

Preprocessing

K. Gibert

Dep. Statistics and Operations Research

*Knowledge Engineering and machine learning group at
Intelligent Data Science and Artificial Intelligence Research Center*

University Research Institute on Science and Technology for Sustainability

Universitat Politècnica de Catalunya-BarcelonaTech

Apunts disponibles (histoIntrodesc)

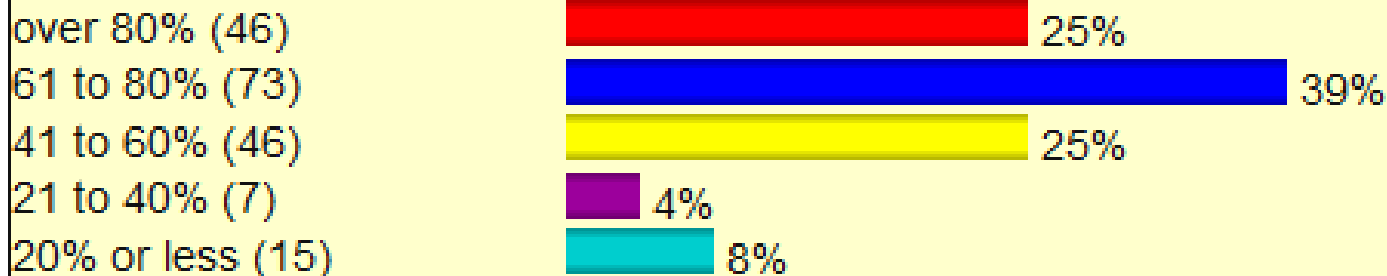
Impact of Preprocessing in real Data Mining projects

Data preparation part in data mining projects



Poll

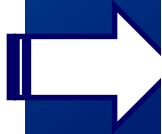
What % of time in your data mining project(s) is spent on data cleaning and preparation [187 votes total]



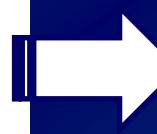
http://www.kdnuggets.com/polls/2003/data_preparation.htm

Preprocessing

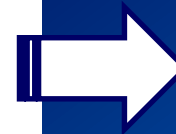
**Data
Quality**



**Quality
of
Analysis**

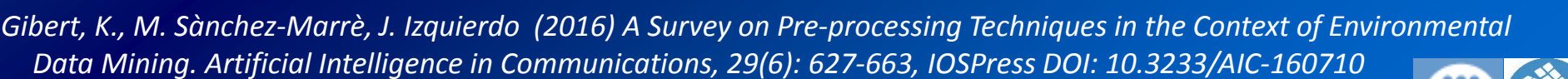


**Quality
of
Results**



**Quality of
Decisions**

[Gibert Aicomm2016]



Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Determine Data Matrix

Which data matrix rows?

Define target population

Objects selection

Which data matrix columns?

Determine objects description

Variables selection

Objects Selection

Inclusion/Exclusion criteria
Filtering

Select from a data base or data warehouse or
from real individuals (costs are different)

- Experimental data (experimental design)
- Observational data (sample theory)



REPORT IN
FINAL DOC

Define the target population

Determines scope of conclusions

Goals' oriented Variables Selection

- Often expert-guided
(highly related with goal of analysis)
- Be maximalists
 - Eliminate irrelevant or redundant information is less risky than detect lack of relevant things to be added in a second wave
- Technically, to complete a final submatrix is highly costly (in both time and resources)

Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Missing data

(empty cells in data matrix)

- ▶ Types and diagnosis
- ▶ Little's test
- ▶ A simple descriptive alternative
- ▶ Some methods
 - ▶ Knn
 - ▶ The MIMMI method
 - ▶ MICE
 - ▶ Interpolation (for time series)

Missing data

(empty cells in data matrix)

- ▶ Random missing
 - non problematic
 - casual
 - follow same distribution as present data
 - imputation is easy: mean, 0
- ▶ Non random missing: absence is informative
 - come from some particular part of population
 - probably correspond to special values
 - difficult to induce from the present data
 - imputation is much difficult
 - very critical
 - very dangerous to ignore those individuals
 - asking religion in israel (muslims do not answer)
 - Asking age to a lady over 45
 - Frequency of observations (microbio tests in water)
- ▶ Non applicable value (non-random, structural)
 - salary of a non-working person
 - number of pregnancies of a man
 - number of cigarretes of a non-smoker person
 - age of menopause

*dangerous to ignore
(specially if non random)*

Missing data

Diagnoses

Little's MCAR test

H_0 : Missings are completely at random (MCAR)

H_1 : Missings are not random

$$d^2 = \sum_{j=1}^J n_j (\bar{X}_j - \bar{X}_j^*)^T \frac{1}{\hat{\Sigma}} (\bar{X}_j - \bar{X}_j^*) \sim \chi^2_{\sum r_j - K}$$

$j=1:J$ missing patterns (subsets of missing variables in a case)

n_j cases in missing pattern j

\bar{X}_j maximum likelihood estimates of the grand means

\bar{X}_j^* means local to cases in missing pattern

$\hat{\Sigma}$ maximum likelihood estimate of the covariance matrix

r_j number of complete variables for pattern j

K total number of variables

Searches significant differences in means conditioned to a certain subset of missing variables (pattern j)

1	Països	Zona	PobTotal	INBzcap	Esper Vida	%Nounat sBaizPes	UsInstal Saneja	VIH	%coneixP resHomes	\$ConeixP resDones	Alfabet Homes	Alfabet Dones	Mortalitat	Natalitat	Economia	despesaSalut/Def ensa	Règim
2	CAMBOYA	AsiaOr+Pacífic	13798	320	57	11	16	170		64	85	64	11	31	enDesenvolupament		Monarquia
3	COREA	AsiaOr+Pacífic	22384		63	7	59						11	16	enDesenvolupament		Republica
4	FIJI	AsiaOr+Pacífic	841	2690	68	10	98	0,6			94	91	6	23	enDesenvolupament	MMesDefensa	Republica
5	FILIPINES	AsiaOr+Pacífic	81617	1170	71	20	73	9	59	44	93	93	5	25	enDesenvolupament	MesSalut	Republica
6	ILLES COOK	AsiaOr+Pacífic	18			3	100								Subdesenvolupat		Monarquia
7	ILLES MARSHALL	AsiaOr+Pacífic	60	2370		12	82								Subdesenvolupat		Republica
8	ILLES SALOMÓ	AsiaOr+Pacífic	466	550	63	13	31						7	33	Subdesenvolupat		Monarquia
9	INDONÈSIA	AsiaOr+Pacífic	220077	1140	67	9	52	110		23	92	83	7	21	enDesenvolupament	MesSalut	Republica
10	KIRIBATI	AsiaOr+Pacífic	97	970		5	39								Subdesenvolupat		Republica
11	LAO, REPÚBLICA DEMOCRÀTICA	AsiaOr+Pacífic	5792	390	55	14	24	1,7			77	61	12	35	enDesenvolupament		Republica
12	MALÀSIA	AsiaOr+Pacífic	24894	4650	73	9		52			92	85	5	22	enDesenvolupament	MesSalut	Monarquia
13	MICRONÈSIA	AsiaOr+Pacífic	110	1990	68	18	28						6	31	Subdesenvolupat		Republica
14	MONGÒLIA	AsiaOr+Pacífic	2614	590	65	7	59	0,5		77	98	98	7	22	enDesenvolupament	MesSalut	Republica
15	MYANMAR	AsiaOr+Pacífic	50004	220	61	15	73	330			94	86	10	20	enDesenvolupament	MesSalut	Dictadura
16	NAURU	AsiaOr+Pacífic	13												Subdesenvolupat		Republica
17	NIUE	AsiaOr+Pacífic	1			0	100								Subdesenvolupat		Monarquia
18	PALAU	AsiaOr+Pacífic	20	6870		9	83								Subdesenvolupat		Republica
19	PAPUA NOVA GUINEA	AsiaOr+Pacífic	5772	580	56	11	45	16			63	51	10	30	enDesenvolupament	MMesDefensa	Monarquia
20	SAMOA	AsiaOr+Pacífic	184	1860	71	4	100				99	98	6	28	Subdesenvolupat		Monarquia
21	SINGAPUR	AsiaOr+Pacífic	4273	24220	79	8		4,1			97	89	5	9	enDesenvolupament	MesSalut	Republica
22	TAILÀNDIA	AsiaOr+Pacífic	63694	2540	70	9	99	0,2			95	91	7	16	enDesenvolupament	MMesDefensa	Monarquia
23	TIMOR-LESTE	AsiaOr+Pacífic	887	550	56	12	33			6			12	50	Subdesenvolupat		Republica
24	TOKELAU	AsiaOr+Pacífic													Subdesenvolupat		Monarquia
25	TONGA	AsiaOr+Pacífic	102	1830	72	0	97				99	99	6	24	Subdesenvolupat		Monarquia
26	TUVALU	AsiaOr+Pacífic	10			5	88								Subdesenvolupat		Monarquia
27	VANUATU	AsiaOr+Pacífic	207	1340	69	6	50						6	31	enDesenvolupament		Republica
28	VIETNAM	AsiaOr+Pacífic	83123	550	71	9	41	220		60	94	87	6	20	enDesenvolupament		Republica
29	XINA	AsiaOr+Pacífic	1307989	1290	72	4	44	840			95	87	7	13	enDesenvolupament	MesSalut	Republica
30	ANGOLA	AfricaOr+Merid	15490	1030	41	12	30	240			82	54	22	48	enDesenvolupament	MesSalut	Republica
31	BOTSWANA	AfricaOr+Merid	1769	4340	35	10	41	350	90	93	76	82	27	26	enDesenvolupament	MesSalut	Republica
32	BURUNDI	AfricaOr+Merid	7282	90	44	16	36	250		47	67	52	19	45	enDesenvolupament	MesSalut	Republica
33	COMORAS	AfricaOr+Merid	777	530	64	25	23			41	63	49	7	36	enDesenvolupament		Republica
34	ERITREA	AfricaOr+Merid	4232	180	54	21	9	60		62			11	39	enDesenvolupament		Republica
35	ETIOPIA	AfricaOr+Merid	75600	110	48	15	6	1500			49	34	16	41	enDesenvolupament	MesSalut	Republica
36	KENYA	AfricaOr+Merid	33467	460	48	10	48	1200	68	59	78	70	15	39	enDesenvolupament	MMesDefensa	Republica
37	LESOTHO	AfricaOr+Merid	1798	740	35	14	37	320		58	74	90	25	28	enDesenvolupament	MMesDefensa	Monarquia
38	MADAGASCAR	AfricaOr+Merid	18113	300	56	17	33	140	54	49	76	65	12	39	enDesenvolupament	MMesDefensa	Republica

Dades UNICEF 2004 <http://www.unicef.org/spanish/infobycountry/index.html> _missing patterns

Missing pattern 1: (ConeixPresHomes)

Missing pattern 2: (ConeixPresHomes+DespesaSalut/Defensa)

Missing pattern 3: (ConeixPresHomes+ConeixPresDones)

Missing pattern 4: (ConeixPresHomes+ConeixPresDones+DespSal/Def)

Missing pattern 5: (ConeixPresHomes+AlfabetDones+DespSal/Def)

Missing pattern 6: (VIH+ConeixPresHomes+ DespSal/Def)

Missing pattern 7: (VIH+ConeixPresHomes+ ConeixPresHomes+ConeixPresDones+DespSal/Def)

Missing pattern 8: (ConeixPresHomes+AlfabetHomes+AlfabetDones+DespSal/Def)

n1=4, x*1=1
n2=2, x*2=2
n3=9, x*3=2
n4=1, x*4=3
n5=1, x*5=3
n6=1, x*6=3
n7=2, x*7=4
n8=1, x*8=4

			PobTotal		Esper	%Nounat	UsInstal		%coneixP	\$ConeixP	Alfabet	Alfabet				despesaSalut/Def		
1	Països	Zona	I	INBxcap	Vida	sBaixPes	Saneja	VIH	resHomes	resDones	Homes	Dones	Mortalitat	Natalitat	Economia	ensa	Règim	
2	CAMBOYA	AsiaOr+Pacífic	13798	320	57	11	16	170		64	85	64		11	31	enDesenvolupament		Monarquia
3	COREA	AsiaOr+Pacífic	22384		63	7	59							11	16	enDesenvolupament		Republica
4	FIJI	AsiaOr+Pacífic	841	2690	68	10	98	0,6			94	91	6	23	enDesenvolupament	MMesDefensa	Republica	
5	FILIPINES	AsiaOr+Pacífic	81617	1170	71	20	73	9	59	44	93	93	5	25	enDesenvolupament	MesSalut	Republica	
6	ILLES COOK	AsiaOr+Pacífic	18			3	100								Subdesenvolupat		Monarquia	
7	ILLES MARSHALL	AsiaOr+Pacífic	60	2370		12	82								Subdesenvolupat		Republica	
8	ILLES SALOMÓ	AsiaOr+Pacífic	466	550	63	13	31						7	33	Subdesenvolupat		Monarquia	
9	INDONÈSIA	AsiaOr+Pacífic	220077	1140	67	9	52	110		23	92	83	7	21	enDesenvolupament	MesSalut	Republica	
10	KIRIBATI	AsiaOr+Pacífic	97	970		5	39								Subdesenvolupat		Republica	
11	LAO, REPÚBLICA DEMOCRÀTICA	AsiaOr+Pacífic	5792	390	55	14	24	1,7			77	61	12	35	enDesenvolupament		Republica	
12	MALÀSIA	AsiaOr+Pacífic	24894	4650	73	9		52			92	85	5	22	enDesenvolupament	MesSalut	Monarquia	
13	MICRONÈSIA	AsiaOr+Pacífic	110	1990	68	18	28						6	31	Subdesenvolupat		Republica	
14	MONGÒLIA	AsiaOr+Pacífic	2614	590	65	7	59	0,5		77	98	98	7	22	enDesenvolupament	MesSalut	Republica	
15	MYANMAR	AsiaOr+Pacífic	50004	220	61	15	73	330			94	86	10	20	enDesenvolupament	MesSalut	Dictadura	
16	NAURU	AsiaOr+Pacífic	13												Subdesenvolupat		Republica	
17	NIUE	AsiaOr+Pacífic	1			0	100								Subdesenvolupat		Monarquia	
18	PALAU	AsiaOr+Pacífic	20	6870		9	83								Subdesenvolupat		Republica	
19	PAPUA NOVA GUINEA	AsiaOr+Pacífic	5772	580	56	11	45	16			63	51	10	30	enDesenvolupament	MMesDefensa	Monarquia	
20	SAMOA	AsiaOr+Pacífic	184	1860	71	4	100				99	98	6	28	Subdesenvolupat		Monarquia	
21	SINGAPUR	AsiaOr+Pacífic	4273	24220	79	8		4,1			97	89	5	9	enDesenvolupament	MesSalut	Republica	
22	TAILÀNDIA	AsiaOr+Pacífic	63694	2540	70	9	99	0,2			95	91	7	16	enDesenvolupament	MMesDefensa	Monarquia	
23	TIMOR-LESTE	AsiaOr+Pacífic	887	550	56	12	33			6			12	50	Subdesenvolupat		Republica	
24	TOKELAU	AsiaOr+Pacífic													Subdesenvolupat		Monarquia	
25	TONGA	AsiaOr+Pacífic	102	1830	72	0	97				99	99	6	24	Subdesenvolupat		Monarquia	
26	TUVALU	AsiaOr+Pacífic	10			5	88								Subdesenvolupat		Monarquia	
27	VANUATU	AsiaOr+Pacífic	207	1340	69	6	50						6	31	enDesenvolupament		Republica	
28	VIETNAM	AsiaOr+Pacífic	83123	550	71	9	41	220		60	94	87	6	20	enDesenvolupament		Republica	
29	XINA	AsiaOr+Pacífic	1307989	1290	72	4	44	840			95	87	7	13	enDesenvolupament	MesSalut	Republica	
30	ANGOLA	AfricaOr+Merid	15490	1030	41	12	30	240			82	54	22	48	enDesenvolupament	MesSalut	Republica	
31	BOTSWANA	AfricaOr+Merid	1769	4340	35	10	41	350	90	93	76	82	27	26	enDesenvolupament	MesSalut	Republica	
32	BURUNDI	AfricaOr+Merid	7282	90	44	16	36	250		47	67	52	19	45	enDesenvolupament	MesSalut	Republica	
33	COMORAS	AfricaOr+Merid	777	530	64	25	23			41	63	49	7	36	enDesenvolupament		Republica	
34	ERITREA	AfricaOr+Merid	4232	180	54	21	9	60		62			11	39	enDesenvolupament		Republica	
35	ETIOPIA	AfricaOr+Merid	75600	110	48	15	6	1500			49	34	16	41	enDesenvolupament	MesSalut	Republica	
36	KENYA	AfricaOr+Merid	33467	460	48	10	48	1200	68	59	78	70	15	39	enDesenvolupament	MMesDefensa	Republica	
37	LESOTHO	AfricaOr+Merid	1798	740	35	14	37	320		58	74	90	25	28	enDesenvolupament	MMesDefensa	Monarquia	
38	MADAGASCAR	AfricaOr+Merid	18113	300	56	17	33	140	54	49	76	65	12	39	enDesenvolupament	MMesDefensa	Republica	

- n1=4, x*1=1
- n2=2, x*2=2
- n3=9, x*3=2
- n4=1, x*4=3
- n5=1, x*5=3
- n6=1, x*6=3
- n7=2, x*7=4
- n8=1, x*8=4

- Missing pattern 9: (VIH+CPH+ CPD+AlfaH+AlfaD+DespSal/Def)
- Missing pattern 10: (INB+VIH+CPH+ CPD+AlfaH+AlfaD+DespSal/Def))
- Missing pattern 1: (EV+ VIH+CPH+ CPD+AlfaH+AlfaD+M+N+D/Def)
- Missing pattern 12: (INB+ EV+ VIH+CPH+ CPD+AlfaH+AlfaD+M+N+D/Def)
- Missing pattern 13: (INB+ EV+NBp+IS+ VIH+CPH+ CPD+AlfaH+AlfaD+M+N+D/Def)
- Missing pattern 14: (Pob+INB+ EV+NBp+IS+ VIH+CPH+ CPD+AlfaH+AlfaD+M+N+D/Def)

- n9=3, x*9=6
- n10=1, x*10=7
- n11=3, x*11=9
- n12=3, x*12=10
- n13=1, x*13=12
- n14=1, x*14=13

1	Països	Zona	PobTotal	INBzcap	Esper Vida	%Nounat sBaizPes	UsInstal Saneja	VIH	%coneixP resHomes	\$ConeixP resDones	Alfabet Homes	Alfabet Dones	Mortalitat	Natalitat	Economia	despesaSalut/Def ensa	Règim
2	CAMBOYA	AsiaOr+Pacífic	13798	320	57	11	16	170		64	85	64	11	31	enDesenvolupament		Monarquia
3	COREA	AsiaOr+Pacífic	22384		63	7	59						11	16	enDesenvolupament		Republica
4	FIJI	AsiaOr+Pacífic	841	2690	68	10	98	0,6			94	91	6	23	enDesenvolupament	MMesDefensa	Republica
5	FILIPINES	AsiaOr+Pacífic	81617	1170	71	20	73	9	59	44	93	93	5	25	enDesenvolupament	MesSalut	Republica
6	ILLES COOK	AsiaOr+Pacífic	18			3	100								Subdesenvolupat		Monarquia
7	ILLES MARSHALL	AsiaOr+Pacífic	60	2370		12	82								Subdesenvolupat		Republica
8	ILLES SALOMÓ	AsiaOr+Pacífic	466	550	63	13	31						7	33	Subdesenvolupat		Monarquia
9	INDONÈSIA	AsiaOr+Pacífic	220077	1140	67	9	52	110		23	92	83	7	21	enDesenvolupament	MesSalut	Republica
10	KIRIBATI	AsiaOr+Pacífic	97	970		5	39								Subdesenvolupat		Republica
11	LAO, REPÚBLICA DEMOCRÀTICA	AsiaOr+Pacífic	5792	390	55	14	24	1,7			77	61	12	35	enDesenvolupament		Republica
12	MALÀSIA	AsiaOr+Pacífic	24894	4650	73	9		52			92	85	5	22	enDesenvolupament	MesSalut	Monarquia
13	MICRONÈSIA	AsiaOr+Pacífic	110	1990	68	18	28						6	31	Subdesenvolupat		Republica
14	MONGÒLIA	AsiaOr+Pacífic	2614	590	65	7	59	0,5		77	98	98	7	22	enDesenvolupament	MesSalut	Republica
15	MYANMAR	AsiaOr+Pacífic	50004	220	61	15	73	330			94	86	10	20	enDesenvolupament	MesSalut	Dictadura
16	NAURU	AsiaOr+Pacífic	13												Subdesenvolupat		Republica
17	NIUE	AsiaOr+Pacífic	1			0	100								Subdesenvolupat		Monarquia
18	PALAU	AsiaOr+Pacífic	20	6870		9	83								Subdesenvolupat		Republica
19	PAPUA NOVA GUINEA	AsiaOr+Pacífic	5772	580	56	11	45	16			63	51	10	30	enDesenvolupament	MMesDefensa	Monarquia
20	SAMOA	AsiaOr+Pacífic	184	1860	71	4	100				99	98	6	28	Subdesenvolupat		Monarquia
21	SINGAPUR	AsiaOr+Pacífic	4273	24220	79	8		4,1			97	89	5	9	enDesenvolupament	MesSalut	Republica
22	TAILÀNDIA	AsiaOr+Pacífic	63694	2540	70	9	99	0,2			95	91	7	16	enDesenvolupament	MMesDefensa	Monarquia
23	TIMOR-LESTE	AsiaOr+Pacífic	887	550	56	12	33			6			12	50	Subdesenvolupat		Republica
24	TOKELAU	AsiaOr+Pacífic													Subdesenvolupat		Monarquia
25	TONGA	AsiaOr+Pacífic	102	1830	72	0	97				99	99	6	24	Subdesenvolupat		Monarquia
26	TUVALU	AsiaOr+Pacífic	10			5	88								Subdesenvolupat		Monarquia
27	VANUATU	AsiaOr+Pacífic	207	1340	69	6	50						6	31	enDesenvolupament		Republica
28	VIETNAM	AsiaOr+Pacífic	83123	550	71	9	41	220		60	94	87	6	20	enDesenvolupament		Republica
29	XINA	AsiaOr+Pacífic	1307989	1290	72	4	44	840			95	87	7	13	enDesenvolupament	MesSalut	Republica
30	ANGOLA	AfricaOr+Merid	15490	1030	41	12	30	240			82	54	22	48	enDesenvolupament	MesSalut	Republica
31	BOTSWANA	AfricaOr+Merid	1769	4340	35	10	41	350	90	93	76	82	27	26	enDesenvolupament	MesSalut	Republica
32	BURUNDI	AfricaOr+Merid	7282	90	44	16	36	250		47	67	52	19	45	enDesenvolupament	MesSalut	Republica
33	COMORAS	AfricaOr+Merid	777	530	64	25	23			41	63	49	7	36	enDesenvolupament		Republica
34	ERITREA	AfricaOr+Merid	4232	180	54	21	9	60		62			11	39	enDesenvolupament		Republica
35	ETIOPIA	AfricaOr+Merid	75600	110	48	15	6	1500			49	34	16	41	enDesenvolupament	MesSalut	Republica
36	KENYA	AfricaOr+Merid	33467	460	48	10	48	1200	68	59	78	70	15	39	enDesenvolupament	MMesDefensa	Republica
37	LESOTHO	AfricaOr+Merid	1798	740	35	14	37	320		58	74	90	25	28	enDesenvolupament	MMesDefensa	Monarquia
38	MADAGASCAR	AfricaOr+Merid	18113	300	56	17	33	140	54	49	76	65	12	39	enDesenvolupament	MMesDefensa	Republica

Dades UNICEF 2004

<http://www.unicef.org/spanish/infobycountry/index.html>

_missing patterns

- n1=4, x*1=1
- n2=2, x*2=2
- n3=9, x*3=2
- n4=1, x*4=3
- n5=1, x*5=3
- n6=1, x*6=3
- n7=2, x*7=4
- n8=1, x*8=4

- n9=3, x*9=6
- n10=1, x*10=7
- n11=3, x*11=9
- n12=3, x*12=10
- n13=1, x*13=12
- n14=1, x*14=13

15 patterns

1	Països	Zona	PobTotal	INBzcap	Esper Vida	%Nounat sBaizPes	UsInstal Saneja	VIH	%coneixP resHomes	\$ConeixP resDones	Alfabet Homes	Alfabet Dones	Mortalitat	Natalitat	Economia	despesaSalut/Def ensa	Règim
2	CAMBOYA	AsiaOr+Pacífic	13798	320	57	11	16	170		64	85	64	11	31	enDesenvolupament		Monarquia
3	COREA	AsiaOr+Pacífic	22384		63	7	59						11	16	enDesenvolupament		Republica
4	FIJI	AsiaOr+Pacífic	841	2690	68	10	98	0,6			94	91	6	23	enDesenvolupament	MMesDefensa	Republica
5	FILIPINES	AsiaOr+Pacífic	81617	1170	71	20	73	9	59	44	93	93	5	25	enDesenvolupament	MesSalut	Republica
6	ILLES COOK	AsiaOr+Pacífic	18			3	100								Subdesenvolupat		Monarquia
7	ILLES MARSHALL	AsiaOr+Pacífic	60	2370		12	82								Subdesenvolupat		Republica
8	ILLES SALOMÓ	AsiaOr+Pacífic	466	550	63	13	31						7	33	Subdesenvolupat		Monarquia
9	INDONÈSIA	AsiaOr+Pacífic	220077	1140	67	9	52	110		23	92	83	7	21	enDesenvolupament	MesSalut	Republica
10	KIRIBATI	AsiaOr+Pacífic	97	970		5	39								Subdesenvolupat		Republica
11	LAO, REPÚBLICA DEMOCRÀTICA	AsiaOr+Pacífic	5792	390	55	14	24	1,7			77	61	12	35	enDesenvolupament		Republica
12	MALÀSIA	AsiaOr+Pacífic	24894	4650	73	9		52			92	85	5	22	enDesenvolupament	MesSalut	Monarquia
13	MICRONÈSIA	AsiaOr+Pacífic	110	1990	68	18	28						6	31	Subdesenvolupat		Republica
14	MONGÒLIA	AsiaOr+Pacífic	2614	590	65	7	59	0,5		77	98	98	7	22	enDesenvolupament	MesSalut	Republica
15	MYANMAR	AsiaOr+Pacífic	50004	220	61	15	73	330			94	86	10	20	enDesenvolupament	MesSalut	Dictadura
16	NAURU	AsiaOr+Pacífic	13												Subdesenvolupat		Republica
17	NIUE	AsiaOr+Pacífic	1			0	100								Subdesenvolupat		Monarquia
18	PALAU	AsiaOr+Pacífic	20	6870		9	83								Subdesenvolupat		Republica
19	PAPUA NOVA GUINEA	AsiaOr+Pacífic	5772	580	56	11	45	16			63	51	10	30	enDesenvolupament	MMesDefensa	Monarquia
20	SAMOA	AsiaOr+Pacífic	184	1860	71	4	100				99	98	6	28	Subdesenvolupat		Monarquia
21	SINGAPUR	AsiaOr+Pacífic	4273	24220	79	8		4,1			97	89	5	9	enDesenvolupament	MesSalut	Republica
22	TAILÀNDIA	AsiaOr+Pacífic	63694	2540	70	9	99	0,2			95	91	7	16	enDesenvolupament	MMesDefensa	Monarquia
23	TIMOR-LESTE	AsiaOr+Pacífic	887	550	56	12	33			6			12	50	Subdesenvolupat		Republica
24	TOKELAU	AsiaOr+Pacífic													Subdesenvolupat		Monarquia
25	TONGA	AsiaOr+Pacífic	102	1830	72	0	97				99	99	6	24	Subdesenvolupat		Monarquia
26	TUVALU	AsiaOr+Pacífic	10			5	88								Subdesenvolupat		Monarquia
27	VANUATU	AsiaOr+Pacífic	207	1340	69	6	50						6	31	enDesenvolupament		Republica
28	VIETNAM	AsiaOr+Pacífic	83123	550	71	9	41	220		60	94	87	6	20	enDesenvolupament		Republica
29	XINA	AsiaOr+Pacífic	1307989	1290	72	4	44	840			95	87	7	13	enDesenvolupament	MesSalut	Republica
30	ANGOLA	AfricaOr+Merid	15490	1030	41	12	30	240			82	54	22	48	enDesenvolupament	MesSalut	Republica
31	BOTSWANA	AfricaOr+Merid	1769	4340	35	10	41	350	90	93	76	82	27	26	enDesenvolupament	MesSalut	Republica
32	BURUNDI	AfricaOr+Merid	7282	90	44	16	36	250		47	67	52	19	45	enDesenvolupament	MesSalut	Republica
33	COMORAS	AfricaOr+Merid	777	530	64	25	23			41	63	49	7	36	enDesenvolupament		Republica
34	ERITREA	AfricaOr+Merid	4232	180	54	21	9	60		62			11	39	enDesenvolupament		Republica
35	ETIOPIA	AfricaOr+Merid	75600	110	48	15	6	1500			49	34	16	41	enDesenvolupament	MesSalut	Republica
36	KENYA	AfricaOr+Merid	33467	460	48	10	48	1200	68	59	78	70	15	39	enDesenvolupament	MMesDefensa	Republica
37	LESOTHO	AfricaOr+Merid	1798	740	35	14	37	320		58	74	90	25	28	enDesenvolupament	MMesDefensa	Monarquia
38	MADAGASCAR	AfricaOr+Merid	18113	300	56	17	33	140	54	49	76	65	12	39	enDesenvolupament	MMesDefensa	Republica

$$n_3=9, x^*_3=2, r_3=13-2=11$$

$$\bar{X}_3 = [\bar{X}_{Pob}, \bar{X}_{INB}, \bar{X}_{EV}, \bar{X}_{\%NbP}, \bar{X}_{AlfaH}, \bar{X}_{AlfaD}, \bar{X}_M, \bar{X}_N]$$

$$\bar{X}_3^* = [\bar{X}_{Pob}^*, \bar{X}_{INB}^*, \bar{X}_{EV}^*,]$$

\sum_3 = Variances and covariances of full variables

$$\bar{X}_{Pob} = 38120$$

$$\bar{X}_{Pob}^* = 172062$$

$$d^2 = \sum_{j=1}^J n_j (\bar{X}_j - \bar{X}_j^*)^T \frac{1}{\sum_j} (\bar{X}_j - \bar{X}_j^*) \sim \chi^2_{\sum_j - K}$$

Type equation here.

		PobTotal	INBxcap	Esper Vida	%Nounat sBaizPes	Uslnstal Saneja	VIH	%coneizP resHomes	\$ConeizP resDones	Alfabet Homes	Alfabet Dones	Mortalitat	Natalitat	Economia	despesaSalut/Def ensa	Règim	
1	Països	I															
2	CAMBOYA	AsiaOr+Pacif	13798	320	57	11	16	170		64	85		11	31	enDesenvolupament	Monarquia	
3	COREA	AsiaOr+Pacif	22384		63	7	59						11	16	enDesenvolupament	Republica	
4	FIJI	AsiaOr+Pacif	841	2690	68	10	98	0,6		94	91	6	23	enDesenvolupament	MMesDefensa	Republica	
5	FILIPINES	AsiaOr+Pacif	81617	1170	71	20	73	3	58	44	93	5	25	enDesenvolupament	MesSalut	Republica	
6	ILLES COOK	AsiaOr+Pacif	18			3	100								Subdesenvolupat	Monarquia	
7	ILLES MARSHALL	AsiaOr+Pacif	60	2370		12	82								Subdesenvolupat	Republica	
8	ILLES SALOMÓ	AsiaOr+Pacif	466	550	63	13	31					7	33	Subdesenvolupat		Monarquia	
9	INDONÈSIA	AsiaOr+Pacif	220077	1140	67	9	52	110		23	92	83	7	21	enDesenvolupament	MesSalut	Republica
10	KIRIBATI	AsiaOr+Pacif	97	970		5	39								Subdesenvolupat	Republica	
11	LAO, REPÚBLICA DEMOCRÀTICA	AsiaOr+Pacif	5792	390	55	14	24	1,7		77	61	12	35	enDesenvolupament		Republica	
12	MALÀSIA	AsiaOr+Pacif	24894	4650	73	9		52		92	85	5	22	enDesenvolupament	MesSalut	Monarquia	
13	MICRONÈSIA	AsiaOr+Pacif	110	1990	68	18	28					6	31	Subdesenvolupat		Republica	
14	MONGÒLIA	AsiaOr+Pacif	2614	590	65	7	59	0,5		77	98	98	7	22	enDesenvolupament	MesSalut	Republica
15	MYANMAR	AsiaOr+Pacif	50004	220	61	15	73	330		94	86	10	20	enDesenvolupament	MesSalut	Dictadura	
16	NAURU	AsiaOr+Pacif	13												Subdesenvolupat	Republica	
17	NIUE	AsiaOr+Pacif	1			0	100								Subdesenvolupat	Monarquia	
18	PALAU	AsiaOr+Pacif	20	6870		9	83								Subdesenvolupat	Republica	
19	PAPUA NOVA GUINEA	AsiaOr+Pacif	5772	580	56	11	45	16		63	51	10	30	enDesenvolupament	MMesDefensa	Monarquia	
20	SAMOA	AsiaOr+Pacif	184	1860	71	4	100			99	98	6	28	Subdesenvolupat		Monarquia	
21	SINGAPUR	AsiaOr+Pacif	4273	24220	79	8		4,1		97	89	5	9	enDesenvolupament	MesSalut	Republica	
22	TAILÀNDIA	AsiaOr+Pacif	63694	2540	70	9	99	0,2		95	91	7	16	enDesenvolupament	MMesDefensa	Monarquia	
23	TIMOR-LESTE	AsiaOr+Pacif	887	550	56	12	33		6			12	50	Subdesenvolupat		Republica	
24	TOKELAU	AsiaOr+Pacif													Subdesenvolupat	Monarquia	
25	TONGA	AsiaOr+Pacif	102	1830	72	0	97			99	99	6	24	Subdesenvolupat		Monarquia	
26	TUVALU	AsiaOr+Pacif	10			5	88								Subdesenvolupat	Monarquia	
27	VANUATU	AsiaOr+Pacif	207	1340	69	6	50					6	31	enDesenvolupament		Republica	
28	VIETNAM	AsiaOr+Pacif	83123	550	71	9	41	220	60	94	87	6	20	enDesenvolupament		Republica	
29	XINA	AsiaOr+Pacif	1307989	1290	72	4	44	840		95	87	7	13	enDesenvolupament	MesSalut	Republica	
30	ANGOLA	AfricaOr+Merid	15490	1030	41	12	30	240		82	54	22	48	enDesenvolupament	MesSalut	Republica	
31	BOTSWANA	AfricaOr+Merid	1769	4340	35	10	41	350	90	93	76	82	27	26	enDesenvolupament	MesSalut	Republica
32	BURUNDI	AfricaOr+Merid	7282	90	44	16	36	250		47	67	52	19	45	enDesenvolupament	MesSalut	Republica
33	COMORAS	AfricaOr+Merid	777	530	64	25	23		41	63	49	7	36	enDesenvolupament		Republica	
34	ERITREA	AfricaOr+Merid	4232	180	54	21	9	60		62		11	39	enDesenvolupament		Republica	
35	ETIOPIA	AfricaOr+Merid	75600	110	48	15	6	1500		49	34	16	41	enDesenvolupament	MesSalut	Republica	
36	KENYA	AfricaOr+Merid	33467	460	48	10	48	1200	68	59	78	70	15	39	enDesenvolupament	MMesDefensa	Republica
37	LESOTHO	AfricaOr+Merid	1798	740	35	14	37	320		58	74	90	25	28	enDesenvolupament	MMesDefensa	Monarquia
38	MADAGASCAR	AfricaOr+Merid	18113	300	56	17	33	140	54	49	76	65	12	39	enDesenvolupament	MMesDefensa	Republica

The Little test in R

- LittleMCAR {BaylorEdPsych}
- **USAGE:** LittleMCAR(x)
x: dataframe, matrix less than 50 variables

- **Returns:**

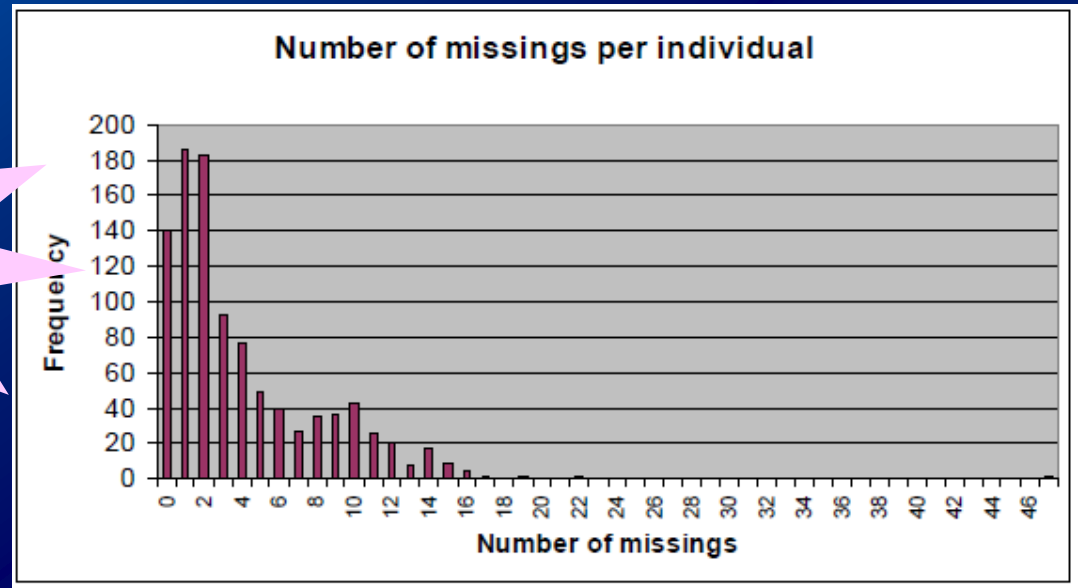
chi.square	Chi-square value
df	Degrees of freedom used for chi-square
missing.patterns	Number of missing data patterns
amount.missing	Amount and percent of mssing data
data	The data, organized my missing data patterns

Missing data

(Simple alternative)

- ▶ Build new variable counting number of missings per individuals.
- ▶ Describe this variable

Reliability of individual



- ▶ Count nr of missing per variable and rank variables
Provides reliability
- ▶ Create indicator of missing/non-missing per variable
and compare both groups of cases

Missing data

(empty cells in data matrix)

- Representation:

- *, ?, “ “, depending on software

- numerical variables: sometimes codified (0, 99999, -1...)

- categorical variables: special modality (Ns/Nc, ...)

- Standardize missing representation

- Causes of missing data:

- voluntary hidden (religion in israel) (always non-random)*

- data non-provided*

- data non-achieveable*

- technical limitations (example anemometers IKE hurrican)*

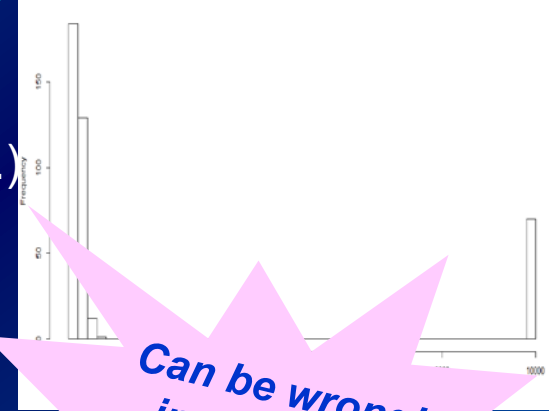
- accessibility (no privileges, sensitive information)*

- data lost*

- data forced to missing (as a result of correction)*

- Identification:

- Numerical indicators (stdev...)*



**Can be wrongly
incorporated
in analysis**

Missing data treatment

Depends on analysis goals!!!!

- ▶ keep it as a missing: only eventually
 - Can significantly reduce the treated observations
- ▶ Inputing: Substituting by a useful value (open problem, difficult)
 - Qualitative variable: Substitute by “Unkown<varName>”
 - Standard way, expert knowledge required
 - use 0
 - use global mean
 - use conditional mean for local groups
 - imputation models (complex)
 - Nearest neighbor (R)
 - Intelligent imputation
 - MIMMI
 - non parametric approach (montecarlo methods, multiple imputation)
 - special software required
 - technical hypothesis about variable distributions required
 - Final models integration required
 - ▶ *Example: French survey, global incomes of household*

For qualitative
variables keep as new
modality: Unknown

Missing data treatment

- ▶ Missing values frequent in real data
- ▶ Imputation before analysis CRITICAL
- ▶ Most statistical packages:
 - ▶ simple imputation by global mean
 - ▶ listwise deletion (dangerous)
- ▶ Specific softwares:
 - ▶ dedicated to sophisticated imputation methods
 - ▶ highly time consuming
 - ▶ non-exportable complete data matrices
- ▶ Find a trade-of between precision and simplicity

Knn method

C_HISTORI	C_TRACTA	DATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de rodes
1569,0	84585,0	09/07/2003 0:00	7	7	6	7	7	6	5	5	7
1642,0	74011,0	20/06/2002 0:00	7	7	7	7	7	7	7	7	7
1645,0	84990,0	21/07/2003 0:00	7		6	6	2	6	6	6	3
1666,0	91980,0	09/03/2004 0:00	7	7	7	7	7	6	6	5	7
1694,0	83561,0	03/06/2003 0:00	7	7	7	7	7	7	6	6	7
1754,0	114451,0	03/02/2006 0:00	7	7	6	7	6	6	6	6	7
1858,0	76281,0	26/09/2002 0:00	7		5	7	7	6	5	5	7
1900,0	84368,0	01/07/2003 0:00	6	6	4	4	3	1	6	4	7
1904,0	82443,0	30/04/2003 0:00	4	7	4	6	5	3	2	3	4
1919,0	74098,0	20/06/2002 0:00	7	7	7	7	7	7	6	6	4
1976,0	80110,0	13/02/2003 0:00	7	5	3	4	3	3	5	5	3
2052,0	81175,0	20/03/2003 0:00	7	7	6	7	6	6	6	6	7
2059,0	82951,0	15/05/2003 0:00	1	1	1	1	1	1	1	1	1
2251,0	76399,0	01/10/2002 0:00	5	5	1	1	1	1	6	5	1
2267,0	86796,0	01/10/2003 0:00	7		7	7	7	7	6	6	7
2524,0	76436,0	02/10/2002 0:00	7	7	6	7	6	6	6	6	7
2533,0	81445,0	28/03/2003 0:00	7	7	7	7	7	7	6	6	7
2604,0	75742,0	06/09/2002 0:00	7	7	6	7	7	7	5	6	7
2646,0	84112,0	20/06/2003 0:00	7	7	7	7	7	7	6	6	7
2685,0	79191,0	15/01/2003 0:00	7	7	7	7	7	7	6	6	7
2694,0	78901,0	02/01/2003 0:00	7	7	7	7	7	7	6	6	7
2726,0	74218,0	27/06/2002 0:00	6	6	4	6	6	5	3	5	6
2765,0	79837,0	05/02/2003 0:00	5	5	2	5	2	1	5	5	4

Original uncomplete data

Knn method

C_HISTORIC	TRACTA/IDATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de ro
1569,0	84585,0	09/07/2003 0:00	7	6	7	7	6	6	5	7
1642,0	74011,0	20/06/2002 0:00	7	7	7	7	7	7	7	7
1645,0	84990,0	21/07/2003 0:00	7	6	6	2	6	6	6	3
1666,0	91980,0	09/03/2004 0:00	7	7	7	7	6	6	5	7
1694,0	83561,0	03/06/2003 0:00	7	7	7	7	7	6	6	7
1754,0	114451,0	03/02/2006 0:00	7	6	7	6	6	6	6	7
1858,0	76281,0	26/09/2002 0:00	7	5	7	7	6	5	5	7
1900,0	84368,0	01/07/2003 0:00	6	4	4	3	1	6	4	7
1904,0	82443,0	30/04/2003 0:00	7	4	6	5	3	2	3	4
1919,0	74098,0	20/06/2002 0:00	7	7	7	7	7	6	6	4
1976,0	80110,0	13/02/2003 0:00	5	3	4	3	3	5	5	3
2052,0	81175,0	20/03/2003 0:00	7	6	7	6	6	6	6	7
2059,0	82951,0	15/05/2003 0:00	1	1	1	1	1	1	1	1
2251,0	76399,0	01/10/2002 0:00	5	1	1	1	1	6	5	1
2267,0	86796,0	01/10/2003 0:00	7	7	7	7	7	6	6	7
2524,0	76436,0	02/10/2002 0:00	7	6	7	6	6	6	6	7
2533,0	81445,0	28/03/2003 0:00	7	7	7	7	7	6	6	7
2604,0	75742,0	06/09/2002 0:00	7	6	7	7	7	5	6	7
2646,0	84112,0	20/06/2003 0:00	7	7	7	7	7	6	6	7
2685,0	79191,0	15/01/2003 0:00	7	7	7	7	7	6	6	7
2694,0	78901,0	02/01/2003 0:00	7	7	7	7	7	6	6	7
2726,0	74218,0	27/06/2002 0:00	6	4	6	6	5	3	5	6
2765,0	79837,0	05/02/2003 0:00	5	2	5	2	1	5	5	4

SPLIT

Cures Full

Cures Missing

C_HISTORIC	TRACTA/IDATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de ro
1569,0	84585,0	09/07/2003 0:00	7	6	7	7	6	6	5	7
1642,0	74011,0	20/06/2002 0:00	7	7	7	7	7	7	7	7
1666,0	91980,0	09/03/2004 0:00	7	7	7	7	6	6	5	7
1694,0	83561,0	03/06/2003 0:00	7	7	7	7	7	6	6	7
1754,0	114451,0	03/02/2006 0:00	7	6	7	6	6	6	6	7
1900,0	84368,0	01/07/2003 0:00	6	4	4	3	1	6	4	7
1904,0	82443,0	30/04/2003 0:00	7	4	6	5	3	2	3	4
1919,0	74098,0	20/06/2002 0:00	7	7	7	7	7	6	6	4
1976,0	80110,0	13/02/2003 0:00	5	3	4	3	3	5	5	3
2052,0	81175,0	20/03/2003 0:00	7	6	7	6	6	6	6	7
2059,0	82951,0	15/05/2003 0:00	1	1	1	1	1	1	1	1
2251,0	76399,0	01/10/2002 0:00	5	1	1	1	1	6	5	1
2524,0	76436,0	02/10/2002 0:00	7	6	7	6	6	6	6	7
2533,0	81445,0	28/03/2003 0:00	7	7	7	7	7	6	6	7
2604,0	75742,0	06/09/2002 0:00	7	6	7	7	7	5	6	7
2646,0	84112,0	20/06/2003 0:00	7	7	7	7	7	6	6	7
2685,0	79191,0	15/01/2003 0:00	7	7	7	7	7	6	6	7
2694,0	78901,0	02/01/2003 0:00	7	7	7	7	7	6	6	7
2726,0	74218,0	27/06/2002 0:00	6	4	6	6	5	3	5	6
2765,0	79837,0	05/02/2003 0:00	5	2	5	2	1	5	5	4

C_HISTORIC	TRACTA/IDATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de ro
1645,0	84990,0	21/07/2003 0:00	7	6	6	2	6	6	6	3
1858,0	76281,0	26/09/2002 0:00	7	5	7	7	6	5	5	7
2267,0	86796,0	01/10/2003 0:00	7	7	7	7	7	6	6	7

Knn method

C_HISTORI.C	TRACTAIDATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de rodes
1645.0	84990.0	21/07/2003 0:00	7	6	2	6	6	6	2	7
1858.0	76281.0	26/09/2002 0:00	7	7	7	6	5	5	7	7
2267.0	86796.0	01/10/2003 0:00	7	7	7	7	6	6	7	7

KNN

C_HISTORI.C	TRACTAIDATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de rodes
1569.0	84585.0	09/07/2003 0:00	7	6	7	7	6	5	5	7
1642.0	74011.0	20/06/2002 0:00	7	7	7	7	7	7	7	7
1666.0	91980.0	09/03/2004 0:00	7	7	7	6	6	5	7	7
1694.0	83561.0	03/06/2003 0:00	7	7	7	7	6	6	7	7
1754.0	114451.0	03/02/2006 0:00	7	6	7	6	6	6	7	7
1900.0	84368.0	01/07/2003 0:00	6	4	4	3	1	6	4	7
1904.0	82443.0	30/04/2003 0:00	4	4	6	5	3	2	3	4
1919.0	74098.0	20/06/2002 0:00	7	7	7	7	7	6	6	4
1976.0	80110.0	13/02/2003 0:00	7	3	4	3	3	5	5	3
2052.0	81175.0	20/03/2003 0:00	7	6	7	6	6	6	6	7
2059.0	82951.0	15/05/2003 0:00	1	1	1	1	1	1	1	1
2251.0	76399.0	01/10/2002 0:00	5	1	1	1	1	6	5	1
2524.0	76436.0	02/10/2002 0:00	7	6	7	6	6	6	6	7
2533.0	81445.0	28/03/2003 0:00	7	7	7	7	7	6	6	7
2604.0	75742.0	06/09/2002 0:00	7	6	7	7	7	5	6	7
2646.0	84112.0	20/06/2003 0:00	7	7	7	7	7	6	6	7
2685.0	79191.0	15/01/2003 0:00	7	7	7	7	7	6	6	7
2694.0	78901.0	02/01/2003 0:00	7	7	7	7	7	6	6	7
2726.0	74210.0	27/06/2002 0:00	6	4	6	6	5	3	5	6
2765.0	79837.0	05/02/2003 0:00	5	2	5	2	1	5	5	4

Knn method

C_HISTORIC	C_TRACTA	DATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de rodes
1645.0	84990.0	21/07/2003 0:00	7		6	2	6	6	6	2	
1858.0	76281.0	26/09/2002 0:00	7		7	7	6	5	5	7	
2267.0	86796.0	01/10/2003 0:00	7	7	7	7	7	6	6	7	

KNN

C_HISTORIC	C_TRACTA	DATA	Alimentació	Cures d'aparença	Higiene	Vestit: part superior	Vestit: part inferior	Utilització del bany	Bufeta	Intestí	Llit, cadira, cadira de rodes
1569.0	84585.0	09/07/2003 0:00	7		6	7	7	6	5	5	7
1642.0	74011.0	20/06/2002 0:00	7		7	7	7	7	7	7	7
1686.0	91980.0	09/03/2004 0:00	7		7	7	7	6	6	5	7
1694.0	83561.0	03/06/2003 0:00	7		7	7	7	7	6	6	7
1754.0	114451.0	03/02/2006 0:00	7		6	7	6	6	6	6	7
1900.0	84368.0	01/07/2003 0:00	6		4	4	3	1	6	4	7
1904.0	82443.0	30/04/2003 0:00	4		4	6	5	3	2	3	4
1919.0	74098.0	20/06/2002 0:00	7		7	7	7	7	6	6	4
1976.0	80110.0	13/02/2003 0:00	7		3	4	3	3	5	5	3
2052.0	81175.0	20/03/2003 0:00	7		6	7	6	6	6	6	7
2059.0	82951.0	15/05/2003 0:00	1		1	1	1	1	1	1	1
2251.0	76399.0	01/10/2002 0:00	5		1	1	1	1	6	5	1
2524.0	76436.0	02/10/2002 0:00	7		6	7	6	6	6	6	7
2533.0	81445.0	28/03/2003 0:00	7		7	7	7	7	6	6	7
2604.0	75742.0	06/09/2002 0:00	7		6	7	7	7	5	6	7
2646.0	84112.0	20/06/2003 0:00	7		7	7	7	7	6	6	7
2685.0	79191.0	15/01/2003 0:00	7		7	7	7	7	6	6	7
2694.0	78901.0	02/01/2003 0:00	7		7	7	7	7	6	6	7
2726.0	74210.0	27/06/2002 0:00	6		4	6	6	5	3	5	6
2765.0	79837.0	05/02/2003 0:00	5		2	5	2	1	5	5	4

MIMMI method *[Gibert 2013]*

- ▶ Select a small number of relevant variables
(with small ratio of missing data)
- ▶ Use intelligent imputation on that reduced data matrix
(expert-based imputation, vertical or horizontal)
- ▶ Multivariate clustering using the imputed variables
- ▶ Determine a partition of the data
- ▶ Impute the missing data of the remaining variables
(use mean local to the group of every individual (conditional means))
- ▶ Good TRADE-OFF between quality improvement vs extra effort

Example OMS

MIMMI *Method* [IJCM Gibert 2013]

Complex process
highly time consuming
rarely applicable in real projects

- ▶ Horizontal imputation:
use the value of other variables of the same individual as predictors of the missing value.

inputing 0 in the income of 4th person if the household has only 1,2 or 3 persons

- ▶ Vertical imputation:
use the value of the same variable in other similar individuals

use the mean of the salary of 4rt persons over 18 years old if the household has more than 4per

MICE method

[vanBuuren1999]

multiple imputation by chained equations

Multiple imputation (MI):

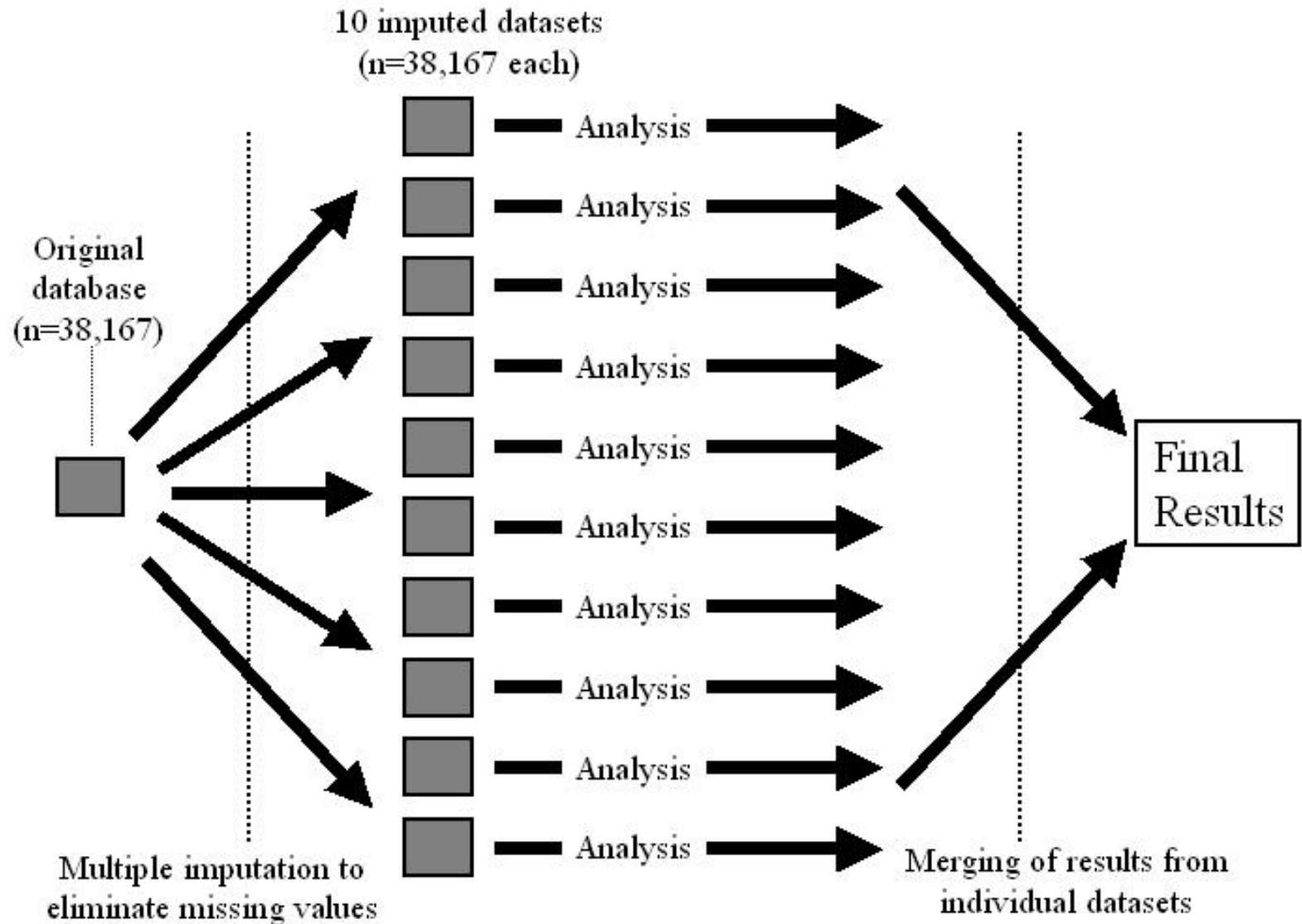
- Replace missing values with *plausible* substitutes
 - Distribution-based
 - maximum-likelihood based Markov-chain Monte Carlo (MCMC)
 - Inject *the right amount of randomness* to reflect uncertainty
- Repeat $m > 1$ times to produce m imputed datasets
- Analyse datasets individually, but identically
- Combine the models, get confidence intervals using Rubin's rules
(`micombine`)

The MICE approach has three components:

- Univariate – implemented in `uviz`
- Multivariate – implemented in `ice`
- Multiple – implemented in `ice`
- `ice` = imputation by chained equations

Context:
Multiple regression
In general

MICE



MICE

The overall *estimate of your parameter* (\bar{Q}) is its mean across the m imputations

$$\bar{Q} = m^{-1} \sum \hat{Q}^{(\ell)}$$

The *within-imputation variance* (\bar{U}) of the Q parameter is the mean of the variances across the m imputations

$$\bar{U} = m^{-1} \sum U^{(\ell)}$$

The *between-imputation variance* (B) of the Q parameter is standard deviation of Q across the m imputations

$$B = (m - 1)^{-1} \sum (\hat{Q}^{(\ell)} - \bar{Q})^2$$

The *total variance* of Q is a function of \bar{U} and B . This total variance is used to calculate the standard error used for test statistics

$$T = (1 + m^{-1})B + \bar{U}$$

$$(\bar{Q} - Q)/\sqrt{T} \sim t_\nu$$

The *degrees of freedom* (ν) are adjusted for the amount of information lost to missing data

$$\nu = (m - 1) \left[1 + \frac{\bar{U}}{(1 + m^{-1})B} \right]^2$$

MICE

- MICE method is very flexible – but demands thought when creating the imputation model
- Strongly recommend mastering the `eq()`, `passive()` and `substitute()` options
- Can deal with interactions using `passive()`
- Choice of m is important
 - may need to be (much) larger than 5
 - See Royston (2004, SJ 4:227-41) for discussion
- available in MICE Rpackage

Interpolation

- Usefull for time-series
- Consider constant between two effective measurements whenever we can assume that the variable evolve with a very slow dynamics

Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Outlier

- ▶ Rare observation (presumed out of range)
- ▶ Multivariate vs univariate outlier



Types of outliers:

- Mistake (Transcription Error or Measurement Error)
 - A person 560 years old
 - FIRST VERIFY *If possible correct.*
If not, substitute by missing
- Informative point
 - A single informative point of a missing part of the population
 - *Complete the sample*
when impossible, restrict scope of analysis
- Extreme value of the population
 - Very old person, 99 years old
 - *Keep*
- Value of another population
 - One swedish in the middle of cannibal tribu, measuring height
 - *Treat apart. CLEARLY REPORT ABOUT IT*
- Missing code
 - *Substitute by missing or inpute*



The danger of suppressions

- *In 1985 British scientists reported a hole in the ozone layer of the earth's atmosphere over the South Pole. This is disturbing, since ozone protects us from cancer-causing ultraviolet radiation. The British report was at first discredited, since it was based on ground instruments looking up. More comprehensive observations from satellite instruments looking down had shown nothing unusual. Then examination of the satellite data revealed that the South Pole ozone readings were so low that the computer software used to analyze the data had automatically suppressed these values as erroneous outliers. Readings dating back to 1979 were reanalyzed and showed a large and growing hole in the ozone layer that is unexplained and possibly dangerous. Computers analyzing large volumes of data are often programmed to suppress outliers as protection against errors in the data. As the example of the hole in the ozone layer illustrated, suppressing an outlier without investigating it can keep valuable information out of the sight*

Moore, McCabe, Introduction to the practice of Statistics, 5th Edition, Freeman

From the paper of John Gleick in New York Times, July 1985

<http://www.nytimes.com/1986/07/29/science/hole-in-ozone-over-south-pole-worries-scientists.html?pagewanted=all>

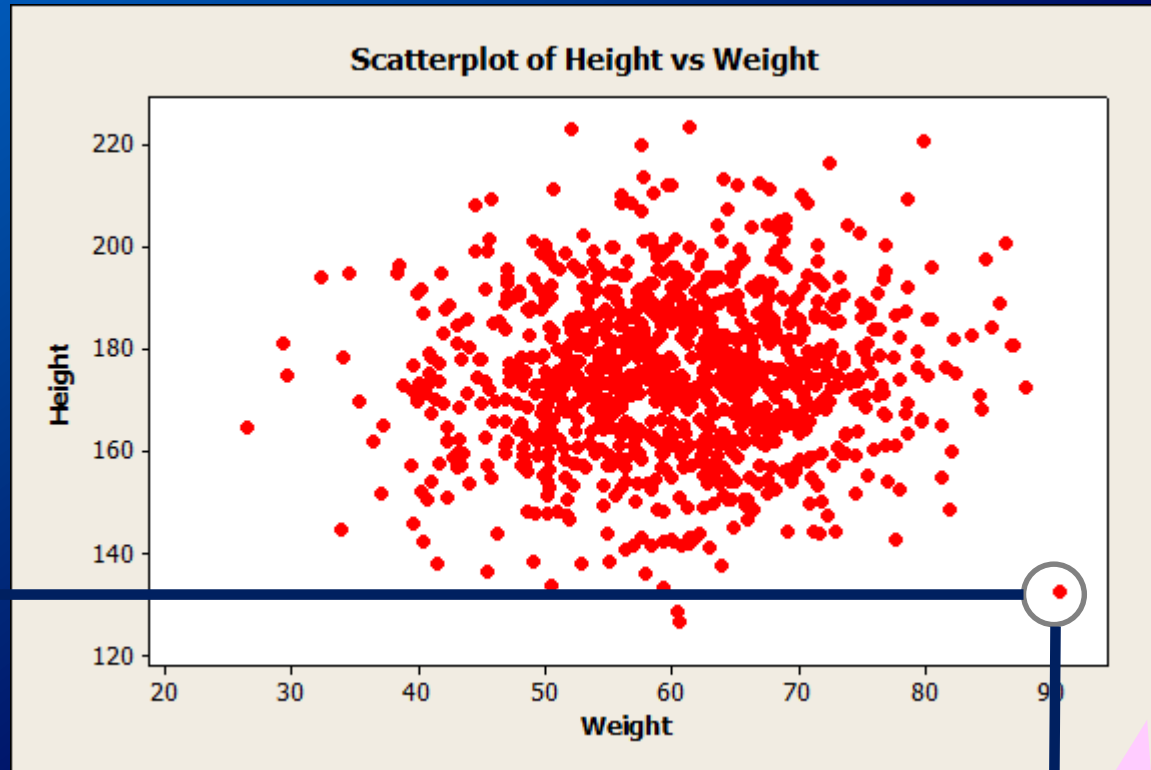
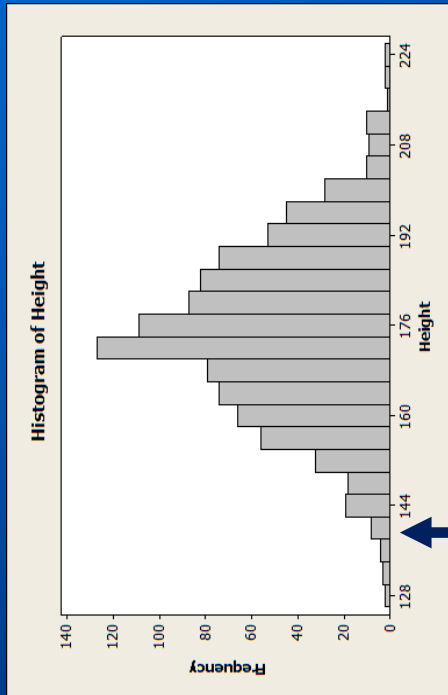
Outlier detection

- ▶ Specific statistical Tests:
 - Depend on the software
 - Usually for specific distributions
- ▶ Graphical representation of the distribution of data
 - Univariate (histogram or boxplot)
 - Bivariate (plots)
 - Clustering for multivariate outliers (singletons)

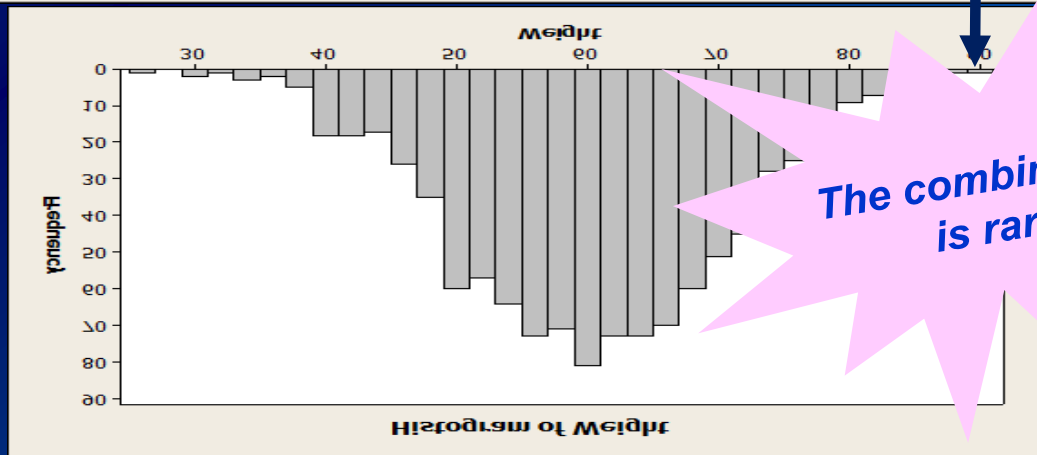


*Care with
boxplot*

Dimensionality of outliers: Bivariate Outlier

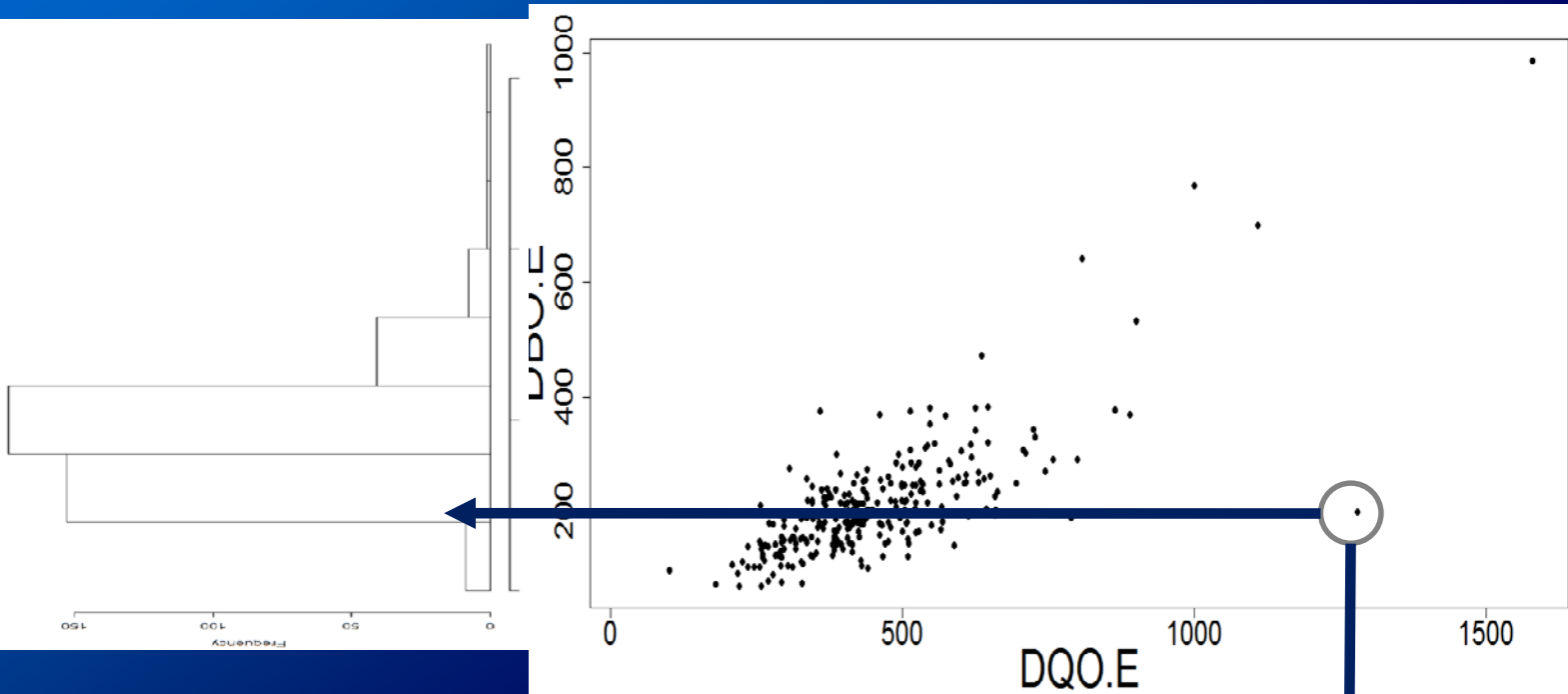


*A person with
90Kg and 1,32 m*

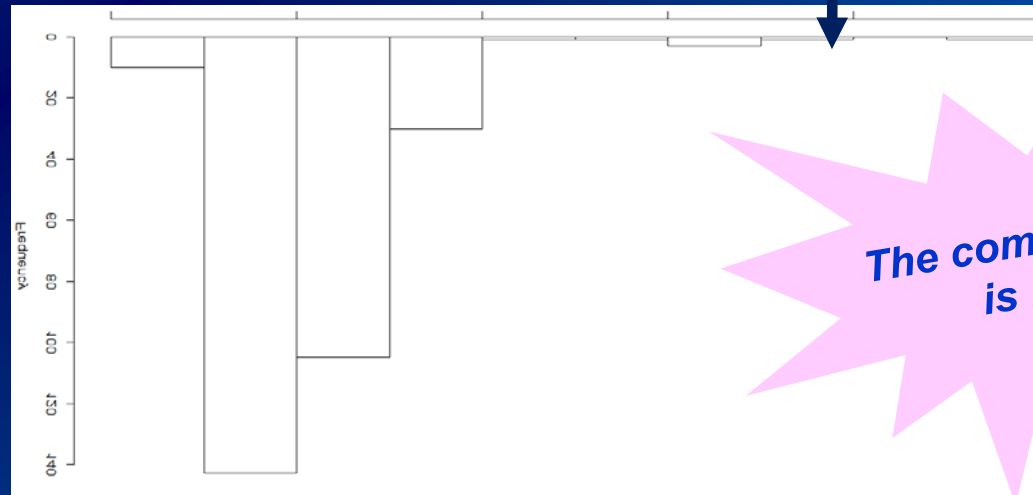


**The combination
is rare**

Dimensionality of outliers: Bivariate Outlier

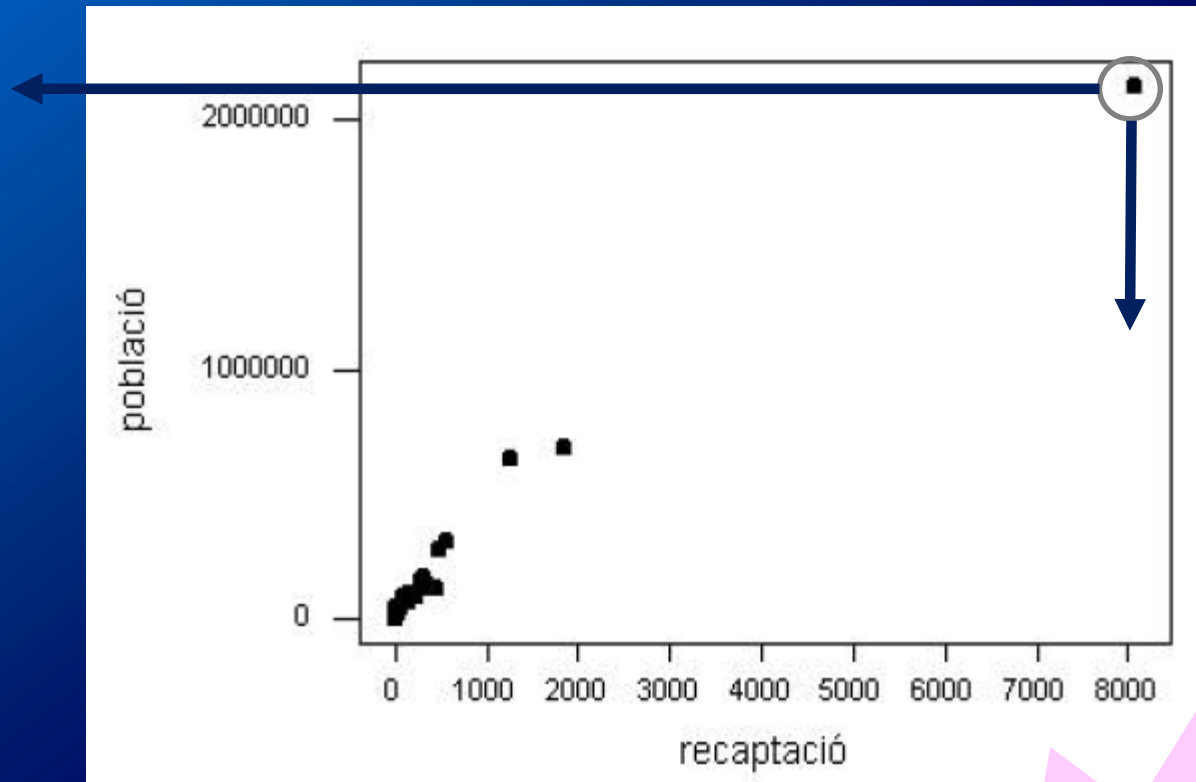


*A day with
DQO.E = 1279 mg/l
DBO.E = 198 mg/l*



**The combination
is rare**

Dimensionality of outliers: Univariate Outlier

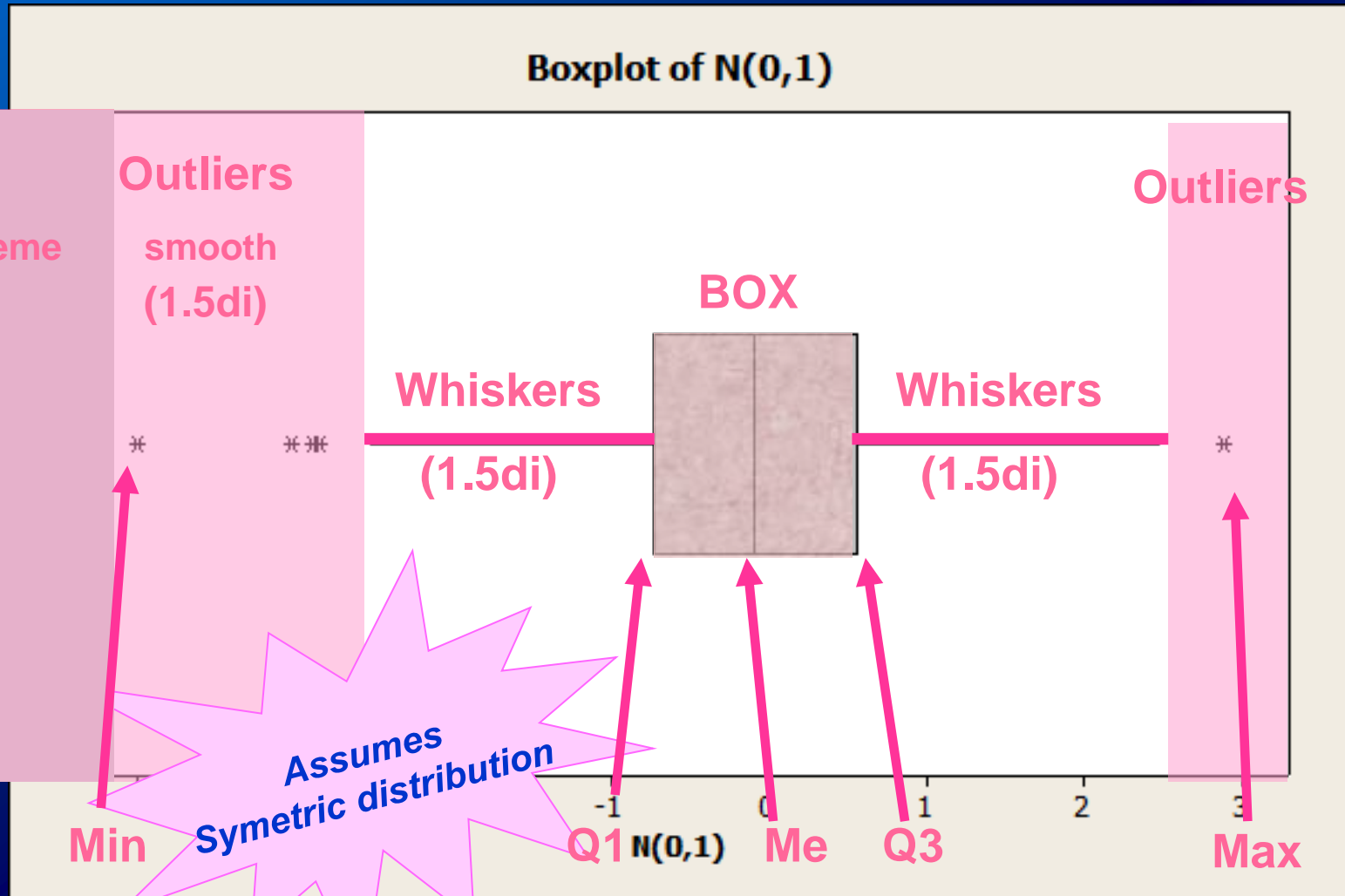


*Aligned with
global model!!!!*

Boxplot

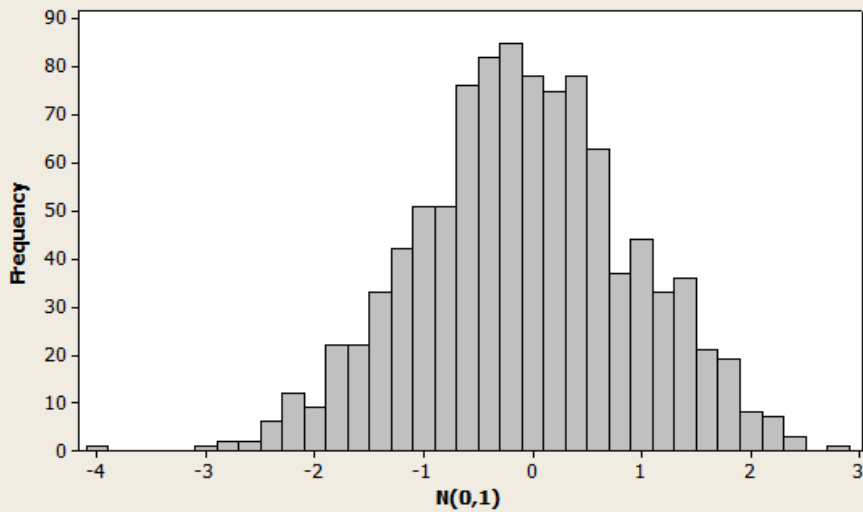
[Tukey 1956]

Symbolic representation of empirical distribution

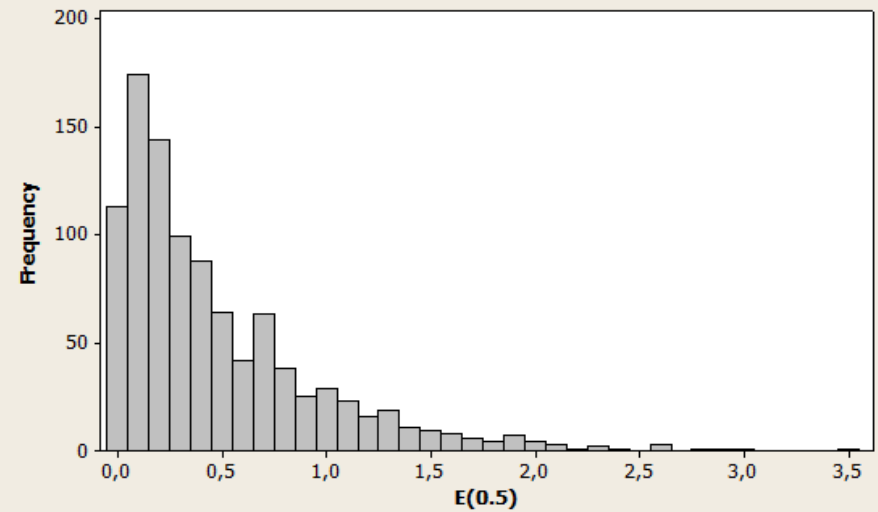


Boxplot [Tukey 1956]

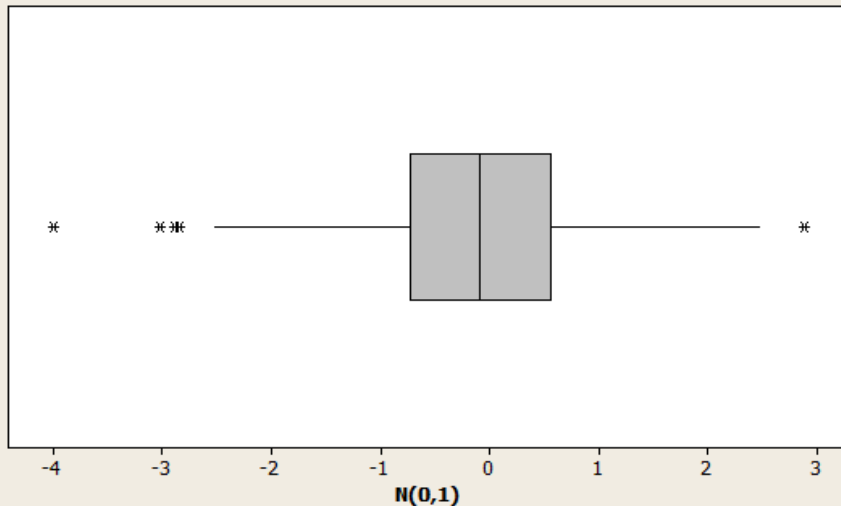
Histogram of $N(0,1)$



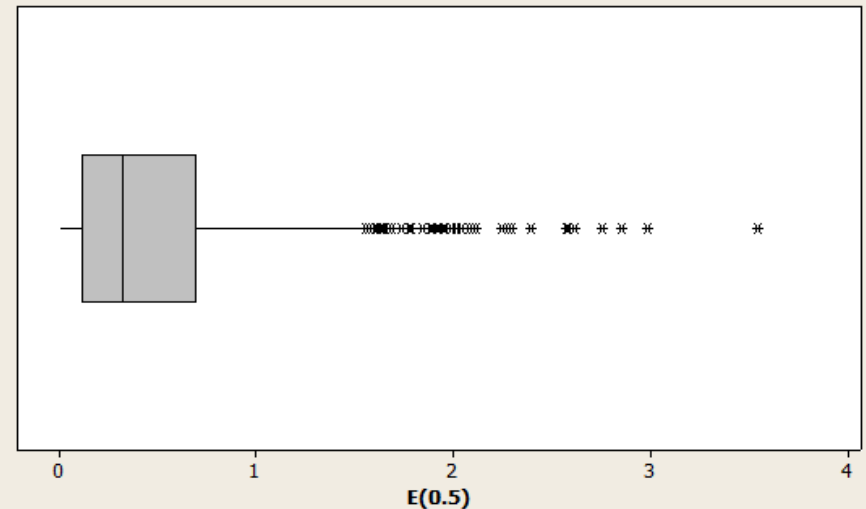
Histogram of $E(0.5)$



Boxplot of $N(0,1)$



Boxplot of $E(0.5)$



Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Instance selection

Evaluation of representative instances in a dataset

- ▶ Elimination of irrelevant instances
- ▶ Sampling
- ▶ Resampling

Reparing unbalanced datasets when required

- ▶ oversampling
- ▶ undersampling

Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Feature selection

Evaluation of relevant variables in a dataset

- ▶ Priorization and ranking under different criteria
 - ▶ Feature weighting (determine weights of variables in the analysis)
- ▶ Elimination of irrelevant variables
 - ▶ Feature selection

Feature selection

- ▶ IA methods
- ▶ Statistical Feature selection: use statistical test for ranking
- ▶ Sometimes just use threshold on feature weighting ranks

Feature selection

- ▶ Goal: discard non-interesting variables
 - ▶ Reduce data dimensionality
 - ▶ Eliminate noise and redundancies
 - ▶ Improve performance of algorithms
 - ▶ Avoid spurious relationships in models
 - ▶ Reduce curse of dimensionality
 - ▶ Requires a response variable to be explained Y
-
- ▶ Rank relevance degree of Y wrt all other variables
 - ▶ Discard less relevant

Statistical Feature selection

Guyon, I. (2008). Practical feature selection: from correlation to causality. NATO science for peace and security, 19, 27-43.

Hypothesis test:

H_0 : There is no relation between the y and x

H_1 : There is a relation

Get p-values for the dependence between Y and X

Lower p-values imply strongest dependence

Rank variables by ascending p-values

Discard irrelevant variables (threshold over p-values)

Specific tests depends on type of variables analyzed

Statistical Feature selection

Hypothesis test:

Y numerical

- ▶ X numerical: Correlations test / Sheffer generalized coefficient
- ▶ X qualitative: F test /Kruskal-Wallis

Y qualitative

- ▶ X numerical: F test/Kruskal-Wallis
- ▶ X qualitative: chi-2 test

**Care with
assumptions**

Feature selection

Evaluation of relevant variables in a dataset

- ▶ Priorization and ranking under different criteria
 - ▶ Feature weighting (determine weights of variables in the analysis)
- ▶ Elimination of irrelevant variables
 - ▶ Feature selection

Feature selection

- ▶ IA methods (based on information theory)
- ▶ Statistical methods (based on statistical tests)
- ▶ Sometimes just use threshold on feature weighting ranks

Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Variables Transformation

- ▶ Homogeneization
- ▶ Approaching to methods hypothesis
- ▶ Getting more interpretability

Variables Transformation

- ▶ Data cleaning reasons

- ▶ Measurement units of Thyroids hormones from different laboratories

1993

Collaboration UPC, Barcelona, Spain

Andrija Stampar School of Public Health, Zagreb, Croatia

Setre Milordsnice Clinical Hospital, Zagreb, Croatia

Find patterns of thyroids dysfunctions 1002 patients, 12 measurements

2013

Collaboration UPC, Atención Primaria ICS

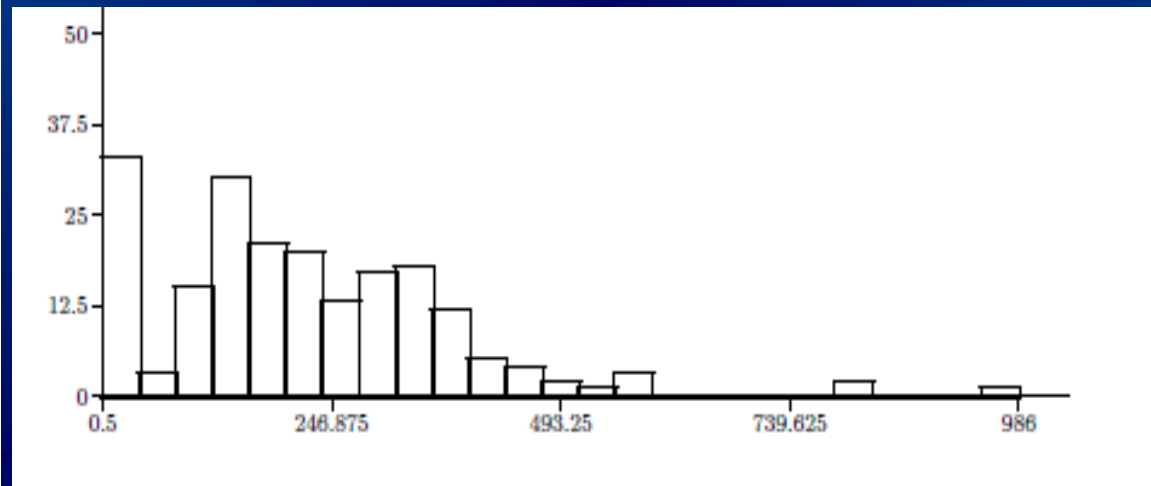
<http://www.sidiap.org/>

Laboratory measuments in TSH

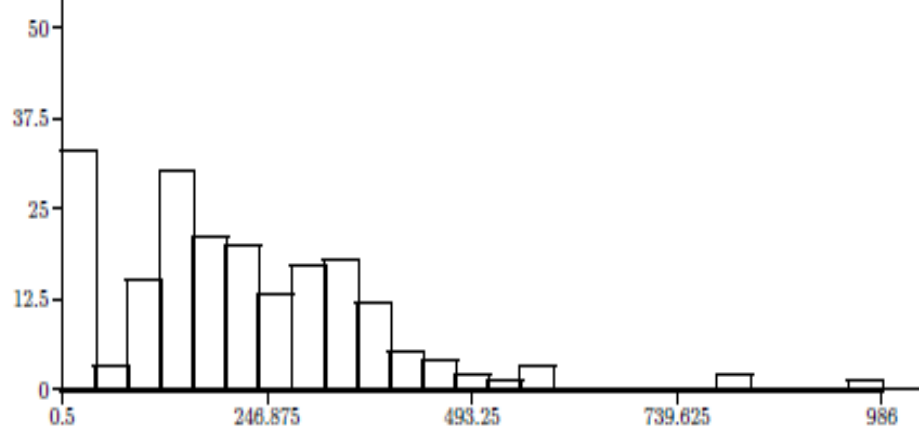
Laboratory Tests measurements

Measurement of Total Cholesterol from 200 pacs from Catalan Public Health System in 2013 (Primary Care)

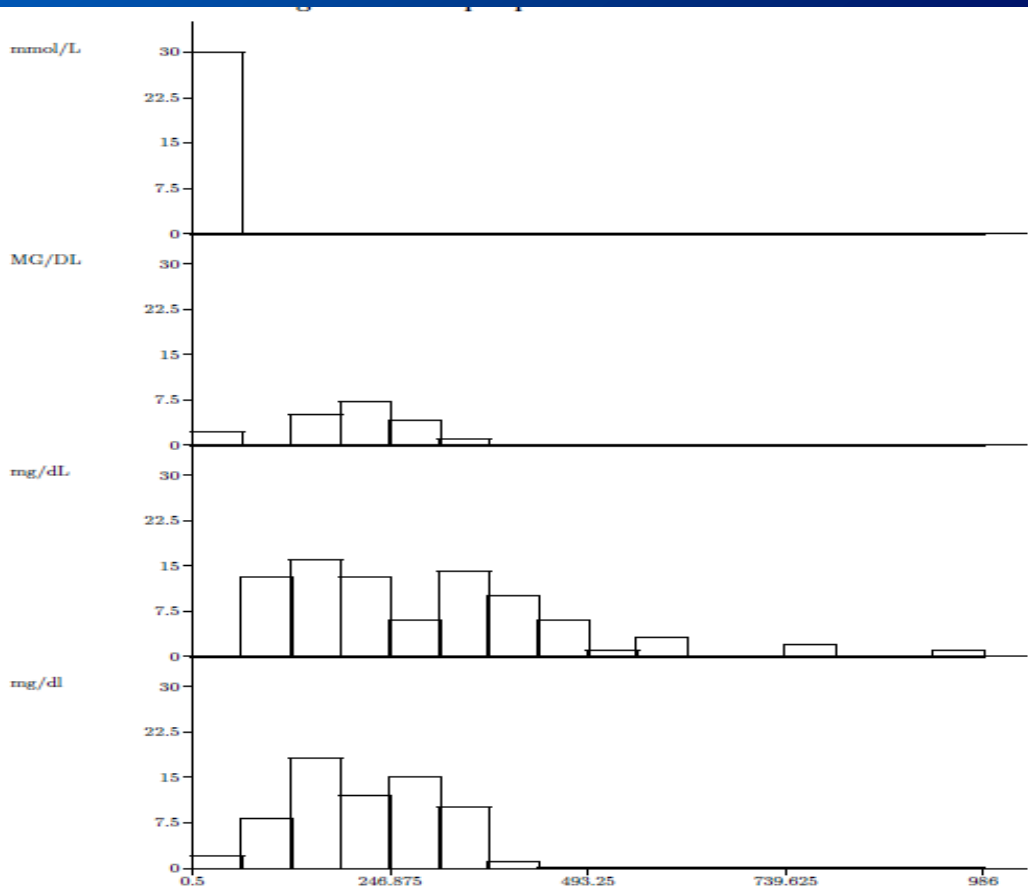
Summary Statistics	
Number of objects	200
Number of missing values	0
Number of useful values	200
Mean	213.2037
Median	193.9
First Quartile (Q1)	115.8
Third Quartile (Q3)	306.2
Minimum	0.5
Maximum	986
Quasi-standard deviation	156.4218
Variation Coefficient	0.7318



Frequency Table		
Modalities	Freq. absol.	Freq. relat.
mg/dl	66	0.33
mg/dL	85	0.425
MG/DL	19	0.095
mmol/L	30	0.15
<i>missing data</i>	0	0



Unitat	Count
(mg/dL)	12
0.0 - 240.0	1
55519	1
G/L	4
mg/dl	50
mg/dL	50
MG/DL	50
mg/dl^MMOL/L.	50
mmol/L	50
Mmol/L	1
MMOL/I	1
MMOL/L	1
mmol/L^mg/dL	50
N=	321



$\text{mmol/l} = 38,669 \text{ mg/dl}$

Variables Transformation

► Data cleaning reasons

- Measurement units of Thyroids hormones from different laboratories
- Refer the whole set of variables to comparable units
all concentration variables in mg/l
proportions instead of absolute numbers,
- Coertions: Information loss.
 - Discretization (h/week working)
 - Categorization (Thiroids levels)
 - Recategorizations (professions)

Better avoid

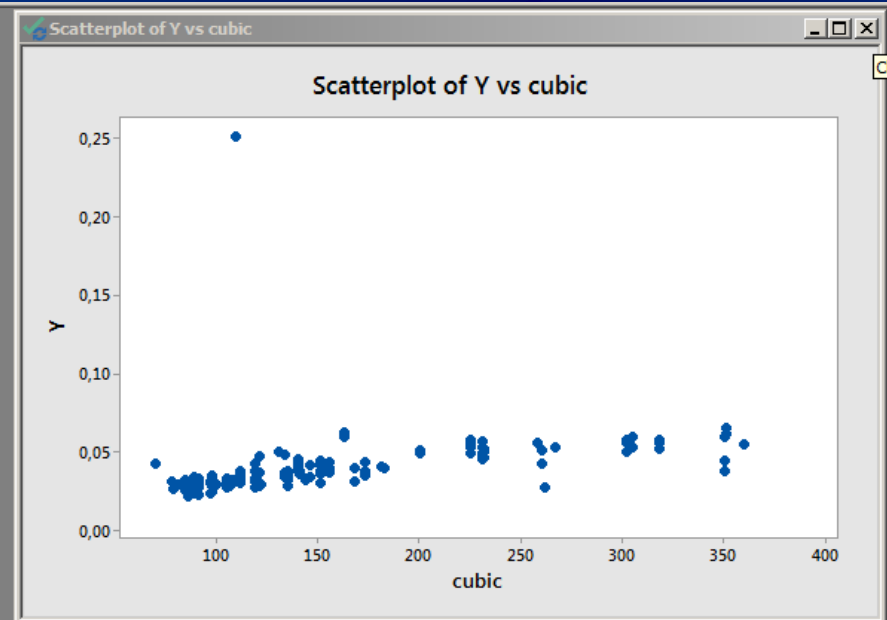
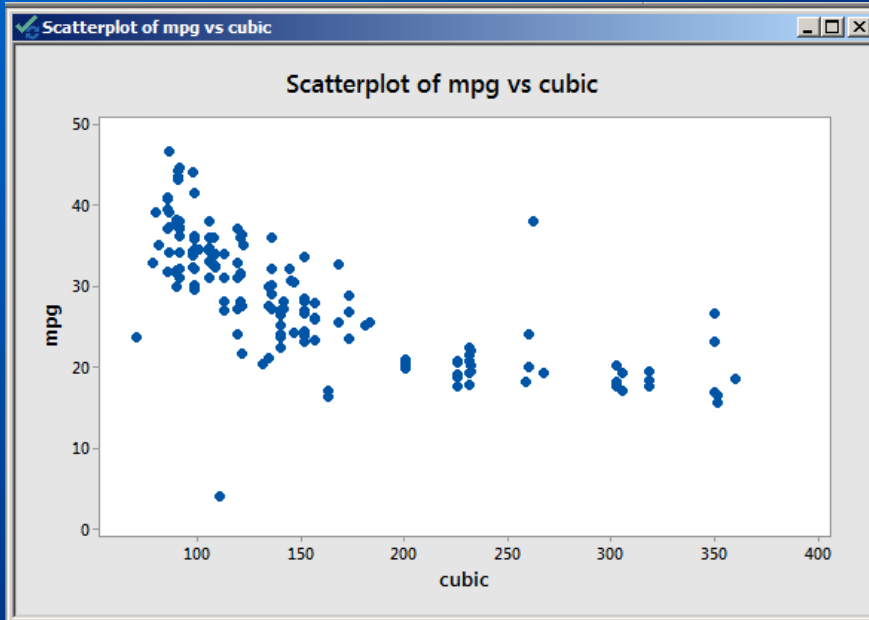
► Technical questions:

- Estandarditization, normalitization o linealirization
- Eventual logarithmic transformation
- Required by data mining technique to apply

**Select a technique
respectfull with original data**

Exceptional situations

where transforms make sense



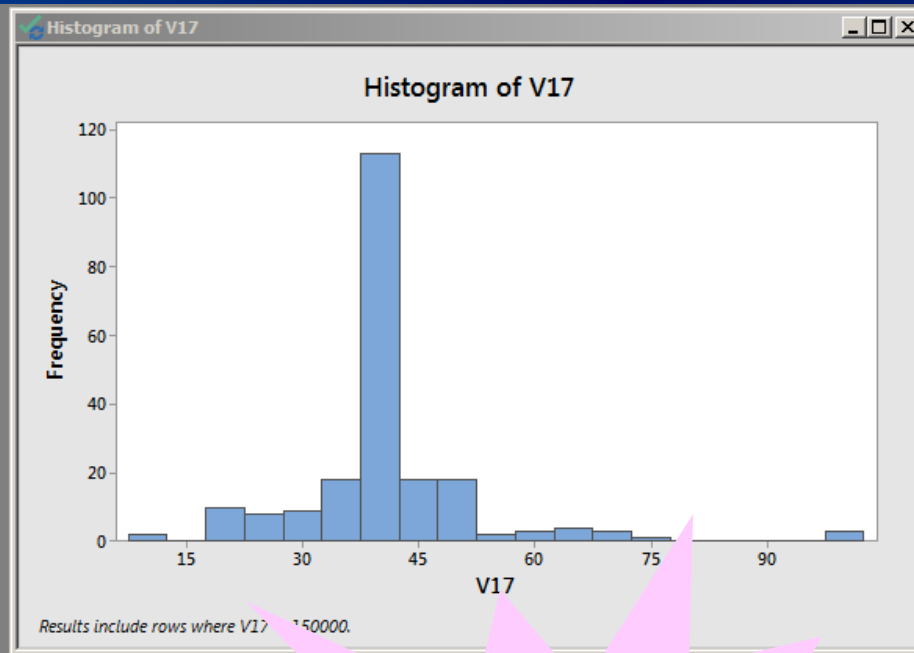
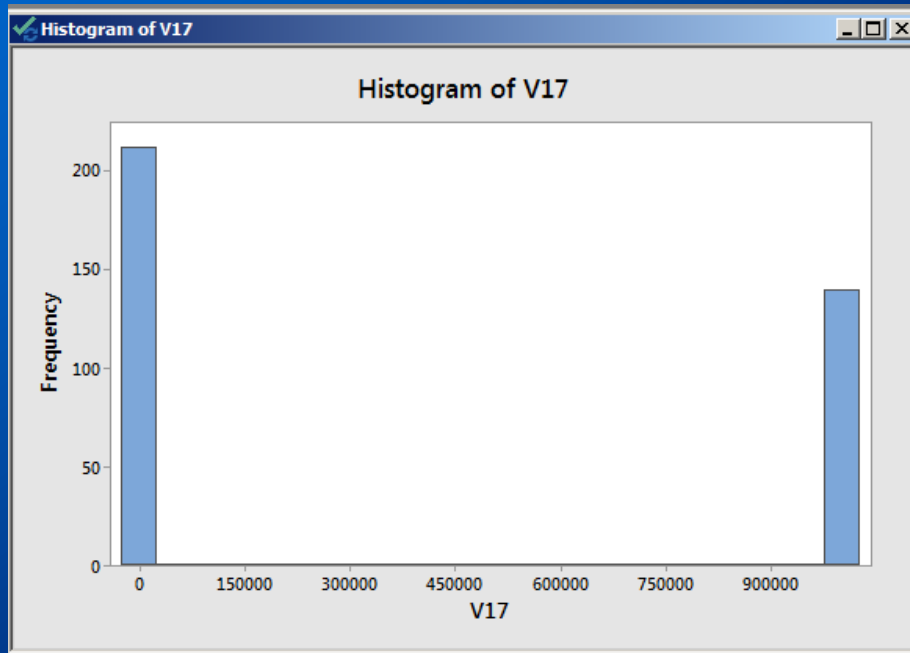
- ▶ Mpg: miles per gallon of a car
- ▶ Cubic: cubic capacity of the car engine
- Non linear relationship (regression non suitable)

- ▶ $Y = 1/\text{mpg}$: Linearizes the relationship

**Y is car
Consumption!!!!**

Exceptional situations

where transforms make sense



- ▶ Hours working per week
- ▶ 3-modal:
 - ▶ Around 20 h/w
 - ▶ Around 40 h/w
 - ▶ Around 65 h/w

- ▶ Correspondence with part-time, full-time, extra turns works

**Build a qualitative
variable:
Type of work
(part-time, full, turn)**

Preprocessing

Data cleaning

Data preparation

Data preprocessing

- ▶ Formatting issues, building software context
- ▶ Determining working matrix, Filtering
- ▶ Identification and treatment of missing data
- ▶ Identification and treatment of outliers
- ▶ Identification and treatment of errors (*correct when possible*)
- ▶ Feature selection/extraction, dimensionality reduction
- ▶ Instance selection
- ▶ Data transformation
- ▶ Derivation of new variables

Derivation of new variables

- ▶ Aggregates (additions of other variables)
 - ▶ Total household income
- ▶ Synthetic indicators
 - ▶ Classical generation of global score in psychometric scales
 - ▶ Indicators

*(Lund parameter = external contacts/days hospital
indicator of “development of a health system”)*

Case Credit Scoring (saving capacity)

**Input missings
According to operation
Previously**

- ▶ Binary indicators
 - ▶ *If condition regarding a combination of values
then indicator=1, else the indicator=0*
- ▶ *Dimensionality reduction techniques*