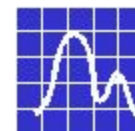


Сегментация в примерах

Полежаев Илья
Владимир Боровиков



StatSoft® Russia

Кластерный анализ

Анализируемый объект представляет собой жилое многоквартирное здание, подвергнутое модернизации с целью экономии потребления ресурсов.



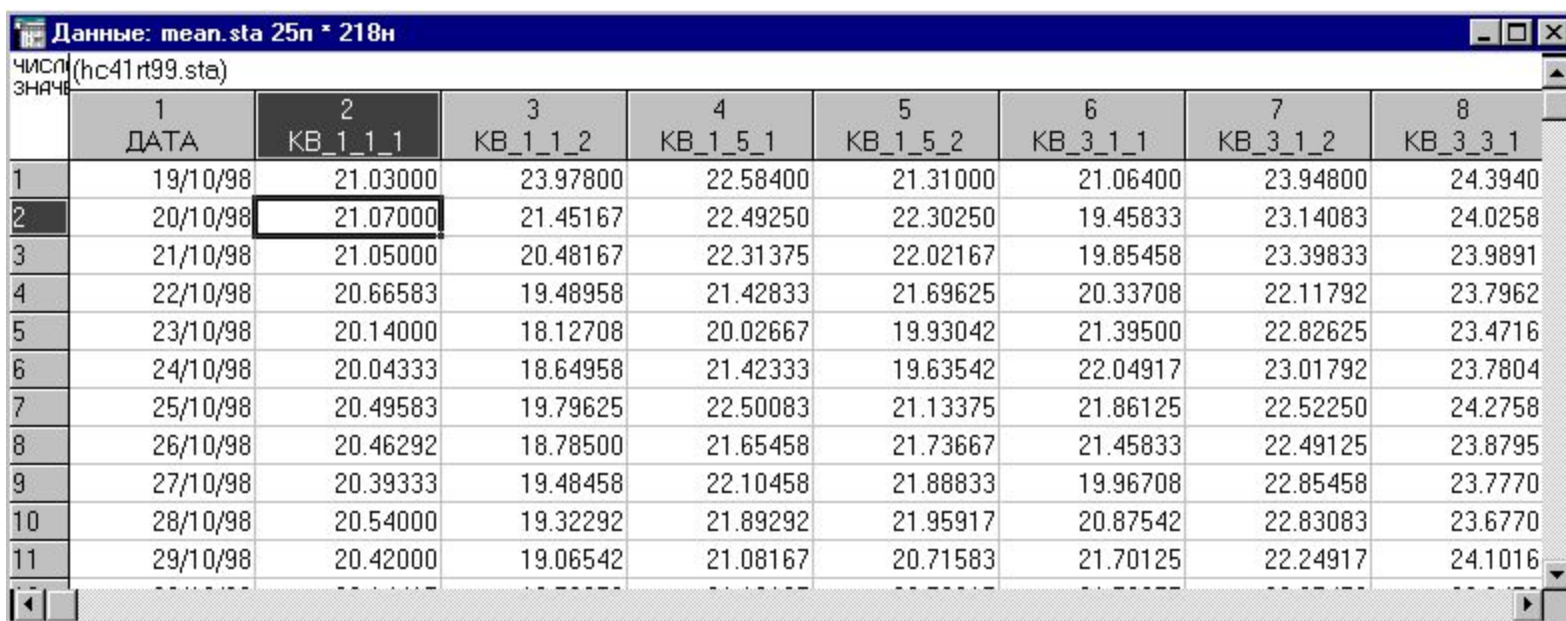
Постановка задачи

Требуется разбить
квартиры на три как можно
больше отличающиеся
друг от друга группы по
признаку - **количество
потребляемого тепла.**



Данные

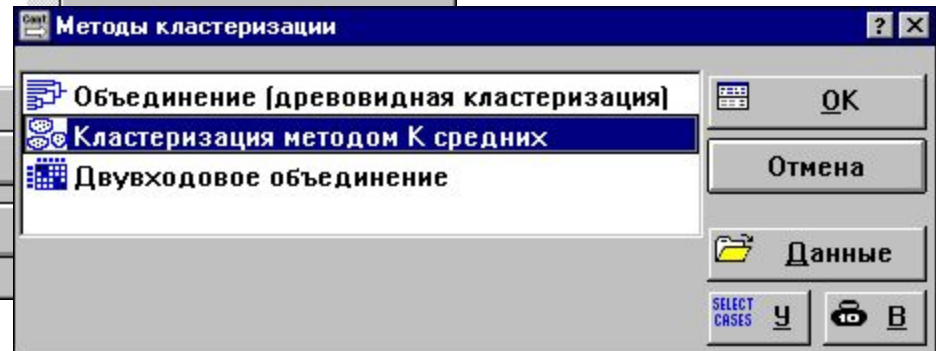
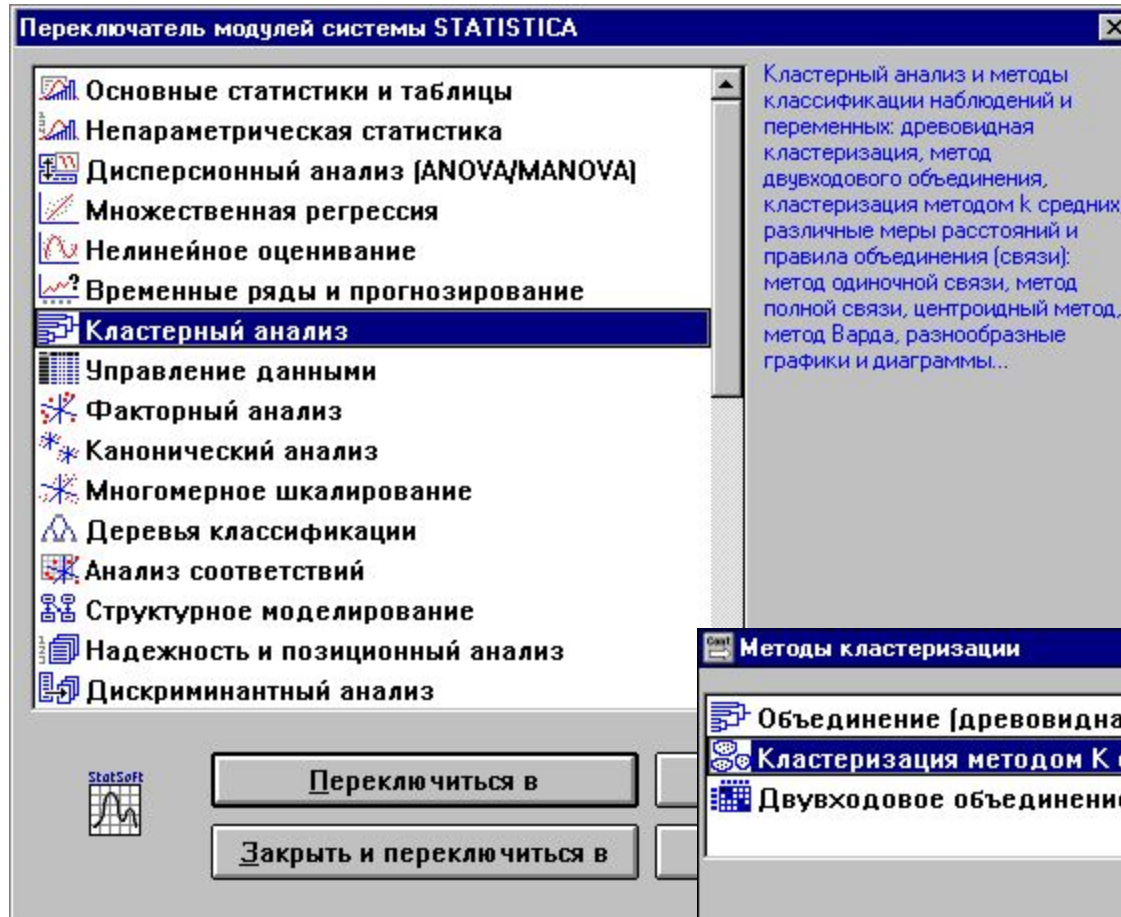
Показания датчиков температуры, размещенных в квартирах. Данные собирались в течение 2-х лет.



Скриншот окна с данными датчиков температуры. Таблица содержит 11 строк дат и 8 столбцов, соответствующих различным квартирам. Значение в ячейке 20/10/98 для квартиры КВ_1_1_1 выделено.

1	2	3	4	5	6	7	8	
ДАТА	КВ_1_1_1	КВ_1_1_2	КВ_1_5_1	КВ_1_5_2	КВ_3_1_1	КВ_3_1_2	КВ_3_3_1	
1	19/10/98	21.03000	23.97800	22.58400	21.31000	21.06400	23.94800	24.3940
2	20/10/98	21.07000	21.45167	22.49250	22.30250	19.45833	23.14083	24.0258
3	21/10/98	21.05000	20.48167	22.31375	22.02167	19.85458	23.39833	23.9891
4	22/10/98	20.66583	19.48958	21.42833	21.69625	20.33708	22.11792	23.7962
5	23/10/98	20.14000	18.12708	20.02667	19.93042	21.39500	22.82625	23.4716
6	24/10/98	20.04333	18.64958	21.42333	19.63542	22.04917	23.01792	23.7804
7	25/10/98	20.49583	19.79625	22.50083	21.13375	21.86125	22.52250	24.2758
8	26/10/98	20.46292	18.78500	21.65458	21.73667	21.45833	22.49125	23.8795
9	27/10/98	20.39333	19.48458	22.10458	21.88833	19.96708	22.85458	23.7770
10	28/10/98	20.54000	19.32292	21.89292	21.95917	20.87542	22.83083	23.6770
11	29/10/98	20.42000	19.06542	21.08167	20.71583	21.70125	22.24917	24.1016

Анализ



Анализ

Кластерный анализ: кластеризация методом К средних

Переменные: KV_1_1_1-KV_6_5_1

Кластер: Переменные (столбцы)

Число кластеров: 3

Число итераций: 2

Выбор переменных для анализа

1-ДАТА	13-KV 6 1
2-KV 1 1 1	14-KV 6 3
3-KV 1 1 2	15-KV 6 3
4-KV 1 5 2	16-KV 6 5
5-KV 1 5 2	17-СРЕДН
6-KV 3 1 1	18-СР1 ЭТ
7-KV 3 1 2	19-СР3 ЭТ
8-KV 3 3 1	20-СР5 ЭТ
9-KV 3 3 2	21-СР1 ПО
10-KV 3 3 2	22-СР3 ПО
11-KV 3 3 2	23-СР6 ПО
12-KV 3 3 2	24-КЛАСТЕ

Результаты метода К средних

Количество переменных: 15
Кол. наблюдений: 212
Класт. перем. методом К средних
ИЦ построчно удалены
Число кластеров: 3
Решение получено после 2 итераций

Дисперсионный анализ

Средние кластеров и евклидовы расстояния

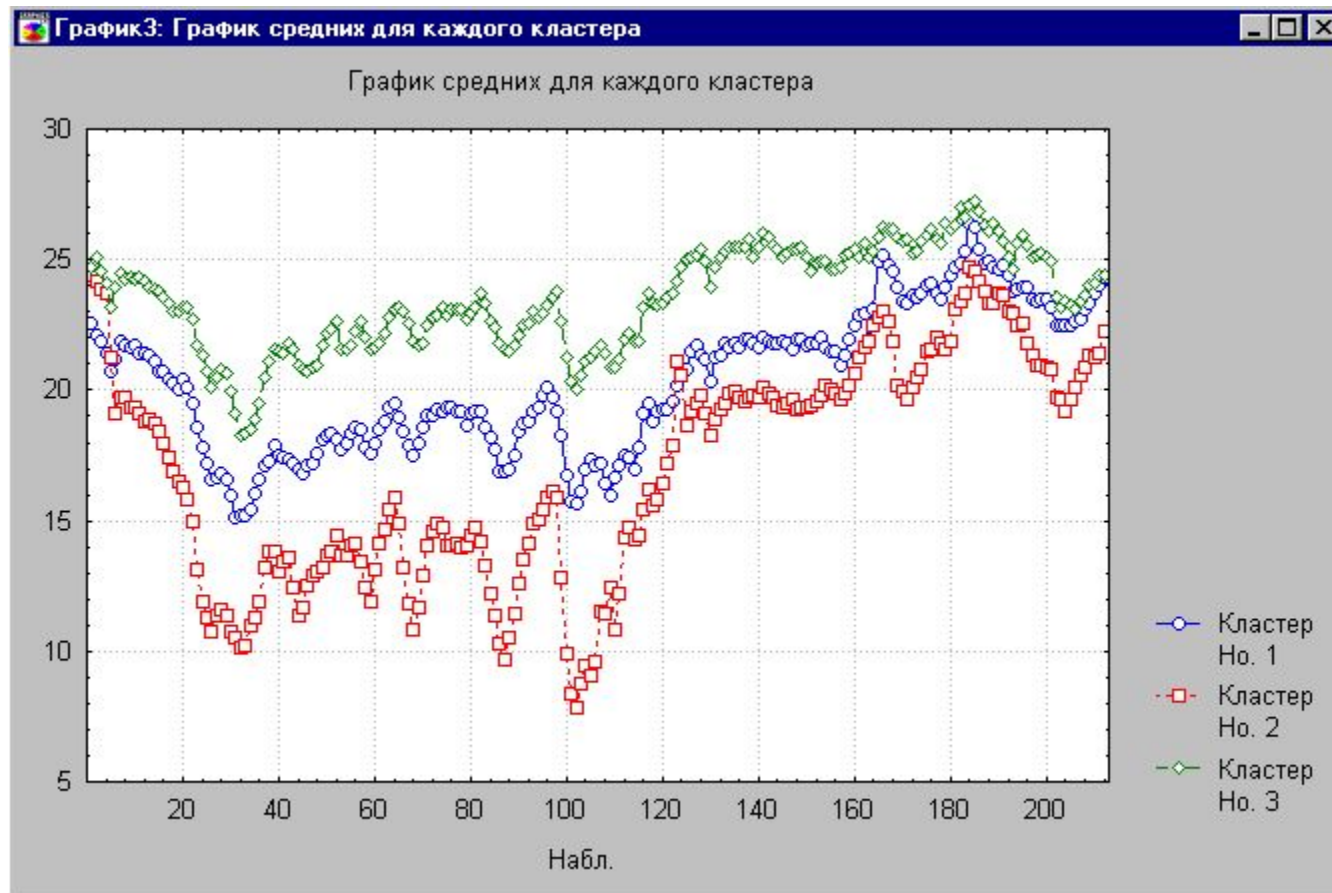
График средних

Описат. статистики для каждого кластера

Элементы кластеров и расстояния

Сохранить классификации и расстояния

Графический анализ



Результаты

Элементы кластера номер 1 (mean.sta)								
КЛАСТЕР. АНАЛИЗ	и расстояния до центра кластера. Кластер содержит 8 переменных							
	Перемен. КВ_1_1_1	Перемен. КВ_1_1_2	Перемен. КВ_1_5_1	Перемен. КВ_1_5_2	Перемен. КВ_3_1_1	Перемен. КВ_3_1_2	Перемен. КВ_6_1_1	Перемен. КВ_6_1_2
Расст-е	1.971874	1.506634	1.089310	1.077315	1.598448	1.637242	1.385841	.954671

Элементы кластера номер 2 (mean.sta)	
КЛАСТЕР. АНАЛИЗ	и расстояния до центра кластера. Кластер содержит 1 переменных
	Перемен. КВ_6_5_1
Расст-е	0.00

Элементы кластера номер 3 (mean.sta)						
Далее...	и расстояния до центра кластера. Кластер содержит 6 переменных					
	Перемен. КВ_3_3_1	Перемен. КВ_3_3_2	Перемен. КВ_3_5_1	Перемен. КВ_3_5_2	Перемен. КВ_6_3_1	Перемен. КВ_6_3_2
Расст-е	1.058891	.824863	1.405683	.976158	.981679	1.529956

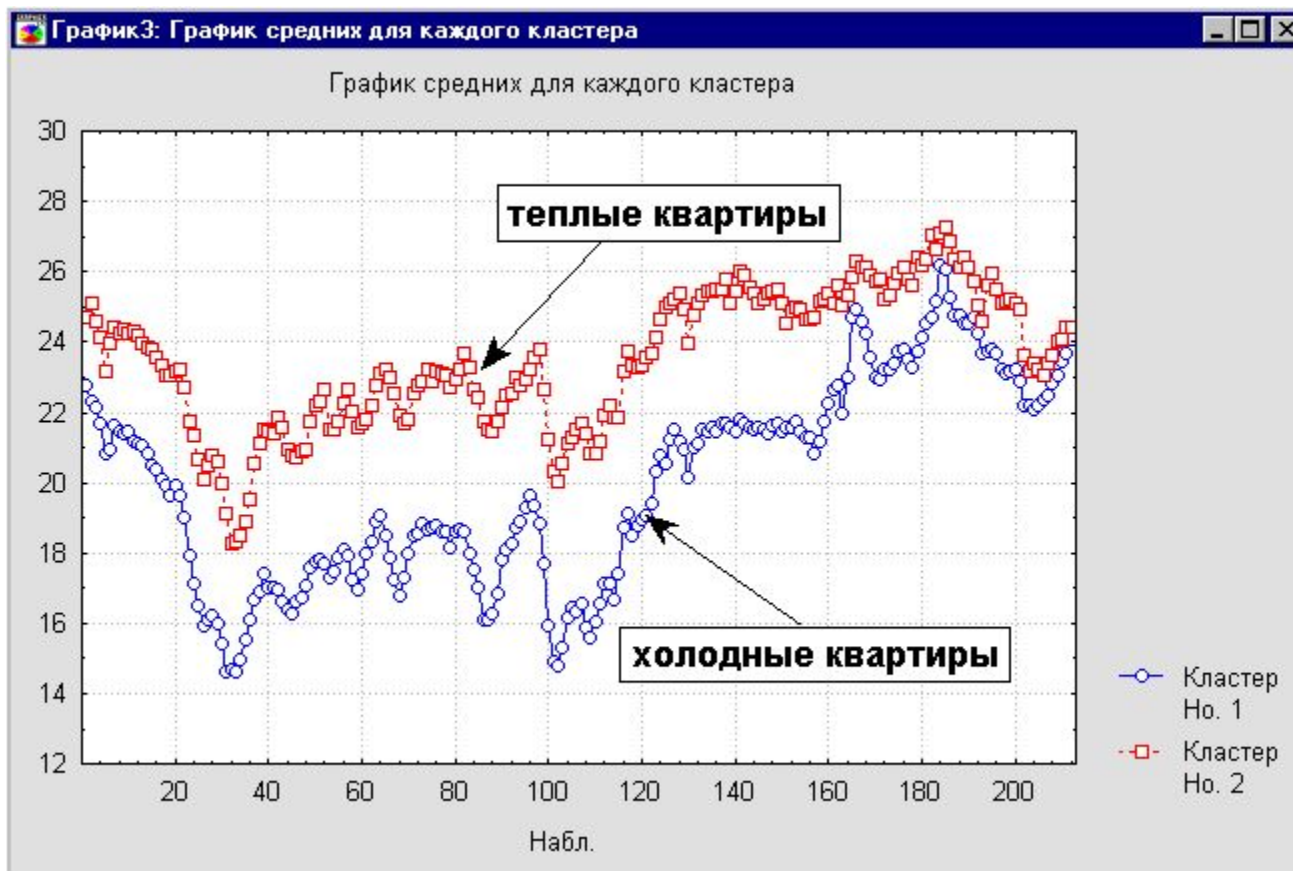
Результаты

КЛАСТЕР. АНАЛИЗ	Кластер Но. 1	Кластер Но. 2	Кластер Но. 3
1	22.57700	24.37400	24.71400
2	22.07531	24.16000	25.10563
3	21.89677	23.90042	24.59660
4	21.43693	23.69958	24.11257
5	20.75703	21.26333	23.17507
6	21.18417	19.14208	23.94160
7	21.86151	19.70167	24.44437
8	21.70307	19.74208	24.24840
9	21.65292	19.36208	24.36618
10	21.69016	19.35708	24.24667
11	21.40417	19.14667	24.29604
12	21.42771	18.82333	24.17430
13	21.30474	18.87083	23.96694
14	21.12156	18.69917	23.83021
15	20.74688	18.43667	23.79194



КЛАСТЕР. АНАЛИЗ Кластер номер	Расстояния по диагоналию Квадраты расстояний над диагоналию		
	Но. 1	Но. 2	Но. 3
Но. 1	0.000000	14.66892	11.26418
Но. 2	3.830003	0.00000	49.23798
Но. 3	3.356215	7.01698	0.00000

А если 2 кластера?

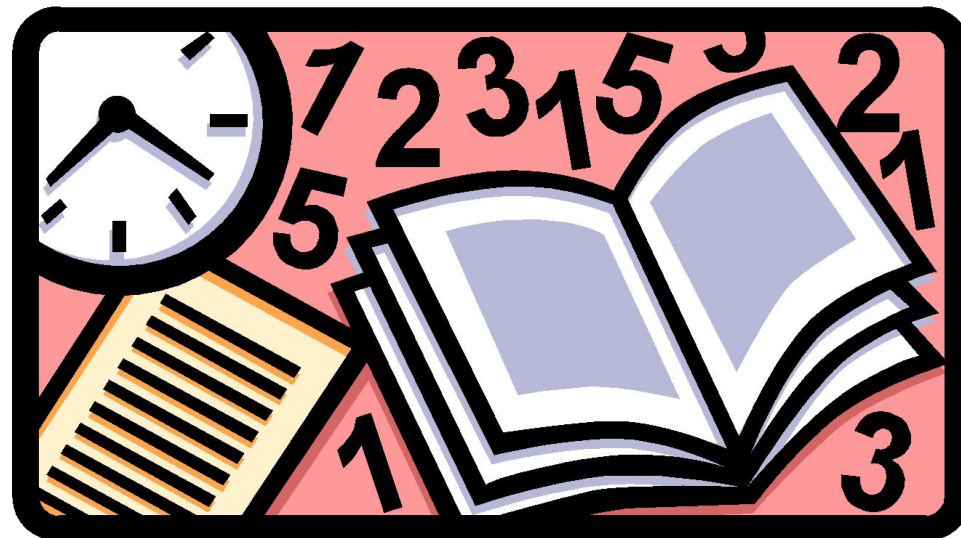


Выводы

- Кластеризация выполнена
- Самые теплые квартиры расположены в середине дома
- Самые холодные квартиры - угловые

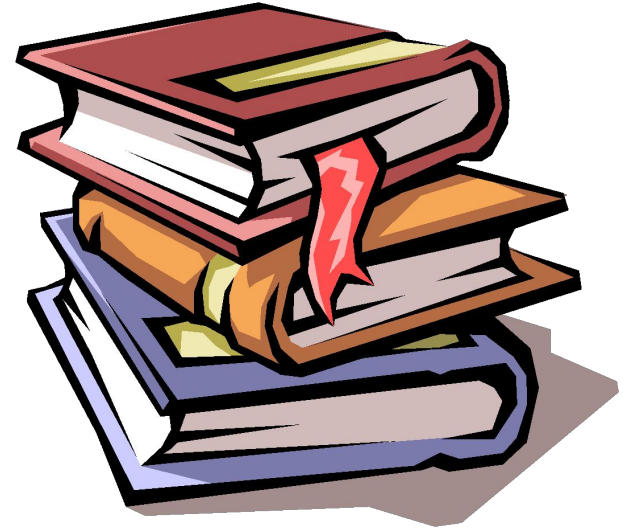


Сегментирование потребителей услуги "Подключение к сети Интернет"



Что необходимо?

- Наличие у фирмы минимум 100 клиентов
- Возможность проведения опроса клиентов фирмы
- Наличие STATISTICA





Этапы методики

- Определение возможных признаков сегментирования
- Проведение опроса
- Определение "пригодных" признаков сегментирования
- Выделение сегментов
- Формулировка сегментов

Признаки сегментирования

- Возраст
- Стаж работы в сети Интернет
- Профессиональная специализация
- Время работы в сети





Опросный лист

- Укажите, пожалуйста, Ваш возраст:
 - до 20 лет
 - 20-35
 - 35-55
 - старше 55
- Как давно Вы работаете в Интернет?
 - менее 1 года
 - 1-2 года
 - 2-3 года
 - более 3-х лет

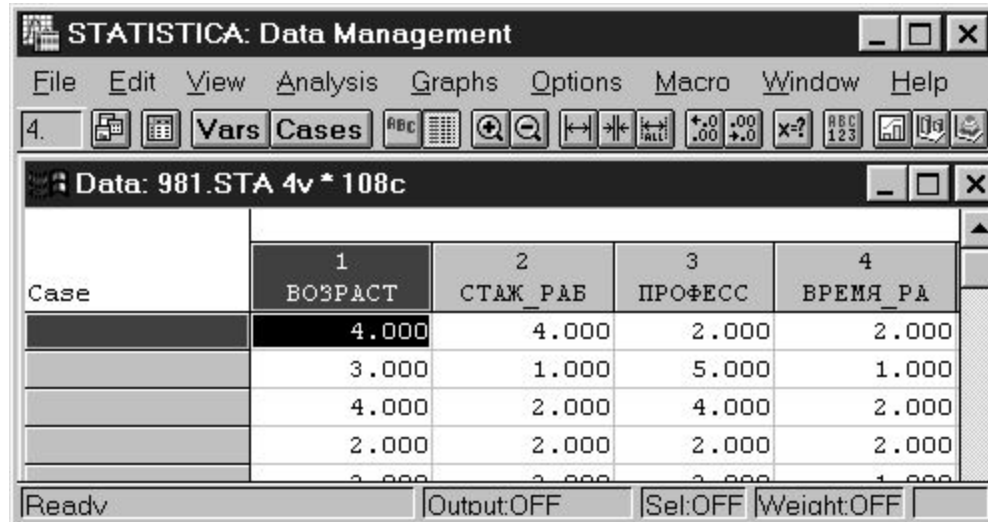


Опросный лист

- В какой степени Вы используете возможности сети в вашей работе?
 - не использую
 - использую крайне редко
 - ежедневно обращаюсь к сети
 - интернет - часть моей работы
- Как часто Вы выходите в Интернет
 - несколько раз в месяц
 - несколько раз в неделю
 - один раз в день
 - несколько раз в день

Ответы кодируются

- Не использую == 1
- Крайне редко == 2
- Ежедневно == 3
- Интернет - моя работа == 4



The screenshot shows the STATISTICA Data Management interface. The main window displays a data table with the following structure:

Case	1 ВОЗРАСТ	2 СТАЖ_РАБ	3 ПРОФЕСС	4 ВРЕМЯ_РА
	4.000	4.000	2.000	2.000
	3.000	1.000	5.000	1.000
	4.000	2.000	4.000	2.000
	2.000	2.000	2.000	2.000
	3.000	3.000	3.000	1.000

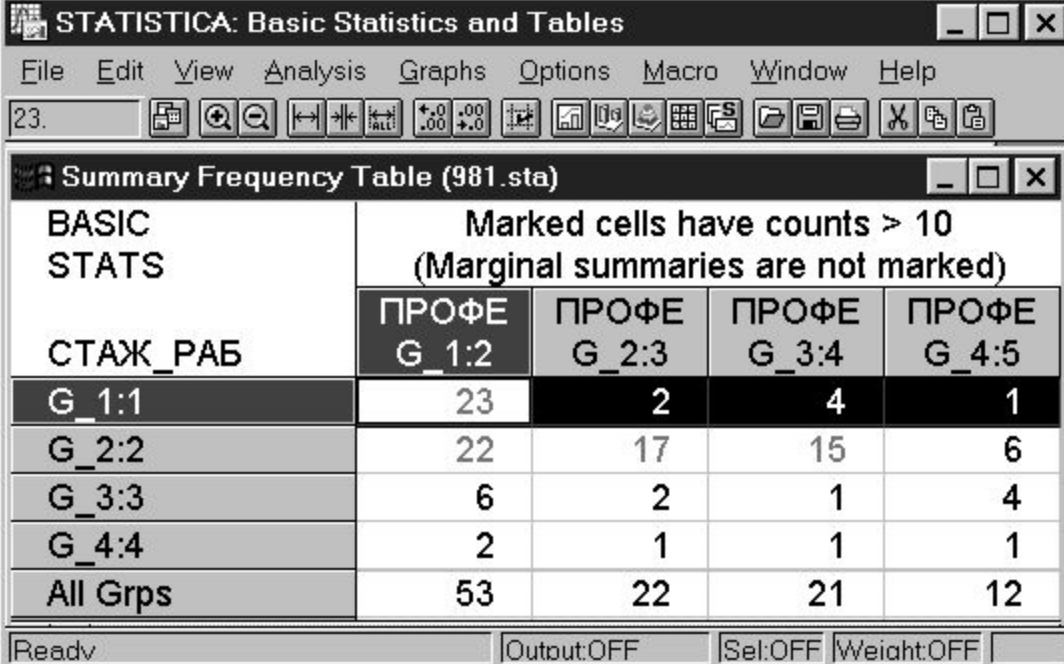
The status bar at the bottom indicates: Ready | Output:OFF | Sel:OFF | Weight:OFF

Выделение признаков

Степенью "пригодности" определенного признака сегментирования можно считать наличие определенной математической корреляции между парой предполагаемых признаков

ЧИСЛОВЫЕ ЗНАЧЕНИЯ	1	2	3	4
	ВОЗРАСТ	СТАЖ_РАБ	ПРОФЕС	ВРЕМЯ_РА
ВОЗРАСТ	1.00	.10	.05	0.00
СТАЖ_РАБ	.10	1.00	.25	.20
ПРОФЕС	.05	.25	1.00	.20
ВРЕМЯ_РА	0.00	.20	.20	1.00

Выделение сегментов



STATISTICA: Basic Statistics and Tables

File Edit View Analysis Graphs Options Macro Window Help

23. [Icons]

Summary Frequency Table (981.sta)

BASIC STATS

Marked cells have counts > 10
(Marginal summaries are not marked)

СТАЖ_РАБ	ПРОФЕ G_1:2	ПРОФЕ G_2:3	ПРОФЕ G_3:4	ПРОФЕ G_4:5
G_1:1	23	2	4	1
G_2:2	22	17	15	6
G_3:3	6	2	1	4
G_4:4	2	1	1	1
All Grps	53	22	21	12

Ready Output:OFF Sel:OFF Weight:OFF

Результат сегментации

Профессиональная специализация"	Стаж работы в сети"	%
Не использую	Менее 1 года	23
Не использую	1-2 года	22
Использую крайне редко	1-2 года	17
Ежедневно обращаюсь к сети	1-2 года	



Формулировка сегментов

Профессиональная специализация	Стаж работы в сети	" Профи Профи"	Описание (характеристики сегмента)
Не использую	Менее 1 года	" Новичок Новичок"	Данный сегмент еще не успел основательно освоиться в среде Интернет, использует Интернет в бытовых целях, не применяет его в своей профессиональной деятельности
Не использую	1-2 года	" Интернет Интернет"	Представители сегмента освоили сеть Интернет и легко в ней ориентируются, но их профессиональная деятельность либо лежит за пределами сети Интернет, либо они не знают возможностей использования Интернет в своей профессиональной деятельности
Использую крайне редко	1-2 года	" Справка на работе Справка на работе"	Представители сегмента представляют себе возможности Интернет и используют его на работе для получения справочной информации в редких случаях. Профессия чаще не связана с информационными потоками
Ежедневно обращаюсь к сети	1-2 года	" Интернет Интернет" в работе	Представители сегмента представляют себе возможности Интернет и активно используют его на работе для получения профессионально ориентированной информации

Кластерный анализ

Разбиение на три
кластера покупателей
продукции со склада



Данные

Сумма, перечисленная
покупателем за месяц

ЧИС ЗНА	1 ДАТА	2 ПОК1	3 ПОК2	4 ПОК3	5 ПОК4	6 ПОК5	7 ПОК6	8 ПОК7
3	3/1/99	3000.000	3250.000	1000.000	40.000	2522.400	50.000	90.000
4	4/1/99	4000.000	4250.000	3000.000	30.000	3353.200	80.000	150.000
5	5/1/99	5000.000	5250.000	4000.000	50.000	4204.000	130.000	250.000
6	6/1/99	5000.000	5250.000	30.000	80.000	4484.000	410.000	810.000
7	7/1/99	1000.000	1250.000	50.000	130.000	1090.800	220.000	430.000
8	8/1/99	3000.000	3250.000	80.000	30.000	2532.400	30.000	110.000
9	9/1/99	4000.000	4250.000	130.000	50.000	3303.200	50.000	50.000
10	10/1/99	5000.000	5250.000	5105.000	80.000	4124.000	80.000	90.000
11	11/1/99	1000.000	5250.000	5105.000	130.000	4154.000	130.000	150.000
12	12/1/99	3000.000	5000.000	1101.000	120.000	1000.800	130.000	250.000
13	1/1/00	4000.000	1000.000	3103.000	400.000	2882.400	410.000	30.000
14	2/1/00	5000.000	3000.000	4104.000	210.000	3493.200	220.000	50.000

Анализ аналогичный

