

*Проекционные
методы.
Основные понятия и
примеры*

Институт химической физики РАН, Москва

Родионова Оксана Евгеньевна

План лекции

1. Введение

- Немного истории
- Природа многомерна
- Пример - многомерный статистический контроль процессов
- Два подхода к анализу данных

2. Идеи, заложенные в проекционном подходе

- Данные – какие они бывают
- Классы решаемых задач

- **3. Метод главных компонент, основные понятия и примеры**

Метод наименьших квадратов (простейший случай)

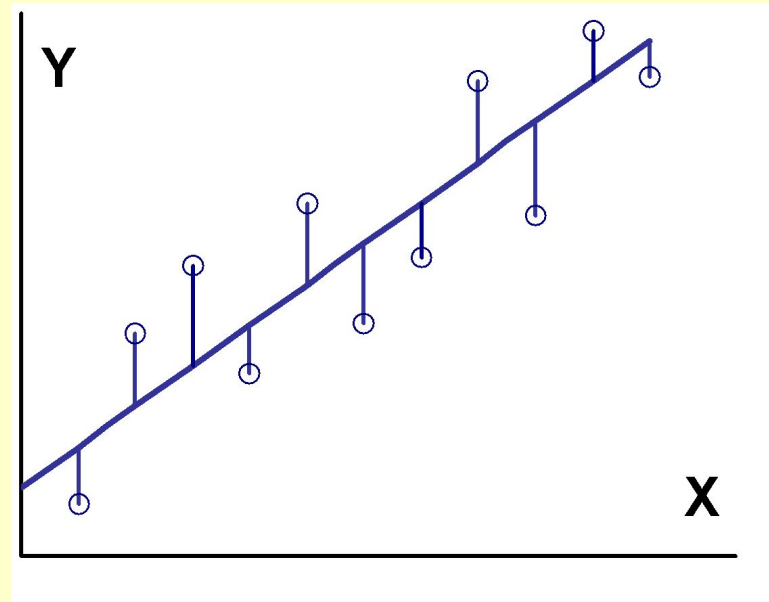
$$y = \alpha + \beta x$$

$$e_i$$

$$a \cong \alpha$$

$$b \cong \beta$$

x_1	y_1
x_2	y_2
.	.
.	.
.	.
...	...
x_n	y_n

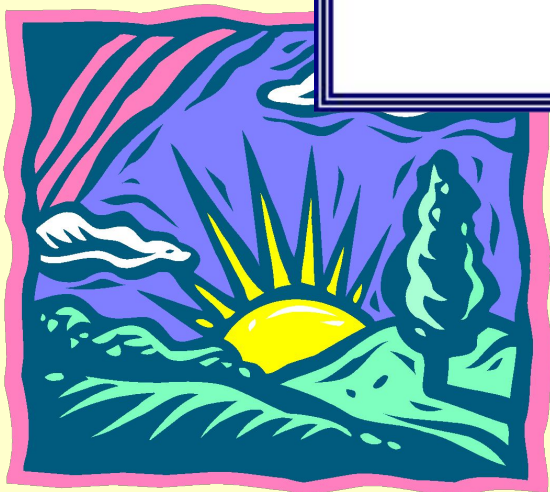


$$\sum_{i=1}^n (y_i - a - bx_i)^2 = \min$$

Многомерные данные



**Исследуемое свойство
очень редко зависит
ТОЛЬКО ОТ ОДНОЙ
переменной**



Контроль производственного процесса

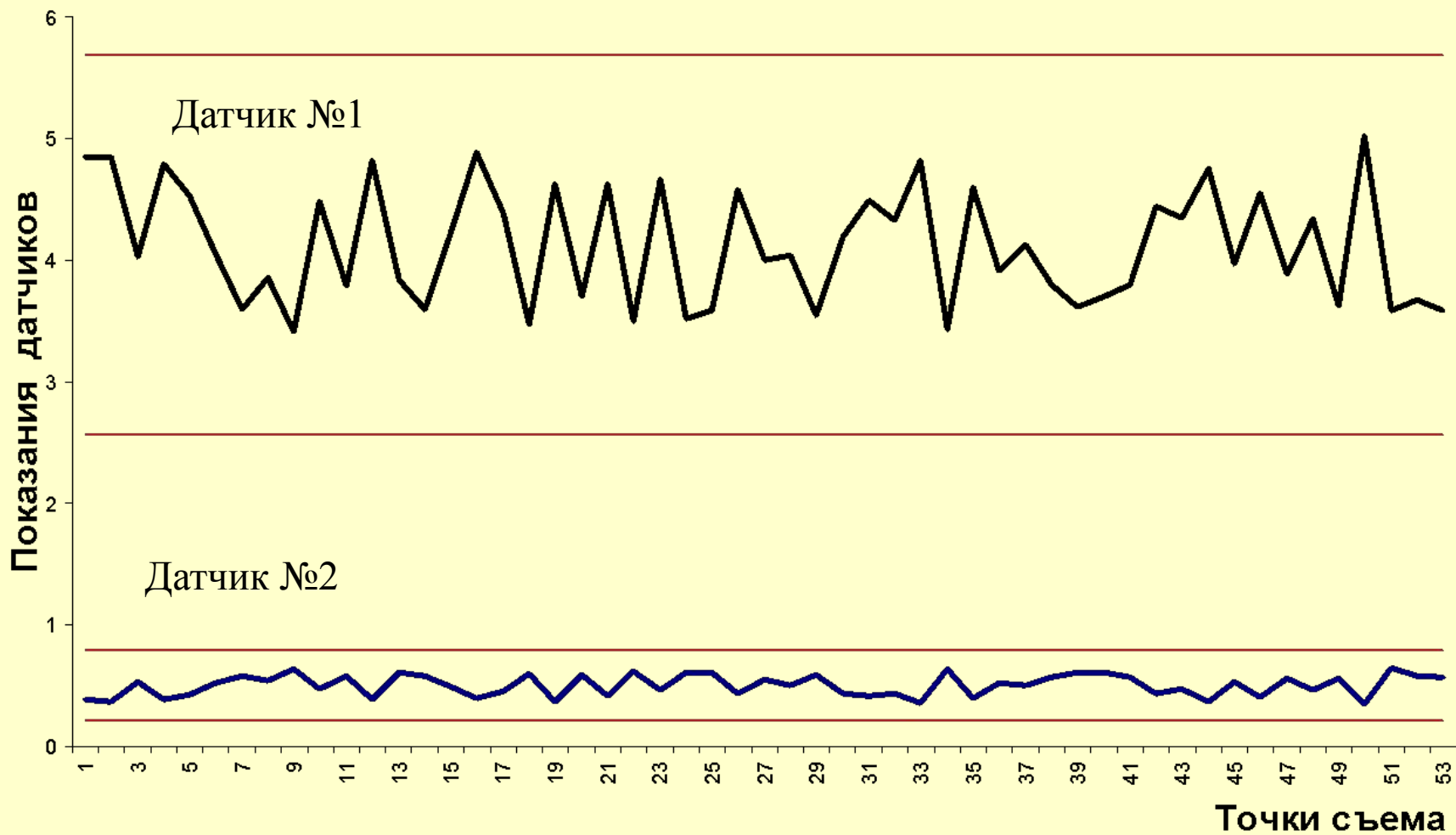
	X1	X2	X3	X4	X5	X6	X7	X8	X9	...	X17
s1	-1.19E-01	7.28E-01	-2.15E-02	5.22E-01	7.06E-04	7.32E-01	3.10E-04	-6.13E-04	-5.92E-05	• • •	9.74E-03
s2	-1.37E-01	7.28E-01	-2.89E-02	6.08E-01	7.09E-04	7.02E-01	6.58E-04	-1.22E-03	-1.49E-04		1.01E-02
s3	2.51E-02	-9.15E-02	6.73E-03	-1.13E-01	-9.07E-05	-7.58E-02	-2.29E-04	4.10E-04	5.65E-05		-1.43E-03
s4	-1.14E-01	6.70E-01	-2.18E-02	5.04E-01	6.50E-04	6.65E-01	3.83E-04	-7.34E-04	-7.96E-05		9.07E-03
s5	-7.93E-02	4.14E-01	-1.69E-02	3.51E-01	4.04E-04	3.98E-01	3.96E-04	-7.35E-04	-9.05E-05		5.78E-03
s6	1.51E-02	-6.38E-02	3.74E-03	-6.75E-02	-6.28E-05	-5.67E-02	-1.15E-04	2.07E-04	2.78E-05		-9.49E-04
s7	7.44E-02	-5.24E-01	1.11E-02	-3.24E-01	-5.06E-04	-5.45E-01	-1.73E-05	7.92E-05	-1.07E-05		-6.79E-03
s8	3.65E-02	-2.66E-01	5.12E-03	-1.59E-01	-2.56E-04	-2.78E-01	1.43E-05	-3.95E-07	-1.14E-05		-3.42E-03
s9	1.36E-01	-7.06E-01	2.89E-02	-6.01E-01	-6.88E-04	-6.77E-01	-6.83E-04	1.26E-03	1.56E-04		-9.86E-03
s10	-2.74E-02	3.60E-01	1.82E-03	1.12E-01	3.42E-04	4.12E-01	-4.31E-04	7.24E-04	1.22E-04		4.18E-03
s11	7.47E-02	-3.31E-01	1.80E-02	-3.34E-01	-3.25E-04	-2.99E-01	-5.30E-04	9.62E-04	1.28E-04		-4.84E-03
s12	-1.17E-01	7.02E-01	-2.16E-02	5.13E-01	6.81E-04	7.03E-01	3.40E-04	-6.63E-04	-6.76E-05		9.44E-03
s13	1.06E-01	-2.82E-01	3.23E-02	-4.82E-01	-2.85E-04	-1.87E-01	-1.25E-03	2.21E-03	3.14E-04		-4.99E-03
s14	7.39E-02	-5.28E-01	1.07E-02	-3.21E-01	-5.09E-04	-5.50E-01	2.49E-06	4.48E-05	-1.59E-05		-6.81E-03
s15	-9.87E-03	1.02E-01	-3.21E-04	4.17E-02	9.75E-05	1.13E-01	-8.29E-05	1.36E-04	2.44E-05		1.23E-03
s16	-1.06E-01	7.68E-01	-1.52E-02	4.62E-01	7.41E-04	8.03E-01	-2.54E-05	-2.68E-05	2.88E-05		9.90E-03
s17	-4.76E-02	2.66E-01	-9.52E-03	2.10E-01	2.59E-04	2.61E-01	1.92E-04	-3.61E-04	-4.19E-05		3.65E-03
	• • •										
s54	6.61E-02	-5.40E-01	7.19E-03	-2.85E-01	-5.19E-04	-5.78E-01	1.81E-04	-2.67E-04	-6.23E-05		-6.78E-03

Цель исследования

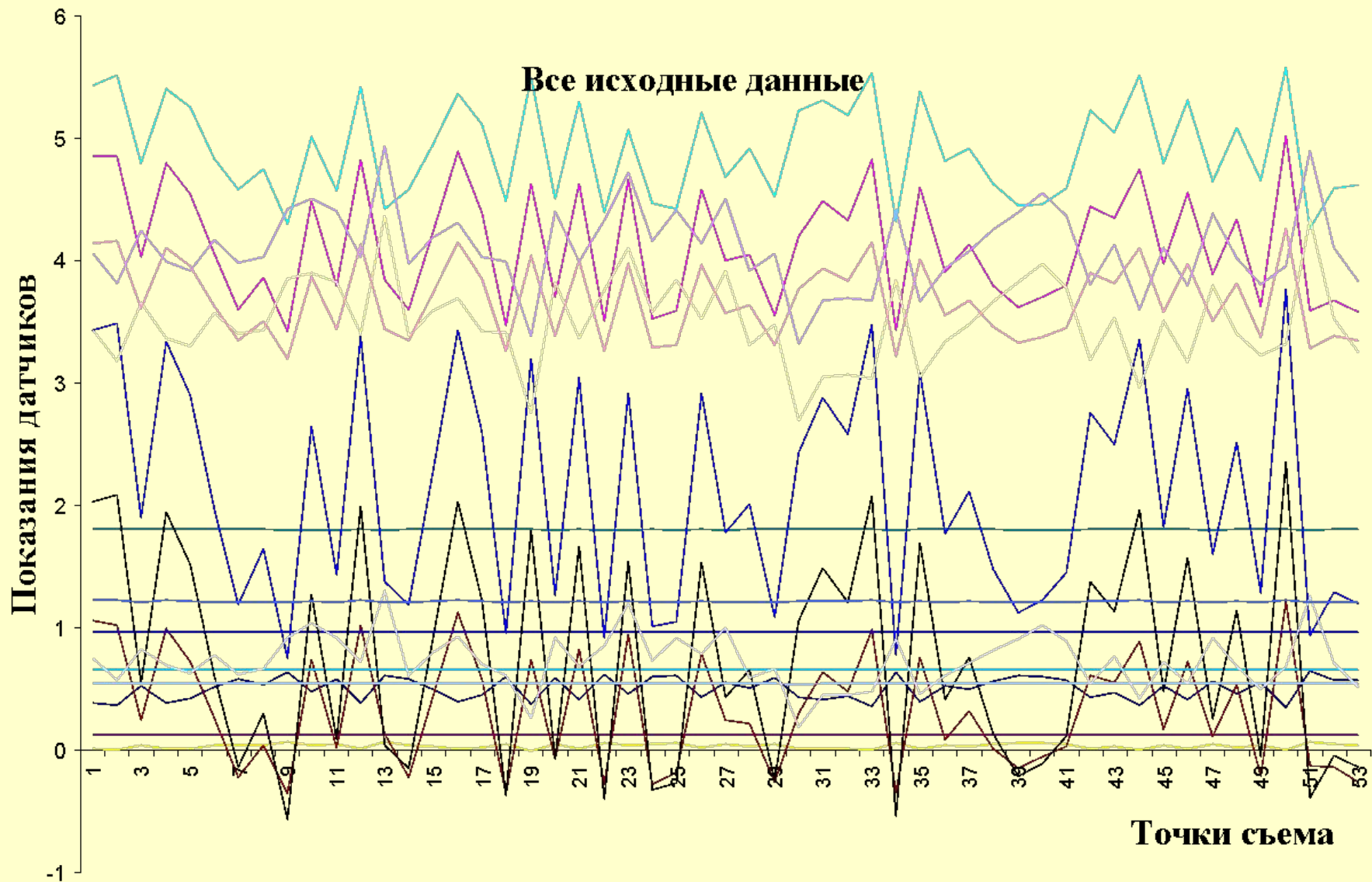
Контроль производства

Возможность воздействовать на процесс для его стабилизации

Контроль производственного процесса

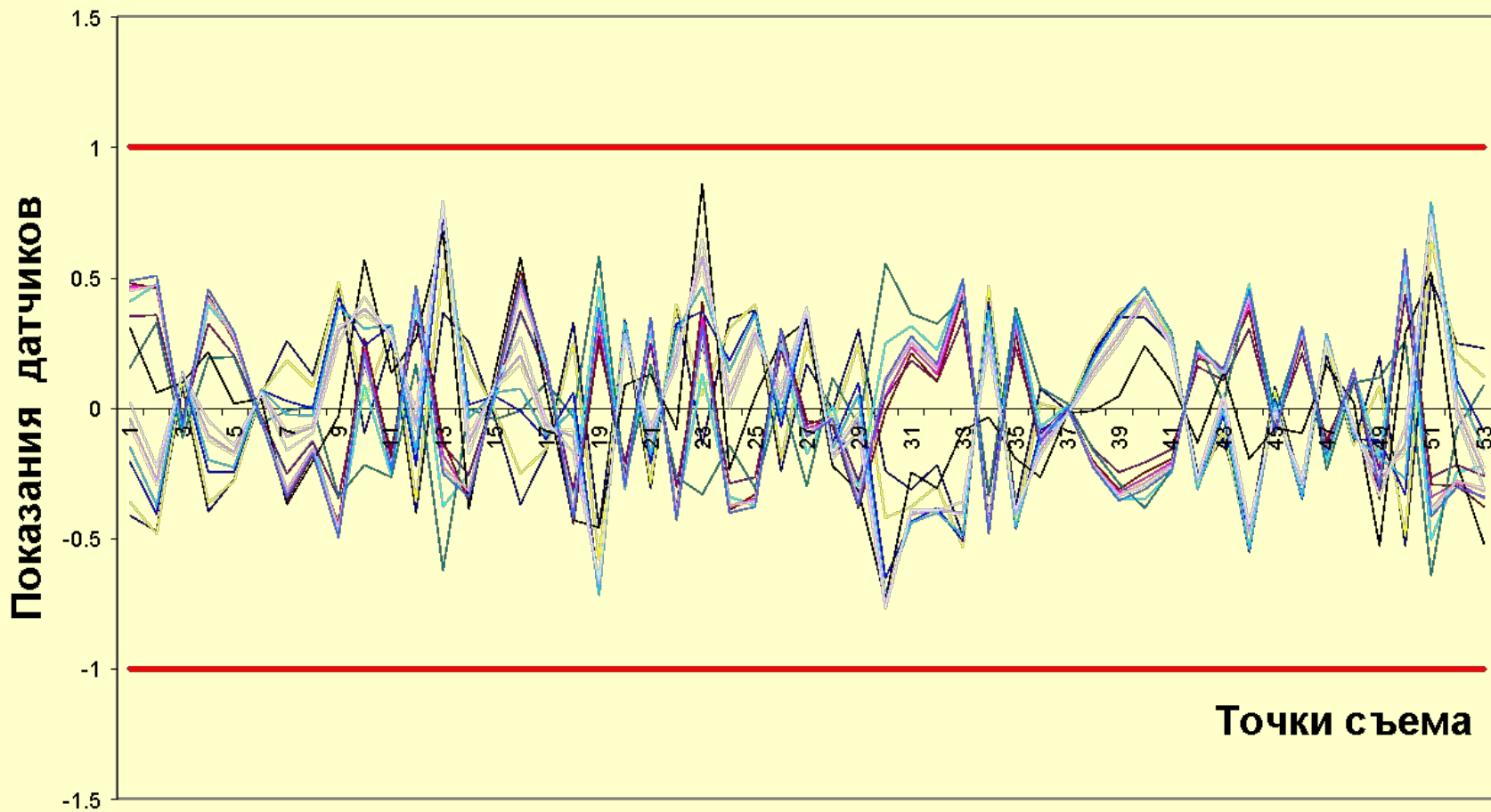


Контроль производственного процесса

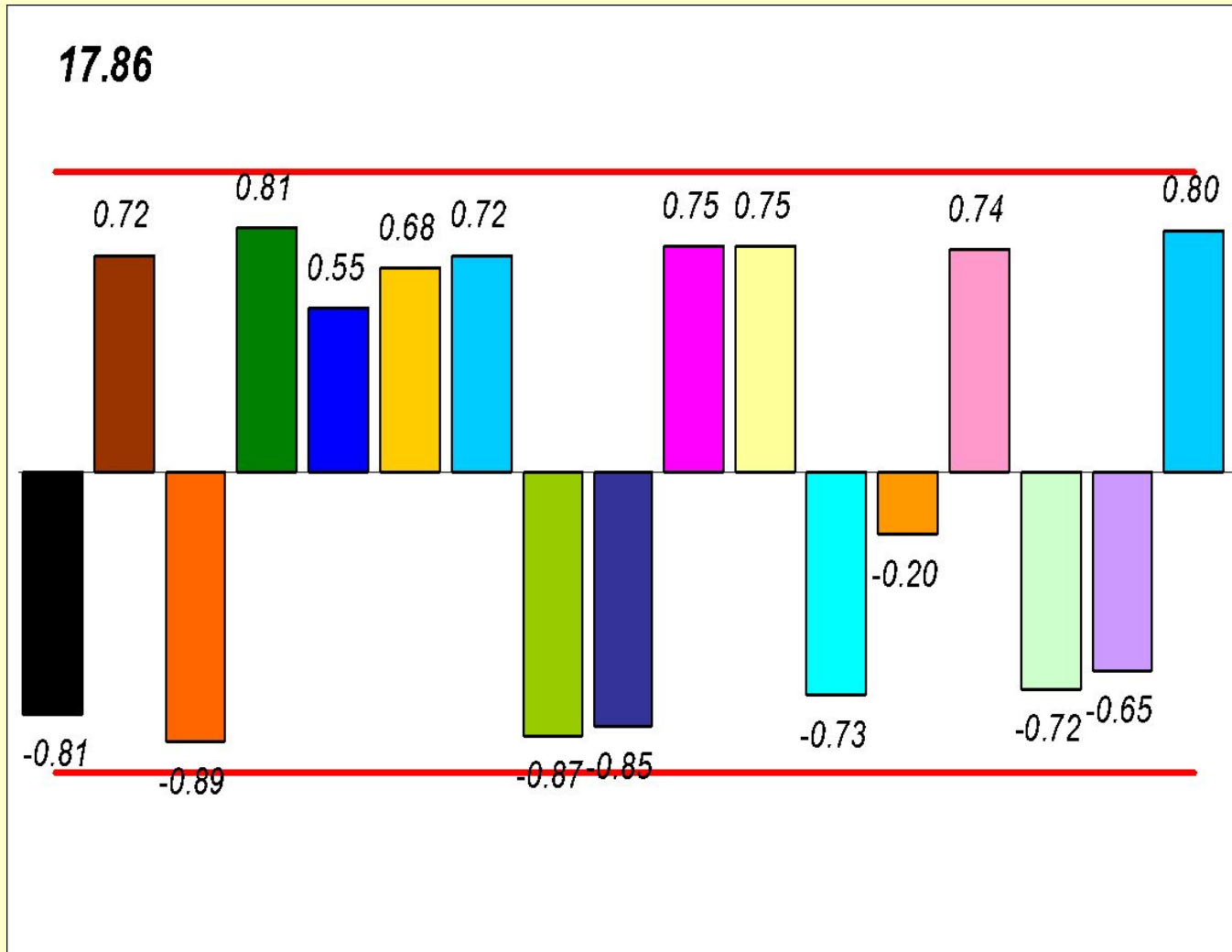


Контроль производственного процесса

Нормированные данные

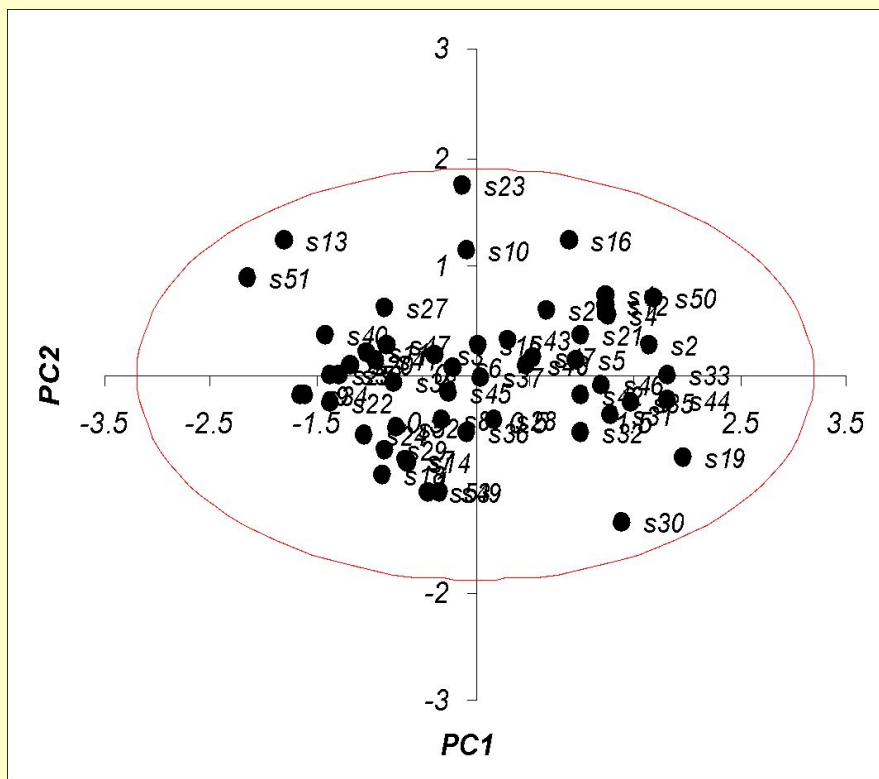


Контроль производственного процесса



Контроль производственного процесса

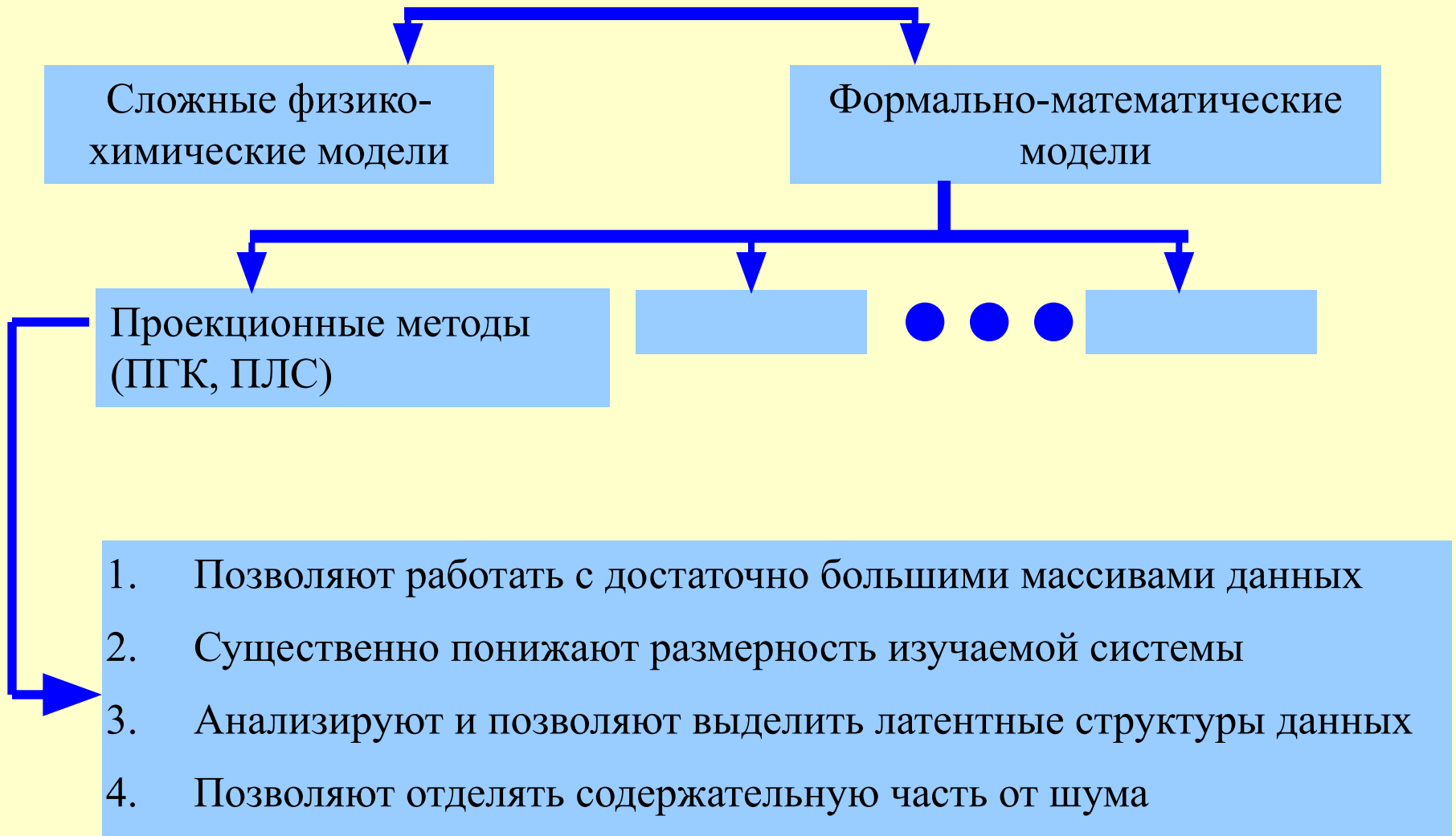
Точки съема



Моделирование
производилось на основе
анализа измерений и
внутренних связей
присущих этому набору
данных

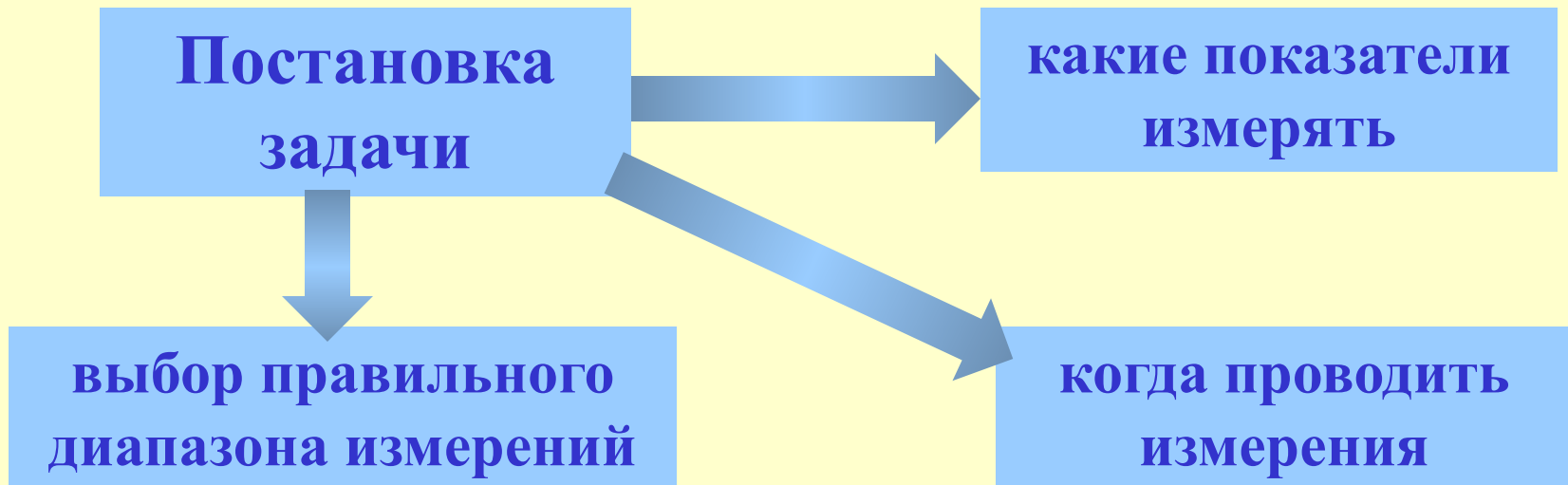
Не применялись
содержательные физико-
химические модели

Моделирование многомерных данных (процессов или явлений)

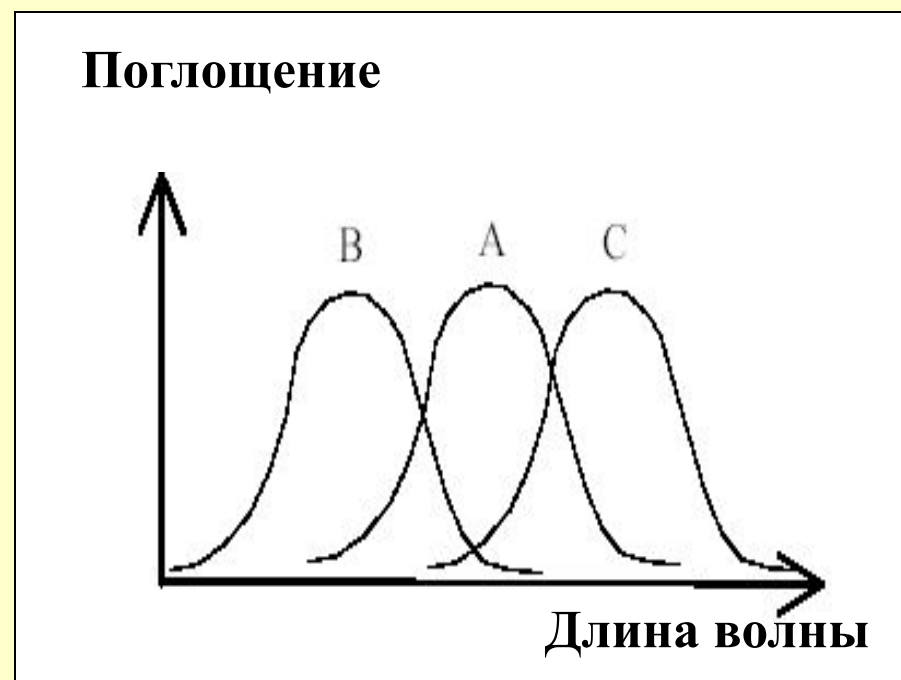
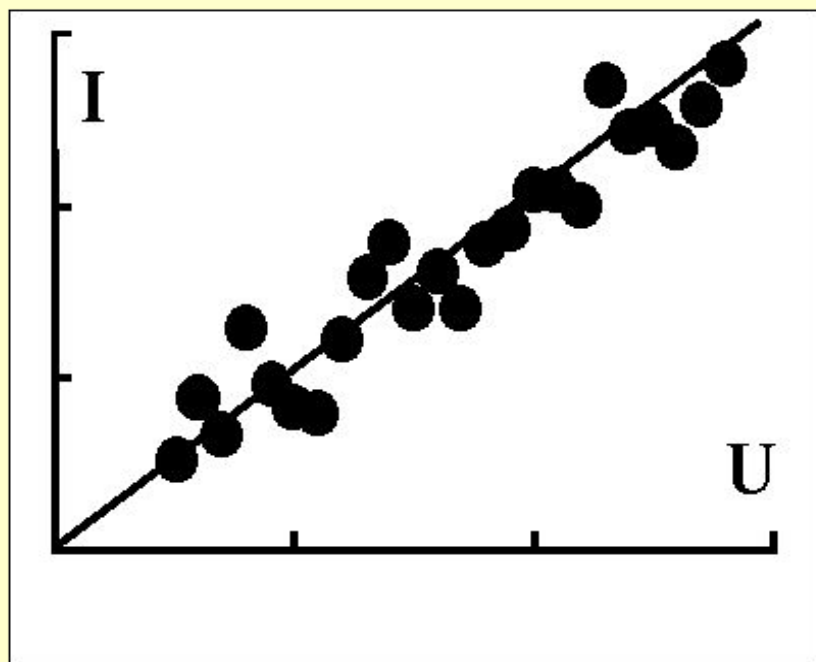
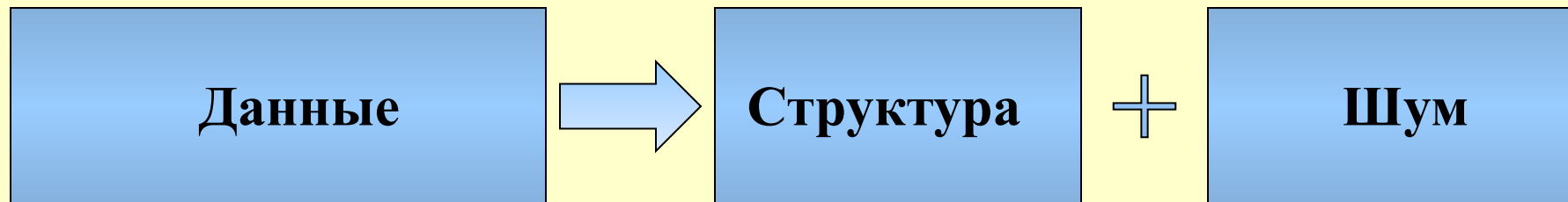


Содержательная составляющая задачи.

Никакие многомерные методы не помогут, если данные не содержат полезной информации об изучаемом свойстве

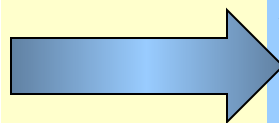


Данные



Два класса решаемых задач

X_{11}	X_{12}		...	X_{1m}
X_{21}	X_{22}		...	
		X		
·	·		·	·
·	·		·	·
·	·		·	·
...
X_{n1}				X_{nm}



Метод- МГК

Задачи

1. Анализ структуры, поиск латентных переменных
2. Классификация и дискриминация

Y_{11}	Y_{12}	...	Y_{1k}
Y_{21}	Y_{22}	...	
		Y	
·	·	·	·
·	·	·	·
·	·	·	·
...
Y_{n1}			Y_{nk}

Методы : РГК, ПЛС

Задачи

1. Построение модели $Y(X)$
2. Прогнозирование

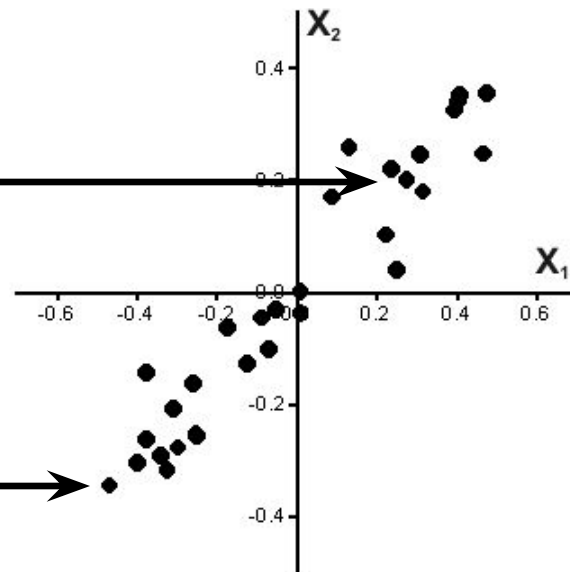
n – количество образцов

m – количество переменных (факторов)

Проекционные методы

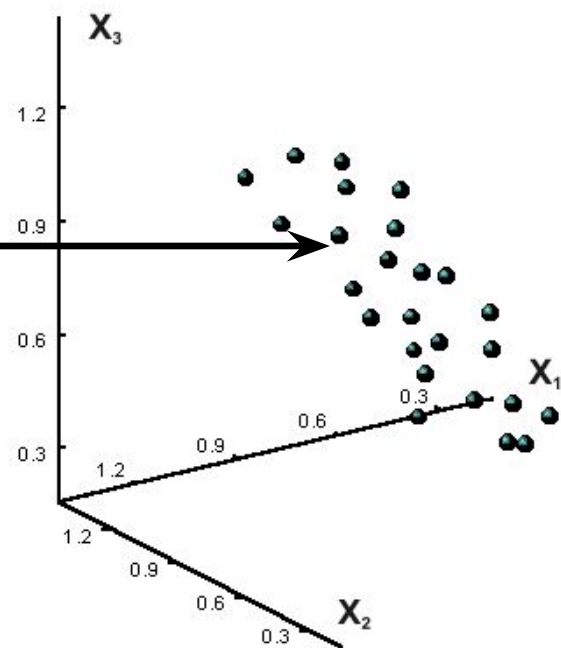
$P=2$

	X_1	X_2
1	0.407	0.353
2	0.475	0.355
3	0.274	0.202
4	0.394	0.325
5	-0.088	-0.045
6	-0.053	-0.031
7	-0.253	-0.253
8	-0.124	-0.128
9	-0.251	-0.255
10	0.088	0.171
11	-0.261	-0.162
12	0.401	0.341
13	-0.469	-0.344
14	-0.376	-0.143



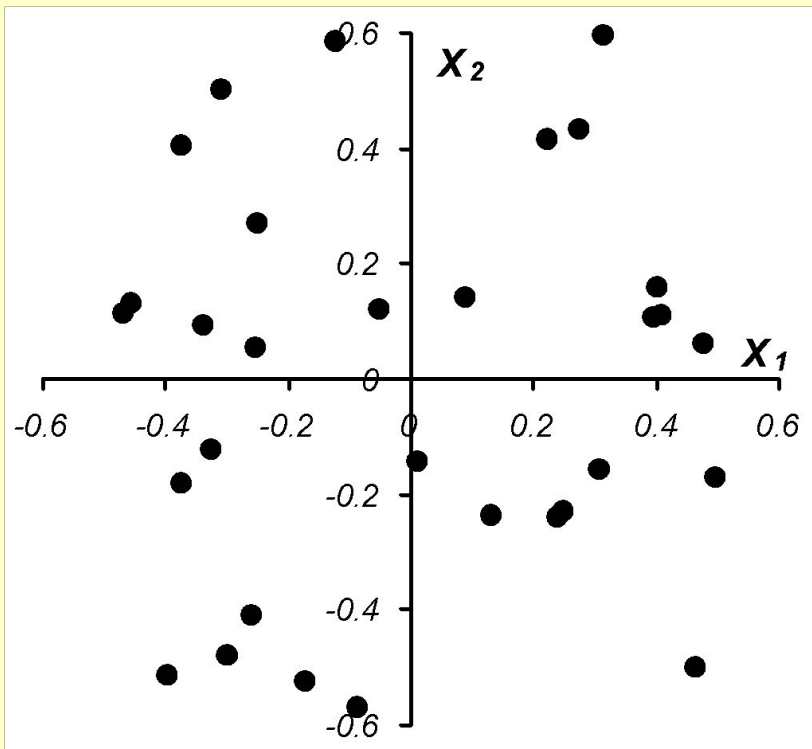
$P=3$

	X_1	X_2	X_3
1	0.631	0.421	0.504
2	0.663	0.537	0.510
3	0.544	0.825	0.637
4	0.662	0.954	0.736
5	0.581	1.178	0.866
6	0.758	0.338	0.482
7	0.679	0.611	0.634
8	0.644	0.870	0.744
9	0.713	1.030	0.756
10	0.748	1.166	0.914
11	0.787	0.372	0.482
12	0.820	0.635	0.678
13	0.773	0.831	0.676
14	0.735	0.964	0.861

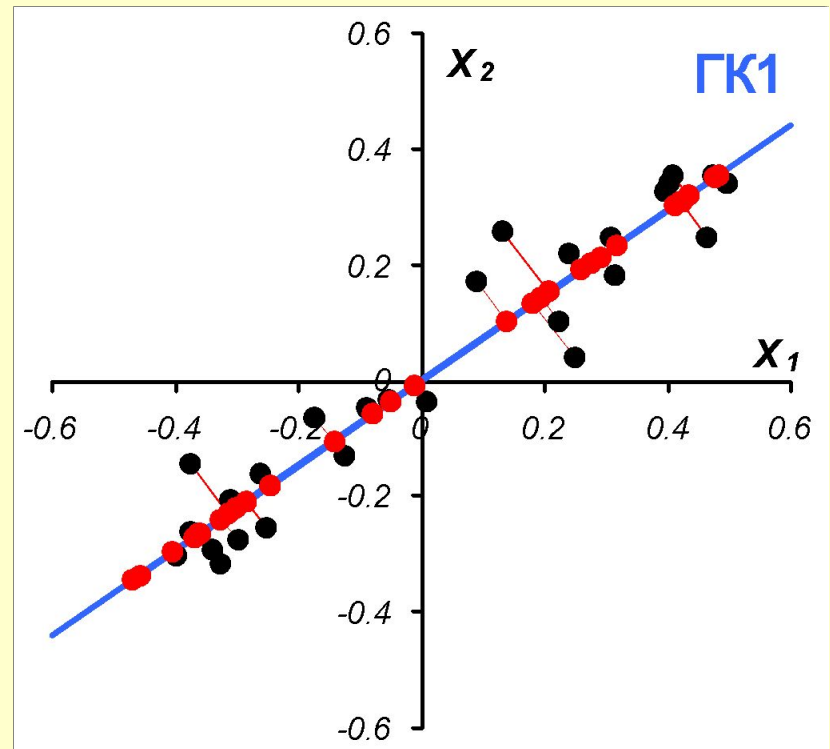


Проекционные методы

Данные без структуры



Данные со скрытой структурой



$$X_2 = aX_1 + E$$

Проекция на подпространство

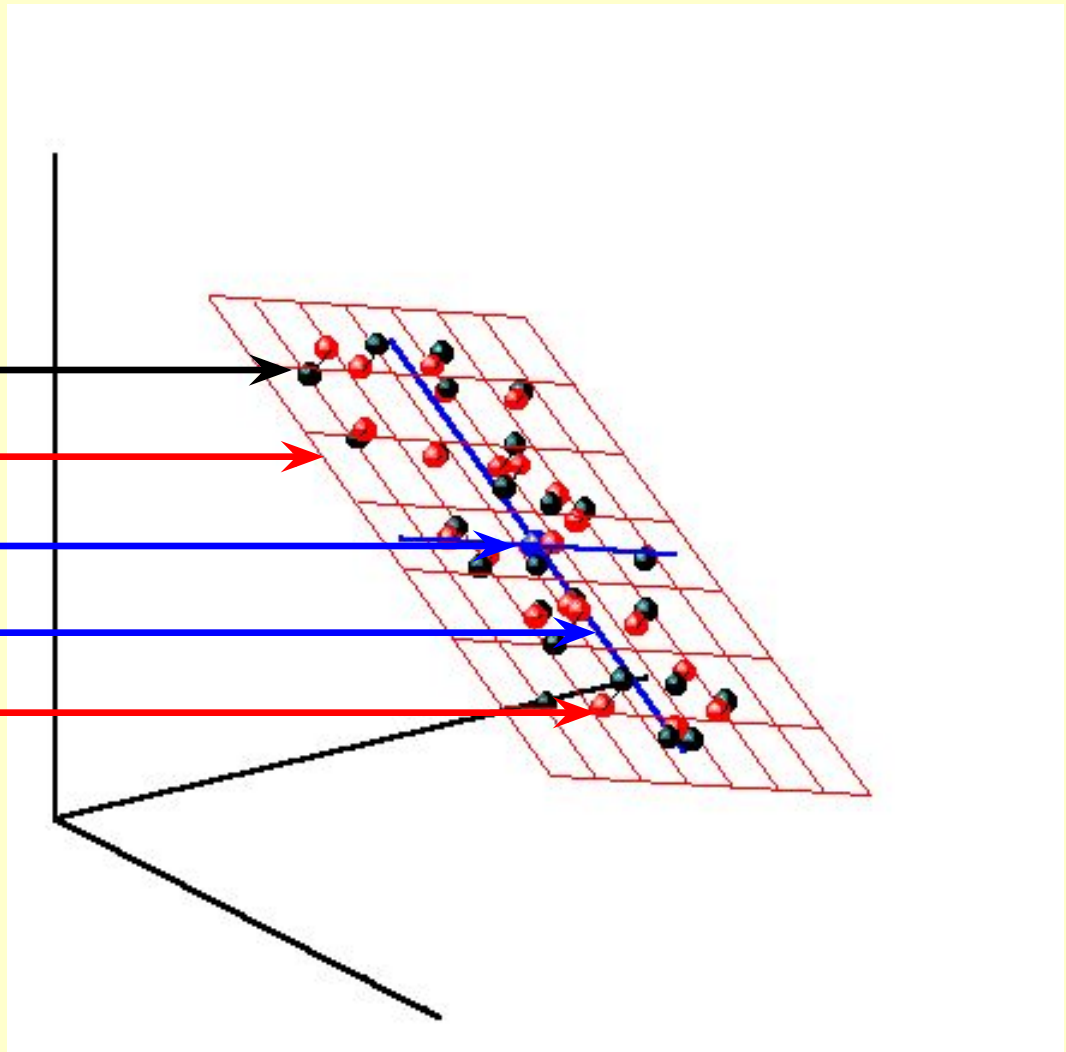
Исходные данные

Подпространство

Центр данных

Главные компоненты

Проекции данных

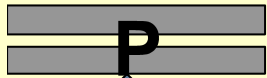
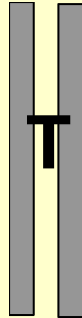
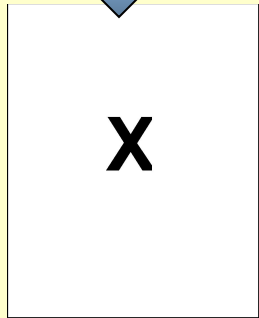


Метод главных компонент (PCA)

Исходные
данные

Матрица счетов
(*Scores*)

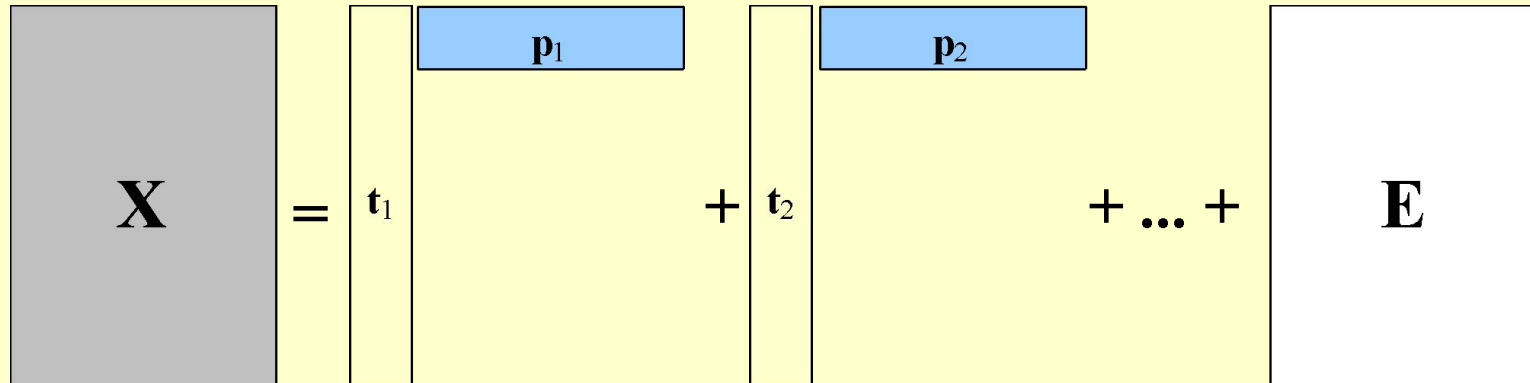
Матрица
ошибок



Матрица нагрузок
(*Loadings*)

$$X = T * P^T + E$$

Метод главных компонент



$$\mathbf{t} = \mathbf{X}\mathbf{p} \Leftarrow \max |\mathbf{X}\mathbf{p}|^2 \text{ при условии } |\mathbf{p}|=1 \Leftrightarrow \mathbf{X}^t \mathbf{X}\mathbf{p} = \lambda \mathbf{p} ; \mathbf{t}^T \mathbf{t} = \lambda$$

\mathbf{X} - матрица данных, \mathbf{E} - матрица ошибок, обе $(n \times p)$

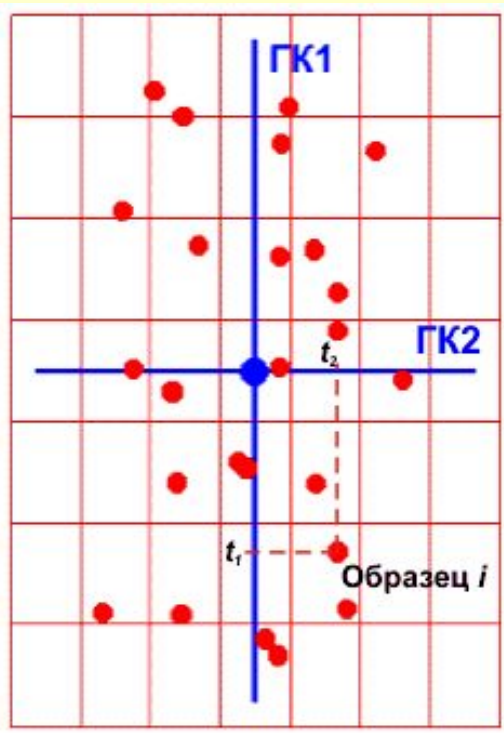
\mathbf{T} - матрица счетов: $(n \times k)$, \mathbf{P} - матрица нагрузок: $(k \times p)$

k - число главных компонент ($k \ll p$)

Karl Pearson, 1901

Матрица счетов T (scores)

$$X = T * P^T + E$$



Строка –
координаты одного
объекта в новой
системе координат



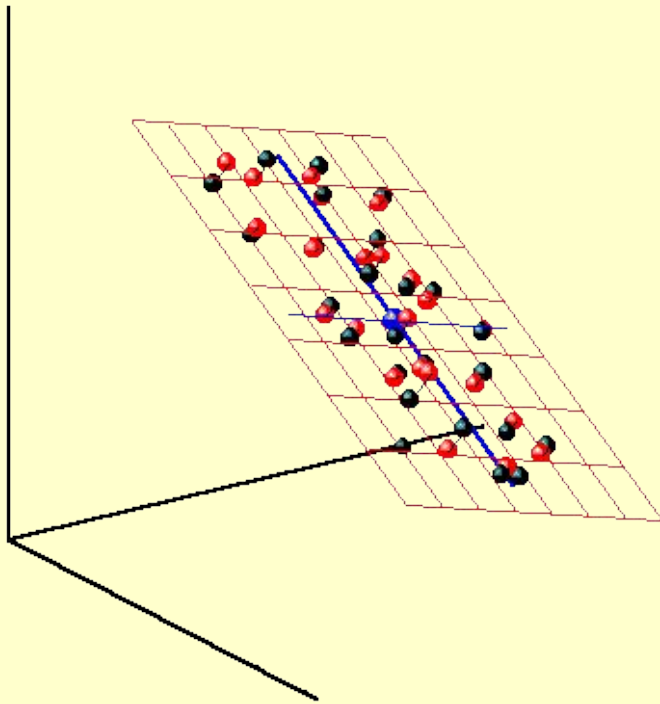
t_{11}	t_{12}
t_{21}	t_{22}
.	.
.	.
.	.
...	...
t_{n1}	t_{n2}

Столбец – проекция
всех объектов на одну
ось главных
компонент

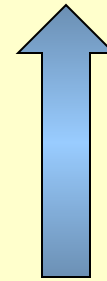


Матрица нагрузок P (loadings)

$$X = T * P^T + E$$



p_{11}	p_{12}		...	p_{1m}
p_{21}	p_{22}		...	p_{2m}



P^T - матрица перехода из пространства X в пространство главных компонент

Остатки E

$$X = T * P^T + E$$

e_{11}	e_{12}		...	e_{1m}
e_{21}	e_{22}		...	
		E		
·	·		·	·
·	·		·	·
·	·		·	·
...
e_{n1}				e_{nm}

$$e_i^2 = \sum_{k=1}^m e_{ik}^2$$

e_i - определяет расстояние от исходного объекта до подпространства главных компонент

$$e_{tot}^2 = \sum_{i=1}^n e_i^2$$

совокупная ошибка для всех объектов

матрица E имеет ту же структуру что и X

E_0, E_1, \dots

E_0 - ошибка при 0 -м ГК, т.е. центрированная матрица X

Математическое обеспечение

Специальные пакеты

UNSCRAMBLER

SIMCA

PLS -ToolBox для MatLab ...

Стандартные статистические пакеты

STATISTICA

SPSS

SAS ...

**Использование стандартного мат. обеспечения
для написания проекционных процедур**

MatLab

Excel+ VBA+.dll

Пример. Демографические данные

	Height	Weight	Hairleng	Shoesize	Age	Income	Beer	Wine	Sex	Swim	A/B	
	1	2	3	4	5	6	7	8	9	10	11	
MA						100e+04	420.0000	115.0000	-1.0000	98.0000	-1.	
MA						100e+04	350.0000	102.0000	-1.0000	92.0000	-1.	
MA						100e+04	320.0000	98.0000	-1.0000	91.0000	-1.	
FA						100e+04	270.0000	78.0000	1.0000	75.0000	-1.	
FA						100e+04	312.0000	99.0000	1.0000	81.0000	-1.	
FA	6	172.0000	64.0000	1.0000	39.0000	24.0000	2.2000e+04	308.0000	91.0000	1.0000	82.0000	-1.
MA	7	182.0000	80.0000	-1.0000	42.0000	35.0000	3.0000e+04	398.0000	65.0000	-1.0000	85.0000	-1.
MA	8	180.0000	80.0000								80	-1.
FA	9	169.0000	51.0000								80	-1.
FA	10	168.0000	52.0000								80	-1.
MA	11	183.0000	81.0000								80	-1.
FA	12	157.0000	47.0000								80	-1.
FA	13	164.0000	50.0000								80	-1.
FA	14	162.0000	49.0000								80	-1.
MA	15	180.0000	82.0000								80	-1.
MA	16	180.0000	81.0000								80	-1.
MB	17	185.0000	82.0000								80	1.
MB	18	187.0000	84.0000								80	1.
FB	19	168.0000	50.0000								80	1.
FB	20	166.0000	49.0000								80	1.
FB	21	158.0000	46.0000								80	1.
MB	22	177.0000	65.0000								80	1.
MB	23	180.0000	72.0000								80	1.
MB	24	181.0000	75.0000								80	1.

Количество объектов (n) = 32

Количество переменных (m) = 12

Рост (*Height*)

в сантиметрах

Вес (*Weight*)

в килограммах

Длина волос (*Hairleng*)

короткие: -1; длинные: +1

Размер обуви (*Shoesize*)

Европейский стандарт

Возраст (*Age*)

в годах

Доход (*Income*)

в евро

Потребление пива (*Beer*)

литров в год

Потребление вина (*Wine*)

литров в год

Пол (*Sex*)

мужской: -1; женский: +1

Способность плавать (*Swim*)

индекс, основанный на 500 м

дистанции

Место жительства (*A/B*)

A: -1 (Скандинавия); B: +1

(Средиземноморье)

Коэффициент интеллекта (*IQ*)

Стандартный евр. тест

Предварительная обработка данных

Цель – преобразование исходных данных в форму, наиболее удобную для анализа.

Автошкалирование

=

Центрирование
относительно
среднего

+

Взвешивание

$$x_{ik}^{scaled} = x_{ik} \frac{1}{SDev}$$

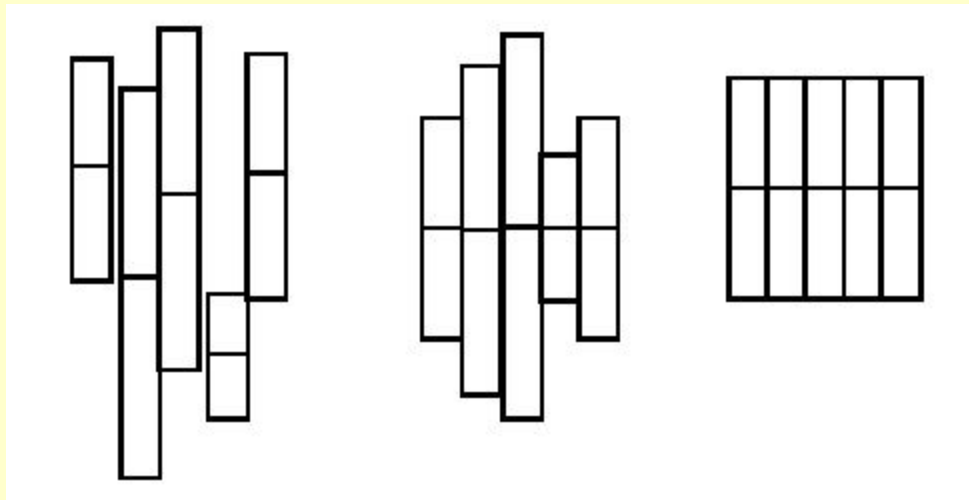
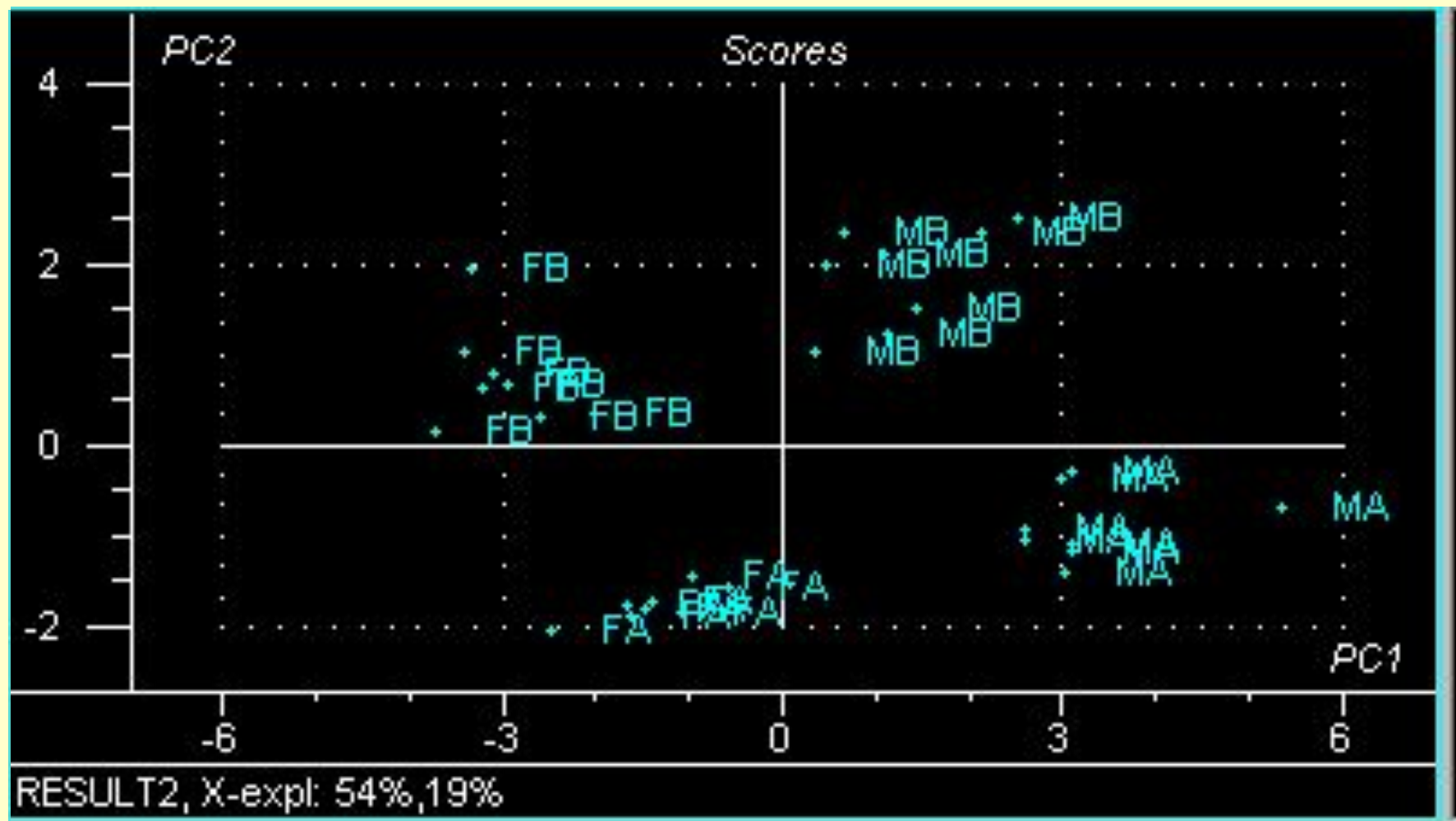
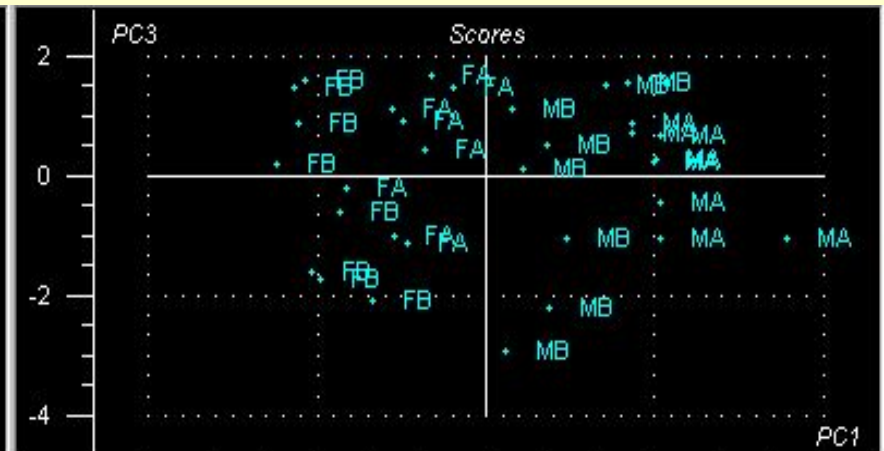
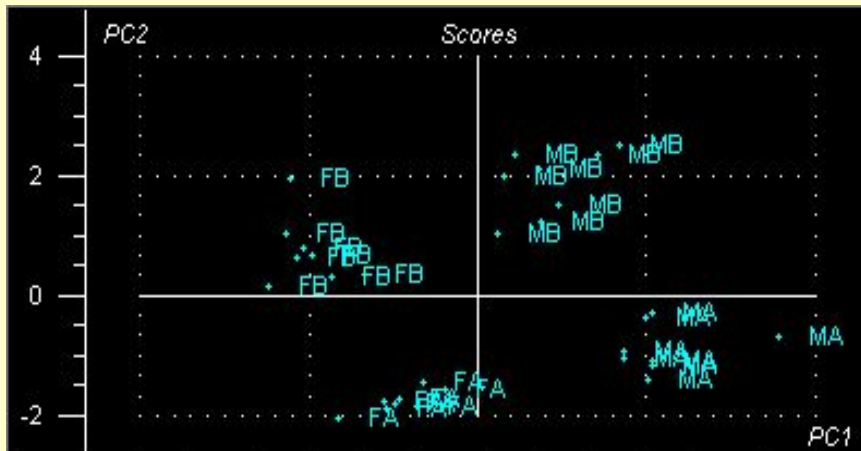


График счетов (ГК1-ГК2)

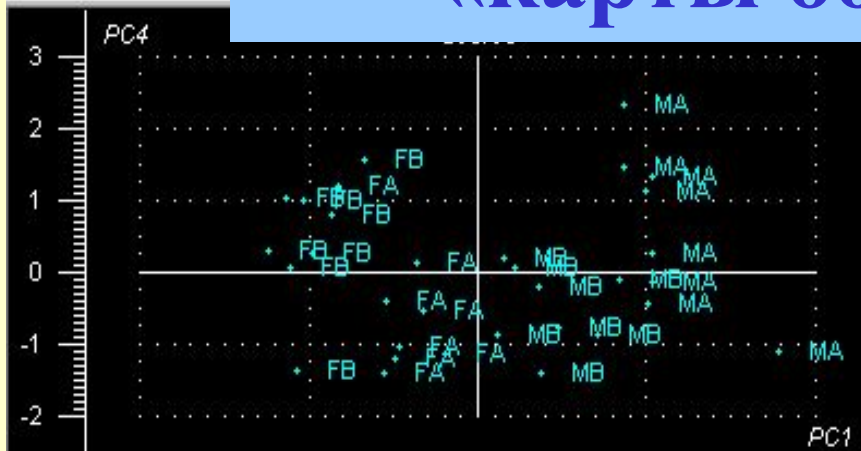


Графики счетов



RESULT2, X-expl: 54%

«карты образцов»

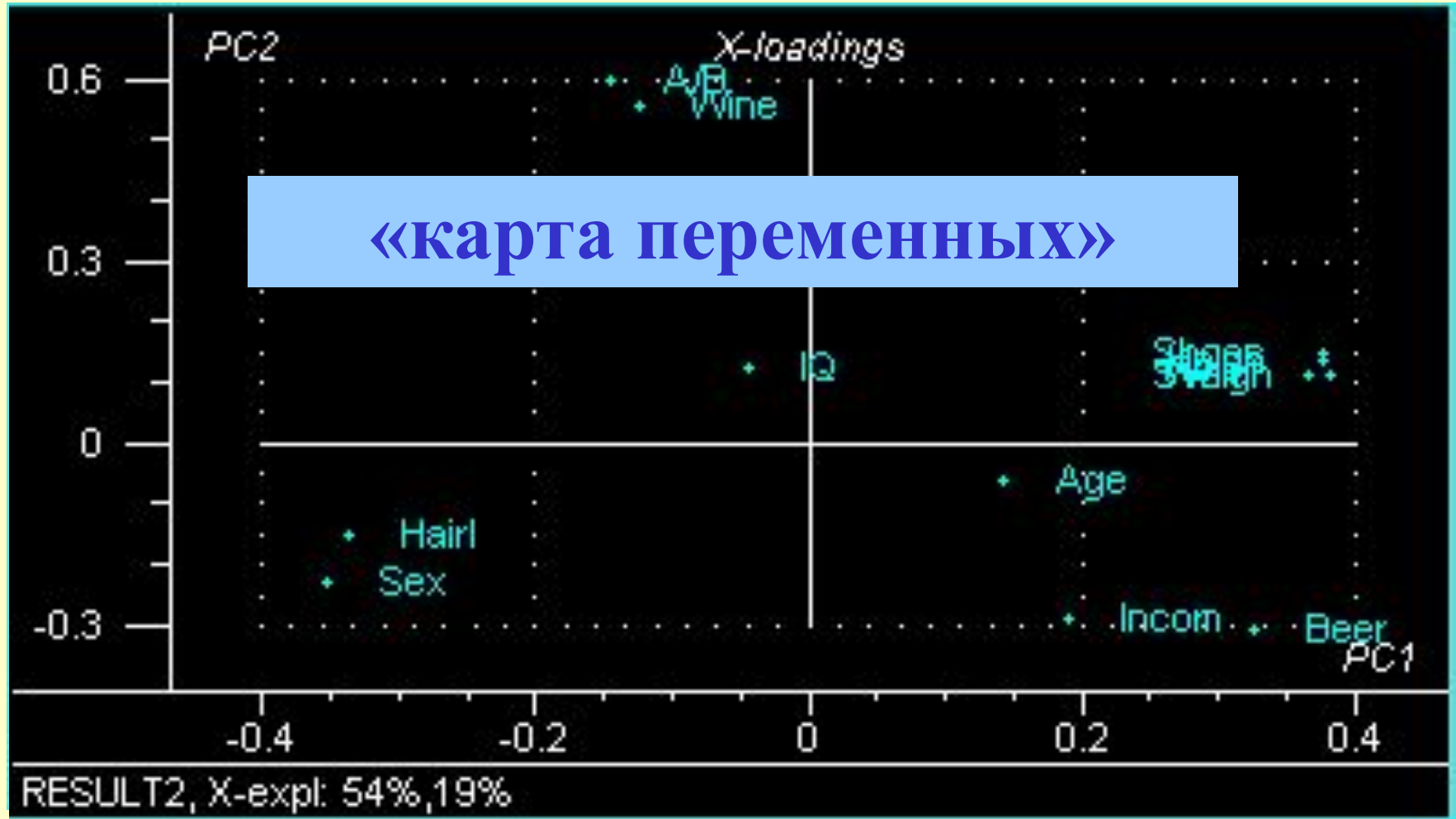


RESULT2, X-expl: 54%,8%

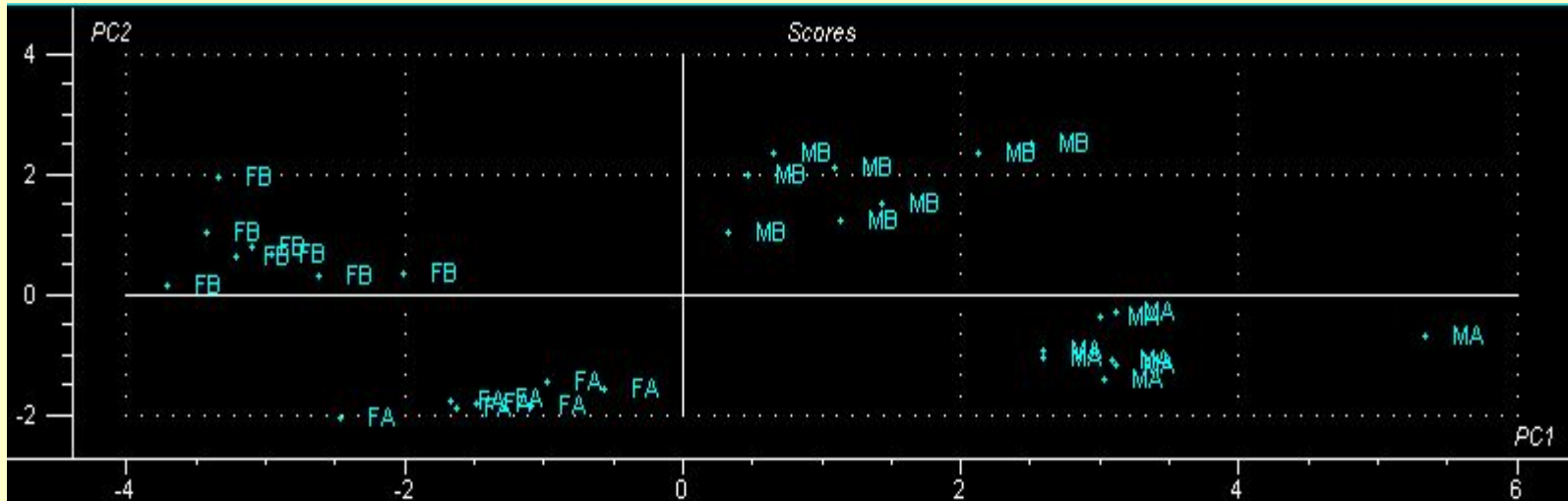


RESULT2, X-expl: 54%,3%

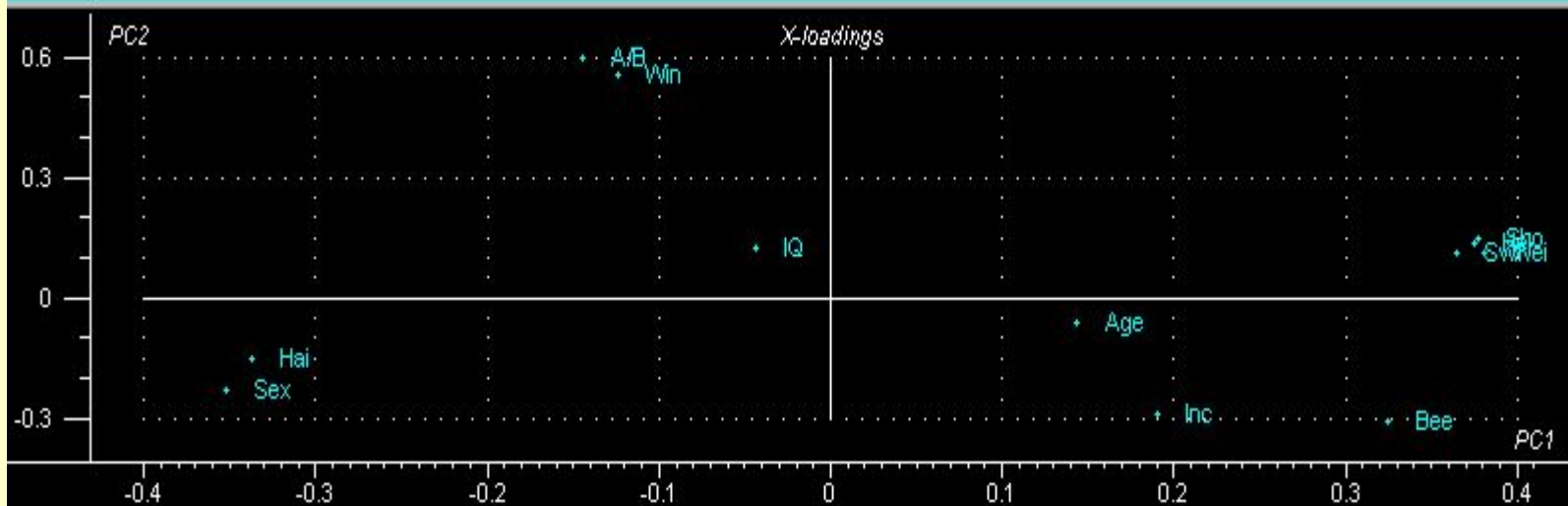
График нагрузок (ГК1-ГК2)



ГК1-ГК2 счета и нагрузки

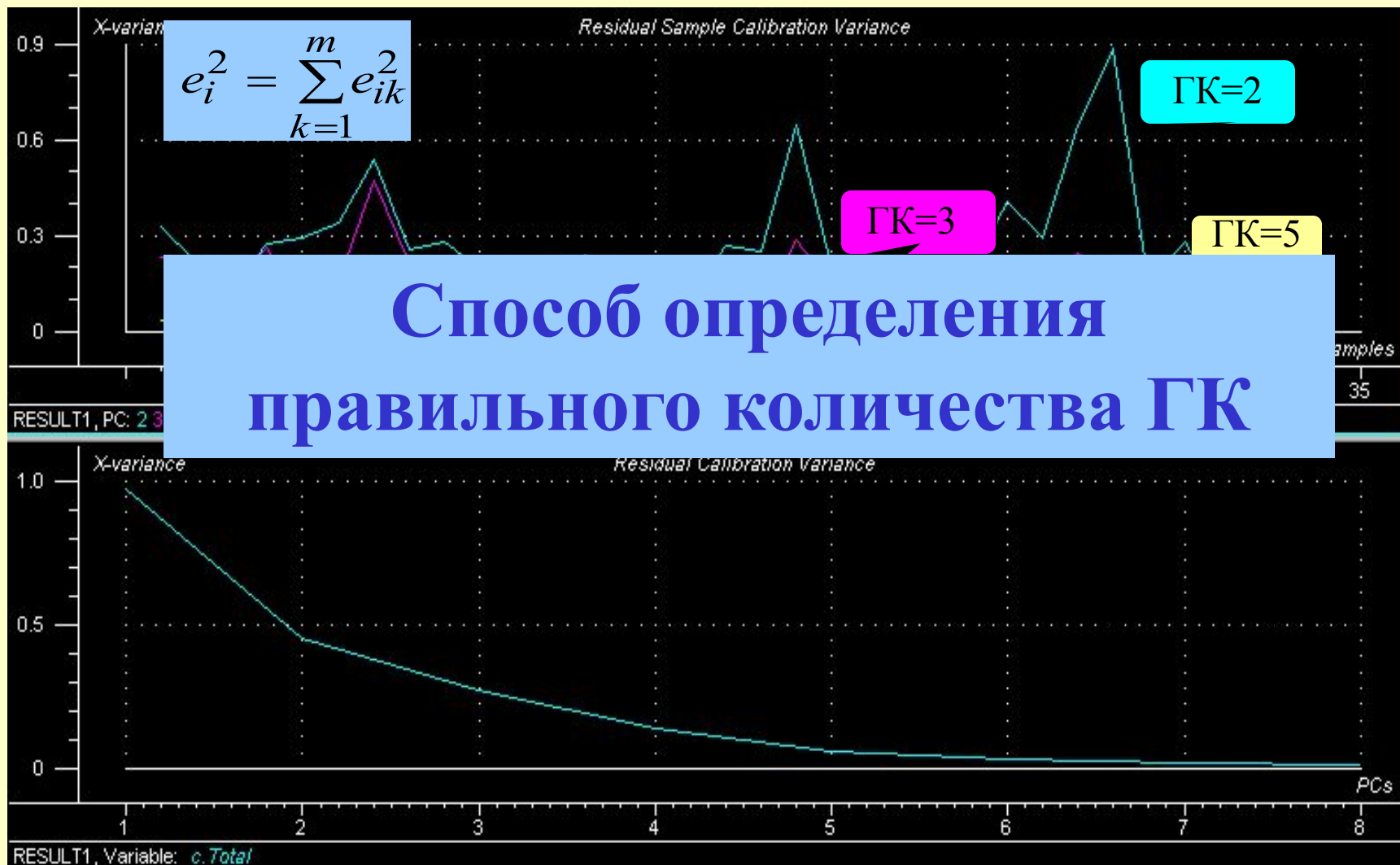


RESULT2, X-expl: 54%,19%



RESULT2, X-expl: 54%,19%

График ошибок



Цели и «инструменты»

Основные цели МГК

1. Представление объектов в пространстве, отражающем внутреннюю структуру изучаемых данных
2. Понижение размерности системы, отделение содержательной части от шума

Основные «инструменты»

1. Графики счетов – «карты образцов»
2. Графики нагрузок – «карты переменных»
3. Графики остатков – способ выбора количества ГК

Что может быть не так?

1. Данные не содержат необходимой информации
2. Использовано недостаточное количество ГК
3. Использовано излишнее количество ГК
4. Не удалены выбросы
5. Удалены точки (псевдовыбросы) содержащие важную информацию
6. Недостаточный анализ графиков счетов/нагрузок
7. Использована только стандартная (машинная) диагностика, без содержательного анализа.
8. Используются неверные методы предварительной обработки данных

Chemometrics

Data Analysis
for the Laboratory
and Chemical Plant

Richard

WILEY

Анализ смеси

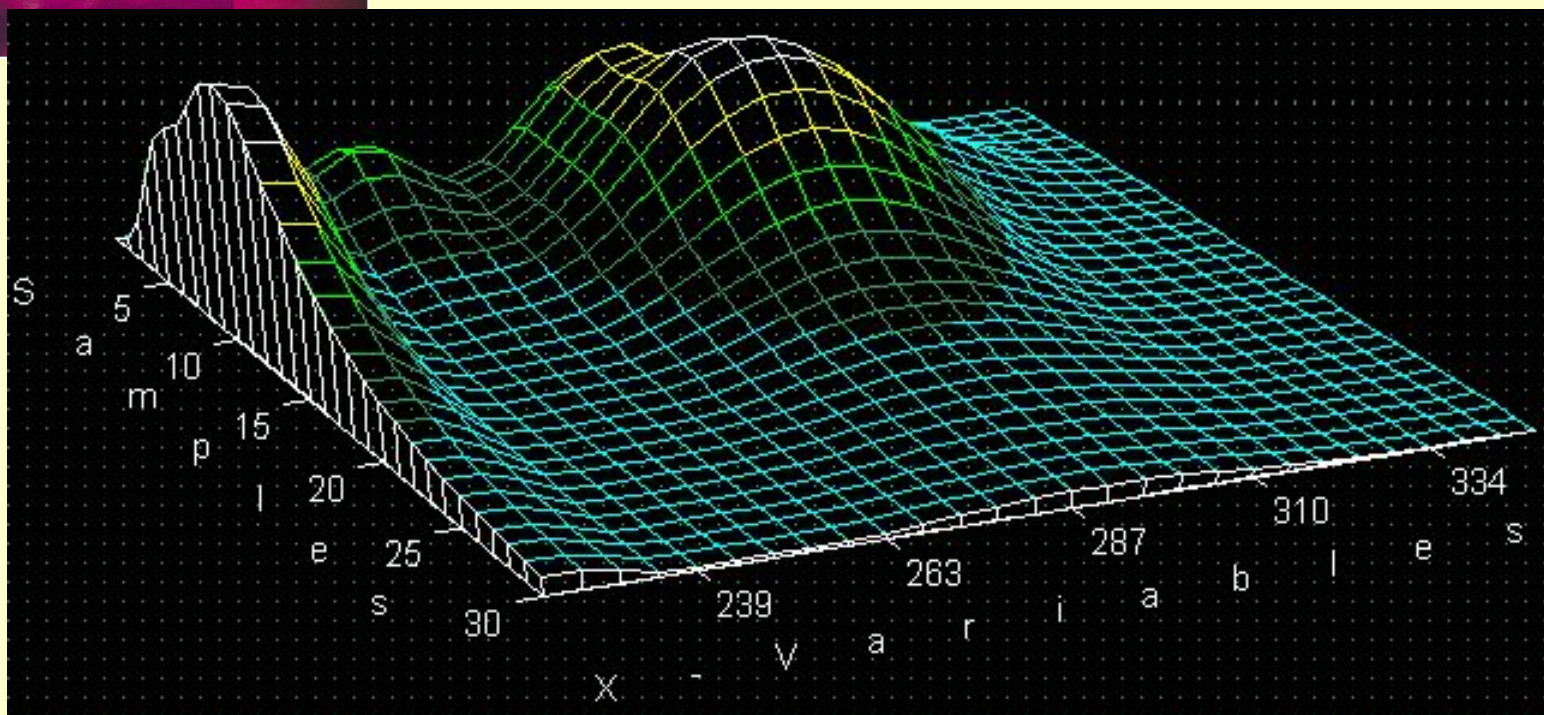
Разделение
перекрывающихся
пикув

n=30
(сек)

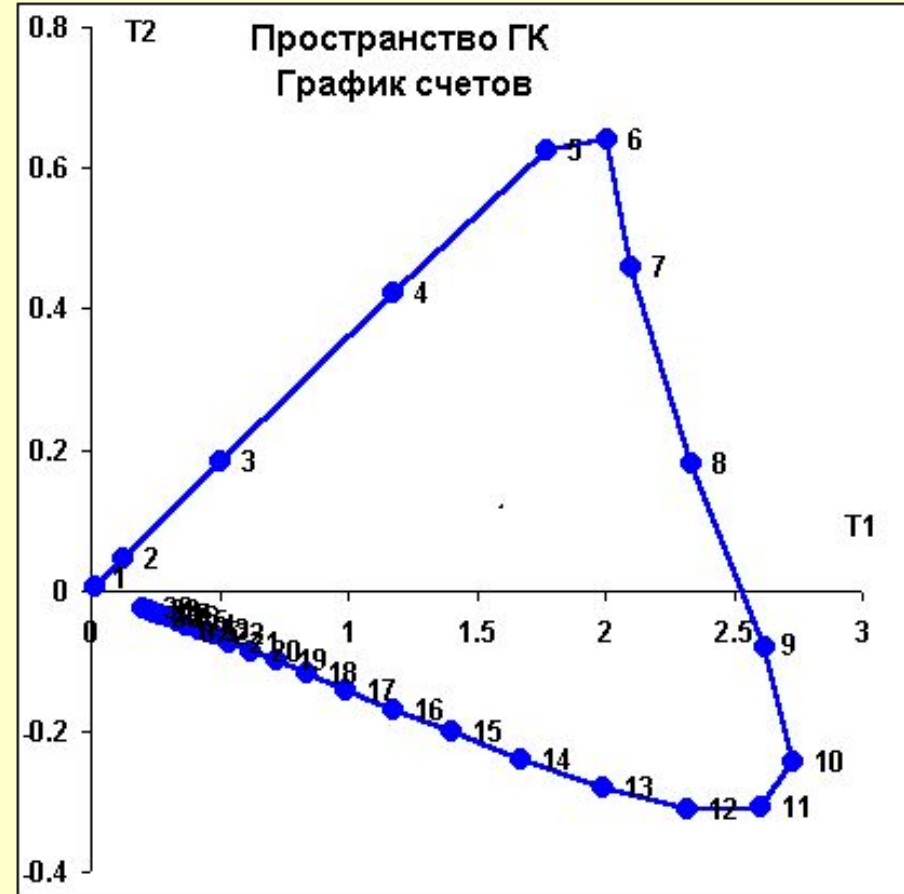
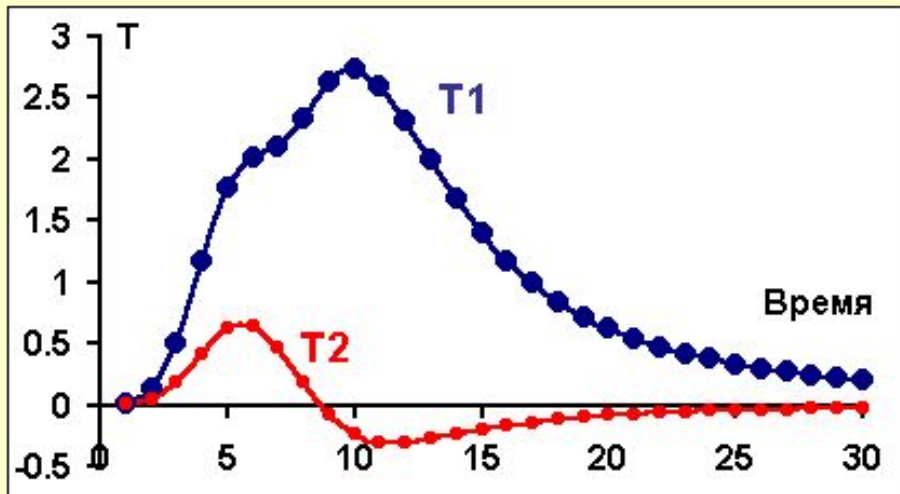
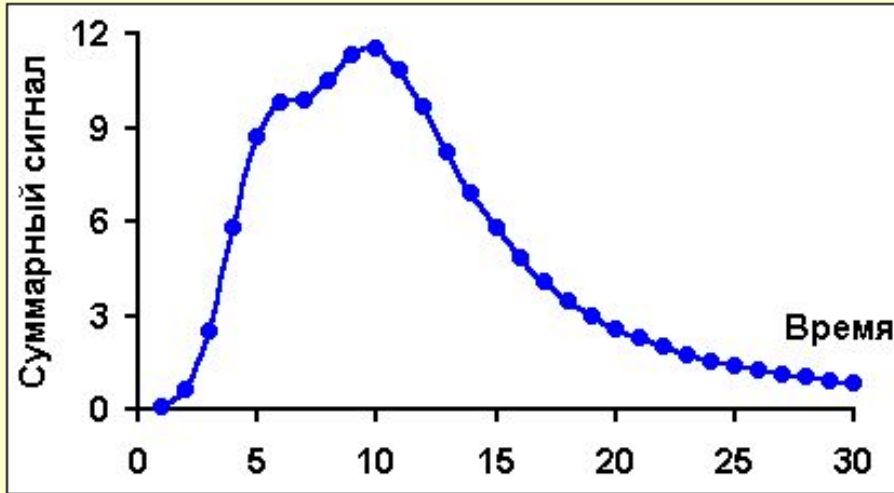
В
р
е
м
я

X
ВЖХ-ДМД

Длина волны
m=28
(нм)



Разделение пиков



Продолжение - за компьютером

