# Data Quality - Edits



IDA 2007 - LJUBLJANA - SLOVENIJA
6-8 September

Hans J. Lenz, Freie Universität Berlin

September 2007

# Edits

**Objective**:
Data Cleansing at the Data Entry
to assert semantic consistency of data with
nominal, ordinal and metric scales.

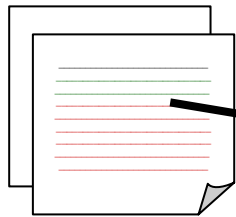„Numbers don't mind where they come from" *(LORD)*

# Agenda

1. Examples
2. Definitions
3. Simple Edits
4. Logical Edits
5. Numerical Edits
6. Statistical Edits
7. Fuzzy Edits
8. Evaluation of Statistical Edits vs. Fuzzy Edits
9. MCMC Simulation

# 1. Examples
## Ex. 1: Triangle Data

x = (2,4,3) consistent?

Frame of discernment:
x is a list of lengths from a triangle with a
right angle at B, and AB=2, AC=4, BC=3.

A

B          C

+

Pythagoras

=

Constraint
$BC^2 + AB^2 = AC^2$
$4 + 9 \neq 16$

# Ex. 2: Relational Data

- Is the tupel

x = (0010, 015, ´elementary´, ´child´, ´single´)

consistent?

- DB Schema as frame of discernment:

questionnaire (<u>oid</u>, age, school_type, household_status, marital_status)
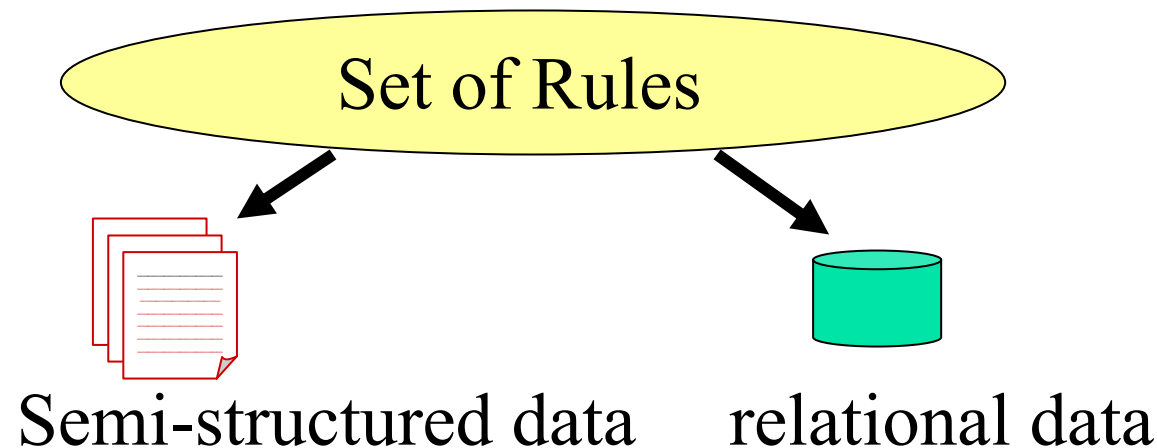
# 2. Definitions
## Edits (Rules, Checks)
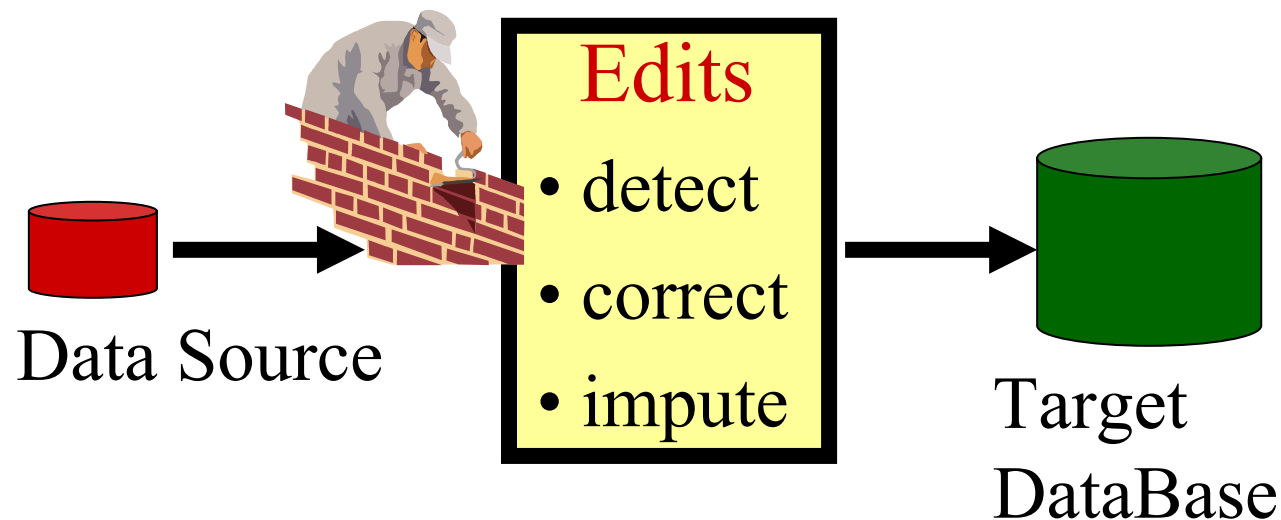
DEF.:

Edits are (normalized) rules (or formulas)
 - applicable to each object of a data set and
 - generated by relationships which exist
   between attributes.



Set of Rules

Semi-structured data    relational data

# Editing Tasks

Given a database **D** on a universe **U**

- *detect* semantic inconsistencies
- *correct* for incoherent data
- *impute* data in case of missing values.



Data Source

**Edits**
- detect
- correct
- impute

Target DataBase

# Objects at Data Entry

- **fix-formatted Data**
  - **Relations**
  - **Sets**
  - **Lists**
  - **Files**
  - **Records**
  - **Fields / Attributes / Variables**
- **semi-structured Data**
  - **HTML**
  - **XML**

# Types of Edits

- simple edits

- logical edits

- numerical edits

- probabilistic edits

- statistical edits
  ( based on Probability Theory )

- fuzzy edits ( based on Fuzzy Logic )

# 3. Simple Edits

☞ Conceptually the simplest edits are those applied to single field or attribute with respect to

- Data type

- Length

- Subset Constraints

- Scale

- Dimension

# 3. Simple Edits
## Examples

Syntax:

<attribute name> <u>predicate</u> <argument>

- Age <u>type</u> `cardinal`

- `code` <u>length</u> `4`

- `date` <u>between</u> `(01.07.06-06.07.06)`

- `Consumption rate` <u>scale</u> `metric`

- `costs` <u>unit</u> `€/year`

# 4. Logical Edits
## (cf. Categorical rules)

Example 1:

$for$ $all$ $u \in \mathbf{U}$

$if$ le(Age,15)

   $or$ school(elementary)

$then$ $not$ household_status(head)

   $and$ marital_status(single)

Generalization:

$if$ $x_1$ $is$ $A_1$ $and$ $x_2$ $is$ $A_2$ $and$...$x_p$ $is$ $A_p$

$then$ y $is$ B

☞ Fellegi and Holt (1973, 1976),
   Mamdani and Assilian (1975)
   Winkler (2004)

# 4. Logical Edits

**Theorem „Normal Form Edit"**
(Fellegi and Holt (1973)

Let f, g,h* $\in$ **C** be compound clauses and A,B,C,..., P,Q,R,... $\in$ $\mathscr{P}$ predicates of 2order-type. Then

$$f(A,B,C,...) => g(P,Q,R,...)$$

$$<=> \bigwedge_s \mathbf{h}^*_s = .false.$$

where s is a subset of the set **S** of attributes

☞ Fellegi and Holt (1973, 1976),
Mamdani and Assilian (1975), Boskovits (2007)

# 4. Normal Forms of Edits
## Examples

- <u>le</u>(Age, 15)∧ household _status(head)=.false.

- <u>le</u>(Age, 15)∧ <u>not</u> marital_status(single)=.false.

- school(elementary) ∧
  household _status(head)=.false.

- school(elementary) ∧
  <u>not</u> marital_status(single)=.false.

# 4. Logical Edits - Algorithms

Algorithms exist for:

- normalising edits into NF

- deciding whether an edit is new

- construct the complete set of essentially different edits

- identifying attributes to be most likely in error.

Sources: Fellegi, Holt (1976), Greenberg and Surdi (1984), Wetherill and Gerson (1986)

# 5. Numerical Edits

- *Imprecise Values* of attributes not allowed

- Defined only for data types *integer, cardinal, real, decimal*

- "Constraint Programming" by LP:

$$A \, x \geq b \qquad \text{(numerical constraints)}$$
$$x \geq 0 \qquad \text{(non negativity)}$$
$$x \in \mathbf{X} \qquad \text{(for all attributes x)}$$

# 5. Numerical Edits
## Example

- Fact: The former Student s
  is now 29 years of age ($x_1$),
  stayed 6 years at elementary school ($x_2$),
  stayed 7 years at high-school ($x_3$),
  studied 5 years at an university ($x_4$), and
  is employed since 2 years ($x_5$).

- *LP: x solves $a'x \geq b$ with*
  $a' = (1,-1,-1,-1,-1)$
  $x' = (x_1,x_2,x_3,x_4,x_5)=(29,6,7,5,2)$
  $b = 6$
  $\mathbf{X} = \prod \text{range}(x_i)$

# 5. Numerical Edits - Objectives

- Detect Redundancy

- Detect Inconsistency

- Error Location (misprint, transcription error, misspelling etc.)

- model-based Imputation

Algorithms for detection & location by:
**Sadiq (1986), ...**

18

# 6. Probabilistic Edits

- Let f $\in$ **C** be a clause over a set of predicates A,B,C,...

- Instead of a numerical edit:
  $$f(A,B,C,...)=.false.$$

- Now a probabilistic edit:

  $$f(A,B,C,...)=.false. \text{ with Prob}\geq 1-\alpha$$

# 6. Probabilistic Edits - Example

- surprise (=low probability) for any $person_A$, $person_B$:

  *if* sex($person_A$ male) *and*
  
  sex($person_B$, female) *and*
  
  married($person_A$,$person_B$) *and*
  
  le(*diff*(age($person_A$),age($person_B$)),-10).

- surprise *for* company x
  
  *if* growth(profit, 30%) *and*
  
  loss(sales, 20%) *and*
  
  *equal*(year(profit), year(sales)).

# 7. Statistical Edits Example

- Data ( six variables ):

Capital = 60 ± 1          Profit = 10 ± 2
Sales = 55 ± 20           Expenditures = 45 ± 20
ROI% = 10 ± 5            Margin% = 50 ± 5

- Constraints ( four equations based on definitions):

Profit = Sales - Expenditures
ROI% = 100 * Profit / Capital
Margin% = 100 * Profit / Sales
Turnover% = 100 * Sales / Capital

*Note*: Margin% = 50% > 100*Profit / Sales=18% !

# 7. Statistical Edits

Syntax:

$$\text{<variable>} = \text{<value>} \pm \quad \text{<abs. error>}/$$
$$\text{< \% error>}/$$
$$\text{<stdv>}$$

Example:

Profit $= 10 \pm 2$  with confidence level 1-$\alpha$

Sources: Schmid (1979), Lenz and Rödel (1991)

# 7. Statistical Edits   Modelling

- **Model Specification**
  - random variables

- Ranges
  - confidence intervals using Gaussian distribution

- **Parameter Estimation / Learning**
  ( ☞ only one value per variable measured! )

- **Inference** (Detection, Correction, Imputation)

# 7. Statistical Edits - Model

- **State Space Model:**

x       observed state vector

$\xi$       unobservable (error-free) state vector

$\zeta$       dependent vector

$\zeta = H \xi$  balance equation system

z       observed vector of $\zeta$

v, w  vectors of measurement errors

*state space equation*: $x = \xi + v$

*observational equation*: $z = H \xi + w$

# 7. Statistical Edits

- **State Space Model:**

u' = (v,w) with E(u) = 0 , i.e. no bias!

E(v,w) = 0 due to partial information
from independent data sources!

$$E(u\,u') = \mathbf{Q} = \begin{pmatrix} \mathbf{P} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix}$$

state space equation: $x = \xi + v$

observational equation: $z = H\,\xi + w$

Source: Schmid (1978), Schneeweiss (1989), Lenz and Rödel (1991, 2000),
Rödel (1999)

# 7. Statistical Edits - Problem

- *Two Problem Classes:*

1. Imputation:

   Given H, P, R and x
   estimate $\xi$, $\zeta = H\xi$ and predict z.

 2. Correction:

   Given H, P, R and x, z
   estimate $\xi$, $\zeta = H\xi$

*state space equation*: $x = \xi + v$

*linear observational equation*: $z = H\xi + w$

# 7. Statistical Edits - Estimation

- General Least Squares – Estimator of $\xi$:

Let $y' = (x,z)$ and $J = (\mathbf{I}, H)$.

$$\min \| y - J\xi \|_{\mathbf{Q}^{-1}}$$

$$\text{s.t. } \zeta = H\xi$$

☞ $\zeta = H\xi$ may be generalized to $\zeta = H(\xi)$

# 7. Statistical Edits - Estimators

1. Imputation:

Given H, P, R and x
estimate $\zeta = H\xi$ and predict z:
$$\hat{\zeta} = Hx \text{ and } \hat{z} = Hx$$

2. Correction:

Given H, P, R and x, z
estimate $\xi$, $\zeta = H\xi$ :

$$\hat{\xi} = x + K(z - Hx) \text{ and } \hat{\zeta} = H\hat{\xi}$$

$$K = PH'(HPH' + R)^{-1}$$

Note: GLS estimators called Kalman Filter Equations and Kalman Gain

# 7 Statistical Edits - Example

- Data: $x_1 = 30 \pm 20$

$$x_2 = 30 \pm 10$$
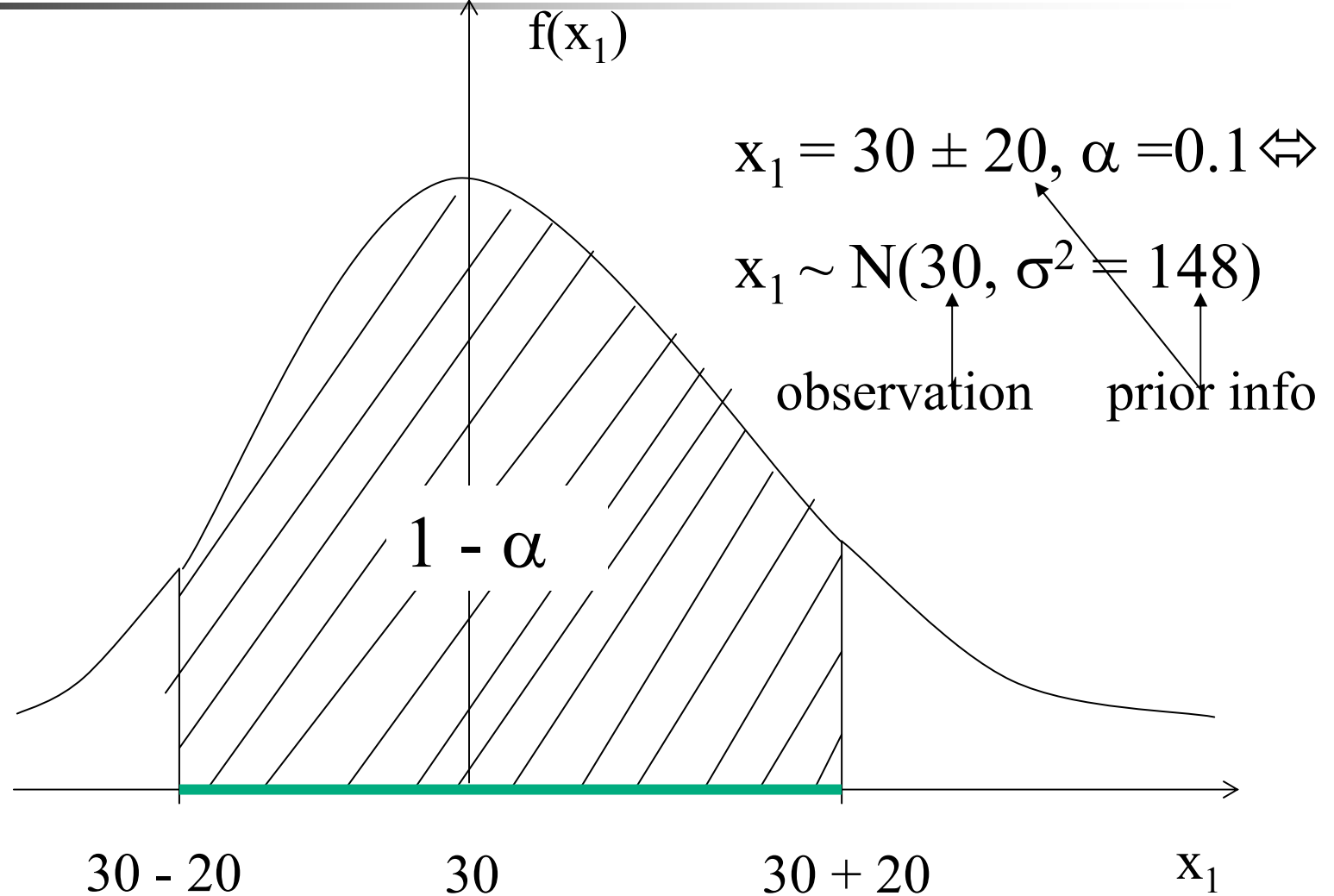
$$z = 50 \pm 10$$

- Model: $\zeta = \xi_1 + \xi_2 = H\,\xi$

Error Distributions ($\alpha = 0,1$):

$$E(vv') = P = \begin{pmatrix} 148 & 0 \\ 0 & 37 \end{pmatrix}$$

$$E(ww') = R^2 = 37$$

# 7 Statistical Edits Transformation
## errors into confidence intervals



$f(x_1)$

$x_1 = 30 \pm 20, \; \alpha = 0.1 \Leftrightarrow$

$x_1 \sim N(30, \sigma^2 = 148)$

observation          prior info

$1 - \alpha$

30 - 20          30          30 + 20          $x_1$

# 7. Statistical Edits - Example

Data: $x_1 = 30 \pm 20$

$x_2 = 30 \pm 10$

$z = 50 \pm 10$

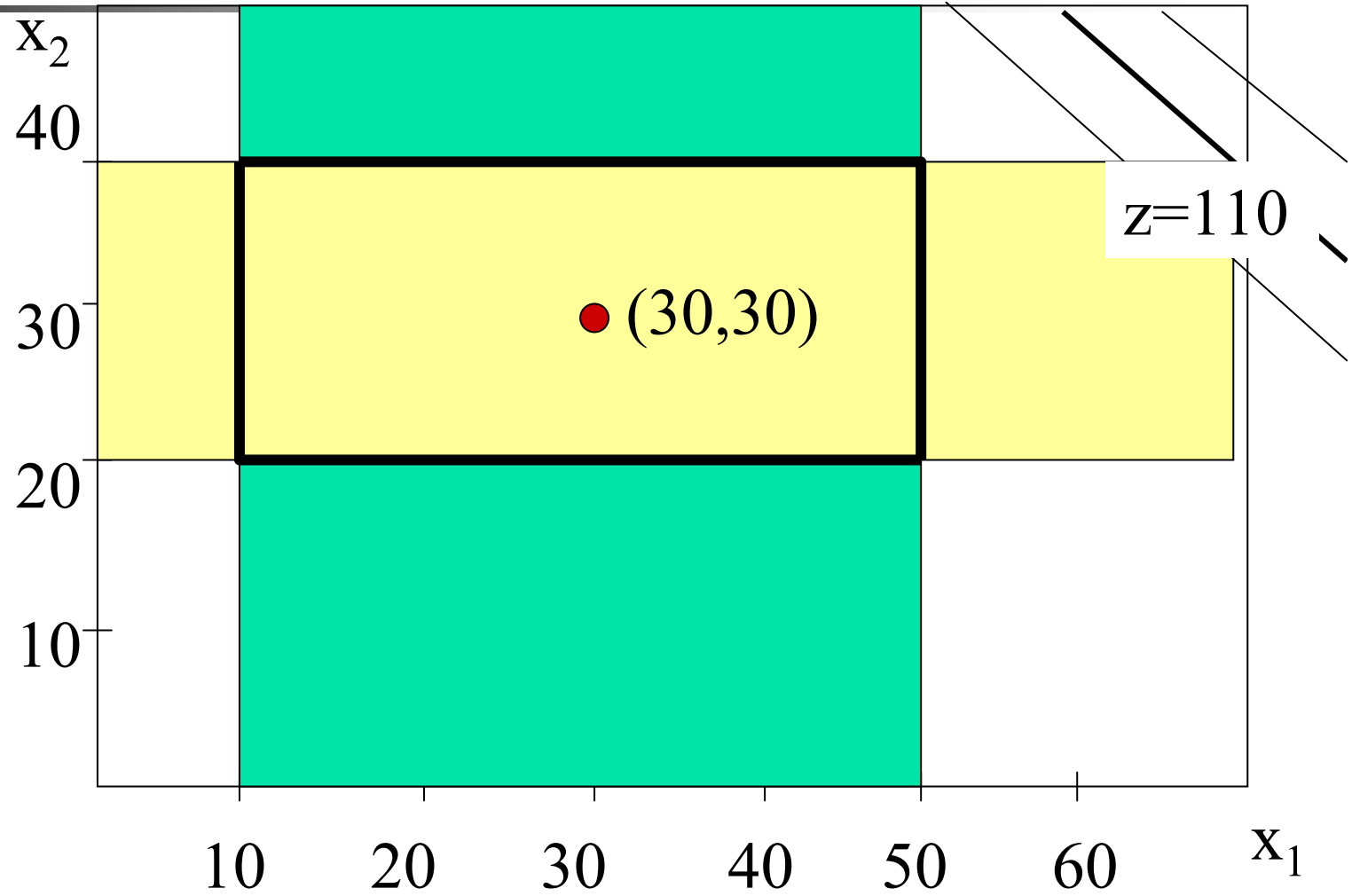Estimates: $\hat{\xi}_1 = 23 \pm 12$

$\hat{\xi}_2 = 28 \pm 9$

$\hat{\zeta} = 51 \pm 9$

Effects: - Constraints satisfied
- Error Reduction ~ error variance
- Corrections right-shifted

$x_2$

z = 50

40

30

$x=(30,30)$

$\hat{\xi} = (23,28)$

20

10

10   20   30   40   50   60   $x_1$

z=110

(30,30)

# 7. Statistical Edits  Algorithms

- PRTI: Schmid (1978) at ETH Zürich
- PRTI II: Schmid and Müller (~1983)
- QUANTOR: Müller and Schürer (~1990),  Daimler-Benz AG
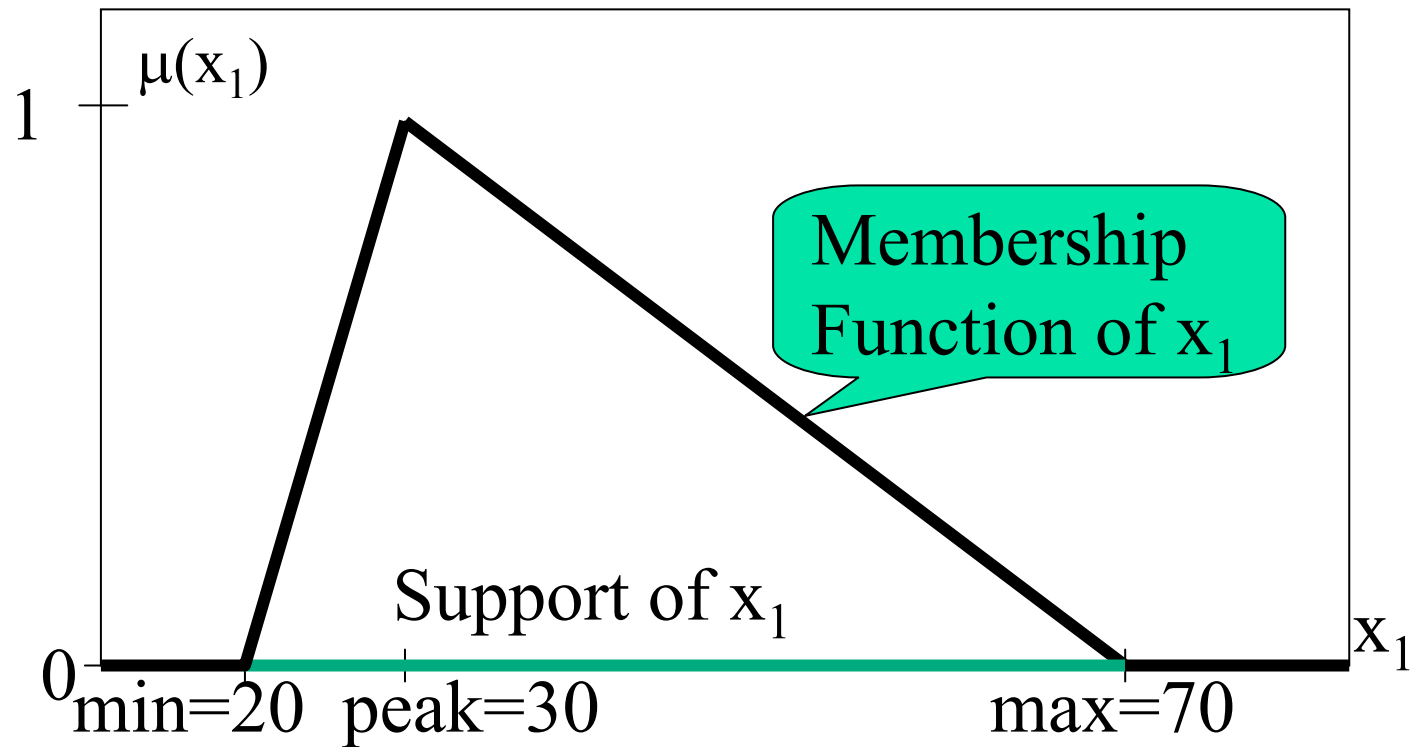
# 8. Fuzzy Edits

Why Fuzzy Edits?

- *Statistical Edits assume*:

  - Gaussian Distribution of errors

  - Linearity (Gaussian Distribution not closed under non-linear transformations!)

  - Correlation almost unknown

- *Fuzzy Edits*
  *closed under non-linear transformations*
  (Extension Principle of Zadeh)

# 8. Fuzzy Edits   Modelling

- Each variable $x_1$, $x_2$,…

  is treated as a fuzzy set (variable).

- A triangle type of membership function is assumed – may be changed.

- Error bounds from expert knowledge is used for specifying the corresponding supports of each variable.

- RHS of balance equations are separable

# 8. Fuzzy Edits  Modelling



μ(x₁)

Membership Function of x₁

Support of x₁

Statement:  $x_1 = 30$   (-10 , +40)

# 8. Fuzzy Edits  FuzzyCalc® Algorithm

$x_0$ vector of variables (observed /missing v.)

**X** product-space spanned by variables x

$G(x) = 0$ system of p nonlinear equations
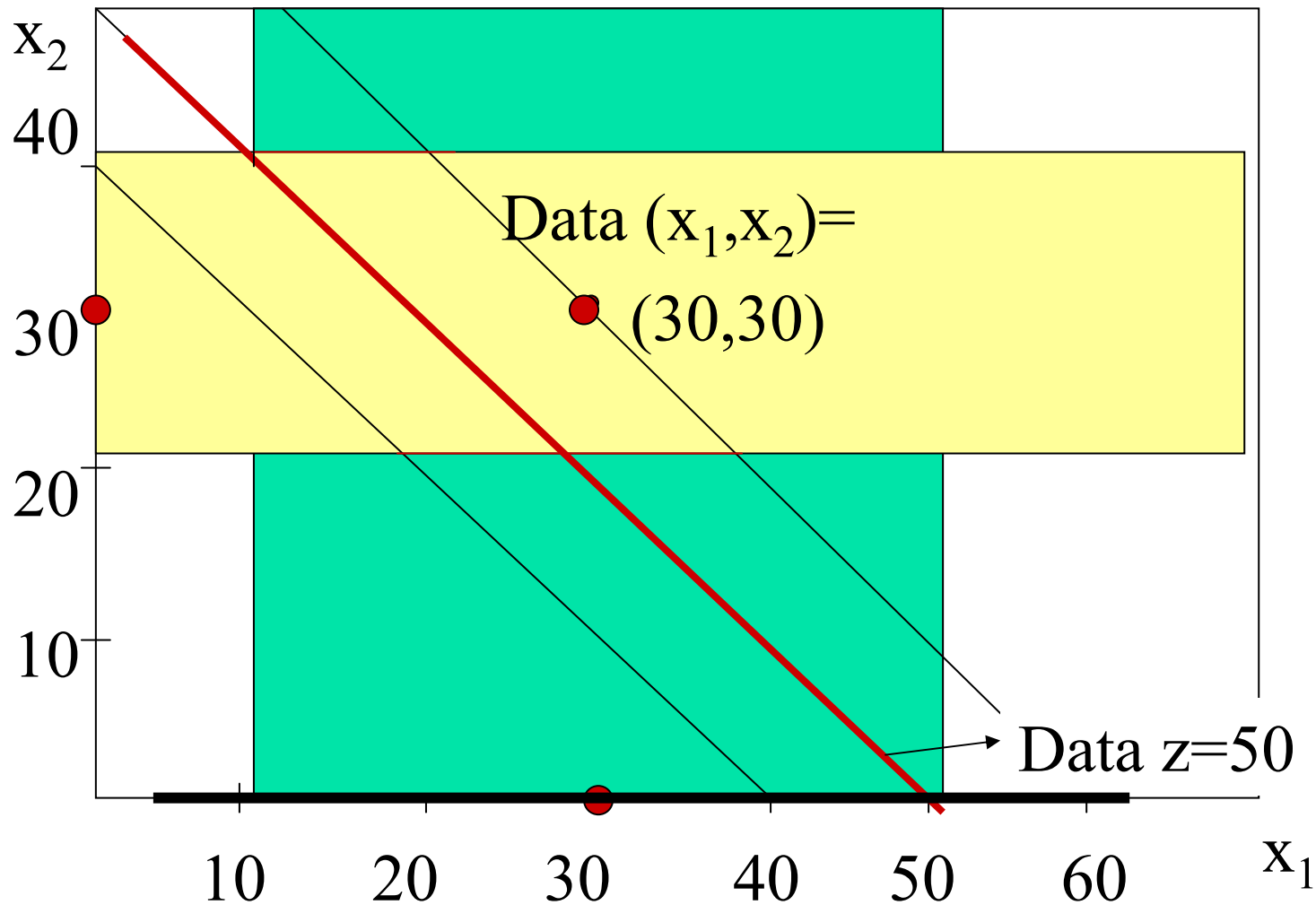
**F** set of p fully specified fuzzy sets on **X**

$$\widetilde{x} \in L = \underset{i=1}{\overset{p}{\times}} support_i$$

$support_i := support_{xi0} \cap support_{x|g1}$

$\cap \, support_{x|g2} \cap \ldots \cap support_{x|gk}$

and $G(\widetilde{x}) = 0$.

SOLD: R. Kruse (1989) Statistics on Linguistic Data;
FuzzyCalc® H.-J. Lenz and R. Müller (1999, 2000)

# 8. Fuzzy Sets Weak Inconsistency



Data $(x_1,x_2)=$ (30,30)

Data z=50

data $(x_1,x_2)=$ (30,30)

Model:
$\zeta = \xi_1 + \xi_2$

$\mu_{x_2}$

$\mu_{x_1}$

$\mu_z$

data z=50

$x_2$

$x_1$

# 8. Fuzzy Edits    Example 1
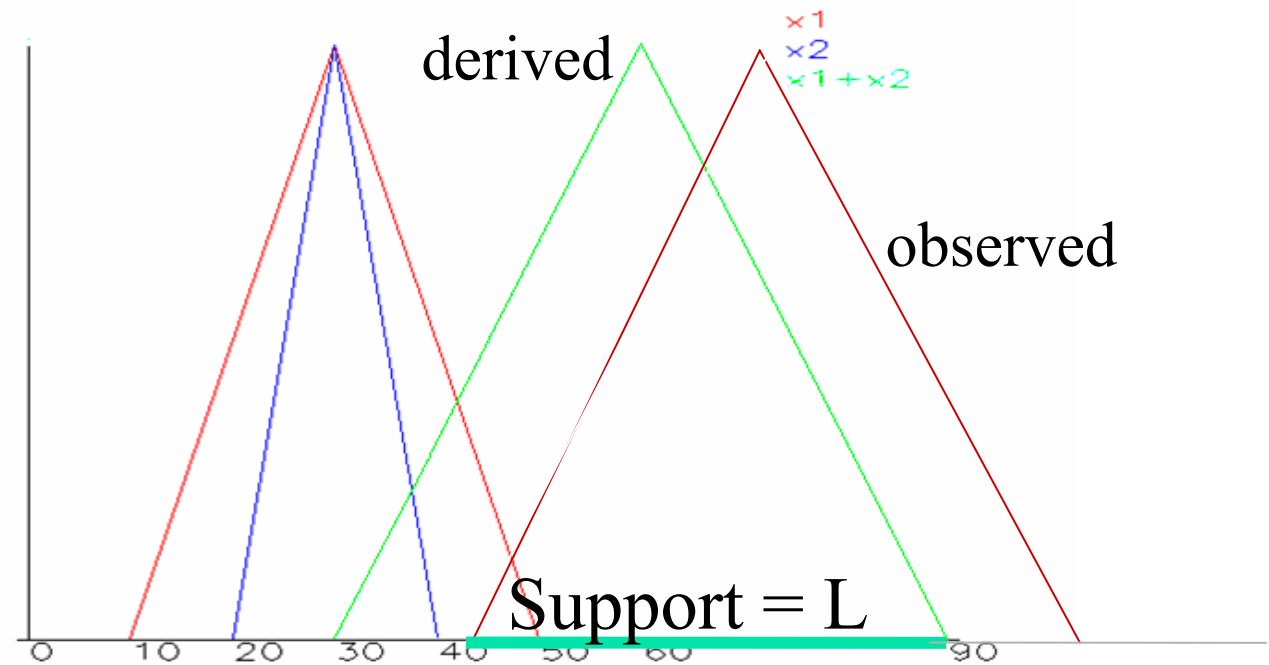
x1
x2
x1 + x2

data:  $x_1 = 30, x_2 = 30$    additive  model: $\zeta = \xi_1 + \xi_2$
derived: $\widetilde{z} = 60$
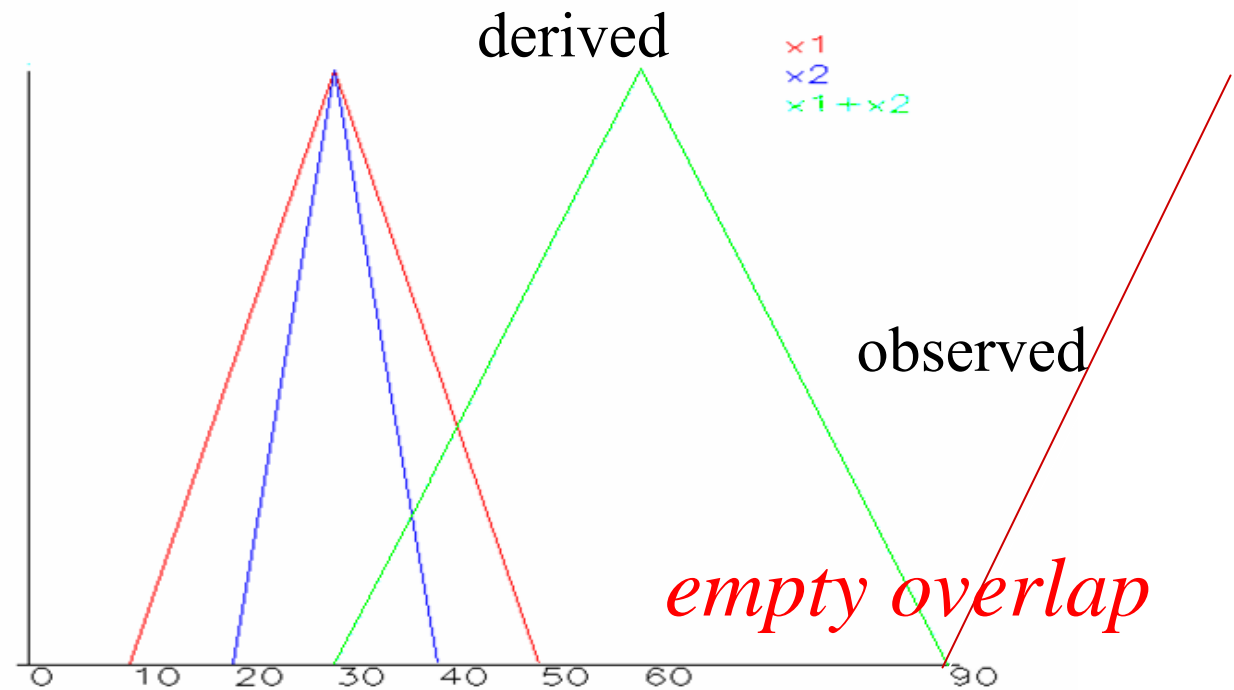
# 8. Fuzzy Edits    Example 2
## Weak Inconsistency of Data



observed: $x_1$, $x_2$, z  derived: $\widetilde{z} = x_1 + x_2$

# 8. Fuzzy Edits    Example 3
## Strong Inconsistency of Data

derived

×1
×2
×1+×2

observed

*empty overlap*

observed: $x_1$, $x_2$, $x_1 + x_2$  derived: $\tilde{z} = x_1 + x_2$

## Zadeh's Extension Principle

**THEOREM 1:** Let $p \in \mathbf{N}$ and
$\mu_1, \mu_2, \ldots, \mu_p, \eta$ (normalized) fully specified membership functions, and
$\Phi: \mathbf{R}^p \to \mathbf{R}$ a mapping with $\Phi \in \mathbf{T}$.

Then

$$\eta(x) := \sup \{\mu(x_1, x_2, \ldots, x_p) \mid \Phi(x_1, x_2, \ldots, x_p) = x,$$

$$(x_1, x_2, \ldots, x_p) \in \mathbf{R}^p \}$$

for all $x \in \mathbf{R}$ with

$$\mu(x_1, x_2, \ldots, x_p) = \min \{\mu_1(x_1), \mu_2(x_2), \ldots, \mu_p(x_p)\}$$

# 8. Fuzzy Set Theory
## Zadeh's Extension Principle     Example

Let p=2 and $\mu_1$, $\mu_2$ normalized & fully specified

Membership functions, and $\Phi: \boldsymbol{R}^2 \rightarrow \boldsymbol{R}$ a mapping

with $\Phi=\oplus$.

Then

$$\mu(y):=\sup_{(x_1, x_2)\in \boldsymbol{R}^2} \{\min \{\mu_1(x_1), \mu_2(x_2)\} \mid x_1+ x_2=y\}$$

for all $x \in \boldsymbol{R}$.

# 8. Fuzzy Sets
## Test Data

| No | Sales | Cost | Capital | Profit | P | | | |
|----|-------|------|---------|--------|---|---|---|---|
| 1 | 100 ± 5 | 80 ± 4 | 80 ± 4 | - | - | - | - | - |
| 2 | 100 ± 10 | 80 ± 8 | 80 ± 8 | - | | | | |
| 3 | 100 ± 50 | 80 ± 40 | 80 ± 40 | - | | | | |
| 4 | - | 80 ± 4 | 80 ± 4 | 20 ± 1 | | | | |
| 5 | - | 80 ± 8 | 80 ± 8 | 20 ± 2 | | | | |
| 6 | - | 80 ± 40 | 80 ± 40 | 20 ±10 | | | | |
| 7 | 100 ± 10 | 80 ± 8 | 80 ± 8 | 20 ± 2 | | | | |
| 8 | 100 ± 10 | 80 ± 8 | 80 ± 8 | 30 ± 3 | | | | |
| 9 | 100 ± 10 | 80 ± 8 | 80 ± 8 | 40 ± 4 | | | | |
| 10 | 100 ± 10 | 80 ± 8 | 80 ± 8 | - | | | | |
| 11 | 100 ± 10 | 80 ± 8 | 80 ± 8 | - | | | | |
| 12 | 100 ± 10 | 80 ± 8 | 80 ± 8 | - | | | 0,5 ± 0,05 | - |
| 13 | 100 ± 5 | 80 ± 4 | 80 ± 4 | 30 ± 1,5 | 0,2 ± 0,01 | 0,4 ± 0,02 | | - |
| 14 | 100 ± 10 | 80 ± 8 | 80 ± 8 | 30 ± 3 | 0,2 ± 0,02 | 0,4 ± 0,04 | | - |
| 15 | 100 ± 50 | 80 ± 40 | 80 ± 40 | 30 ± 15 | 0,2 ± 0,10 | 0,4 ± 0,20 | | - |

Base Model: DuPont Model

# *Quantor - FuzzyCalc*
## *Benchmark : DuPont-Model*

| No | Sales | Costs | Cap | Profit | Profitab | ROI | Turnover | M |
|----|-------|-------|-----|--------|----------|-----|----------|---|
| 1 | 100±5 | 80±4 | 80±4 | 20 ±9 | 0.2 ±0.1 | 0.25 ±0.1 | 1.25 ±0.1 | *FC* |
|  |  |  |  | 20 ±6.4 | .2±0.1 | 0.3 ±0.1 | 1.3 ±0.1 | *Q* |
| 12 | 100±10 | 80±8 | 80±8 | 35 ±3 | 0.32 ±0.04 | 0.5 ±0.05 | 1.4 ±0.2/0.1 | *FC* |
|  |  |  |  | 37 ±5.2 | 0.33 ± 0.04 | 0.48 ±0.05 | 1.44 ±0.16 | *Q* |

Interval Length!
Asymmetry!

Legend: FC FuzzyCalc
Q  Quantor
Items: 100 ± 5
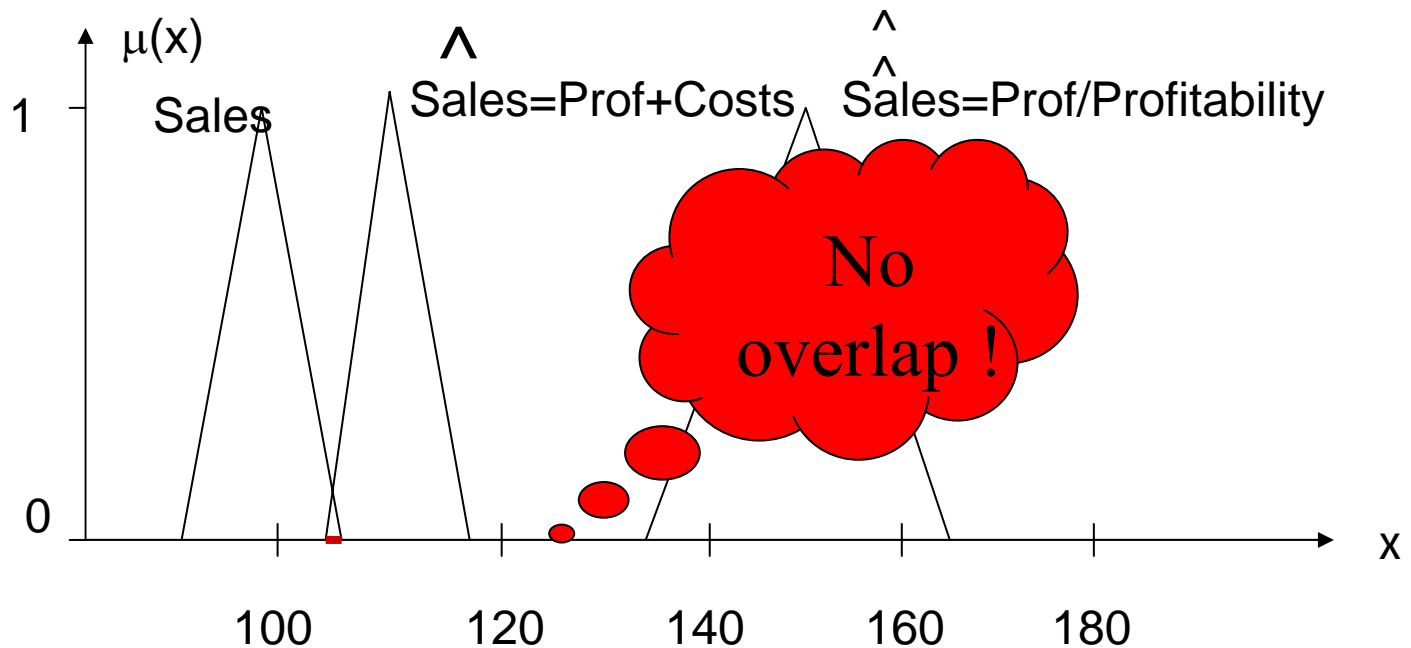
# Measurement Errors
## Effects I: Precision

| No | Sales | Costs | Capital | Profitability (**F**uzzy) | (**Q**uantor) |
|---|---|---|---|---|---|
| 1 | 100±5 | 80±4 | 80±4 | 0.2 ± 0.1 | 0.2 ± 0.06 |
| 2 | 100±10 | 80±8 | 80±8 | 0.2 ± [0.18, 0.22] | 0.2 ± 0.11 |
| 3 | 100±50 | 80±40 | 80±40 | 0.2 ± [1.6, 2.0] | 0.2 ± 0.44 |

Error Rate (posterior)

2.0

1.0

0.1

FuzzyCalc

Quantor

Error Rate (prior)

±5%  ±10%  ±20%  ±30%  ±40%  ±50%

# Measurement Errors
## Effects II: Strong Inconsistency

| No | Sales | Costs | Capital | Profit | Profitability | ROI | Turnover | |
|----|-------|-------|---------|--------|---------------|-----|----------|---|
| 13 | 100±5 | 80±4 | 80±4 | 30±1.5 | 0,2±0,01 | 0,4±0,02 | ? | **Q** |
| | 110±3 | 85±2,8 | 72±3 | 26±0,9 | 0,2±0,01 | 0,36±0,02 | 1,5 ±0.07 | |
| | - | - | - | - | - | - | - | **FQ** |



μ(x)

1

Sales

∧
Sales=Prof+Costs

∧
∧
Sales=Prof/Profitability

No overlap !

0

x

100   120   140   160   180

49

# *Computational Aspects*

- *Contraction* of prior supports
- *Intuitive correct shifts* of peaks
- *Shift Distance* ~ length of prior support =>no shift of singletons
- *Invariance Property* w.r.t. model consistent data
- *Strong Inconsistency* of data <=>
  *empty intersection* of supports
- Corrected (=estimated) data fulfill balance equations (=model)
- Bad Scalability ➜ Use edits at data entry!

# 9  MCMC Simulation

# *Example: Business Figures*

*Annual Report*

| Sales | Profit | Costs | ROI | Capital |
|-------|--------|-------|-----|---------|
| 55 | 10 | 45 | 10 | 60 |

| ± 20 | ± 2 | ± 20 | ± 5 | ± 1 |
|------|-----|------|-----|-----|

- *Balance Equations*
- Sales = Profit + Costs          (linear equation)
- ROI = 100 * Profit / Capital     (nonlinear equation)
- *Measurement Errors Model*

$$x = \xi + u \quad \text{with } u \sim N(0, \Sigma_{uu})$$

# MCMC Simulation & Estimation

**Given**

M  structural equation system (balance equations)
  of fully separable variables

x  observed state vector (missing values allowed)

$f_x$  fully specified density function of x
  (given mean and standard deviation)

$C_{xx}$ correlation matrix

**Wanted**

$\hat{f}, \hat{\xi}, \hat{\sigma}$  density, mean and standard deviation for
  each state variable

# SamPro-Algorithm
## *Resoving, Sampling, Estimation and Projection*

**begin**

**resolve** (set) LHS $\equiv$ RHS (x) for all variables of all equations

**sample** from the joint density function of all RHS variables for each LHS (z)

**estimate** f, $\mu$, $\sigma$ of all LHS variables

**estimate** $\alpha/2$, $1-\alpha/2$-quantiles $\underline{q}_{max}, \overline{q}_{min}$ for each variable, i=1,2,…,p .

if $\overline{q}_{min} < \underline{q}_{max}$ flag "M -inconsistency" else

**compute** the distribution $\hat{f}_{xz}$ restricted to the subspace x-z = 0, and $\mu_x$ and $\sigma^2$ (in case of k= 2 equations).

**end**

54

# M-Inconsistency (k = 2)

$f_z$                                                                $f_x$

Example: Sales = Profit + Costs

LHS                          RHS

$\overline{q}_{min}$                    $\underline{q}_{max}$

$\overline{q}_{min} = \{\overline{q}_1, \overline{q}_2, ..., \overline{q}_k\}$                    $\underline{q}_{max} = \{\underline{q}_1, \underline{q}_2, ..., \underline{q}_k\}$

inconsistent

weak consistent

$I_q$

$$I_q = [\underline{q}_{max}, \overline{q}_{min}]$$

# 6  Experiments and Analysis

**Experimental Group A**
   (missing values allowed)
1.  Single Effect of skewness
2.  Effect of correlation
3.  Interaction Effect of Skewness and Correlation

**Experimental Group B**
   (complete data sets, Gaussian distributions, correlation between variables)
1.  Effect of M-consistency
2.  Effect of M-inconsistency

# Exp.-Group 1: Scenario 1
## Normality - no Correlation

*Specification of the distributions:*
*Profit* $\sim N(20, 2^2)$; *Costs* $\sim N(80, 8^2)$; *Capital* $\sim N(60, 6^2)$
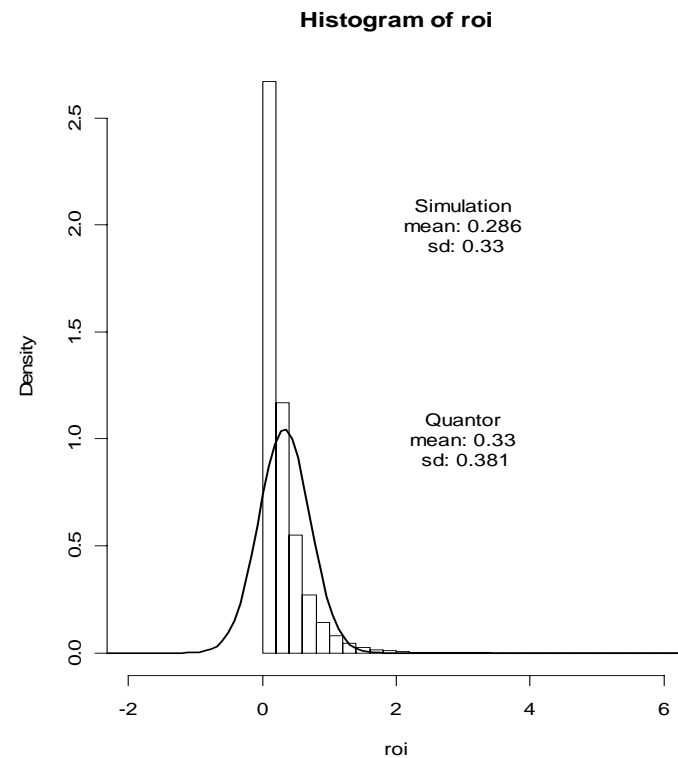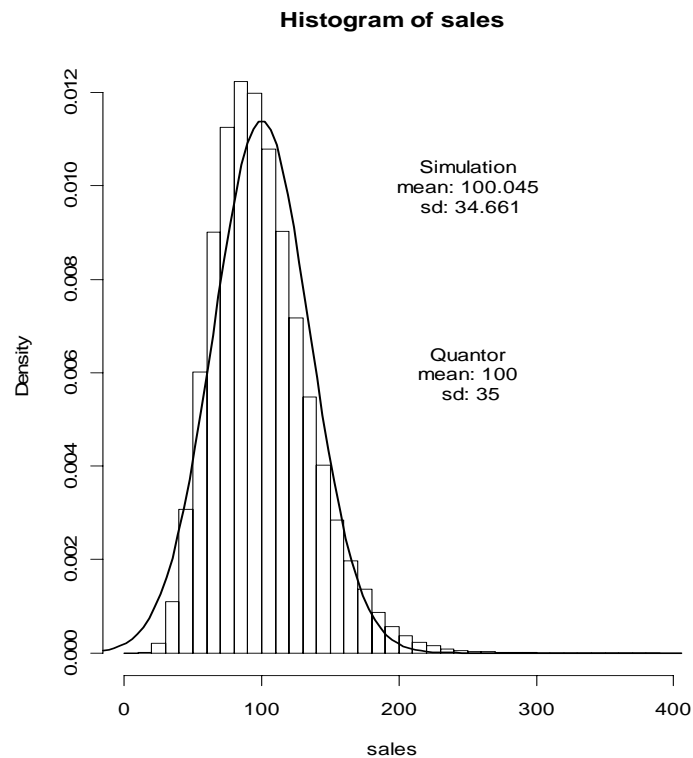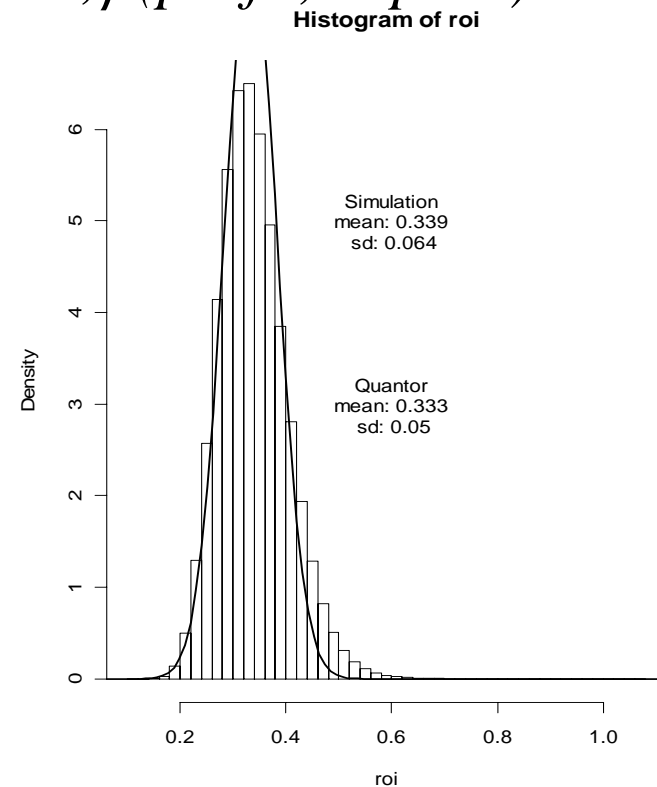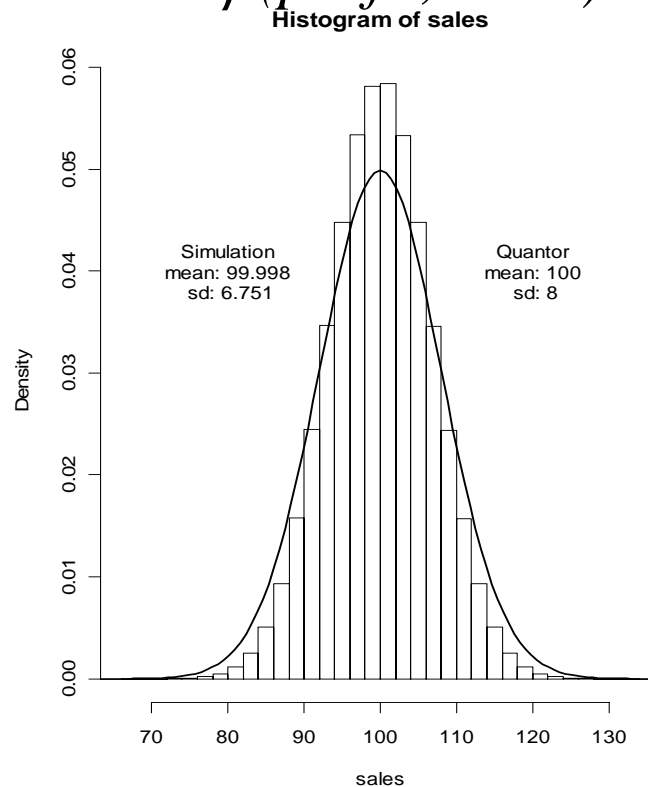
*Missing values: Sales, ROI*

**Histogram of sales**

Simulation
mean: 99.998
sd: 21.552

Quantor
mean: 100
sd: 22

Density

sales

**Histogram of roi**

Simulation
mean: 0.337
sd: 0.34

Quantor
mean: 0.333
sd: 0.34

Density

roi

57

# Scenario 2
## Skewness - no correlation

***Specification of the distributions****: Profit  ~ Exp(1/20) vs. N(20, 20²); Costs  ~ N(80, 8²); Capital  ~ N(60, 6²)*

***Missing values****: Sales, ROI*



**Histogram of sales**

Simulation
mean: 100.009
sd: 21.537

Quantor
mean: 100
sd: 22

**Histogram of roi**

Simulation
mean: 0.253
sd: 0.255

Quantor
mean: 0.333
sd: 0.34

# Scenario 3
## Skewness - no Correlation

***Specification of the distributions****: Profit ~ Exp(1/20) vs. N(20, 20²); Costs ~ Gamma(8,0.1) vs. N(80, 28²); Capital ~ Gamma(15, 0.25) vs. N(60, 8.6²)*

***Missing values****: Sales, ROI*



**Histogram of sales**

Simulation
mean: 100.045
sd: 34.661

Quantor
mean: 100
sd: 35

**Histogram of roi**

Simulation
mean: 0.286
sd: 0.33

Quantor
mean: 0.33
sd: 0.381

# Scenario 4
## Normality - neg. Correlation

**Specification of the distributions:**

*Profit ~ N(20, 2²); Costs ~ N(80, 8²); Capital ~ N(60, 6²)*

**Missing values:** *Sales, ROI*

**Correlation:** $\rho$(profit, costs) = -0.7; $\rho$(profit, capital) = -0.7



Histogram of sales

Simulation
mean: 99.998
sd: 6.751

Quantor
mean: 100
sd: 8

Histogram of roi

Simulation
mean: 0.339
sd: 0.064

Quantor
mean: 0.333
sd: 0.05

# Scenario 5
## Normality - pos. Correlation

*Specification of the distributions*:

*Profit ~ N(20, 2²); Costs ~ N(80, 8²); Capital ~ N(60, 6²)*

**Missing values***: Sales, ROI*

**Correlation:** *$\rho$(profit, costs) = 0.7; $\rho$(profit, capital) = 0.7*



**Histogram of sales**

Simulation
mean: 100.001
sd: 9.506

Quantor
mean: 100
sd: 8

**Histogram of roi**

Simulation
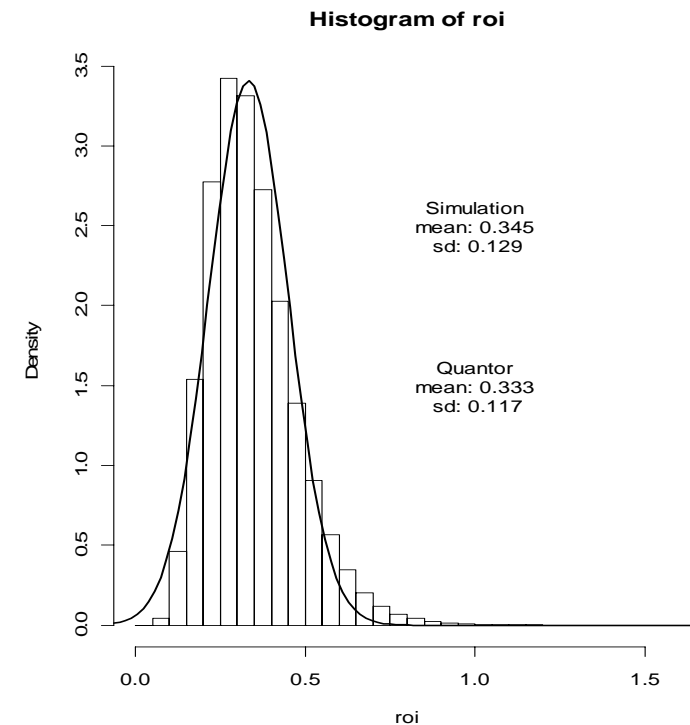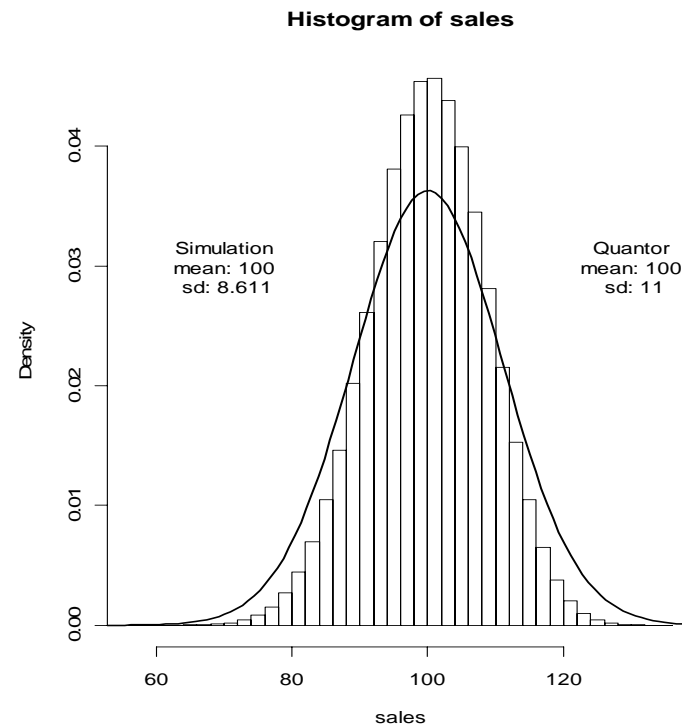mean: 0.334
sd: 0.026

Quantor
mean: 0.333
sd: 0.05

# Scenario 6
## Skewness - neg. Correlation

***Specification of the distributions****: (Profit, Costs, Capital) ~ Dir(10, 40, 30) vs. Profit ~ N(20, 5.9²); Costs ~ N(80, 8.9²); Capital ~ N(60, 8.6²)*

***Missing values****: Sales, ROI*

***Correlation****: $\rho$(profit, costs) = -0.4; $\rho$(profit, capital) = -0.3*



**Histogram of sales**

Simulation
mean: 100
sd: 8.611

Quantor
mean: 100
sd: 11

**Histogram of roi**

Simulation
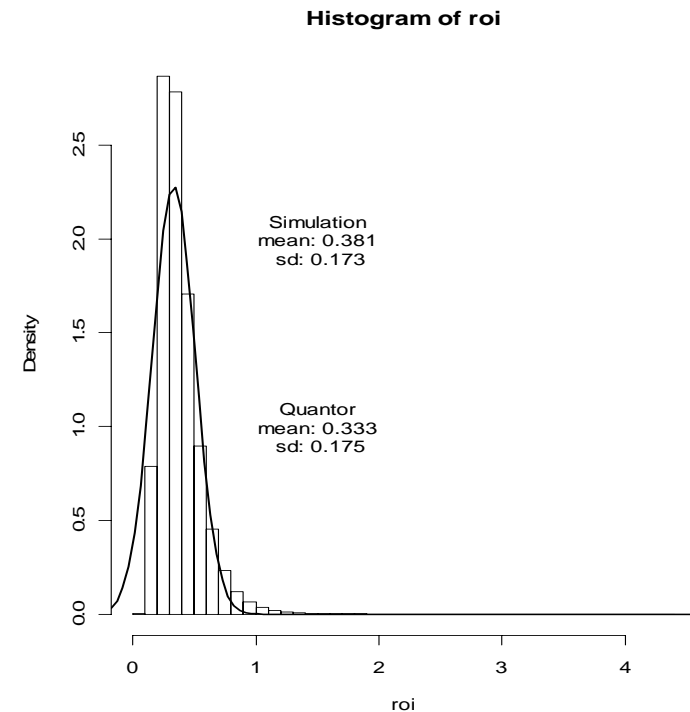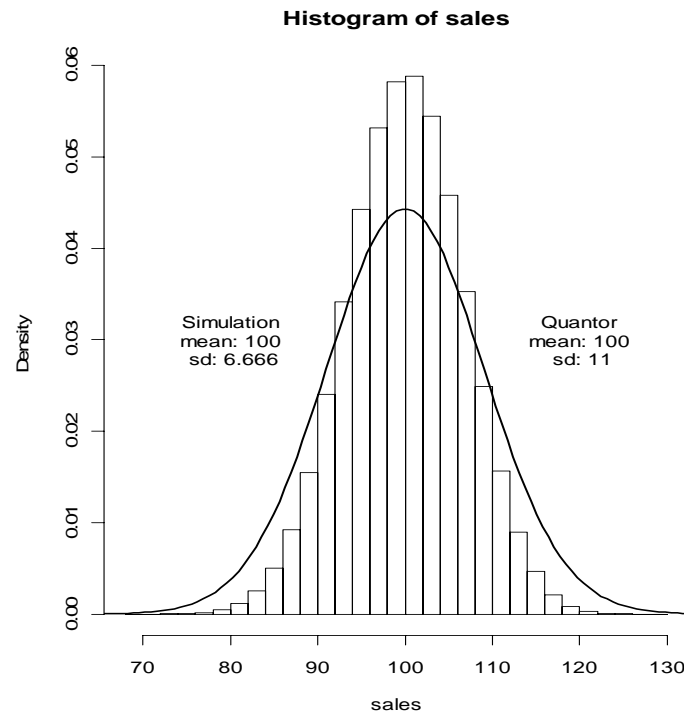mean: 0.345
sd: 0.129

Quantor
mean: 0.333
sd: 0.117

# Scenario 7
## Skewness - pos. Correlation

***Specification of the distributions****:(Profit, Costs, Capital) ~ Dir(30, 40, 8) vs. Profit ~ N(20, 2.8²); Costs ~ N(80, 8.8²); Capital ~ N(60, 20²)*
***Missing values****: Sales, ROI*
***Correlation****:* $\rho(profit, costs) = -0.8$; $\rho(profit, capital) = -0.3$



**Histogram of sales**

Simulation
mean: 100
sd: 6.666

Quantor
mean: 100
sd: 11

**Histogram of roi**

Simulation
mean: 0.381
sd: 0.173

Quantor
mean: 0.333
sd: 0.175

# Exp.-Group B:
## Normality, Correlation

*Prior Information:*

Sales and ROI M-consistent vs. M-inconsistent

*Missing values: no*

*Correlation Matrix* $(\rho \in \{-0.4, 0.0, +0.4\})$

$$R = \begin{bmatrix} 1 & \rho & 0 & -\rho & -\rho \\ \rho & 1 & 0 & -\rho & \rho \\ 0 & -\rho & 1 & 0 & \rho \\ -\rho & -\rho & 0 & 1 & -\rho \\ -\rho & \rho & \rho & -\rho & 1 \end{bmatrix}$$

# Scenario 1
## Normality, Correlation

***Specification of the distributions:*** *Profit $\sim N(20, 2^2)$; Costs $\sim N(80, 8^2)$; Capital $\sim N(60, 6^2)$; Sales $\sim N(100, 10^2)$; ROI $\sim N(0.333, 0.0333^2)$*
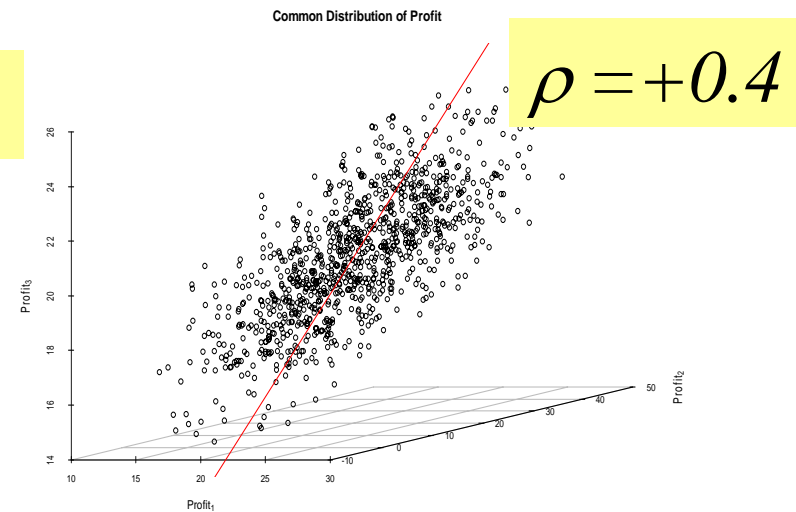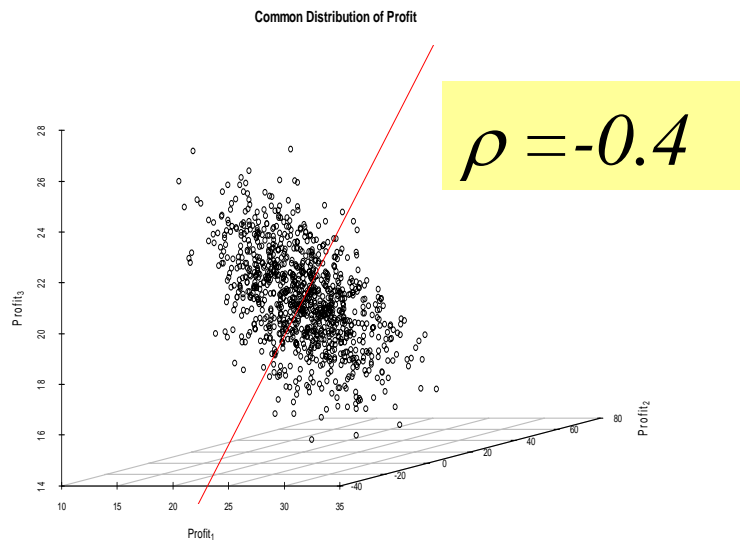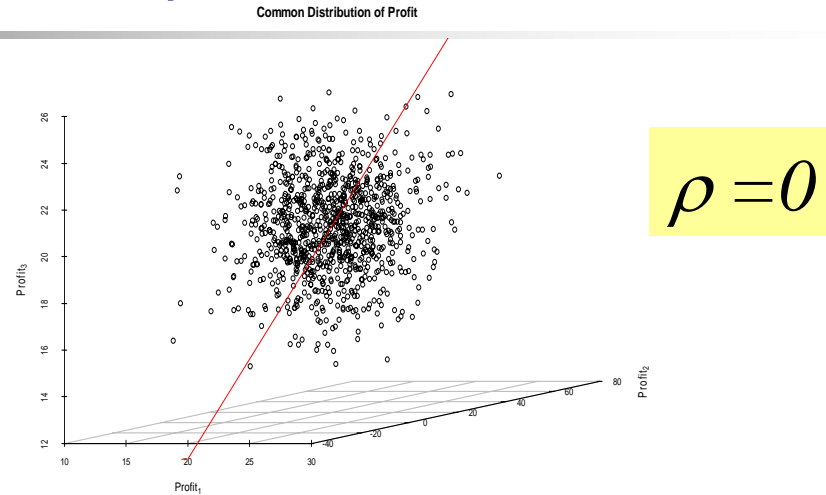***Missing values****: no*

**Special Set-up**:
- M-consistent observations of Sales, ROI
- $\rho = 0$

| Variable | Mean | Stdv |
|----------|------|------|
| **profit** | 19.97 | 1.58 |
| **costs** | 79.95 | 6.29 |
| **capital** | 59.81 | 4.89 |
| **sales** | 99.96 | 6.37 |
| **ROI** | 0.33 | 0.03 |

# Scenario 1
## 3D scatter plots of Profit

1. Profit $\sim N(20, 2^2)$
2. Profit = Sales -Costs
3. Profit = Capital * ROI



Common Distribution of Profit

$\rho = 0$



Common Distribution of Profit

$\rho = -0.4$



Common Distribution of Profit

$\rho = +0.4$

# Scenario 2
## Normality, Correlation

***Specification of the distributions:*** Profit $\sim N(30, 3^2)$;
Costs $\sim N(80, 8^2)$; Capital $\sim N(60, 6^2)$; Sales $\sim N(100, 10^2)$;
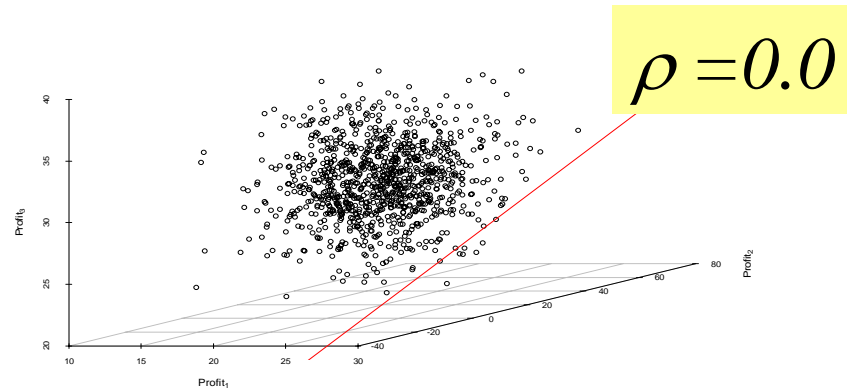ROI $\sim N(0.333, 0.333^2)$

***Missing values***: No

**Special Set-up**:
- M-inconsistent observations of Sales, ROI
- $\rho = 0$

| Variable | Mean | Stdv |
|----------|--------|------|
| **profit** | 24.84 | 2.18 |
| **costs** | 76.27 | 6.36 |
| **capital** | 67.09 | 5.05 |
| **sales** | 105.74 | 6.50 |
| **ROI** | 0.37 | 0.03 |

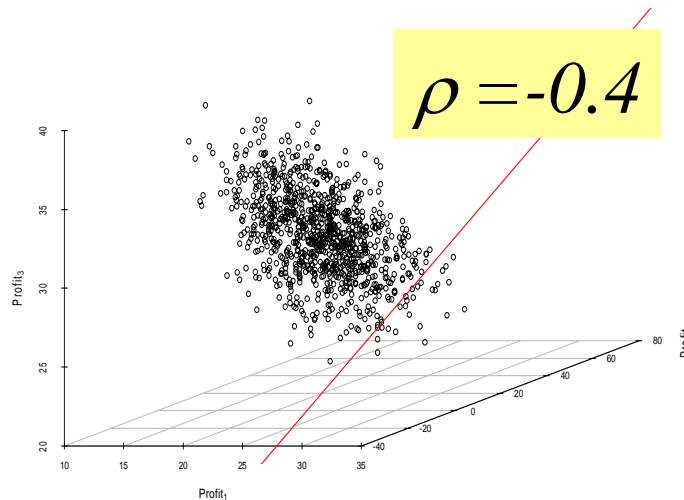# 3D Scatterplots of Profit

1. Profit $\sim N(30,3^2)$
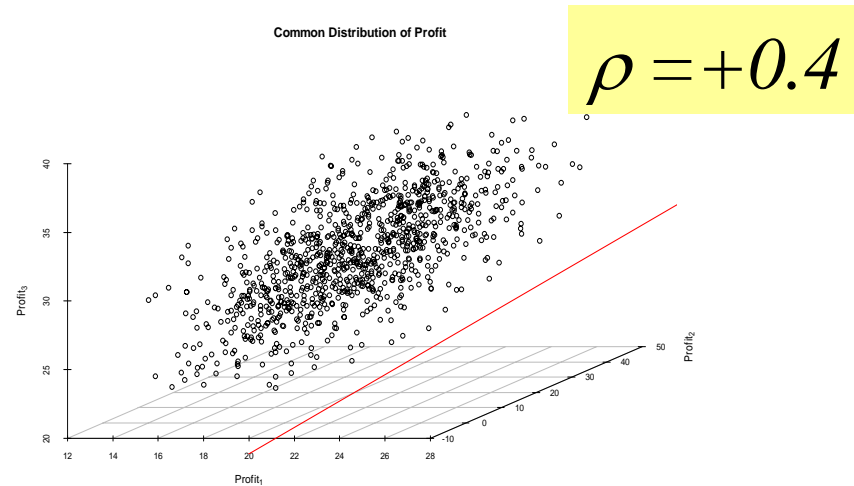2. Profit = Sales - Costs
3. Profit = Capital * ROI



Common Distribution of Profit

$\rho = 0.0$



Common Distribution of Profit

$\rho = -0.4$



Common Distribution of Profit

$\rho = +0.4$

strong M-inconsistency !

# Summary

1. No correlation:
   (the means of ) the simulated quantities are about the same as the GLS estimates under a Gaussian hypothesis.

2. Skewness of distributions:
   mostly has only a small effect on the estimates.

3. Positive cross-correlations of the variables: can lead to severe problems: The balance equation system may become M-inconsistent !

# Basic Literature

1. **Batini, C., Scannapieco, M.: Data Quality: Concepts, Methods and Techniques. Heidelberg: Springer Verlag (2006)**

2. **Dombrowski, Erik und Lechtenböger, Jens: Evaluation objektorientierter Ansätze zur Data-Warehouse-Modellierung, Datenbank-Spektrum 15/2005**

3. **Naumann, Felix: Datenqualität, Informatik-Spektrum_30_1_2007**

4. Naumann, Felix: Quality-driven Query Answering for integrated Information Systems, LNCS 2261, Springer, Heidelberg, 2002

# *Edits* - A Firewall against inconsistent Data Entry

- **Thank you!**

- **hjlenz@wiwiss.fu-berlin.de**
- **http://www.wiwiss.fu-berlin.de**