

Анализ и прогнозирование гидрологических данных Дипломная работа

Александр Сергеевич Павлов Научный руководитель: Цеховая Татьяна Вячеславовна

Факультет прикладной математики и информатики
Кафедра теории вероятностей и математической статистики

Минск, 2015

Содержание



Постановка задачи

Обзор реализованного программного обеспечения Модуль предварительного анализа Модуль анализа остатков Модуль вариограммного анализа

Детерминированный подход Проверка на нормальность Корреляционный анализ Регрессионный анализ Анализ остатков

Геостатистический подход Вариограммный анализ Автоматический подход

Заключение

Постановка задачи



- 1. Предварительный статистический анализ гидроэкологических данных озера Баторино;
- 2. Вариограммный анализ временного ряда: построение оценок семивариограммы, подбор моделей семивариограммы;
- 3. Исследование статистических свойств оценки семивариограммы гауссовского случайного процесса;
- 4. Прогнозирование значений временного ряда с помощью интерполяционного метода Кригинг;
- 5. Исследование точности прогноза в зависимости от оценки семивариограммы и модели семивариограммы, лежащих в основе метода Кригинг.

Обзор реализованного программного обеспечения Особенности

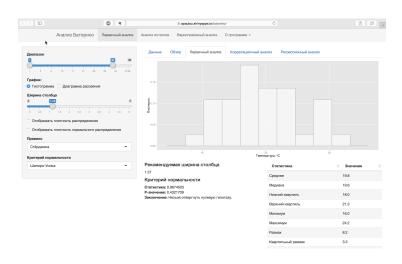


- Доступно с любого устройства, имеющего доступ в Интернет, по адресу apaulau.shinyapps.io/batorino;
- Реализовано на языке программирования R;
- Логически разделено на три модуля;
- Имеет простой, быстро расширяемый гибкий интерфейс;
- Широкие графические возможности;
- Проверка тестов и критериев;
- Мгновенный отклик на изменение параметров;
- Быстрая проверка различных моделей.

Модуль предварительного анализа



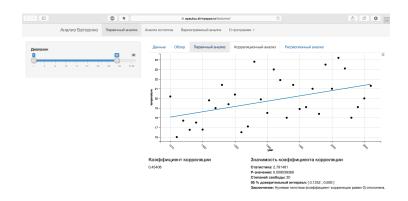
Первичный анализ и описательные статистики



Модуль предварительного анализа



Корреляционный анализ

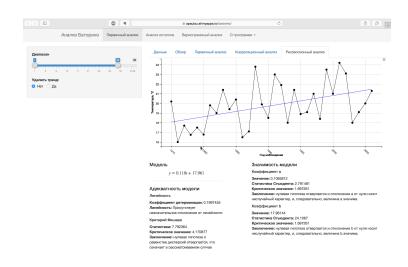


Минск, 2015

Модуль предварительного анализа



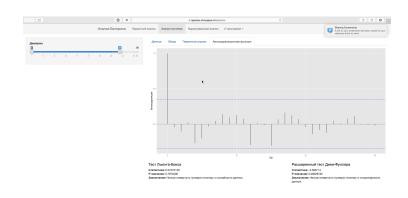
Регрессионный анализ



Модуль анализа остатков



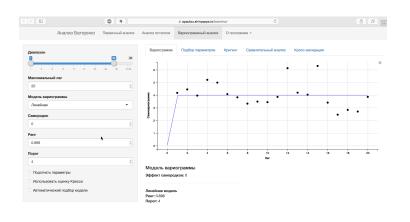
Автокорреляционная функция



Модуль вариограммного анализа



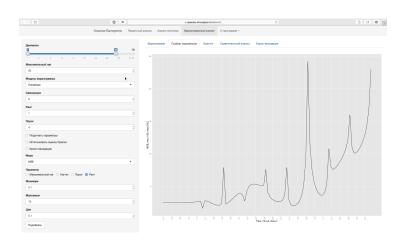
Возможности по подбору модели вариограммы



Модуль вариограммного анализа



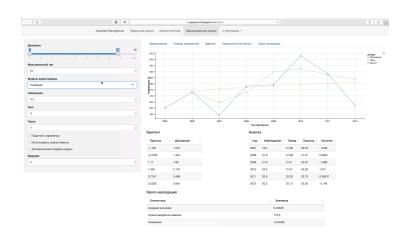
Подбор параметров модели вариограммы



Модуль вариограммного анализа



Сравнение прогнозных значений





Исходные данные

Данные получены от учебно-научного центра «Нарочанская биологическая станция им. Г.Г.Винберга».

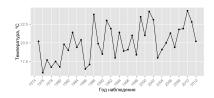


Рис.: Исходные данные

Исходные данные: выборка X(t), $t=\overline{1,n}$, n=38, X(t) — значение средней температуры воды оз. Баторино в июле месяце для каждого года в период с 1975 по 2012 годы.



Проверка на нормальность

- Коэффициент асимметрии 0.30 ⇔ распределение скошено вправо;
- Коэффициент эксцесса $-0.746 \Leftrightarrow$ пик кривой распределения пологий относительного нормального.

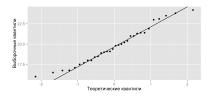


Рис.: График квантилей

Выборочное распределение близко к нормальному $\mathcal{N}(19.77,5.12)$ (визуально, критерии Шапиро-Уилка, χ^2 -Пирсона и Колмогорова-Смирнова).



Корреляционный анализ

- Выбросы в исходных данных отсутствуют (критерий Граббса);
- Выборочный коэффициент корреляции $r_{xt} = 0.454$ при уровне значимости $\alpha = 0.05$ является значимым.

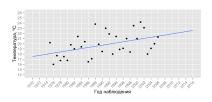


Рис.: Диаграмма рассеяния



Регрессионный анализ: регрессионная модель

Исследуемый временной ряд является аддитивным,

$$X(t) = y(t) + \varepsilon(t);$$
 (1)

y(t) — тренд, $\varepsilon(t)$ — нерегулярная составляющая.

Найденная модель тренда:

$$y(t) = at + b = 0.1014t + 18.0521.$$

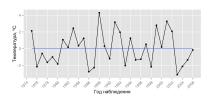


Рис.: Ряд остатков $\varepsilon(t)$



Регрессионный анализ: качество регрессионной модели

	X(t)	<i>y</i> (<i>t</i>)	X(t) - y(t)
2007	19.400	18.071	1.329
2008	21.800	18.181	3.619
2009	21.900	18.290	3.610
2010	24.300	18.400	5.900
2011	22.800	18.509	4.291
2012	20.200	18.619	1.581

Таблица: Сравнение прогнозных значений (модель y(t))

- Коэффициенты регрессионной модели значимы (критерий Стьюдента, $\alpha=0.05$);
- Модель адекватна (F-критерий Фишера, $\alpha = 0.05$);
- Точность модели невысока (поскольку коэффициент детерминации $\eta_{\mathbf{x}(t)}^2 = 0.275$).



Анализ остатков

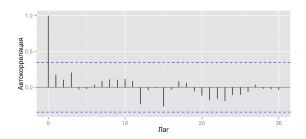


Рис.: Автокорреляционная функция

- Выборочное распределение близко к $\mathcal{N}(0.00, 4.07)$;
- Значимые автокорреляции отсутствуют;
- Значения имеют небольшую амплитуду и имеют тенденцию к затуханию ⇔ ряд стационарен в широком смысле.



Оценка вариограммы

Пусть $X(t),\ t\in\mathbb{Z}$ — стационарный в широком смысле гауссовский случайный процесс с дискретным временем, нулевым математическим ожиданием и постоянной дисперсией.

Определение 1

Вариограмма случайного процесса $X(t), t \in \mathbb{Z}$:

$$2\gamma(h) = V\{X(t+h) - X(t)\}, t, h \in \mathbb{Z}.$$
 (2)

При этом функция $\gamma(h), h \in \mathbb{Z}$, называется семивариограммой.

Оценка вариограммы (Матерон):

$$2\tilde{\gamma}(h) = \frac{1}{n-h} \sum_{t=1}^{n-h} (X(t+h) - X(t))^2, \quad h = \overline{0, n-1}.$$
 (3)



Первые два момента оценки вариограммы

Теорема 2

Для оценки $2\tilde{\gamma}(\mathbf{h})$ имеют место следующие соотношения:

$$\mathbf{E}\{2\tilde{\gamma}(\mathbf{h})\} = 2\gamma(\mathbf{h}),$$

$$\begin{split} \cos\!v(2\tilde{\gamma}(\textit{h}_1),2\tilde{\gamma}(\textit{h}_2)) &= \frac{2}{(\textit{n}-\textit{h}_1)(\textit{n}-\textit{h}_2)} \sum_{t=1}^{\textit{n}-\textit{h}_1} \sum_{\textit{s}=1}^{\textit{n}-\textit{h}_2} (\gamma(\textit{t}-\textit{h}_2-\textit{s}) + \\ &+ \gamma(\textit{t}+\textit{h}_1-\textit{s}) - \gamma(\textit{t}-\textit{s}) - \gamma(\textit{t}+\textit{h}_1-\textit{s}-\textit{h}_2))^2, \\ V\{2\tilde{\gamma}(\textit{h})\} &= \frac{2}{(\textit{n}-\textit{h})^2} \sum_{\textit{t},\textit{s}=1}^{\textit{n}-\textit{h}} (\gamma(\textit{t}-\textit{h}-\textit{s}) + \gamma(\textit{t}+\textit{h}-\textit{s}) - 2\gamma(\textit{t}-\textit{s}))^2, \end{split}$$

где $\gamma(h)$, — семивариограмма процесса X(t), $h,h_1,h_2=\overline{0,n-1}$.



Асимптотическое поведение оценки вариограммы

Теорема 3

Если имеет место соотношение $\sum_{h=-\infty}^{+\infty} |\gamma(h)| < +\infty$, то

$$\lim_{n\to\infty}(\mathbf{n}-\min\{\mathbf{h}_1,\mathbf{h}_2\})\mathbf{cov}\{2\tilde{\gamma}(\mathbf{h}_1),2\tilde{\gamma}(\mathbf{h}_2)\}=2\sum_{\mathbf{m}=-\infty}^{+\infty}\gamma(\mathbf{m}-\mathbf{h}_2)+$$

$$\gamma(\textit{\textit{m}}+\textit{\textit{h}}_1) - \gamma(\textit{\textit{m}}) - \gamma(\textit{\textit{m}}+\textit{\textit{h}}_1-\textit{\textit{h}}_2))^2, \ +\infty$$

$$\lim_{\mathbf{n}\to\infty}(\mathbf{n}-\mathbf{h})\mathsf{V}\{2\tilde{\gamma}(\mathbf{h})\}=2\sum_{\mathbf{m}=-\infty}^{+\infty}\gamma(\mathbf{m}-\mathbf{h})+\gamma(\mathbf{m}+\mathbf{h})-2\gamma(\mathbf{m}))^2.$$

20 / 32



Асимптотическое поведение оценки вариограммы

Следствие 4

Из теоремы 2 следует соотношение

$$\lim_{{\bf n}\to\infty} {\bf V}\{2\tilde{\gamma}({\bf h})\}=0, \quad {\bf h}=\overline{0,{\bf n}-1}.$$

Следствие 5

В силу показанной в теореме 1 несмещённости оценки и вышеприведённого следствия получаем, что оценка вариограммы $2\tilde{\gamma}(h)$ является состоятельной в среднеквадратическом смысле для вариограммы $\gamma(h), h \in \mathbb{Z}.$



График экспериментальной вариограммы

Прогнозные значения $X^*(t)$ вычисляются по формуле:

$$X^*(t) = y(t) + \varepsilon^*(t),$$

где y(t) — тренд, $\varepsilon^*(t)$ — значения, вычисленные с помощью кригинга.

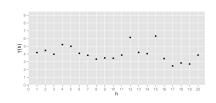


Рис.: Оценка семивариограммы Матерона

Для оценки качества модели используются

- коэффициент корреляции $r_{\varepsilon\varepsilon^*}$;
- Среднеквадратическая ошибка (*n* объём выборки):

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (\varepsilon(t_i) - \varepsilon^*(t_i))^2.$$
 (4)



Линейная модель

Общий вид модели:

$$\widehat{\gamma}(h) = \mathbf{c}_0 + \mathbf{Lin}(h) =$$

$$= \begin{cases} \mathbf{c}_0 + \mathbf{b} \cdot \mathbf{h}, & h > 0, \\ \mathbf{c}_0, & h \le 0, \end{cases}$$
 (5)

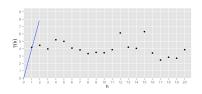
где b – параметр, отвечающий за угол наклона, c_0 — эффект самородков.

Подобранная модель:

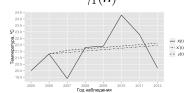
$$\widehat{\gamma}_1(\mathbf{h}) = \mathbf{Lin}(\mathbf{h}), \quad \mathbf{b} = 4$$
 (6)

Показатели качества

$$r_{\varepsilon\varepsilon^*} = -0.09129$$
, $MSE = 6.324$



Модель семивариограммы $\widehat{\gamma}_1(h)$



Прогноз по модели $\widehat{\gamma}_1(h)$



Чистый эффект самородков

Общий вид модели:

$$\widehat{\gamma}(h) = \mathbf{c} \cdot \mathsf{Nug}(h) = = \begin{cases} 0, & h = 0, \\ \mathbf{c}, & h \neq 0, \end{cases}$$
(7)

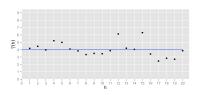
где b – параметр, отвечающий за угол наклона, c_0 — эффект самородков.

Подобранная модель:

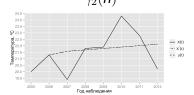
$$\widehat{\gamma}_2(\mathbf{h}) = 4.04 \cdot \mathbf{Nug}(\mathbf{h}).$$
 (8)

Показатели качества

$$r_{\varepsilon \varepsilon^*} = -1$$
, $MSE = 4.199$



Модель семивариограммы $\widehat{\gamma}_2(h)$



Прогноз по модели $\widehat{\gamma}_2(h)$

БЕЛАРУСКІ ДЗЯРЖАЎНЫ ЎНІВЕРСІТЭТ

Линейная модель с порогом

Общий вид модели:

$$\widehat{\gamma}(h) = \mathbf{c}_0 + \mathbf{c} \cdot Lin(h, \mathbf{a}) =$$

$$= \begin{cases} \mathbf{c}_0 + \mathbf{c} \cdot \frac{h}{\mathbf{a}}, & 0 \le h \le \mathbf{a}, \\ \mathbf{c}_0 + \mathbf{c}, & h > \mathbf{a}, \end{cases}$$
(9)

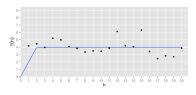
где c_0 — эффект самородков, c — порог, a — ранг.

Подобранная модель:

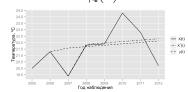
$$\widehat{\gamma}_4(\mathbf{h}) = 4 \cdot \mathbf{Lin}(\mathbf{h}, 2).$$
 (10)

Показатели качества

$$r_{\varepsilon\varepsilon^*} = 0.152$$
, $MSE = 18.69$



Модель семивариограммы $\widehat{\gamma}_4(h)$



Прогноз по модели $\widehat{\gamma}_4(h)$



Сферическая модель

Общий вид модели:

$$\begin{split} \widehat{\gamma}(\textbf{\textit{h}}) &= \textbf{\textit{c}}_0 + \textbf{\textit{c}} \cdot \textbf{\textit{Sph}}(\textbf{\textit{h}}, \textbf{\textit{a}}) = \\ &= \left\{ \begin{array}{ll} \textbf{\textit{c}}_0 + \textbf{\textit{c}} \cdot (\frac{3}{2}\frac{\textbf{\textit{h}}}{\textbf{\textit{a}}} - \frac{1}{2}(\frac{\textbf{\textit{h}}}{\textbf{\textit{a}}})^3), & \textbf{\textit{h}} \leq \textbf{\textit{a}}, \\ \textbf{\textit{c}}_0 + \textbf{\textit{c}}, & \textbf{\textit{h}} \geq \textbf{\textit{a}}, \\ & \textbf{\textit{(11)}} \end{array} \right. \end{split}$$

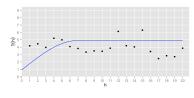
где c_0 – эффект самородков, c – порог, a – ранг.

Подобранная модель:

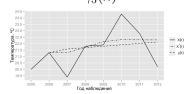
$$\widehat{\gamma}_{5}(\mathbf{h}) = 0.9 + 4 \mathbf{Sph}(\mathbf{h}, 6.9),$$
 (12)

Показатели качества

$$r_{\varepsilon\varepsilon^*} = -0.009$$
, $MSE = 5.396$



Модель семивариограммы $\widehat{\gamma}_5(h)$



Прогноз по модели $\widehat{\gamma}_5(h)$



Периодическая модель

Общий вид модели:

$$\widehat{\gamma}(\mathbf{h}) = \mathbf{c}_0 + \mathbf{c} \cdot \mathbf{Per}(\mathbf{h}, \mathbf{a}) = \mathbf{cos}(\frac{2\pi\mathbf{h}}{\mathbf{a}}),$$

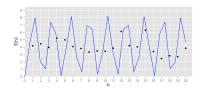
где c_0 — эффект самородков, c — порог, a — ранг.

Подобранная модель:

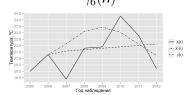
$$\widehat{\gamma}_{6}(\mathbf{h}) = 4 \cdot \textit{Per}(\mathbf{h}, 0.898),$$
 (14)

Показатели качества

$$r_{\varepsilon \varepsilon^*} = 0.404$$
, $MSE = 4.369$



Модель семивариограммы $\widehat{\gamma}_6(h)$



Прогноз по модели $\widehat{\gamma}_6(\mathbf{h})$

Автоматический подход

БЕЛАРУСКІ ДЗЯРЖАЎНЫ ЎНІВЕРСІТЭТ

Волновая модель

Общий вид модели:

$$\widehat{\gamma}(\mathbf{h}) = \mathbf{c}_0 + \mathbf{c} \cdot \mathbf{Wav}(\mathbf{h}, \mathbf{a}) = \mathbf{1} - \frac{\mathbf{a}}{\mathbf{h}} \cdot \mathbf{sin}(\frac{\mathbf{h}}{\mathbf{a}}),$$

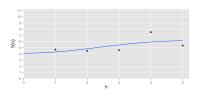
где c_0 — эффект самородков, c — порог, a — ранг.

Подобранная модель:

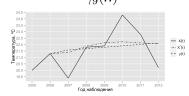
$$\widehat{\gamma}_{9}(\mathbf{h}) = 4.11 + 1.65 \cdot \mathbf{Wav}(\mathbf{h}, 3.59),$$
 (16)

Показатели качества

$$r_{\varepsilon\varepsilon^*} = -1$$
, $MSE = 4.20$



Модель семивариограммы $\widehat{\gamma}_9(h)$



Прогноз по модели $\widehat{\gamma}_9(\mathbf{h})$

Автоматический подход



Периодическая модель

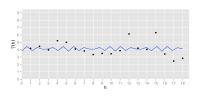
Модель семивариограммы вида (13).

Подобранная модель:

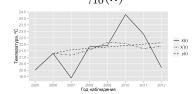
$$\widehat{\gamma}_{10}(\emph{h}) = 3.8 + 0.32 \cdot \emph{Per}(\emph{h}, 1.3)$$
 (17)

Показатели качества

$$r_{\varepsilon \varepsilon^*} = -0.15, \quad \textit{MSE} = 5.22$$



Модель семивариограммы $\widehat{\gamma}_{10}(h)$



Прогноз по модели $\widehat{\gamma}_{10}(\emph{h})$

Заключение



Список использованных источников





Cressie N.

Statistics for Spatial Data.

New York. — Wiley, 1991.



А.А. Савельев, С.С. Мухарамова, А.Г. Пилюгин, Н.А. Чижикова Геостатистический анализ данных в экологии и природопользовании (с применением пакета R)

Казань: Казанский университет, 2012.



Robert H. Shumway, David S. Stoffer Time series and Its Applications: With R Examples (Springer Texts in Statistics). Springer Science+Business Media, LLC 2011, 3d edition, 2011.



Paul Teetor

R Cookbook (O'Reilly Cookbooks)).

O'Reilly Media, 1 edition, 2011.

Спасибо за внимание!



Анализ и прогнозирование гидрологических данных

Александр Сергеевич Павлов Научный руководитель: Цеховая Татьяна Вячеславовна

Факультет прикладной математики и информатики Кафедра теории вероятностей и математической статистики

Минск, 2015