

Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Владимирский государственный университет имени Александра
Григорьевича и Николая Григорьевича Столетовых»
(ВлГУ)
Кафедра ИЗИ

ОТЧЕТ
по лабораторной работе №3
по дисциплине «МАД»
ТЕМА: «ПРЕДОБРАБОТКА ДАННЫХ»

Студент гр. ИСБ-121 —

Ходыкин Л.Ю.

Преподаватель —

Артюшина Л.А.

Владимир
2024

Оглавление

Задание.	3
Ход работы 1 задание.....	4
Ход работы 2 задание.....	14

Задание.

Задание:

- 1.1 Используя текстовый редактор Word создать файл
- 1.2 Импортируйте данные из созданного файла в Deductor.
- 1.3 Восстановите пропущенные данные в столбце «Синус»
- 1.4 Удалите аномалии в данных
- 1.5 Сгладьте данные методом спектральной обработки
- 1.6 Удалите шумы из данных
- 2.1 импортировать данные из файла в Deductor;
- 2.2 используя инструмент Филт, выделите из получившейся таблицы страны, согласно индивидуальному варианту (по 2 столбца таблицы);
- 2.3 визуализируйте данные;
- 2.4 при наличии заполните пропуски в данных. Обоснуйте выбранный способ заполнения;
- 2.5 визуализируйте данные после заполнения пропусков;
- 2.6 удалите аномалии в данных;
- 2.7 сгладьте данные одним из методов спектральной обработки. Обоснуйте выбор метода;
- 2.8 визуализируйте данные после удаления аномалий;
- 2.9 обоснуйте наличие шумов в данных. Удалите шумы из данных (при наличии);
- 2.10 Создайте многомерный отчет по количеству заболевших в виде OLAP-куба, показатели выберите самостоятельно (например, общее число заболевших по странам и т.п.).

Ход работы 1 задание.

Создаем файл со столбцами – Аргумент, Синус, Аномалии, Большие Шумы, Средние Шумы, Малые Шумы.

1	Argument	Sinus	Anomaly	Large Noise	Medium Noise	Small Noise
2	0.000000000	0.000000000	0.079914694	0.023834934	-0.049147772	-0.008032414
3	0.020000000	0.019998667	0.019998667	0.035346857	-0.023802290	0.010271700
4	0.040000000	0.039989334	0.059964006	0.134215134	0.101732203	0.056760332
5	0.060000000	0.059964006	0.059964006	0.064336934	0.036772284	0.040865458
6	0.080000000	0.079914694	0.079914694	-0.064865690	0.080033216	0.071900237
7	0.100000000	0.099833417	0.099833417	0.102934196	0.030902416	0.106550496
8	0.120000000	0.119712207	0.119712207	0.076163088	0.179586020	0.119247206
9	0.140000000	0.139543115	0.139543115	0.146402760	0.151334605	0.132139887
10	0.160000000	0.159318207	0.159318207	0.038418294	0.213413657	0.170795312
11	0.180000000	0.179029573	0.179029573	0.292217395	0.138959784	0.172577694
12	0.200000000	0.198669331	0.198669331	0.104018445	0.127732180	0.207973849
13	0.220000000	0.218229623	0.198669331	0.282018452	0.262107705	0.217532104
14	0.240000000	0.237702626	0.237702626	0.414769226	0.227303021	0.241287057
15	0.260000000	0.257080552	0.257080552	0.285008637	0.226875528	0.274392074
16	0.280000000	0.276355649	0.276355649	0.232482808	0.247364849	0.272814433
17	0.300000000	0.295520207	0.295520207	0.200311542	0.221370401	0.295651797
18	0.320000000	0.314566561	0.314566561	0.298577344	0.304446289	0.318777560
19	0.340000000	0.333487092	0.333487092	0.457663415	0.357547962	0.329305261
20	0.360000000	0.352274233	0.352274233	0.248535588	0.379479480	0.353890904
21	0.380000000	0.370920469	0.370920469	0.371209442	0.298687001	0.374706208
22	0.400000000	0.389418342	0.352274233	0.482408129	0.403232394	0.379481870
23	0.420000000	0.407760453	0.407760453	0.426426145	0.420541009	0.408829429
24	0.440000000	0.425939465	0.425939465	0.357989956	0.451032228	0.436068216
25	0.460000000	0.443948107	0.461779176	0.571466869	0.565075088	0.430901859
26	0.480000000	0.461779176	0.461779176	0.407071183	0.496868321	0.476819297
27	0.500000000	0.479425539	0.479425539	0.551837898	0.490657268	0.486050577
28	0.520000000	0.496880138	0.496880138	0.465588559	0.485270633	0.498147726
29	0.540000000	0.514135992	0.514135992	0.278569566	0.478830734	0.517984159
30	0.560000000	0.531186198	0.531186198	0.555910496	0.573211965	0.527753621
31	0.580000000	0.548023937	0.548023937	0.500921765	0.590447079	0.538078755
32	0.600000000	0.564642473	0.564642473	0.514768949	0.554359703	0.560216576
33	0.620000000	0.581035161	0.581035161	0.608518813	0.635938262	0.575216029
34	0.640000000	0.597195441	0.597195441	0.381441989	0.615273238	0.593728733
35	0.660000000	0.613116852	0.613116852	0.810053369	0.644584710	0.611191840
36	0.680000000	0.628793024	0.564642473	0.607654046	0.572206817	0.643697066
37	0.700000000	0.644217687	0.644217687	0.719469146	0.624530389	0.637676888
38	0.720000000	0.659384672	0.659384672	0.514273191	0.604166317	0.664693432
39	0.740000000	0.674287912	0.674287912	0.674762174	0.654361139	0.685770217
40	0.760000000	0.688921445	0.688921445	0.719196803	0.676934425	0.706751512
41	0.780000000	0.703279419	0.703279419	0.633862940	0.736795667	0.691378332
42	0.800000000	0.717356091	0.717356091	0.701852718	0.694907810	0.703719691
43	0.820000000	0.731145830	0.731145830	0.693102708	0.708464384	0.730554219
44	0.840000000	0.744643120	0.744643120	0.753082864	0.786260804	0.737966239
45	0.860000000	0.757842563	0.757842563	0.795911547	0.864086706	0.747049789
46	0.880000000	0.770738879	0.770738879	0.753916689	0.727164734	0.764832918
47	0.900000000	0.783238666	0.783238666	0.700000000	0.700000000	0.700000000

Далее загружаем данный файл и визуализируем по синусу и аргументу



Дальше заполняем пропуски с помощью заполнения пропусков в мастере обработки

Мастер обработки - Заполнение пропусков (1 из 5)

Заполнение пропущенных данных
Общие настройки набора данных

☐ Использовать информацию узла оценки качества данных

☒ Обработать как упорядоченный набор

Максимально допустимый процент пропущенных данных: 50

< Назад Далее > Отмена

Нам необходимо только синус и аргумент интерполировать. Используется интерполяция потому что данные упорядочены.

Мастер обработки - Заполнение пропусков (2 из 5)

Заполнение пропущенных данных
Настройка отдельных полей (упорядоченный набор данных)

✓ Argument
✓ Sinus
i Anomaly
i Large Noise
i Medium Noise
i Small Noise

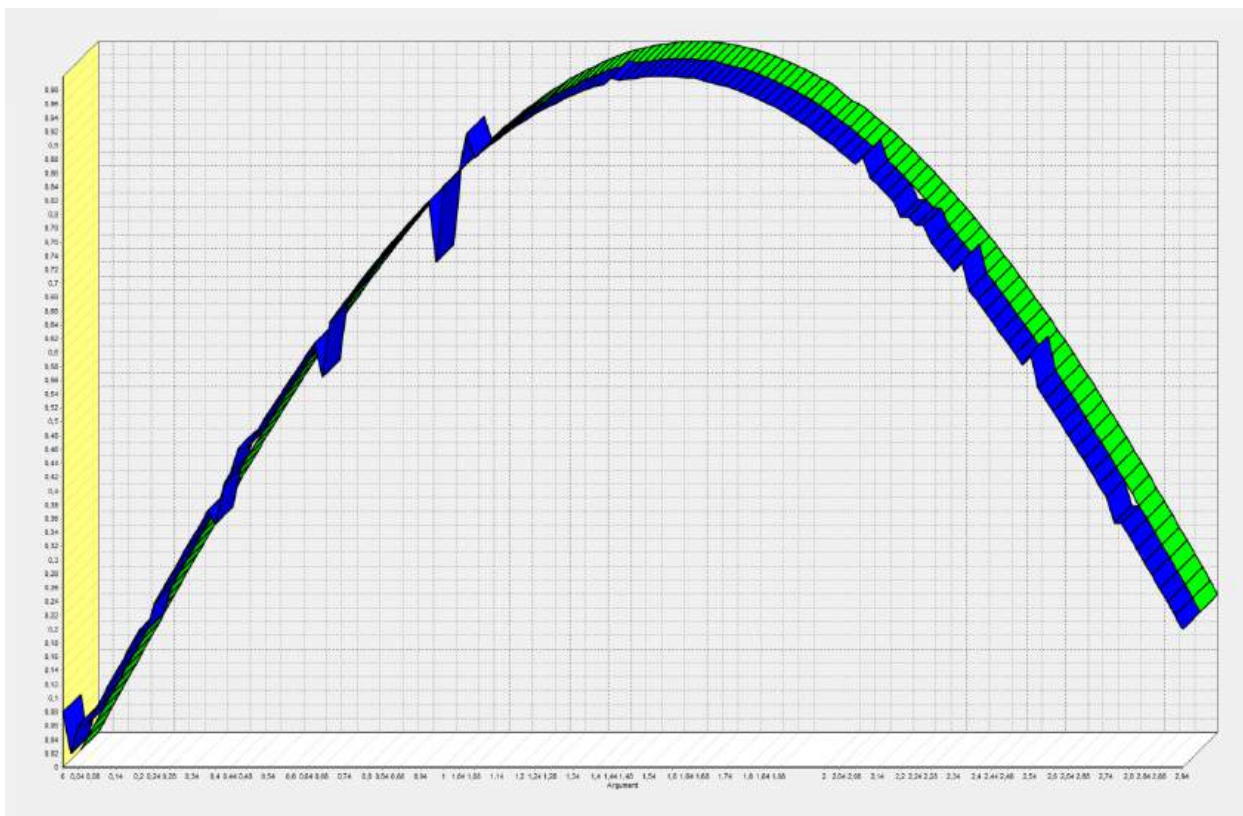
Тип данных: Вещественный
Вид данных: Непрерывный
Назначение: ☒ Используемое

Выбор метода обработки пропущенных данных:
Интерполировать

Описание метода
Вместо пропущенных данных будут вставляться значения, полученные интерполяцией исходного временного ряда (только для вещественных непрерывных данных в упорядоченных наборах)

< Назад Далее > Отмена

Получаем вот такую обновленную диаграмму для синуса (зеленая) и диаграмму для аномалий(синяя)



Далее удаляем аномалии. В данной версии дедуктора это делается через редактирование выбросов и экстремальных значений в мастере обработки

Мастер обработки - Редактирование выбросов (1 из 5)

Редактирование выбросов и экстремальных значений
Общие настройки набора данных

☐ Использовать информацию узла оценки качества данных

☒ Обращивать как упорядоченный набор данных

Метод определения выбросов и экстремальных значений

☒ Стандартное отклонение от среднего

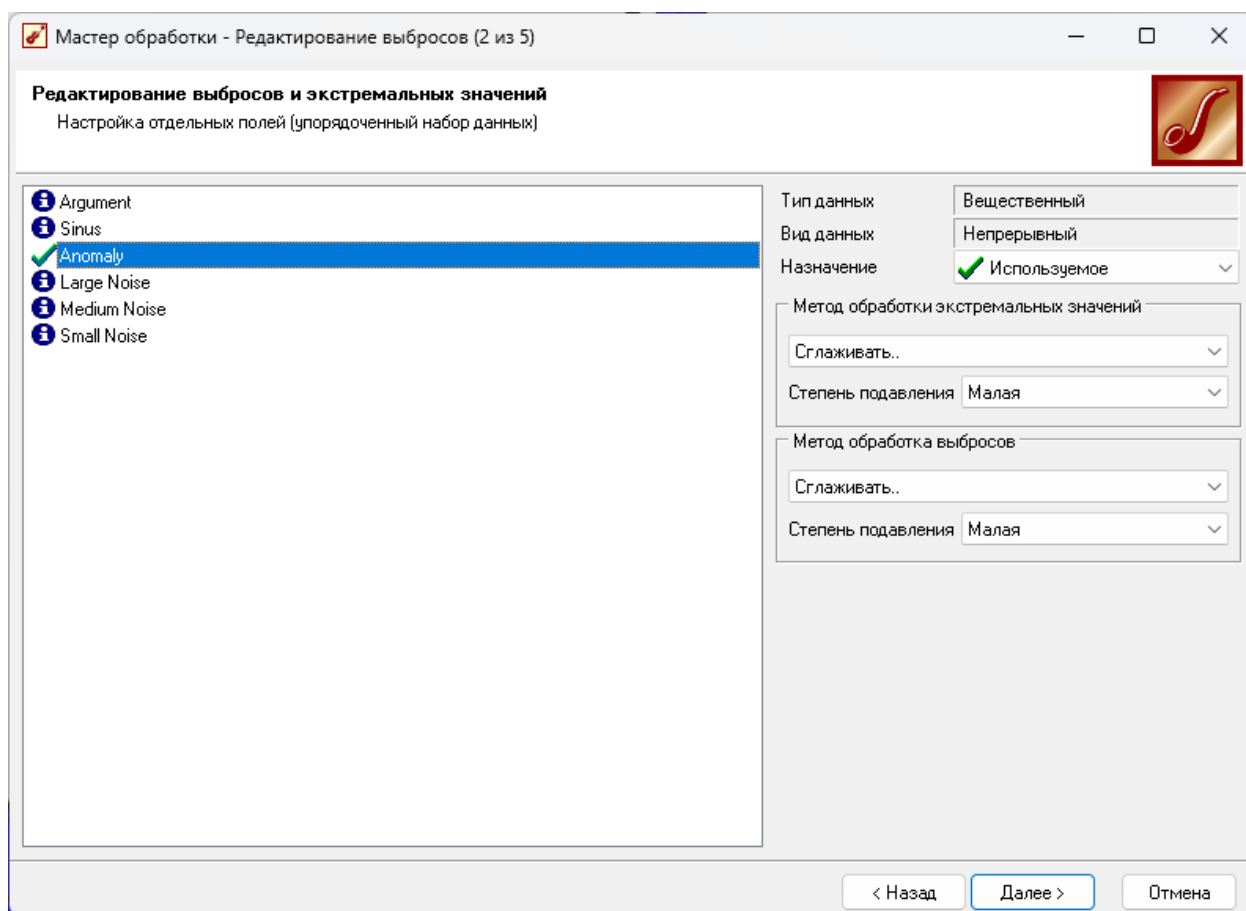
Для выбросов: 3,0 Для экстремальных значений: 5,0

☐ Интерквартильная ширина

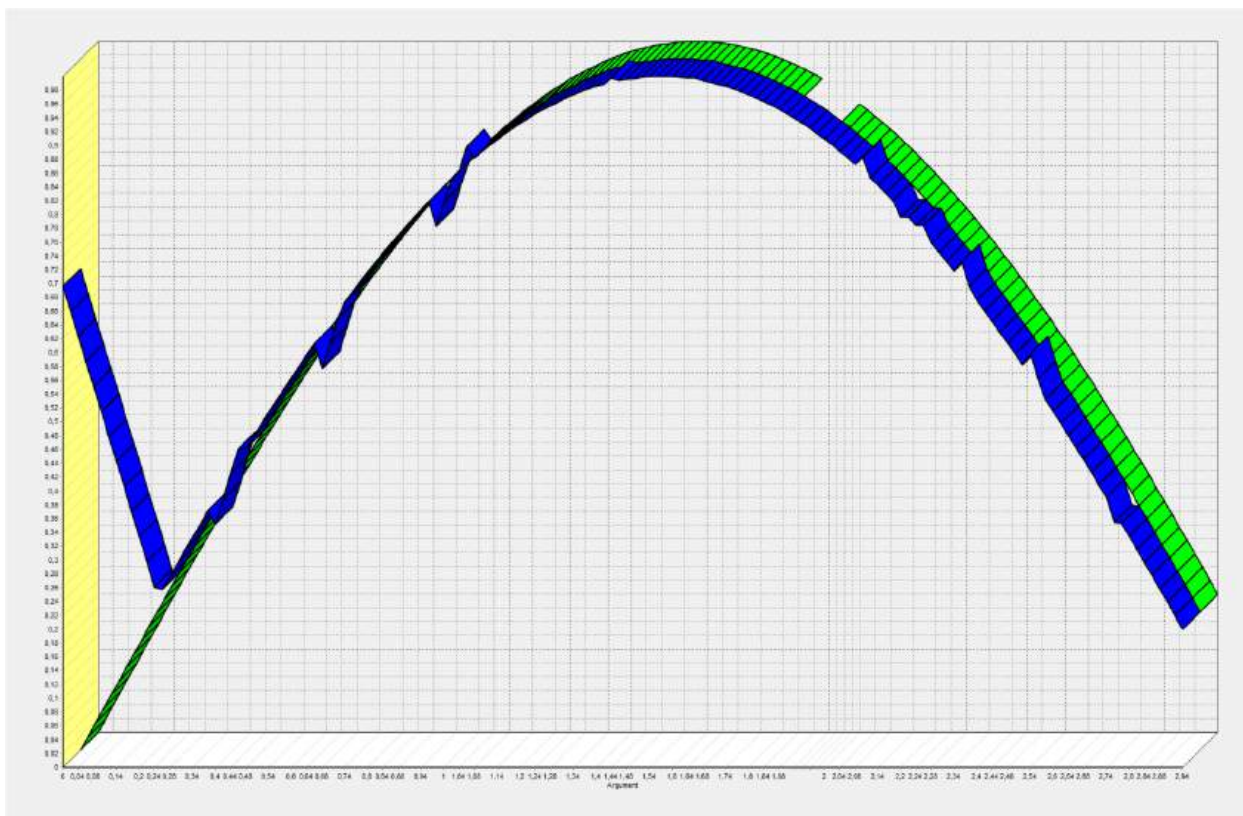
Для выбросов: 1,5 Для экстремальных значений: 3,0

< Назад Далее > Отмена

Далее выбираем нужный столбец в нашем случае аномалии выбираем метод обработки – сглаживать и степень малая

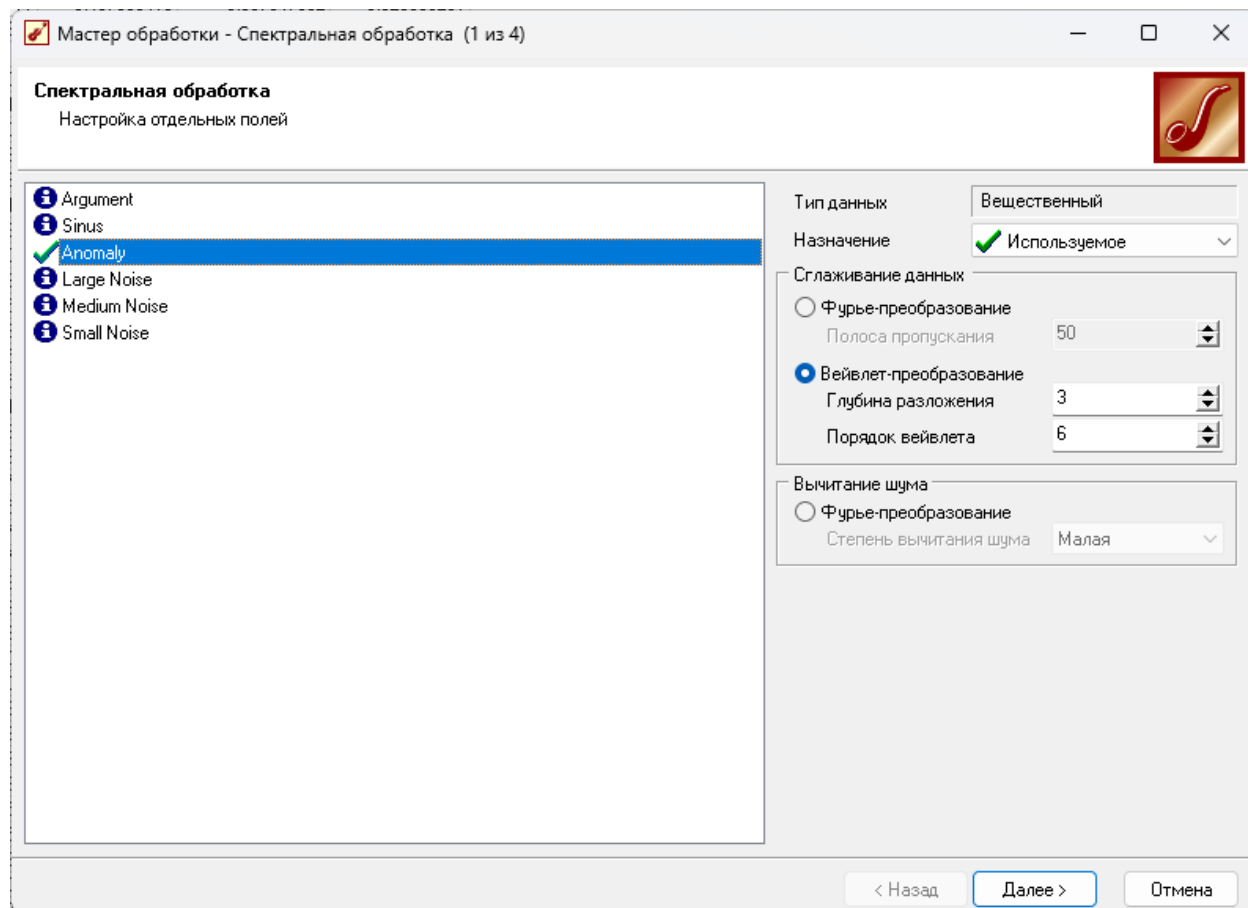


После чего получаем вот такую диаграмму со сглаженными экстремальными значениями по столбцу аномалии(синий)

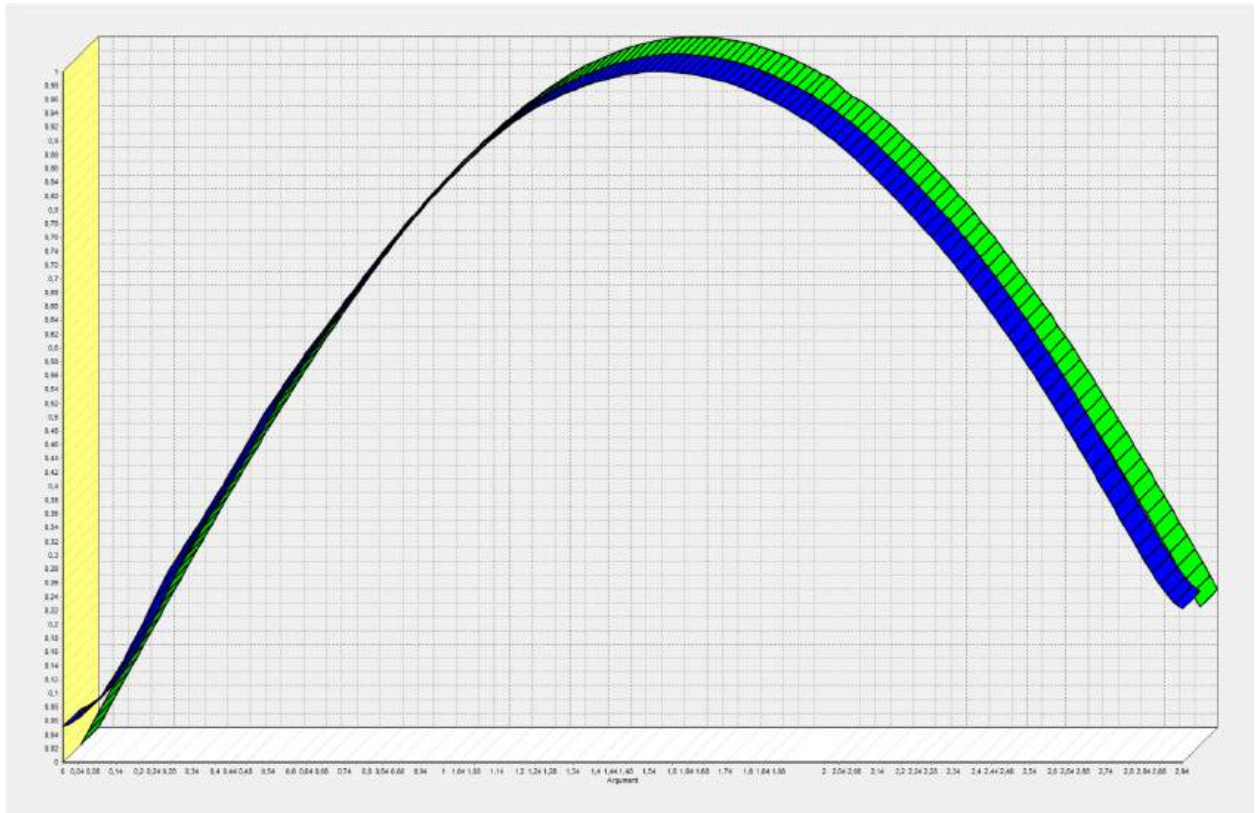


Далее необходимо удалить сгладить данные для аномалий. В данной версии дедуктора используется метод спектральной обработки в мастере обработки.

Выбираем нужный нам столбец и выбираем вейвлет-преобразование.



После чего получаем вот такую диаграмму



Далее необходимо удалить шумы. Шумы начинаем удалять от больших к меньшим. В данной версии дедуктора используется метод спектральной обработки в мастере обработки. Выбираем нужную колонку и метод вычитание шума, фурье-преобразование со степенью Большая.

Мастер обработки - Спектральная обработка (1 из 4)

Спектральная обработка

Настройка отдельных полей

- Argument
- Sinus
- Anomaly**
- Large Noise
- Medium Noise
- Small Noise

Тип данных: Вещественный

Назначение: ☒ Используемое

Сглаживание данных:

☐ Фурье-преобразование
Полоса пропускания: 50

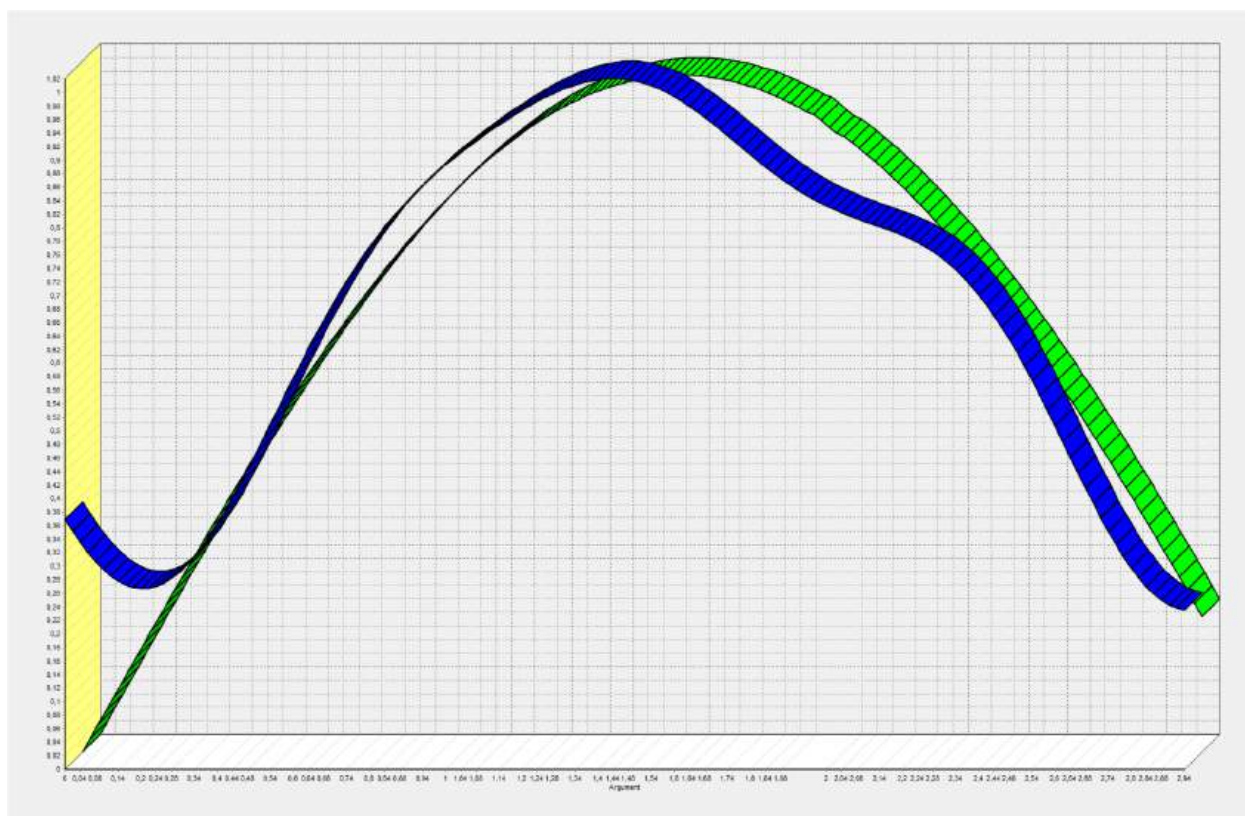
☐ Вейвлет-преобразование
Глубина разложения: 3
Порядок вейвлета: 6

Вычитание шума:

☒ Фурье-преобразование
Степень вычитания шума: Большая

< Назад Далее > Отмена

Получаем вот такую диаграмму



Дальше проделываем те же самые действия только ставим степень средняя.

Мастер обработки - Спектральная обработка (1 из 4)

Спектральная обработка

Настройка отдельных полей

- Argument
- Sinus
- Anomaly**
- Large Noise
- Medium Noise
- Small Noise

Тип данных: Вещественный

Назначение: ☒ Используемое

Сглаживание данных

☐ Фурье-преобразование
Полоса пропускания: 50

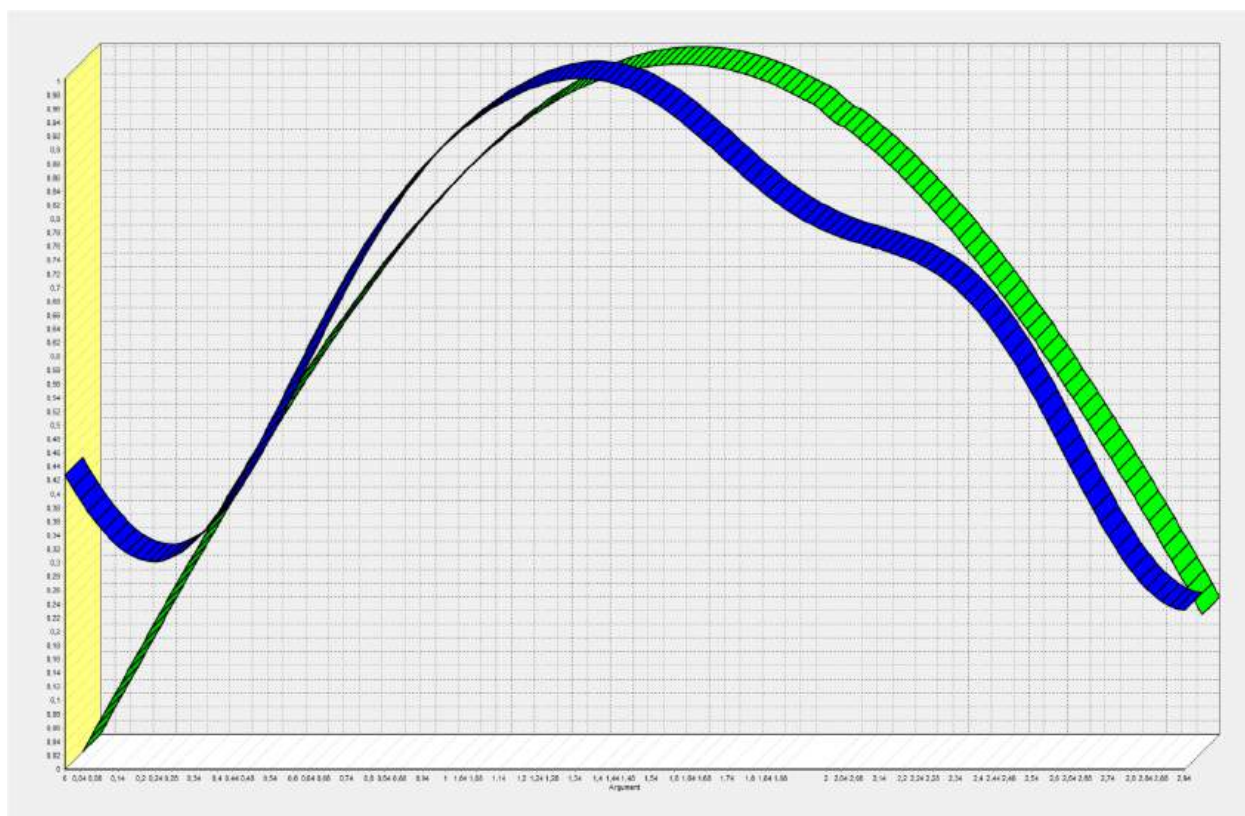
☐ Вейвлет-преобразование
Глубина разложения: 3
Порядок вейвлета: 6

Вычитание шума

☒ Фурье-преобразование
Степень вычитания шума: Средняя

< Назад Далее > Отмена

Получаем вот такую диаграмму



Дальше проделываем те же самые действия только ставим степень малая.

Мастер обработки - Спектральная обработка (1 из 4)

Спектральная обработка

Настройка отдельных полей

- Argument
- Sinus
- Anomaly**
- Large Noise
- Medium Noise
- Small Noise

Тип данных: Вещественный

Назначение: ☒ Используемое

Сглаживание данных:

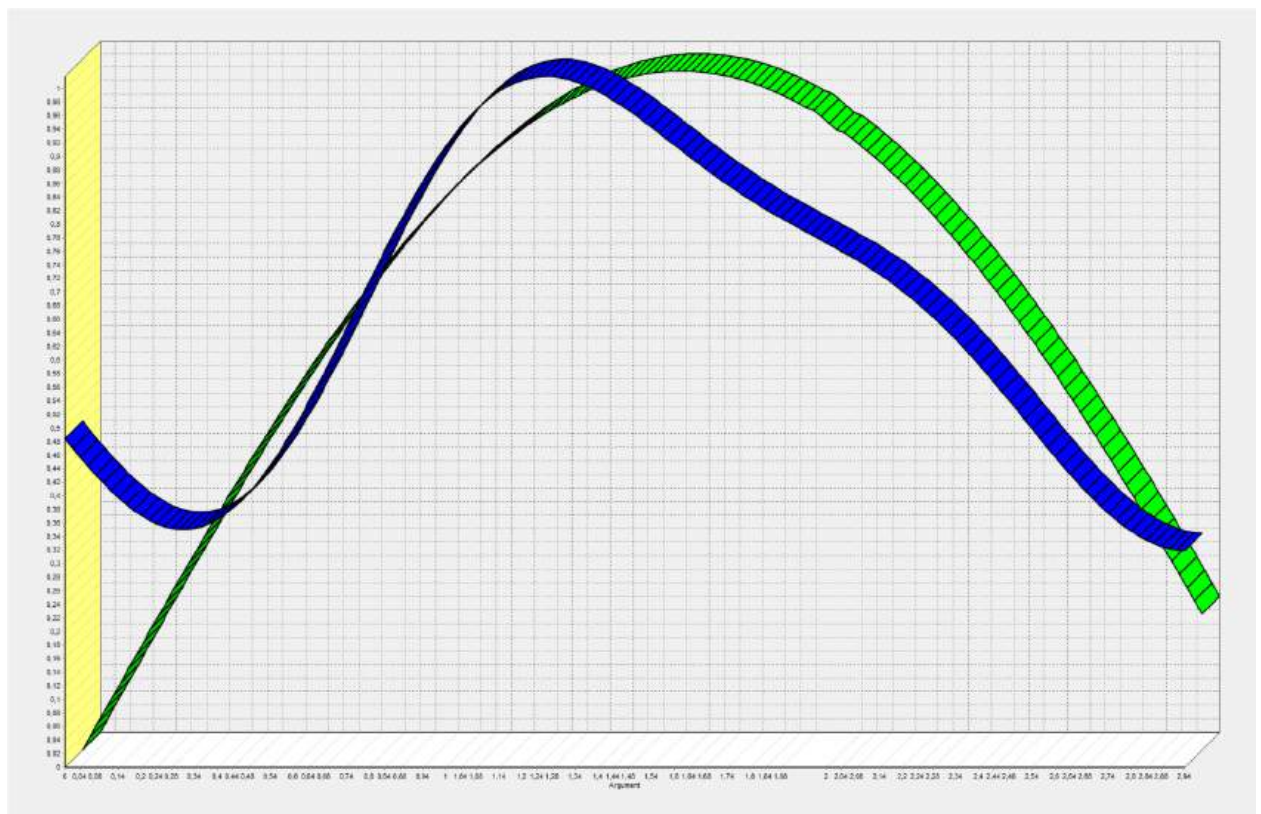
- ☐ Фурье-преобразование
 - Полоса пропускания: 50
- ☐ Вейвлет-преобразование
 - Глубина разложения: 3
 - Порядок вейвлета: 6

Вычитание шума:

- ☒ Фурье-преобразование
 - Степень вычитания шума: Малая

< Назад Далее > Отмена

Получаем вот такую диаграмму

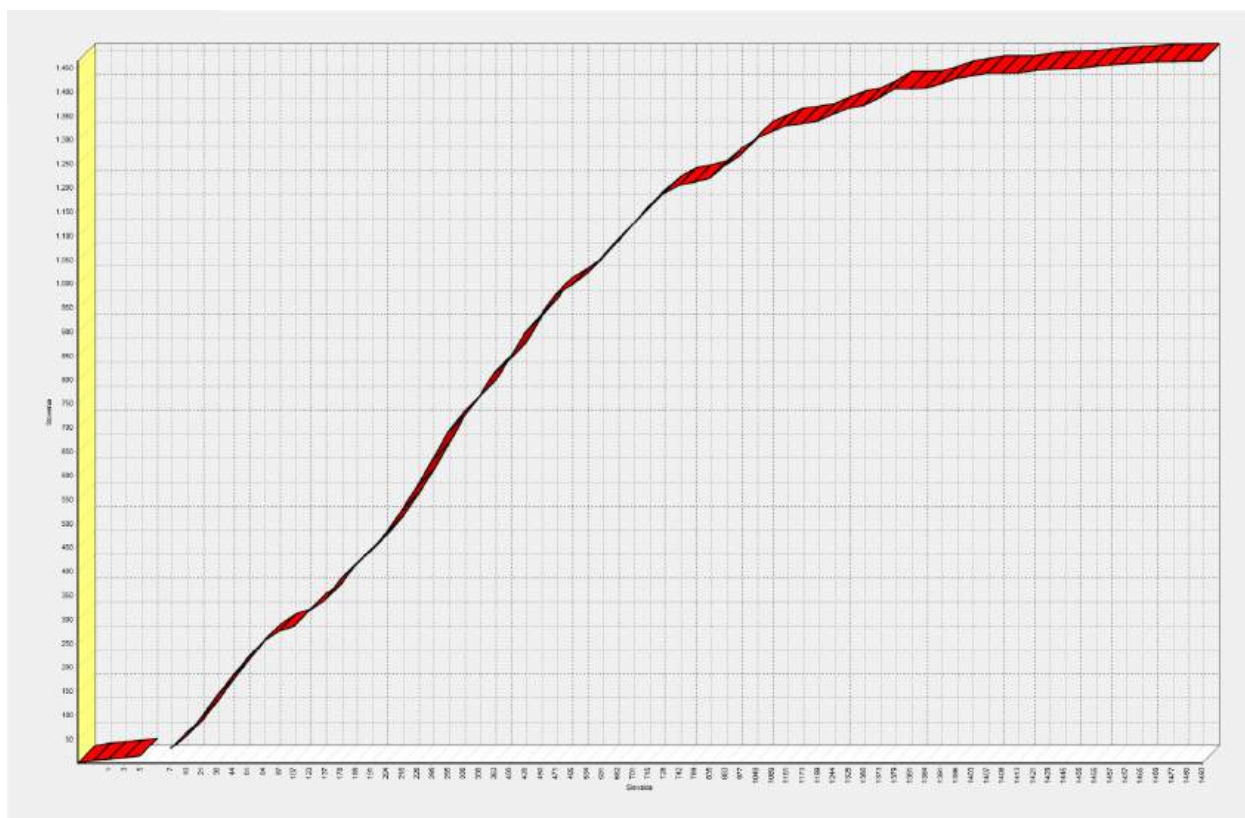


Ход работы 2 задание.

Для начала необходимо выбрать данные из excel таблицы по своему варианту и перенести их в txt файл.

1	Slovakia	Slovenia
2	nan	1
3	nan	6
4	1	9
5	3	12
6	5	16
7	nan	nan
8	7	31
9	10	57
10	21	96
11	30	141
12	44	181
13	61	219
14	84	253
15	97	275
16	107	286
17	123	319
18	137	341
19	178	383
20	185	414
21	191	442
22	204	480
23	216	528
24	226	577
25	295	632
26	295	691
27	336	730
28	336	763
29	363	814
30	400	841
31	426	897
32	450	934
33	471	977
34	485	997
35	534	1021
36	581	1055
37	682	1091
38	701	1124
39	715	1160
40	728	1188
41	742	1205
42	769	1212
43	835	1220
44	863	1248
45	977	1268
46	1049	1304
47	1089	1317
48	1161	1330
49	1173	1335
50	1199	1340
51	1244	1353
52	1325	1366
53	1360	1373
54	1373	1388

После чего загружаем эти данные в deductor, строим диаграмму и смотрим на пропуски.



Для заполнения пропусков буду использовать метод заполнения пропущенных данных в мастере обработки. Данный метод используется, потому что по графику видно, что у данных есть некоторый порядок и мы можем заменить на наиболее вероятные значения.

Мастер обработки - Заполнение пропусков (1 из 5)

Заполнение пропущенных данных
Общие настройки набора данных

☐ Использовать информацию узла оценки качества данных

☒ Обращивать как упорядоченный набор

Максимально допустимый процент пропущенных данных

< Назад **Далее >** Отмена

Мастер обработки - Заполнение пропусков (2 из 5)

Заполнение пропущенных данных
Настройка отдельных полей (упорядоченный набор данных)

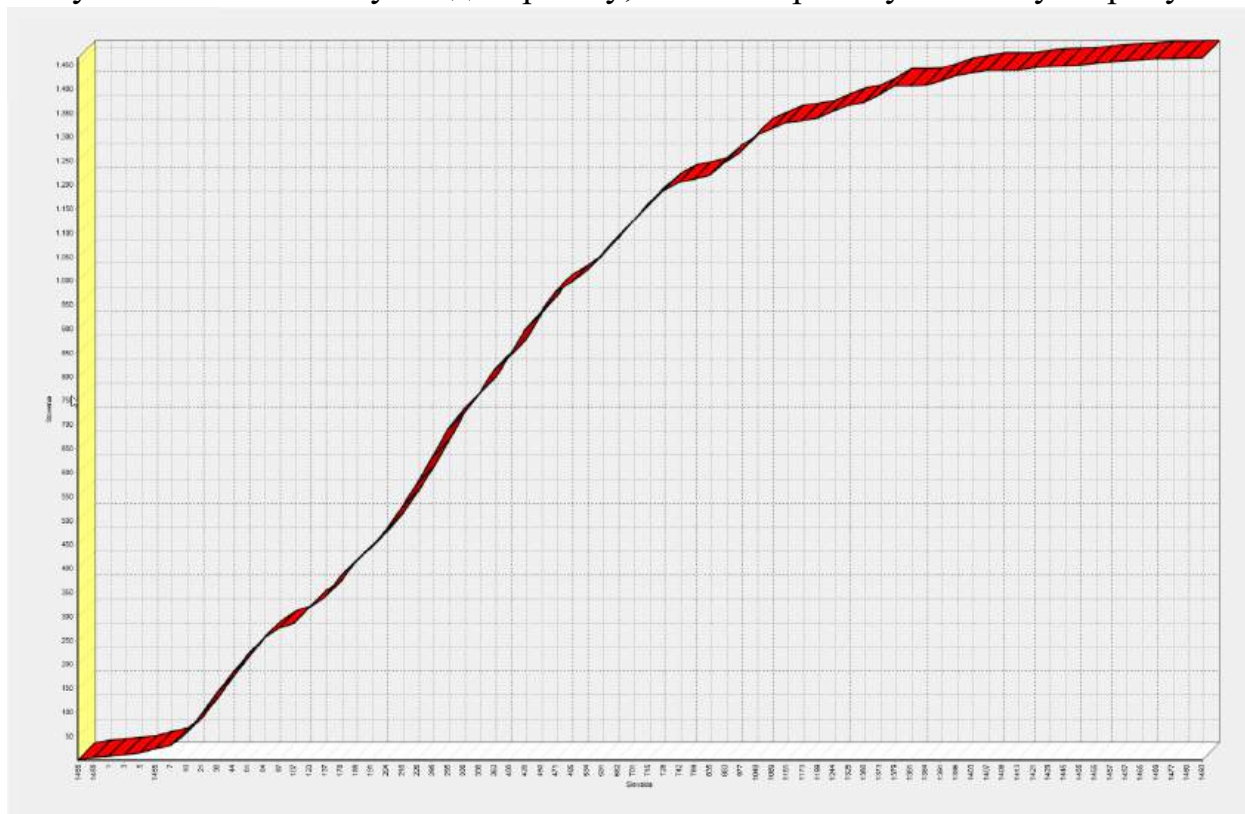
<input checked="" type="checkbox"/> Slovakia	Тип данных	Строковый
<input checked="" type="checkbox"/> Slovenia	Вид данных	Дискретный
	Назначение	<input checked="" type="checkbox"/> Используемое

Выбор метода обработки пропущенных данных:

Описание метода
Вместо пропущенных данных будет вставляться значение, наиболее часто встречающееся во входных данных (мода)

< Назад **Далее >** Отмена

Получаем обновленную диаграмму, в которой уже нету пропусков



Т.к. значения являются строго упорядоченными можно не редактировать аномальные выбросы. Чтобы подтвердить данную теорию в мастере обработки выберем редактирование выбросов и экстремальных значений в мастере обработки и зададим стандартные параметры с заменой на наиболее вероятное значение.

Мастер обработки - Редактирование выбросов (1 из 5)

Редактирование выбросов и экстремальных значений
Общие настройки набора данных

☐ Использовать информацию узла оценки качества данных

☐ Обращивать как упорядоченный набор данных

Метод определения выбросов и экстремальных значений

☒ Стандартное отклонение от среднего

Для выбросов Для экстремальных значений

☐ Интерквартильная ширина

Для выбросов Для экстремальных значений

< Назад **Далее >** Отмена

Мастер обработки - Редактирование выбросов (2 из 5)

Редактирование выбросов и экстремальных значений
Настройка отдельных полей (неупорядоченный набор данных)

☒ Slovakia
☒ Slovenia

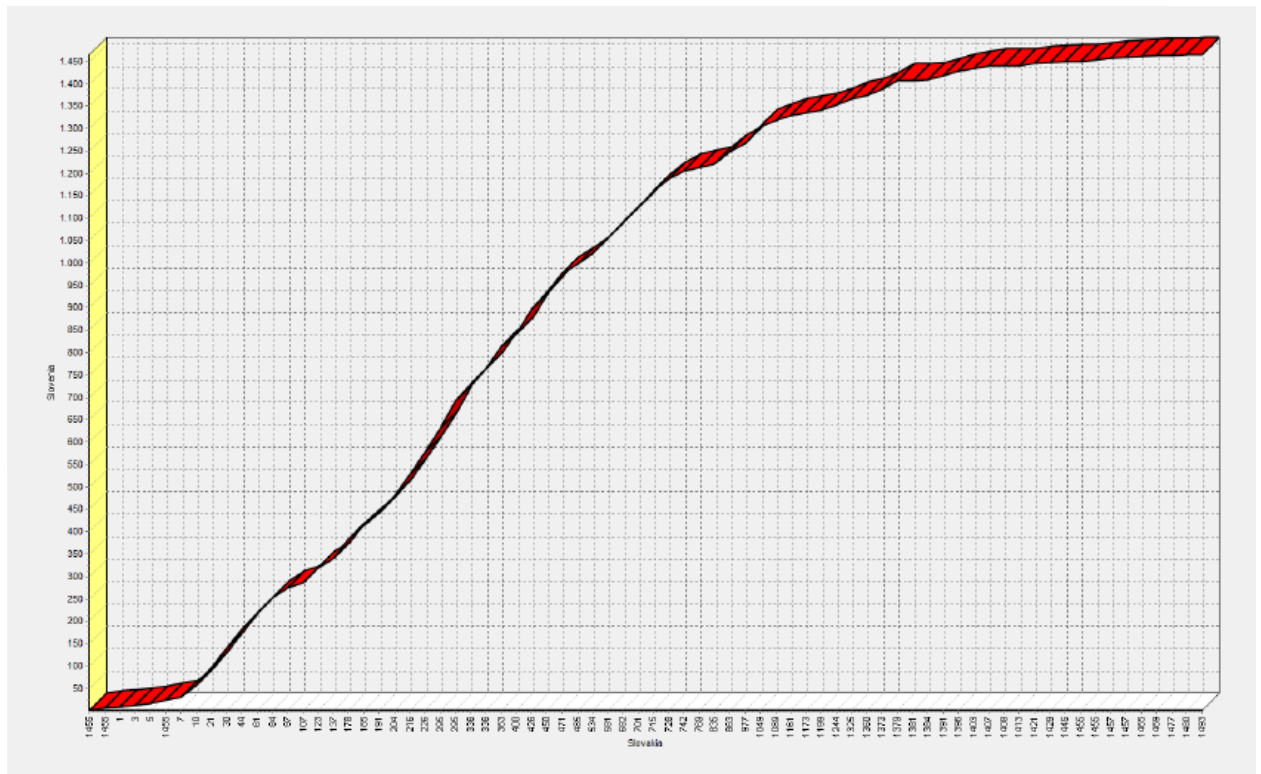
Тип данных
Вид данных
Назначение ☒ Используемое

Метод обработки экстремальных значений

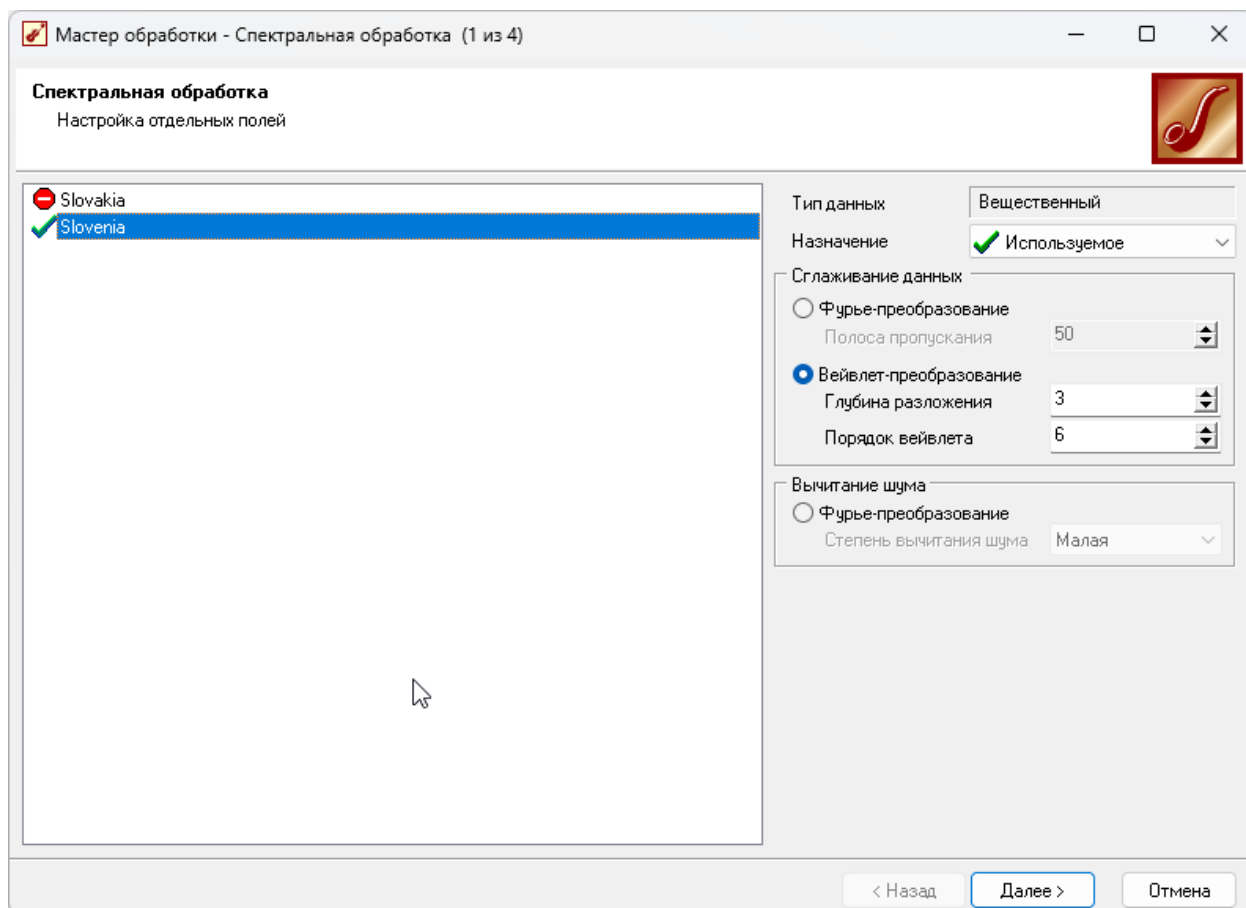
Метод обработки выбросов

< Назад **Далее >** Отмена

После построения диаграммы мы видим, что она никак не изменилась, т.е. у нас нету выбросов аномалий.



Далее сглаживаем данные с метода спектральной обработки в мастере обработки, с помощью вейвлет преобразования. Данный метод позволяет эффективно отделить сигнал от шума и так же он сохраняет важные составляющие – резкие изменения, выбросы и пр.



Далее удаляем шумы с больших до меньших. С помощью метода спектральная обработка в мастере обработке – вычитание шума. Для начала ставим степень Большая.

Мастер обработки - Спектральная обработка (2 из 5)

Спектральная обработка

Настройка отдельных полей

Slovakia

Slovenia

Тип данных: Вещественный

Назначение: Используемое

Сглаживание данных:

☐ Фурье-преобразование
Полоса пропускания: 50

☐ Вейвлет-преобразование
Глубина разложения: 3
Порядок вейвлета: 6

Вычитание шума:

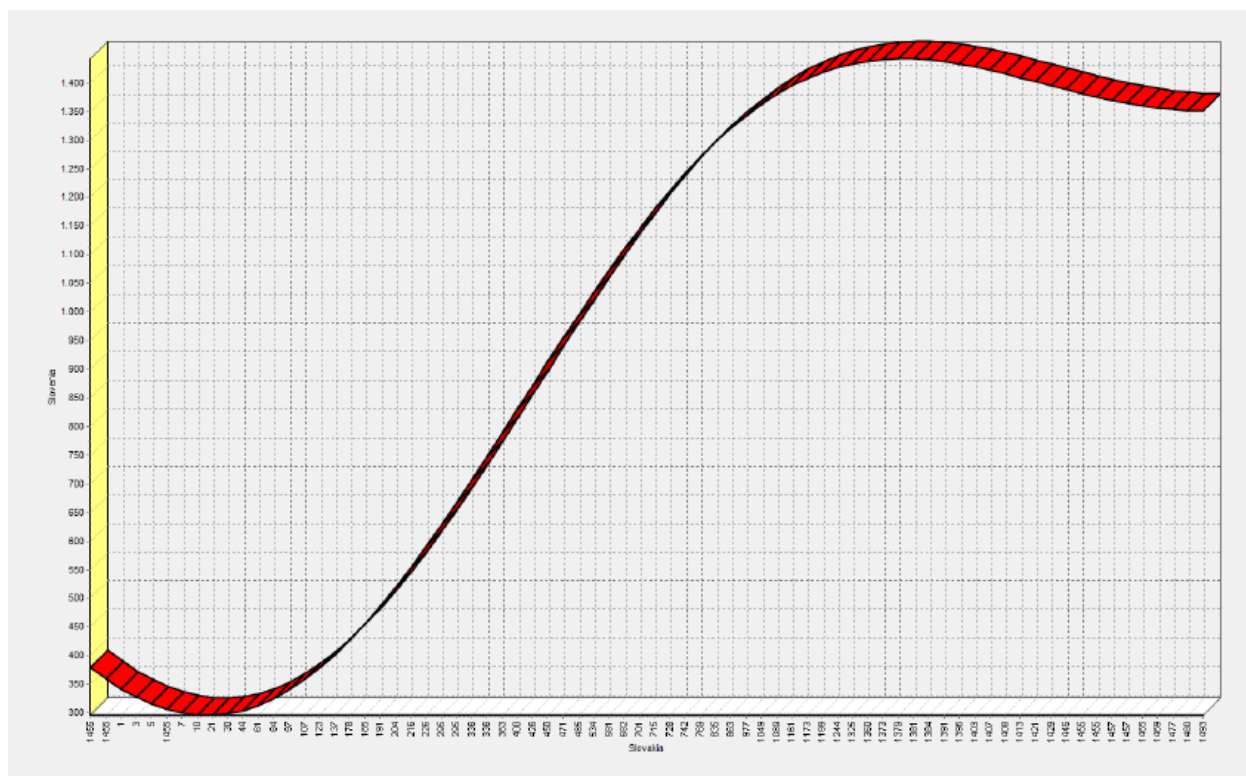
☒ Фурье-преобразование
Степень вычитания шума: Большая

< Назад

Далее >

Отмена

Получаем вот такую диаграмму



Теперь ставим степень Средняя.

Мастер обработки - Спектральная обработка (2 из 5)

Спектральная обработка

Настройка отдельных полей

Slovakia

Slovenia

Тип данных: Вещественный

Назначение: Используемое

Сглаживание данных:

☐ Фурье-преобразование
Полоса пропускания: 50

☐ Вейвлет-преобразование
Глубина разложения: 3
Порядок вейвлета: 6

Вычитание шума:

☒ Фурье-преобразование
Степень вычитания шума: Средняя

< Назад

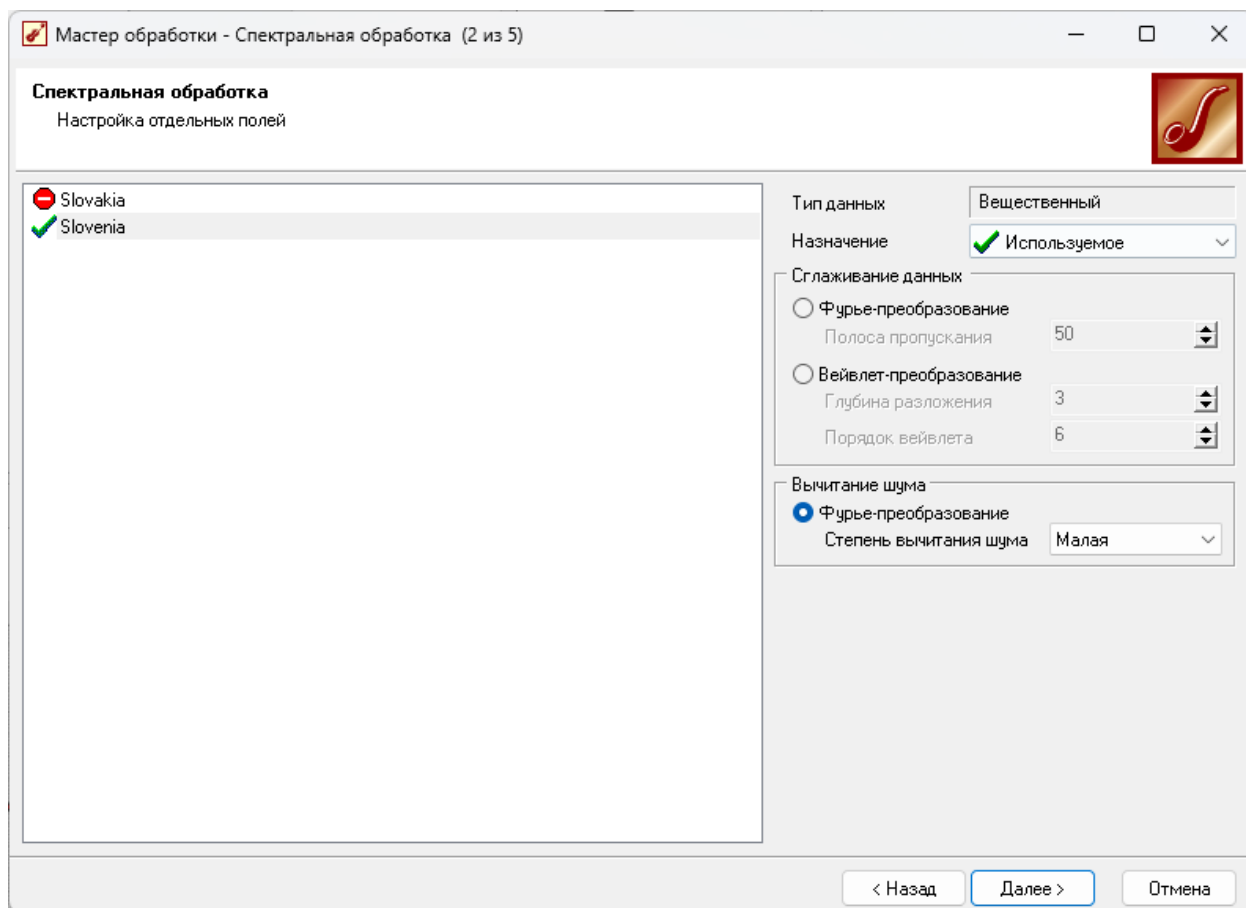
Далее >

Отмена

Получаем вот такую диаграмму.



Теперь ставим степень Малая.

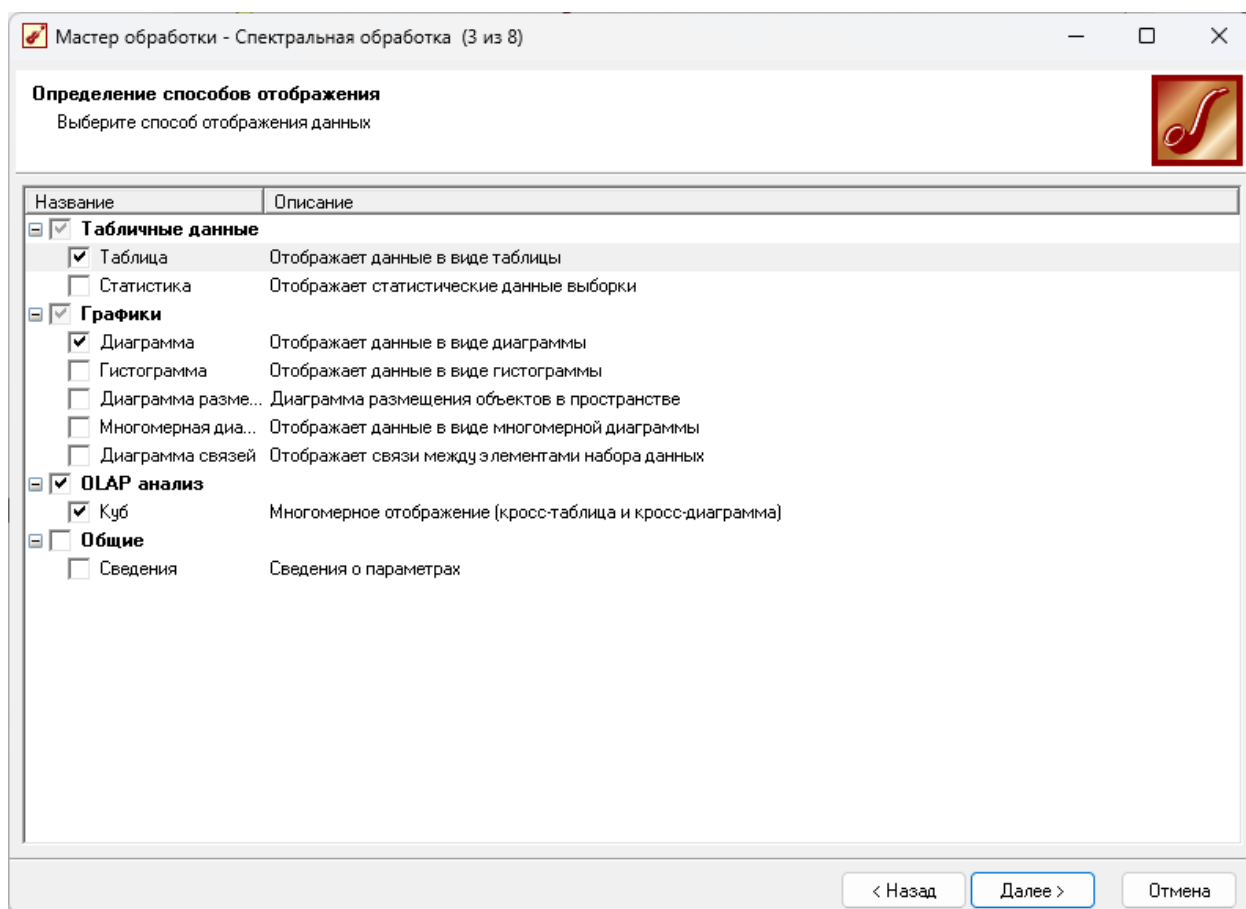


Получаем вот такую диаграмму.



В результате проделанной работы мы получаем качественные отфильтрованные данные.

Во время построения последней диаграммы, нужно указать OLAP анализ – куб.



Указываем 2 измерения

Мастер обработки - Спектральная обработка (5 из 8)

Настройка назначений полей куба

Задайте назначения столбцов для многомерного отображения (кросс-таблицы и кросс-диаграммы)

↗

Slovakia

↗

Slovenia

Имя столбца: Slovenia
Метка столбца: Slovenia
Тип данных: 9.0 Вещественный
Вид данных: — Непрерывный
Назначение: ↗ Измерение

☒ Кэшировать исходные значения фактов
☒ Использовать инкрементный расчет фактов

< Назад
Далее >
Отмена

Ставим эти измерения в строки

Мастер обработки - Спектральная обработка (6 из 8)

Настройка измерений

Настройка размещения измерений

Доступные измерения

↑

↓

↔

↔

<

<<

Выбранные измерения

Колонки

Строки

9.0 Slovenia

ab Slovakia

< Назад
Далее >
Отмена

В фактах указываем количество

Мастер обработки - Спектральная обработка (7 из 8)

Настройка фактов

Выбор отображаемых фактов, агрегации и вариантов их отображения

↑ ↓ [иконки]

Факты и варианты агрегации

☑ # Количество

☑ # Количество

↑ ↓ [иконки]

Варианты отображения

☑ Σ Значение

☐ % Процент по горизонтали

☐ % Процент по вертикали

< Назад

Далее >

Отмена

Получаем вот такой куб

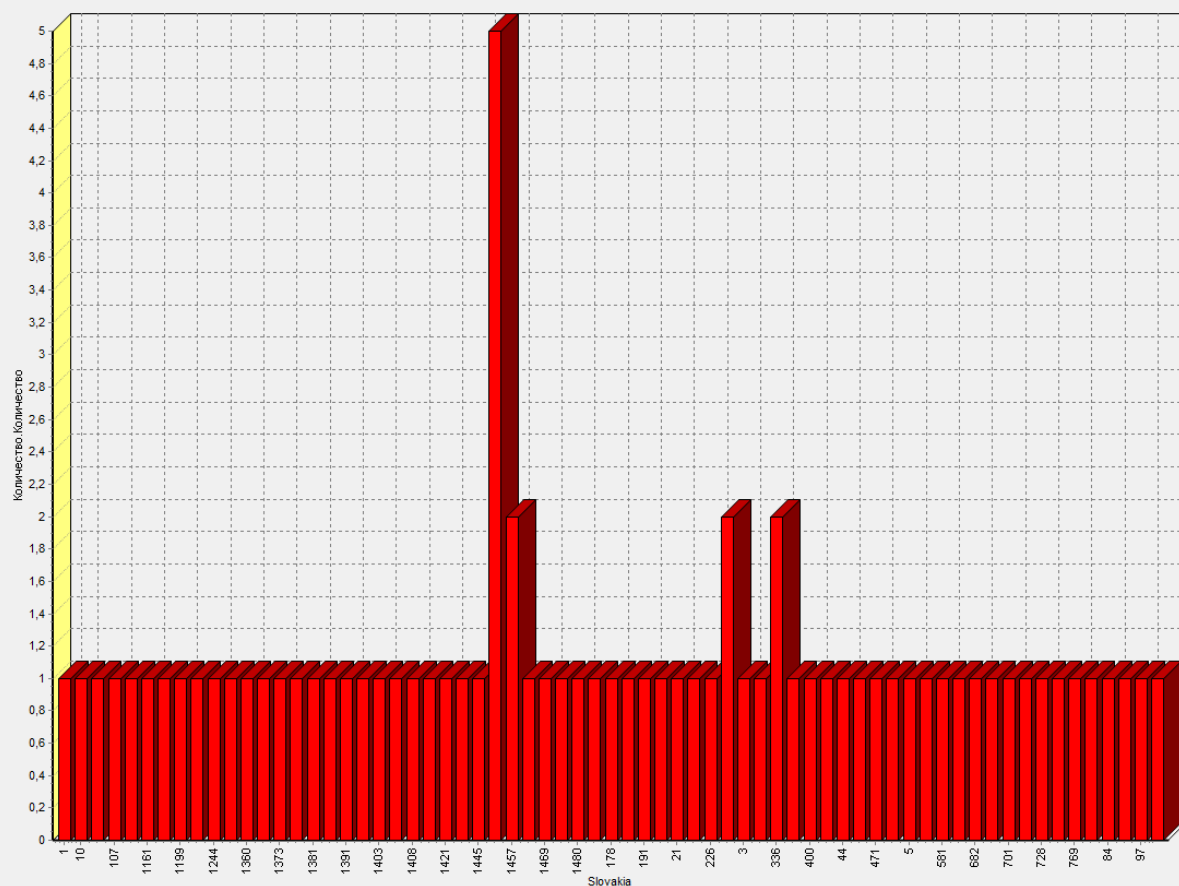
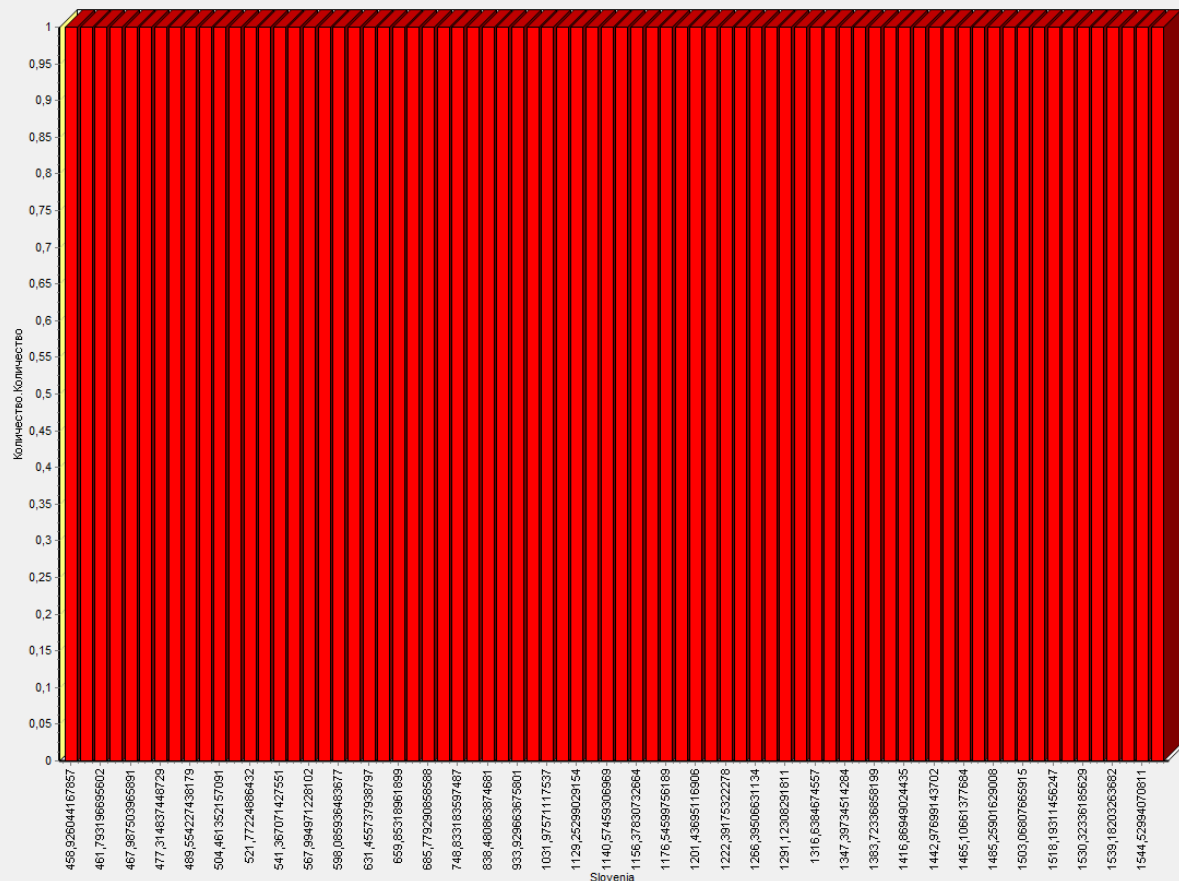
Только словакия

Slovakia	# Количество
1	1
10	1
1049	1
107	1
1089	1
1101	1
1173	1
1199	1
123	1
1244	1
1325	1
1360	1
137	1
1373	1
1379	1
1381	1
1384	1
1391	1
1396	1
1403	1
1407	1
1408	1
1413	1
1421	1
1426	1
1445	1
1455	1
1457	2
1465	1
1469	1
1477	1
1480	1
1492	1
178	1
185	1
191	1
204	1
21	1
216	1
226	1
295	2
3	1

Только Словения

Словения	#	Количество
498,926044167617		1
499,55012021895		1
461,79319069002		1
463,795555555277		1
467,907503905951		1
471,76162095908		1
477,35463748729		1
483,50165395387		1
489,584227438179		1
495,037344918183		1
504,461352152091		1
518,34664641552		1
521,77224886432		1
541,207191124077		1
541,267071427511		1
552,474674979066		1
567,994971238162		1
585,275490202272		1
595,085930483677		1
609,305381085307		1
631,405737928797		1
634,265354340085		1
659,85318961899		1
667,881801513495		1
685,770380858588		1
707,40520898077		1
748,833181907487		1
792,742697923501		1
838,480863874681		1
885,674771063061		1
933,926491675801		1
982,8364251412626		1
1021,97371117327		1
1080,82287952483		1
1128,252999281154		1
1157,205991217516		1
1140,57459305069		1
1146,9438084034		1
1156,37830732664		1
1168,7362013222		1
1176,54590756189		1
1183,83165371248		1

Теперь построим диаграммы по кубу.



Ответы на вопросы

1. Для чего следует проводить подготовку данных для анализа?

Подготовка данных обеспечивает качество и точность анализа, устраняя ошибки, пропуски и несоответствия, а также преобразуя данные в подходящий формат для эффективного использования аналитическими инструментами.

2. Что такое шумы и аномалии в данных?

- **Шумы** — случайные и нерегулярные колебания или ошибки в данных, не несущие полезной информации.
- **Аномалии** — значительные отклонения от нормального поведения данных, которые могут указывать на ошибки, необычные события или важные изменения.

3. Какими методами можно убрать шумы в системе Deductor?

- **Скользящее среднее:** сглаживание данных путем усреднения соседних точек.
- **Вейвлет-преобразование:** разделение сигнала и удаление высокочастотных шумов.
- **Фильтры низких частот:** устранение высокочастотных компонентов шума.

4. Какими методами можно убрать аномалии данных в системе Deductor?

- **Статистические методы:** использование Z-оценки или межквартильного размаха для выявления выбросов.
- **Методы кластеризации:** обнаружение точек, не принадлежащих ни одному кластеру.
- **Визуальный анализ:** применение графических методов, таких как ящики с усами или диаграммы рассеяния.

5. Для чего используется парциальная предобработка?

- Парциальная предобработка применяется для частичной очистки и трансформации данных на отдельных этапах обработки, улучшая качество данных перед основным анализом и позволяя более эффективно управлять большими объемами информации.

6. Для чего используется спектральная обработка?

Спектральная обработка анализирует частотные компоненты данных, позволяя выявлять скрытые паттерны, удалять шумы и аномалии, а также

сглаживать данные для улучшения их качества и подготовки к дальнейшему анализу.

7. Какие виды спектральной обработки имеются в системе Deductor?

- **Быстрое преобразование Фурье (БПФ):** анализ частотных компонентов сигнала.
- **Вейвлет-преобразование:** многошкальный анализ данных.
- **Фильтрация частот:** удаление нежелательных частотных составляющих.

8. Что такое аппроксимация данных?

Аппроксимация данных — процесс приближения сложных данных более простыми моделями или функциями, сохраняя основные тенденции и структуры, что облегчает анализ и визуализацию.