

Многомерный статистический анализ

Многмерный статистический анализ

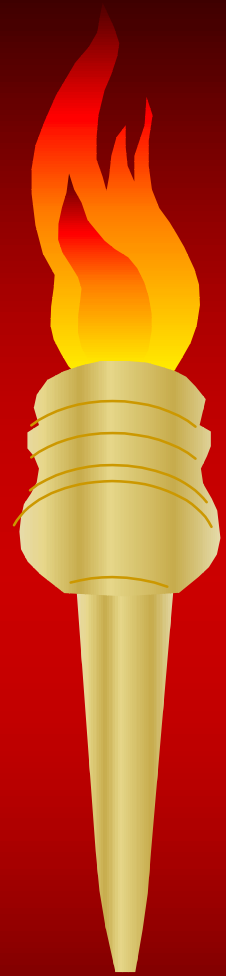
■ Зачем он нужен?

– При тестировании моделей

- Необходимо исключить влияние некоторых показателей
- Установить наличие совместного действия факторов
- Установить наличие скрытых связей

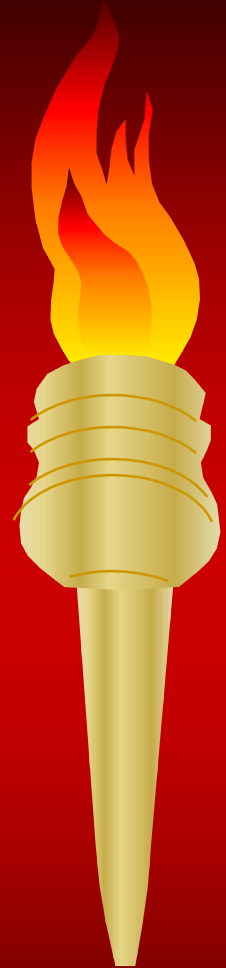
– При исследовательском анализ данных

- Снизить размерность используемых данных
- Сегментировать данные



Разработка данных (Data Mining)

- Относительно новый подход, появление которого связано с накоплением больших объемов данных в компьютеризированной форме (data warehouses)
- Нахождение внутренних закономерностей в данных, а не тестирование статистических гипотез



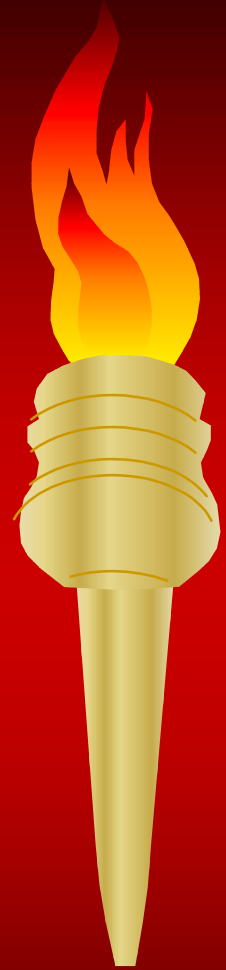
Разработка данных (Data Mining)

■ Наиболее популярные методы

- OLAP (On-Line Analytical Data Processing) - построение многомерных таблиц
- Нейронные сети (neural networks)
- Сегментирование данных (древовидное моделирование - classification trees)

■ Более распространенные методы

- Факторный анализ
- Корреспондентский анализ

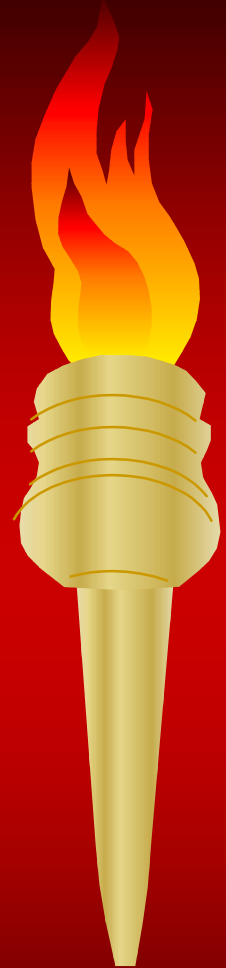


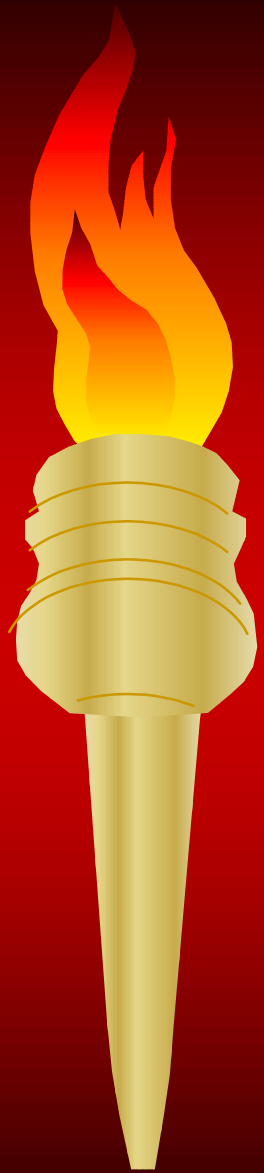
Разработка данных (Data Mining)

- Требуют очень больших массивов данных

- не менее десяти наблюдений на переменную в факторном анализе (но более ста наблюдений)
- не менее десяти наблюдений на связь при нейросетевом моделировании

- В общем сотни, лучше тысячи наблюдений

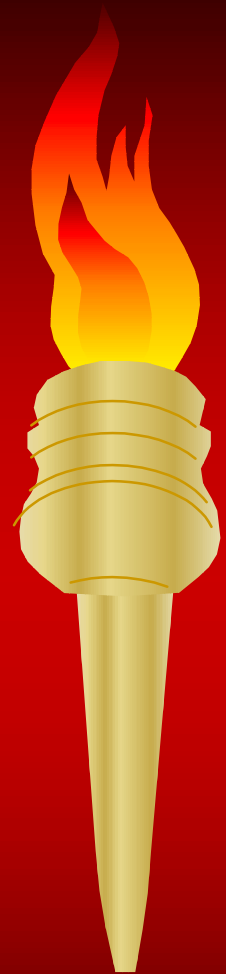




Тестирование моделей

Необходимо определить тип зависимой переменной

- Качественная
- Количественная



Количественная зависимая переменная

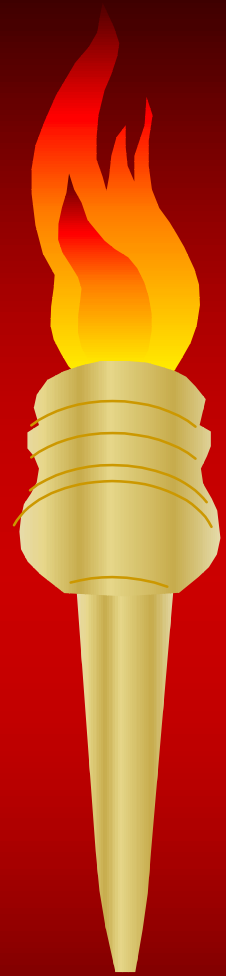
■ Общая линейная модель (GLM)

– Какие независимые переменные участвуют в анализе?

■ Только качественные - MANOVA

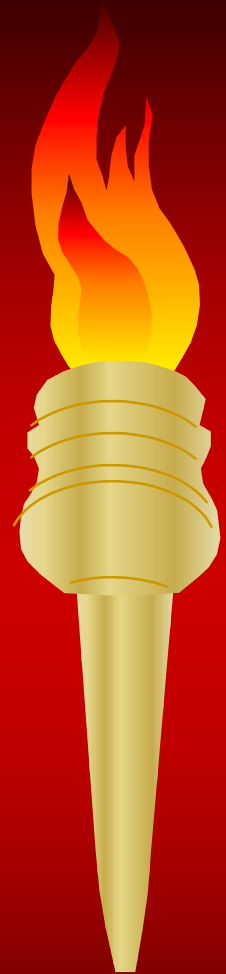
■ Количественные и качественные - MANCOVA

■ Только количественные - MLR




MANOVA

- Многофакторный дисперсионный анализ (не путать с факторным анализом!)
- Изучение влияния нескольких качественных факторов и их сочетаний на значения количественных переменных
- На самом деле, сравнение средних значений в таблицах



MANOVA

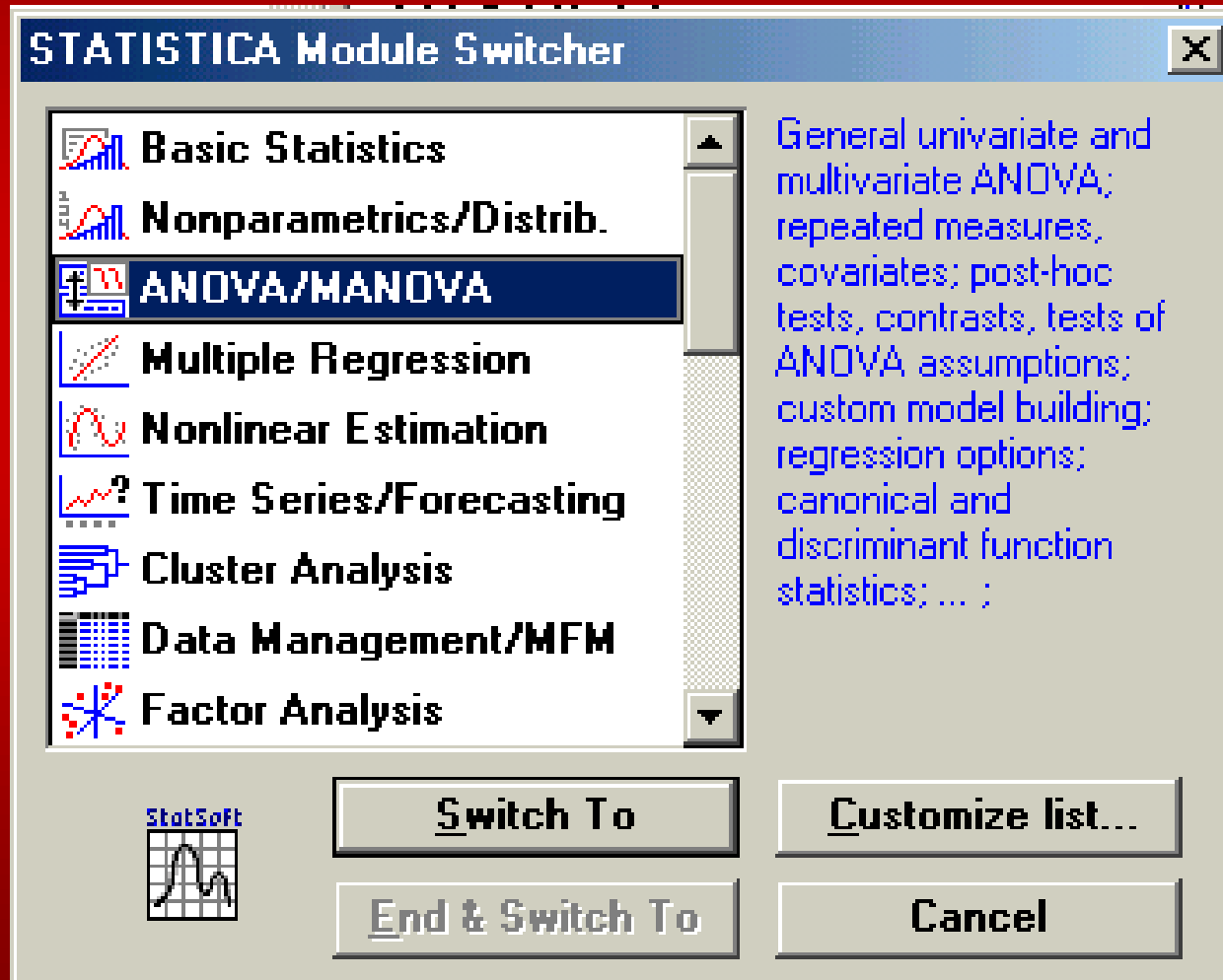


A1						A2					
B1		B2		B3		B1		B2		B3	
C1	C2	C1	C2	C1	C2	C1	C2	C1	C2	C1	C2
x1	x5	x9	x13	x17	x21	x25	x29	x33	x37	x41	x45
x2	x6	x10	x14	x18	x22	x26	x30	x34	x38	x42	x46
x3	x7	x11	x15	x19	x23	x27	x31	x35	x39	x43	x47
x4	x8	x12	x16	x20	x24	x28	x32	x36	x40	x44	x48


■ Ответит на вопросы:

- Одинаковы ли средние всех 12 групп
- Ести ли различия между средними групп, образованными фактором А, В или С по отдельности
- Имеется ли суммарное влияние факторов А, В и С

MANOVA



MANOVA



General ANOVA/MANOVA

Variables **Covariates** **OK**

Independent (factors): **SEX FAM-ED**
Dependent: **PAR**
Covariates: **none**

Cancel

Codes for between-groups factors: **Selected**

Repeated measures (within SS) design: **none**

Nested design: **none**

Random factors: **none**


Isolated control group: **none**

Open Data

Regression approach (Type I, II, III SS) **SELECT CASES** **10 W**

For large main effect and non full-factorial designs, hierarchically nested models or designs with unbalanced nesting, and mixed-model (random effect) designs, see also the Variance Components or Experimental Design modules.

MANOVA



ANOVA Results

DESIGN: 3 - way ANOVA, fixed effects

DEPENDENT: 1 variable: PAR

BETWEEN: 1-SEX (2): 1 2
2-FAM (2): 1 2
3-ED (4): 1 2 3 4

WITHIN: none

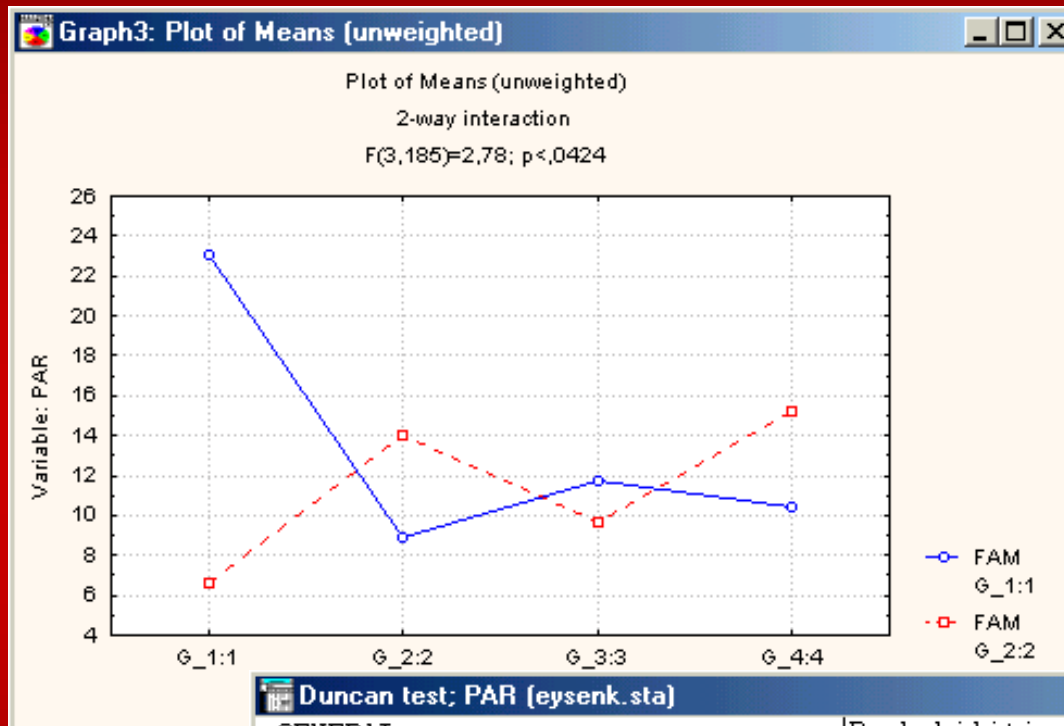
Buttons: All effects, Means/Graphs, Post hoc comparisons, Cancel

Summary of all Effects; design: [eyenk.sta]

Continue... 1-SEX, 2-FAM, 3-ED

Effect	df Effect	MS Effect	df Error	MS Error	F	p-level
1	1	294,1616	185	156,7267	1,876908	,172346
2	1	87,7365	185	156,7267	,559805	,455289
3	3	80,3911	185	156,7267	,512938	,673854
12	1	39,1748	185	156,7267	,249956	,617700
13	3	195,1141	185	156,7267	1,244932	,294824
23	3	435,8036	185	156,7267	2,780660	,042410
123	3	177,6508	185	156,7267	1,133507	,336836

MANOVA

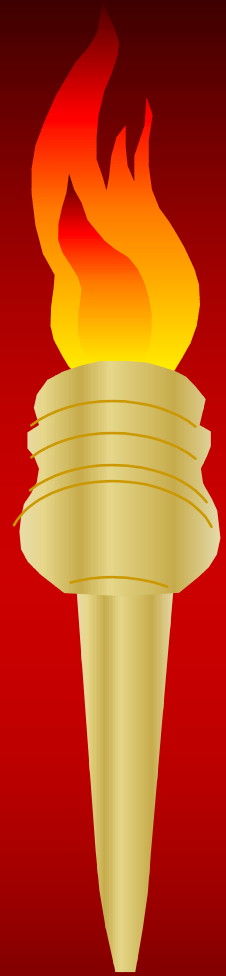


Duncan test; PAR (eysenk. sta)

GENERAL MANOVA				Probabilities for Post Hoc Tests						
				INTERACTION: 2 x 3						
SEX	FAM	ED		{1}	{2}	{3}	{4}	{5}	{6}	{7}
				23,100	8,8750	11,700	10,421	6,5714	13,971	9,6490
....	1	1	{1}		.01	.03	.02	.00	.08	.01
....	1	2	{2}	.01		.61	.77	.64	.36	.87
....	1	3	{3}	.03	.61		.79	.36	.64	.70
....	1	4	{4}	.02	.77	.79		.48	.50	.87
....	2	1	{5}	.00	.64	.36	.48		.19	.56
....	2	2	{6}	.08	.36	.64	.50	.19		.43
....	2	3	{7}	.01	.87	.70	.87	.56	.43	
....	2	4	{8}	.11	.27	.50	.38	.13	.80	.32

MLR

- Множественный линейный регрессионный анализ
- Оценка зависимости количественной переменной от одной или более количественных переменных



MLR (ОСНОВЫ)

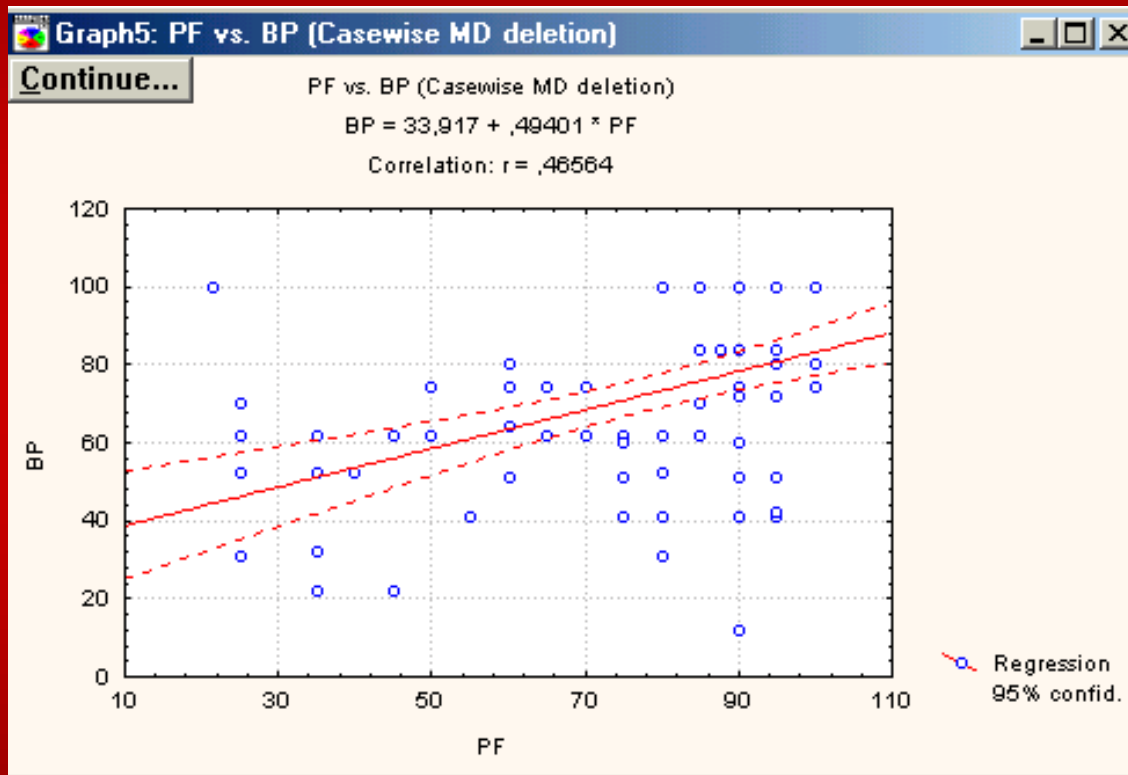


Диаграмма рассеяния + регрессионная линия
пытаемся провести ее так, чтобы расстояния от нее
до наблюдаемых точек на графике были минимальными

MLR



Multiple Regression

Variables:
 Independent: BP
 Dependent: PF

Input file: Raw Data

MD deletion: Casewise

Mode: Standard

Perform default (non-stepwise)
 Review descr. stats, corr. matr
 Extended precision computation
 Batch processing/printing
 Print residual analysis

Specify all variables for the analysis; add later. For stepwise regression etc. deselection

Buttons: OK, Cancel

Multiple Regression Results

Multiple Regression Results

Dep. Var. : PF Multiple R : ,46564420 F = 25,19363
 R²: ,21682452 df = 1, 91
 No. of cases: 93 adjusted R²: ,20821819 p = ,000003
 Standard error of estimate: 19,603387779
 Intercept: 45,863772341 Std.Error: 6,635755 t(91) = 6,9116 p < ,0000

BP beta = ,466

(significant beta's are highlighted)

Buttons: Regression summary, Predict dependent var., Compute confidence limits

Regression Summary for Dependent Variable: PF (sf36l.sta)

Continue... R = ,46564420 R² = ,21682452 Adjusted R² = ,20821819
 F(1,91) = 25,194 p < ,00000 Std. Error of estimate: 19,606

N=93	BETA	St. Err. of BETA	B	St. Err. of B	t(91)	p-level
Intercept			45,86377	6,635755	6,911613	,000000
BP	,465644	,092770	,43891	,087444	5,019325	,000003

Buttons: OK, Cancel

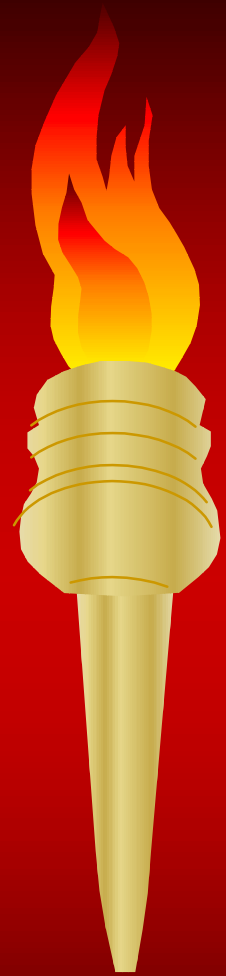
Residual analysis

Correlations & desc. stats

Alpha (display): .05 Apply

MLR (анализ остатков)

- Остатки - это значения, которые не могут быть объяснены независимой переменной
- Если у нас несколько переменных, мы берем первую, рассчитываем остатки и повторяем анализу уже на остатках (исключив влияние первой переменной) и т.д.
- Поэтому коэффициенты, получаемые в результате регрессионного анализа показывают *независимое* влияние переменных
- Если модель адекватна, остатки будут распределены по нормальному закону и разброс, соответственно, будет не очень сильным



MLR (анализ остатков)



Residual Analysis

Dep. Var. : PF Multiple R : ,46564420 F = 25,19363
R² : ,21682452 df = 1,91
No. of cases: 93 adjusted R² : ,20821819 p = ,000003
Standard error of estimate: 19,605567779
Intercept: 45,863772341 Std.Error: 6,635755 t(91) = 6,9116 p < ,0000

Statistics

- Correlations & descr. (1)
- Regression summary (2)
- Display residuals & pred. (3)
- Durbin-Watson stat (4)
- Save residuals & pred. (5)

Scatter Plots

- Pred. & residuals (D)
- Pred. & squared resid (E)
- Pred. & observed (F)

Probability Plots

- Normal plot of resid (M)

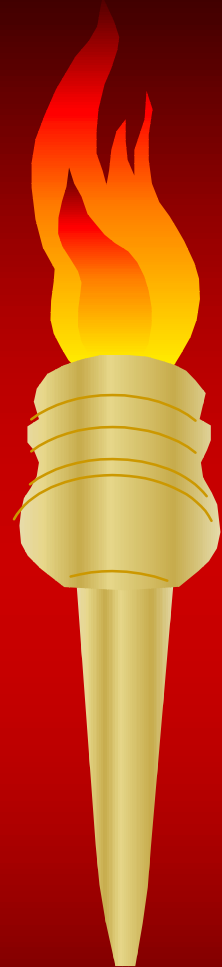
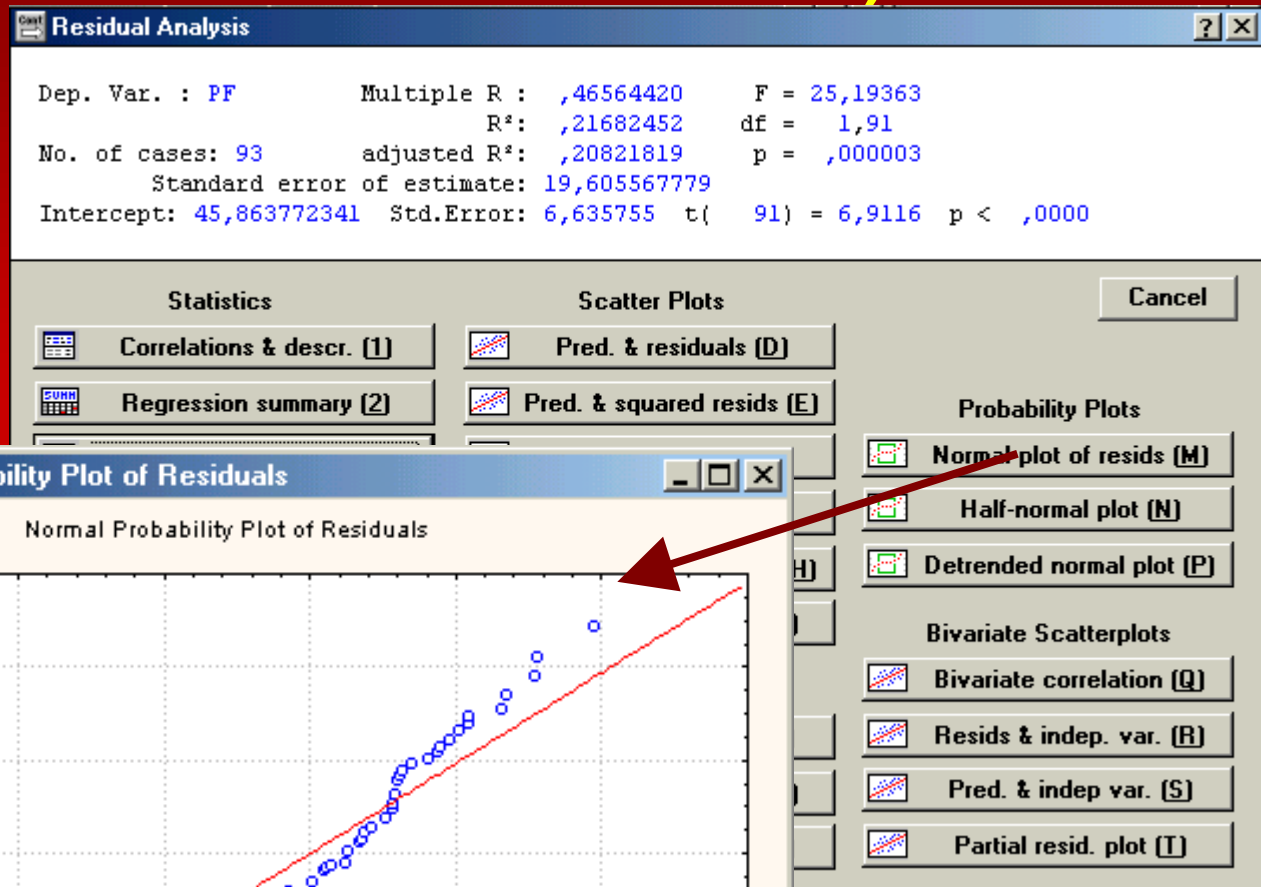
Casewise Plots

- Plots of residuals (A)
- Plots of outliers (B)
- Plots of predicted (C)

Standard Residual (sf361.sta)

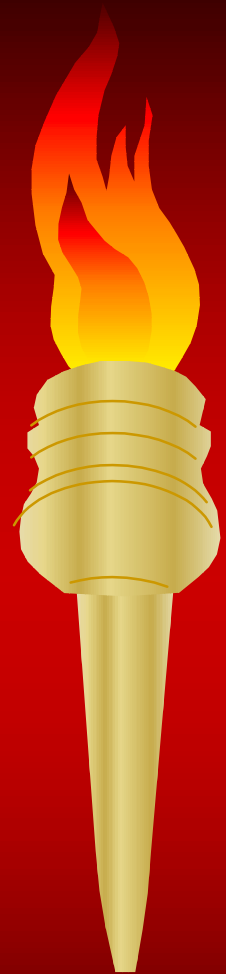
Case	Standard Residuals	Observed Value	Predictd Value	Residual
17	-3,5	90,0000	89,75478	,2452
18	*	21,4286	89,75478	-68,3262
19	*	90,0000	72,19838	17,8016
20	*	90,0000	89,75478	,2452
21	*	75,0000	63,85909	11,1409
22	*	95,0000	89,75478	5,2452
23	*	25,0000	68,68710	-43,6871
24	*	65,0000	73,07619	-8,0762
25	*	70,0000	73,07619	-3,0762
26	*	60,0000	80,97658	-20,9766
27	*	100,0000	89,75478	10,2452
28	*	50,0000	78,34312	-28,3431
29	*	95,0000	89,75478	5,2452
30	*	95,0000	68,24818	26,7518
31	*	95,0000	89,75478	5,2452

MLR (анализ остатков)



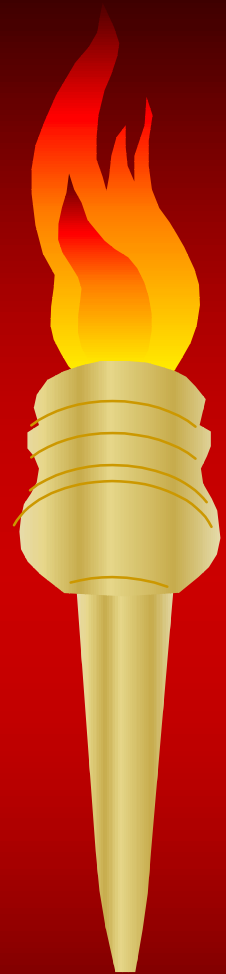
MLR (модели)

- Полная (включаются все переменные)
- Пошаговый отбор переменных (stepwise)
 - Forward
 - Backward



MANCOVA

- Многофакторный дисперсионный анализ с ковариантами
 - Комбинация из множественного линейного регрессионного анализа и дисперсионного анализа. Этапы
 - Построение модели с учетом только количественных переменных при помощи MLR
 - Дисперсионный анализ на остатках



MANCOVA

General ANOVA/MANOVA

Variables **Covariates**

Independent (factors): **SEX FAM-ED**
Dependent: **PAR**
Covariates: **AGE VNP**

Codes for between-groups factors: Selected

Repeated measures (within SS) design: none

OK **Cancel**

Nested design:

Random factors:

Isolated control group:

Regression approach (C)
For large main effect and non f...
with unbalanced nesting, and r...
Components or Experimental D...

ANOVA Results

DESIGN: 3 - way ANCOVA, fixed effects

DEPENDENT: 1 variable: **PAR**
COVARIATE: 2 variables: **AGE VNP**

BETWEEN: 1-SEX (2): 1 2
 2-FAM (2): 1 2
 3-ED (4): 1 2 3 4

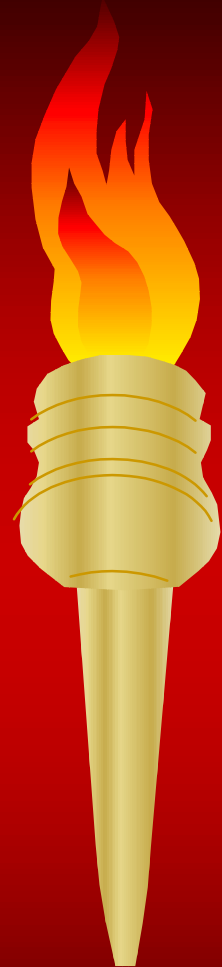
WITHIN: none

All effects **Means/Graphs** **Post hoc comparisons** **Cancel**


Specific effect/Means/Graphs **Descriptive stats & graphs** **Print design**

Planned comparisons **Options** **Within-cell regression**

Pooled effect/error term **Output options**



MANCOVA



Summary of all Effects; design: (eysenk.sta)

Continue... 1-SEX, 2-FAM, 3-ED

Effect	df Effect	MS Effect	df Error	MS Error	F	p-level
1	1	85,1137	183	122,7064	,693637	,406016
2	1	55,9069	183	122,7064	,455615	,500532
3	3	15,3134	183	122,7064	,124797	,945351
12	1	137,5578	183	122,7064	1,121033	,291092
13	3	239,7315	183	122,7064	1,953701	,122558
23	3	317,7355	183	122,7064	2,589397	,054347
123	3	214,2954	183	122,7064	1,746408	,159103

После коррекции по значениям переменных VNP и AGE эффект FAM*ED исчез (на самом деле, это эффект VNP)

Summary of all Effects; design: (eysenk.sta)

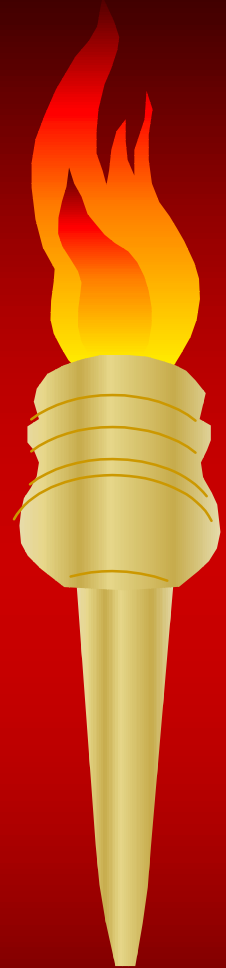
Continue... 1-SEX, 2-FAM, 3-ED

Effect	df Effect	MS Effect	df Error	MS Error	F	p-level
1	1	572,1018	184	143,0188	4,000185	,046966
2	1	91,7552	184	143,0188	,641560	,424179
3	3	201,7854	184	143,0188	1,410901	,241055
12	1	43,1038	184	143,0188	,301386	,583681
13	3	260,6784	184	143,0188	1,822686	,144560
23	3	481,9720	184	143,0188	3,369990	,019717
123	3	232,0732	184	143,0188	1,622676	,185640

Без коррекции по VNP

MANCOVA

- Наиболее общая форма анализа данных, однако в системе Statistica в соответствующем модуле невозможно выполнять MLR
- SAS позволяет выполнять все виды анализа в одном модуле



MANCOVA (SAS)

PROGRAM EDITOR

Command ==>

```
00001 PROC GLM DATA=mydat.mu99;
00002 CLASS SMPR SMPA EDUC MORT99;
00003 MODEL CH= SMPR SMPA|EDUC|MORT99 HDL;
00004 LSMEANS MORT99 EDUC /STDERR;
00005 MEANS MORT99|EDUC / DUNCAN;
00006 RUN; QUIT;
```

SAS


18:41 Monday, June 6,

General Linear Models Procedure Class Level Information

Class	Levels	Values
SMPR	2	0 1
SMPA	2	0 1
EDUC	3	1 2 3
MORT99	2	0 1

Number of observations in data set = 3907

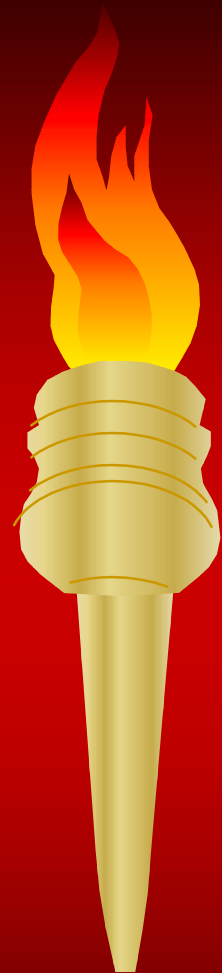
MANCOVA (SAS)



General Linear Models Procedure					
Dependent Variable: CH V1 TOTAL CHOLESTEROL					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	13	97326.42633	7486.64818	4.80	0.0001
Error	3110	4847719.72060	1558.75232		
Corrected Total	3123	4945046.14693			
	R-Square	C.V.	Root MSE	CH Mean	
	0.019682	18.05328	39.48104	218.691741	

Dependent Variable: CH V1 TOTAL CHOLESTEROL					
Source	DF	Type I SS	Mean Square	F Value	Pr > F
SMPR	1	516.03680	516.03680	0.33	0.5651
SMPA	1	1265.23144	1265.23144	0.81	0.3677
EDUC	2	29229.53245	14614.76623	9.38	0.0001
SMPA*EDUC	2	8520.60084	4260.30042	2.73	0.0652
MORT99	1	32009.07352	32009.07352	20.54	0.0001
SMPA*MORT99	1	385.45248	385.45248	0.25	0.6190
EDUC*MORT99	2	6965.28601	3482.64301	2.23	0.1072
SMPA*EDUC*MORT99	2	5992.86857	2996.43428	1.92	0.1464
HDL	1	12442.34421	12442.34421	7.98	0.0048

MANCOVA (SAS)



General Linear Models Procedure


Duncan Grouping	Mean	N	MORT99
A	221.279	1520	1
B	216.240	1604	0

General Linear Models Procedure

Duncan Grouping	Mean	N	EDUC
A	222.773	984	3
B	218.652	830	2
B	215.651	1310	1

Level of EDUC	Level of MORT99	N	-----CH----- Mean	SD
1	0	522	214.013410	38.1368610
1	1	788	216.736041	39.7042289
2	0	440	213.586364	36.8341329
2	1	390	224.366667	43.2912230
3	0	642	219.869159	38.5875225
3	1	342	228.225146	41.8455358

MANCOVA (SAS)



MORT99	CH LSMEAN	Std Err LSMEAN	Pr > T H0:LSMEAN=0
0	215.933556	1.207323	0.0
1	222.908094	1.612393	0.0

Качественная зависимая переменная

■ Независимые переменные

– Качественные

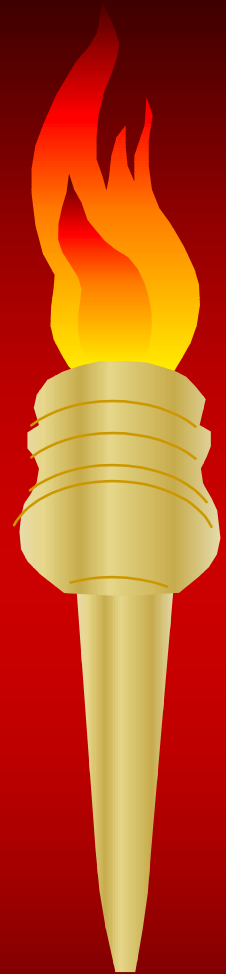
- Стратифицированный анализ Mantel-Haenszel
- Логлинейный анализ

– Качественные и количественные

- Логистическая регрессия

– Количественные

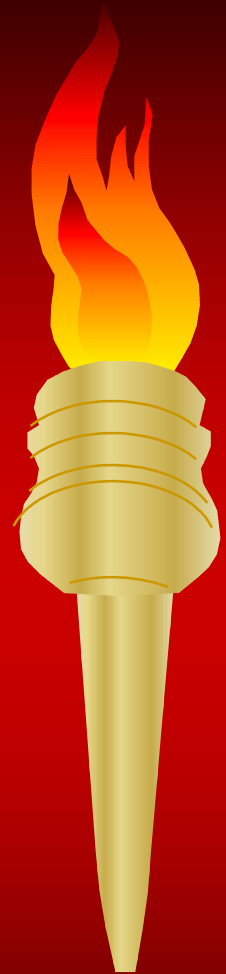
- Дискриминантный анализ



Связь нескольких качественных переменных с качественной

■ Метод Mantel-Haenszel

- Используется в рамках т.н. стратифицированного анализа (например, влияние курение на частоту развития ИМ с учетом пола, возраста - 5-летние группы, социального статуса)
- Часто используется в мультицентровых исследованиях
- Может использоваться для анализа выживаемости (год дожития - отдельная переменная)



Метод Mantel-Haenszel

- Выбираем пару базовых переменных (например, ИМ и курение)
- Строим отдельные четырехпольные таблицы для каждого уровня остальных факторов:

ПОЛ=мужской		
	Курит	Не курит
ИМ+	A1	B1
ИМ-	C1	D1

ПОЛ=женский		
	Курит	Не курит
ИМ+	A2	B2
ИМ-	C2	D2

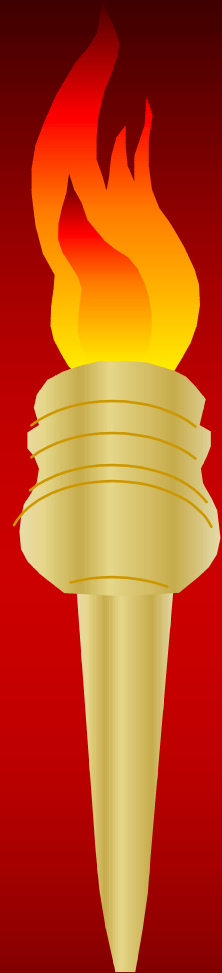
$$OR_{MH} = \frac{\sum A_i \cdot D_i / N_i}{\sum B_i \cdot C_i / N_i}$$

$$x = \sum A_i$$

$$E(x) = \sum \frac{R_{1i} \cdot C_{1i}}{N_i}$$


$$Var(x) = \sum \frac{R_{1i} \cdot C_{1i} \cdot R_{2i} \cdot C_{2i}}{N_i^2 \cdot (N_i - 1)}$$

$$\chi^2 = \frac{(|x - E(x)| - 0.5)^2}{Var(x)}$$



Метод Mantel-Haenszel

В системе Statistica можно использовать только для данных по выживаемости (требуется времена дожития)
В SAS для любых таблиц



```
Command ==>

      1  PROC FREQ DATA=MYDAT.MU99;
      2  TABLES EDUC*MORT99*SMPA/ CMH;
      3  RUN;
NOTE: The PROCEDURE FREQ used 1.00 seconds.
PROGRAM EDITOR
Command ==> ■

00001 PROC FREQ DATA=MYDAT.MU99;
00002 TABLES EDUC*MORT99*SMPA/ CMH;
00003 RUN;
```

Метод Mantel-Haenszel (SAS)

TABLE 3 OF MORT99 BY SMPA
CONTROLLING FOR EDUC=3

MORT99	SMPA(SMOKED IN THE PAST)		
	0	1	Total
Frequency			
Percent			
Row Pct			
Col Pct			
0	247	413	660
	24.46	40.89	65.35
	37.42	62.58	
	78.16	59.51	
Total	316	694	1010
	31.29	68.71	100.00

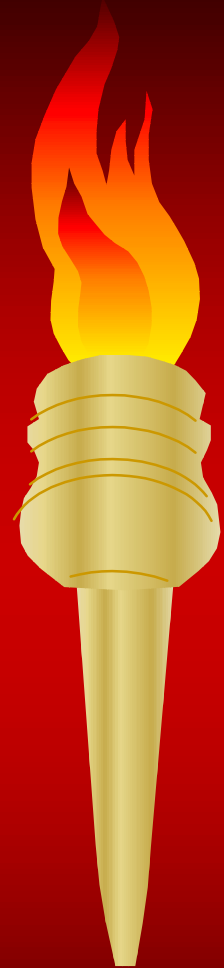
(Continued)

MORT99	SMPA(SMOKED IN THE PAST)		
	0	1	Total
Frequency			
Percent			
Row Pct			
Col Pct			
1	69	281	350
	6.83	27.82	34.65
	19.71	80.29	
	21.84	40.49	
Total	316	694	1010
	31.29	68.71	100.00

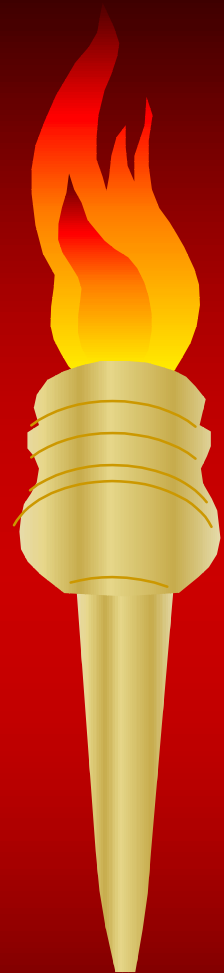
SUMMARY STATISTICS FOR MORT99 BY SMPA
CONTROLLING FOR EDUC

Cochran-Mantel-Haenszel Statistics (Based on Table Scores)

Statistic	Alternative Hypothesis	DF	Value	Prob
1	Nonzero Correlation	1	83.951	0.000
2	Row Mean Scores Differ	1	83.951	0.000
3	General Association	1	83.951	0.000



Метод Mantel-Haenszel (SAS)



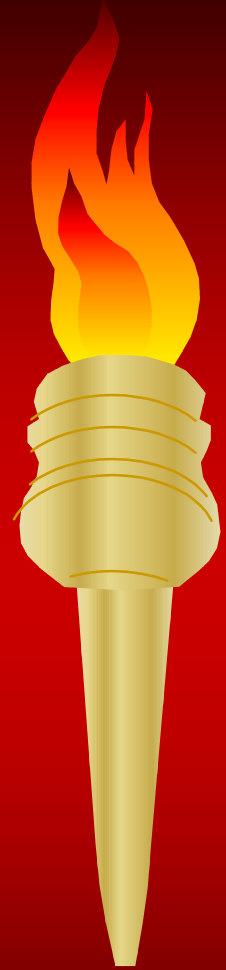
Estimates of the Common Relative Risk (Row1/Row2)
95%

Type of Study	Method	Value	Confidence Bounds	
Case-Control (Odds Ratio)	Mantel-Haenszel	2.371	1.971	2.851
	Logit	2.372	1.966	2.862
Cohort (Col1 Risk)	Mantel-Haenszel	1.978	1.709	2.289
	Logit	1.980	1.700	2.306
Cohort (Col2 Risk)	Mantel-Haenszel	0.850	0.820	0.880
	Logit	0.857	0.828	0.887

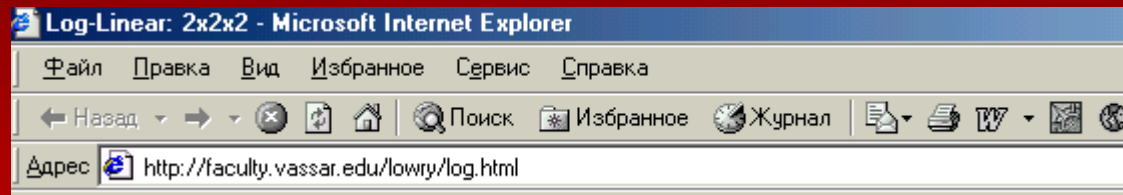
The confidence bounds for the M-H estimates are test-based.

Логлинейный анализ

- Аналог дисперсионного анализа для многомерных таблиц
- Анализируется совместное влияние нескольких переменных на частоты в многомерной таблице
- Если модель соответствует данным (fit), то переменные не влияют на частоты в таблице - модель можно подогнать зная только краевые частоты (R_1, R_2, C_1, C_2)



Логлинейный анализ



Data Entry

	A ₁		A ₂	
	C ₁	C ₂	C ₁	C ₂
B ₁	<input type="text" value="1"/>	<input type="text" value="1"/>	<input type="text" value="5"/>	<input type="text" value="5"/>
B ₂	<input type="text" value="2"/>	<input type="text" value="2"/>	<input type="text" value="10"/>	<input type="text" value="10"/>

df [ABC] = 4

G² [ABC] =

p =

Note that G² is distributed approximately as chi-square. p values are derived by referring calculated values of G² to the appropriate sampling distributions of chi-square.

Data Entry


	A ₁		A ₂	
	C ₁	C ₂	C ₁	C ₂
B ₁	<input type="text" value="10"/>	<input type="text" value="1"/>	<input type="text" value="25"/>	<input type="text" value="5"/>
B ₂	<input type="text" value="2"/>	<input type="text" value="2"/>	<input type="text" value="10"/>	<input type="text" value="10"/>

df [ABC] = 4

G² [ABC] =

Note that G² is distributed approximately as chi-square. p values are derived by referring calculated values of G² to the appropriate sampling distributions of

Логлинейный анализ




AB Table			AC Table			BC Table					
	B ₁	B ₂		C ₁	C ₂		C ₁	C ₂			
A ₁	11	4	15	A ₁	12	3	15	B ₁	35	6	41
A ₂	30	20	50	A ₂	35	15	50	B ₂	12	12	24
	41	24	65		47	18	65		47	18	65

df [AB] = 1		df [AC] = 1		df [BC] = 1	
G ² [AB] =	0.91	G ² [AC] =	0.6	G ² [BC] =	9.29
p =	ns	p =	ns	p =	p<.005

df [ABC-AB] = 3		df [ABC-AC] = 3		df [ABC-BC] = 3	
G ² [ABC-AB] =	9.7	G ² [ABC-AC] =	10.01	G ² [ABC-BC] =	1.32
p =	p<.025	p =	p<.025	p =	ns

df [ABC-AC-BC] = 2		df [ABC-AB-BC] = 2		df [ABC-AB-AC] = 2	
G ² [ABC-AC-BC] =	0.72	G ² [ABC-AB-BC] =	0.41	G ² [ABC-AB-AC] =	9.1
p =	ns	p =	ns	p =	p<.025

Логлинейный анализ



Data: NEW.STA 10v * 10c

NUM VAL	1 A	2 B	3 C	4 WEIG
1	1,000	1,000	1,000	10,000
2	1,000	1,000	2,000	1,000
3	1,000	2,000	1,000	2,000
4	1,000	2,000	2,000	2,000
5	2,000	1,000	1,000	25,000
6	2,000	1,000	2,000	5,000
7	2,000	2,000	1,000	10,000
8	2,000	2,000	2,000	10,000
9				

Analysis

Case analyzed:

A 2 B 2 C 2
x x x


Input file: Frequencies with coding variables

Variables: A-C

Variable containing frequencies: WEIG

Select codes Selected

Логлинейный анализ



Log-Linear Model Specification

Table to be analyzed:

(1)	A	(2)	B	(3)	C
2	x	2	x	2	

Minimum cell frequency: 1, Maximum: 25, Sum: 65,

Review complete observed table

Display slices of tables
 One by one As many as fit at a time

Specify model to be tested

Test all marginal & partial association models

Automatic selection of best model

Save the table


Structural zeros:

Delta:

Max. no. of iterations:

Convergence criterion:

Логлинейный анализ



Results

Table to be analyzed:

(1)	A	(2)	B	(3)	C
2	x	2	x	2	

Minimum cell frequency: 1, Maximum: 25, Sum: 65,

Model to be tested: 32

Delta: ,5000 ; Maximum iterations: 50 ; Conv. criterion: ,0100
Convergence reached after 2 iterations

	Chi-square	p
Maximum Likelihood Chi-square:	19,46161	,0006392
Pearson Chi-square:	18,38187	,0010413

Observed table Graph of observed vs. fitted

Fitted table Graph of fitted vs. residuals

Residuals (observed-fitted) Graph of fitted vs. std. residuals

Residuals standardized Graph of fitted vs. m-L components


Components of max-Likelihood Chi Graph of fitted vs. F-T deviates

Freeman-Tukey deviates

Marginal tables

OK Cancel

Логлинейный анализ (SAS)



```
PROGRAM EDITOR
Command ==>

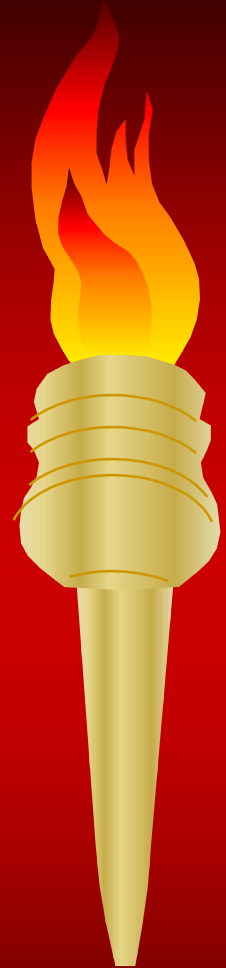
00001 DATA loglin; INPUT A B C W; CARDS;
00002 1 1 1 10
00003 1 1 2 1
00004 1 2 1 2
00005 1 2 2 2
00006 2 1 1 25
00007 2 1 2 5
00008 2 2 1 10
00009 2 2 2 10
00010 ;
00011 RUN;
00012 PROC CATMOD DATA=loglin;
00013 MODEL C*A*B =_RESPONSE_ /ML;
00014 LOGLIN A|B|C;
00015 WEIGHT W;
00016 RUN;
```

MAXIMUM LIKELIHOOD ANALYSIS OF VARIANCE TABLE

Source	DF	Chi-Square	Prob
A	1	12.99	0.0003
B	1	0.51	0.4746
A*B	1	0.19	0.6636
C	1	6.03	0.0141
C*A	1	0.19	0.6636
C*B	1	6.03	0.0141
C*A*B	1	0.19	0.6636
LIKELIHOOD RATIO	0	.	.

Логлинейный анализ (SAS)

- На самом деле необходимо проанализировать несколько моделей:
 1. «Независимую» (LOGLIN A B C)
 2. «С исключением верхнего уровня» (LOGLIN A B C B*C)
 3. «С двойными взаимодействиями» (LOGLIN A B C A*B A*C B*C)
 4. «Полную» (LOGLIN A|B|C)
- Записать остатки (LIKLEHOOD RATIO)
- Посчитать выигрыш при усложнении модели (χ^2 и df)



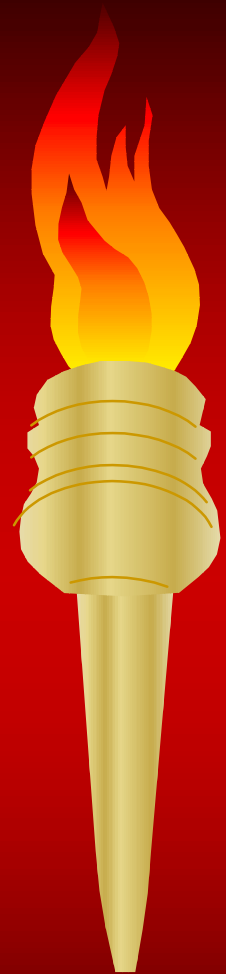
Логлинейный анализ (SAS)

■ Составить таблицу

модель	df	LR	p	dChi	DeltaDF
1	4	10.61	0.0313		
2	3	1.31	0.7258	9.30	1
3	1	0.19	0.6568	1.12	1
4	0	0	0	0.19	1

■ Видно, что наибольший выигрыш при переходе от модели 1 к модели 2, все остальные χ^2 не достоверны

■ Значит выбираем модель 2



Логлинейный анализ (SAS)

MAXIMUM LIKELIHOOD ANALYSIS OF VARIANCE TABLE

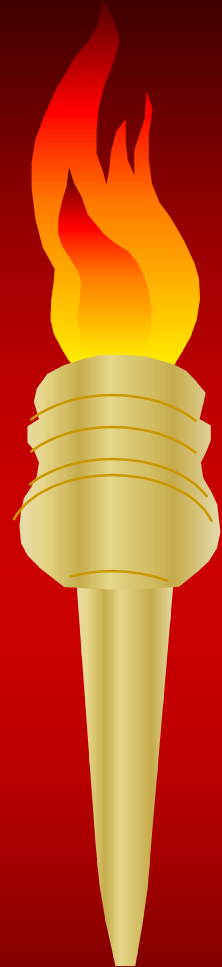
Source	DF	Chi-Square	Prob
A	1	16.73	0.0000
B	1	0.39	0.5305
C	1	8.59	0.0034
C*B	1	8.59	0.0034
LIKELIHOOD RATIO	3	1.31	0.7258

ANALYSIS OF MAXIMUM LIKELIHOOD ESTIMATES

Effect	Parameter	Estimate	Standard Error	Chi-Square	Prob
RESPONSE	1	-0.6020	0.1472	16.73	0.0000
	2	0.0943	0.1504	0.39	0.5305
	3	0.4409	0.1504	8.59	0.0034
	4	0.4409	0.1504	8.59	0.0034

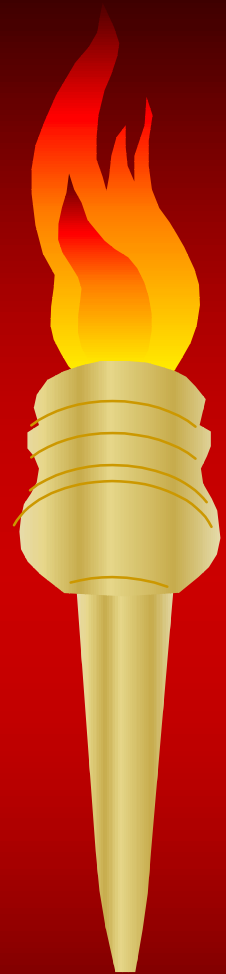
NOTE: _RESPONSE_ = A B C C*B

Мы можем оценить эффекты. Например, коэффициент для А -0,6. Беря антилогарифм ($e^{-0.6}=0,55$) имеем, что значений первого уровня А почти в 2 раза меньше, чем второго. Глядя на 4 строку, можно сказать, что С1 значительно больше при В1, и т.д.

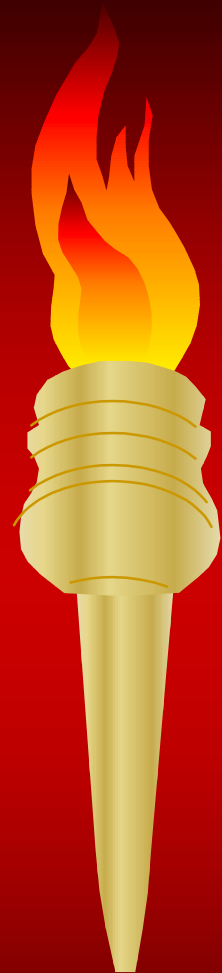
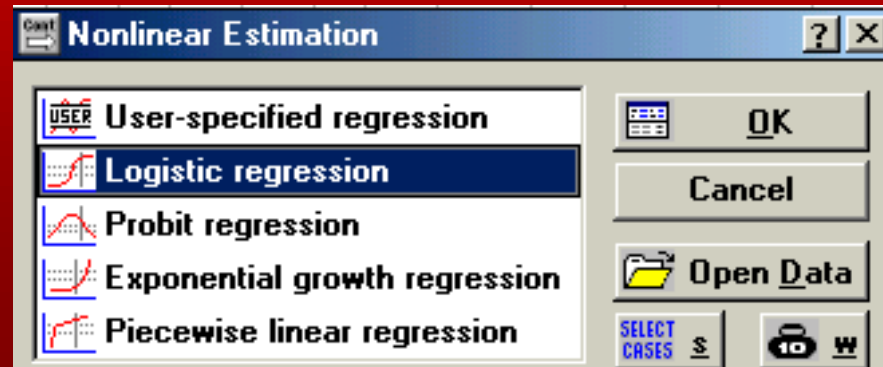
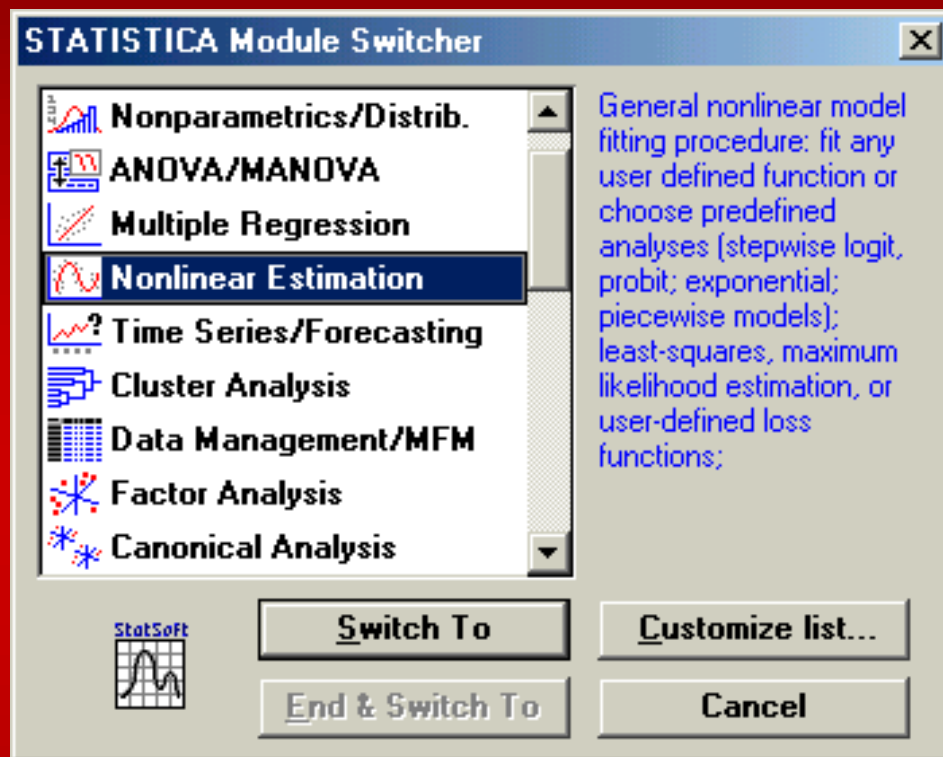


Логистическая регрессия

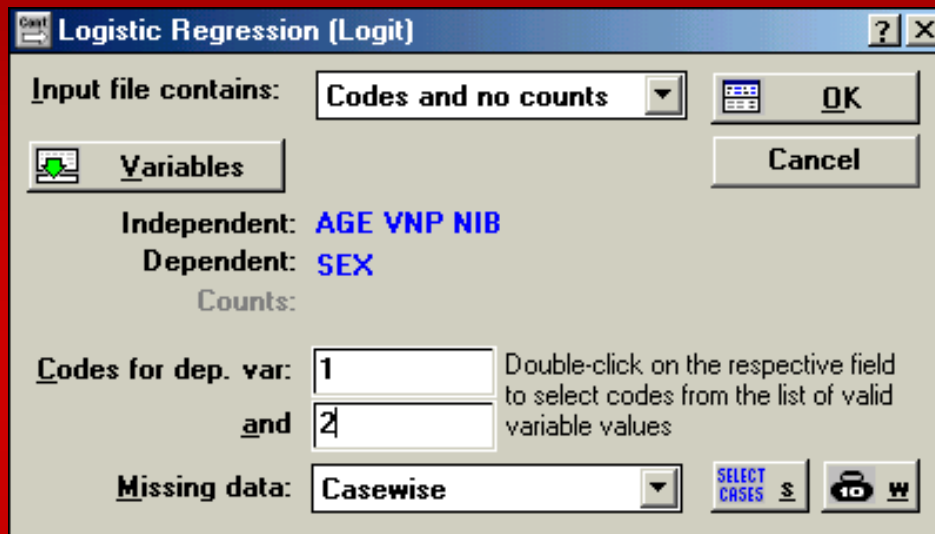
- На самом деле, наиболее общий случай (более общий, чем логлинейный анализ), правда оптимизирована для работы с бинарными зависимыми переменными
- Может принимать для анализа как качественные, так и количественные переменные
- В принципе, если логлинейный анализ аналог ANOVA, то логистическая регрессия - аналог ANCOVA (в SAS можно даже использовать одну и ту же процедуру)



Логистическая регрессия




Логистическая регрессия



- При проведении логистического регрессионного анализа качественные переменные с более чем двумя уровнями должны быть превращены в т.н. *dummy* переменные (пустышки) с двумя уровнями. Например, образование с 3 будет превращено в 2 переменных, например высшее (есть/нет) и ниже среднего (да/нет).

Логистическая регрессия



Results

Model is: **logistic regression (logit)** No. of 0's: 69,00000 (34,32836%)
No. of 1's: 132,0000 (65,67164%)
Dependent variable: **SEX** Independent variables: 3
Loss function is: **maximum likelihood** Final value: 119,61091998
-2*log(Likelihood): for this model = 39,34033, intercept only: 258,5622
Chi-square = 19,34033, df = 3, p = ,0002331

Parameters & standard errors Fitted 2D function & observed vals OK
Cov./corr. of parameters Fitted 3D function & observed vals Cancel
 Scale MSE to 1 Cnf. interval: 95,0 % Distribution of residuals
Residuals Predicted values Normal probability plot of residuals
Observed values Half-normal probability plot
Classification of cases
Means & standard deviations
Difference (previous - current)
Save predicted and residuals

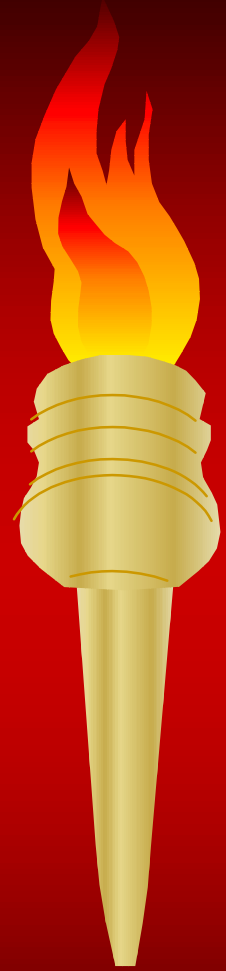
Model: Logistic regression (logit) N of 0's:69 1's:132

Dep. var: SEX Loss: Max likelihood (MS-err. scaled to 1)
Final loss: 119,61091998 ChiI(3)=19,340 p=,00023

	Const. B0	AGE	VNP	NIB
N=201				
Estimate	-5,24838	,00834	,3242	-,17996
Standard Error	1,54970	,03663	,0843	,45949
t(197)	-3,38671	,22763	3,8445	-,39165
p-level	,00085	,82017	,0002	,69574
-95%CL	-8,30450	-,06390	,1579	-1,08611
+95%CL	-2,19225	,08058	,4905	,72619
Wald's Chi-square	11,46980	,05181	14,7799	,15339
p-level	,00071	,81994	,0001	,69532
Odds ratio (unit ch)	,00526	1,00837	1,3829	,83530
-95%CL	,00025	,93810	1,1711	,33753
+95%CL	,11167	1,08391	1,6332	2,06719
Odds ratio (range)		1,27355	35,3878	,83530
-95%CL		,15675	5,6800	,33753
+95%CL		10,34733	220,4760	2,06719

Логистическая регрессия

- Для качественных бинарных показателей мы можем взять OR прямо из таблицы (это будет OR, откорректированный на влияние других факторов риска)
- Для количественных показателей необходимо взять регрессионный коэффициент (estimate), умножить его на стандартное отклонение для этого показателя и взять от произведения антилогарифм.
- Например, для возраста в описанном ранее примере $estimate=0,008$, $SD=5,4$. Соответственно $OR=e^{0,008*5,4}=1,04$



Логистическая регрессия (SAS)

- Можно использовать несколько процедур - CATMOD и LOGISTIC. Первая может также использоваться для мультиномиальных (не бинарных) типов анализа
- CATMOD очень требовательна к памяти

```
00001 PROC LOGISTIC DATA=mydat.mu99;  
00002 MODEL MORT99=SBP EDHIGH EDLOW SMPA SMPR;  
00003 RUN;  
00004 PROC CATMOD DATA=mydat.mu99;  
00005 DIRECT SBP;  
00006 MODEL MORT99=SBP EDHIGH EDLOW SMPA SMPR/ NOGLS ML NOPROFILE;  
00007 RUN;  
00008 QUIT;
```

Логистическая регрессия (SAS)

The LOGISTIC Procedure

Analysis of Maximum Likelihood Estimates

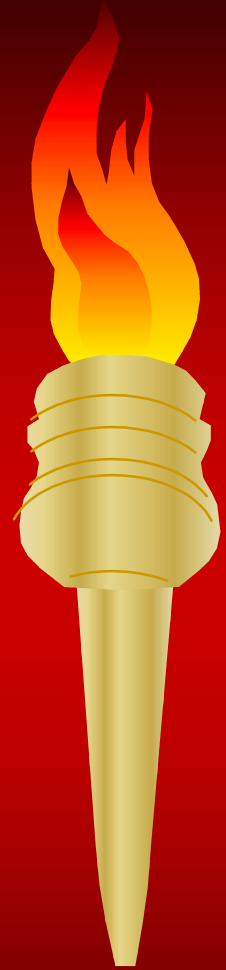
Variable	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square	Standardized Estimate
INTERCPT	5.2032	0.2993	302.2534	0.0001	.
SBP	-0.0313	0.00197	251.6423	0.0001	-0.385224
EDHIGH	0.3162	0.1025	9.5254	0.0020	0.080750
EDLOW	-0.3936	0.0938	17.6021	0.0001	-0.107249
SMPA	-0.4710	0.1208	15.1963	0.0001	-0.104752
SMPR	-0.6347	0.0951	44.5831	0.0001	-0.173424

ANALYSIS OF MAXIMUM LIKELIHOOD ESTIMATES

Effect	Parameter	Estimate	Standard Error	Chi-Square	Prob
INTERCEPT	1	4.6116	0.2787	273.88	0.0000
SBP	2	-0.0313	0.00197	251.64	0.0000
EDHIGH	3	-0.1581	0.0512	9.53	0.0020
EDLOW	4	0.1968	0.0469	17.60	0.0000
SMPA	5	0.2355	0.0604	15.20	0.0001
SMPR	6	0.3174	0.0475	44.58	0.0000

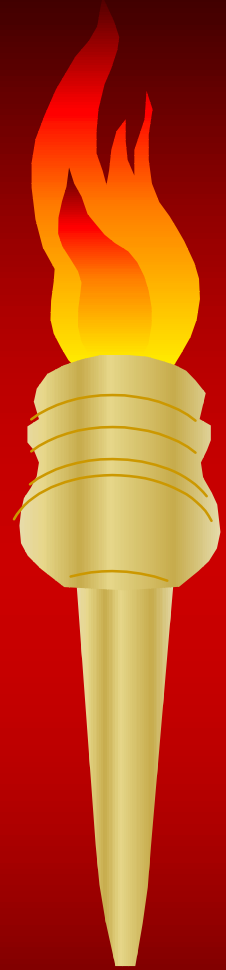
Дискриминантный анализ

- Методика создания классифицирующих правил
- Пытается создать серию уравнений, которые позволили бы правильно классифицировать группы
- В принципе, является аналогом ANOVA, но только «перевернутым» (рост позволяет отличить мужчин от женщин)



Особый случай - анализ выживаемости

- Нечто среднее между анализом качественных и количественных показателей (качественный показатель - цензурирование, количественный - время под наблюдением)
- Основа унивариантного анализа - построение таблиц дожития (life tables) и кривых дожития (Kaplan-Meier)
- Многомерные методики
 - Только качественные независимые переменные
 - Стратифицированный анализ Mantel-Haenszel
 - Коррекция кривых дожития
 - Качественные и количественные независимые переменные
 - Моделирование выживаемости (Weibull и др.)
 - Модель пропорционального риска (Cox)



Снижение размерности

■ Количественные переменные

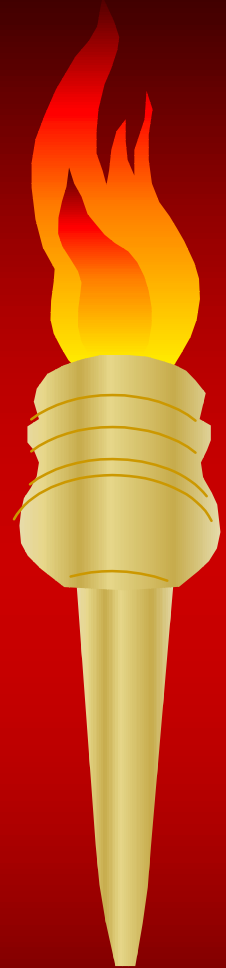
- Анализ главных компонент
- Факторный анализ
- Многомерное шкалирование
- Иерархический кластерный анализ

■ Качественные переменные

- Корреспондентский анализ

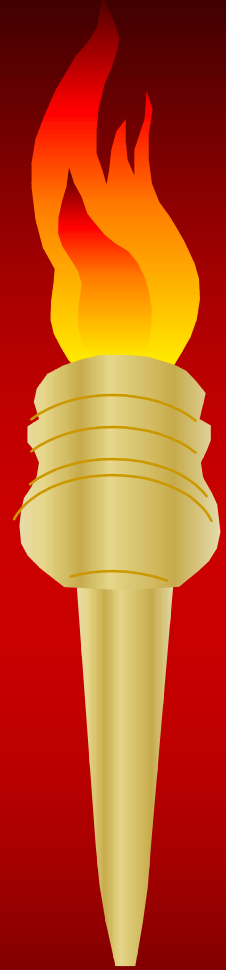
■ Объединение наблюдений

- Кластерный анализ



Факторный анализ

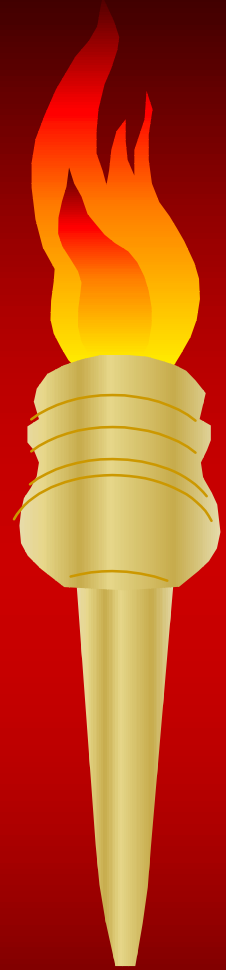
- Вместе с анализом главных компонент - основной метод снижения размерности
- Пытается найти новые переменные, которые бы комбинировали в себе несколько старых (рост, измеренный в метрах, футах и локтях)
- Отличие от метода главных компонент - отсутствие “уникальных” факторов - факторов, представленных одной переменной
- Единицей анализа является корреляционная матрица



Факторный анализ

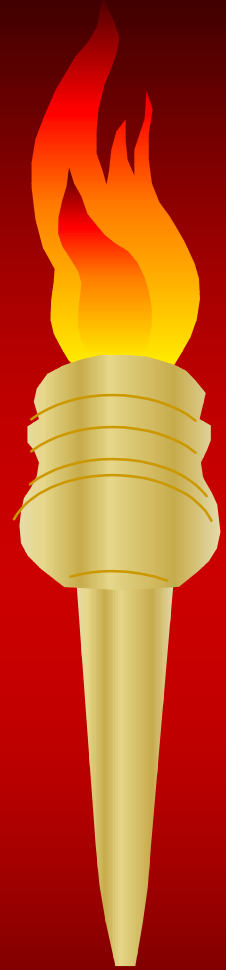
■ Выделяют несколько вариантов

- Стандартный (R) - факторный анализ переменных
- Обратный (Q) - факторный анализ наблюдений
- Временной (T) - факторный анализ временных тенденций (у разных людей по одной переменной)
- Обратный временной (S) - факторный анализ наблюдений в зависимости от времени
- Одного наблюдения (O) - факторный анализ временных тенденций по разным показателям у одного субъекта

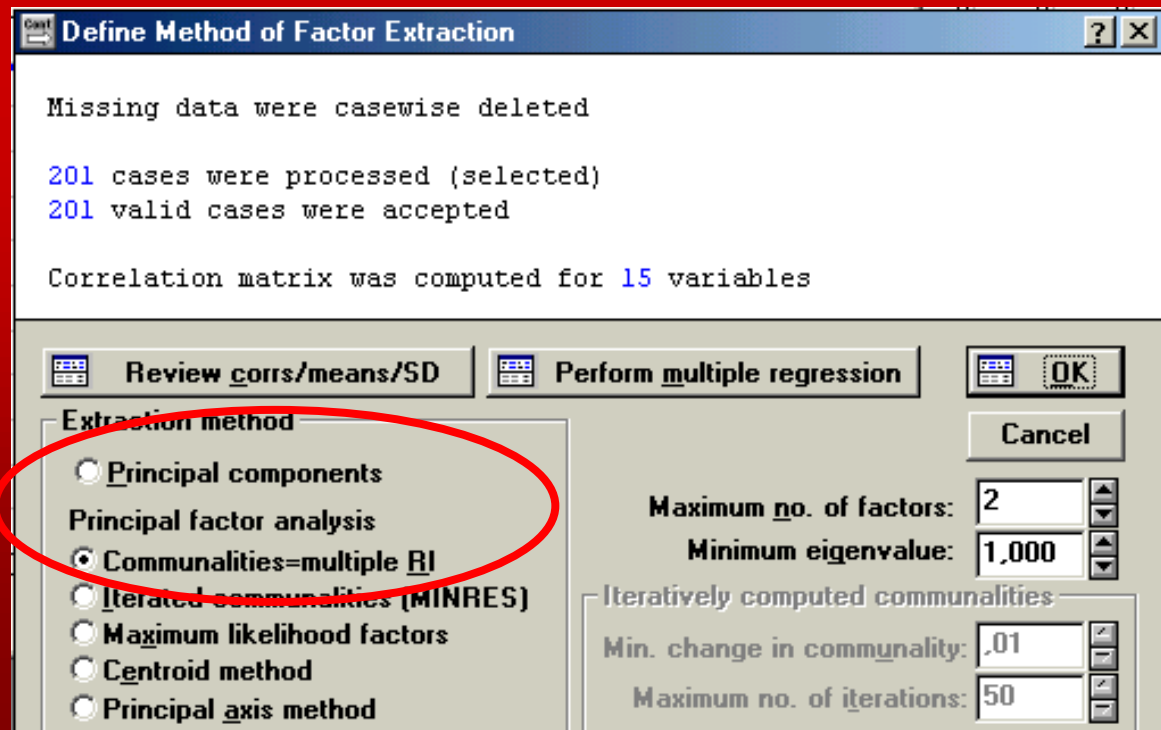
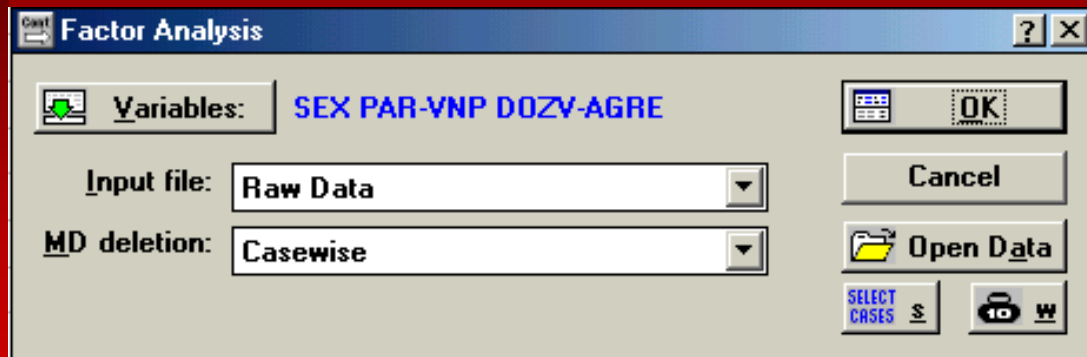


Факторный анализ

- Требуется не менее 10 (5) наблюдений на каждую включенную в анализ переменную, но общее число наблюдений должно быть больше 100
- Факторный анализ отличается от компонентного тем, что находится в диагональных элементах корреляционной матрицы (1 в компонентном, иное число - максимальное значение r в строке, множ. Коэффициент корреляции и т.п. - в факторном)



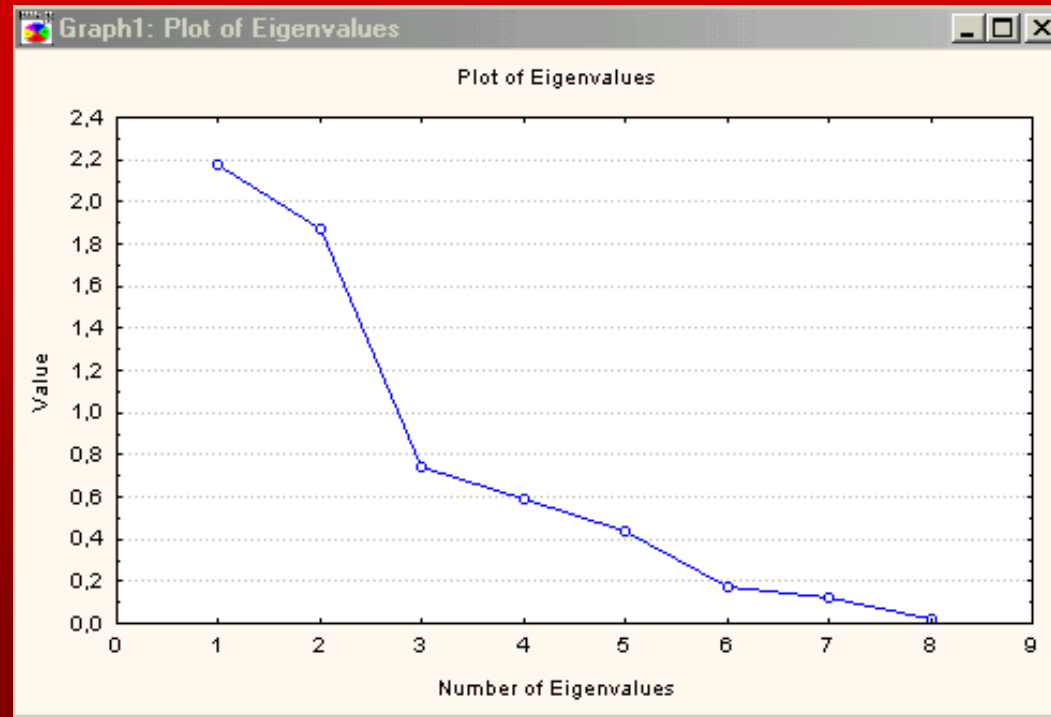
Факторный анализ



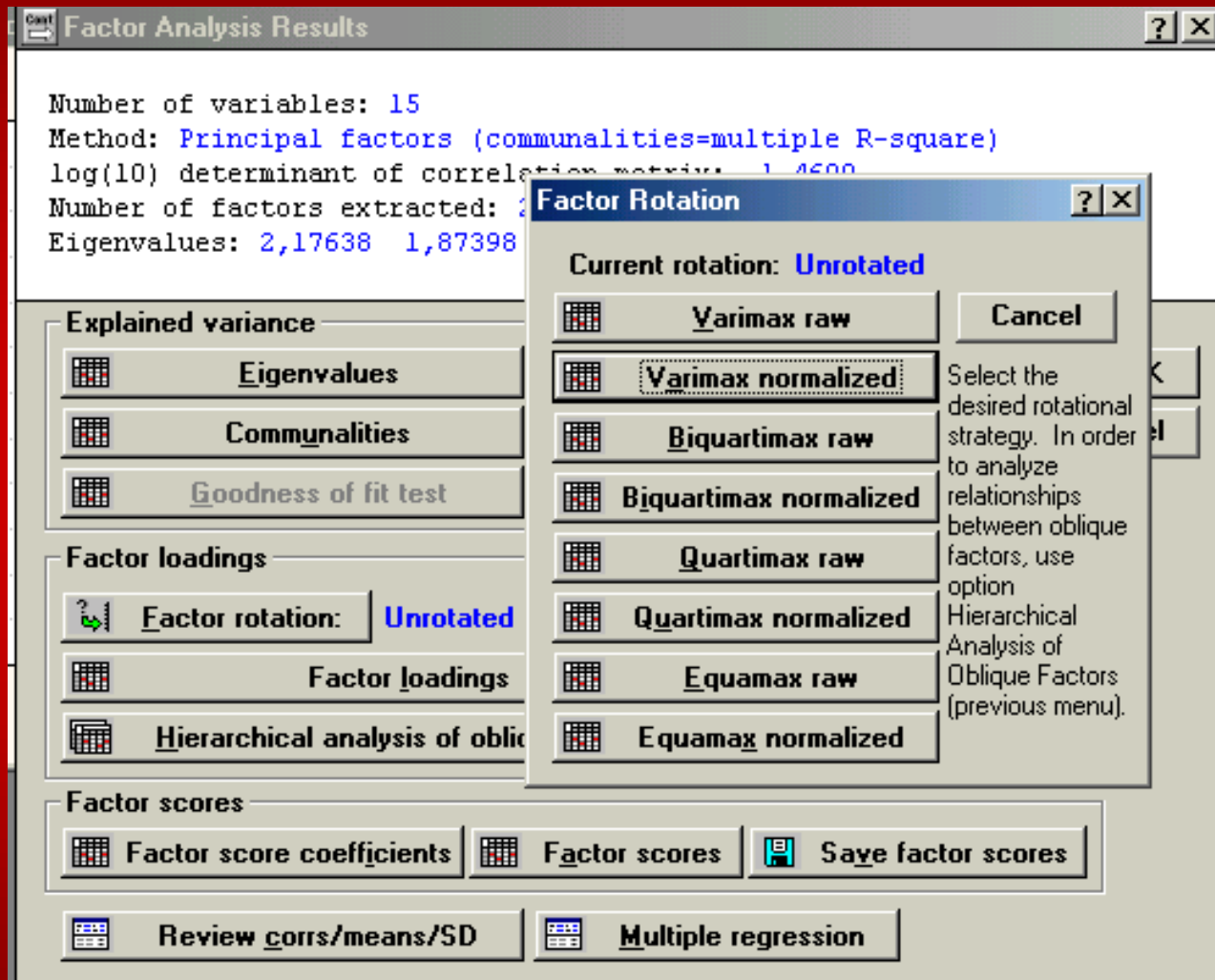
Факторный анализ

■ Методы определения числа факторов

- Критерий Кайзера (собственное значение фактора - Eigenvalue больше 1)
- Критерий Жоффе (Eigenvalue > 0.7)
- Процент объясненной дисперсии (90% и более)
- Scree plot



Факторный анализ



Factor Analysis Results

Number of variables: 15
Method: Principal factors (communalities=multiple R-square)
log(10) determinant of correlation matrix: 1,4500
Number of factors extracted: 2
Eigenvalues: 2,17638 1,87398

Explained variance

- Eigenvalues
- Communalities
- Goodness of fit test

Factor loadings

- Factor rotation: Unrotated
- Factor loadings
- Hierarchical analysis of oblique factors

Factor scores

- Factor score coefficients
- Factor scores
- Save factor scores

Review corrs/means/SD Multiple regression

Factor Rotation

Current rotation: Unrotated


- Varimax raw
- Varimax normalized
- Biquartimax raw
- Biquartimax normalized
- Quartimax raw
- Quartimax normalized
- Equamax raw
- Equamax normalized

Cancel

Select the desired rotational strategy. In order to analyze relationships between oblique factors, use option Hierarchical Analysis of Oblique Factors (previous menu).

Вращение факторов - Varimax или Oblimin - необходимо для интерпретируемости результатов - простая структура

Факторный анализ



Factor Loadings (Varimax normalized) [eysenk.sta]

Continue... Extraction: Principal factors (comm.=multiple R-square)
(Marked loadings are > ,400000)

Variable	Factor 1	Factor 2
SEX	-,434748	,091442
PAR	,414157	-,033518
VNP	-,442744	,101523
DOZV	,437745	-,147297
UDOV	,266659	-,660534
KONF	,080965	,742981
BEZL	,557509	,391287
PORN	,308654	,161431
ZAST	-,249146	,612586
STID	-,184641	,081979
DOMI	,035320	-,543752
OTVR	-,492575	-,038430
VOZB	,668622	,055499
BIOL	,195436	-,062284
AGRE	,020808	,312492
Expl. Var	2,057647	1,992715
Prp. Totl	,137176	,132848

Факторный анализ (SAS)

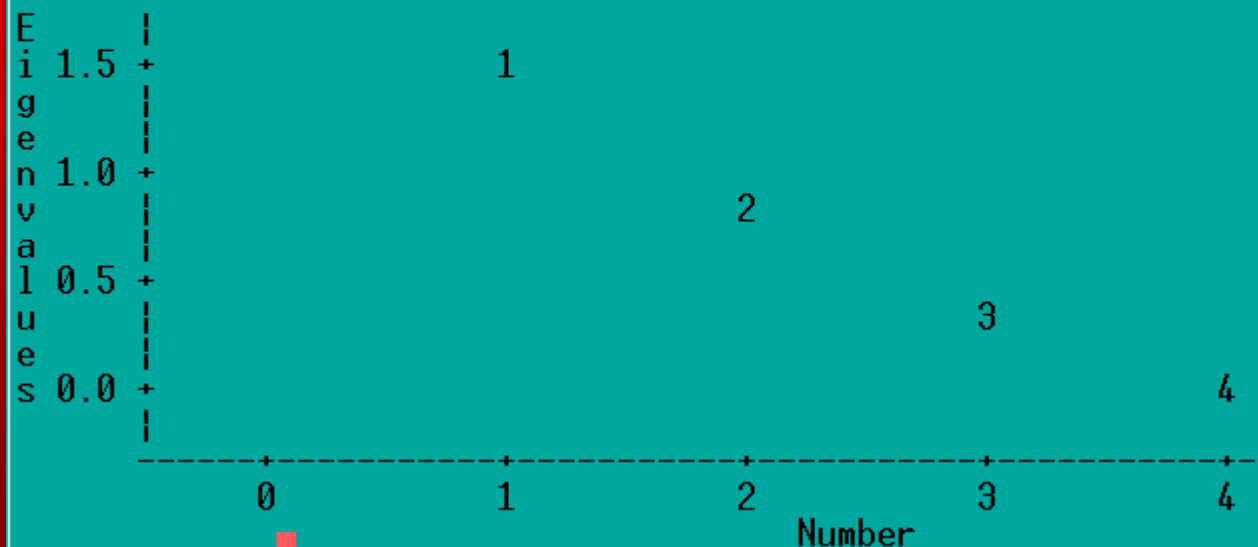
PROGRAM EDITOR

Command ==> ■

```
00001 PROC FACTOR DATA=mydat.mu99 PRIORS=MAX SCREE ROTATE=V;  
00002 VAR CH HDL TG SBP DBP;  
00003 RUN;
```

Initial Factor Method: Principal Factors

Scree Plot of Eigenvalues



Факторный анализ (SAS)

Rotation Method: Varimax

Rotated Factor Pattern

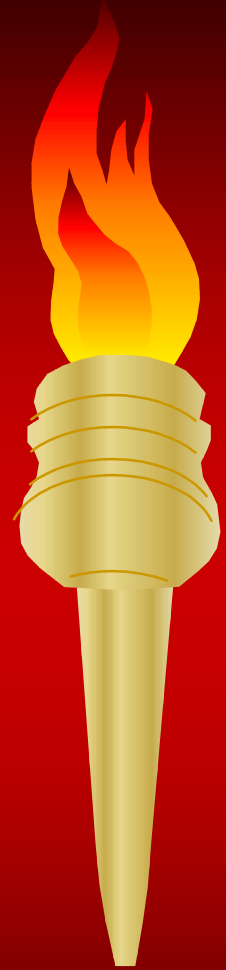
	FACTOR1	FACTOR2	FACTOR3	
CH	0.08055	0.63287	0.04089	V1 TOTAL CHOLESTEROL
HDL	0.06757	-0.04088	0.52926	V1 HDL CHOLESTEROL
TG	0.08077	0.53919	-0.36173	TRIG1 MINUS TRIGBLK1
SBP	0.85181	0.08656	0.11457	
DBP	0.85557	0.10789	0.00277	

Variance explained by each factor

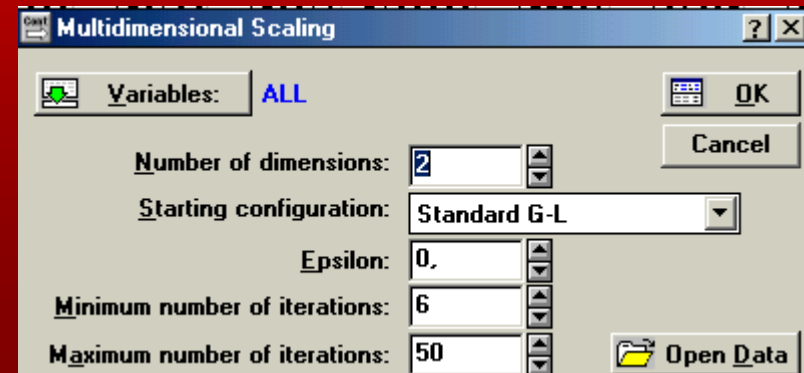
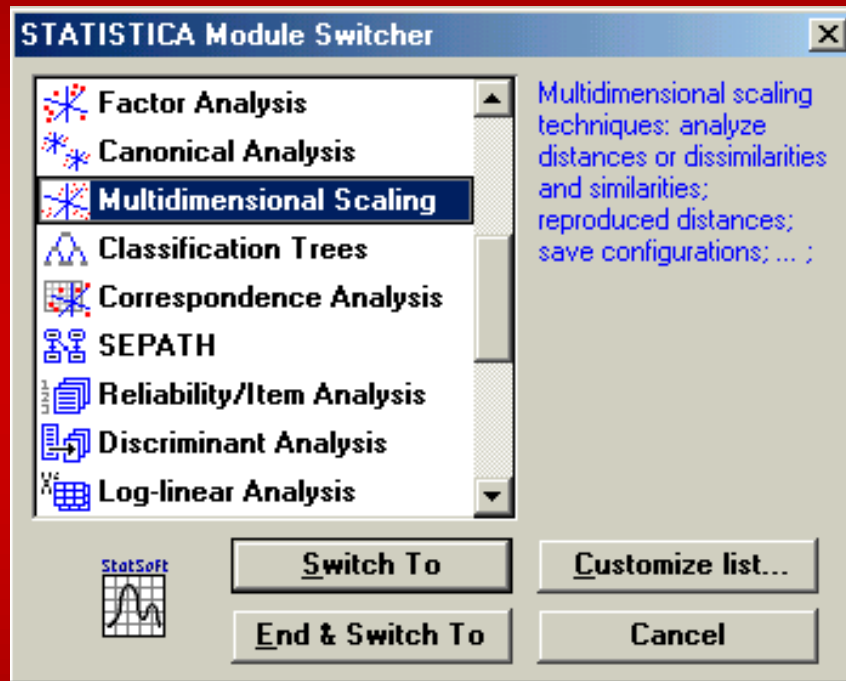
FACTOR1	FACTOR2	FACTOR3
1.475154	0.712047	0.425780

Многомерное шкалирование

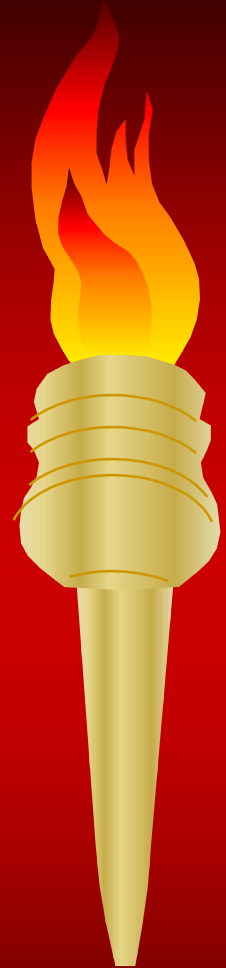
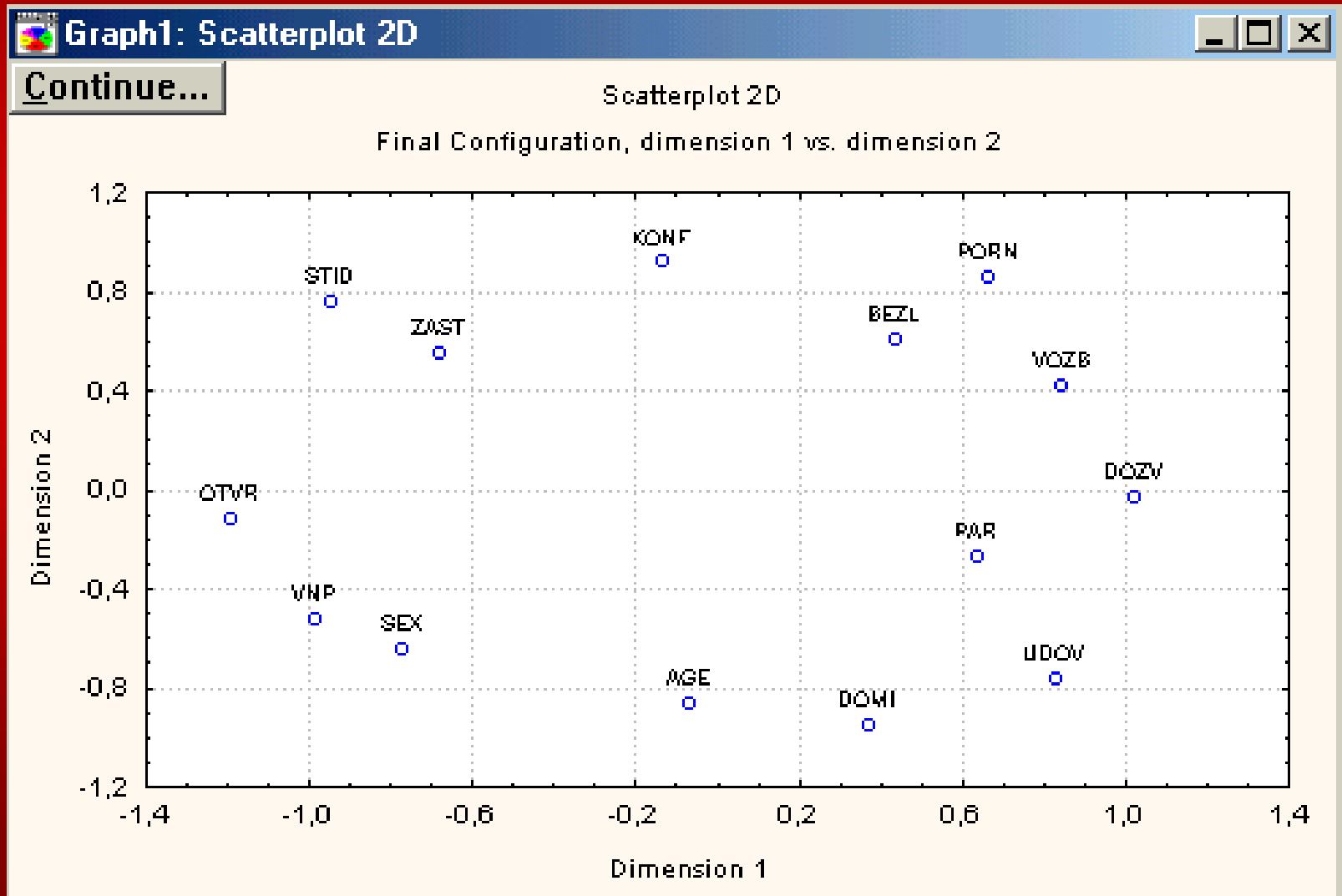
- Метод анализа матрицы дистанций
 - Истинные дистанции
 - Коэффициенты корреляции ($1-r$)
 - Другие индексы
- Аналог географической карты



Многомерное шкалирование

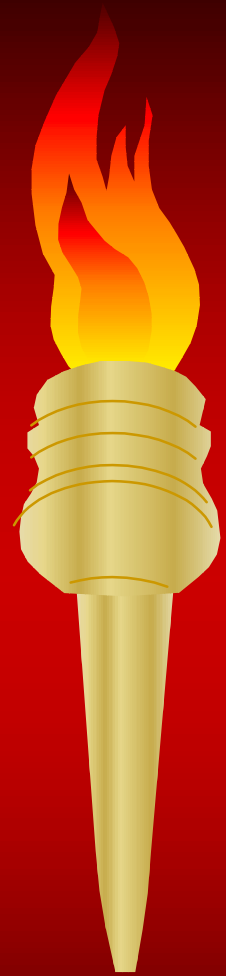


Многомерное шкалирование

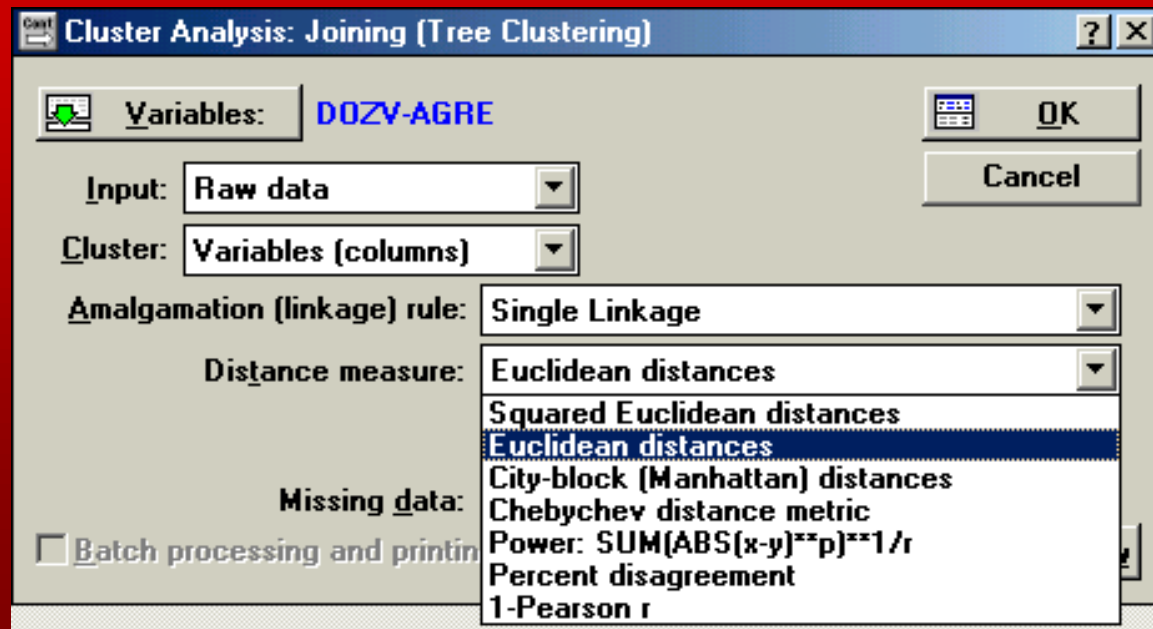
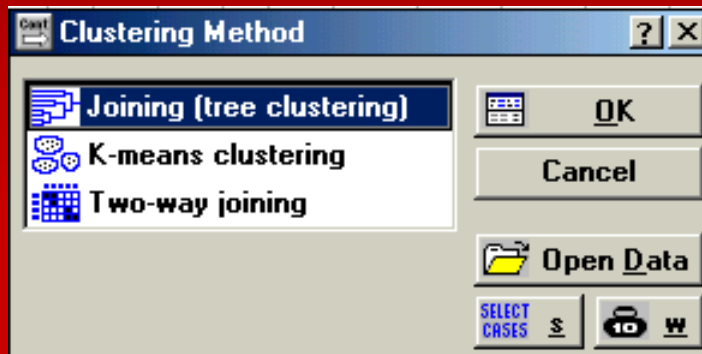


Кластерный анализ

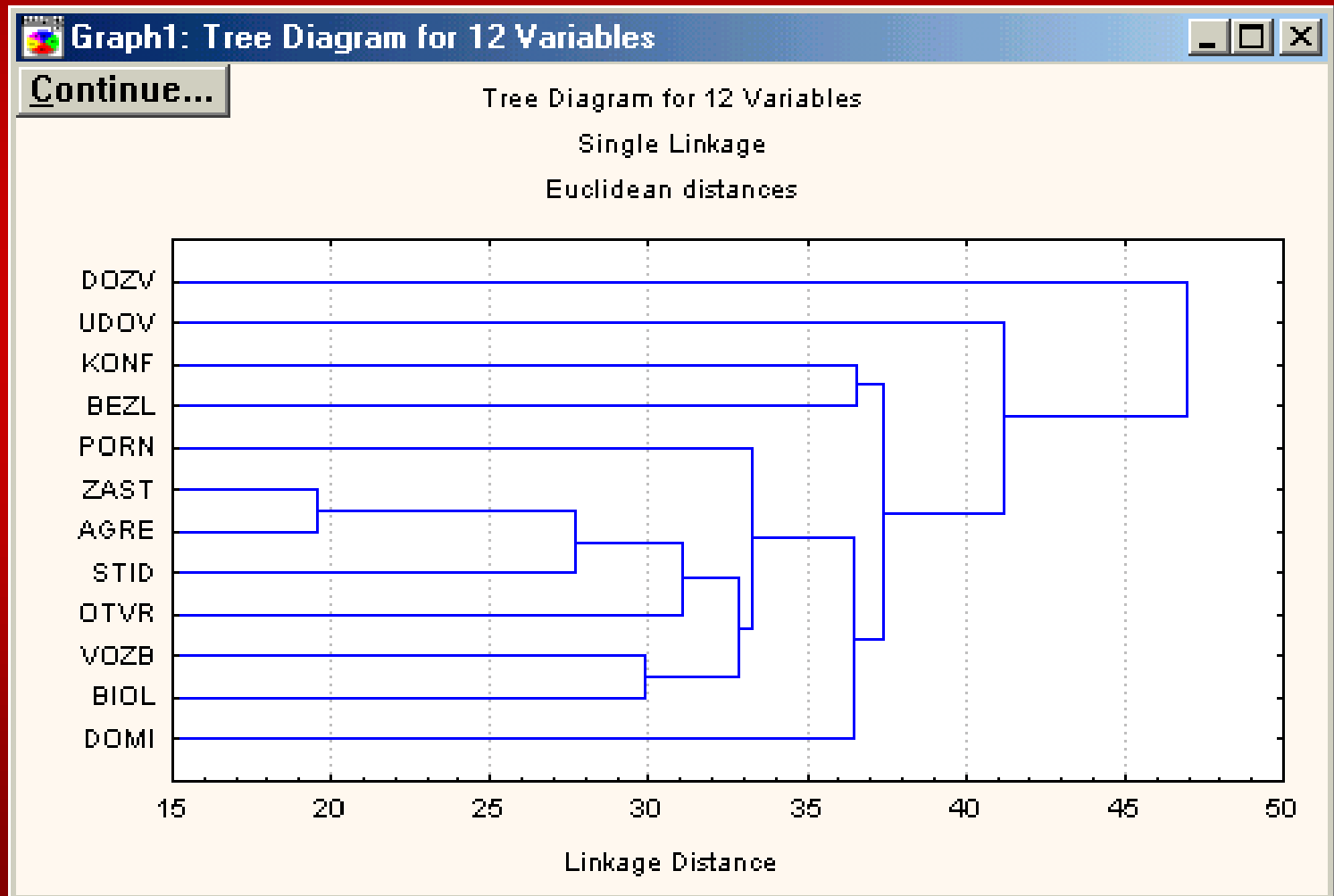
- Может использоваться для снижения размерности/классификации переменных (иерархический анализ)
- Может использоваться для объединения наблюдений на основании ряда показателей (групп - кластеров-пациентов) k-means



Кластерный анализ

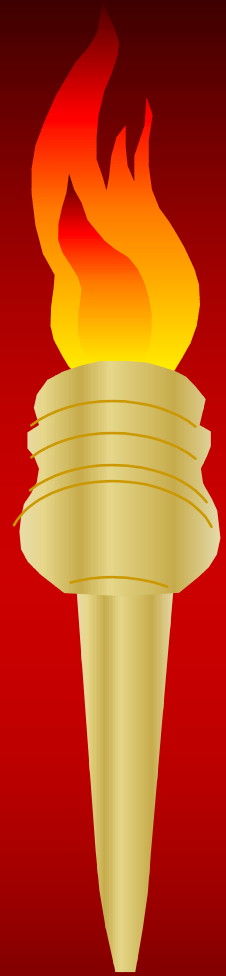


Кластерный анализ



Кластерный анализ (k-means)

- Необходимо заранее решить, какое количество кластеров следует выделить
- Работает как ANOVA наоборот:
 - наблюдения объединяются таким образом, чтобы расстояния между группами (между средними групп) были максимальными
 - поэтому группы почти всегда будут статистически достоверно отличаться по включенным параметрам



Кластерный анализ (k-means)

Cluster Analysis: K-Means Clustering

Variables: DOZV-AGRE

Cluster: Cases (rows)

Number of clusters: 2

Number of iterations: 10

Missing data: Casewise deleted

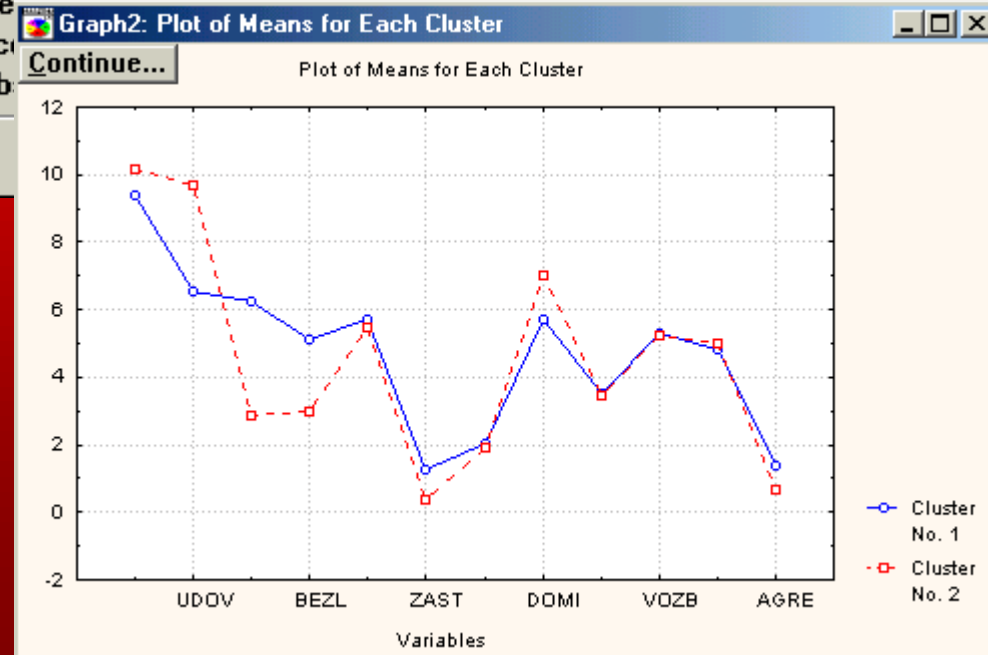
Initial cluster centers

- Choose observations to maximize initial between-cluster distances
- Sort distances and take observations at each step
- Choose the first N (Number of clusters) observations

Batch processing and printing

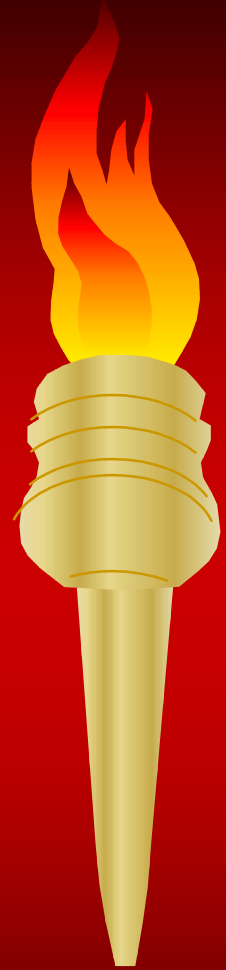
OK

Cancel




Корреспондентский анализ

- Аналог факторного анализа для качественных переменных
- Результаты также сходны с многомерным шкалированием
- Был разработан для лингвистического анализа



Корреспондентский анализ



Multiple Correspondence Analysis (MCA): Table Specifications

Method

Correspondence Analysis (CA) **Multiple Correspondence Analysis (MCA)**

Input

Raw data (requires tabulation)
 Frequencies with grouping variables
 Frequencies w/out grouping vars
(Data must specify a Burt table for MCA)

You can tabulate variables with codes, or input a [stacked] table of frequencies, with/out coding variables.

Variables (Factors in Burt Table)

Codes for grouping variables

Supplementary columns (variables)

OK Cancel Open Data

Multiple Correspondence Analysis Results

No. of columns in table to be analyzed: 16
Variables and number of categories:
SEX(2) VNP(12) NIB(2)
Eigenvalues: ,4678 ,3857 ,3333 ,3333 ,3333 ,3333 ,3333 ,3333 ,3333 ...
Total chi-square=2685,41 df=225 p=0,000
(Chi-square (df, p) only valid if table were an ordinary two-way table)

Column coordinates

Eigenvalues **Plot**

Unstandardized matrices

Number of dimensions

Number of dimensions: 2

Cumulative contribution to inertia:
90,0 % (or more)

Plots of coordinates

Column, 1D 2D 3D

Plot selected dimensions only
 Truncate labels to 2 chars
 Use identical X/Y/Z scales

Review tables

Include supplementary points

<input type="checkbox"/> Observed frequenc.	<input type="checkbox"/> Expected (chi-square)
<input type="checkbox"/> Row percentages	<input type="checkbox"/> Obs. minus expected
<input type="checkbox"/> Column percentages	<input type="checkbox"/> Contrib. to chi-square
<input type="checkbox"/> Total percentages	<input type="checkbox"/> Standardized deviates

To produce 3D histograms of frequencies, percentages, etc., use the default Quick Stats Graphs option in these Scrollsheets.

Print summary

Корреспондентский анализ

