

# Extraire les données d'un article en vue d'une méta-analyse: trucs et astuces statistiques



Combescure Christophe, chargé de cours

Centre de recherche clinique & Service d'épidémiologie clinique

# Intérêts d'une revue systématique

---

- Faire un état des lieux des connaissances scientifiques sur une certaine question de recherche
  - quelles recherches ont été menées
  - quelle qualité des études
- Proposer une synthèse de ces connaissances
  - quels ont été les résultats obtenus
- Combiner les résultats des études: méta-analyse
  - généralisabilité de l'effet estimé
  - gain de puissance et précision statistique
- Rendre accessible les résultats de la revue

# Étapes d'une revue systématique

---

## METHODES

Protocole

Recherche de la littérature

Sélection des études

Lecture critique

Extraction des données des études

Synthèse

Communication des résultats de la revue

# Ce que devrait rapporter un article



## CONSORT 2010 checklist of information to include when reporting a randomised trial\*

Section/Topic	Item No	Checklist item	Reported on page No
<b>Results</b>			
Numbers analysed	16	For each group, number of participants (denominator) included in each analysis and whether the analysis was by original assigned groups	
Outcomes and estimation	17a	For each primary and secondary outcome, results for each group, and the estimated effect size and its precision (such as 95% confidence interval)	

- Selon CONSORT 2010, un RCT devrait rapporter les outcomes de manière détaillée
- Mais:
  - certaines études incluses dans la revue peuvent être antérieures à CONSORT 2010
  - l'adhésion à CONSORT 2010 n'est pas de 100%
  - les outcomes secondaires peuvent être rapportés avec moins de soins que l'outcome principal

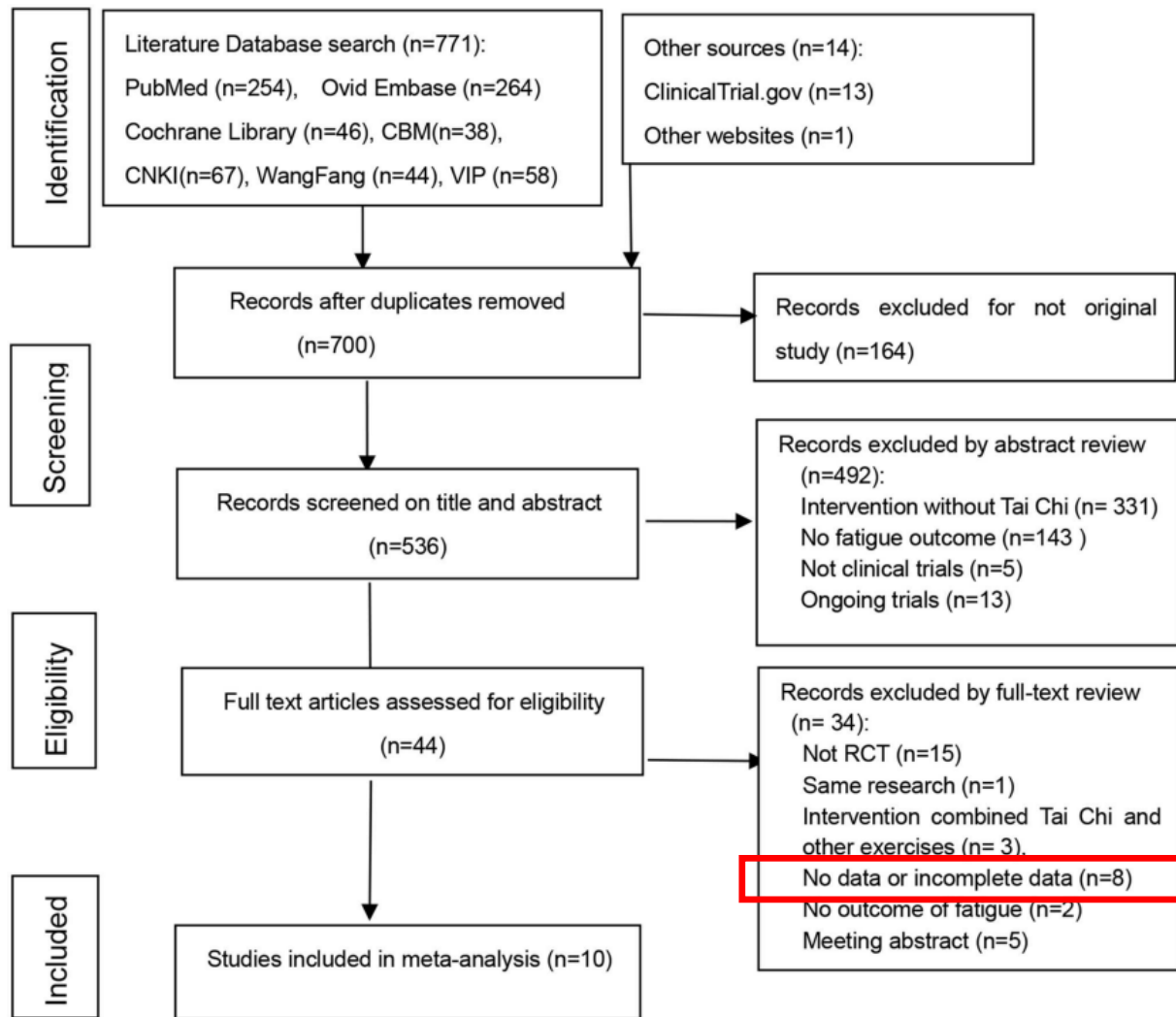


Fig 1. Flow diagram of study selection and identification.

There was insufficient data on fatigue (n = 8); 5) There was no outcome of fatigue (n = 2); 6) Conference abstracts (n = 5). Of the excluded 8 studies which there were insufficient data, two described fatigue outcome in words without data, one provided only median, four provided data that cannot be used to calculate SMD. One study we didn't get full-text. For these studies, we sent e-mails to the authors, two authors had replied to us. One author sent the full text to us that we didn't find before, but the data was useless. Both of the two authors said the original data has been destroyed. While others did not reply within three months.

# Conséquences potentielles des données incomplètes sur une méta-analyse

---

- Exclusion d'études de la méta-analyse (mais pas de la revue systématique)
- Diminution de la puissance/précision statistique de la méta-analyse
- Complique l'investigation d'un biais de publication et des sources d'hétérogénéité (on considère qu'il faut 10 études ou plus)
- Risque de sélection biaisée des études (moins de soins à rapporter un outcome secondaire non statistiquement significatif?)

# Solutions pour limiter le problème des données incomplètes

---

- Contacter les auteurs
- Accorder un peu plus d'attention à l'extraction des données
  - Parfois, les données rapportées ne sont pas directement exploitables pour la méta-analyse mais on peut en déduire l'information nécessaire à la méta-analyse
  - Ça ne coûte pas grand-chose
- En pratique, dans cette présentation:
  - Méthodes statistiques pour compléter les données d'un outcome continu
  - Extraction de données d'un graphique



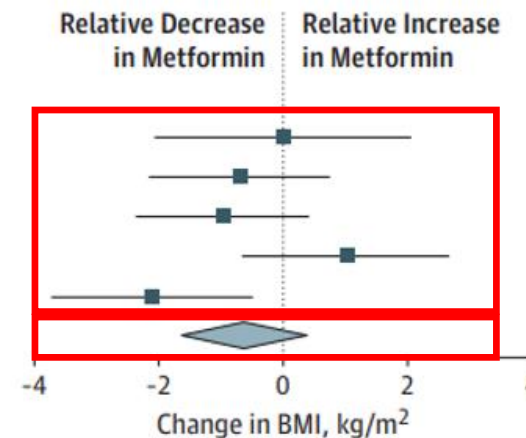
# Principe d'une méta-analyse d'outcomes continus

Données entrées par l'analyste

Etudes

Source	Metformin		Control		MD
	No.	Mean (SD)	No.	Mean (SD)	
Low cumulative metformin dose					
Evia-Viscarra et al, <sup>18</sup> 2012	12	-0.7 (2.5)	14	-0.7 (2.9)	0.0
Gómez-Díaz et al, <sup>6</sup> 2012	28	-2.2 (2.7)	24	-1.5 (2.6)	-0.7
Mauras et al, <sup>20</sup> 2012	23	-2.1 (2.6)	19	-1.2 (2.0)	-1.0
Casteels et al, <sup>21</sup> 2010	19	1.4 (2.8)	23	0.4 (2.7)	1.0
Burgert et al, <sup>5</sup> 2008	15	-0.9 (2.5)	14	1.2 (1.9)	-2.1
Subtotal	97		94		-0.6

Heterogeneity:  $\chi^2 = 7.52$ ;  $P = .11$ ;  $I^2 = 47\%$ , Overall effect:  $z = 1.22$ ;  $P = .22$



Représentation graphique des différences de moyenne avec l'intervalle de confiance à 95%

Effet commun  
Largeur du losange = intervalle de confiance à 95%

Habituellement, les données extraites des études pour chaque bras sont:

- le nombre d'observations
- la moyenne et l'écart type de l'outcome



# Principe d'une méta-analyse d'outcomes continus

L'effet commun est une moyenne pondérée des effets observés dans les études:

$$\text{Effet commun} = \frac{\sum_{k=1}^K \text{Poids}_k \times \text{Différence de moyenne dans l'étude } k}{\sum_{k=1}^K \text{Poids}_k}$$

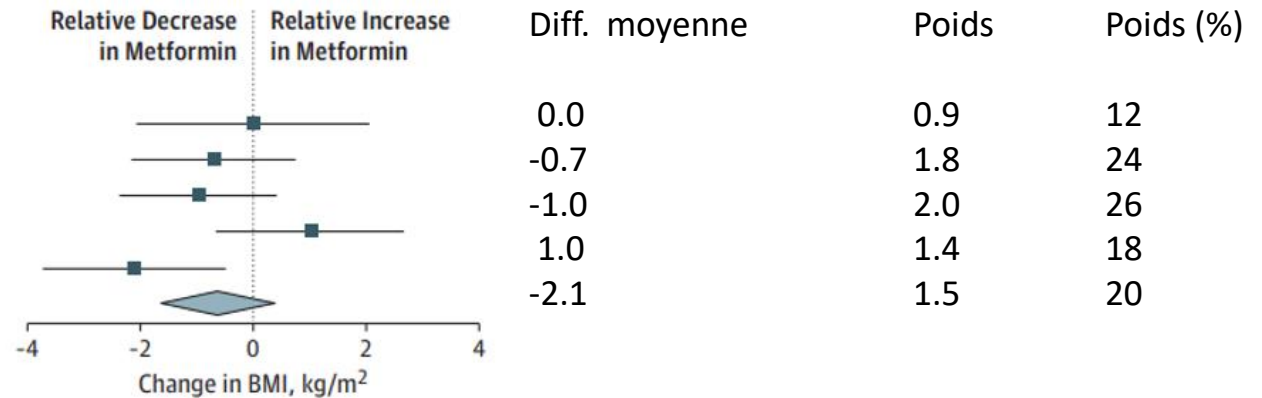
Le poids d'une étude dans la méta-analyse reflète la précision statistique de l'effet avec les données de cette étude:

- plus l'intervalle de confiance à 95% de la différence de moyenne est étroit, plus le poids est important
- le poids tend à augmenter lorsque la taille d'échantillon augmente

# Principe d'une méta-analyse d'outcomes continus

C Effect of metformin use on body mass index

Source	Metformin		Control		MD
	No.	Mean (SD)	No.	Mean (SD)	
Low cumulative metformin dose					
Evia-Viscarra et al, <sup>18</sup> 2012	12	-0.7 (2.5)	14	-0.7 (2.9)	0.0
Gómez-Díaz et al, <sup>6</sup> 2012	28	-2.2 (2.7)	24	-1.5 (2.6)	-0.7
Mauras et al, <sup>20</sup> 2012	23	-2.1 (2.6)	19	-1.2 (2.0)	-1.0
Casteels et al, <sup>21</sup> 2010	19	1.4 (2.8)	23	0.4 (2.7)	1.0
Burgert et al, <sup>5</sup> 2008	15	-0.9 (2.5)	14	1.2 (1.9)	-2.1
Subtotal	97		94		-0.6
Heterogeneity: $\chi^2 = 7.52$ ; $P = .11$ ; $I^2 = 47\%$ , Overall effect: $z = 1.22$ ; $P = .22$					



Différence de moyenne commune

$$\begin{aligned}
 &= [ 0.9 \times 0.0 + 1.8 \times (-0.7) + 2.0 \times (-1.0) + 1.4 \times 1.0 + 1.5 \times (-2.1) ] / 7.6 \\
 &= 0.20 \times 0.0 + 0.24 \times (-0.7) + 0.26 \times (-1.0) + 0.18 \times 1.0 + 0.20 \times (-2.1) \\
 &= -0.6
 \end{aligned}$$

# Principe d'une méta-analyse d'outcomes continus

Le poids d'une étude est égal à l'inverse de la variance de la différence de moyenne estimée dans cette étude (modèle à effet fixe):

Différence de moyenne =  $m_2 - m_1$

Variance de la différence de moyenne =  $\frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2}$

Poids =  $1 / \text{Variance} = 1 / \left( \frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2} \right)$

Erreur type de la différence de moyenne =  $\sqrt{\frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2}}$

IC95% de la différence de moyenne  
=  $(m_2 - m_1) \pm 1.96$  erreur type de la différence de moyenne

**Les données nécessaires à la méta-analyse sont:**

- 1) la différence de moyenne
- 2) la variance (ou l'erreur type) de la différence de moyennes

Gómez-Díaz et al,<sup>6</sup> 2012

Metformin		Control	
No.	Mean (SD)	No.	Mean (SD)
28	-2.2 (2.7)	24	-1.5 (2.6)

$$m_2 - m_1 = -2.2 - (-1.5) = -0.7$$

$$\frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2} = 2.7^2/28 + 2.6^2/24 = 0.54$$

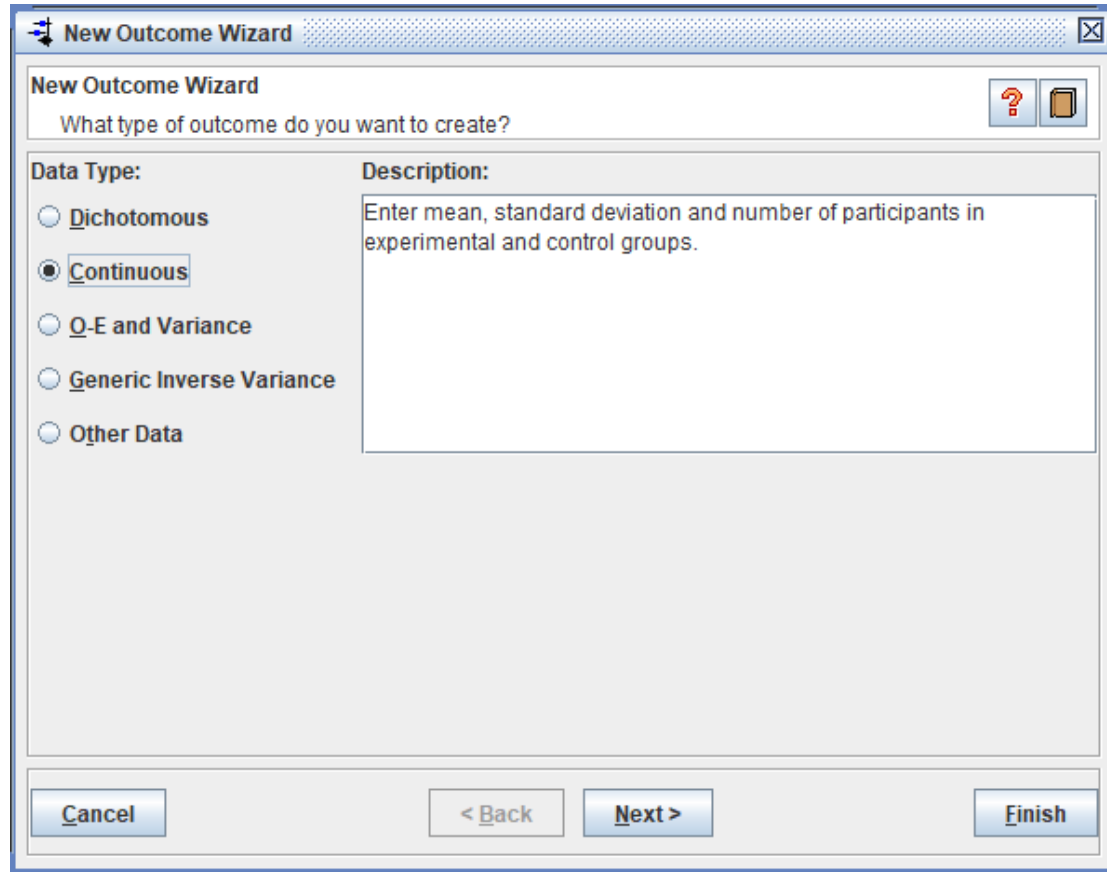
$$\text{Poids} = 1/0.54 = 1.8$$

$$\text{Erreur type} = \sqrt{0.54} = 0.74$$

$$\text{IC95\%} = -0.7 \pm 1.96 \times 0.74 = -2.2 \text{ à } +0.8$$

## Approche habituelle:

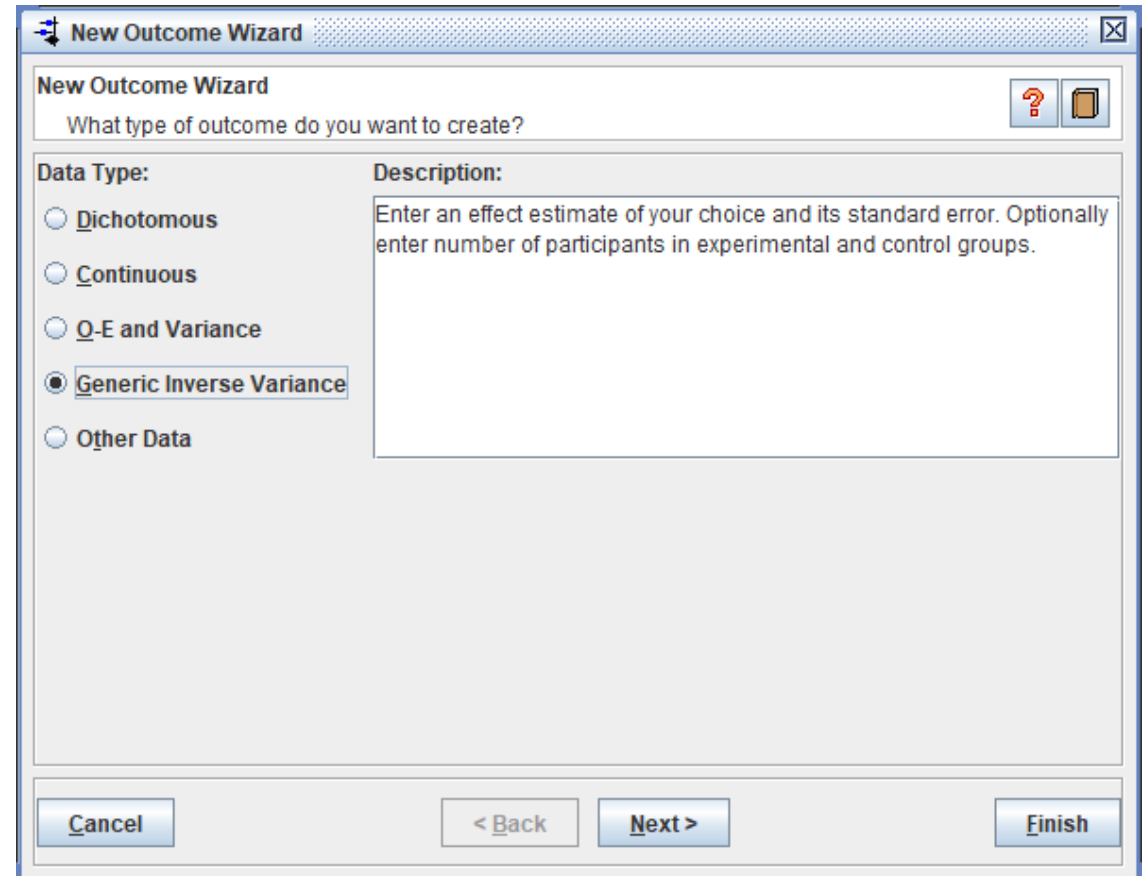
Entrer  $n_1$ ,  $m_1$ ,  $sd_1$  et  $n_2$ ,  $m_2$ ,  $sd_2$



The screenshot shows the 'New Outcome Wizard' dialog box. The title bar reads 'New Outcome Wizard'. Below the title bar, there is a question mark icon and a folder icon. The main text asks 'What type of outcome do you want to create?'. Under 'Data Type:', there are five radio button options: 'Dichotomous', 'Continuous' (which is selected), 'O-E and Variance', 'Generic Inverse Variance', and 'Other Data'. To the right, under 'Description:', there is a text box containing the text: 'Enter mean, standard deviation and number of participants in experimental and control groups.' At the bottom, there are four buttons: 'Cancel', '< Back', 'Next >', and 'Finish'.

## Méthode générique de l'inverse de la variance :

Entrer  $m_2 - m_1$  et **l'erreur type de la différence de moyenne**



The screenshot shows the 'New Outcome Wizard' dialog box. The title bar reads 'New Outcome Wizard'. Below the title bar, there is a question mark icon and a folder icon. The main text asks 'What type of outcome do you want to create?'. Under 'Data Type:', there are five radio button options: 'Dichotomous', 'Continuous', 'O-E and Variance', 'Generic Inverse Variance' (which is selected), and 'Other Data'. To the right, under 'Description:', there is a text box containing the text: 'Enter an effect estimate of your choice and its standard error. Optionally enter number of participants in experimental and control groups.' At the bottom, there are four buttons: 'Cancel', '< Back', 'Next >', and 'Finish'.

Les 2 approches sont disponibles dans les logiciels de méta-analyses  
L'estimation de l'effet commun est exactement identique

# Données incomplètes: situation 1

- Informations rapportées dans un article:
  - Nombre d'observations dans chaque bras:  $n_1, n_2$
  - Moyenne et intervalle de confiance à 95% (ou erreur type de la moyenne) dans chaque bras:  $m_1$  (IC95%  $l_1$  à  $u_1$ ) et  $m_2$  (IC95%  $l_2$  à  $u_2$ )
  - Les écarts types sont manquants

IC95% de  $m_1 = m_1 \pm 1.96 \times$  erreur type de la moyenne

Erreur type de la moyenne = largeur de l'IC95% / (2x1.96)  
 $= (u_1 - l_1) / (2 \times 1.96)$

Erreur type de la moyenne =  $\frac{sd_1}{\sqrt{n_1}}$

$sd_1 =$  Erreur type de la moyenne  $\times \sqrt{n_1}$   
 $= \sqrt{n_1} (u_1 - l_1) / (2 \times 1.96)$

Groupe Metformin de l'étude Gomes-Diaz et al

$n_1=28, m_1= -2.2$  (IC95% -3.2 à -1.2)

Erreur type =  $(-1.2 - (-3.2)) / (2 \times 1.96) = 0.50$

$sd_1 = 0.50 \times \sqrt{28} = 2.6$

# Données incomplètes: situation 2

- Certaines études rapportent les données complètes
  - $n_1, m_1, sd_1$  et  $n_2, m_2, sd_2$
- D'autres rapportent la différence de moyenne avec l'IC95%
  - $m_2 - m_1$  (IC95%  $l_{diff}$  à  $u_{diff}$ )

Pour chaque étude, on calcule :

- la différence de moyenne
- l'erreur type de la différence de moyenne

Puis on applique la méthode générique de l'inverse de la variance pour combiner les études

# Données incomplètes: situation 2

- Etudes rapportant les données complètes  $n_1, m_1, sd_1$  et  $n_2, m_2, sd_2$

Différence de moyenne =  $m_2 - m_1$

Erreur type de la différence de moyenne =  $\sqrt{\frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2}}$

The screenshot shows a statistical calculator window titled "Calculator - [New Outcome]". It is divided into two main sections: "Experimental" and "Control". Each section has input fields for Mean, N, and SD. Below these are summary statistics for each group, including SE, CI Start, CI End, tTest, and P value. At the bottom, there are summary statistics for the combined data, including N, MD, SE, CI Start, CI End, tTest, and P value. The MD and SE fields for the combined data are highlighted with a red box. The confidence interval is set to 95%.

Group	Mean	N	SD	SE	CI Start	CI End	tTest	P value	
Experimental	-2.2	28	2.7	0.5103	-3.2470	-1.1530	-4.3116	0.0002	
Control	-1.5	24	2.6	0.5307	-2.5979	-0.4021	-2.8263	0.0096	
Summary		52							
				MD	SE	CI Start	CI End	tTest	P value
				-0.7000	0.7362	-2.1430	0.7430	-0.9480	0.3477

Calculateur dans le logiciel Revman



# Données incomplètes: situation 2

- Etudes rapportant  $m_2 - m_1$  (IC95%  $l_{\text{diff}}$  à  $u_{\text{diff}}$ )

IC95% de la différence de moyenne =  $(m_2 - m_1) \pm 1.96$  erreur type de la différence de moyenne

Erreur type de la différence de moyenne =  $(u_{\text{diff}} - l_{\text{diff}}) / (2 \times 1.96)$

Etude Mauras et al

$m_2 - m_1 = -0.9$  (IC95% -2.3 à +0.5)

Erreur type de la différence de moyenne =  $(0.5 - (-2.3)) / (2 \times 1.96) = 0.71$

# Données incomplètes: situation 3

- Certaines études rapportent
  - Le nombre d'observations dans chaque bras:  $n_1, n_2$
  - la moyenne dans chaque bras:  $m_1, m_2$  (ou la différence  $m_2 - m_1$ )
  - la valeur p d'un test de Student comparant les 2 bras

Différence de moyenne =  $m_2 - m_1$

Statistique de test  $t = \frac{m_2 - m_1}{\text{erreur type de la différence de moyenne}}$

La statistique de test  $t$  se déduit de la valeur p rapportée

Erreur type de la différence de moyenne =  $\frac{m_2 - m_1}{\text{statistique de test } t}$

# Données incomplètes: situation 3

Etude Mauras et al:

$n_1=23$  et  $n_2=19$

$m_2 - m_1 = -0.9$ ,  $p=0.22$

<https://goodcalculators.com/student-t-value-calculator/>

Online T-Value Calculator

Degrees of Freedom (df):   $n_2+n_1-2$

Significance Level ( $\alpha$ ):  Valeur p

## Results

T-Value (right-tailed): 0.7799693

T-Value (two-tailed): +/- 1.246036

Statistique de test

$$\text{Erreur type de la différence de moyenne} = \frac{-0.9}{-1.25} = 0.72$$

Une erreur type est toujours positive

Pour ne pas se tromper, on peut diviser la valeur absolue de la différence de moyenne par la valeur absolue de la statistique de test

# Données incomplètes: situation 3

---

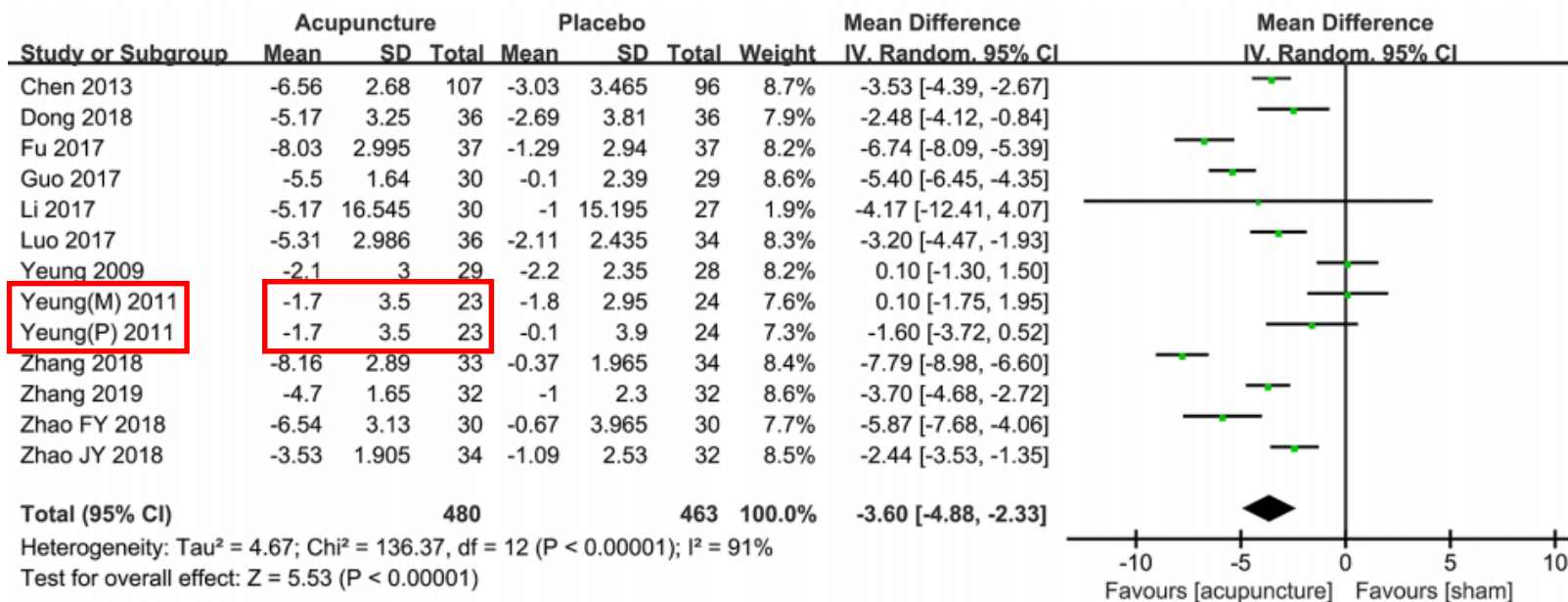
- Cette approche est inutilisable lorsque:
  - la différence de moyenne est égale à 0: la statistique de test vaut 0 quelle que soit l'erreur type (valeur  $p = 1$ )
  - la valeur  $p$  est rapportée de manière imprécise (e.g.  $p < 0.001$ ,  $p < 0.05$  ou NS)
  - la valeur  $p$  est celle d'un test de Mann-Whitney
- La méthode présentée suppose que le test de Student est bilatéral (two-sided)
- Bien s'assurer que le test utiliser est un test de Student bilatéral

# Données incomplètes: situation 4

- Certaines études sont des essais cliniques randomisés à 3 bras dont l'intervention donnée dans 2 bras correspond approximativement à l'intervention évaluée dans la revue systématique (par ex. doses assez proches)
- Données rapportées:
  - groupe comparateur:  $n_{\text{comp}}, m_{\text{comp}}, sd_{\text{comp}}$
  - groupe expérimental 1:  $n_1, m_1, sd_1$
  - groupe expérimental 2:  $n_2, m_2, sd_2$
- L'analyste souhaite agréger les 2 groupes expérimentaux (si cela fait sens scientifiquement)

# Données incomplètes: situation 4

Duplication des données du groupe intervention



Yeung(M): contrôle = acupuncture minimale  
 Yeung(P): contrôle = placebo

# Données incomplètes: situation 4

- Si cela fait sens de regrouper deux bras de l'essai, les données des 2 groupes peuvent être agrégées:
  - $n_{1+2} = n_1 + n_2$
  - $m_{1+2} = (n_1 \times m_1 + n_2 \times m_2) / (n_1 + n_2)$
  - $sd_{1+2} = ?$

Décomposition de la variance (conséquence du théorème de König-Huygens)

La variance du groupe agrégé est égale à la somme de

- 1) la moyenne des variances dans les sous-groupes
- 2) la variance des moyennes des sous-groupes

$$sd_{1+2} = \sqrt{\frac{n_1 sd_1^2 + n_2 sd_2^2}{n_{1+2}} + \frac{n_1}{n_{1+2}} (m_1 - m_{1+2})^2 + \frac{n_2}{n_{1+2}} (m_2 - m_{1+2})^2}$$

Moyenne des variances  
des 2 groupes

Variances des moyennes des 2 groupes (par  
rapport à la moyenne du groupe agrégé)



# Données incomplètes: situation 4

Agrégation des groupes contrôles de l'étude Yeung 2011:

Study or Subgroup	Acupuncture			Placebo			Weight	Mean Difference	
	Mean	SD	Total	Mean	SD	Total		IV, Random	95% CI
Yeung(M) 2011	-1.7	3.5	23	-1.8	2.95	24	7.6%	0.10	[-1.75, 1.95]
Yeung(P) 2011	-1.7	3.5	23	-0.1	3.9	24	7.3%	-1.60	[-3.72, 0.52]

$$n_1 = 24, m_1 = -1.8, sd_1 = 2.95$$

$$n_2 = 24, m_2 = -0.1, sd_2 = 3.90$$

$$n_{1+2} = 24 + 24 = 48$$

$$m_{1+2} = (24 \times (-1.8) + 24 \times (-0.1)) / 48 = -0.95$$

$$\text{Moyenne des variances} = (24 \times 2.95^2 + 24 \times 3.90^2) / 48 = 11.96$$

$$\text{Variance des moyennes} = 24 \times (-1.8 - (-0.95))^2 / 48 + 24 \times (-0.1 - (-0.95))^2 / 48 = 0.72$$

$$sd_{1+2} = \sqrt{11.96 + 0.72} = 3.56$$

Study or Subgroup	Acupuncture			Placebo			Weight	Mean Difference	
	Mean	SD	Total	Mean	SD	Total		IV, Random	95% CI
Yeung 2011	-1.7	3.5	23	-0.95	3.56	48			

# Données incomplètes: situation 5

---

- Certaines études rapportent le nombre d'observations, la médiane, les quartiles et/ou l'étendue (min-max) dans chaque groupe
- Principe:
  - encadrer la moyenne (inconnue) par une borne inférieure et une borne supérieure qui peuvent être calculées avec les données à disposition (la médiane, quartiles et/ou min-max)
  - approximer la moyenne par le milieu de ces bornes
  - similaire pour l'écart type

Hozo et al. Estimating the mean and variance from the median, range, and the size of a sample. *BMC Medical Research Methodology*. 5:13 (2005).

Wan et al. Estimating the sample mean and standard deviation from the sample size, median, range and/or interquartile range. *BMC Medical Research Methodology*. 14:135 (2014).

Bland M. Estimating mean and standard deviation from the sample size, three quartiles, minimum, and maximum. *International Journal of Statistics in Medical Research*, 2015,4:57-64.

# Données incomplètes: situation 5

Exemple lorsque seules les valeurs médiane, minimale et maximale sont rapportées dans un bras de l'essai

$$\text{min} \leq x_{(1)} \leq \text{min}$$

$$\text{min} \leq x_{(2)} \leq \text{médiane}$$

$$\text{min} \leq x_{(3)} \leq \text{médiane}$$

....

$$\text{min} \leq x_{((n-1)/2)} \leq \text{médiane}$$

$$\text{médiane} \leq x_{((n+1)/2)} \leq \text{médiane}$$

$$\text{médiane} \leq x_{((n+3)/2)} \leq \text{max}$$

....

$$\text{médiane} \leq x_{(n-1)} \leq \text{max}$$

$$\text{max} \leq x_{(n)} \leq \text{max}$$

Echantillon de taille  $n$  ( $n$  est supposé impair pour simplifier)

$x_{(i)}$  sont les observations ordonnées par ordre croissant

$x_{(1)}$  est la plus petite des  $n$  observations

$x_{(2)}$  la 2<sup>ème</sup> plus petite observation

....

$x_{(n)}$  est la plus grande des  $n$  observations

$x_{((n+1)/2)}$  est l'observation du «milieu» = médiane

Si on connaît les valeurs médiane, minimale et maximale des observations, on peut encadrer chacune des observations

# Données incomplètes: situation 5

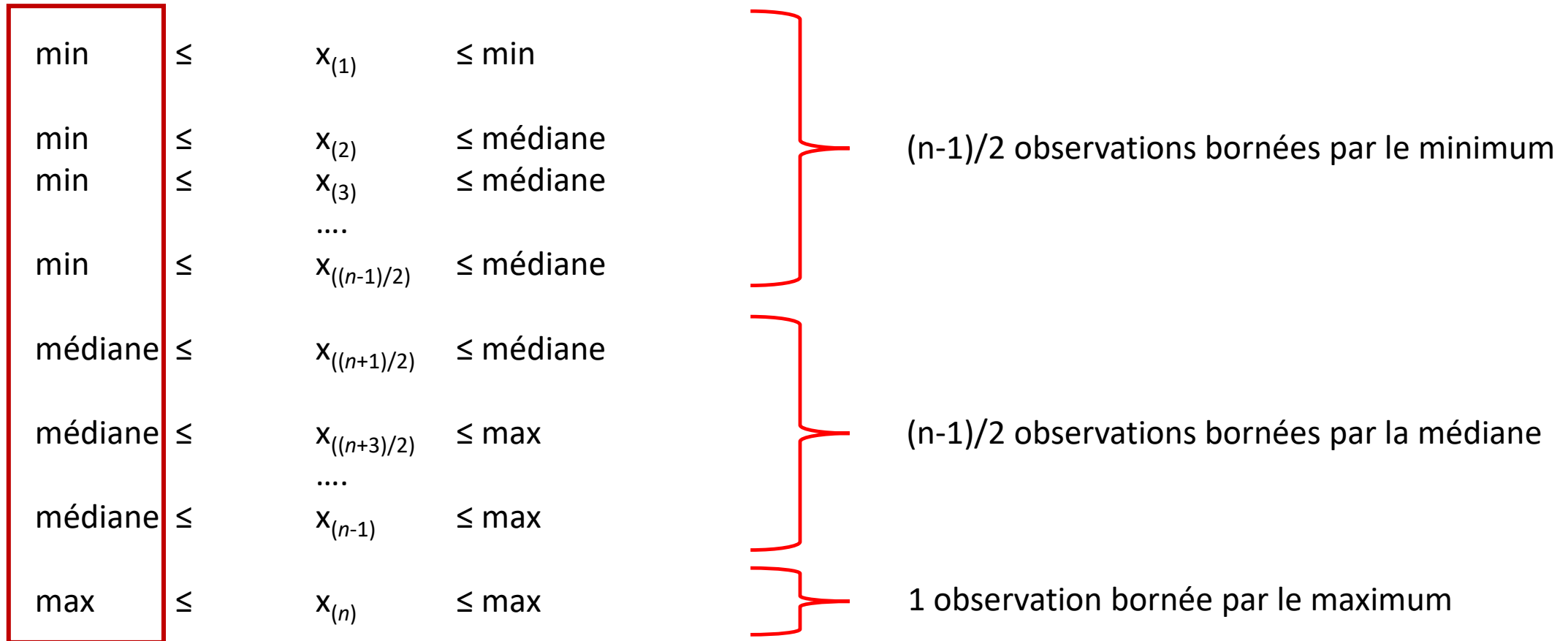
min	≤	$x_{(1)}$	≤ min
min	≤	$x_{(2)}$	≤ médiane
min	≤	$x_{(3)}$	≤ médiane
		....	
min	≤	$x_{((n-1)/2)}$	≤ médiane
médiane	≤	$x_{((n+1)/2)}$	≤ médiane
médiane	≤	$x_{((n+3)/2)}$	≤ max
		....	
médiane	≤	$x_{(n-1)}$	≤ max
max	≤	$x_{(n)}$	≤ max

Moyenne

Somme terme à terme des  $n$   
inégalités puis division par le  
nombre d'observations  $n$

$$\frac{\sum_{i=1}^n x_{(i)}}{n} = \text{moyenne}$$

# Données incomplètes: situation 5



LB ≤ Moyenne

$$LB = \frac{\frac{n-1}{2}min + \frac{n-1}{2}médiane + max}{n}$$

# Données incomplètes: situation 5

min	≤	$x_{(1)}$	≤ min	
min	≤	$x_{(2)}$	≤ médiane	
min	≤	$x_{(3)}$	≤ médiane	
....		....		
min	≤	$x_{((n-1)/2)}$	≤ médiane	
médiane	≤	$x_{((n+1)/2)}$	≤ médiane	
médiane	≤	$x_{((n+3)/2)}$	≤ max	
....		....		
médiane	≤	$x_{(n-1)}$	≤ max	
max	≤	$x_{(n)}$	≤ max	

---

LB ≤ Moyenne ≤ UB

$$UB = \frac{\frac{n-1}{2} \text{médiane} + \frac{n-1}{2} \text{max} + \text{min}}{n}$$

# Données incomplètes: situation 5

Finalement la moyenne peut être encadrée :

$$\frac{\frac{n-1}{2} \min + \frac{n-1}{2} \text{médiane} + \max}{n} \leq \text{Moyenne} \leq \frac{\frac{n-1}{2} \text{médiane} + \frac{n-1}{2} \max + \min}{n}$$

Estimation de la moyenne par:

$$(\text{LB} + \text{UP})/2 = \frac{\min + 2 \text{médiane} + \max}{4} + \frac{\min - 2 \text{médiane} + \max}{4n}$$



Terme négligeable lorsque  
 $n$  est grand



# Données incomplètes: situation 5

	Scenario 1	Scenario 2	Scenario 3
Données rapportées	Nombre d'observations ( $n$ ) min médiane ( $méd$ ) max	Nombre d'observations ( $n$ ) min 1 <sup>er</sup> quartile ( $q_1$ ) médiane ( $méd$ ) 3 <sup>ème</sup> quartile ( $q_3$ ) max	Nombre d'observations ( $n$ ) 1 <sup>er</sup> quartile ( $q_1$ ) médiane ( $méd$ ) 3 <sup>ème</sup> quartile ( $q_3$ )
Moyenne $\approx$	$\frac{min + 2méd + max}{4}$	$\frac{min + 2q_1 + 2méd + 2q_3 + max}{8}$	$\frac{q_1 + méd + q_3}{3}$
Ecart type $\approx$	$\frac{max - min}{2\Phi^{-1}\left(\frac{n - 0.375}{n + 0.25}\right)}$	$\frac{max - min}{4\Phi^{-1}\left(\frac{n - 0.375}{n + 0.25}\right)} + \frac{q_3 - q_1}{4\Phi^{-1}\left(\frac{0.75n - 0.125}{n + 0.25}\right)}$	$\frac{q_3 - q_1}{2\Phi^{-1}\left(\frac{0.75n - 0.125}{n + 0.25}\right)}$

Où  $\Phi(\cdot)$  est la distribution cumulée d'une loi normale standard

NB: lorsque  $n$  est grand,  $2\Phi^{-1}\left(\frac{0.75n-0.125}{n+0.25}\right) \approx 1.35$

# Données incomplètes: situation 5

Une étude rapporte pour le bras intervention (n=50) une médiane de 15 et un intervalle interquartile de 10 à 22 (scenario 3)

$$\text{Moyenne} = (10 + 15 + 22) / 3 = 15.67$$

$$\text{Ecart type} = \frac{q_3 - q_1}{2\Phi^{-1}\left(\frac{0.75n - 0.125}{n + 0.25}\right)}$$

$$q_3 - q_1 = 22 - 10 = 12$$

$$\frac{0.75 \times 50 - 0.125}{50 + 0.25} = 0.744$$

$$\text{Ecart type} = \frac{12}{2 \times 0.656}$$

$$= 9.15$$

<https://planetcalc.com/4987/>

Probability  
0.744

Variance  
1

Mean  
0

Calculation precision  
Digits after the decimal point: 3

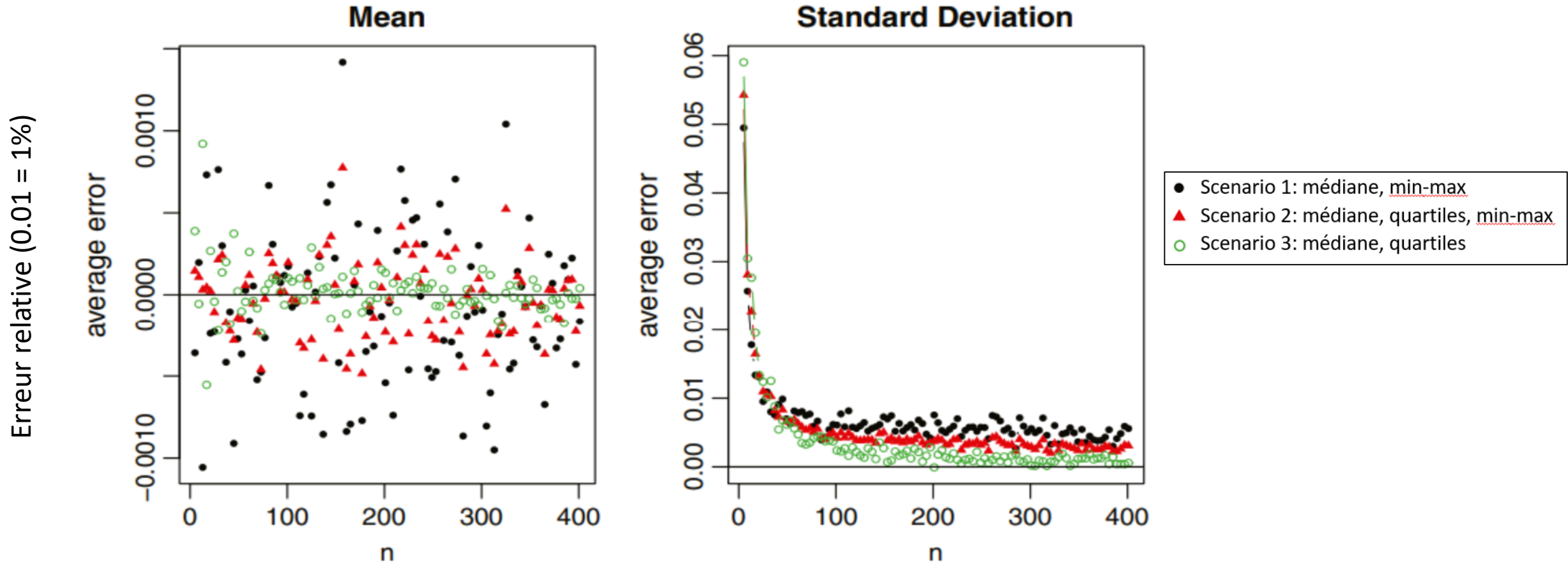
Quantile  
0.656 =  $\Phi^{-1}(0.744)$

CALCULATE

Paramètres  
d'une loi  
normale  
standard

# Données incomplètes: situation 5

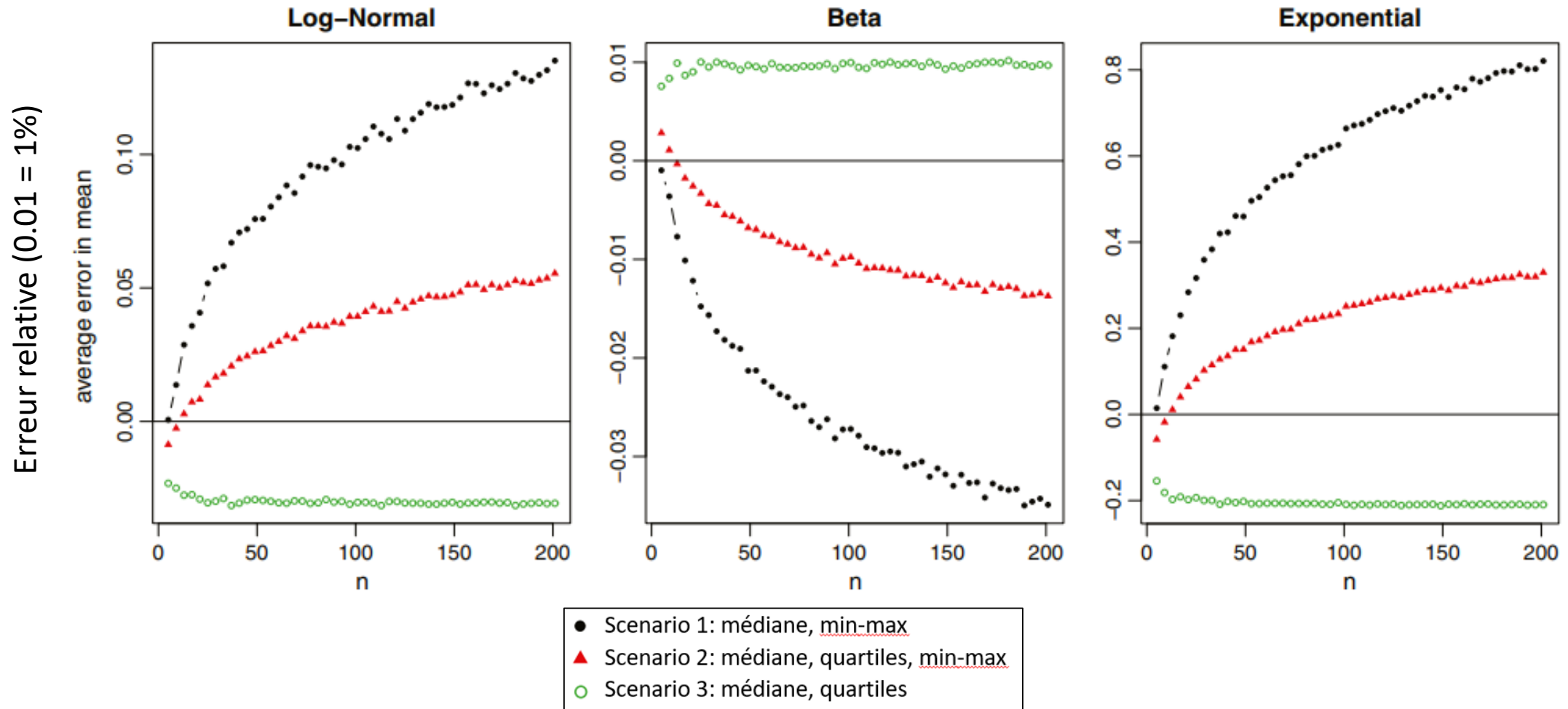
- Performances des 3 scénarios lorsque les données sont normalement distribuées



Toutes les méthodes donnent de bons résultats (erreur relative < 1% ou 2% globalement)

# Données incomplètes: situation 5

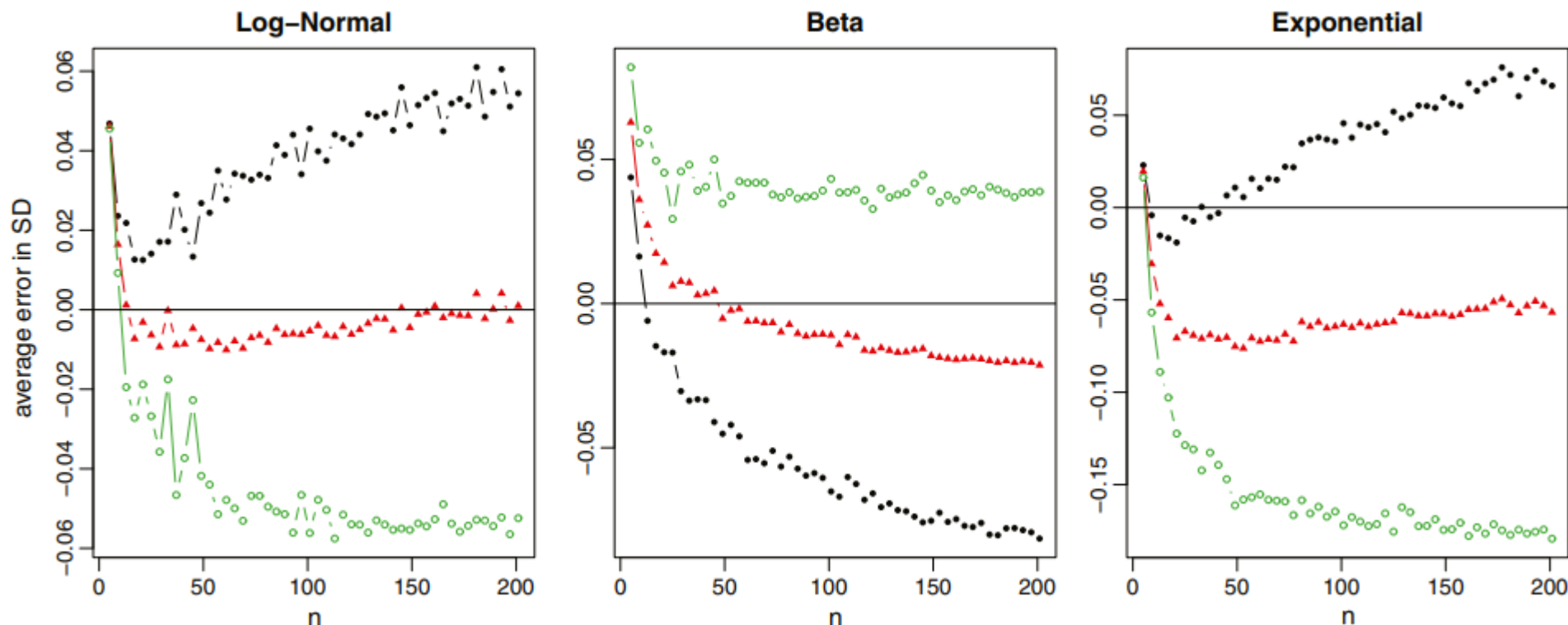
- Performances de l'estimation de la moyenne lorsque les données ne sont pas normalement distribuées



# Données incomplètes: situation 5

- Performances de l'estimation de l'écart type lorsque les données ne sont pas normalement distribuées

Erreur relative (0.01 = 1%)



- Scenario 1: médiane, min-max
- ▲ Scenario 2: médiane, quartiles, min-max
- Scenario 3: médiane, quartiles

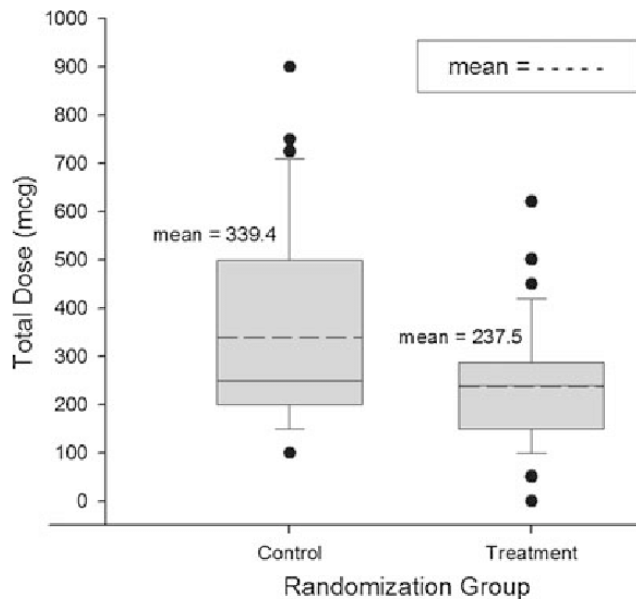
# Données incomplètes: situation 5

- Attention: les erreurs présentées sont des erreurs moyennes sur un grand nombre d'études simulées aléatoirement
- En général, les données ont une distribution approximativement normale ou log-normale
- Recommandations par Wan et al:
  - Données normalement distribuées:
    - Toutes les méthodes donnent de bon résultats (erreur relative  $< 1$  ou 2%)
    - Scenario 3 si possible (un peu meilleur)
  - Données non normalement distribuées
    - Scenario 2 si possible car scenario 3 sensible à l'asymétrie pour l'estimation de l'écart type
- Recommandations personnelles:
  - Données non normalement distribuées et on dispose de médiane, quartiles, min-max
    - Scenario 2 pour l'estimation de l'écart type
    - Scenario 2 pour l'estimation de la moyenne lorsque  $n < 100$
    - Scenario 3 pour l'estimation de la moyenne lorsque  $n \geq 100$

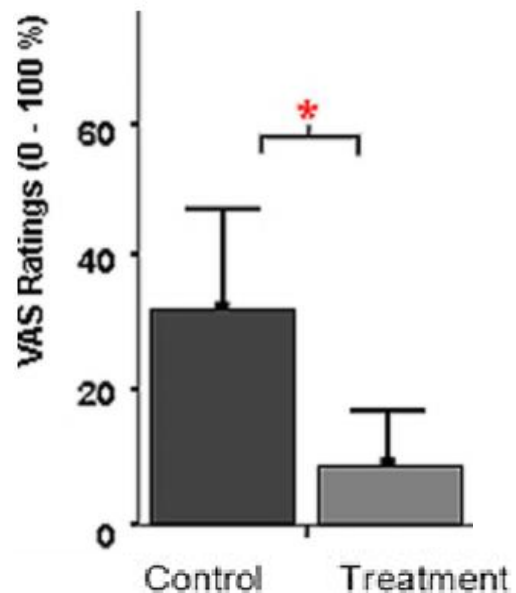
# Données incomplètes: situation 6

- Certaines études rapportent les données seulement sous forme graphique

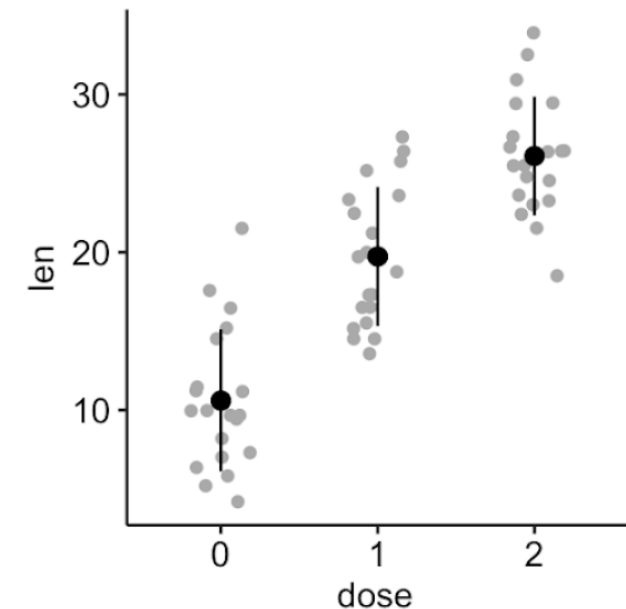
Box plot: médiane, quartiles, minimum, maximum



«Dynamite» plot: moyenne avec sd ou sem ou IC95%



Dot plot: valeurs individuelles avec éventuellement moyenne / médiane, ...

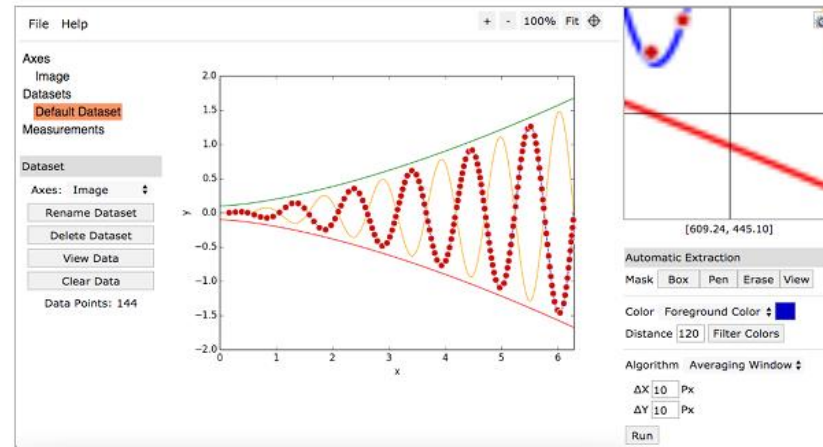




# Données incomplètes: situation 6

---

- Utilisation d'un logiciel / application permettant d'extraire les coordonnées de points sélectionnés sur la figure
  - Enregistrer la figure dans un fichier au format image (par ex. jpg)
    - par ex. capture d'écran
    - puis coller la capture d'écran dans le logiciel Paint3D (Microsoft)
    - sélectionner la figure et l'enregistrer au format .jpg
  - Importer la figure dans le logiciel / application
  - Définir un repère (axes x et y)
  - Sélectionner les points en cliquant dessus



### Web Application

English ▼

Launch Now!

### Desktop Version



### View Source

GitHub

It is often necessary to reverse engineer images of data visualizations to extract the underlying numerical data. WebPlotDigitizer is a semi-automated tool that makes this process extremely easy:

- Works with a wide variety of charts (XY, bar, polar, ternary, maps etc.)
- Automatic extraction algorithms make it easy to extract a large number of data points
- Free to use, opensource and cross-platform (web and desktop)
- Used in hundreds of published works by thousands of users
- Also useful for measuring distances or angles between various features
- More to come soon...

Version 4.4 Released (November 28, 2020)

[\[ Release Notes \]](#)

Load Image File(s)

Sélect. fichiers

Aucun fichier choisi

Load

Cancel

Ouvrir



<< COMB... > PresentationCRC10mai2021

Rechercher dans : Presentati...

Organiser

Nouveau dossier

Accès rapide

Bureau

Téléchargements

Documents

Nom

Modifié le

693.full.pdf

08.05.2021 12:31

FigureDotPlot.jpg

06.05.2021 15:56

PresentationExtractionDataMetaAnalysev...

08.05.2021 17:16

[739.44, 234.51]

WebPlotDigitizer 4.4

Load Image

Tutorial Video

Visit Website

Visit GitHub

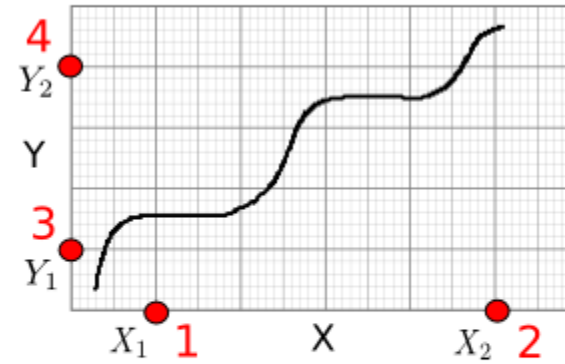
### Choose Plot Type

2D (X-Y) Plot

- 2D Bar Plot
- Polar Diagram
- Ternary Diagram
- Map With Scale Bar
- Image

Align Axes Cancel

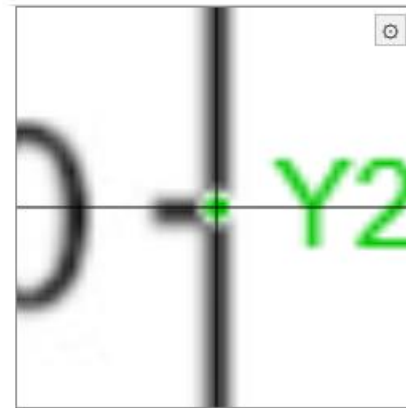
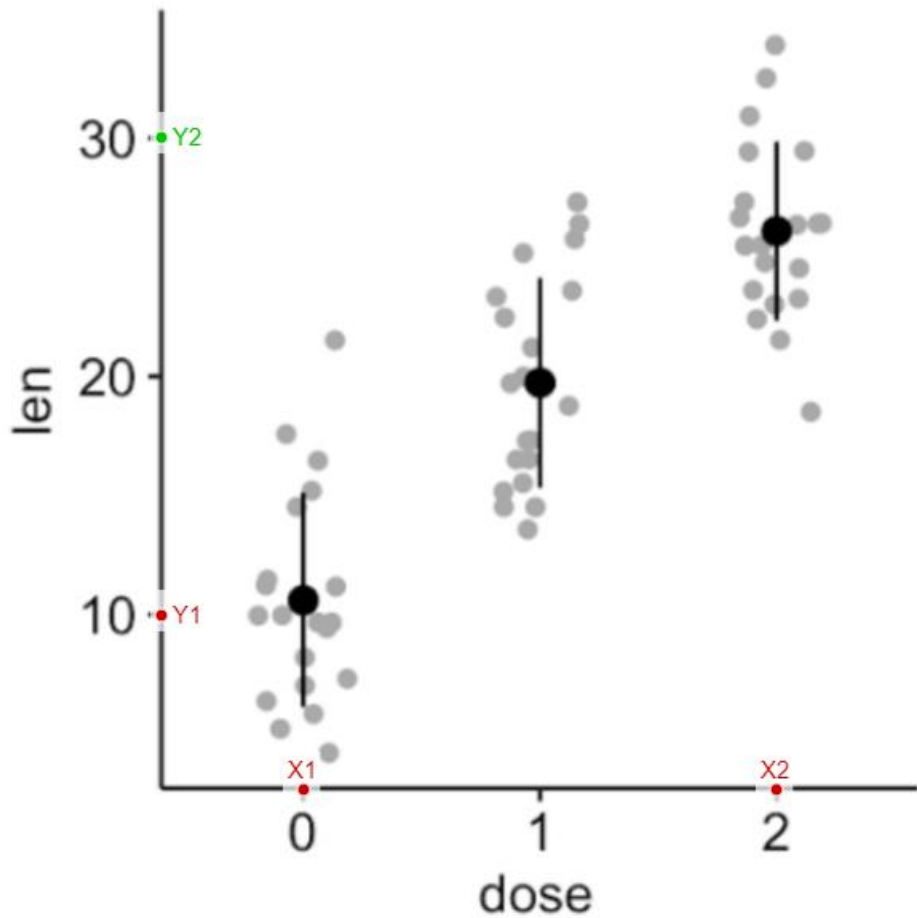
### Align X-Y Axes



Click four known points on the axes in the **order shown in red**. Two on the X axis ( $X_1$ ,  $X_2$ ) and two on the Y axis ( $Y_1$ ,  $Y_2$ ).

Proceed

## Détermination du repère



[525.72, 140.07]

#### Axes Calibration

Click points to select and use cursor keys to adjust positions. Use Shift+Arrow for faster movement. Click complete when finished.

Complete!

#### X and Y Axes Calibration

Enter X-values of the two points clicked on X-axis and Y-values of the two points clicked on Y-axes

	Point 1	Point 2	Log Scale
X-Axis:	0	2	<input type="checkbox"/>
Y-Axis:	10	30	<input type="checkbox"/>

Assume axes are perfectly aligned with image coordinates (skip rotation correction)

\*For dates, use yyyy/mm/dd hh:ii:ss format, where ii denotes minutes (e.g. 2013/10/23 or 2013/10 or 2013/10/23 10:15 or just 10:15). For exponents, enter values as 1e-3 for 10^-3.

OK

Cliquer sur les points «repère» et renseigner les valeurs correspondantes

Image  
Axes  
XY  
Datasets  
■ Default Dataset  
Measurements

Dataset

Axes: XY

Display Color

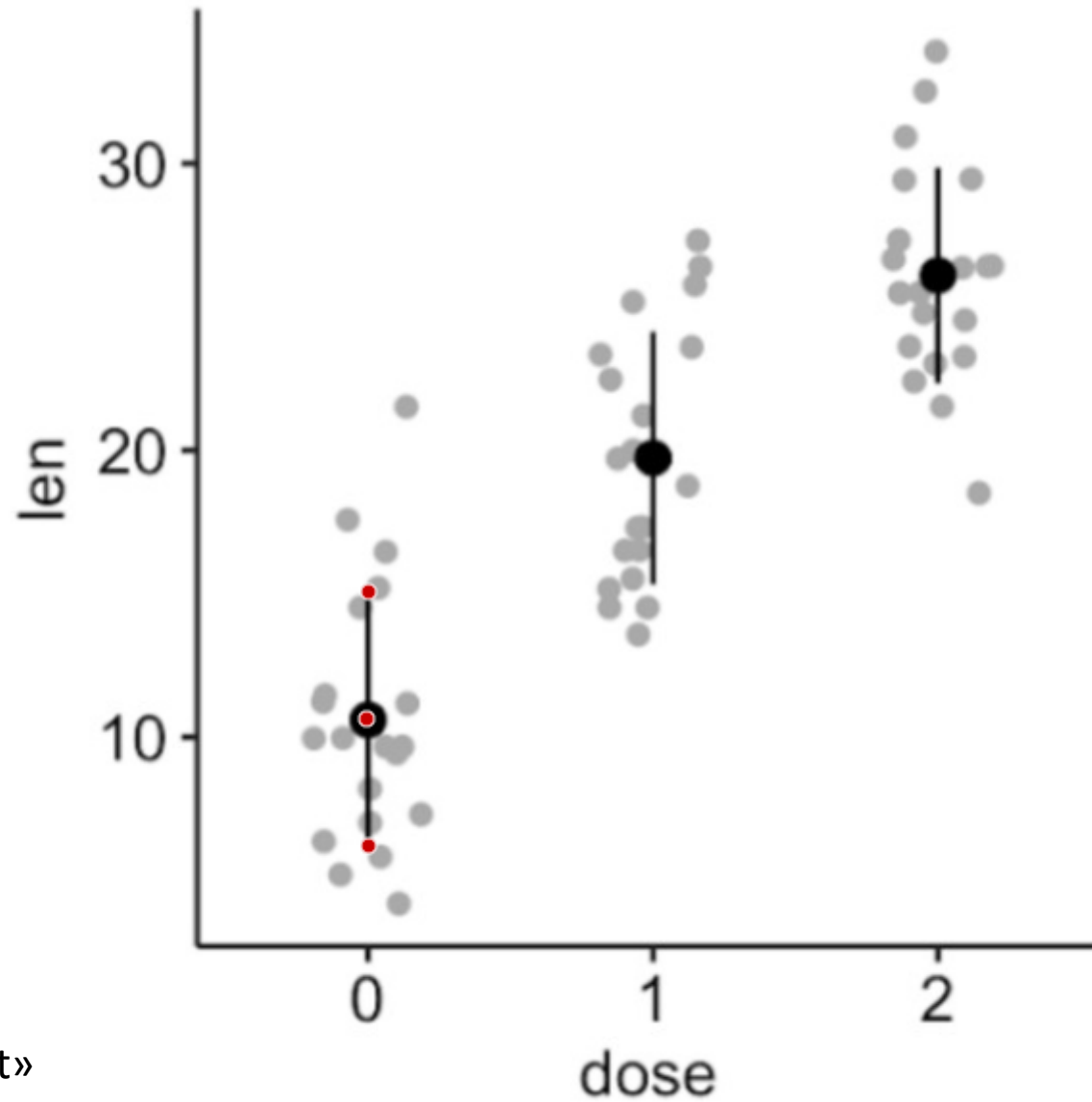
Rename Dataset

Delete Dataset

View Data

Clear Data

Data Points: 3



[ -6.7568e-3, 1.5084e+1 ]

Manual Extraction

Add Point (A) Adjust Point (S)

Delete Point (D)

Automatic Extraction

Mask Box Pen Erase View

Color Foreground Color

Distance 120 Filter Colors

Algorithm Averaging Window

$\Delta X$  10 Px

$\Delta Y$  10 Py

Run

- 1) Sélectionner «Add point»
- 2) Cliquer sur les points souhaités
- 3) Sélectionner «View data»

## Acquired Data

Dataset: Default Dataset ▾

Variables: X, Y

```
-0,006756756756756133; 10,668896321070239  
0; 6,2541806020066915  
0; 15,083612040133787
```

### Sort

Sort by: Raw ▾

Order: Ascending ▾

### Format

Number Formatting:

Digits:  Ignore ▾

Column Separator:

Format

\*Plotly is a secure data analysis and graphing site with data sharing and access controls.

Visit <http://plot.ly> for details.

- Si le point central sur le graphe représente la moyenne et la barre horizontale l'IC95%, j'obtiens:
  - moyenne = 10.7
  - IC95%: 6.3 à 15.1

Très facile à utiliser

# Données incomplètes: situation 6

- Très bonnes performances (concordance correlation coefficients  $> 0.90$ )<sup>1</sup>:
  - précision des valeurs extraites
  - reproductibilité inter- et intra-extracteurs
- Précision des valeurs extraites similaire à celle d'une extraction manuelle mais extraction computer-based 2 fois plus rapide<sup>2</sup>
- WebPlotDigitizer (gratuit) plus ergonomique que d'autres logiciels (payants)<sup>3</sup>
- Tutoriels:
  - <https://automeris.io/WebPlotDigitizer/tutorial.html>
  - Vidéos sur YouTube

<sup>1</sup>Van der Mierden et al. Extracting data from graphs: a case-study on animal research with implications for meta-analysis. Research Synthesis Methods 2021; 1-10.

<sup>2</sup>Kadic et al. Extracting data from figures with software was faster, with higher interrater reliability than manual extraction. Journal of Clinical Epidemiology 2016; 74:119-123.

<sup>3</sup>Drevon et al. Intercoder reliability and validity of WebPlotDigitizer in extracting data. Behaviour Modification 2016; 41(2)



# Conclusion

---

- Situations couvertes
  - $n_1, m_1, \text{IC95\% } l_1 \text{ à } u_1$  et  $n_2, m_2, \text{IC95\% } l_2 \text{ à } u_2$
  - $m_2 - m_1$  et  $\text{IC95\% } l_{\text{diff}} \text{ à } u_{\text{diff}}$
  - $n_1, n_2, m_1$  et  $m_2$  (ou  $m_2 - m_1$ ), valeur p du test de Student
  - Essai à 3 bras:  $n_0, m_0, \text{sd}_0, n_1, m_1, \text{sd}_1, n_2, m_2, \text{sd}_2$ ,
  - $n$ , médiane, quartiles et/ou min-max dans chaque bras
  - représentation graphique des données
  
- Groupes supposés indépendants

# Conclusion

---

- Certaines méthodes sont plus «exactes» que d'autres
  - Risque d'erreur lorsque la moyenne et l'écart type sont estimées à partir de la médiane et des quartiles et/ou min-max, surtout la distribution des observations est asymétrique
- Faire une analyse de sensibilité pour tester l'influence des études dont le résultat a été reconstruit (une méta-analyse sans ces études et une méta-analyse)
  - Si le résultat change sensiblement et que la distribution des observations est asymétrique, on ne peut pas savoir si la différence provient des nouvelles études elles-mêmes ou d'un problème méthodologique
- Extraction à partir de figures: il faut être certain de ce que représente la figure
- Pour ces méthodes, préférer les réponses des auteurs à la reconstruction de valeurs d'outcomes

# Conclusion

---

- Les méthodes peuvent être enchaînées
  - Par ex. extraction d'un intervalle de confiance à 95% à partir d'une figure puis calcul de l'écart type à partir des bornes de l'intervalle de confiance
- Autres outcomes
  - outcomes binaires
    - problème des données incomplètes moins fréquent (selon moi)
  - outcomes de survie
    - plus complexe
    - idéalement, il faut combiner les hazards ratios
    - lorsqu'un hazard ratio n'est pas rapporté, on peut l'estimer en reconstruisant les données à partir des courbes de survie publiées