

Н.В. Комарова

*Днепропетровский национальный университет им. Олеся Гончара***МЕТОД БЛИЖАЙШИХ СОСЕДЕЙ С НОРМАЛИЗОВАННЫМ ЕВКЛИДОВЫМ РАССТОЯНИЕМ**

Пропонується використання нормалізованої евклідової відстані у методі найближчих сусідів для прогнозування часових рядів, приводиться порівняльний аналіз застосування методу до прогнозу цін на метали та металопродукцію.

Предлагается использование нормализованного евклидоваго расстояния в методе ближайших соседей для прогнозирования временных рядов, приводится сравнительный анализ результатов применения метода для прогноза цен на металлы и металлопродукцию.

In the article the using of normalized Euclidian distance is considered in the nearest neighbor method for time series forecasting; the results of the method application to metal and steel product prices forecasting are analyzed in comparison with the results of other methods.

Ключевые слова: нормализованное евклидово расстояние, метод ближайших соседей, прогнозирование, временные ряды.

Вступлення. Аналіз і прогнозування часових рядів – важлива задача во многих областях исследований. Основными инструментами ее решения являются статистические методы, среди которых наиболее широкое распространение получили методы авторегрессии и скользящего среднего. Однако их применение для прогноза многих временных рядов, описывающих сложные системы, не дает удовлетворительных результатов.

Новый подход предложила нелинейная динамика, изучающая структуру и свойства эволюционных процессов в нелинейных динамических системах.

Постановка задачи. Нестационарные временные ряды встречаются довольно часто. Известно, что классические методы авторегрессии, скользящего среднего и индексов сезонности часто не являются адекватным инструментом их прогнозирования. Поэтому возникает задача поиска метода прогнозирования, который можно было бы использовать для широкого класса нестационарных временных рядов. Одним из методов, используемых для прогнозирования нестационарных временных рядов, является метод ближайших соседей, предложенный нелинейной динамикой. Зачастую в методе используется евклидово расстояние, однако результаты прогноза при этом не всегда удовлетворительны, в частности, когда ряд имеет ярко выраженный тренд. Решению этой задачи посвящена настоящая работа.

Подход к решению

Метод ближайших соседей прогнозирования временных рядов. В последнее время при решении задачи прогнозирования временных рядов широко используется метод ближайших соседей, предложенный нелинейной динамикой (нелинейная динамика изучает структуру и свойства эволюционных процессов в нелинейных динамических системах).

Пусть $\{x_i, i = 1, 2, \dots, N\}$ – временной ряд. По значениям $x_i, i = 1, 2, \dots, M-1$, временного ряда необходимо построить прогноз \tilde{x}_M его значения x_M в момент времени M . Для прогнозирования воспользуемся классическим методом ближайших соседей в его простейшем варианте одного ближайшего соседа нулевого порядка [1]. Строим фазовую траекторию временного ряда $\{x_i, i = 1, 2, \dots, N\}$ – последовательность точек с координатами $\mathbf{x}_i = (x_i, x_{i+\tau}, \dots, x_{i+(m-1)\tau})$, $i = 1, 2, \dots, N - (m-1)\tau$; ($m, \tau \ll N$) в пространстве \mathbb{R}^m . Для точки $\mathbf{x}_{M-1-(m-1)\tau} = (x_{M-1-(m-1)\tau}, x_{M-1-(m-2)\tau}, \dots, x_{M-1})$ находим точку $\mathbf{x}_{k-1-(m-1)\tau} = (x_{k-1-(m-1)\tau}, x_{k-1-(m-2)\tau}, \dots, x_{k-1})$ фазовой траектории, ближайшую к точке $\mathbf{x}_{M-1-(m-1)\tau}$. Следующей за точкой $\mathbf{x}_{k-1-(m-1)\tau}$ фазовой траектории является точка $\mathbf{x}_{k-(m-1)\tau} = (x_{k-(m-1)\tau}, x_{k-(m-1)\tau+1}, \dots, x_k)$. Согласно методу ближайших соседей в качестве прогноза \tilde{x}_M значения x_M временного ряда в момент M принимается

значение x_k . Если необходимо построить прогноз на несколько точек $x_M, x_{M+1}, \dots, x_{M+L}$, то, как правило, в качестве прогнозных значений $\tilde{x}_M, \tilde{x}_{M+1}, \dots, \tilde{x}_{M+L}$ рассматривают значения $x_k, x_{k+1}, \dots, x_{k+L}$ временного ряда.

Прогноз можно строить по двум и большему числу ближайших соседей, и в этом случае в качестве прогнозных значений рассматривается среднее или средневзвешенное значение прогнозов по каждому из ближайших соседей.

Выбор параметров m фазового пространства \mathbb{R}^m и параметра τ является отдельной задачей. На практике не худшим вариантом оказывается подбор этих параметров.

Нормализованное евклидово расстояние. В методе ближайших соседей в качестве меры близости между точками фазового пространства, как правило, используют евклидово расстояние. В настоящей работе в качестве меры близости предложено использовать так называемое нормализованное евклидово расстояние.

Введем обозначения:

$$\mathbf{x} = (x_1, x_2, \dots, x_n), \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{\mathbf{x}} = (\bar{x}, \bar{x}, \dots, \bar{x}), \quad \|\mathbf{x}\|^2 = \frac{1}{n} \sum_{i=1}^n x_i^2, \quad d_E - \text{евклидова метрика в } \mathbb{R}^n.$$

Определим нормализованное евклидово расстояние d_{NE} между точками $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ равенством

$$d_{NE}(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\|\mathbf{x} - \bar{\mathbf{x}}\|} - \frac{y_i - \bar{y}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right)^2 \right)^{1/2} = d_E \left(\frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|}, \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right). \quad (1)$$

Нормализованное евклидово расстояние, хотя и определяется через евклидову метрику, но само на \mathbb{R}^n метрикой не является, уже хотя бы потому, что из $d_{NE}(\mathbf{x}, \mathbf{y}) = 0$ не следует равенство $\mathbf{x} = \mathbf{y}$. Например, если $\mathbf{y} = \mathbf{a}\mathbf{x} + \mathbf{b}$, то $d_{NE}(\mathbf{x}, \mathbf{y}) = 0$. Заметим также, что d_{NE} не определено для пары точек, среди которых хотя бы одна имеет вид $\mathbf{c} = (c, c, \dots, c)$. Поэтому далее, когда мы будем говорить о нормализованном евклидовом расстоянии между точками \mathbf{x} и \mathbf{y} , всегда будем предполагать, что они не равны константам, т. е. \mathbf{x} и \mathbf{y} отличны от $\mathbf{c} = (c, c, \dots, c)$.

Нормализованное евклидово расстояние d_{NE} обладает следующими свойствами:

1. $d_{NE}(\mathbf{x}, \mathbf{y}) \geq 0$.
2. $d_{NE}(\mathbf{x}, \mathbf{y}) = 0$ тогда и только тогда, когда $\mathbf{y} = \mathbf{a}\mathbf{x} + \mathbf{b}$.
3. $d_{NE}(\mathbf{x}, \mathbf{y}) = d_{NE}(\mathbf{y}, \mathbf{x})$.
4. $d_{NE}(\mathbf{x}, \mathbf{z}) \leq d_{NE}(\mathbf{x}, \mathbf{y}) + d_{NE}(\mathbf{y}, \mathbf{z})$.

Свойство 1 очевидно.

Свойство 2. Из $\mathbf{x} = \mathbf{y}$ очевидно следует $d_{NE}(\mathbf{x}, \mathbf{y}) = 0$.

Пусть \mathbf{x}, \mathbf{y} таковы, что

$$d_{NE}(\mathbf{x}, \mathbf{y}) = d_E \left(\frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|}, \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right) = 0,$$

тогда $\frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|} = \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|}$. Из последнего равенства для \mathbf{y} получаем представление

$$\mathbf{y} = \frac{\|\mathbf{y} - \bar{\mathbf{y}}\|}{\|\mathbf{x} - \bar{\mathbf{x}}\|} \cdot \mathbf{x} + \left(\bar{\mathbf{y}} - \frac{\|\mathbf{y} - \bar{\mathbf{y}}\|}{\|\mathbf{x} - \bar{\mathbf{x}}\|} \cdot \bar{\mathbf{x}} \right) = a\mathbf{x} + \mathbf{b},$$

где $a = \frac{\|\mathbf{y} - \bar{\mathbf{y}}\|}{\|\mathbf{x} - \bar{\mathbf{x}}\|}$, $\mathbf{b} = \bar{\mathbf{y}} - \frac{\|\mathbf{y} - \bar{\mathbf{y}}\|}{\|\mathbf{x} - \bar{\mathbf{x}}\|} \cdot \bar{\mathbf{x}}$.

Свойство 3 следует из определения (1) d_{NE} .

Свойство 4 следует из соотношений

$$d_{NE}(\mathbf{x}, \mathbf{y}) = d_E \left(\frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|}, \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right) \leq d_E \left(\frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|}, \frac{\mathbf{z} - \bar{\mathbf{z}}}{\|\mathbf{z} - \bar{\mathbf{z}}\|} \right) + d_E \left(\frac{\mathbf{z} - \bar{\mathbf{z}}}{\|\mathbf{z} - \bar{\mathbf{z}}\|}, \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right) = d_{NE}(\mathbf{x}, \mathbf{z}) + d_{NE}(\mathbf{z}, \mathbf{y}).$$

Метод ближайших соседей с нормализованным евклидовым расстоянием. Для построения прогноза \tilde{x}_M значения x_M временного ряда $\{x_i, i = 1, 2, \dots, N\}$ сначала для точки $\mathbf{x}_{M-1-(m-1)\tau} = (x_{M-1-(m-1)\tau}, x_{M-1-(m-2)\tau}, \dots, x_{M-1})$ находим ближайшего соседа $\mathbf{x}_{k-1-(m-1)\tau}$, принадлежащего фазовой траектории в \mathbb{R}^m (ближайшего в смысле нормализованного евклидова расстояния). Ближайший сосед в смысле нормализованного евклидова расстояния существует (выбор производится из конечного множества), но, вообще говоря, не единственный.

Точки вида $\mathbf{y}_{k-1-(m-1)\tau} = a\mathbf{x}_{k-1-(m-1)\tau} + \mathbf{b}$ из пространства \mathbb{R}^m имеют такое же расстояние (в смысле нормализованного евклидова расстояния) до точки $\mathbf{x}_{M-1-(m-1)\tau}$, как и точка $\mathbf{x}_{k-1-(m-1)\tau}$. Действительно,

$$\begin{aligned} d_{NE}(a\mathbf{x}_{k-1-(m-1)\tau} + \mathbf{b}, \mathbf{x}_{M-1-(m-1)\tau}) &= d_E \left(\frac{a\mathbf{x}_{k-1-(m-1)\tau} + \mathbf{b} - a\bar{\mathbf{x}}_{k-1-(m-1)\tau} - \mathbf{b}}{\|a\mathbf{x}_{k-1-(m-1)\tau} + \mathbf{b} - a\bar{\mathbf{x}}_{k-1-(m-1)\tau} - \mathbf{b}\|}, \frac{\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}}{\|\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}\|} \right) = \\ &= d_E \left(\frac{a(\mathbf{x}_{k-1-(m-1)\tau} - \bar{\mathbf{x}}_{k-1-(m-1)\tau})}{\|a\mathbf{x}_{k-1-(m-1)\tau} - a\bar{\mathbf{x}}_{k-1-(m-1)\tau}\|}, \frac{\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}}{\|\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}\|} \right) = \\ &= d_E \left(\frac{\mathbf{x}_{k-1-(m-1)\tau} - \bar{\mathbf{x}}_{k-1-(m-1)\tau}}{\|\mathbf{x}_{k-1-(m-1)\tau} - \bar{\mathbf{x}}_{k-1-(m-1)\tau}\|}, \frac{\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}}{\|\mathbf{x}_{M-1-(m-1)\tau} - \bar{\mathbf{x}}_{M-1-(m-1)\tau}\|} \right) = d_{NE}(\mathbf{x}_{k-1-(m-1)\tau}, \mathbf{x}_{M-1-(m-1)\tau}). \end{aligned}$$

Обозначим $L(\mathbf{x}_{M-1-(m-1)\tau}) = \{\mathbf{y}_{k-1-(m-1)\tau} = a\mathbf{x}_{k-1-(m-1)\tau} + \mathbf{b} : a \in \mathbb{R} \setminus \{0\}, \mathbf{b} = (b, b, \dots, b) \in \mathbb{R}^m\}$.

Из множества $L(\mathbf{x}_{M-1-(m-1)\tau})$ выбираем точку $\tilde{\mathbf{x}}_M$, ближайшую точку к $\mathbf{x}_{M-1-(m-1)\tau}$ (в смысле евклидова расстояния):

$$\tilde{\mathbf{x}}_M = \tilde{a}\mathbf{x}_{k-1-(m-1)\tau} + \tilde{\mathbf{b}}.$$

$$\tilde{a} = \frac{n \left(\sum_{i=1}^L x_{M-1-(m-i)\tau} \cdot x_{k-1-(m-i)\tau} - \sum_{i=1}^L x_{M-1-(m-i)\tau} \sum_{i=1}^L x_{k-1-(m-i)\tau} \right)}{n \sum_{i=1}^L x_{k-1-(m-i)\tau}^2 - \left(\sum_{i=1}^L x_{k-1-(m-i)\tau} \right)^2};$$

$$\tilde{b} = \frac{\left(\sum_{i=1}^L x_{M-1-(m-i)\tau} \cdot \sum_{i=1}^L x_{k-1-(m-i)\tau}^2 - \sum_{i=1}^L x_{M-1-(m-i)\tau} \cdot x_{k-1-(m-i)\tau} \sum_{i=1}^L x_{k-1-(m-i)\tau} \right)}{n \sum_{i=1}^L x_{k-1-(m-i)\tau}^2 - \left(\sum_{i=1}^L x_{k-1-(m-i)\tau} \right)^2}.$$

Значение $\tilde{x}_M = \tilde{a}x_k + \tilde{b}$

будем рассматривать в качестве прогноза значения x_M временного ряда $\{x_i, i = 1, 2, \dots, N\}$.

Сравнение различных методов прогнозирования цен. Метод ближайших соседей с нормализованным евклидовым расстоянием был применен для прогноза цен на металлы. Прогноз строился на 6 месяцев – июль-декабрь 2010 г. В качестве прогноза значений $x_M, x_{M+1}, \dots, x_{M+5}$ брались значения $\tilde{x}_M, \tilde{x}_{M+1}, \dots, \tilde{x}_{M+5}$.

Интерес представляет сравнение прогнозов, построенных методом ближайших соседей с нормализованным евклидовым расстоянием и другими методами (в частности, методом ближайших соседей с евклидовым расстоянием).

На рис. 1 приведены исторические (март 2001 г – июнь 2010 г) значения очищенных от инфляции цен (в \$/тонну) на стальной г/к рулон, алюминий и цинк. Для этих показателей далее разными методами строились прогнозы на 6 месяцев (июль – декабрь 2010 г).

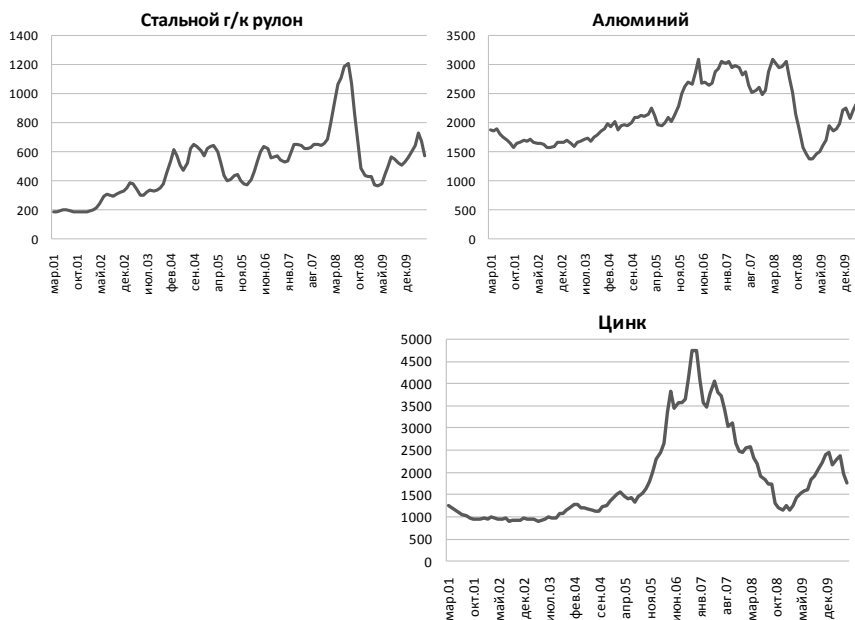


Рис. 1. Динамика цен на различные виды металлов (очищенных от инфляции), март 2001 – июнь 2010, \$/т

Были применены следующие методы прогноза: метод ближайших соседей с использованием евклидова расстояния (БСЕ), и метод ближайших соседей с использованием нормализованного евклидова расстояния (БСНЕ); метод авторегрессии и скользящего среднего (АРСС); метод индексов сезонности (ИС). В случае метода ближайших соседей были построены прогнозы при $m=9$ и $\tau=3$.

Для каждого показателя с помощью перечисленных методов прогнозирования построены по два прогноза: первый, когда значения цен в 2008 г, соответствующие периоду кризиса, не использовались для построения прогноза (с исключением кризиса), и второй, когда использовались (без исключения кризиса). Для метода АРСС в первом случае использовалась модель с интервенциями.

Результаты прогнозирования показаны на графиках (см. рис. 2 – 4).

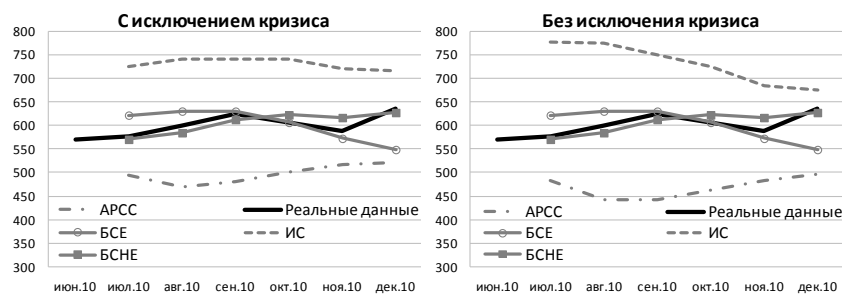


Рис. 2. Реальные значения цены на г/к рулон и прогноз на июль–декабрь 2010 г., полученный методами APCC, BCE, BCNE, IS

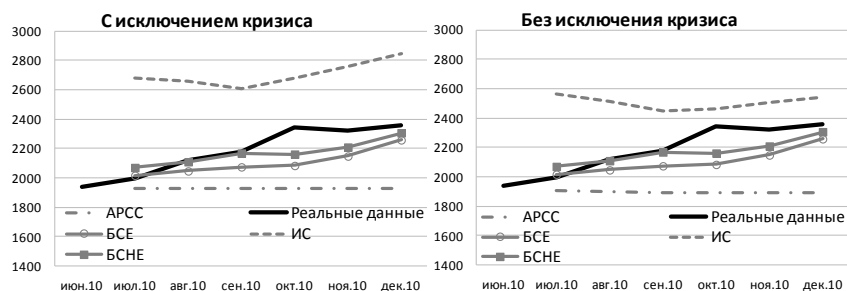


Рис. 3. Реальные значения цены на алюминий и прогноз на июль–декабрь 2010 г., полученный методами APCC, BCE, BCNE, IS

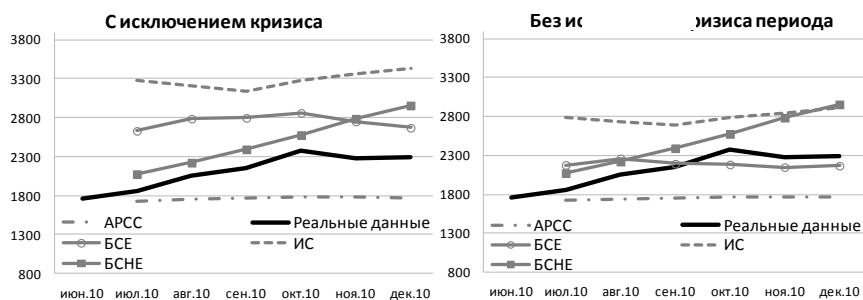


Рис. 4. Реальные значения цены на цинк и прогноз на июль–декабрь 2010 г., полученный методами APCC, BCE, BCNE, IS

В табл. 1 приведены средние погрешности прогнозирования цен:

$$\Delta_1 = 100\% \cdot \frac{1}{6} \sum_{i=0}^5 (\tilde{x}_{M+i} - x_{M+i}) / x_{M+i};$$

и средние абсолютные погрешности:

$$\Delta_2 = 100\% \cdot \frac{1}{6} \sum_{i=0}^5 |\tilde{x}_{M+i} - x_{M+i}| / x_{M+i}.$$

Таблица 1

Погрешности прогнозов, полученных различными методами

Показатель	Метод	Погрешность			
		С исключением кризиса		Без исключения кризиса	
		Δ_1	Δ_2	Δ_1	Δ_2
Цена на г/к рулон	BCE	-0,4 %	5 %	-0,4 %	5 %
	BCNE	0,1 %	2 %	0,1 %	2 %
	ИС	21 %	21 %	21 %	21 %
	APCC	-18 %	18 %	-23 %	23 %
Цена на алюминий	BCE	-5 %	5 %	-5 %	5 %
	BCNE	-2 %	3 %	-2 %	3 %
	ИС	22 %	22 %	13 %	13 %
	APCC	-13 %	13 %	-14 %	14 %
Цена на цинк	BCE	28 %	28 %	2 %	8 %

	БСНЕ	15 %	15 %	15 %	15 %
	ИС	52 %	52 %	30 %	30 %
	АРСС	- 18 %	18 %	- 19 %	19 %

Во всех рассмотренных случаях лучшие прогнозы давал метод ближайших соседей с использованием нормализованного евклидова расстояния.

Замечание. Предложенный метод также апробировался для прогноза цен на г/к рулон на периоды продолжительностью в 6 мес. начиная с периода [июнь 09, июль 09,..., дек. 09] и заканчивая периодом [июнь 10, июль 10,..., дек. 10]. При этом погрешности имеют такой же порядок, как и приведенные в таблице.

Выводы. В статье показаны преимущества метода ближайших соседей с нормализованным евклидовым расстоянием при прогнозировании нестационарных временных рядов на примерах нестационарных рядов цен на металлы и металлопродукцию.

Библиографические ссылки

1. **Зульпукаров М.-Г.М.** Метод русел и джокеров на примере исследования системы Розенцвейга–Макартура. / М.-Г.М. Зульпукаров, Г.Г. Малинецкий, А.В. Подлазов // Институт прикладной математики им. М.В. Келдыша РАН (препринт) – 2006. – № 21.
2. **Лоскутов А.Ю.** Временные ряды: анализ и прогноз. / А.Ю. Лоскутов, О.Л. Котляров, Д.И. Журавлев // Математика. Компьютер. Образование. Сборник трудов XI международной конференции. – Ижевск, 2004. – С. 9-46.
3. **Малинецкий Г.Г.** Современные проблемы нелинейной динамики. / Г.Г. Малинецкий, А.Б. Потапов. – М., 2000.
4. **Малинецкий Г.Г.** Руслу и джокеры: о новых методах прогноза поведения сложных систем. / Г.Г. Малинецкий, А.Б. Потапов. // ИПМ им. М.В.Келдыша РАН (препринт). – 1998. – № 32.
5. **Комарова Н.В.** Некоторые инструменты нелинейной динамики в анализе и прогнозировании временных рядов. / Н.В. Комарова // Питання прикладної математики і математичного моделювання. Зб. наукових праць. – Д., 2011. – С. 191-199.

Надійшла до редколегії 08.06.2012