

РЕГРЕССИОННЫЙ АНАЛИЗ. ПРОЦЕССНЫЙ ПОДХОД

*Национальный технический университет Украины «КПИ», Киев, Украина

Анотація. Побудова регресійних моделей, починаючи з формалізації і закінчуючи їх використанням, розглядається з точки зору керування процесами, тобто процесного підходу. Його застосування дає можливість розглядати процес побудови емпіричних математичних моделей як єдине ціле і обґрунтовано визначити критерії якості як кожного процесу, так і технології в цілому.

Ключові слова: регресійний аналіз (РА), процесний підхід, планування експерименту.

Аннотация. Построение регрессионных моделей, начиная с формализации и заканчивая их использованием, рассматривается с точки зрения управления процессами, то есть процессного подхода. Его использование позволяет рассматривать процесс построения эмпирических математических моделей как единое целое и обоснованно определить критерии качества каждого процесса как отдельно, так и технологии в целом.

Ключевые слова: регрессионный анализ (РА), процессный подход, планирование эксперимента.

Abstract. Building of regression models, from the formalization to their use is considered in terms of process control, in other words, the process approach. Its use makes it possible to consider the process of constructing empirical mathematical models as a whole and reasonably determine quality criteria both each process and technology in general.

Keywords: regression analysis (РА), process approach, design of experiment.

1. Введение. Проблема и цель работы

Многие проблемы в использовании планирования эксперимента и регрессионного анализа связаны с тем, что они представляют собой набор отдельных теорий и методов, не имеющих единства [1, 2]. Разумеется, со стороны практиков попытки представить их в виде некоторой технологии построения эмпирических моделей существуют, например, [3, 4]. Но эта технология не имеет теоретического обоснования. Подтверждением ее служит только длительное практическое использование. Предлагается с целью разрешения указанного противоречия рассматривать построение регрессионных моделей как часть выборочного метода [2]. Тогда требования к матрице (выборке) и методу обработки получаются естественным образом. К сожалению, это требует теоретического обоснования высокого уровня теоретиками. Вряд ли кто-то этим займется: кто может – тому неинтересно, кому нужно – недоступно по уровню теоретической подготовки.

Рассмотрим построение регрессионной модели как некоторую деятельность, используя процессный подход. Отметим, что в реальности отдельные подпроцессы этого процесса выполняются разными людьми, да и руководство ими зачастую тоже разное, принадлежащее даже к разным организациям и ведомствам. Использование данного подхода позволит получить обоснованные решения, не проводя вышеупомянутых теоретических изысканий.

Принципиальным в процессном подходе есть то, что на каждом этапе должна обеспечиваться конечная цель [4]. По У.Е. Демингу функционирование любого подпроцесса системы должно оцениваться по его вкладу в исполнение цели всей системы, а не по индивидуальной эффективности.

По сути, это положение заимствовано из принципа динамического программирования Беллмана [5]. Исходя из этого положения, локальные (на отдельном этапе) цели и критерии должны обеспечивать наилучший результат не данного этапа, а конечной цели. Та-

кой подход позволяет выполнить объединение всех ранее разрозненных методов в единый процесс построения модели.

2. Классификация процессов

Процесс – совокупность связанных или взаимодействующих видов деятельности, которые превращают входы в выходы. Деятельность рассматривается как совокупность процессов, каждый из которых характеризуется показателями, описывающими его выполнение, результат и влияние на деятельность в целом. С этой точки зрения деятельность по построению регрессионной модели можно представить в виде следующей совокупности процессов (табл. 1). К процессам, создающим ценность, относятся формализация, проведение экспериментов, определение структуры, идентификация. Прочие процессы создают возможности для создания ценностей.

Таблица 1. Этапы построения регрессионной модели

Процессы	Подпроцессы			
Предпланирование эксперимента	Постановка задачи в предметной области			
	Формализация задачи			
Формирование выборки	Пассивный эксперимент	Сбор данных	Планирование эксперимента	Активный эксперимент
		Выделение обучающей подвыборки	Проведение эксперимента	
Предварительный статистический анализ выборки	Проверка гетероскедастичности			
	Проверка соотношения сигнал/шум			
	Проверка неразрывности факторного пространства			
	Работа с выделяющимися наблюдениями			
Построение модели	Спецификация общая			
	Преобразования переменных специальные			
	Преобразования переменных стандартные			
	Преобразования специальные пространства			
	Спецификация частная			
	Выбор способа оценки коэффициентов модели			
	Идентификация модели			
Анализ качества модели и принятие решений	Анализ статистических характеристик модели			
	Проверка на контрольной последовательности			
	Исследование наличия нарушений и предпосылок РА			
	Анализ адекватности			
	Корректировка модели			
Использование модели (вычислительный эксперимент)	Графический анализ			
	Прогнозирование			
	Оптимизация			

Любой процесс описывается циклом Деминга-Шухарта: План (Plan) – Реализация (Do) – Проверка (Check) – Исправление (Action). При этом для завершения процесса может потребоваться несколько витков.

Рассмотрим, как будут выглядеть некоторые процессы построения модели регрессии. Например, конструирование плана может быть представлено следующим образом.

Plan: Сформировать требования к плану в виде задания на конструирование.

Do: Построить план.

Check: Проверить план на соответствие желаемым статистическим свойствам и реализуемость.

Action: Внести необходимые исправления.

На самом деле циклы более сложные. Например, описанный выше цикл может замыкаться на цикл, включающий его. А именно, при невозможности воспользоваться данным планом переходим в начало для переформулирования задания на план. Если и это невозможно, то переход осуществляется на предыдущий процесс – формализацию, который тоже теперь выполняется в новом цикле.

В общем случае описываемая деятельность представляет собой совокупность вложенных и переплетающихся циклов, разворачивающихся в сложных ситуациях. Рассмотрим её более подробно (в упрощенном виде).

1.1. Определение цели построения математической модели, сформулированной в терминах предметной области (возможно построение дерева целей).

1.2. Определение средств и методов построения модели. Здесь определяется класс математических моделей, в котором возможно построение модели исследуемого процесса или объекта.

1.3. Формулирование требований к модели, то есть перечень требований, которым должна удовлетворять модель, чтобы соответствовать цели (сформулированных в терминах предметной области).

2.1. Определение формализованной цели исследования (исходя из 1.2).

2.2 Анализ и структурирование объекта исследования, основываясь на 2.1.

2.3. Определение необходимых ресурсов на проведение исследований. При невозможности выделения требуемых ресурсов или удовлетворения условия проведения исследования, выдвигаемых выбранными методами, – переход в п. 1.2 или 1.1.

2.4. Установление характеристик проверки адекватности модели (формальные + прикладной области).

2.5. Верификация соответствию постановки задачи. Если ожидаемые полученные результаты не соответствуют цели (1.1), то переход в п. 1.2.

3.1. Определение способа получения данных: активный или пассивный эксперимент.

3.2а. При пассивном эксперименте планирования способа и условий сбора информации.

3.3а. Сбор данных.

3.4а. Выделение обучающей выборки.

3.5а. Определение возможности построения модели на существующей выборке. В случае невозможности – определение необходимых недостающих данных и переход в 3.3а.

3.2б. Формулирование требований к плану эксперимента.

3.3б. Построение плана эксперимента.

3.4б. Проверка свойств плана. В случае неудовлетворения требованиям – переход в п. 3.2б.

3.5б. Построение рабочей матрицы и определение условий эксперимента.

3.6б. Проведение эксперимента.

3.7б. При наличии пропущенных, ошибочных экспериментов или выбросов, обнаруженных в экспериментальном отделе, – проведение дополнительных экспериментов для исправления обнаруженных ошибок.

4.1. Проверка неразрывности факторного пространства.

4.2а. В случае разрывности – разделение на неразрывные подобласти.

4.3а. Определение возможности построить модели в каждой подобласти отдельно.

Если да, то продолжение деятельности по ветке «б».

4.4а. Если нет, то выяснение возможности проведения дополнительных экспериментов или сбора данных для обеспечения построения модели. Если это возможно, то переход в п. 3.2 с дальнейшим выполнением процесса параллельно для каждой подобласти.

4.5а. Если получение дополнительных данных невозможно, то переход в 2.2.

4.2б. Проверка соотношения сигнал/шум в результатах. Если уровень сигнала неотличим от уровня шума, то переход в 2.2.

4.3б. Выявление выделяющихся наблюдений.

4.3.1б. Принятие решений, считать ли их выбросами. Если да, то переход в 3.7б. Если нет – считать, что ошибки распределены не по закону Гаусса.

4.4б. Проверка гетероскедастичности. Если обнаружена гетероскедастичность, то выполняется расчет матрицы весовых коэффициентов для последующей корректировки модели.

5.1. При необходимости следует проведение специальных преобразований, связанных с особенностями поведения отдельных факторов.

5.2. При необходимости – выполнение преобразований пространства, если оно нестандартное.

5.3. Выполнение стандартных преобразований рабочей матрицы (ортогонализация, нормировка, построение взаимодействий).

5.4. Проверка правильности стандартных преобразований. При наличии ошибок – устранение причин и возврат в п. 5.3.

5.5. Определение частной структуры уравнения регрессии.

5.6. Идентификация уравнения регрессии. Производится с учетом наличия гетероскедастичности и отличия распределения ошибок от гауссовского (при необходимости).

6.1. Анализ качества уравнения регрессии.

6.2. Если показатели качества неудовлетворительные, то переход в п.5.5, если возможности по выбору элементов структуры не исчерпаны, или в п.5.3, если необходимо расширение списка регрессоров.

6.3. Проверка нарушений предпосылок и допущений и при необходимости корректировка или переход в п. 5.5.

6.4. Проверка модели на контрольной последовательности опытов и при необходимости переход на п. 5.5 или 5.1.

При более детальном рассмотрении мы имеем множество циклов, позволяющих за конечное число шагов достичь цели. При этом можно описать принятие решения в виде алгоритма или таблицы теории игр с определением потерь.

3. Планирование

Под планированием в процессном подходе понимается установление целей, определение требований, разработка способов осуществления процессов и определение ресурсов. В данном случае целью (удовлетворением потребности) является построение модели. Построить модель означает выполнить ее спецификацию и идентификацию. Расшифровка приведена в виде дерева целей (табл. 2).

Таблица 2. Построение модели (дерево целей)

Модель			
Спецификация		Идентификация	
Общая	Частная	Оценки коэффициентов	Статистические характеристики коэффициентов

В классическом РА структура полагается известной до планирования, но в реальности это не так. Во многих случаях целью построения модели как раз и является выяснение структуры взаимосвязей между независимыми и зависимой переменными. Установление этой структуры и есть одна из задач в регрессионном анализе [4]. Ранее ее формулировали как формирование наиболее информативного множества регрессоров, что не совсем верно, так как в такой формулировке происходит подмена исходной цели из предметной области целью более низкого уровня из математической статистики.

Таким образом, для построения модели нам необходимо определить общую и частную структуры уравнения регрессии, а затем получить оценки коэффициентов регрессии и их статистические характеристики.

Исходя из целей построения модели, мы должны определить требования к процессам.

Для линейной регрессии общая структура представляет собой алгебраическую сумму произвольных функций. Обычно это полином или ряд Фурье (для периодических). Мы же будем рассматривать её как полином Чебышева с возможным предварительным преобразованием отдельных переменных. Эти преобразования необходимы в некоторых особых случаях: при быстро изменяемых функциях [7, 8] и при асимптотических процессах [9], а также при описании процессов со скачками. Обоснование такого выбора: теорема Вееерштрассе, с одной стороны, и ограниченные возможности по количеству экспериментов и необходимость экстраполяции с другой.

Для определения частной структуры требуются:

- правильная общая спецификация;
- неразрывность пространства;
- независимость регрессоров;
- достаточное количество уровней;
- достаточное количество опытов;
- известная дисперсия воспроизводимости;
- эффективные алгоритмы определения частной структуры.

Поскольку первые два требования объяснений не требуют, то следует сказать несколько слов об остальных.

Вне зависимости от алгоритмов определения частной структуры, чем ближе матрица регрессоров к ортогональной, тем выше вероятность правильного определения структуры [2, 4, 11]. Чем больше количество уровней, тем выше возможная степень полинома и сложность модели. Чем больше число экспериментов, тем больше может быть членов в модели. Кроме того, увеличение числа уровней и числа опытов приводит к снижению общей закоррелированности матрицы регрессоров. Дисперсия воспроизводимости нужна как критерий остановки при последовательном включении членов в модель.

Для получения оценок коэффициентов регрессии и обеспечения требуемых их свойств необходимы:

- известная частная структура;
- независимость регрессоров;
- хорошая обусловленность матрицы регрессоров;
- детерминированность рабочей матрицы;

- гомоскедастичность дисперсий ошибок;
- независимость ошибок.

Соответствие требований, вытекающих из цели, отдельным процессам приведено в табл. 3. При этом следует иметь в виду, что часть требований обеспечивается в процессе, с которого они требуются в таблице, лишь частично. Например, «хорошая обусловленность матрицы регрессоров» в процессе построения модели обеспечивается за счет ортогонализации, нормировки и алгоритмов формирования структуры. Но она принципиально не может быть достигнута, если в процессе формализации не обеспечены независимость факторов, в процессе планирования – независимость регрессоров, а при проведении эксперимента – соответствие сгенерированному плану. В тех случаях, когда требования по каким-то причинам не выполнены или не могут быть выполнены, то необходимы соответствующие корректирующие действия. Эти действия увеличивают затраты ресурсов на деятельность по построению модели в целом, в основном вычислительных ресурсов и ресурсов персонала по формированию заданий на корректировки. Иногда эти затраты в несколько раз могут превышать собственно затраты на получение модели в идеальных условиях.

Так, если невозможно достичь физической независимости факторов (формализация задачи), то необходимо преобразование области эксперимента к стандартной форме, например, гиперкуб, обеспечивающий независимость. Случайная рабочая матрица требует обработки методами конфлюэнтного анализа. В ряде случаев необходим возврат на предыдущие этапы, как описано ранее в алгоритме.

Таблица 3. Соответствие общих требований отдельным процессам

Процесс	Требования
Формализация задачи	Правильная общая спецификация Независимость факторов Неразрывность пространства
Планирование эксперимента	Независимость регрессоров Достаточное количество уровней Достаточное количество опытов Устойчивость к незначительным отклонениям
Проведение эксперимента	Известная дисперсия воспроизводимости Независимость ошибок Детерминированность рабочей матрицы Гомоскедастичность дисперсий ошибок Соответствие плану эксперимента
Построение модели	Хорошая обусловленность матрицы регрессоров Эффективные алгоритмы определения частной структуры

Анализ требований приводит нас к выводу, что наилучшим образом удовлетворяют им робастные планы на основе ЛП_т чисел.

4. Результативность и эффективность

Возникает вопрос, как нам оценивать качество полученного результата? В процессном подходе для этого используется результативность. Это уровень реализации цели и эффективность, то есть соотношение между результатами и затраченными ресурсами.

Для оценки уровня реализации цели необходимо множество показателей (рис. 1). Использование только одного показателя, как это часто делается, совершенно неправиль-

но. Дело в том, что модель, имеющая один хороший показатель, может быть, тем не менее, совершенно непригодна для использования. Например, информативная и адекватная модель может быть вычислительно неустойчивой, и её использование бессмысленно [10]. Адекватная модель может быть неинформативной, то есть не нести никакой полезной информации [11]. Любимое занятие прикладников: включить в модель значимые факторы и больше ничего не делать. Но такая модель может быть неадекватной и неинформативной. С другой стороны, адекватной и информативной может быть модель, состоящая из незначимых факторов [12].

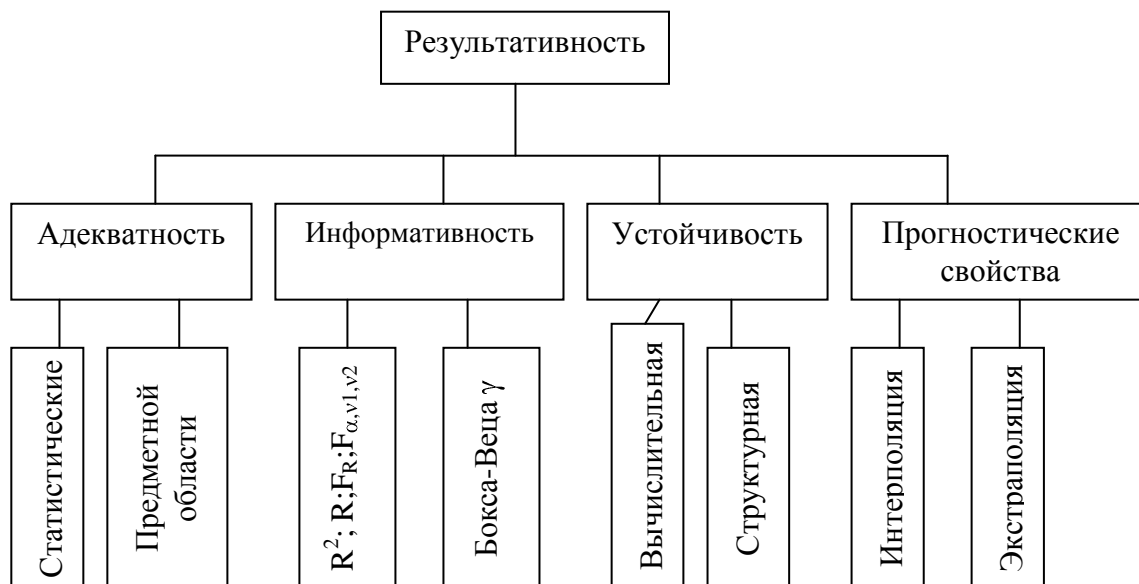


Рис. 1. Дерево оценки результативности

От цели построения модели зависит относительная значимость конкретной ветки, что позволяет выполнять настройку соответствующих алгоритмов построения модели и влиять на конечный результат. Этой же целью определяются конкретные требования к модели в ветке «предметной области». При этом следует помнить, что достижение наилучшего результата по любой ветви может конфликтовать с другой (и, возможно, не одной) ветвью. Например, изменение структуры модели для достижения адекватности или лучших описывающих свойств может приводить к снижению её информативности и устойчивости. Возможные конфликты отражены в табл. 4.

Таблица 4. Возможные конфликты в дереве оценки результативности

	Адекватность статистическая	Адекватность предметной области	R^2	Бокса-Веца γ	Устойчивость вычислительная	Устойчивость структурная	Интерполяция	Экстраполяция
Адекватность статистическая		+		+	+	+		+
Адекватность предметной области	+		+	+			+	
R^2		+			+	+		+
Бокса-Веца γ	+	+			+	+	+	+

Устойчивость вычислительная	+		+	+			+	+
Устойчивость структурная	+	+	+	+			+	+
Интерполяция		+		+	+	+		+
Экстраполяция	+	+	+	+	+	+	+	

Таким образом, для достижения цели построения регрессионной модели необходимо решать задачу многокритериальной оптимизации, позволяющей удовлетворять противоречивым требованиям.

Эффективность определяется затраченными ресурсами по отношению к полученным результатам. То есть необходимо соотносить результативность и затраченные на её достижение ресурсы. Ресурсы при построении модели можно разделить на материальные, человеческие и вычислительные. Относительная стоимость этих ресурсов зависит как от конкретной задачи, так и от условий в организации, её решающей, в частности, от наличия специалистов и программного обеспечения.

Как некий образец рассматривают наивысшую результативность и соответствующие ей затраты ресурсов. Обычно таким образом служит план полного факторного эксперимента (ПФЕ), обеспечивающий наилучшие условия для достижения максимальной результативности. Недостатком такого образца служат запредельные ресурсы для её достижения: огромное количество необходимых экспериментов для реальной многофакторной задачи. Например, план полного факторного эксперимента $3^6/729$ требует 729 опытов. А если учитывать необходимость дублирования опытов, то 1458. А ведь шесть факторов – это достаточно маленькая задача. Поэтому более рациональным является сравнение возможных вариантов друг с другом по результативности и затратам.

Если мы рассмотрим описанный в параграфе 2 упрощенный алгоритм действий при построении модели, то становится ясно, что деятельность по построению модели – это решение задачи в условиях неопределенности. То есть мы вынуждены принимать какие-то решения, не имея достаточной информации об условиях, в которых наши действия будут выполняться. В этих случаях для принятия решения и оценки возможных потерь используется теория игр. В табл. 5 приведена платёжная матрица для принятия решения о выборе вида плана.

Таблица 5. Платёжная матрица для выбора вида плана

Действия (выбранный вид плана)	Проблемная ситуация			
	Недостаточное количество опы- тов	Несоответствие части значений факторов плану эксперимента	Пропущенные эксперименты	Недостаточная степень полинома
План МФРП	Проведение всех эксперимен- тальных иссле- дований по новому плану	Или выполнение «испорченных» экспериментов наново, или рез- кое снижение ре- зультативности	Резкое сниже- ние результа- тивности	Проведение всех эксперименталь- ных исследова- ний по новой матрице
План ЛП _т	Проведение до- полнительных экспериментов	Не требуется	Не требуется	Выполнение рас- четов с уточнен- ными степенями

Конкурирующими планами являются план на основе многофакторных регулярных планов (МФРП) и план на основе равномерно распределенных псевдослучайных чисел (ЛП_т). При этом первые в общем случае имеют лучшие статистические свойства, то есть должны обеспечить более высокую результативность, но только в случае отсутствия проблемных ситуаций. Как видно из табл. 5, при наличии проблемных ситуаций эффективность и надежность выше при использовании ЛП_т плана. В этом случае ресурсы, затрачиваемые в случае ЛП_т плана, меньше и/или результативность выше. Так, при недостаточном количестве опытов или недостаточной степени полинома для МФРП затраты на экспериментальную часть увеличиваются вдвое, а для ЛП_т плана в первом случае увеличиваются незначительно, а во втором увеличиваются затраты вычислительных ресурсов вдвое, что не является существенным в общем случае. В остальных случаях для МФРП происходит снижение результативности, а для ЛП_т плана «убытки» отсутствуют.

Любые дополнительные действия и циклы снижают эффективность, поэтому отсюда следует необходимость выполнять требования теории планирования к проведению эксперимента. Они снижают возможность нарушения предпосылок и допущений и повышают эффективность и результативность. Например, рандомизация последовательности опытов, которую очень не любят экспериментаторы, препятствует образованию зависимости между ошибками экспериментов и заниженному определению дисперсии воспроизводимости, а, следовательно, повышает результативность.

5. Результаты и выводы

Построение регрессионной модели является сложной деятельностью, состоящей из множества взаимосвязанных процессов. В общем случае для её успешного выполнения необходимо задействовать не только знания математической статистики и предметной области, в которой строится модель, но других областей знаний: многокритериальной оптимизации, теории игр и пр. Применение процессного подхода позволяет осознанно рассматривать эту деятельность в единстве, независимо от области применимых в каждом процессе знаний. Следовательно, возможно обоснованное выдвижение требований на каждом этапе, позволяющее обеспечить достижение требуемой цели. Также становится возможным обоснованная оценка качества полученного результата, то есть модели. Таким образом, становится возможным устранение недостатков планирования экспериментов и регрессионного анализа, препятствующих их эффективному использованию.

Побочным эффектом становится облегчение обучению и использованию методов научными работниками, для которых математическая статистика не является специальностью. Им гораздо проще воспринимать не как математические методы, а как технологические процессы обработки данных.

Направление дальнейших работ

Детальное описание деятельности по построению регрессионной модели с поддержкой принятия решения находится в стадии разработки.

СПИСОК ЛИТЕРАТУРЫ

1. Налимов В.В. Логические основания планирования эксперимента / В.В. Налимов, Т.И. Голикова. – [2-е изд., перераб. и доп.]. – М.: Металлургия, 1980. – 152 с.
2. Лапач С.Н. Основные проблемы построения регрессионных моделей / С.Н. Лапач, С.Г. Радченко // Математичні машини і системи. – 2012. – № 4. – С. 125 – 133.
3. Лапач С.Н. Планирование, регрессия и анализ моделей PRIAM (ПРИАМ) / С.Н. Лапач, С.Г. Радченко, П.Н. Бабич // Каталог программные продукты Украины. – К., 1993. – С. 24 – 27.
4. Лапач С.Н. Статистические методы в фармакологии и маркетинге фармацевтического рынка / С.Н. Лапач, М.Ф. Пасечник, А.В. Чубенко. – К.: ЗАТ “Укрспецмонтаж”, 1999. – 312 с.

5. Галямина И.Г. Управление процессами / Галямина И.Г. – СПб.: Питер, 2013. – 304 с.
6. Беллман Р. Динамическое программирование / Беллман Р. – М.: Иностранная литература, 1960. – 400 с.
7. Лагутин М.В. Наглядная математическая статистика / Лагутин М.В. – М.: Бином. Лаборатория знаний, 2007. – 472 с.
8. Калиткин Н.Н. Численные методы / Калиткин Н.Н. – М.: Наука, ГРФМЛ, 1978. – 512 с.
9. Лапач С.М. Лінійна регресія при прогнозуванні асимптотичних залежностей / С.М. Лапач // Вестник Херсонского национального технического университета. – 2010. – № 3 (39). – С. 257 – 260.
10. Петрович М.Л. Регрессионный анализ и его математическое обеспечение на ЕС ЭВМ: Практическое руководство / Петрович М.Л. – М.: Финансы и статистика, 1982. – 199 с.
11. Лапач С.Н. Проблемы построения математических моделей экспериментально-статистическими методами / С.Н. Лапач // Прогресивна техніка і технологія машинобудування, приладобудування і зварювального виробництва. Праці НТУУ “КПІ”. – Т. 2. – К.: НТУУ “КПІ”, – 1998. – С. 25 – 29.
12. Pardoux C. Sur la selection de variables en regression multiple / C. Pardoux // Cah. Bur. Univ. rech. oper. – 1982. – № 39–40. – P. 101 – 133.

Стаття надійшла до редакції 26.08.2015