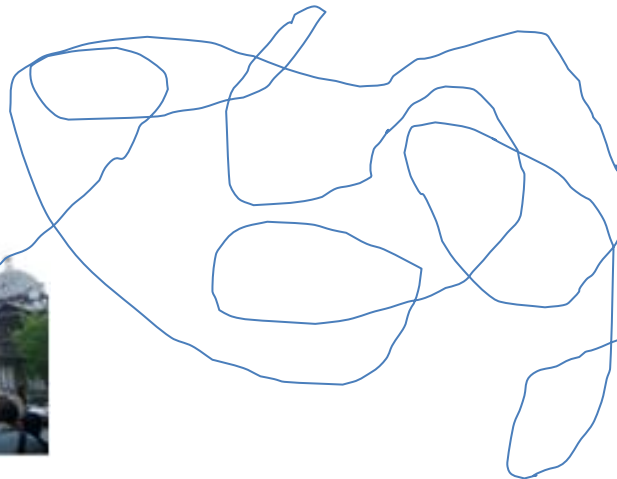


# *on doing fine work in statistics*



Old Huxley Building



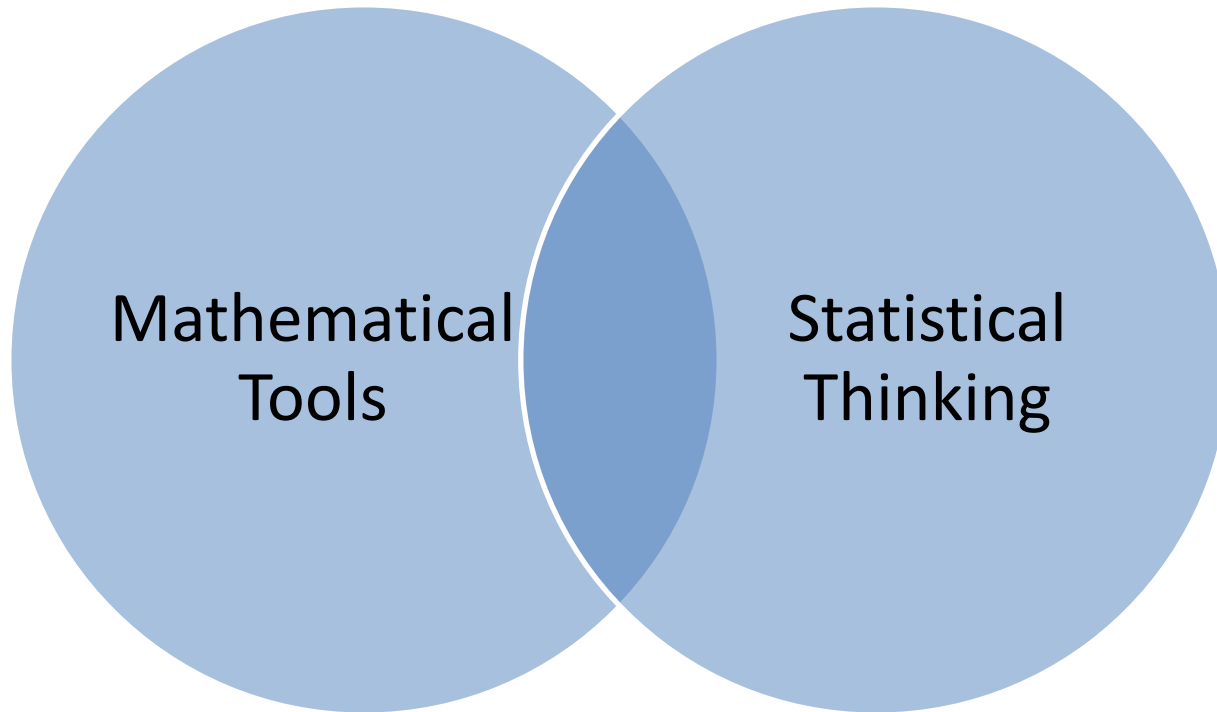
RSS 2013

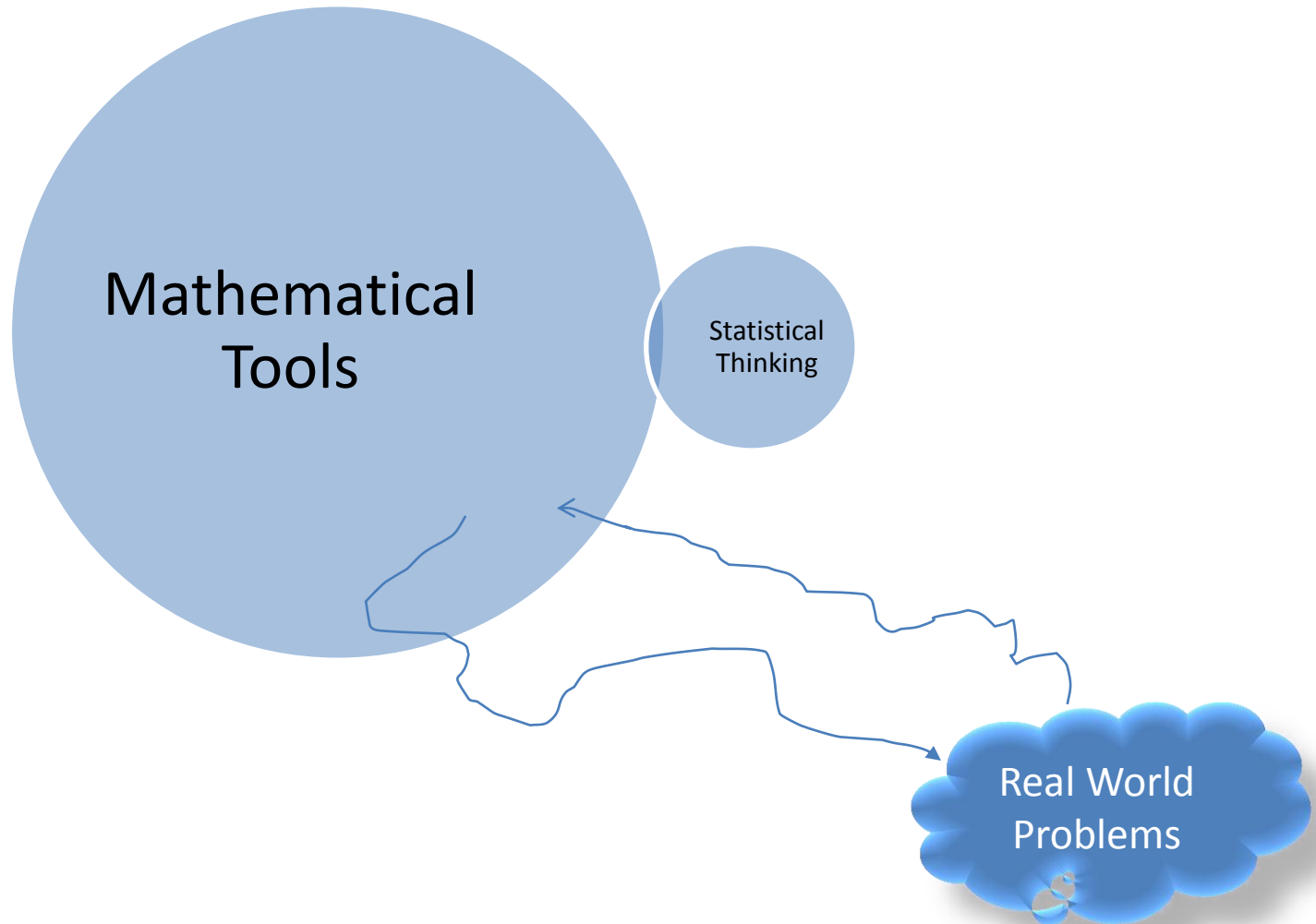
*"Much fine work in statistics involves minimal mathematics; some bad work in statistics gets by because of its apparent mathematical content."*

David Cox (1981),  
Theory and general principle in statistics, JRSS(A), 144, pp. 289-297.

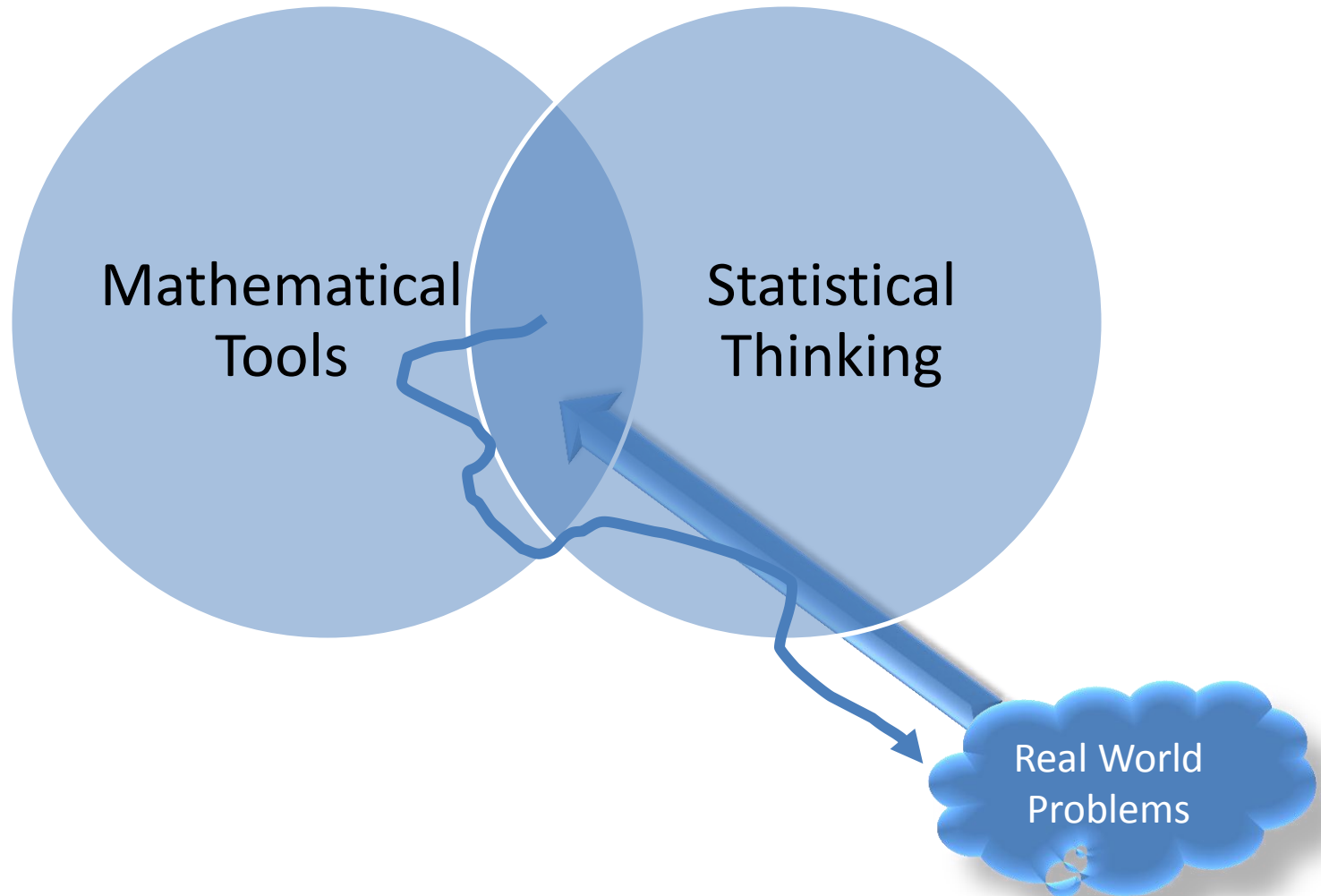


**The challenge of solving  
real world problems  
with mathematical tools  
and statistical thinking**





*The mathematical statistician*



**Applications of Bayesian Networks to Big Data, with open source implications**

Ron S. Kenett      Michael Ashcroft



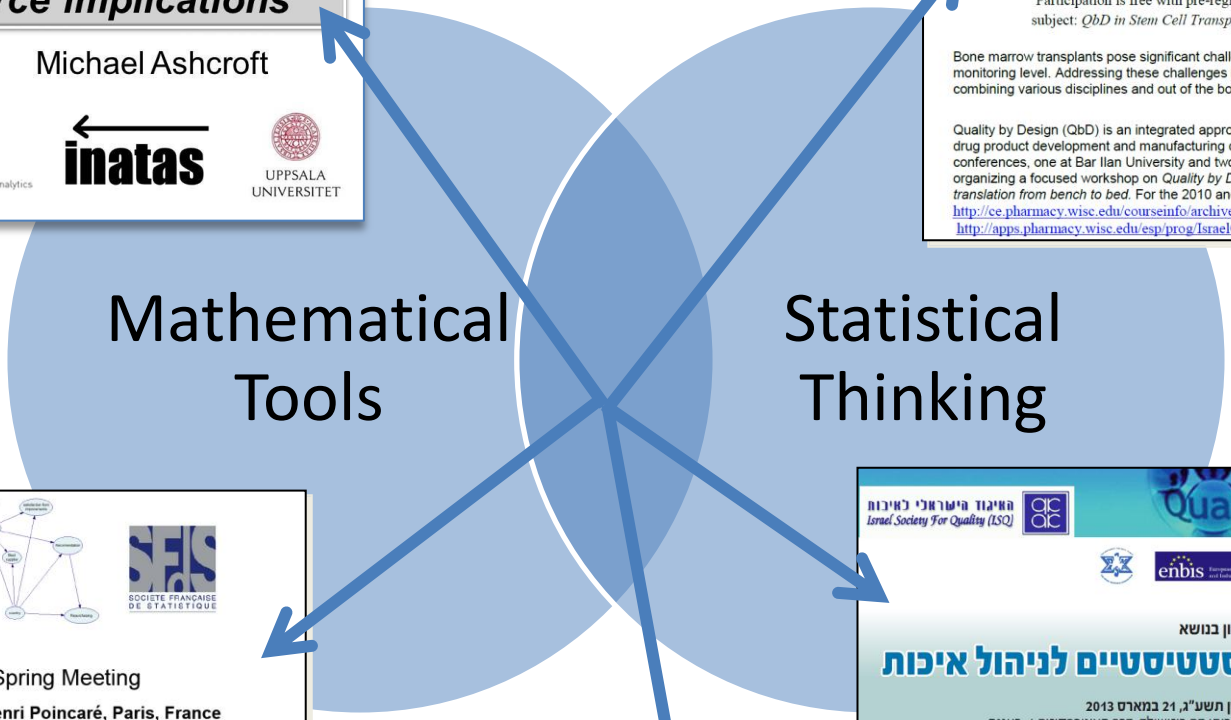
**The Hebrew University of Jerusalem**  
**The School of Pharmacy Institute for Drug Research**

**Quality by Design in Stem Cell Transplantation: Rational translation from bench to bed**  
 Wednesday, May 1<sup>st</sup>, 2013  
 Faculty of Medicine, Butnar Building, Butnar Small Hall, Ein Kerem, Jerusalem  
 13:00 – 18:00

A Half Day Conference for Multidisciplinary Brainstorming  
 Participation is free with pre-registration by sending an email with subject: *QbD in Stem Cell Transplantation* to [avrir@ekmd.huji.ac.il](mailto:avrir@ekmd.huji.ac.il)

Bone marrow transplants pose significant challenges at the therapeutic, clinical and monitoring level. Addressing these challenges requires an interdisciplinary perspective combining various disciplines and out of the box thinking.

Quality by Design (QbD) is an integrated approach that applies to the clinical, analytical and drug product development and manufacturing domains. Following 3 successful QbD conferences, one at Bar Ilan University and two at the Hebrew University in Jerusalem, we are organizing a focused workshop on *Quality by Design in Stem Cell Transplantation: Rational translation from bench to bed*. For the 2010 and 2012 QbD conferences see: <http://ce.pharmacy.wisc.edu/courseinfo/archive/2012Israel> <http://apps.pharmacy.wisc.edu/esp/prog/IsraelQBD/>



**ENBIS-SFdS 2014 Spring Meeting**  
**9-10-11<sup>th</sup> of April, 2014 - Institut Henri Poincaré, Paris, France**

The European Network for Business and Industrial Statistics (ENBIS) and the French Statistical Society (SFdS) are planning an ENBIS-SFdS spring meeting at the Institut Henri Poincaré (IHP) in the center of Paris in April 2014. The topic of the meeting is:

**Graphical causality models: Trees, Bayesian Networks and Big Data**

The objective of the meeting is to investigate new methodological advances in the analytics, causality mechanisms and robust modelling. These areas of theoretical and research concern a large spectrum of statisticians in various fields of specialisation. The scientific format of this special topic spring meeting is:

- On 9/4/2014: A general talk on causality networks, open to the public plus several case studies (½ day, free admission, Amphitheatre Hermite). An open meeting with presentations from researchers, practitioners and software companies (½ day, free admission, Amphitheatre Hermite).
- On 10-11/4/2014: A research workshop, by invitation, to specialists interested in discussing case studies and new developments in an informal setting (2 days, Amphitheatre Darboux). The workshop will consist of technical presentations of methods and tools effectively applied to a collection of well-documented set of examples and case studies. The plan is to organize this event with free registration but without financial support for transportation costs and living expenses.

**קישורים חשובים**  
 LinkedIn Group – Israel  
 Statistical Association

**התפתחויות ביישומים סטטיסטיים לניהול איכות**  
**יום עיון בנושא**  
**יום חמישי, 1 ביוני 2013, 21 במרץ 2013**  
 קריית האוניברסיטה הפתוחה ע"ש דורותי דה רוטשילד, דרך האוניברסיטה 1, רעננה (הכניסה לתניה דרך צומת רעננה צפון; התניה ללא תשלום)

**סדנת הביוסטטיסטיקה – סדר יום ותקצירי הרצאות**  
 מאת זיסלי | בקסגוריה 2013  
 יום רביעי 20 במרץ, 2013

הסדנה השנתית של האיגוד התקיים בתאריך 10.4.2013 בבית התפוזות בתל אביב, בשעות 9.00-16.00. ההרשמה לסדנה בלינק <https://tixwise.co.il/he/biostat>

תכנית הסדנה	שעות
התכנסות וכבוד קל	8.15
ד"ר יוסי טל – אסטרונומיה וסטטיסטיקה בניסויים קליניים	9.00
פרופ' איילה כהן – מילופון לנרגיסיה הלינארית וישעון מחקר רפואי	10.00
ד"ר יוסי לוי – גל מה שרציתם לדעת על ה-P-value הפקסת זהירים*	11.00
ד"ר חות הלר – בדיקת השעויות מרובות על ידי FDR	14.00
פרופ' רון קנת – תכנון ניסויים ו-Quality by Design	15.00
סיום	16.00



*J. R. Statist. Soc. A*, (1980),  
143, Part 4, pp. 383–430

## Sampling and Bayes' Inference in Scientific Modelling and Robustness

By GEORGE E. P. BOX

*University of Wisconsin–Madison*

[Read before the ROYAL STATISTICAL SOCIETY at a meeting organized by the South West Group on Thursday, May 15th, 1980, the President SIR CLAUS MOSER in the C

# Warning



**We do not teach tools and methods for doing that**

**“This is not as the analyst of a single set of data, nor even as the designer and analyzer of a single experiment, but rather as a colleague working with an investigator throughout the whole course of iterative deductive-inductive investigation.”**

Problem elicitation

RLD

CED

PSE

InfoQ

1

2

3

4

5

**Are we generating knowledge? (InfoQ)**      **Are we making an impact? (PSE)**

Data integration  
(ETL, data fusion)

Statistical education  
(concept science)

Problem elicitation  
(cognitive science)

**Statistical Thinking**

Confidentiality

Visualization

Integrated Models

and data exploration

Unstructured data

(semantic data, networks)

(HMI) Causality  
(CS, BN)

Reproducible research  
(Sweave, CFR Part11)

**Other Disciplines**

# Problem Elicitation



# Problem Elicitation

‘The teacher acts as a mentor in training the student in **unstructured problem solving** while the computer stores and retrieves information. This sets the mind free to do what it does best – **be inductively creative**’ (Box, 1997)

‘Most iterative investigations involve **multi-dimensional learning**: in particular, how to study the problem is a necessary precursor to how to solve the problem’ (Box, 2000)

**The Theory of Applied Statistics**

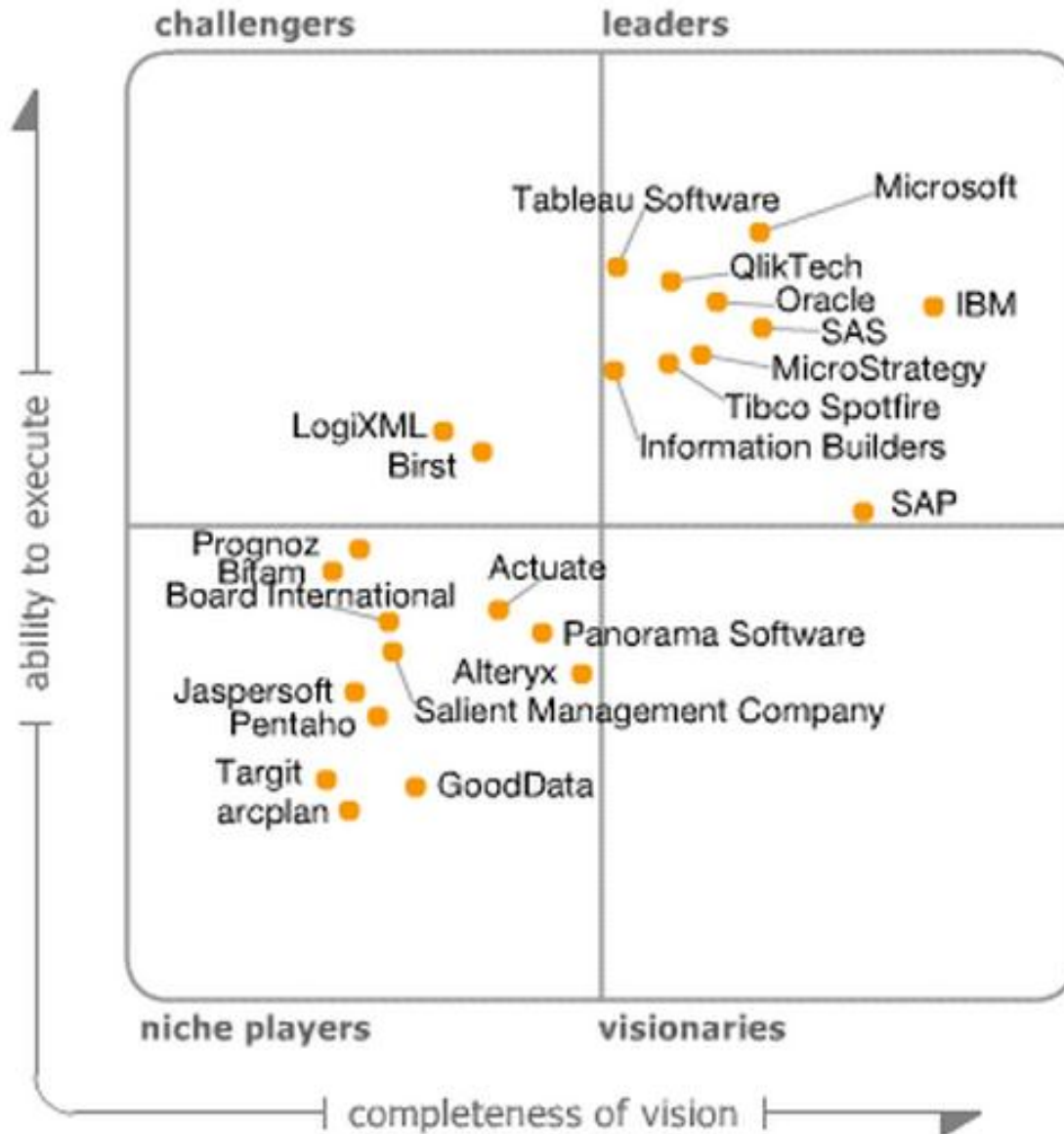
# Visualization: Some (old) references

- **The Visual Display of Quantitative Information**, E. Tufte, Graphics Press, 1983.
- **Graphical Perception: Theory, Experimentation, and Application to the Development of Graphical Methods**, W. Cleveland and R. McGill, *Journal of the American Statistical Association*, 79 (387):531-554, 1984.
- **Graphs in Scientific Publications**, W. Cleveland, *The American Statistician*, 38 (4):261-269, 1984.
- **How to Display Data Badly**, H. Wainer, *The American Statistician* 38(2):137-147, 1984.
- **Data-Based Graphics: Visual Display in the Decades to Come**, J. Tukey, *Statistical Science*, 5(3): 327-339 , 1990.
- **Communicating Statistics**, T. Greenfield, *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 156 (2): 287-297, 1993.

# Business Intelligence and Analytics Platforms Market Segment

“The dominant theme of the market in 2012 was that **data discovery** became a mainstream BI and analytic architecture. The market also saw increased activity in **real time, content and predictive analytics.**” Gartner, 2013

**Figure 1.** Magic Quadrant for Business Intelligence and Analytics Platforms



As of February 2013

Gartner's Magic Quadrant for Business Intelligence and Analytics Platforms  
5 February, 2013.  
Analyst(s): K. Schlegel, R. Sallam, D/ Yuen, J. Tapadinhas

## Web Search Interest: obama

Worldwide, 2004 - present

Categories: [Arts & Current Events \(50-75%\)](#), [Entertainment \(0-10%\)](#), [Society \(0-10%\)](#), [more...](#)

2

Totals [?](#)

obama  10

### Interest over time

Forecast [?](#)  News headlines

[Learn what these numbers mean](#)



\* The last value on the graph is based on partial data and may change. [Learn more](#)

 [Embed this chart](#)


# Do they interact?

## Web Search Interest: depression

Worldwide, 2004 - present

Categories: [Health \(50-75%\)](#), [Arts & Humanities \(10-25%\)](#), [Society \(0-10%\)](#), [Lifestyles \(0-10%\)](#), [more...](#)

Totals [?](#)

depression  63

### Interest over time

Forecast  News headlines

[Learn what these numbers mean](#)



\* The last value on the graph is based on partial data and may change. [Learn more](#)

 [Embed this chart](#)

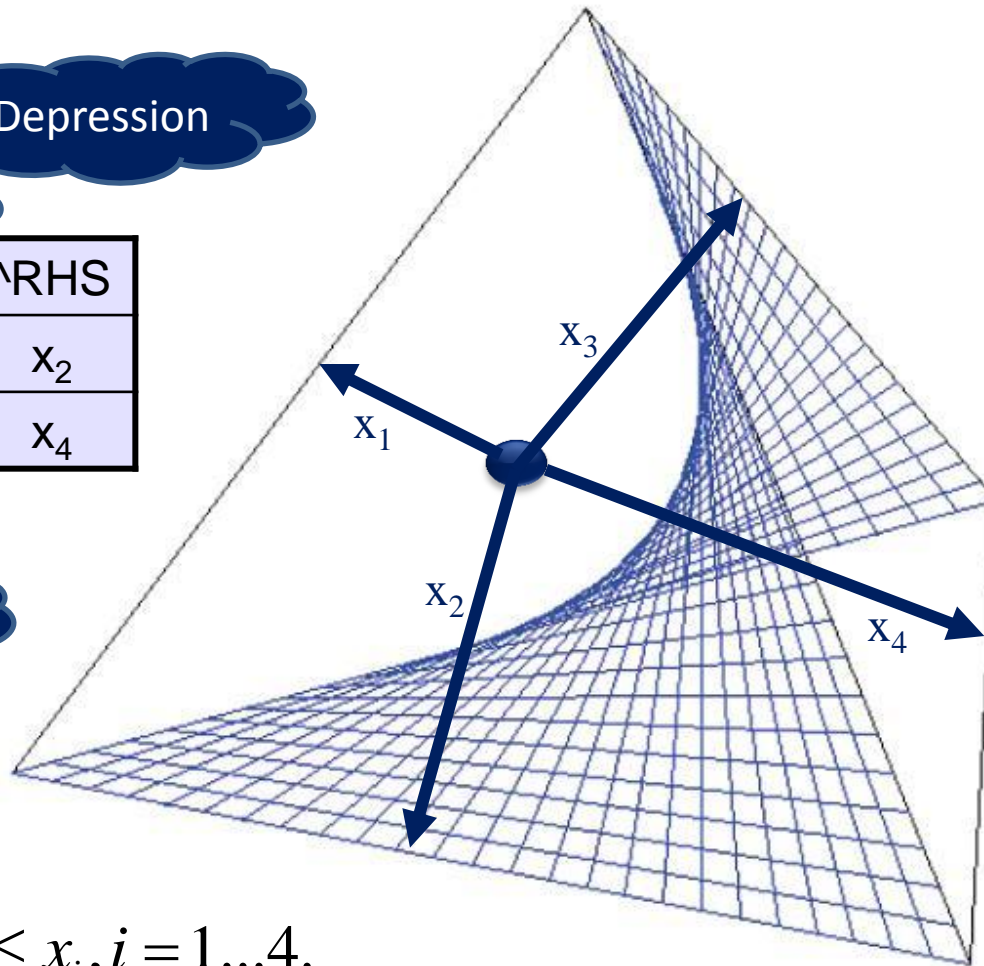


# The Simplex

Depression

	RHS	^RHS
LHS	$x_1$	$x_2$
^LHS	$x_3$	$x_4$

Obama



Bishop, Fienberg,  
& Holland, 1975

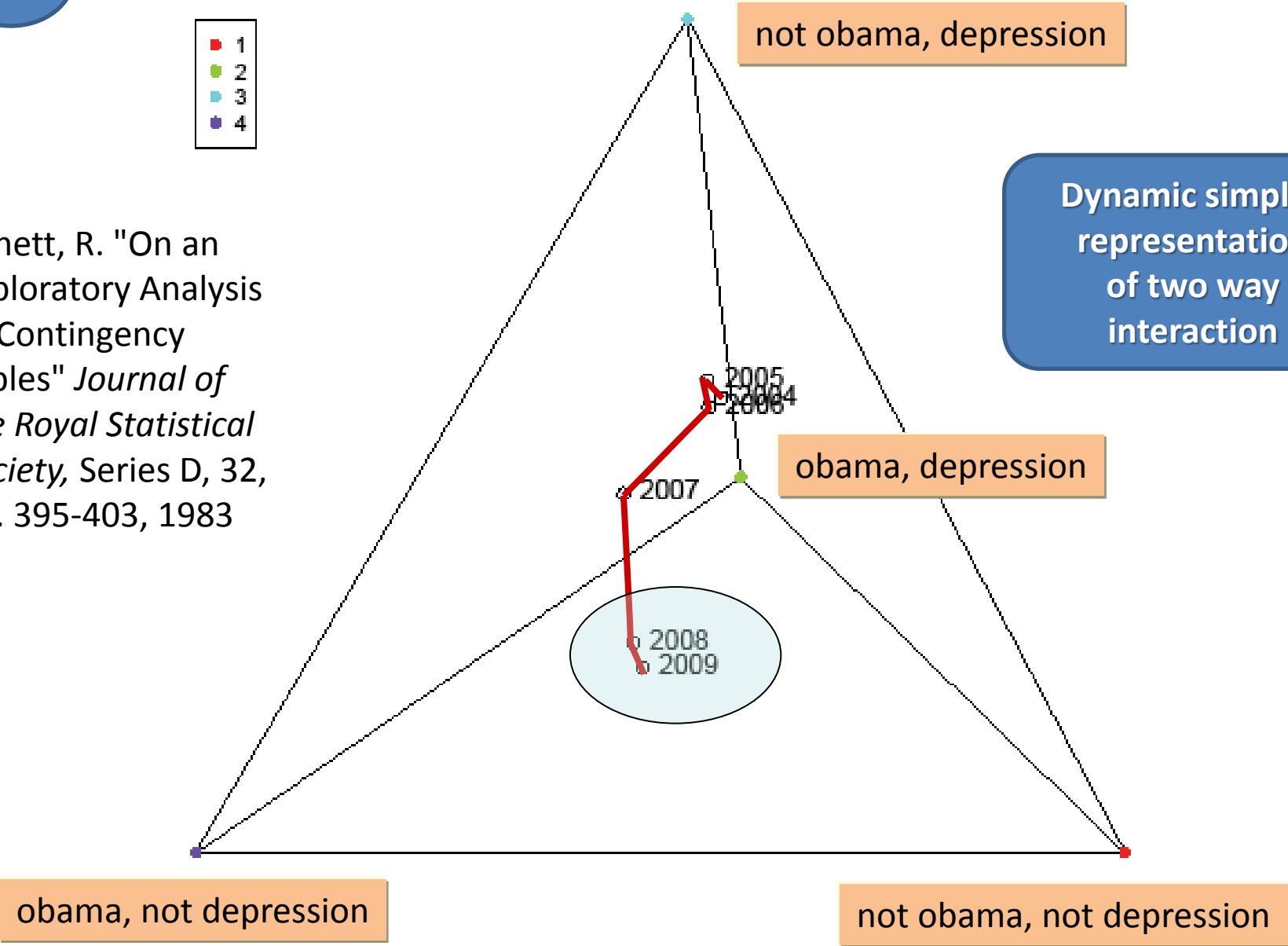
$$D = x_1 x_4 - x_2 x_3 = 0$$

$$\sum_{i=1}^4 x_i = 1, 0 \leq x_i, i = 1 \dots 4.$$

### Depression AND Obama - 3D

- 1
- 2
- 3
- 4

Kenett, R. "On an Exploratory Analysis of Contingency Tables" *Journal of the Royal Statistical Society, Series D*, 32, pp. 395-403, 1983



# Linkage Disequilibrium

$$D = x_1 x_4 - x_2 x_3$$

$$x_1 = fg + D$$

$$x_2 = (1-f)g - D$$

$$x_3 = f(1-g) - D$$

$$x_4 = (1-f)(1-g) + D$$

	RHS	^RHS
LHS	$x_1$	$x_2$
^LHS	$x_3$	$x_4$

$$g = x_1 + x_2$$

$$f = x_1 + x_3$$

$$\sum_{i=1}^4 x_i = 1, \quad 0 \leq x_i, i = 1 \dots 4.$$

*D can be extended to more dimensions...*

# Linkage Disequilibrium

$$\underline{X} = \underline{f} \otimes \underline{g} + D \underline{e} \otimes \underline{e}$$

where

$$\underline{X} = (x_1, x_2, x_3, x_4)$$

$$\underline{f} = (f, 1-f)$$

$$\underline{g} = (g, 1-g)$$

$$\underline{e} = (1, -1)$$

$$\sum_{i=1}^4 x_i = 1, \quad 0 \leq x_i, i = 1 \dots 4.$$

$$f = x_1 + x_3$$

$$g = x_1 + x_2$$

An algebraic observation...

# Relative Linkage Disequilibrium

$D$  is the distance from the point corresponding to the contingency table in the simplex, to the surface  $D=0$  in the  $e \otimes e$  direction.

$$RLD = \frac{D}{D_M}$$

$D_M$  is the distance from the point corresponding to the contingency table on the surface  $D=0$  in the  $e \otimes e$  direction, to the surface of the simplex, in that direction.

*If  $D > 0$   
then  
if  $x_3 < x_2$*

*then*  $RLD = \frac{D}{D + x_3}$

*else*  $RLD = \frac{D}{D + x_2}$

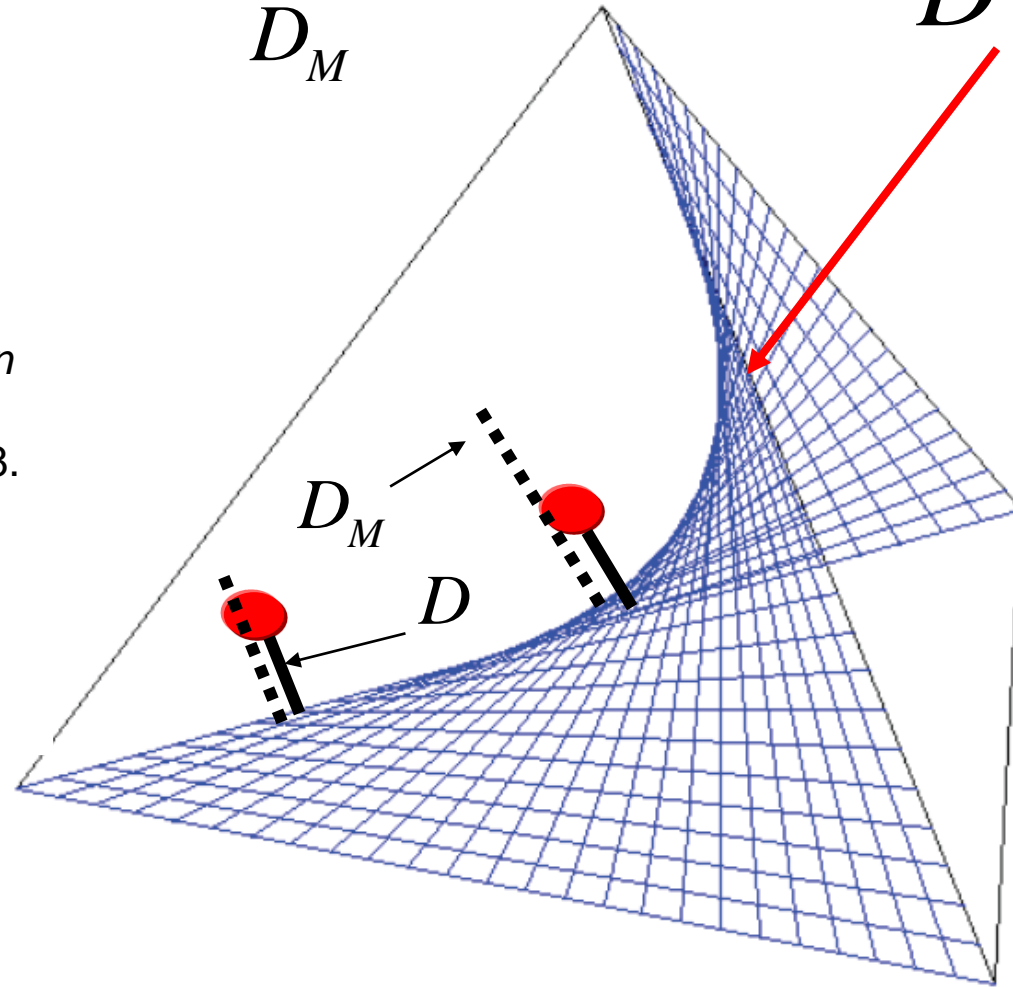
*else  
if  $x_1 < x_4$*

*then*  $RLD = \frac{D}{D - x_1}$

*else*  $RLD = \frac{D}{D - x_4}$

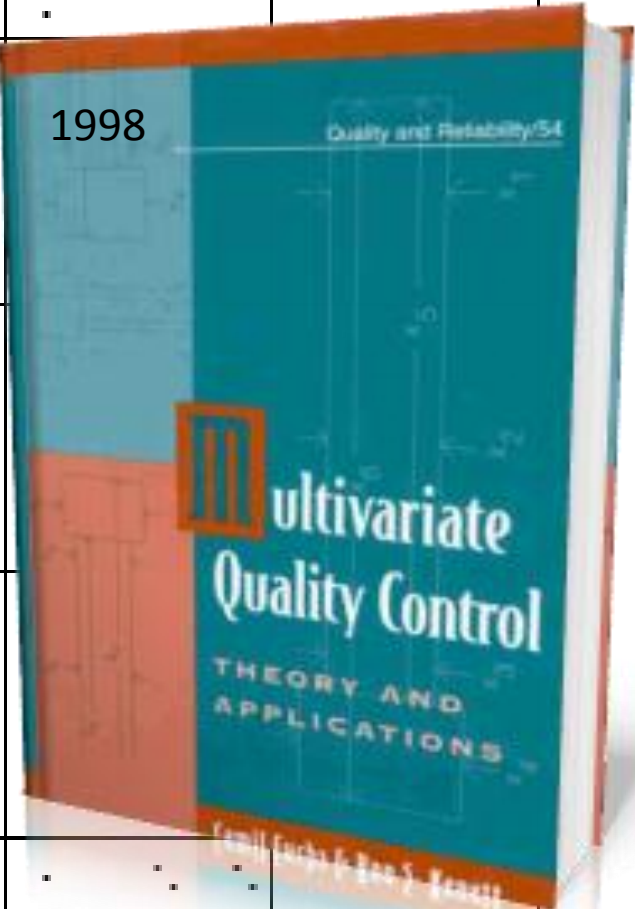
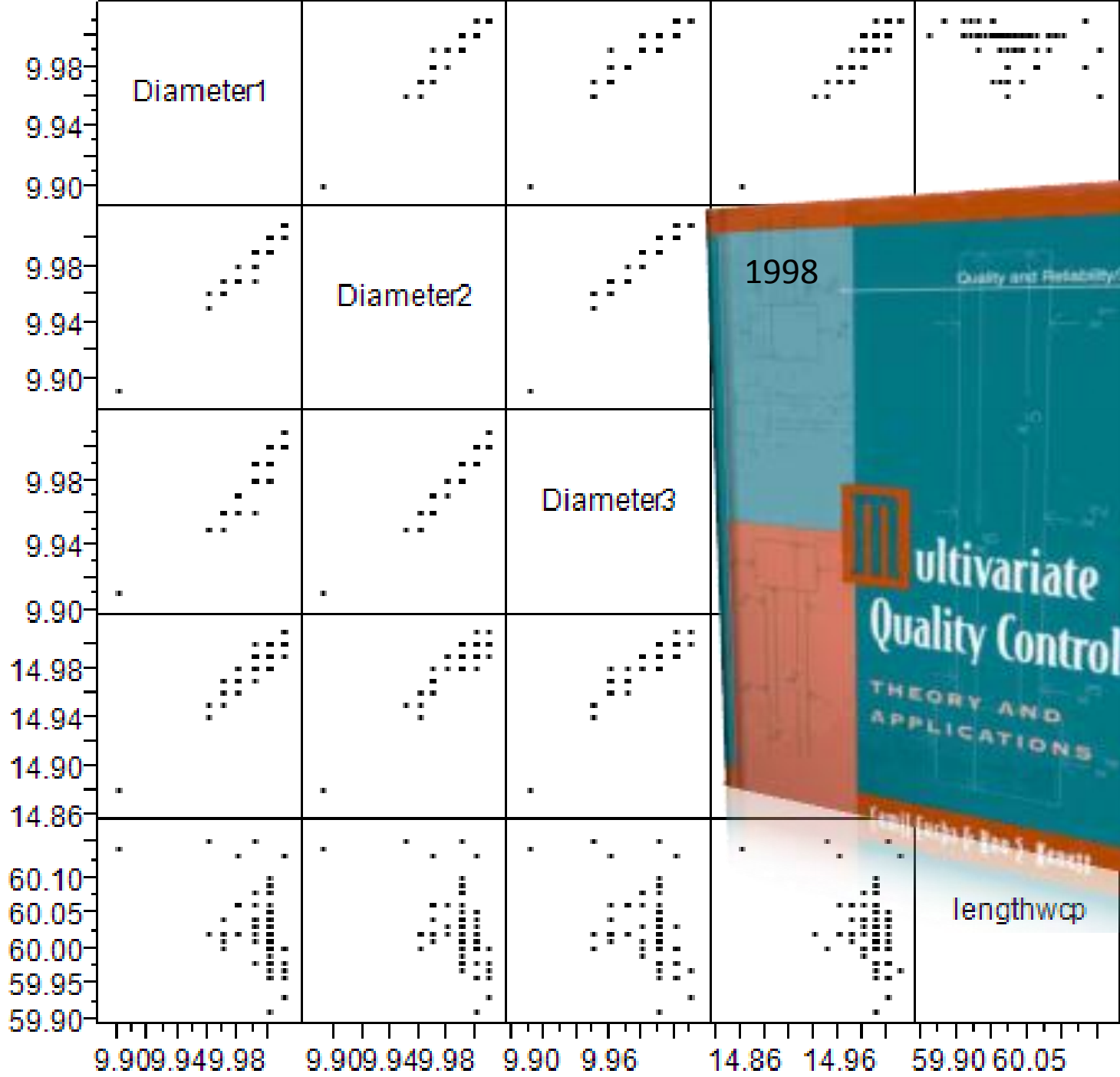
$$RLD = \frac{D}{D_M}$$

$$D = 0$$

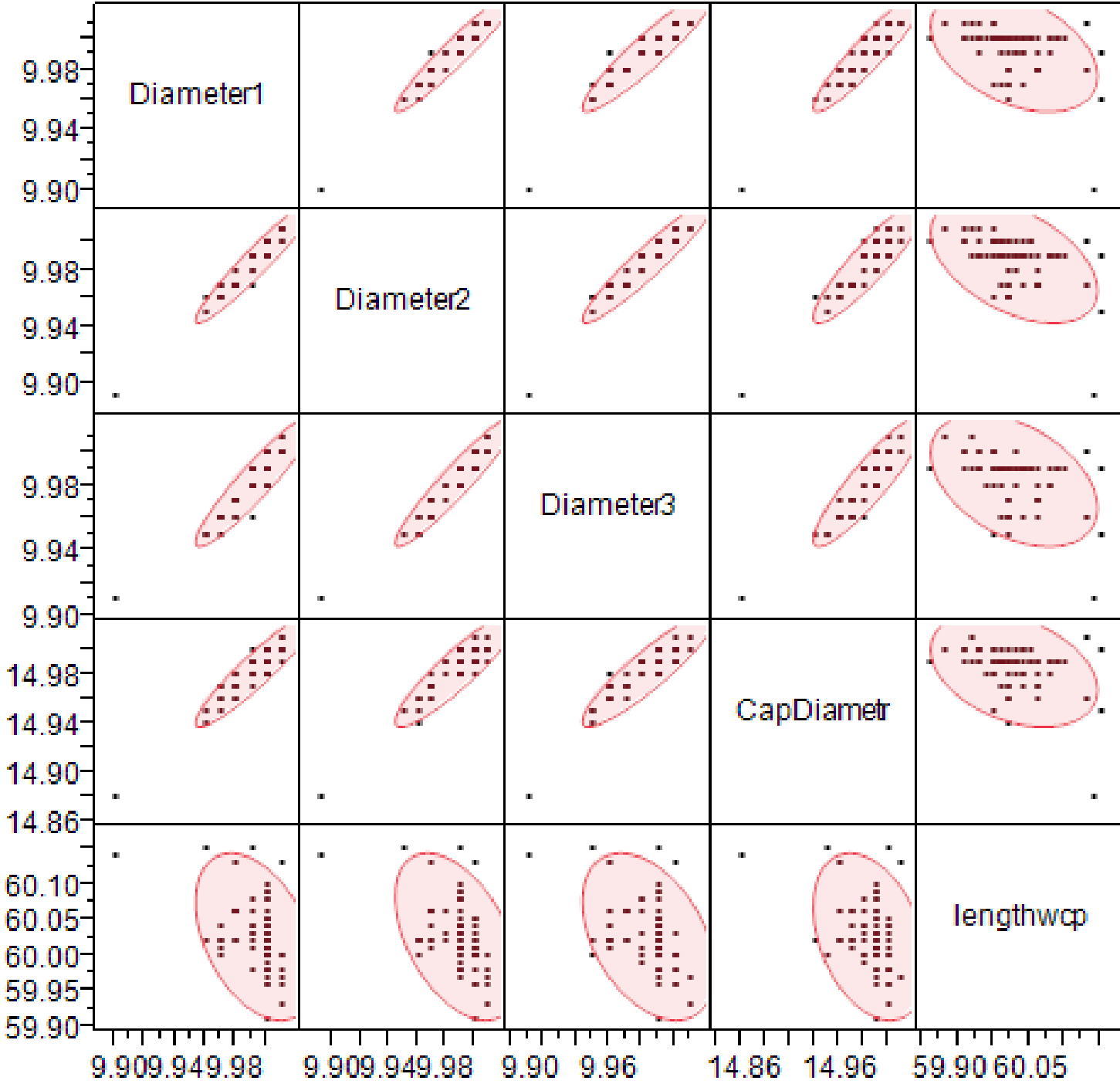


Kenett, R. and Salini, S.,  
 "Relative Linkage  
 Disequilibrium Applications to  
 Aircraft Accidents and  
 Operational Risks" . *Trans. on  
 Machine Learning and Data  
 Mining*, 1,(2), pp. 83-96, 2008.

**Implemented in arules R  
 Package. Version 0.6-6,  
 Mining Association Rules  
 and Frequent Itemsets**

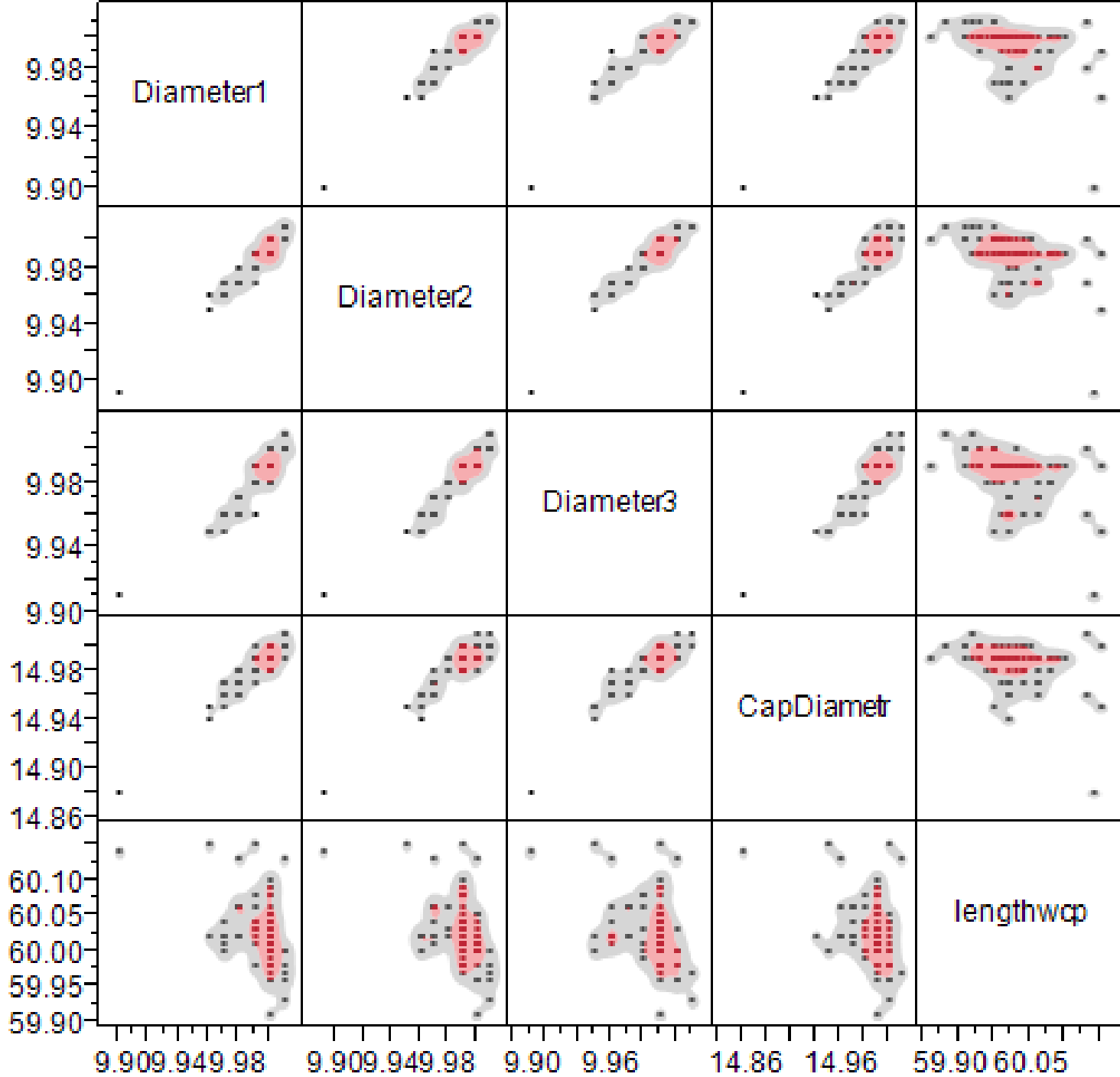


2



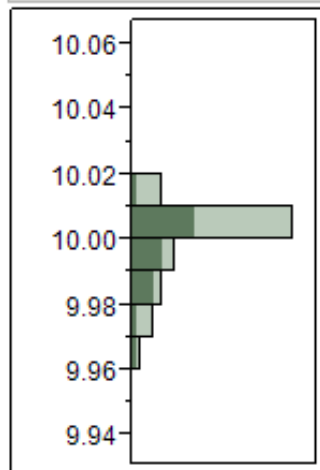


2

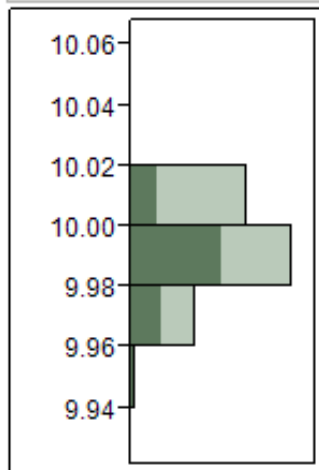


## Distributions

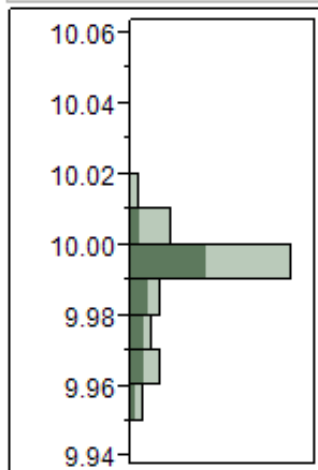
Diameter1



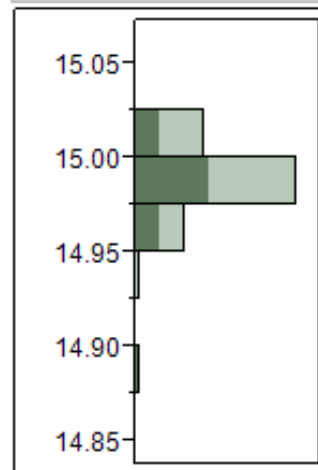
Diameter2



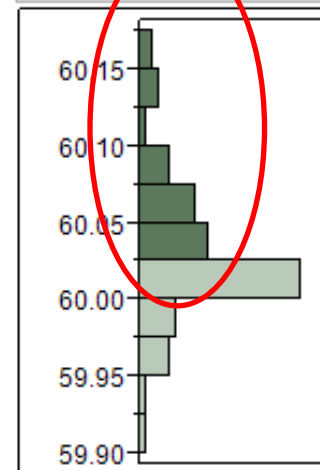
Diameter3



CapDiametr



lengthwcp



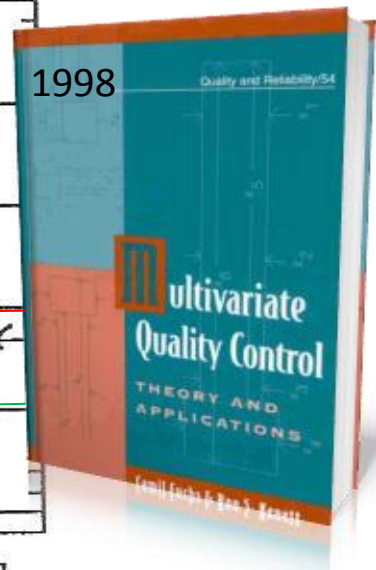
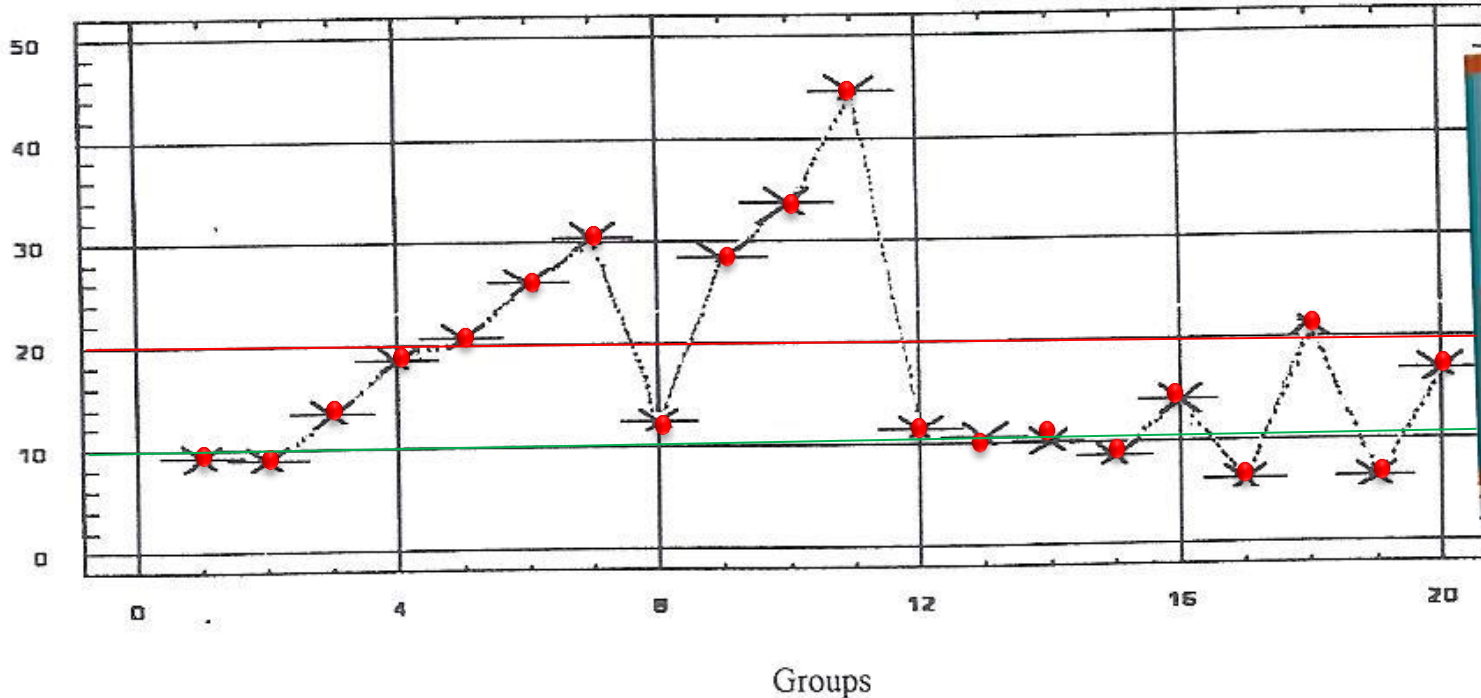
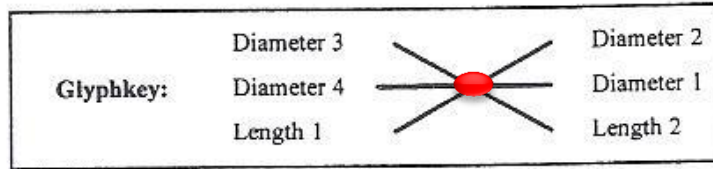


Figure 11.2: Starplot of the 6 variables from Case Study 1.

## Starplots

The starplot was apparently first developed at the SCS corporation as an enhancement to the Multivariate Control Charts available through

# ( $ARL_0$ and $ARL_1$ ) or (PFA and CED)?

## Special Issue Article

Quality and  
Reliability  
Engineering  
International

(wileyonlinelibrary.com) DOI: 10.1002/qre.1436

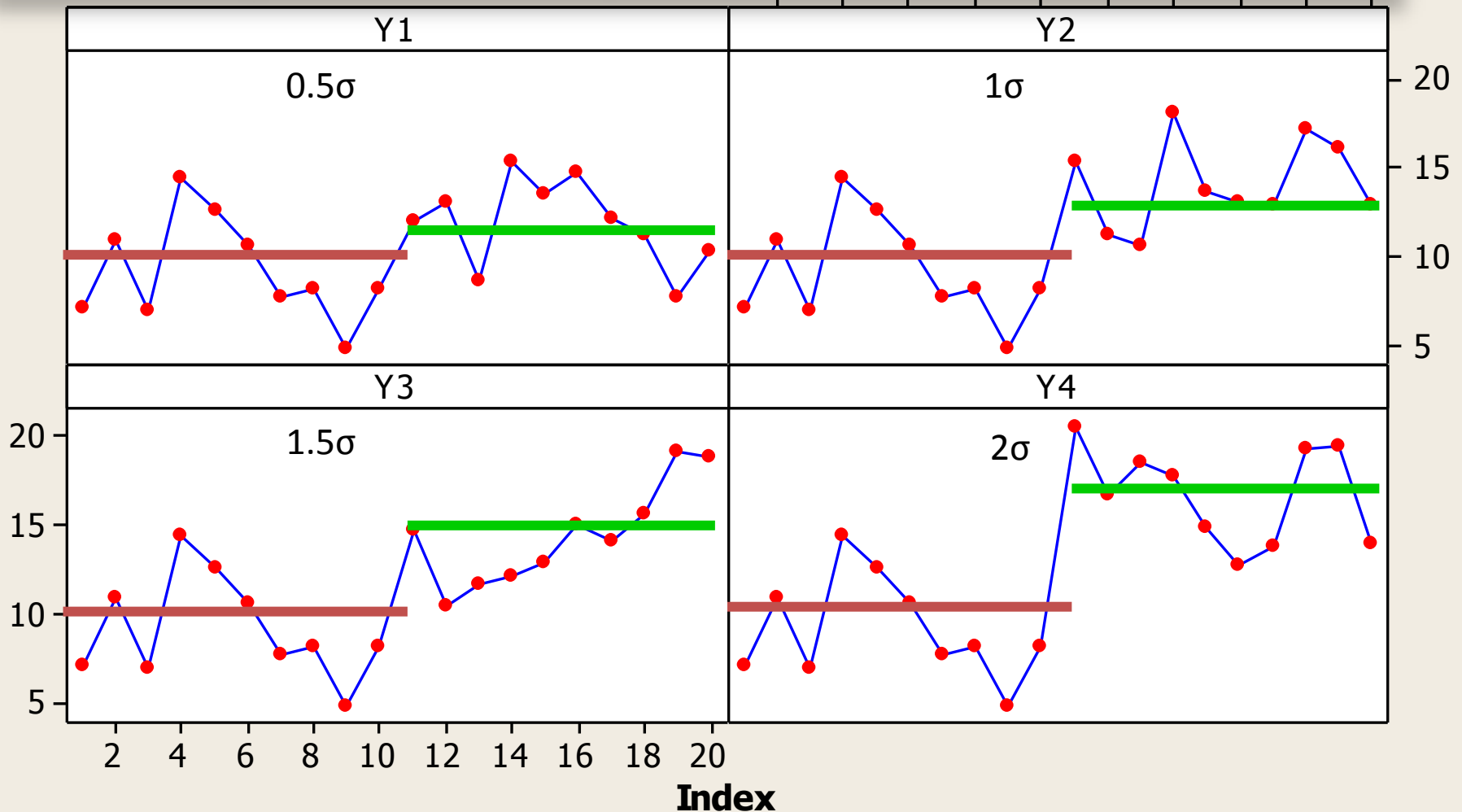
Published online in Wiley Online Library

# On Assessing the Performance of Sequential Procedures for Detecting a Change

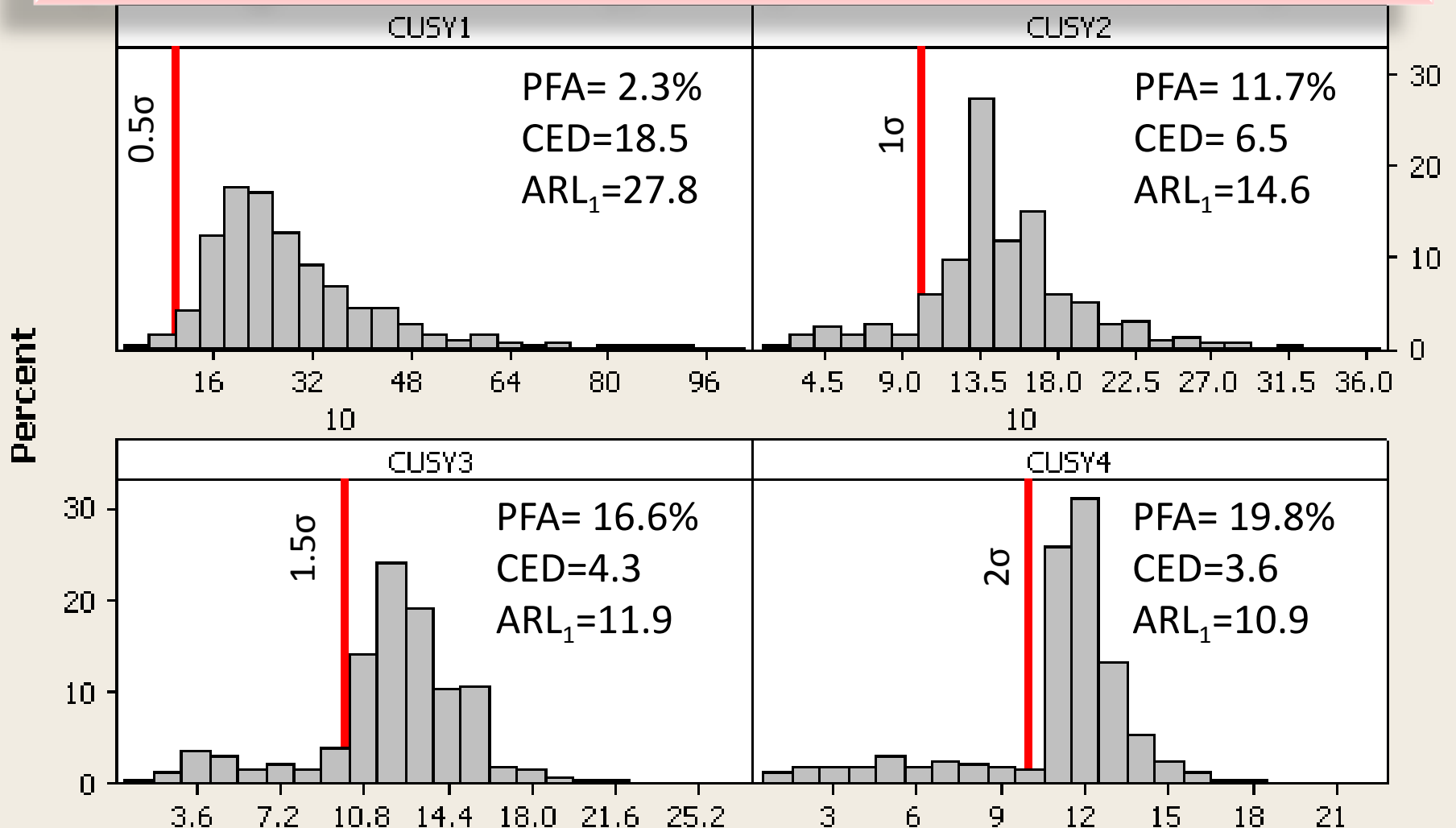
Ron S. Kenett<sup>a,b\*†</sup> and Moshe Pollak<sup>c</sup>

The literature on statistical process control has focused mostly on the average run length (ARL) to an alarm, as a performance criterion of sequential schemes. When the process is in control, this is the ARL to false alarm, generally denoted by  $ARL_0$ , and represents the in-control operating characteristic of the procedure. The ARL from the occurrence of a change to its detection represents an out-of-control operating characteristic and is typically embodied by  $ARL_1$ , the ARL to detection assuming that the change occurs at the very start of surveillance. However, these indices do not tell the whole story, and at times they are not defined well by a single number. We review the role of various operating characteristics in assessing performance of sequential procedures in comparison with  $ARL_0$  and  $ARL_1$ . Copyright © 2012 John Wiley & Sons, Ltd.

# ( $ARL_0$ and $ARL_1$ ) or (PFA and CED)?



# ( $ARL_0$ and $ARL_1$ ) or (PFA and CED)?



# Are we making an impact?

## Practical Statistical Efficiency (PSE)

$$\text{PSE} = \mathbf{E}\{\mathbf{R}\} \times \mathbf{T}\{\mathbf{I}\} \times \mathbf{P}\{\mathbf{I}\} \times \mathbf{V}\{\mathbf{PS}\} \times \mathbf{P}\{\mathbf{S}\} \times \mathbf{V}\{\mathbf{P}\} \times \mathbf{V}\{\mathbf{M}\} \times \mathbf{V}\{\mathbf{D}\}$$

$V\{D\}$  = value of the data actually collected

Kenett, Coleman,  
Stewardson

$V\{M\}$  = value of the statistical method employed

E.J.G.  
Pitman

$V\{P\}$  = value of the problem to be solved

$P\{S\}$  = probability that the problem actually gets solved

$V\{PS\}$  = value of the problem being solved

$P\{I\}$  = probability the solution is actually implemented

$T\{I\}$  = time the solution stays implemented

$E\{R\}$  = expected number of replications

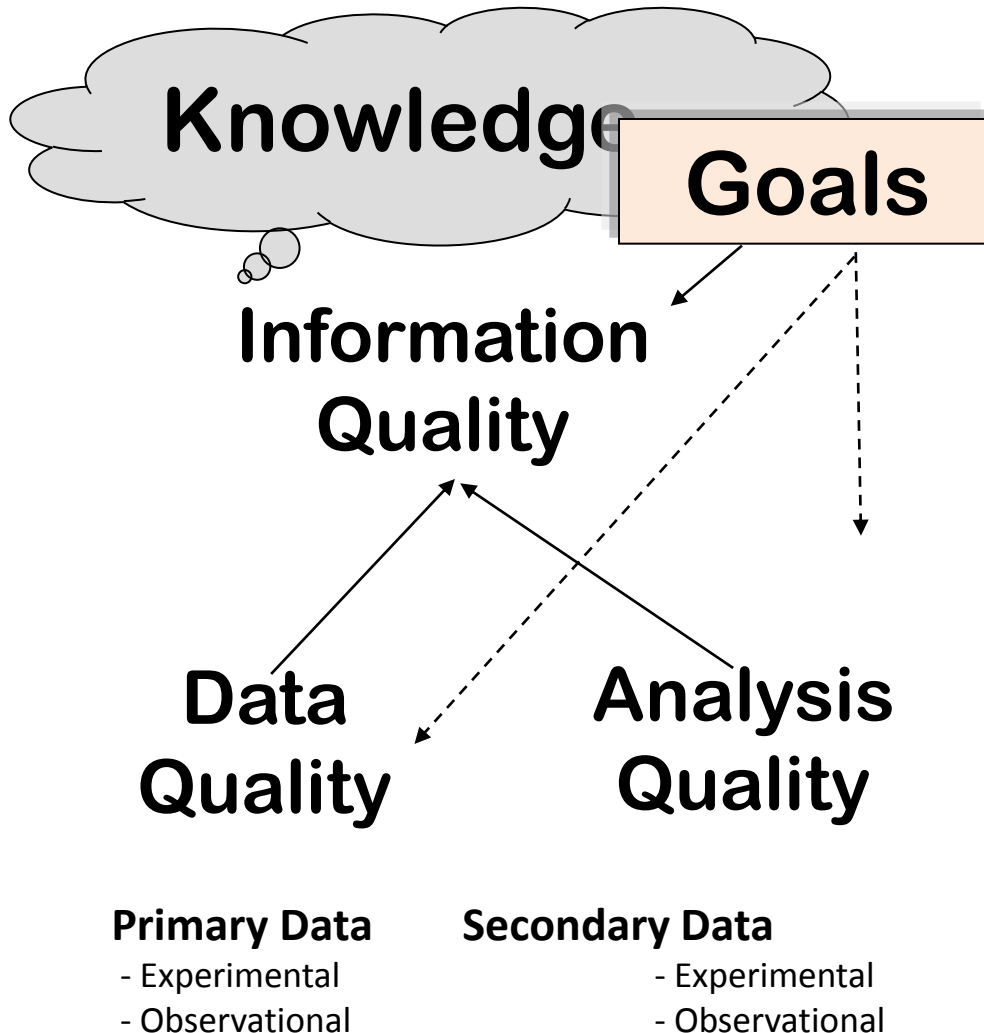
C. Eisenhart

B. Hoadley

A.B. Godfrey

# Are we generating knowledge?

## Information Quality (InfoQ)



$$InfoQ(f, X, g) = U(f(X|g))$$

<i>g</i>	A specific analysis goal	
<i>X</i>	The available dataset	
<i>f</i>	An empirical analysis method	
<i>U</i>	A utility measure	<b>What</b>

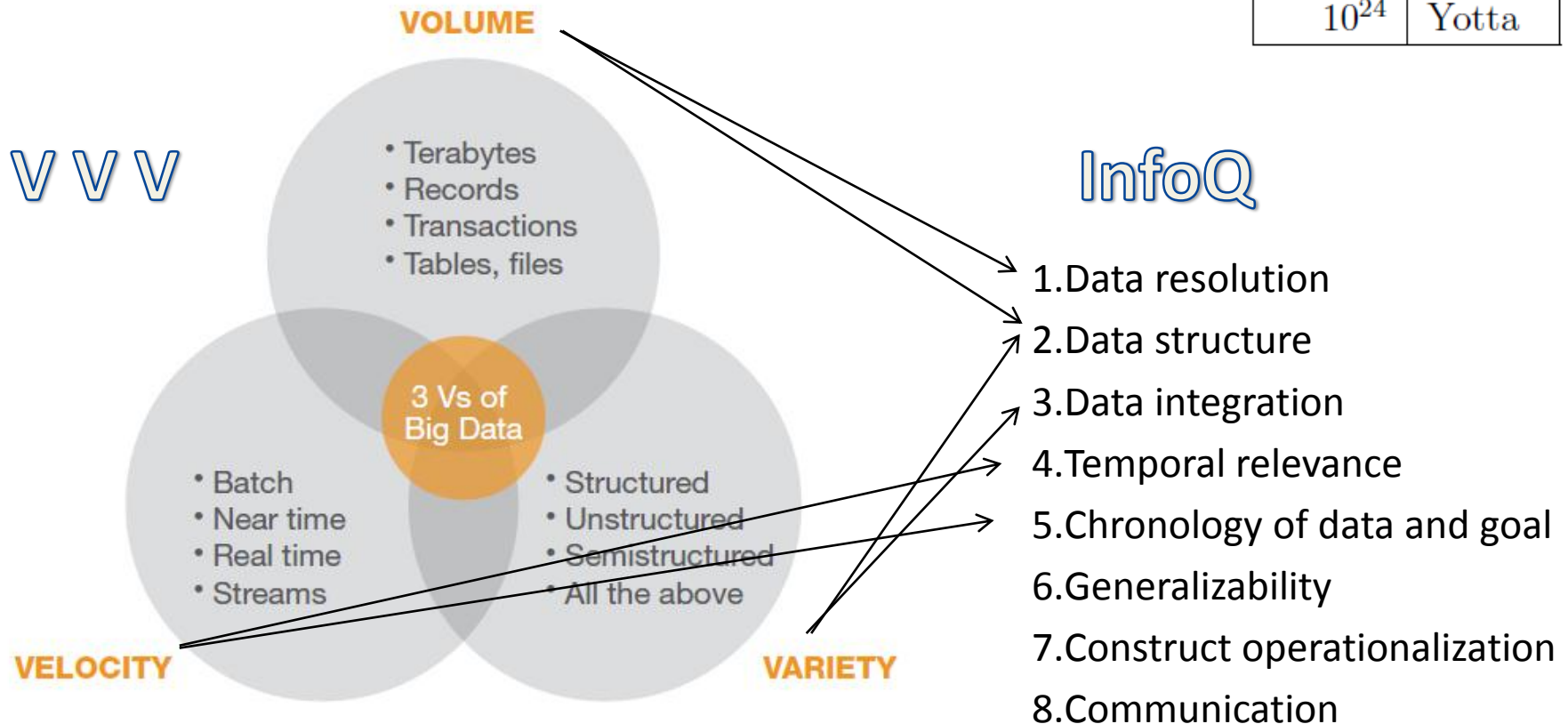
1. Data resolution
  2. Data structure
  3. Data integration
  4. Temporal relevance
  5. Chronology of data and goal
  6. Generalizability
  7. Construct operationalization
  8. Communication
- How**



# Big Data

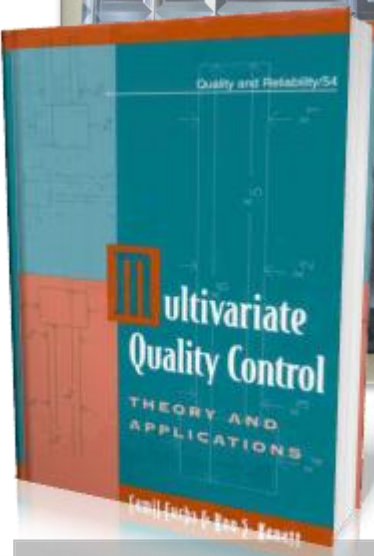
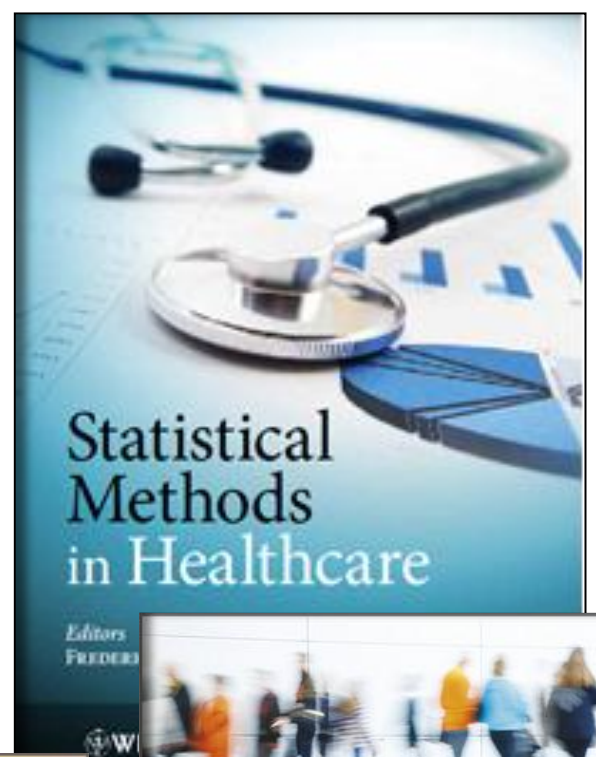
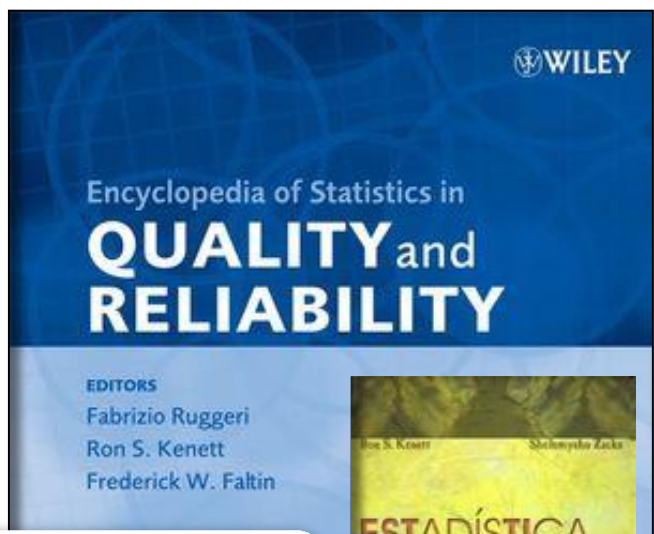
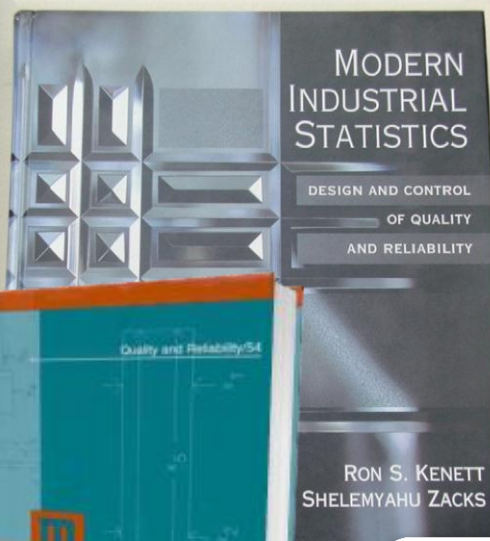
## VVV and InfoQ

Power	Prefix
$10^9$	Giga
$10^{12}$	Tera
$10^{15}$	Peta
$10^{18}$	Exa
$10^{21}$	Zetta
$10^{24}$	Yotta

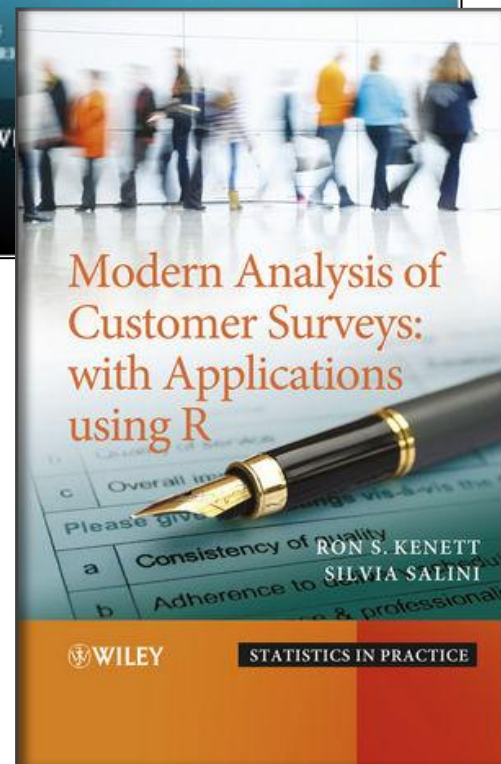
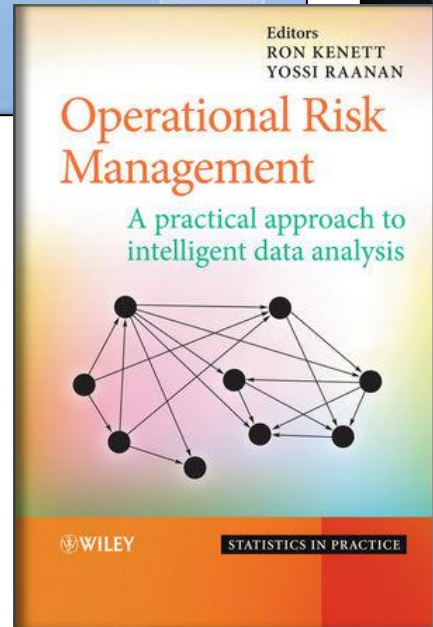
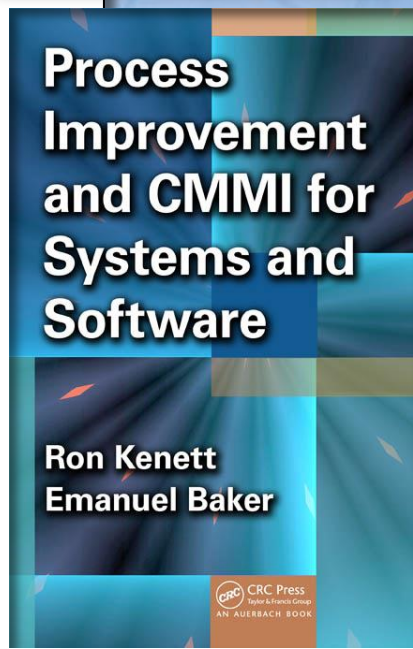
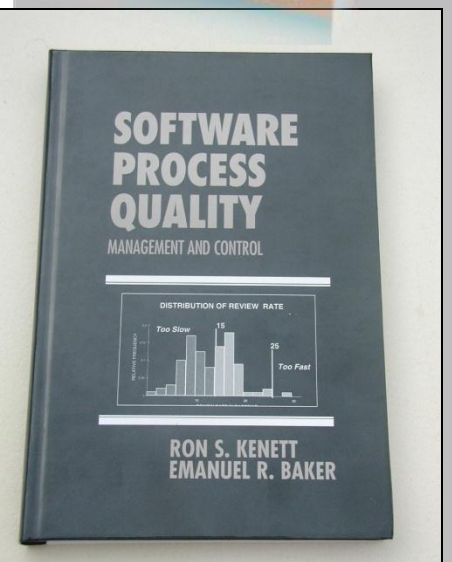
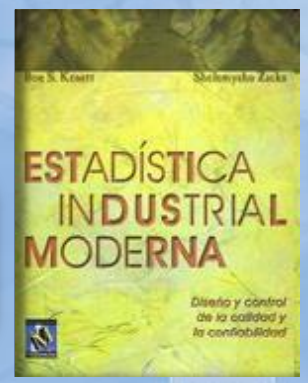


# Dimensions of Information Quality (InfoQ)

1. Data resolution — Mallows
2. Data structure — Hand
3. Data integration — Huber
4. Temporal relevance — Cox
5. Chronology of data and goal — Juran
6. Generalizability — Box
7. Construct operationalization — Deming
8. Communication — Greenfield



Modern Industrial Statistics with R, MINITAB and JMP



איגוד ישראלי  
לאבטחת איכות



# מדריך לשיטות מתקדמות בהנדסת איכות

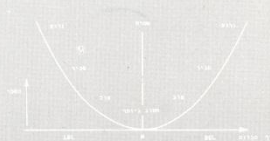
רון קנת ודוד שטיינברג



אורגון  
שרותי כתיבה טכנית

מהדורה 1.0

מאי 1989



# RSS Greenfield Medalists

Year	Name	Year	Name
1991	D Price	1992	T P Davis & D M Grove
1995	A Racine-Poon	1998	R Caulcutt & M Gerson
2000	L Furlong	2005	Susan Lewis
2007	S J Morrison	2010	D Montgomery



Imperial College  
London

מכון ויצמן למדע  
WEIZMANN INSTITUTE OF SCIENCE



UNIVERSITÀ  
DEGLI STUDI  
DI TORINO  
ALMA UNIVERSITAS  
TAURINENSIS



BINGHAMTON  
UNIVERSITY  
STATE UNIVERSITY OF NEW YORK



# CREATE IMPACT

# GENERATE KNOWLEDGE

Problem elicitation



RLD



CED



PSE



InfoQ

