

# Reliability and agreement

Cristina Costa Santos

MEDCIDS - Departamento Medicina da Comunidade, Informação e Decisão em Saúde  
CINTESIS - Centro de Investigação em Tecnologias e Serviços de Saúde  
Faculdade de Medicina da Universidade do Porto



► The quality of the measurements taken by health professionals or by measurement devices is fundamental not only for clinical care but also for research



- ▶ Measurement of variables always implies some degree of error.

- ▶ When an observer takes a measurement, the value obtained depends on several things such us:
  - ▶ the skills of the observer,
  - ▶ observer experience,
  - ▶ the measurement instrument,
  - ▶ observer's expectations
  - ▶ ...
- ▶ Also, natural continuous variation in a biological quantity can be presente.

- ▶ When natural continuous variation in a biological quantity is present, it is outside the control of the observer.
- ▶ It is, however, possible to minimize the observer variability by:
  - ▶ training of observers,
  - ▶ use of guidelines
  - ▶ automation, ...

# Reliability and agreement studies

before

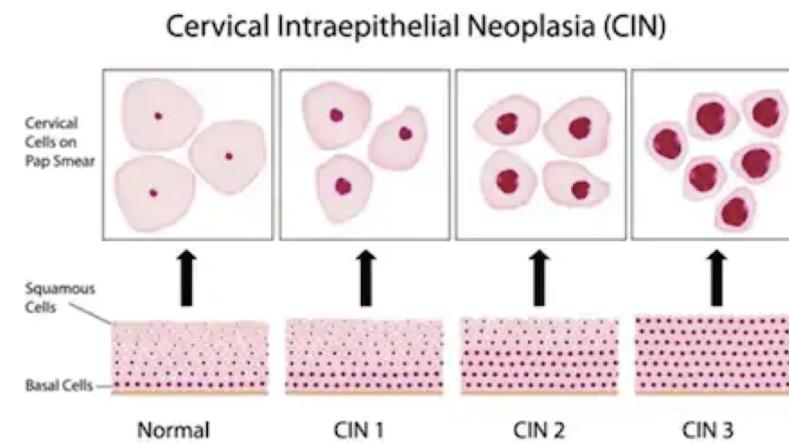
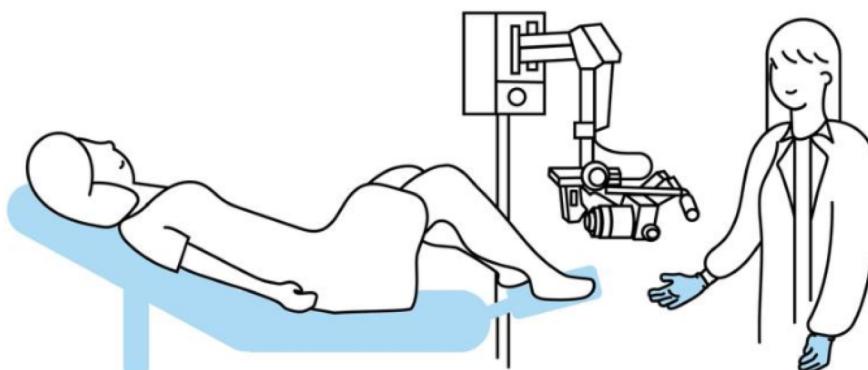
Validity studies / RCT / ...

Agreement studies are very importante

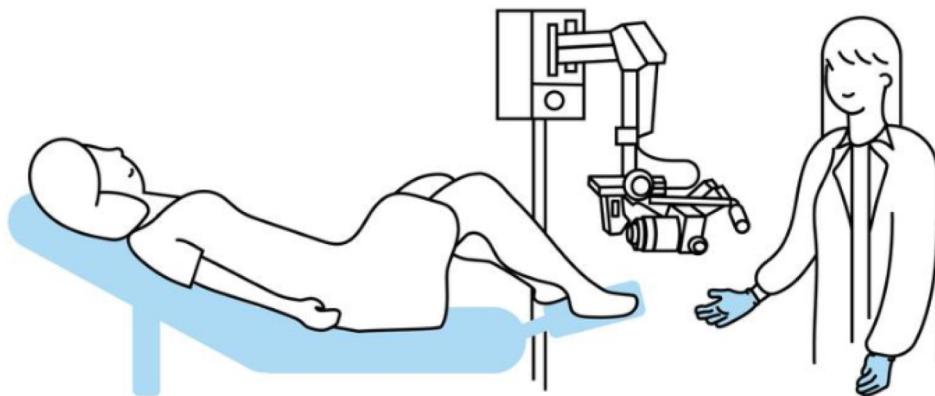
however

neglected in medical literature

► What is, for example, the sensitivity and specificity of colposcopy and of cervical cytology in detection of Cervical Neoplasia?

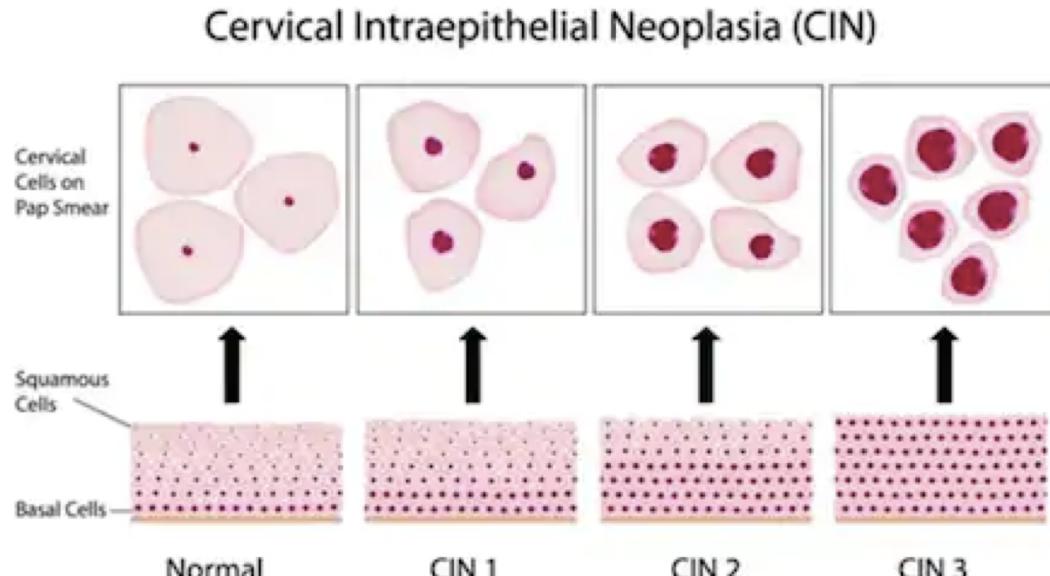


- Mitchell et al. found a variation in sensitivity and specificity of colposcopy ranging from 87% to 100% and 26 to 87%, respectively, in nine published studies using similar methods.

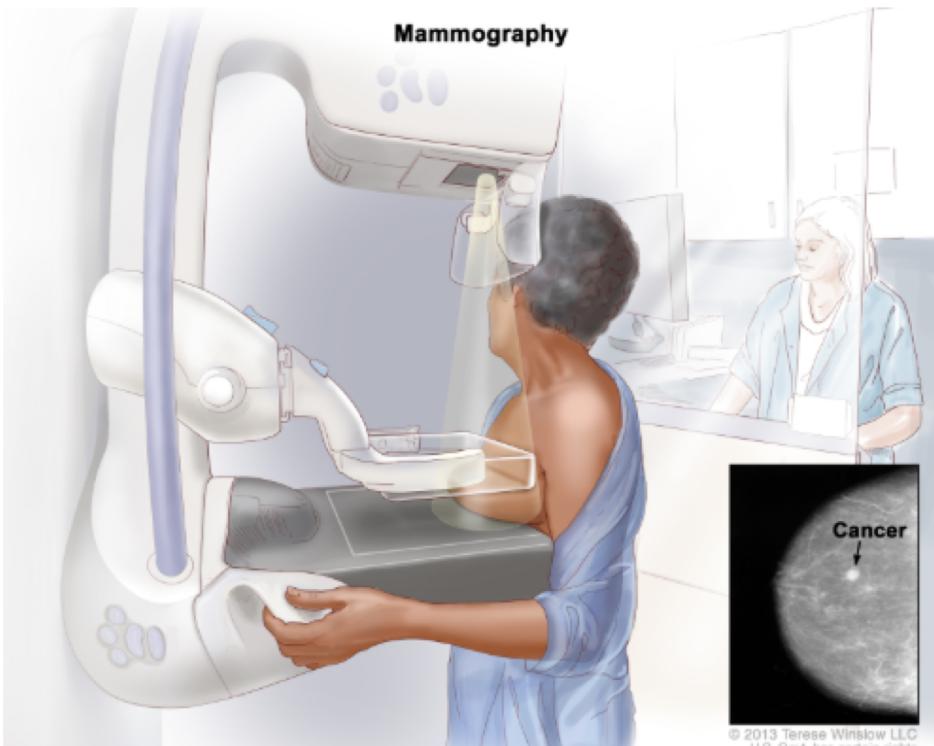


*Mitchell et al.  
Flourescence spectroscopy for diagnosis of squamous intraepithelial  
lesions of the cervix. Obstet Gynecol 1999;93:462- 470.*

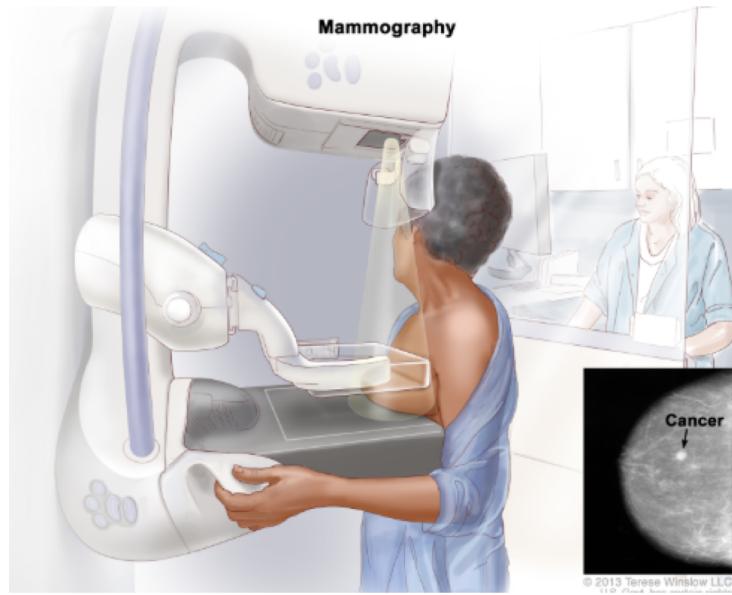
- The same authors found a variation in sensitivity and specificity of cervical cytology from 0 to 93% and 0 to 100%, respectively, in 26 published studies using similar methods.



► And, what are the relative risks of poor outcomes generally supposed to be avoided by interventions based on cardiotocography and mammography?



► Are they really close to 1, as suggested by most metaanalyses or are they just the mean result of discrepant results from the different randomised clinical trial included in meta-analyses?

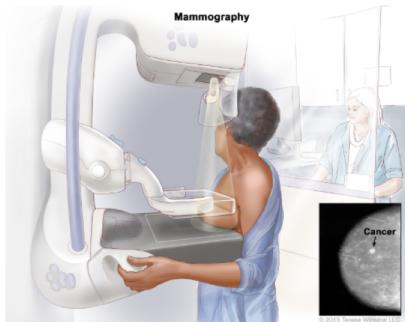
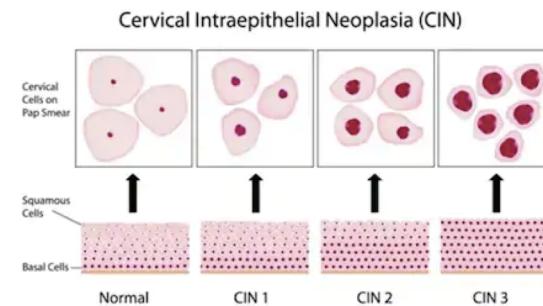
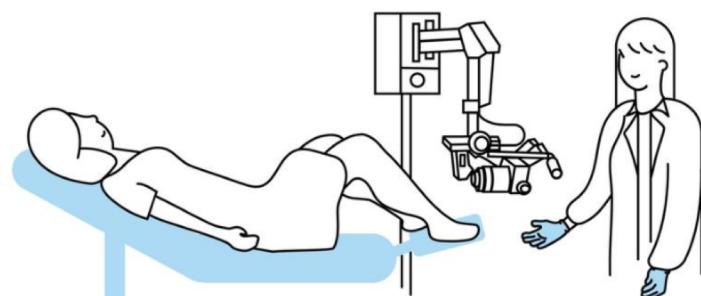


Thacker SB, Stroup DF, Peterson HB. Efficacy and safety of intrapartum electronic fetal monitoring: an update. *Obstet Gynecol* 1995;86:613- 620. 8. Olsen O, Gotzsche PC. Cochrane review on screening for breast cancer with mammography. *Lancet* 2001;358(9290):1340- 1342.

In fact, several studies revealed that all these widely used methods

(colposcopy, cervical cytology, cardiotocography and mammography)

have very low inter-observer agreement among among the doctors who perform them or who analyse them.



- The terms “reliability” and “agreement” are often used interchangeably.

However, the two concepts are conceptually distinct

- ▶ **Reliability** can be defined as the ability of a measurement to differentiate between subjects.
- ▶ **Agreement** is the degree to which scores or ratings are identical.

- ▶ Both concepts are important, because they provide information about the quality of measurements.
- ▶ The study designs for examining the two concepts are similar.

# Guidelines for Reporting Reliability and Agreement Studies (GRRAS) were proposed

Jan Kottner<sup>a,\*</sup>, Laurent Audigé<sup>b</sup>, Stig Brorson<sup>c</sup>, Allan Donner<sup>d</sup>, Byron J. Gajewski<sup>e</sup>,  
Asbjørn Hróbjartsson<sup>f</sup>, Chris Roberts<sup>g</sup>, Mohamed Shoukri<sup>h</sup>, David L. Streiner<sup>i</sup>

<sup>a</sup>*Department of Nursing Science, Centre for Humanities and Health Sciences, Charité-Universitätsmedizin Berlin, Berlin, Germany*

<sup>b</sup>*AO Clinical Investigation and Documentation, Dübendorf, Switzerland*

<sup>c</sup>*Department of Orthopaedic Surgery, Herlev University Hospital, Herlev, Denmark*

<sup>d</sup>*Department of Epidemiology and Biostatistics, Schulich School of Medicine and Dentistry, The University of Western Ontario, London, Ontario, Canada*

<sup>e</sup>*Department of Biostatistics, University of Kansas Schools of Medicine & Nursing, Kansas City, KS, USA*

<sup>f</sup>*The Nordic Cochrane Centre, Rigshospitalet, Copenhagen, Denmark*

<sup>g</sup>*School of Community Based Medicine, The University of Manchester, Manchester, UK*

<sup>h</sup>*Department of Biostatistics, Epidemiology and Scientific Computing, King Faisal Specialist Hospital and Research Center,  
The University of Western Ontario, London, Ontario, Canada*

<sup>i</sup>*Department of Psychiatry, University of Toronto, Toronto, Ontario, Canada*

# Measures of agreement

Statistical methods for analyzing interrater/intrarater reliability and agreement studies

Level of measurement	Reliability measures	Agreement measures
Nominal	Kappa statistics	Proportions of agreement Proportions of specific agreement
Ordinal	Ranked intraclass correlation Matrix of kappa coefficients Weighted kappa	Proportions of agreement Proportions of specific agreement
Continuous	Intraclass correlation coefficients	Proportions of agreement (ranges) Proportions of specific agreement (ranges) Standard errors of measurement Coefficients of variation Bland–Altman plots and limits of agreement

# Cohen's Kappa coefficient

Two raters are asked to classify objects into categories 1 and 2.  
The table below contains cell probabilities for a 2 by 2 table.

		Rater #1		Total
		1	2	
Rater #2	1	$p_{11}$	$p_{12}$	$p_{1.}$
	2	$p_{21}$	$p_{22}$	$p_{2.}$
Total		$p_{.1}$	$p_{.2}$	1

$$\text{Proportions of agreement (Po)} = P_{11} + P_{22}$$

$$\text{Proportions of expected by chance (Pe)} = P_{.1} P_{1.} + P_{.2} P_{2.}$$

$$K = \frac{Po - Pe}{1 - Pe}$$

Chance-Corrected Agreement? Or measure of reliability?

Cohen J.  
A coefficient of agreement for  
nominal scales.  
Educational and  
Psychological Measurement,  
1960:37-46, 1960.

# High agreement but low kappa

		Observer 1			
		yes	no	total	
Observer 2	yes	40	5	45	
	no	3	2	5	
	total	43	7	50	

$$PA = 42/50 = 0.84$$

$$Pe = (43/50)(45/50) + (7/50)(5/50) = 0.79$$

$$K = (0.84 - 0.79) / (1 - 0.79) = 0.24$$

# Specific agreement

Summary of binary ratings by two raters

		Rater 2		total
Rater 1	+	a	b	a + b
	-	c	d	c + d
total		a + c	b + d	N

$$\text{Proportion of overall agreement} = (a+d)/N$$

$$PA = \frac{2a}{2a + b + c}; \quad NA = \frac{2d}{2d + b + c}.$$

PA estimates the conditional probability, given that one of the raters, randomly selected, makes a positive rating, the other rater will also do so

NA estimates the conditional probability, given that one of the raters, randomly selected, makes a negative rating, the other rater will also do so

The values a, b, c and d here denote the observed frequencies for each possible combination of ratings by Rater 1 and Rater 2.

Cicchetti DV.  
Feinstein AR.  
*High agreement but low kappa: II. Resolving the paradoxes.*  
*Journal of Clinical Epidemiology*, 1990,  
43, 551-558.

Spitzer R, Fleiss J.  
*A re-analysis of the reliability of psychiatric diagnosis.*  
*British Journal on Psychiatry*, 1974,  
341-47.

# High agreement but low kappa

		Observer 1			
		yes	no	total	
Observer 2	yes	40	5	45	
	no	3	2	5	
	total	43	7	50	

$$PA = 42/50 = 0.84$$

$$PA \text{ yes} = 40 \times 2 / (40 \times 2 + 3 + 5) = 0.92$$

$$K = (0.84 - 0.79) / (1 - 0.79) = 0.24$$

$$PA \text{ no} = 2 \times 2 / (2 \times 2 + 3 + 5) = 0.33$$

# Intraclass Correlation Coefficient

- ▶ The Intraclass Correlation (ICC) assesses rating reliability by comparing the variability of different ratings of the same subject to the total variation across all ratings and all subjects.

*Shrout PE, Fleiss JL. Intraclass correlations: Uses in assessing rater reliability. Psychol Bulletin 1979;86:420-427.*

# Information Based Measure of Disagreement

- ▶ The sum over all logarithms of possible outcomes of the variable is a valid measure of the amount of information, or uncertainty, contained in a variable.

Consider that we aim to measure disagreement between measurements obtained by Observer Y and Observer X .

The disagreement between Y and X is related to the differences between them. So, we consider

$$\sum_{i=1}^n \log_2 |a_i - b_i|$$

the amount of information contained in the differences between observers.

By adding 1 to the differences, we avoid the behavior of the logarithmic function between 0 and 1.

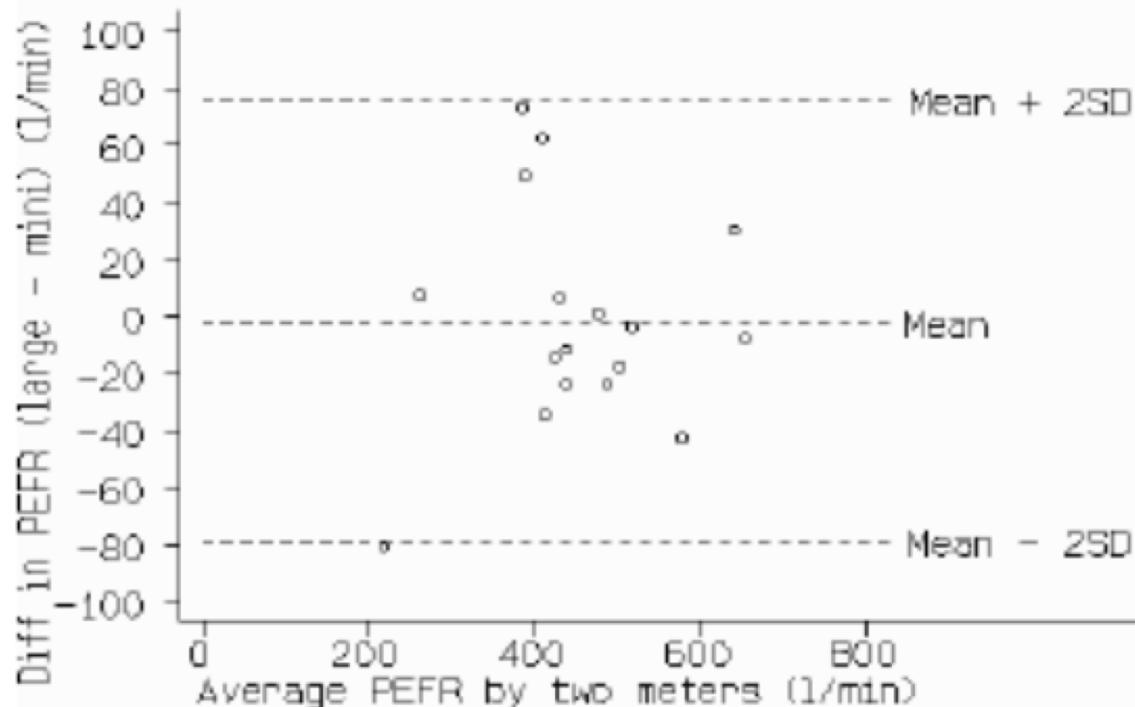
To get a value between 0 and 1 we normalize the amount of information contained in the differences to obtain the following measure of information-based measure of disagreement (IBMD):

$$\frac{1}{n} \sum_{i=1}^n \log_2 \left( \frac{|a_i - b_i|}{\max\{a_i, b_i\}} + 1 \right)$$

0 - no disagreement  
Tends to 1 (total disagreement)

# Bland and Altman limits of agreement

Difference against mean



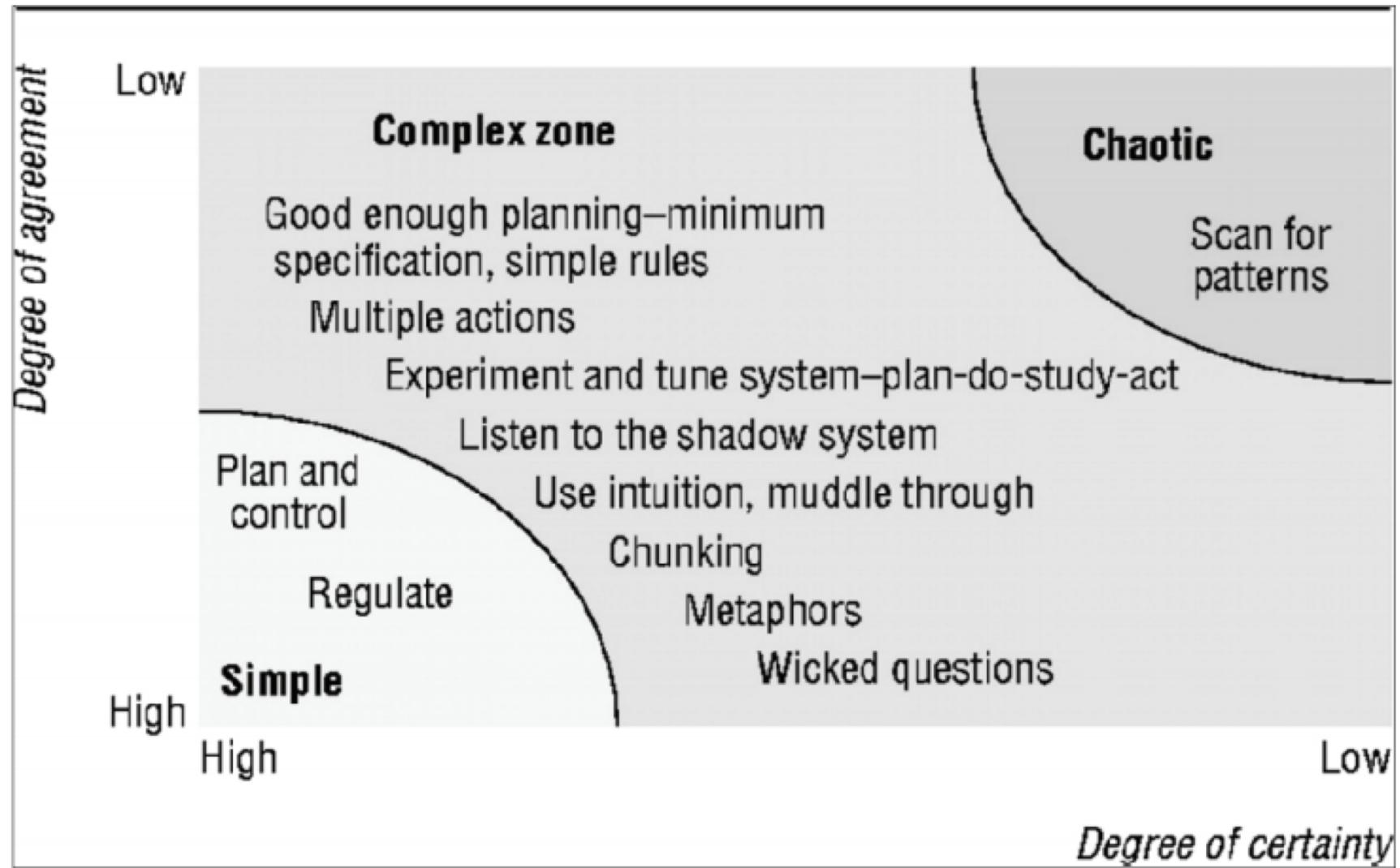
Provided differences within *mean differences*  $\pm 2SD$  would not be clinically important, we could use the two measurement methods interchangeably.

IS A CLINICAL (NOT  
STATISTICAL)  
INTERPRETATION



## Cardiotocography CTG guidelines:

- FIGO, International Federation of Gynecology and Obstetrics;
- ACOG, the American College of Obstetricians and Gynecologists,
- NICE, National Institute for Health and Care Excellence.



Reprinted from BMJ, 323, Suneet Wilson T, Holt T, Greenhalgh T, Complexity science: complexity and clinical care, 685-6888, Copyright (2001), with permission from BMJ Publishing Group Ltd.



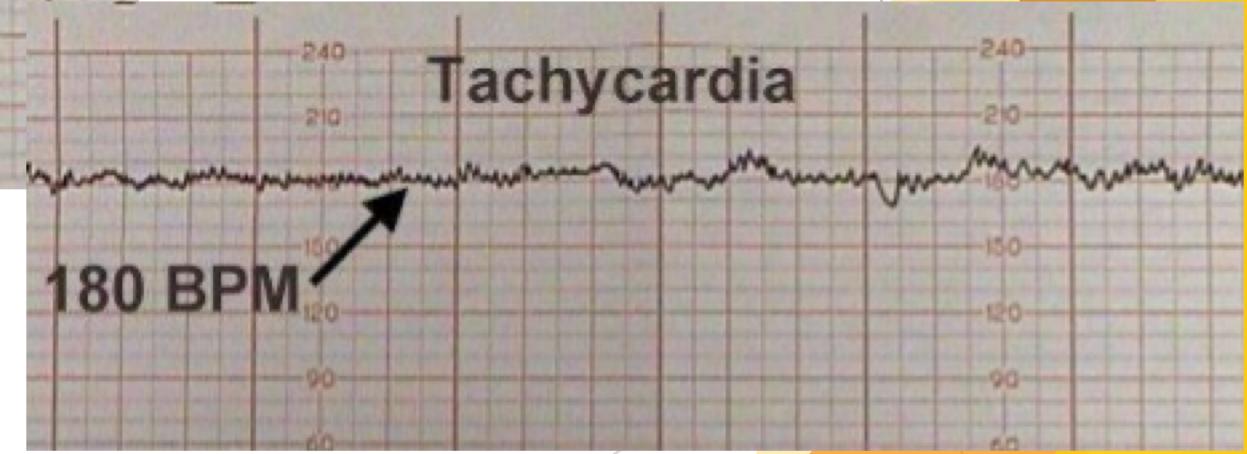
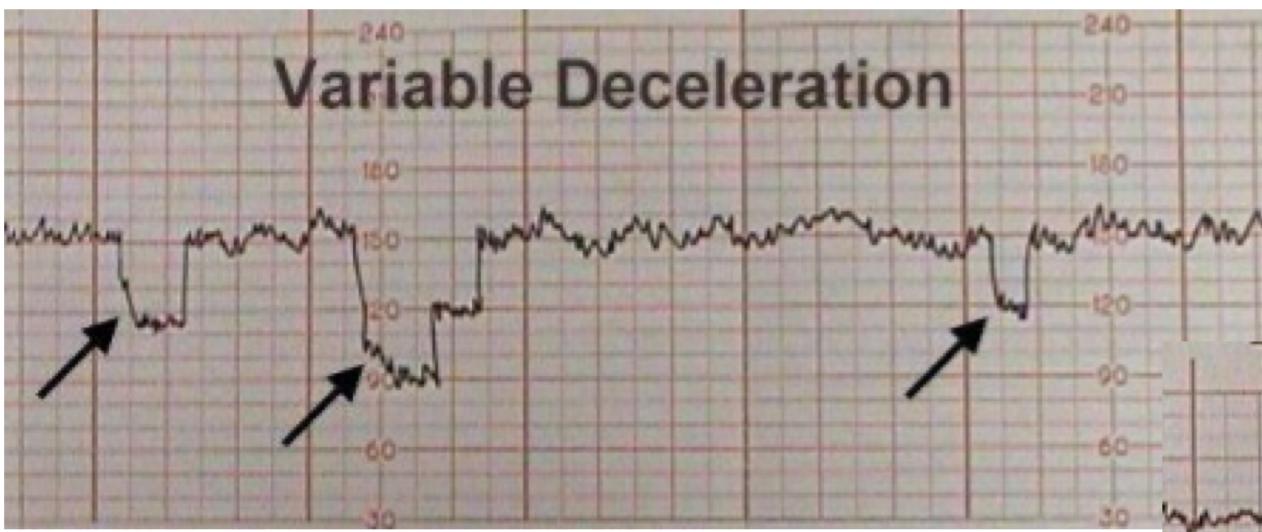
**AOGS** ORIGINAL RESEARCH ARTICLE

## **Agreement and accuracy using the FIGO, ACOG and NICE cardiotocography interpretation guidelines**

SUSANA SANTO<sup>1</sup>, DIOGO AYRES-DE-CAMPOS<sup>2</sup>, CRISTINA COSTA-SANTOS<sup>3</sup>, WILLIAM SCHNETTLER<sup>4</sup>, AUSTIN UGWUMADU<sup>5</sup> & LUÍS M. DA GRAÇA<sup>1</sup> FOR THE FM-COMPARE COLLABORATION\*

<sup>1</sup>*Department of Obstetrics and Gynecology, Santa Maria Hospital, Faculty of Medicine of Lisbon University, Lisbon,*

<sup>2</sup>*Department of Obstetrics and Gynecology, Medical School, University of Porto, S. João Hospital, Institute of Biomedical Engineering, Porto,* <sup>3</sup>*Department of Medical Informatics, Medical School, University of Porto, Porto, Portugal,* <sup>4</sup>*Center for Maternal Cardiac Care, TriHealth, Good Samaritan Hospital, Cincinnati, OH, USA, and* <sup>5</sup>*Department of Obstetrics & Gynecology, St George's Hospital, University of London, London, UK*



**Table 2.** Comparison of cardiotocography (CTG) classification criteria in the FIGO, NICE and ACOG guidelines.

FIGO	NICE	ACOG
Normal pattern <ul style="list-style-type: none"><li>• Baseline heart rate between 110 and 150 bpm</li><li>• Amplitude of heart rate variability between 5 and 25 bpm</li></ul>	Normal (a CTG where all of the following four reassuring features are present) <ul style="list-style-type: none"><li>• Baseline rate: 110–160 bpm</li><li>• Variability: <math>\geq 5</math> bpm</li><li>• No decelerations</li><li>• Accelerations: present</li></ul>	Category I (category I FHR tracings include all of the following) <ul style="list-style-type: none"><li>• Baseline rate: 110–160 bpm</li><li>• Baseline variability: 6–25 bpm</li><li>• Late or variable decelerations: absent</li><li>• Early decelerations: present or absent</li><li>• Accelerations: present or absent</li></ul>
Suspicious pattern <ul style="list-style-type: none"><li>• Baseline heart rate between 150 and 170 bpm or between 100 and 110 bpm</li><li>• Amplitude of variability between 5 and 10 bpm for more than 40 min</li><li>• Increased variability above 25 bpm</li><li>• Variable decelerations</li></ul>	Suspicious (a CTG where one of the following features is present and all others fall into the reassuring category) <ul style="list-style-type: none"><li>• Baseline rate<ul style="list-style-type: none"><li>- 100–109 bpm</li><li>- 161–180 bpm</li></ul></li><li>• Baseline variability<ul style="list-style-type: none"><li>- <math>&lt;5</math> bpm for 40–90 min</li></ul></li><li>• Decelerations<ul style="list-style-type: none"><li>- Typical variable decelerations with <math>&gt;50\%</math> of contractions occurring for <math>&gt;90</math> min</li><li>- Single prolonged deceleration for up to 3 min</li></ul></li></ul>	Category II (Category II FHR tracings include all FHR tracings not categorised as Category I or Category III. Examples of Category II FHR tracings include any of the following) <ul style="list-style-type: none"><li>• Baseline rate<ul style="list-style-type: none"><li>- Bradycardia not accompanied by absent baseline variability</li><li>- Tachycardia</li></ul></li><li>• Baseline variability<ul style="list-style-type: none"><li>- Minimal variability</li><li>- Absent variability with no recurrent decelerations</li><li>- Marked variability</li></ul></li><li>• Accelerations</li></ul>

	FIGO 1987		ACOG 2010		NICE 2007	
	PA (95% CI)	$\kappa$ (95% CI)	PA (95% CI)	$\kappa$ (95% CI)	PA (95% CI)	$\kappa$ (95% CI)
FHR baseline	0.81 (0.78–0.85)	0.63 (0.57–0.70)	0.88 (0.84–0.91)	0.59 (0.49–0.69)	0.88 (0.85–0.91)	0.65 (0.58–0.72)
Normal	0.86 (0.83–0.89)		0.93 (0.90–0.95)		0.93 (0.91–0.96)	
Tachycardia	0.80 (0.73–0.85)		0.67 (0.57–0.77)		0.73 (0.66–0.80)	
Bradycardia	0.40 (0.25–0.54)		0.49 (0.07–0.71)		0.42 (0.00–1.00)	
Variability	0.83 (0.80–0.86)	0.51 (0.42–0.61)	0.85 (0.82–0.88)	0.49 (0.39–0.59)	0.83 (0.80–0.86)	0.38 (0.29–0.50)
Normal	0.89 (0.87–0.92)		0.91 (0.89–0.93)		0.90 (0.88–0.92)	
Abnormal	0.61 (0.51–0.69)		0.57 (0.45–0.66)		0.44 (0.32–0.53)	
Accelerations	0.67 (0.64–0.71)	0.34 (0.29–0.41)	0.67 (0.63–0.70)	0.34 (0.28–0.40)	0.71 (0.68–0.75)	0.41 (0.35–0.48)
Yes	0.73 (0.68–0.77)		0.72 (0.67–0.76)		0.61 (0.57–0.68)	
No	0.59 (0.54–0.64)		0.59 (0.53–0.64)		0.76 (0.72–0.80)	
Decelerations	0.92 (0.89–0.95)	0.53 (0.43–0.66)	0.85 (0.82–0.88)	0.28 (0.18–0.46)	0.89 (0.85–0.91)	0.47 (0.35–0.59)
Yes	0.96 (0.94–0.97)		0.92 (0.89–0.93)		0.94 (0.92–0.95)	
No	0.59 (0.48–0.69)		0.35 (0.23–0.45)		0.49 (0.36–0.59)	
Classification	0.64 (0.61–0.67)	0.37 (0.31–0.43)	0.73 (0.70–0.76)	0.15 (0.10–0.21)	0.55 (0.51–0.58)	0.33 (0.28–0.39)
Cat. I/Normal	0.54 (0.39–0.64)		0.26 (0.18–0.33)		0.55 (0.48–0.62)	
Cat. II/Suspicious	0.67 (0.62–0.70)		0.83 (0.81–0.86)		0.42 (0.38–0.47)	
Cat. III/Pathological	0.63 (0.57–0.68)		0.26 (0.18–0.34)		0.66 (0.59–0.71)	



“Fine, tell him he’s going to make it.  
We’ll just have to agree to disagree.”