

Тема 7

# Статистические методы снижения размерности признакового пространства

Снижение размерности

# <u>Снижение размерности признакового пространства</u> <u>обусловлено</u>:

- □ необходимость наглядного представления данных
- □ стремление к лаконизму исследуемых моделей:
  - отбор наиболее информативных показателей
- ☐ необходимость существенного сжатия хранимой статистических данных



Классификация

Снижение размерности

# Инструментарий снижения размерности признакового пространства:

- □ метод главных компонент
- □ методы факторного анализа
- □ многомерное шкалирование
- □ отбор существенных предикторов в регрессионном анализе
- □ отбор типообразующих признаков в дискриминантном анализе
- ☐ целенаправленное проецирование и отбор типообразующих признаков в кластер-анализе

Снижение размерности

Факторный анализ - совокупность статистических методов, которые на основе реально существующих связей признаков (или объектов) позволяют выявить <u>латентные</u> обобщающие характеристики организационной структуры и механизма развития изучаемых явлений и процессов.

**Снижение размерности** 

$$Z_{j}=a_{j1}F_{1}+a_{j2}F_{2}+\ldots+a_{jn}F_{n},$$

где  $Z_i$  – стандартизированные переменные

$$F_1, F_2, ... F_n$$
 – общие факторы,

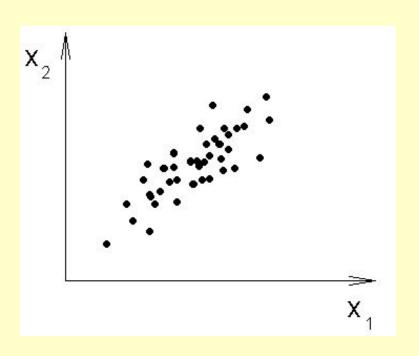
$$Var F_1 \ge Var F_2 \ge ... \ge Var F_n$$

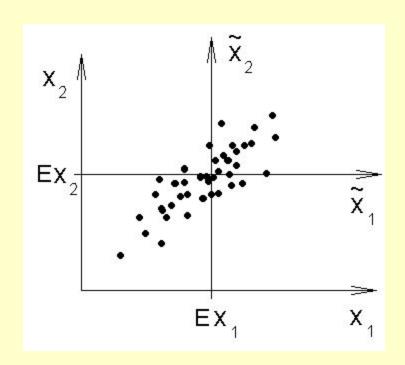
$$a_{i1}, a_{i2}, ..., a_{in}$$
 – факторные нагрузки



Снижение размерности

### Метод главных компонент

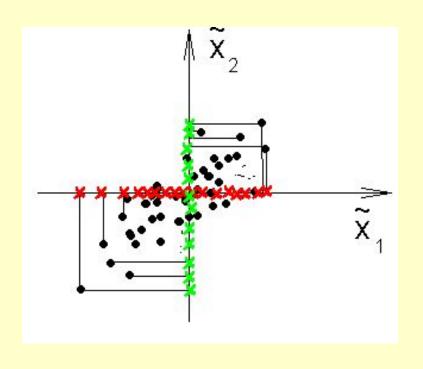


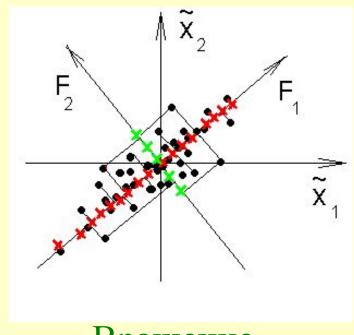


Стандартизация

Снижение размерности

# Метод главных компонент

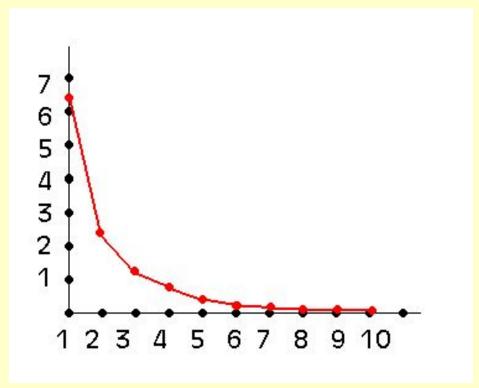




Вращение

Снижение размерности

# Объяснение общей дисперсии



$$\frac{\lambda_1 + \lambda_2 + \lambda_3}{\lambda_1 + \lambda_2 + \lambda_3 + \ldots + \lambda_{10}} = \frac{\text{Var } F_1 + \text{Var } F_2 + \text{Var } F_3}{\text{общая ддисперси}}$$

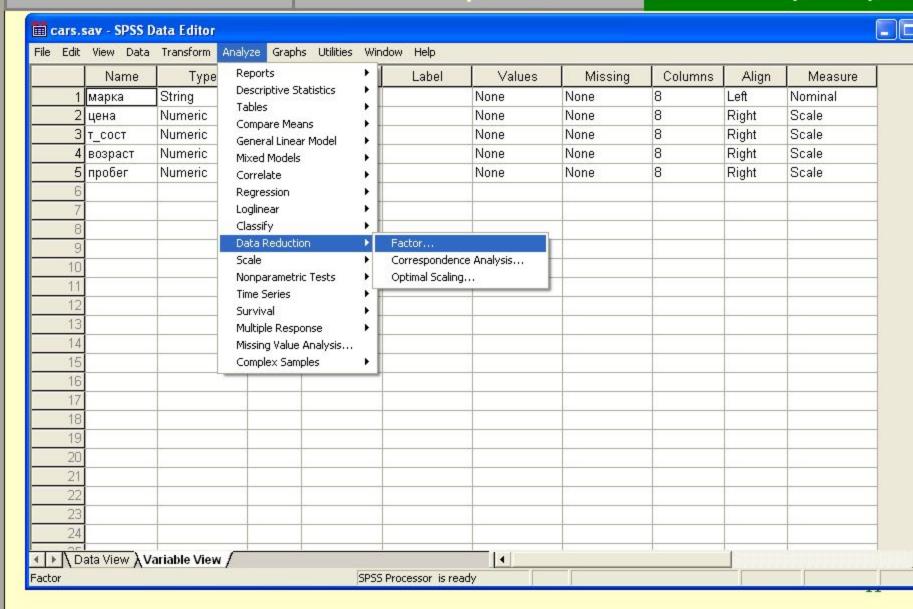
#### Анализ данных в маркетинговых исследовани

Взаимосвязи

# Классификация

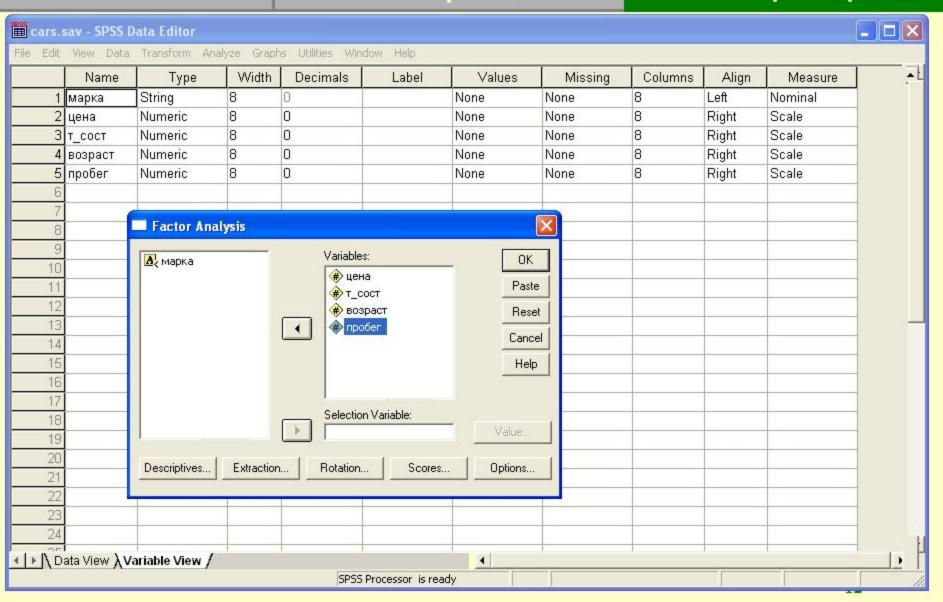
: возраст											
марка		т_сост	возраст	пробег	var						
1 Audi	16350	7	3	66000							
2 BMW	14500	8	4	92500							
3 Buick	8950	9	8	92500		2					
4 Chrysler	8950	6	6	92500							
5 Dodge	8450	5	6	92500							
6 Honda	9850	7	4	118500							
7 Mazda	12650	7	8	58000							
8 Mercede:	17250	8	5	92500							
9 Mitsub.	8950	3	5	136000							
10 Nissan	9850	6	4	150000							
11 Pontiac	8950	7	6	110000							
12 Saab	14950	7	4	101000							
13 Toyota	11700	8	7	57500							
14 VW	8450	5	6	110000							
15 Volvo	13100	8	4	75000							
16											
17											
18											
19											
20											
21				-							

#### Классификация



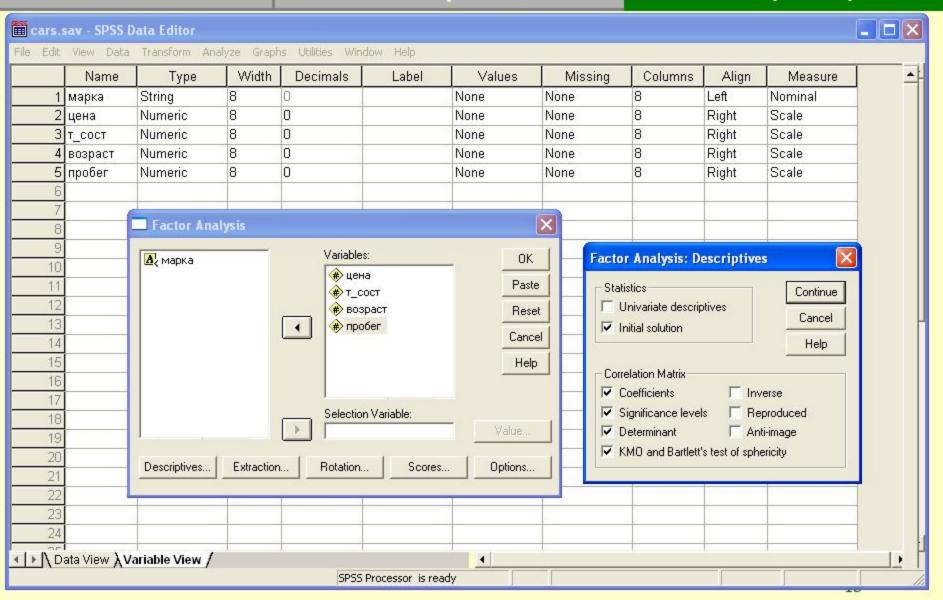


#### Классификация



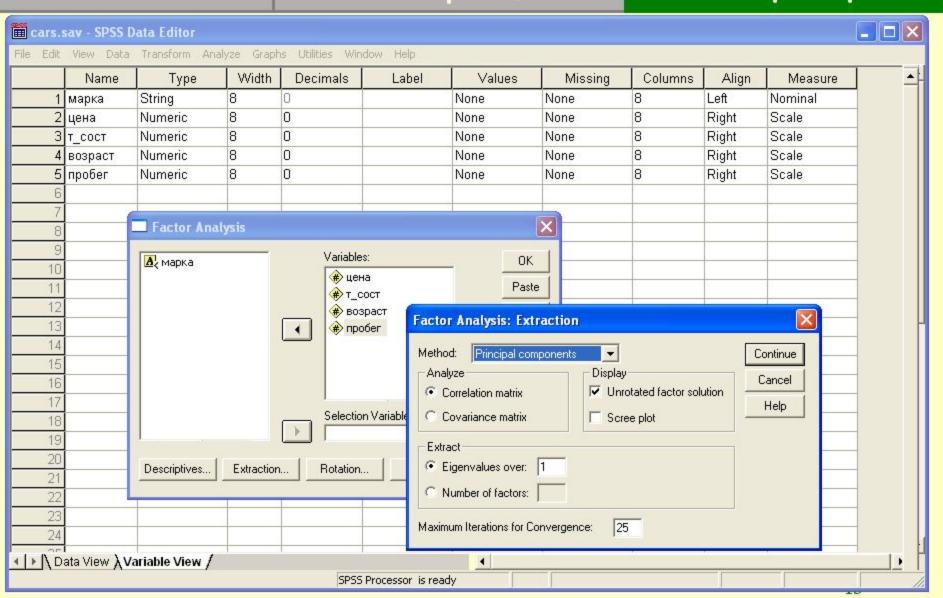


#### Классификация



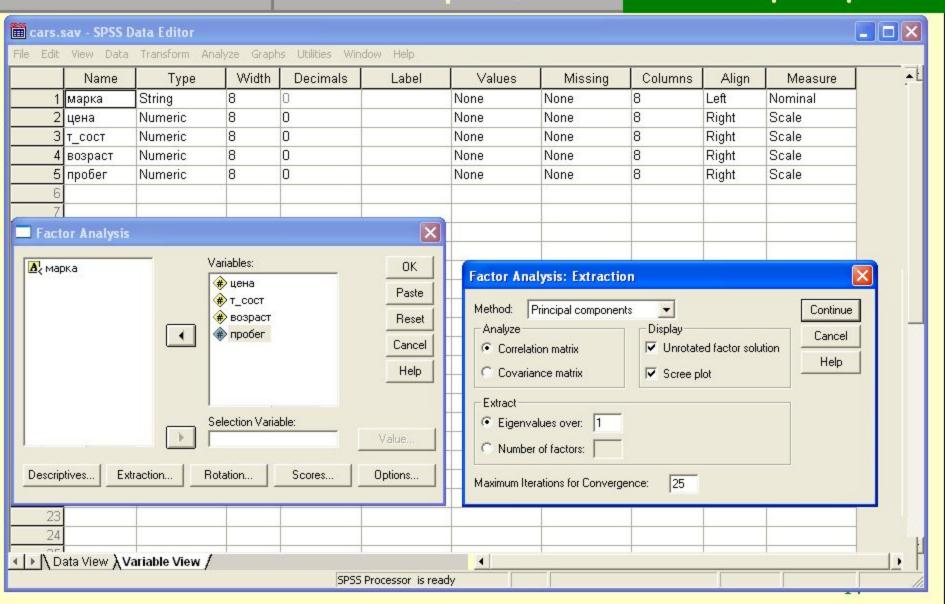


# Классификация



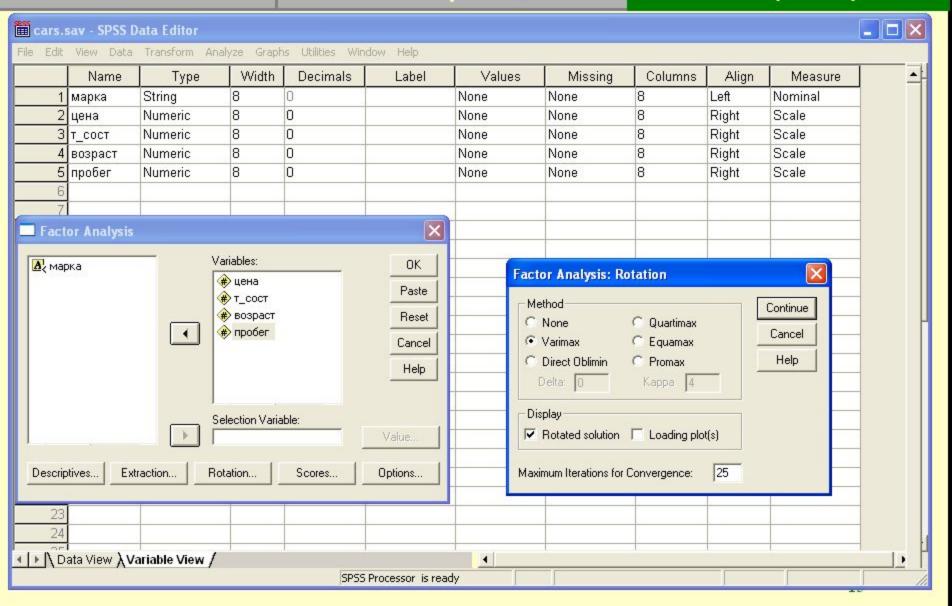


#### Классификация



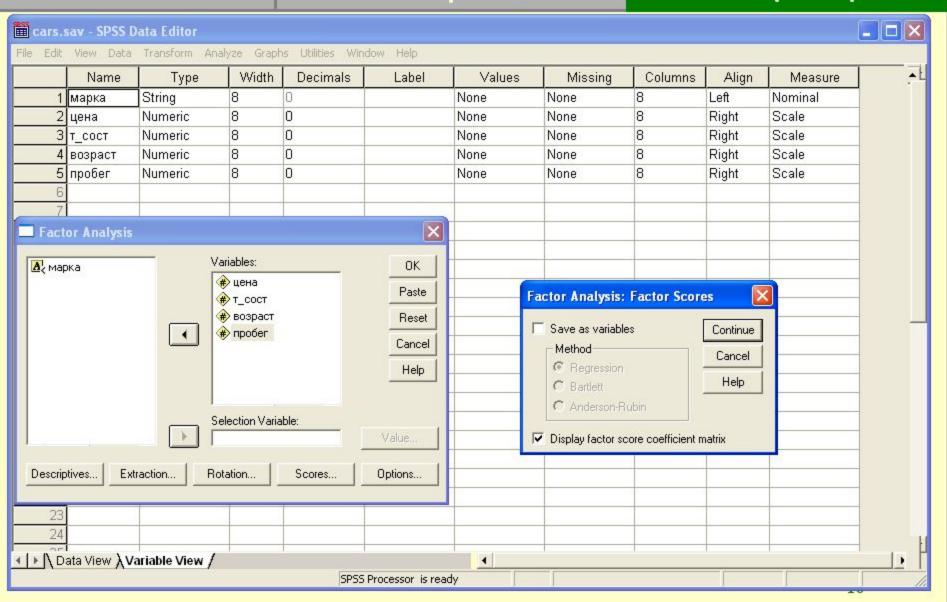


#### Классификация



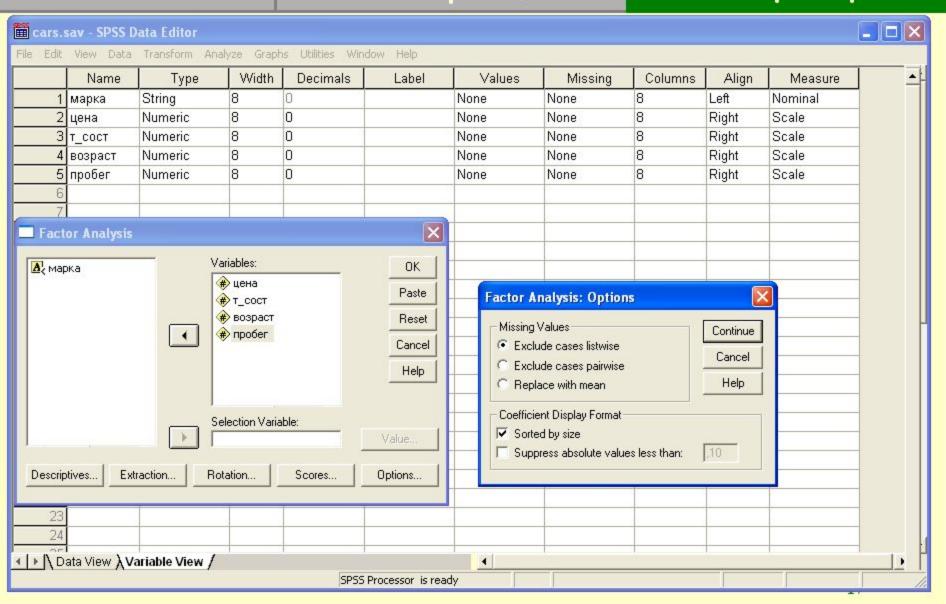


#### Классификация



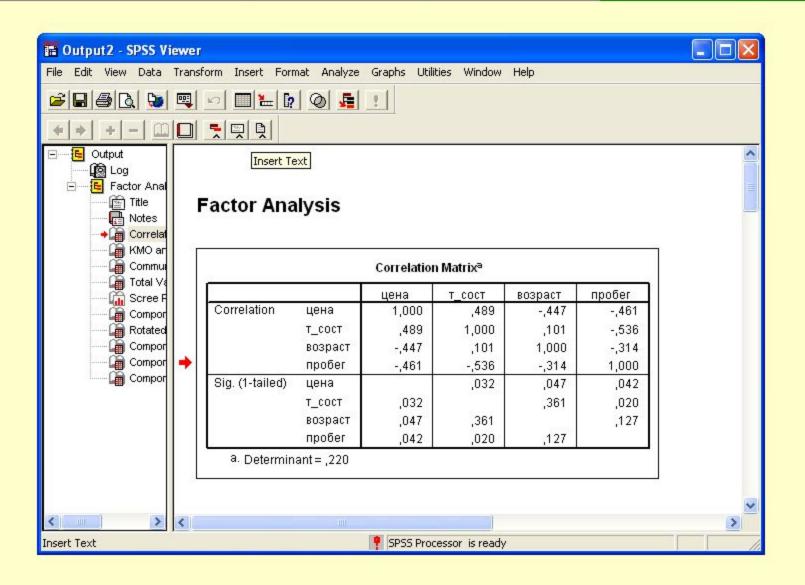


#### Классификация



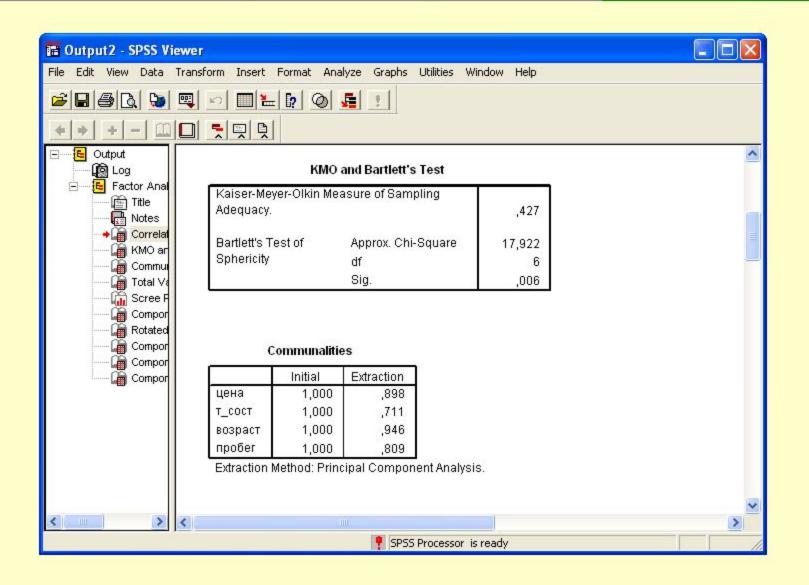


#### Классификация



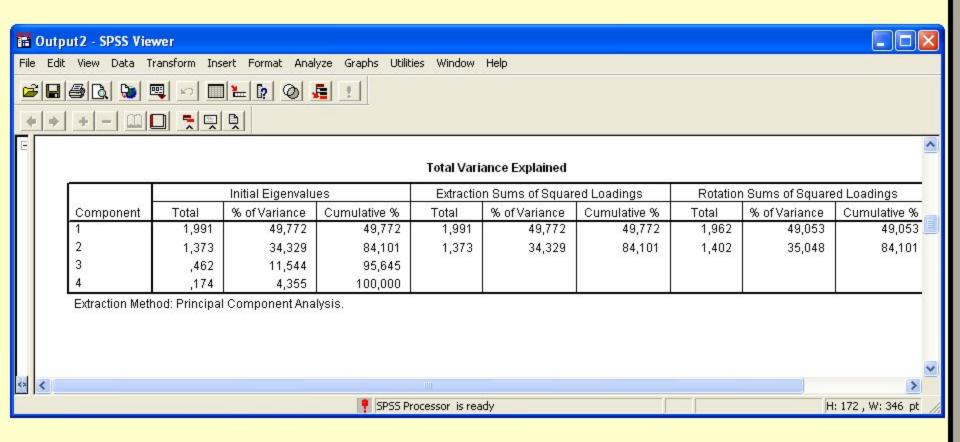


#### Классификация



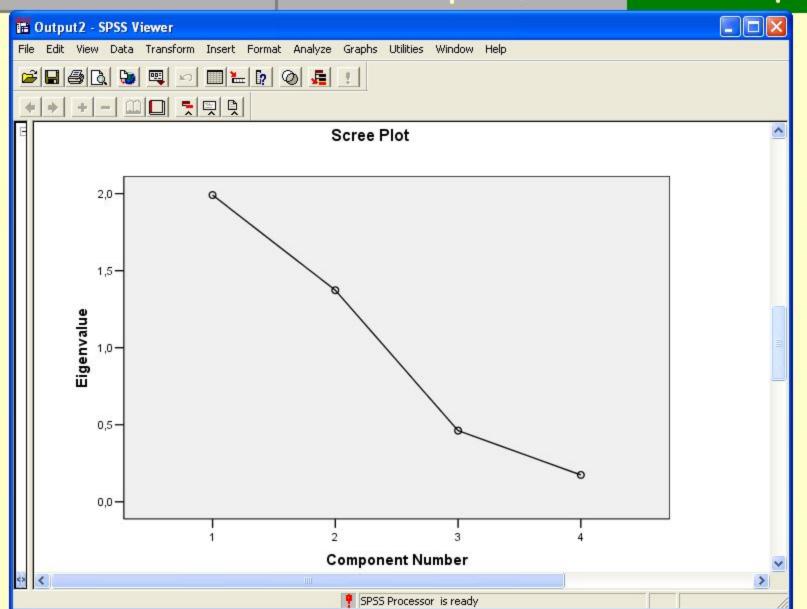


#### Классификация



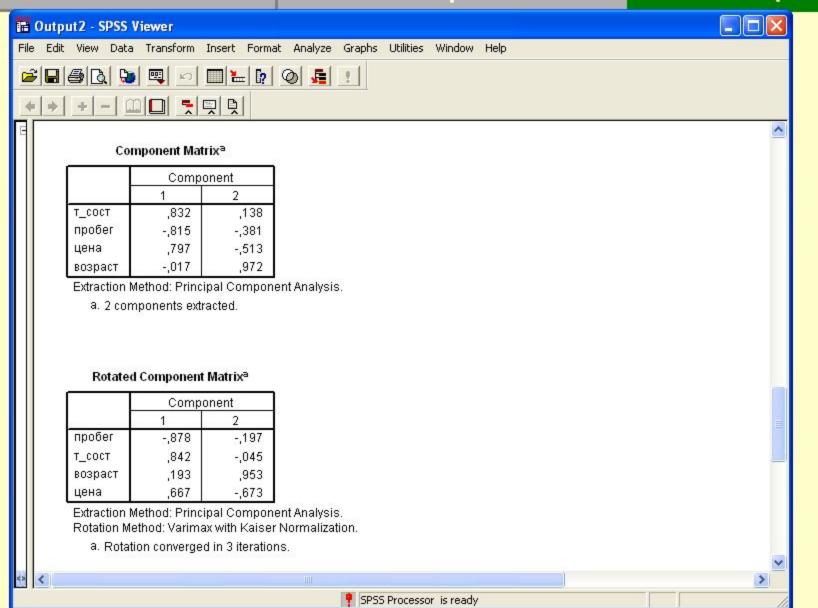


#### Классификация





# Классификация



Rotation Method: Varimax with Kaiser Normalization. Extraction Method: Principal Component Analysis.

1	1,000	,000
Component	1,000	3,000

#### **Component Score Covariance Matrix**

Rotation Method: Varimax with Kaiser Normalization. Extraction Method: Principal Component Analysis.

1	,976	-,216
Component	<sub>1</sub> ,216	2,976

**Component Transformation Matrix** 

Rotation Method: Varimax with Kaiser Normalization. Extraction Method: Principal Component Analysis.

цена	. 310 Comp	onent -,451
пробег	-,460	-,183
возраст	,144	,693
т_сост	<sub>1</sub> ,430	,008

**Component Score Coefficient Matrix** 

#### Классификация

Снижение размерности

# Простые методы факторного анализа:

- □ однофакторная модель Ч. Спирмена: выделяется один латентный и один характерный факторы. Для возможно существующих других латентных факторов делается предположение об их незначительности
- □ **бифакторная модель Г.Хользингера:** допускается влияние на вариацию элементарных признаков не одного, а нескольких латентных факторов (обычно двух) и одного характерного фактора

□ **центроидный метод Л. Тэрстоуна:** корреляция между переменными рассматривается как пучок векторов, а латентный фактор геометрически представляется как уравновешивающий вектор, проходящий через центр этого пучка



#### Классификация

Снижение размерности

# Аппроксимирующие методы факторного анализа:

□ главных факторов
□ групповой
□ минимальных остатков
□ канонический
□ распознавания образов
□ максимального правдоподобия



Классификация

Снижение размерности

# Основные стадии факторного анализа:

- □ Вычисление корреляционной матрицы для всех переменных, участвующих в анализе; если переменная слабо коррелированна с остальными, следует подумать о ее исключении при следующем запуске программы (учитывая при этом ее общность и нагрузки).
- ☐ Извлечение факторов; оцениваются нагрузки факторов; выбирается либо метод главных компонент, либо один из методов факторного анализа.
- □ Вращение факторов для создания упрощенной структуры; вращение увеличивает, либо уменьшает нагрузки на каждый фактор; просмотр результатов позволяет уменьшить первоначальное число факторов.
- □ Интерпретация факторов.