прямых и величине их интерсептов (уровень относительно горизонтальной оси).

Рассматривая наш вымышленный эксперимент, мы можем проверить, например, такую гипотезу: чем больше вклад факультета в учебный процесс в рамках всего вуза (т. е. на макроуровне), тем выше связь мотивации достижения с уровнем ФС студента перед сессией. Или такую гипотезу: особенности модульной системы обучения на данном факультете снижают для «сильных» студентов связь между их научной мотивацией и оптимальным уровнем их ФС.

Далее остановимся на основных идеях, положенных в основу обобщенной линейной модели как нового статистического подхода к анализу неметрических данных.

1. Обобщенная линейная модель (*GenLin-model*) – это расширение общей линейной модели, описывающее зависимую переменную как линейно связанную с факторами и ковариатами посредством соответствующей *функции связи*. *GenLin-model* позволяет эмпирически оцененной зависимой переменной

иметь ненормальное распределение и различный уровень измерения, что покрывает широкий ряд статистических моделей:

- линейная регрессия для нормально распределенных данных;
- логистические модели для бинарных данных;
- логлинейные модели для частотных данных;
- лог-лог-модели для интервальных данных и др.

Для *каждого типа* распределения эмпирических данных и для *каждого уровня* измерения подбирается соответствующая ему функция связи. Например, функция связи вида $f(x) = \log(-\log(1-x))$ соответствует биномиальному распределению, а функция $f(x) = \tan(\pi(x - 0,5))$ соответствует только мультиномиальному распределению. Таким образом, согласно специфике полученных данных в рамках обобщенной линейной модели подбирается соответствующее преобразование, делающее их линейными.

Кратко опишем те четыре процедуры, которые стали доступны исследователям в статистической системе IBM SPSS Statistics (начиная с 2010 г. версия 19).

1. *Обобщенные линейные модели*: позволяют проводить многофакторный анализ влияния *межгрупповых* факторов и ковариат на количественные, порядковые, номинальные и мультиномиальные зависимые переменные.

2. *Обобщенные уравнения оценки*: позволяют проводить многофакторный анализ совокупного влияния *внутригрупповых и межгрупповых* факторов и ковариат на количественные, порядковые, номинальные и мультиномиальные зависимые переменные, т. е. включать в обработку данные *повторных измерений*.

3. *Смешанные линейные модели*: позволяют проводить многофакторный анализ совокупного влияния *внутригрупповых и межгрупповых* факторов и ковариат, а также *случайных* факторов на количественные, порядковые, номинальные и мультиномиальные зависимые переменные. Специфика данной процедуры – возможность реализации экспериментальных планов со *случайными факторами*.

4. *Обобщенные смешанные модели*: самая универсальная и многофункциональная процедура. Позволяет проводить многофакторный анализ совокупного влияния *внутригрупповых и межгрупповых* факторов (фиксированных и/или случайных) и ковариат на порядковые, номинальные и мультиномиальные зависимые переменные. Специфика данной процедуры – возможность проведения *многоуровневого* анализа данных, т. е. оценки вклада в общую дисперсию факторов высшего порядка.

Для знакомства с указанными процедурами мы рекомендуем обратиться к соответствующей документации статистической системы IBM SPSS Statistics, свободно представленной на сайте компании IBM. Для детального освоения этих процедур с использованием хороших примеров из психолого-педагогических исследований и детальным разбором сути каждой процедуры следует обратиться к блестящей книге Multilevel Modeling of Categorical Outcomes Using IBM SPSS Рональда Хека, Скотта Томаса и Лин Табата [4].

Новые процедуры регрессионного анализа

Психологи, обрабатывающие данные эмпирических исследований, часто применяют регрессионный анализ как способ проверить статистические гипотезы о взаимодействии двух переменных, исходя из предположения о влиянии одной переменной на другую, т. е. рассматривая одну переменную как предиктор, объясняющий, предсказывающий вариацию другой – зависимой переменной (ЗП). Статистически – это попытка объяснить дисперсию ЗП вкладом влияния одной или нескольких независимых переменных (НЗП) как ее предикторов. Традиционно в психологических исследованиях используют процедуры линейного *одномерного* (один предиктор) или *многомерного* (несколько предикторов) регрессионного анализа [3]. В соответствии с требованиями этих процедур в обработку могут быть включены лишь метрические данные. Как правило, большинство исследователей используют линейные модели регрессионного анализа, априори предполагая, что связь анализируемых переменных описывается линейной функцией: $y = ax + b$.

Кратко изложим назначение достаточно новых процедур регрессионного анализа, появившихся в статистической системе IBM SPSS Statistics в последние годы [6].

Регрессия частично наименьших квадратов (PLS) представляет собой метод для предсказания изменения количественной ЗП, который является альтернативой обычной регрессии, основанной на использовании метода наименьших квадратов, каноническим корреляциям или построению линейных моделей с помощью структурных уравнений. Процедура PLS соединяет свойства метода главных компонент и множественной регрессии. Она особенно полезна, когда предикторные переменные сильно коррелированы (т. е. их нельзя рассматривать как независимые факторы) или когда число предикторов превышает число наблюдений. Особенностью использования в SPSS данной процедуры является необходимость дополнительной установки модуля для языка Питон (Integration Plug-in for Python).

Логистическая процедура специально предназначена для построения одномерных и многомерных регрессионных линейных моделей для номинальных *дихотомических* ЗП. Для каждой НЗП вычисляется отношение правдоподобия. НЗП должны быть порядковыми или интервальными.

Мультиномиальная логистическая процедура специально предназначена для построения одномерных и многомерных регрессионных линейных моделей для *мультиномиальных* ЗП. НЗП может быть категориальной переменной или количественной (ковариата).

Порядковая – специфическая процедура для построения одномерных и многомерных регрессионных линейных моделей для *порядковых* ЗП. НЗП может быть категориальной или ковариатой.

Пробит – процедура для построения одномерных и многомерных регрессионных линейных моделей для ЗП, представляющих собой *частоты* ответов испытуемого. НЗП может быть категориальной или ковариатой. Частоты ответов нередко оцениваются в психологических исследованиях, например: количество пересдач студентом экзаменов за учебный год или весь период обу-

чения, количество оригинальных решений, число обращений за помощью к экспериментатору, число пропусков целевого стимула.

Нелинейная регрессия работает только с количественными данными. Позволяет построить одномерную или многомерную регрессионную модель, используя различные – как линейные, так и нелинейные – преобразования НЗП. Эти комбинации задаются в виде формулы самим исследователем.

Взвешенное оценивание предназначено для построения одномерных и многомерных регрессионных линейных моделей только для количественных ЗП и НЗП. Характерной особенностью данной процедуры является предположение о том, что дисперсия ЗП непостоянна для всех уровней НЗП. Интересно, что данные, полученные от разных испытуемых, – отдельные наблюдения, могут иметь разный вес, т. е. в этой процедуре предусмотрена количественная переменная, *взвешивающая* наблюдения, придающая им особый статус в регрессионном уравнении при сочетании с определенным уровнем НЗП.

Общий логлинейный анализ – еще одна процедура для построения многомерных регрессионных линейных моделей для *номинальных* ЗП. Фактически проводится анализ закономерности распределения частот в таблице кросс-табуляции по нескольким НЗП (до 10). В модель могут включаться ковариаты.

Логит-логлинейный анализ – процедура для построения одномерных и многомерных регрессионных линейных моделей для одной или *нескольких номинальных* или *мультиномиальных* ЗП. НЗП может быть категориальной переменной (т. е. номинативной) или ковариатой (т. е. количественной). При построении моделей число ЗП и НЗП ограничено 10. Фактически по назначению – это неметрический аналог MANOVA.

Подбор модели – еще одна процедура для построения многомерных регрессионных линейных моделей для *номинальных* ЗП. Так же как и в процедуре *общий логлинейный анализ*, здесь анализируются закономерности распределения частот в таблице кросс-табуляции по множеству факторов, число которых может быть больше 10. Включение в модель ковариат не предусмотрено.

Представленный выше краткий обзор процедур регрессионного анализа показывает, что они позволяют исследователю работать с неметрическими данными – номинальными, мультиномиальными и порядковыми, частотами ответов, данными, имеющими самые разные распределения, весьма отличные от привычного для многих психологов нормального распределения. Таким образом, для исследователя раскрываются новые горизонты по проверке гипотез не только корреляционного, но и каузального типа по данным всех уровней измерений – от шкалы наименований до шкалы отношений.

Список литературы

1. Гусев А. Н. Дисперсионный анализ в экспериментальной психологии / А. Н. Гусев. – М. : Психология, 2000. – 136 с.
2. Гусев А. Н. Психологические измерения: Теория. Методы: Общепсихологический практикум / А. Н. Гусев, И. С. Уточкин. – М. : Аспект Пресс, 2011. – 317 с.
3. Наследов А. IBM SPSS 20 Statistics и AMOS: профессиональный статистический анализ данных / А. Наследов. – СПб. : Питер, 2013. – 416 с.

4. Heck R. H. Multilevel Modeling of Categorical Outcomes Using IBM SPSS / R. H. Heck, S. L. Thomas, L. N. Tabata. – N. Y. : Routledge, Naylor & Francis Group, 2012. – 448 p.
5. Raudenbush S. W. Hierarchical linear models: Applications and data analysis method. Second Edition / S. W. Raudenbush, A. S. Bryk. – Thousand Oaks, CA: Sage Publications, 2002. – 510 p.
6. Tabachnick B. R. Using Multivariate Statistics / B. R. Tabachnick, L. S. Fidell. – 6-th ed. – Boston : Allyn and Bacon, 2013. – 1024 p.

New Capacities of IBM SPSS Statistics for Handling Psychological Research Data

A. N. Gusev

Moscow State University M. V. Lomonosov, Moscow

Abstract. The paper presents a review of new data handling procedures recently included in IBM SPSS Statistics (beginning from version 19) that may be helpful for psychologists when carrying out theoretical and applied research. Data statistical analysis of a new type is given. It is a multilevel data analysis that presents higher level factors as mediators and moderators. Multilevel analysis enables to assess the effect of these factors on the nature of influence of independent variables under study. The function of four new statistical procedures based on a generalized linear model is described. Specific features of a generalized linear model as a new universal way of empirical data handling are described. A brief review of modern regression models handling various data is given.

Keywords: IBM SPSS Statistics, metric and nonmetric data, generalized linear model, multilevel data analysis, regression models.

*Гусев Алексей Николаевич
доктор психологических наук, профессор
МГУ им. М. В. Ломоносова
119991, г. Москва, Ленинские горы, ГСП-1
e-mail: angusev@mail.ru*

*Gusev Alexey Nikolaevich
Doctor of Sciences (Psychology), Professor
Moscow State University M. V. Lomonosov
GSP-1, Leninskie Gory, Moscow, 119991
e-mail: angusev@mail.ru*