

# ISTA Statistics Committee open meeting: overview of some recent projects



Presenter: Kirk Remund & Jean-Louis Laffont  
Location: Verona, Italy  
Date: May 31, 2023

# ISTA Statistics Committee



<b>Chair:</b>	<b>Kirk Remund</b>	<b>USA</b>
<b>Vice:</b>	<b>Jean-Louis Laffont</b>	<b>France</b>
<b>Members:</b>	<b>Gabriel Carré</b>	<b>France</b>
	<b>Mustapha El Yakhlifi</b>	<b>France</b>
	<b>Zhou Fang</b>	<b>USA</b>
	<b>Bonnie Hong</b>	<b>USA</b>
	<b>Bo-Jein Kuo</b>	<b>Separate Custom Territory of Taiwan, Penghu, Kinmen and Matsu</b>
	<b>Ray Shillito</b>	<b>USA</b>
	<b>Thomas Michelon</b>	<b>Brasil</b>
	<b>Oluseyi Odubote</b>	<b>USA</b>

1. New statistical tool for determining working sample weight to amend Table 2C of ISTA Rules
2. Number of sub-lots for which an OIC established for the lot is still valid
3. Group testing: number of groups to ensure that estimation is possible
4. Opportunities...

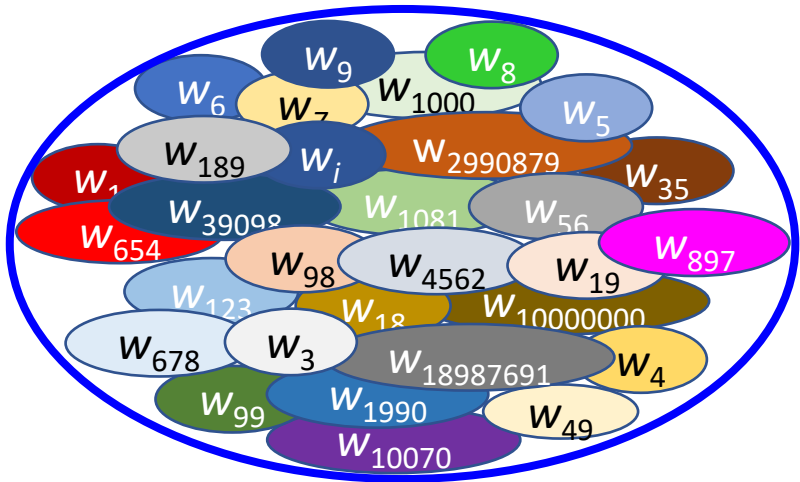


# 1. New statistical tool for determining working sample weight to amend Table 2C of ISTA Rules

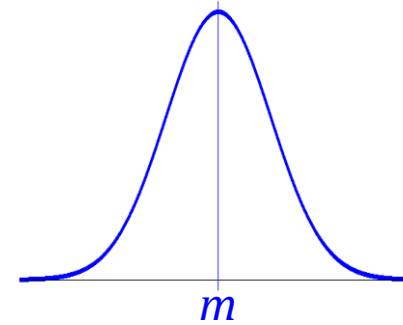
# 1. Principle

Population (all possible varieties, lots, labs and 100 seeds samples) of 100 seed weights

$Y \sim$  Normal distribution with mean  $m$  and variance  $\sigma^2$

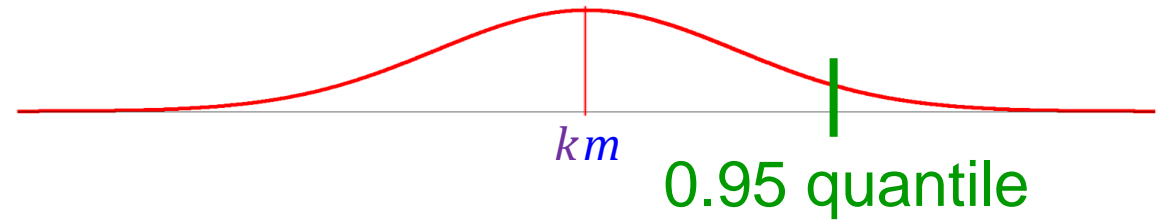


Assumption



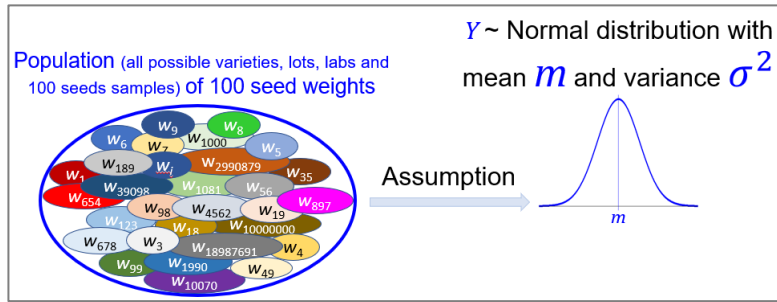
$X$ : “2500 ( $k = 25$ ) or 25000 ( $k = 250$ ) seeds weight”

If  $Y \sim N(m, \sigma^2)$ , then  $X = kY$  is  $\sim N(km, k^2\sigma^2)$



95% confident to have at least 2500 or 25000 seeds in a random sample with the 0.95 quantile weight

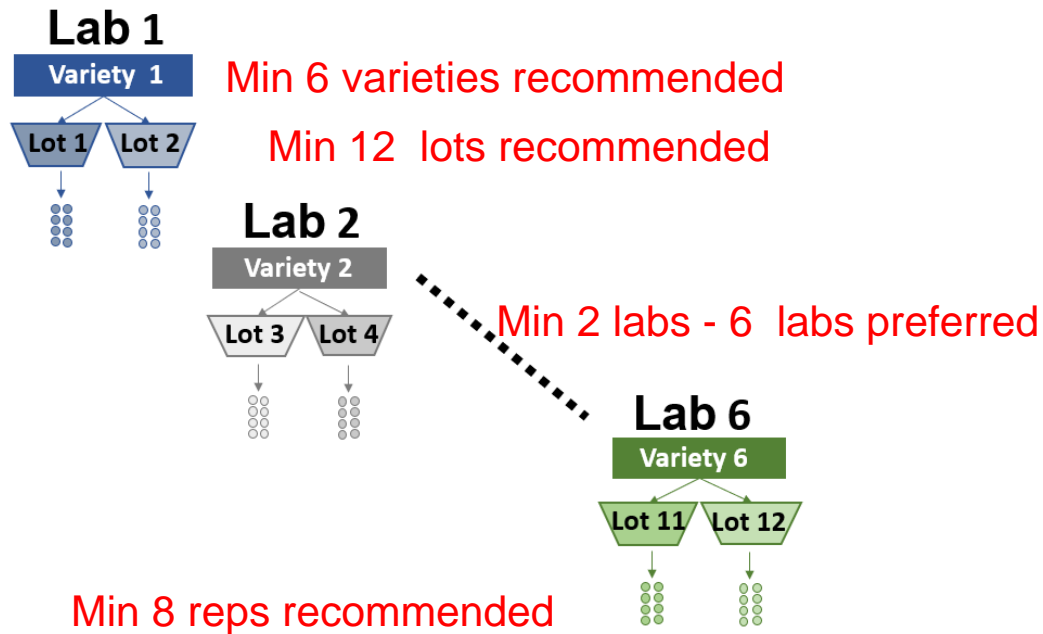
# 1. Principle



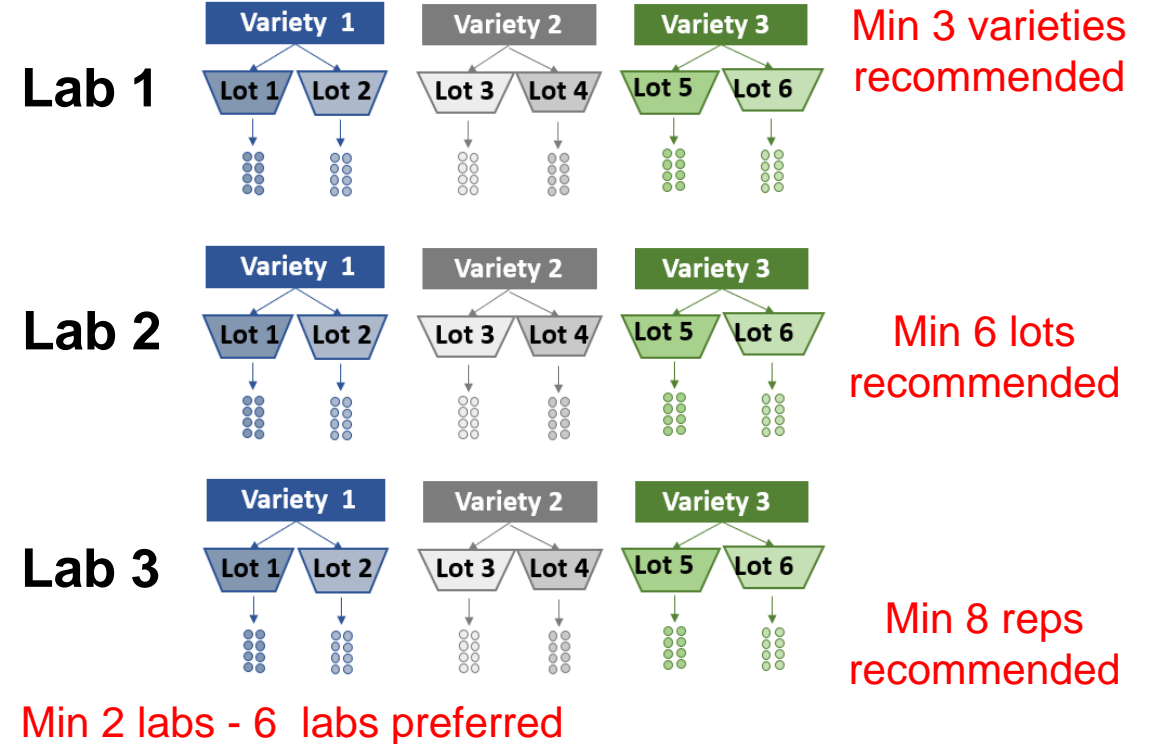
# Estimating $m$ and $\sigma^2$

2 experiment designs to capture all the possible sources of variation at its best

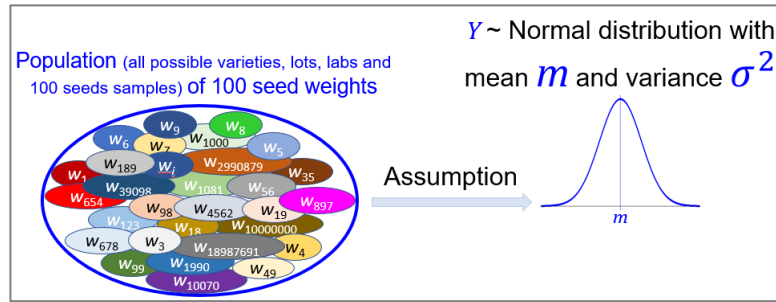
## Experiment design 1: 2-way nested design



## Experiment design 2: 2-way crossed design



# 1. Principle



## Estimating $m$ and $\sigma^2$

### Fitting linear random effects model

#### Experiment design 1 (2-way nested design)

$100\_seeds\_weight = general\_mean$

+  $Lab\_effect \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Lab}^2)$

+  $Lot(within\ Lab)\_effect \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Lot}^2)$

+  $Residual \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Res}^2)$

#### Experiment design 2 (2-way crossed design)

$100\_seeds\_weight = general\_mean$

+  $Lab\_effect \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Lab}^2)$

+  $Lot\_effect \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Lot}^2)$

+  $Lab \times Lot\_effect \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Lab \times Lot}^2)$

+  $Residual \rightarrow \sim$  i.i.d.  $N(0, \sigma_{Res}^2)$

$$\hat{m} = \widehat{general\_mean}$$

$$\widehat{\sigma^2} = \widehat{\sigma_{Lab}^2} + \widehat{\sigma_{Lot}^2} + \widehat{\sigma_{Res}^2}$$

$$\widehat{\sigma^2} = \widehat{\sigma_{Lab}^2} + \widehat{\sigma_{Lot}^2} + \widehat{\sigma_{Lab \times Lot}^2} + \widehat{\sigma_{Res}^2}$$

(estimates are denoted with a "hat")



## 1. Calculations details – Outliers detection

Prior to estimation, reps outliers are detected using **Grubbs's method**

$i$	$y_i$	$T_i$
1	0.4460	0.7600748
2	0.4190	0.3052932
3	0.4000	0.0147383
4	0.4270	0.4400433
5	0.2600	2.3728652
6	0.4100	0.1536993
7	0.4420	0.6926998
8	0.4030	0.035793

1. Calculate the mean  $\bar{y}$  and the standard-deviation  $s$  :

$$\bar{y} = 0.4009$$

$$s = 0.0594$$

2. For each value  $y_i$  in the dataset, calculate: 
$$T_i = \frac{|y_i - \bar{y}|}{s}$$

3. If  $T_i$  is greater than a **critical value** corresponding to a given significance probability (usually 5%), then identify  $y_i$  as an outlier



# 1. Calculations details – Outliers detection

$i$	$y_i$	$T_i$
1	0.4460	0.7600748
2	0.4190	0.3052932
3	0.4000	0.0147383
4	0.4270	0.4400433
5	0.2600	2.3728652
6	0.4100	0.1536993
7	0.4420	0.6926998
8	0.4030	0.035793

Critical values for 5% level of significance

Sample size	Critical value	Sample size	Critical value	Sample size	Critical value	Sample size	Critical value
3	1.15	15	2.55	27	2.86	39	3.03
4	1.48	16	2.59	28	2.88	40	3.04
5	1.71	17	2.62	29	2.89	50	3.13
6	1.89	18	2.65	30	2.91	60	3.20
7	2.02	19	2.68	31	2.92	70	3.26
8	2.13	20	2.71	32	2.94	80	3.31
9	2.21	21	2.73	33	2.95	90	3.35
10	2.29	22	2.76	34	2.97	100	3.38
11	2.34	23	2.78	35	2.98	110	3.42
12	2.41	24	2.80	36	2.99	120	3.44
13	2.46	25	2.82	37	3.00	130	3.47
14	2.51	26	2.84	38	3.01	140	3.49



$y_5$  is identified as an outlier

# 1. Calculations details – Outliers detection

Grubbs's method **critical values** can be calculated from the **student distribution** as follows:

$$\sqrt{\frac{\left[ (n - 1) t_{1 - \frac{\alpha}{2n}, n - 2} \right]^2}{n \left[ (n - 2) + \left( t_{1 - \frac{\alpha}{2n}, n - 2} \right)^2 \right]}}$$

where: .  $n$  : sample size

.  $\alpha$  : level of significance

.  $t_{1 - \frac{\alpha}{2n}, n - 2}$  :  $1 - \frac{\alpha}{2n}$  critical point of a  $t$ -distribution with  $n - 2$  degrees of freedom

Can be easily implemented into Excel:

	A	B	C	D	E	F	G
1	<b>Grubbs' method: critical values</b>						
2							
3	Level of significance:	5%					
4	Sample size:	15					
5							
6	Grubbs' critical value:	2.55					
7							

# 1. Calculations details – Variance components estimation

## 2-way nested design

$$\begin{aligned}
 100\_seeds\_weight &= general\_mean \\
 &+ Lab\_effect \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Lab}^2) \\
 &+ Lot(\text{within } Lab)\_effect \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Lot}^2) \\
 &+ Residual \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Res}^2)
 \end{aligned}$$

## 2-way crossed design

$$\begin{aligned}
 100\_seeds\_weight &= general\_mean \\
 &+ Lab\_effect \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Lab}^2) \\
 &+ Lot\_effect \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Lot}^2) \\
 &+ Lab \times Lot\_effect \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Lab \times Lot}^2) \\
 &+ Residual \quad \rightarrow \sim \text{i.i.d. } N(0, \sigma_{Res}^2)
 \end{aligned}$$

- There are several methods to get estimates of variance components:
  - . ANOVA based methods
  - . Maximum Likelihood (ML) methods
  - . REstricted Maximum Likelihood (REML) methods
  - ...
- Today, the preferred method is REML ... but it requires heavy computations

➡ Selected **Henderson Method I** (ANOVA based method) for its ease of implementation in Excel.

This method works for **unbalanced data**; for **balanced data**, it provides identical estimates as REML method

# 1. Calculations details – Variance components estimation

## Henderson Method I

Searle, S.R., Casella, G. and C.E. McCulloch (1992).  
 In Variance components (pp. 429, 434-435). *Wiley-Interscience*, New York.

This is what is implemented in the calculator

[F.2] THE 2-WAY NESTED CLASSIFICATION 429

Then

$$\text{var}(\hat{\theta}_1^2) = 2\sigma_1^4(\sum_i w_i^2)/D,$$

$$\text{var}(\hat{\theta}_2^2) = 2\sigma_2^4(N - a + \sum_i w_i^2/n_i^2)/D$$

and

$$\text{cov}(\hat{\theta}_1^2, \hat{\theta}_2^2) = -2\sigma_1^2\sigma_2^2(\sum_i w_i^2/n_i)/D$$

(Crump, 1951; Searle, 1956).

F.2. THE 2-WAY NESTED CLASSIFICATION

a. Model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$$

$$i = 1, 2, \dots, a, \quad j = 1, 2, \dots, b, \quad \text{and} \quad k = 1, 2, \dots, n_{ij}$$

with

$$b = \sum_i b_i \quad \text{and} \quad N = \sum_i \sum_j n_{ij}$$

b. Analysis of variance estimators

Calculate

$$k_1 = \sum_i n_i^2/N, \quad k_2 = \sum_i \sum_j n_{ij}^2/N, \quad k_{12} = \sum_i (\sum_j n_{ij}^2/n_i) = \sum_i \hat{\theta}_i^2$$

$$T_A = \sum_i y_i^2/n_i, \quad T_{AB} = \sum_i \sum_j y_{ij}^2/n_{ij}$$

$$T_0 = \sum_i \sum_j \sum_k y_{ijk}^2 \quad \text{and} \quad T_B = y^2/N$$

Then

$$\hat{\theta}_1^2 = (T_0 - T_{AB})/(N - b),$$

$$\hat{\theta}_2^2 = [T_{AB} - T_A - (b - a)\hat{\theta}_1^2]/(N - k_{12}),$$

$$\hat{\theta}_3^2 = [T_A - T_B - (k_{12} - k_2)\hat{\theta}_1^2 - (a - 1)\hat{\theta}_2^2]/(N - k_1)$$

(Searle, 1961).

c. Variances of analysis of variance estimators (under normality)

$$\text{var}(\hat{\theta}_1^2) = 2\sigma_1^4/(N - b).$$

Calculate

$$k_4 = \sum_i \sum_j n_{ij}^2, \quad k_5 = \sum_i (\sum_j n_{ij}^2/n_i),$$

$$k_6 = \sum_i (\sum_j n_{ij}^2)^2/n_i, \quad k_7 = \sum_i (\sum_j n_{ij}^2)^2/n_i^2,$$

$$k_8 = \sum_i n_i (\sum_j n_{ij}^2), \quad k_9 = \sum_i n_i^2$$

434 APPENDIX F [F.4]

PART II. THE 2-WAY CROSSED CLASSIFICATION

F.4. WITH INTERACTION, RANDOM MODEL

a. Model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$$

$$i = 1, 2, \dots, a, \quad j = 1, 2, \dots, b \quad \text{and} \quad k = 1, 2, \dots, n_{ij}$$

with

$$n_{ij} > 0 \quad \text{for } i(j)\text{-cells and } \sum_i \sum_j n_{ij} = N$$

b. Henderson Method I estimators

Calculate Table F.1 and, for  $n_{ij} > 0$ ,

$$T_0 = \sum_i \sum_j \sum_k y_{ijk}^2, \quad T_A = \sum_i y_i^2/n_i, \quad T_B = \sum_j y_j^2/n_j,$$

$$T_{AB} = \sum_i \sum_j y_{ij}^2/n_{ij} \quad \text{and} \quad T_C = y^2/N;$$

$$\text{SSA} = T_A - T_C, \quad \text{SSB} = T_B - T_C,$$

$$\text{SSAB}^* = T_{AB} - T_A - T_B + T_C,$$

$$\text{SSE} = T_0 - T_{AB}$$

TABLE F.1. ANALYSIS OF VARIANCE ESTIMATION OF VARIANCE COMPONENTS IN THE 2-WAY CROSSED CLASSIFICATION, INTERACTION, RANDOM MODEL

Terms needed for calculating estimators and their variances.

For estimators only, calculate  $k_1, k_2, k_3, k_4$  and  $k_{12}$ .

$\sqrt{k_1} = \sum_i n_i$	$\sqrt{k_2} = \sum_j n_j$
$\sqrt{k_3} = \sum_i (\sum_j n_{ij})/n_i$	$\sqrt{k_4} = \sum_j (\sum_i n_{ij})/n_j$
$k_5 = \sum_i n_i^2$	$k_6 = \sum_j n_j^2$
$k_7 = \sum_i (\sum_j n_{ij}^2)/n_i$	$k_8 = \sum_j (\sum_i n_{ij}^2)/n_j$
$k_9 = \sum_i (\sum_j n_{ij}^2)^2/n_i^2$	$k_{10} = \sum_j (\sum_i n_{ij}^2)^2/n_j^2$
$k_{11} = \sum_i (\sum_j n_{ij}^2)/n_i$	$k_{12} = \sum_j (\sum_i n_{ij}^2)/n_j$
$k_{13} = \sum_i (\sum_j n_{ij}^2)(\sum_k n_{ik})/n_i$	$k_{14} = \sum_j (\sum_i n_{ij}^2)(\sum_k n_{jk})/n_j$
$k_{15} = \sum_i (\sum_j n_{ij}^2)^2/n_i$	$k_{16} = \sum_j (\sum_i n_{ij}^2)^2/n_j$
$k_{17} = \sum_i (\sum_j n_{ij}^2)/n_i$	$k_{18} = \sum_j (\sum_i n_{ij}^2)/n_j$
$k_{19} = \sum_i (\sum_j n_{ij}^2)/n_i$	$k_{20} = \sum_j (\sum_i n_{ij}^2)/n_j$
$k_{21} = \sum_i \sum_j n_{ij}^2/n_i n_j$	$k_{22} = \sum_j \sum_i n_{ij}^2/n_i n_j$

$\sqrt{k_{13}} = \sum_i \sum_j n_{ij}$

$\sqrt{k_{14}} = \sum_j \sum_i n_{ij}$

$k_{23} = \sum_i \sum_j n_{ij}^2/n_i n_j$

$k_{24} = \sum_j \sum_i n_{ij}^2/n_i n_j$

$\sqrt{k_1} = \sum_i n_i$  for all  $i$ .

[F.4] WITH INTERACTION, RANDOM MODEL 435

Then

$$\hat{\theta}_1^2 = \text{SSE}/(N - s) = \text{MSE}$$

and with

$$P = \begin{bmatrix} N - k_1 & k_2 - k_1 & k_3 - k_{12} \\ k_4 - k_1 & N - k_2 & k_4 - k_{12} \\ k_5 - k_4 & k_6 - k_4 & N - k_5 - k_4 + k_{12} \end{bmatrix}$$

$$\hat{\theta}^2 = \begin{bmatrix} \hat{\theta}_1^2 \\ \hat{\theta}_2^2 \\ \hat{\theta}_3^2 \end{bmatrix} = P^{-1} \begin{bmatrix} \text{SSA} - (a - 1)\text{MSE} \\ \text{SSB} - (b - 1)\text{MSE} \\ \text{SSAB}^* - (s - a - b + 1)\text{MSE} \end{bmatrix}$$

as in (32) of Section 5.3b. This is equivalent to calculating

$$\delta_A = [\text{SSB} + \text{SSAB}^* - (s - a)\text{MSE}]/(N - k_1)$$

and

$$\delta_B = [\text{SSA} + \text{SSAB}^* - (s - b)\text{MSE}]/(N - k_2)$$

with which

$$\hat{\theta}_1^2 = [(N - k_1)\delta_A + (k_2 - k_1)\delta_B - \{\text{SSA} - (a - 1)\text{MSE}\}]/(N - k_1 - k_2 + k_{12}),$$

$$\hat{\theta}_2^2 = \delta_A - \hat{\theta}_1^2 \quad \text{and} \quad \hat{\theta}_3^2 = \delta_B - \hat{\theta}_1^2$$

(Searle, 1958).

c. Variances of Henderson Method I estimators (under normality)

$$\text{var}(\hat{\theta}_1^2) = 2\sigma_1^4/(N - s)$$

For P given above and for H and I being

$$H = \begin{bmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ -1 & -1 & 1 & 1 \end{bmatrix} \quad \text{and} \quad I = \begin{bmatrix} a - 1 \\ b - 1 \\ s - a - b + 1 \end{bmatrix}$$

$$\text{var}(\hat{\theta}^2) = P^{-1} [\text{H var}(t) \text{H} + \text{var}(\hat{\theta}_1^2) \text{I}] P^{-1}$$

and

$$\text{cov}(\hat{\theta}_1^2, \hat{\theta}_2^2) = -P^{-1} [\text{var}(\hat{\theta}_1^2) \text{I}]$$

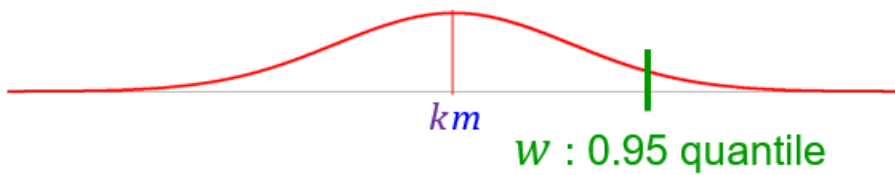
where

$$\text{var}(t) = \text{var}[T_A \quad T_B \quad T_{AB} \quad T_C]$$

Var(t) has 10 different elements, each element is a function of the 10 squares and products of  $\sigma_1^2, \sigma_2^2, \sigma_3^2$  and  $\sigma_1^2$ . The  $10 \times 10$  matrix of these coefficients is shown in Table F.2. Apart from  $N, a, b, s$  and unity, Table F.2 involves only

# 1. Calculations details – 0.95 quantile weight

$X$ : “2500 ( $k = 25$ ) or 25000 ( $k = 250$ ) seeds weight”



**Reporting** →

- If  $w < 1 \text{ g}$ ,  $w$  is rounded up to the nearest multiple of 0.01
- If  $1 \text{ g} \leq w < 5 \text{ g}$ ,  $w$  is rounded up to the nearest multiple of 0.1
- If  $w \geq 5 \text{ g}$ ,  $w$  is rounded up to the nearest integer

Examples:

$w$	Value reported
0.34567	0.35
0.96781	0.97
0.99001	1.00
1.08962	1.1
4.45687	4.5
5.00768	6
9.76981	10



# 1. Overview of the calculator

### Calculator for adding working weights to Table 2C of the ISTA Rules

THE CALCULATOR IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND. IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY ARISING IN CONNECTION WITH THE CALCULATOR.

#### Experiment designs

Two types of experiment designs are considered in the calculator:

#### Experiment design 1: 2-way nested design

A minimum of 12 lots are considered across a minimum of six varieties represented in the experiment as a general rule. These 12 lots will be evaluated by a minimum of two labs however six labs are preferred. A minimum of eight 100 seed reps are weighed per lot.

#### Experiment design 2: 2-way crossed design

A minimum of six lots are considered across a minimum of three varieties represented in the experiment as a general rule. These six lots will be evaluated each by a minimum of two labs however six labs are preferred. A minimum of eight 100 seed reps are weighed per lot.

#### Calculations

The rep weights are entered into the unprotected yellow cells of the calculator. If experiment design 1 is used, data for the different lots from each lab are entered in different columns. In order to avoid conditional formatting conflicts, always copy/paste data in the calculator using Paste Special -> Values.

- For each lot in a given laboratory, outliers are highlighted in red using Grubbs's method at the 5% significance level (Grubbs, 1969). These outliers are then excluded manually from the computations.
- The linear random effects models used for the analysis of the two experiment designs are:
  - Experiment design 1:**

$$y_{ijk} = \mu + \alpha_i + \beta_{ij} + \epsilon_{ijk}$$

in which:

    - $y_{ijk}$  is the observed 100-seeds weight of lot  $j$  ( $j = 1, 2, \dots, b_j$ ) in lab  $i$  ( $i = 1, 2, \dots, a$ ) and replication  $k$  ( $k = 1, 2, \dots, n_{ij}$ );
    - $\mu$  is the intercept;
    - $\alpha_i$  is the random effect of lab  $i$  ( $\alpha_i \sim \text{i.i.d. } N(0, \sigma_{lab}^2)$ );
    - $\beta_{ij}$  is the random effect of lot  $j$  within lab  $i$  ( $\beta_{ij} \sim \text{i.i.d. } N(0, \sigma_{lot}^2)$ );
    - $\epsilon_{ijk}$  is the residual ( $\epsilon_{ijk} \sim \text{i.i.d. } N(0, \sigma_{res}^2)$ ).
  - Experiment design 2:**

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$$

in which:

    - $y_{ijk}$  is the observed 100-seeds weight of lot  $j$  ( $j = 1, 2, \dots, b$ ) in lab  $i$  ( $i = 1, 2, \dots, a$ ) and replication  $k$  ( $k = 1, 2, \dots, n_{ij}$ );
    - $\mu$  is the intercept;
    - $\alpha_i$  is the random effect of lab  $i$  ( $\alpha_i \sim \text{i.i.d. } N(0, \sigma_{lab}^2)$ );
    - $\beta_j$  is the random effect of lot  $j$  ( $\beta_j \sim \text{i.i.d. } N(0, \sigma_{lot}^2)$ );
    - $(\alpha\beta)_{ij}$  is the random interaction effect between lab  $i$  and lot  $j$  ( $(\alpha\beta)_{ij} \sim \text{i.i.d. } N(0, \sigma_{lab \times lot}^2)$ );
    - $\epsilon_{ijk}$  is the residual ( $\epsilon_{ijk} \sim \text{i.i.d. } N(0, \sigma_{res}^2)$ ).

The calculator automatically selects which model to fit according to the dataset structure.

Variance components for the two models are estimated from the data by the Henderson Method I (Searle et al., 1992, Appendix F). When an estimate is negative, this estimate is reported as zero. Let  $\hat{\sigma}_{lab}^2$ ,  $\hat{\sigma}_{lot}^2$ ,  $\hat{\sigma}_{lab \times lot}^2$  and  $\hat{\sigma}_{res}^2$  be these

Instructions Calculator

### Supporting Data of New Species Proposal to ISTA Rules Table 2C

Submitter Name:		Lab Full Name:	
Scientific Name of the Crop kind: Genus Species		ISTA Member Code:	
		Contact Email:	

Change any value in a yellow cell

Number of observations	0
Number of labs	0
Number of lots	0
General mean	
Lab variance	
Lot variance	
Lab x Lot variance	
Residual variance	
2500 seed weight*	
25000 seed weight*	

\* 95% Confidence

Decision

Rep weights in red are identified as outliers by Grubbs's method at the 5% significance level and needs to be suppressed (removed) manually

Lab \ Seed lot	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
Lab 1	Rep1																
	Rep2																
	Rep3																
	Rep4																
	Rep5																
	Rep6																
	Rep7																
	Rep8																
	Rep9																
	Rep10																
	Rep11																
	Rep12																
Mean																	
St. Dev.																	
Number of reps	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Grubbs critical values																	
Lab 2	Rep1																
	Rep2																
	Rep3																
	Rep4																
	Rep5																
	Rep6																
	Rep7																
	Rep8																
	Rep9																
	Rep10																
	Rep11																

Instructions Calculator

# 1. Overview of the calculator

- The spreadsheet is protected (no password): entering data is only possible in **yellow cells**
- In order to avoid **conditional formatting conflicts**, always copy/paste data in the calculator using Paste Special → Values

When needed, some warnings are displayed in **red**

Supporting Data of New Species Proposal to ISTA Rules Table 2C																	
Submitter Name: XXX			Lab Full Name: YYY			Number of observations: 232		Number of labs: 7		Number of lots: 5		General mean: 3.2584		Lab variance: 0.0020037		Lot variance: 0.1453077	
Scientific Name of the Crop kind: Basella B. alba			ISTA Member Code: ZZZ			Residual variance: 0.0101622		Lab x Lot variance: 0.0963083		Decision		2500 seed weight*: 103		25000 seed weight*: 1022		1050	
Contact Email: AAA			Change any value in a yellow cell			* 95% Confidence		6 lots are preferred for an accurate estimation		Decision value should be greater than or equal to 103							
Rep weights in red are identified as outliers by Grubbs's method at the 5% significance level and needs to be suppressed (removed) manually																	
Lab \ Seed lot	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	
Lab 2	Rep1	2.3672	3.4036	2.3585	3.1927	3.7473											
	Rep2	2.2734	3.4207	2.3530	3.0972	3.7309											
	Rep3	2.3198	3.5878	2.4268	3.2861	3.7818											
	Rep4	2.3866	3.4322	2.3827	3.2858	3.7380											
	Rep5	2.3600	3.3296	2.2917	3.2861	3.7908											
	Rep6	2.3720	3.3873	2.2663	3.1889	3.8542											
	Rep7		3.4601	2.2407	3.1103												
	Rep8			2.7779	3.2201												
	Rep9																
	Rep10																
	Rep11																
	Rep12																
Mean	2.3465	3.4316	2.3872	3.2084	3.7738												
St. Dev.	0.04225	0.08008	0.16963	0.07636	0.04613												
Number of reps	6	7	8	8	6	0	0	0	0	0	0	0	0	0	0		
Grubbs critical values	1.89	2.02	2.13	2.13	1.89												

Outliers are automatically identified in **red**



# 1. Overview of the calculator

Although not recommended, the calculator can provide estimates of the variance components when there is only **1 lab**

Random effects model:  $100\_seeds\_weight = general\_mean + Lot\_effect + Residual$

$\sim \text{i.i.d. } N(0, \sigma_{Lot}^2)$   
 $\sim \text{i.i.d. } N(0, \sigma_{Res}^2)$

Balanced dataset example:

Number of observations	224	
Number of labs	1	<b>6 labs are preferred for an accurate estimation</b>
Number of lots	28	
General mean	0.4261	
<del>Lab variance</del>		
Lot variance	0.0040998	←
<del>Lab x Lot variance</del>		
Residual variance	0.0003288	← <b>Decision</b>
2500 seed weight*	14	
25000 seed weight*	134	

REML estimates (R package lme4)

Number of obs: 224, groups: Lot, 28

Random effects:

Groups	Variance
Lot	0.0040998
Residual	0.0003288

\* 95% Confidence



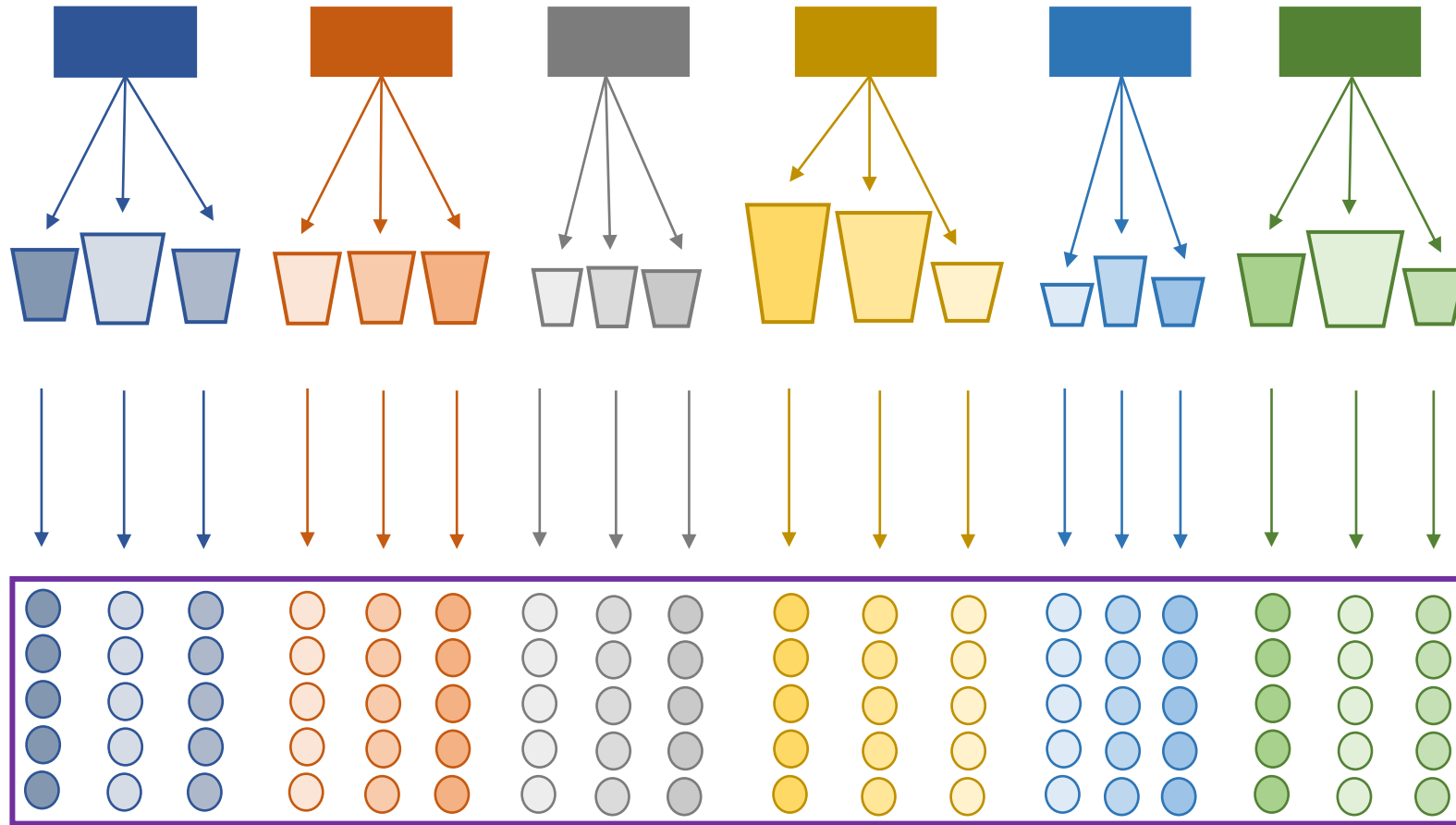
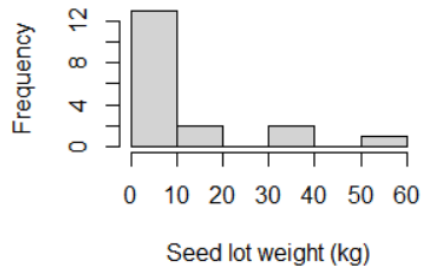


2. Number of sub-lots for which an OIC established for the lot is still valid

## 2. 2021 Tomato experiment – Experiment design

6 companies

3 seed lots/company  
produced in 11  
different countries  
Different sizes:

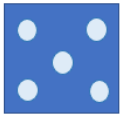


90 samples

**Question:**  
*are the 5  
sub-samples  
results  
homogeneous?*

**Measurements:**

- Purity test
- Germination test: 1<sup>st</sup> count at day 6,  
final count at days 8 to 14





Purity % are all equal to 100%



1. Homogeneity of the test replications for normal seedlings, final count:

→ for each of the **90** samples, the **4** reps are within ISTA tolerances

2. Heterogeneity of the **5** samples from each lot

→ Use of the  $H$  statistic:

$$H = \frac{\overset{400}{\#\_of\_seeds\_in\_the\_sample} \times (\overset{5}{\#\_of\_subsamples} - 1) \times \text{observed\_subsample\_variance}}{\text{mean} \times (100 - \text{mean})}$$

→  $H$  has a chi-squared distribution

→ Statistical test: p-value that all the sample values are equal

→ The lower the p-value, the greater the statistical evidence for **heterogeneity**

## 2. 2021 Tomato experiment – Analysis – Sub-samples homogeneity for germination



Company	Lot weight	Normal seedlings %, 1 <sup>st</sup> count			Normal seedlings %, final count			Abnormal seedlings %			Dead seeds %		
		Mean	H	p-value	Mean	H	p-value	Mean	H	p-value	Mean	H	p-value
A	3kg	85.8	8.80	0.0663	96.0	4.17	0.3839	3.6	1.38	0.8471	0.4	12.05	0.0170
A	5.9kg	76.6	18.57	0.0010	92.0	2.17	0.7038	3.0	2.75	0.6006	5.0	1.68	0.7936
A	7.7kg	84.8	17.01	0.0019	89.0	10.62	0.0311	8.2	2.55	0.6356	2.8	15.87	0.0032
B	1.5kg	87.2	5.30	0.2575	89.6	1.37	0.8488	3.0	0.00	1.0000	7.4	1.87	0.7600
B	20kg	78.2	16.61	0.0023	98.2	1.81	0.7706	1.8	1.81	0.7706	0.0		
B	6kg	63.8	12.26	0.0155	92.0	5.43	0.2455	4.6	8.39	0.0784	3.4	3.90	0.4201
C	5.9kg	42.0	46.63	0.0000	98.6	3.48	0.4813	1.2	2.70	0.6094	0.2	16.03	0.0030
C	6.4kg	34.8	18.48	0.0010	89.6	4.81	0.3076	7.6	2.96	0.5642	2.8	4.12	0.3906
C	7.8kg	60.6	41.08	0.0000	98.4	8.13	0.0869	1.0	8.08	0.0887	0.6	8.05	0.0898
D	2.1kg	71.4	18.26	0.0011	96.4	3.69	0.4498	1.6	3.05	0.5497	2.0	4.08	0.3951
D	3.3kg	96.0	27.08	0.0000	98.2	6.34	0.1754	1.0	8.08	0.0887	0.8	4.03	0.4017
D	3.7kg	17.8	29.74	0.0000	96.8	6.20	0.1848	1.8	6.34	0.1754	1.4	3.48	0.4813
E	13kg	37.8	29.74	0.0000	97.4	5.05	0.2818	1.6	8.13	0.0869	1.0	8.08	0.0887
E	6.8kg	90.2	8.51	0.0747	96.4	1.38	0.8471	1.4	3.48	0.4813	2.2	5.21	0.2669
E	7.8kg	74.6	30.23	0.0000	98.4	3.05	0.5497	1.2	2.70	0.6094	0.4	12.05	0.0170
F	32kg	84.6	18.79	0.0009	94.2	18.16	0.0012	1.2	2.70	0.6094	3.4	23.38	0.0001
F	36kg	90.2	3.98	0.4084	95.8	6.76	0.1491	2.0	4.08	0.3951	2.2	5.21	0.2669
F	51kg	97.0	5.50	0.2399	97.0	5.50	0.2399	1.6	3.05	0.5497	1.4	3.48	0.4813

Evidence for  
**heterogeneity**

Evidence for  
**homogeneity**

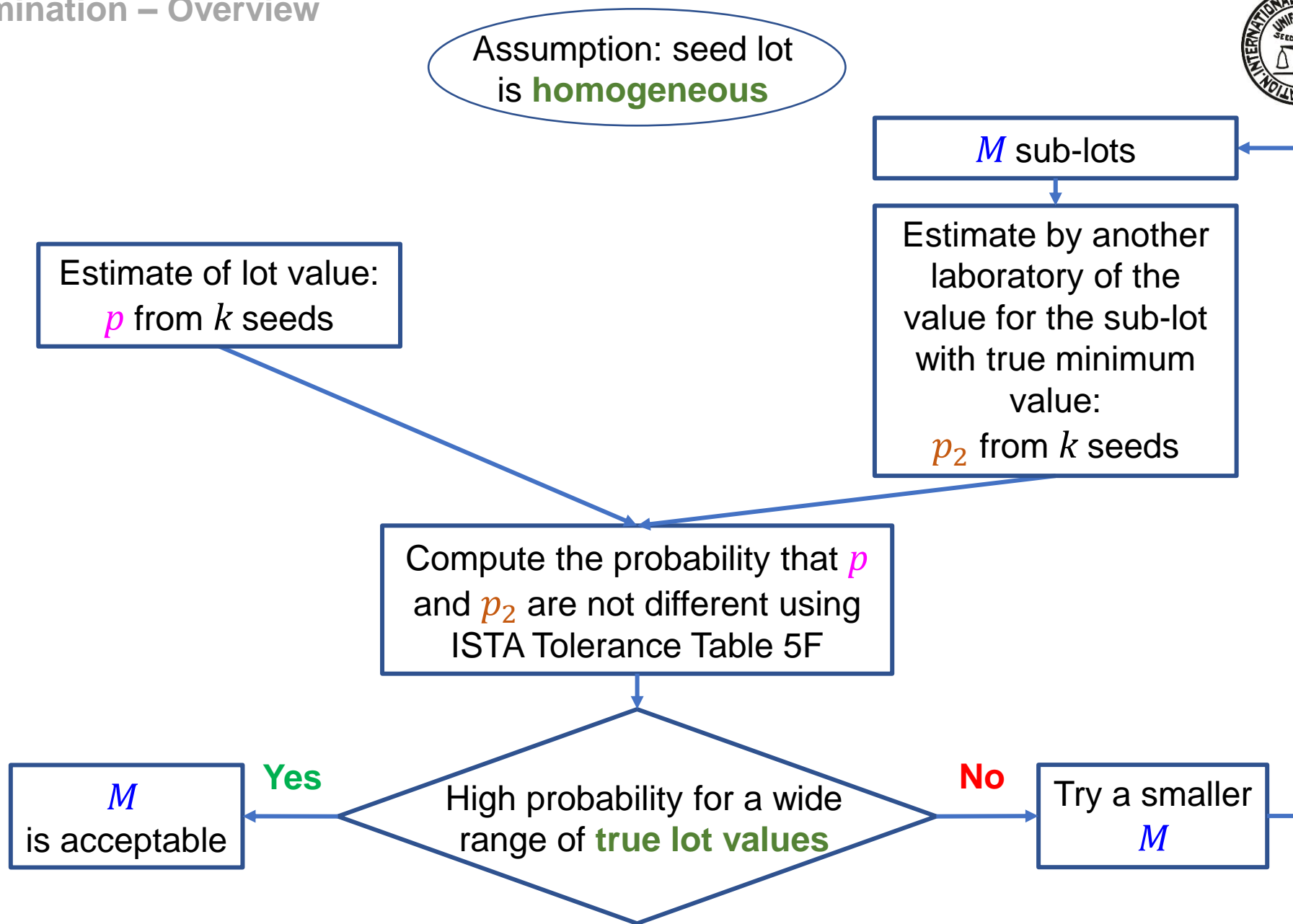
## 2. 2021 Tomato experiment – Analysis – Sub-samples homogeneity for germination



- Homogeneity for the germination final count is reinforced by the **R test**:

Company	Lot weight	Normal seedlings %, final count				
		Mean	H	p-value H test	Range	p-value R test
A	3kg	96.0	4.17	0.3839	2.0	0.5995
A	5.9kg	92.0	2.17	0.7038	1.5	0.8355
A	7.7kg	89.0	10.62	0.0311	3.8	0.0522
B	1.5kg	89.6	1.37	0.8488	1.3	0.8867
B	20kg	98.2	1.81	0.7706	1.5	0.8251
B	6kg	92.0	5.43	0.2455	2.9	0.2264
C	5.9kg	98.6	3.48	0.4813	1.7	0.7493
C	6.4kg	89.6	4.81	0.3076	2.6	0.3430
C	7.8kg	98.4	8.13	0.0869	3.2	0.1601
D	2.1kg	96.4	3.69	0.4498	2.1	0.5505
D	3.3kg	98.2	6.34	0.1754	3.0	0.2083
D	3.7kg	96.8	6.20	0.1848	3.4	0.1124
E	13kg	97.4	5.05	0.2818	2.5	0.3868
E	6.8kg	96.4	1.38	0.8471	1.1	0.9422
E	7.8kg	98.4	3.05	0.5497	1.6	0.7922
F	32kg	94.2	18.16	0.0012	6.0	0.0002
F	36kg	95.8	6.76	0.1491	3.0	0.2135
F	51kg	97.0	5.50	0.2399	2.3	0.4602

## 2. Number of sub-lots determination – Overview



**Table 5F.** Tolerances between results of two tests made in different laboratories on the same or different samples from the same seed lot (two-way test at 5 % significance level) on 400 seed tests. Updated by ISTA Statistics Technical Committee, based on Miles (1963) Table G5, column C, 400 seed tests.

Average germination percentage of 2 tests		Tolerance
51–100 %	0–50 %	
99	2	2
98	3	3
96–97	4–5	4
94–95	6–7	5
91–93	8–10	6
88–90	11–13	7
84–87	14–17	8
79–83	18–22	9
74–78	23–27	10
68–73	28–33	11
60–67	34–41	12
51–59	42–50	13

## 2. Number of sub-lots determination – Some details

Assumption: seed lot is **homogeneous**

True lot value:  $\pi$

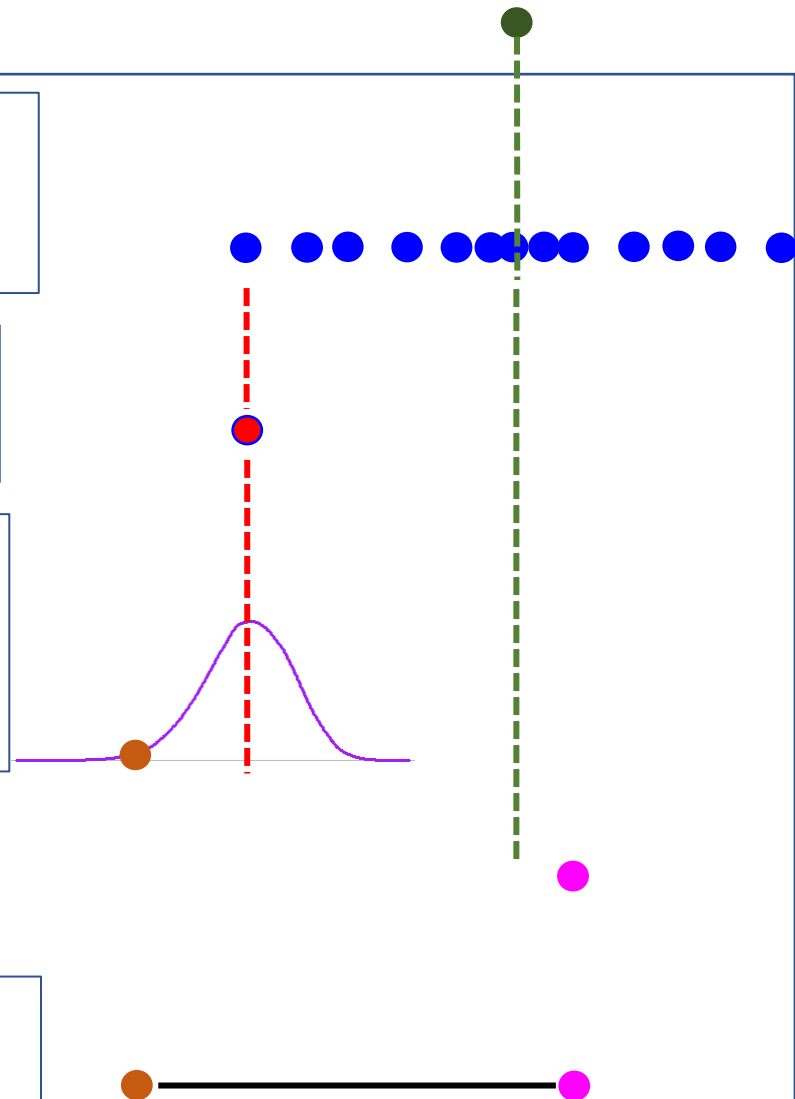
$M$  true sub-lot values  $\pi_i$  from multivariate hypergeometric distribution \*

Minimum of the true  $M$  sub-lot values:  $\pi_m = \min(\{\pi_i\}_{i=1}^M)$

Test value from another lab on the sub-lot with minimum value from a **Beta-binomial** distribution:  $p_2$

Estimate of lot value:  $p$  from  $k$  seeds

Determine if  $p$  and  $p_2$  are within tolerance according to ISTA Tolerance Table 5F



Repeated 10,000 times

Compute the probability that  $p$  is not different from  $p_2$  as  $\frac{T}{10000}$  where  $T$  is the number of times  $p$  and  $p_2$  are within tolerance

Table 5F. Tolerances between results of two tests made in different laboratories on the same or different samples from the same seed lot (two-way test at 5% significance level) on 400 seed tests. Updated by ISTA Statistics Technical Committee, based on Miles (1963) Table G5, column C, 400 seed tests.

Average germination percentage of 2 tests		Tolerance
51–100 %	0–50 %	
99	2	2
98	3	3
96–97	4–5	4
94–95	6–7	5
91–93	8–10	6
88–90	11–13	7
84–87	14–17	8
79–83	18–22	9
74–78	23–27	10
68–73	28–33	11
60–67	34–41	12
51–59	42–50	13

\* Laffont, J-L., Hong, B., Kuo, B-J. and K.M. Remund (2019). Exact theoretical distributions around the replicate results of a germination test. *Seed Science Research* **29**, 64-72.



## 2. Number of sub-lots determination – Some details

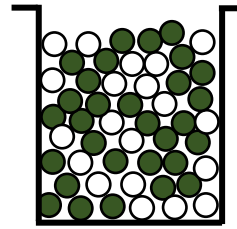
Seed Science Research

Exact theoretical distributions around the replicate results of a germination test

cambridge.org/ssr

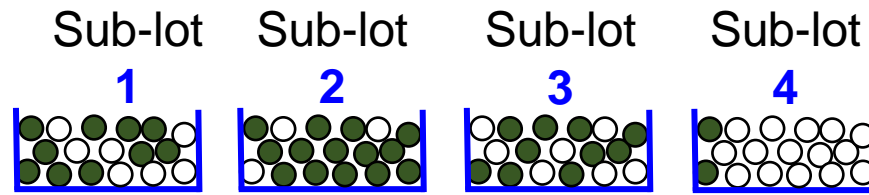
Jean-Louis Laffont<sup>1</sup>, Bonnie Hong<sup>2</sup>, Bo-Jein Kuo<sup>3</sup> and Kirk M. Remund<sup>4</sup>

Seed lot:  $N$  seeds with  
 $G = N\pi$  seeds to germinate



$N = 1,600,000$   
 $\pi = 90\%$   
 $G = 1,440,000$

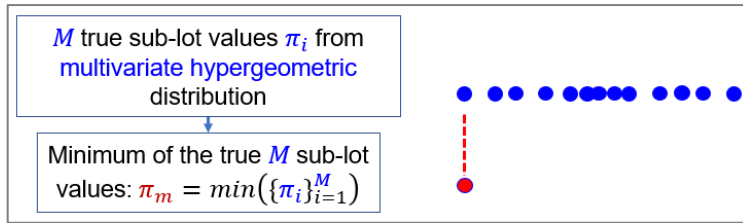
Random assignment  
of the seeds into  
 $M = 4$  sub-lots  
( $n = 400,000$  seeds  
per sub-lot)



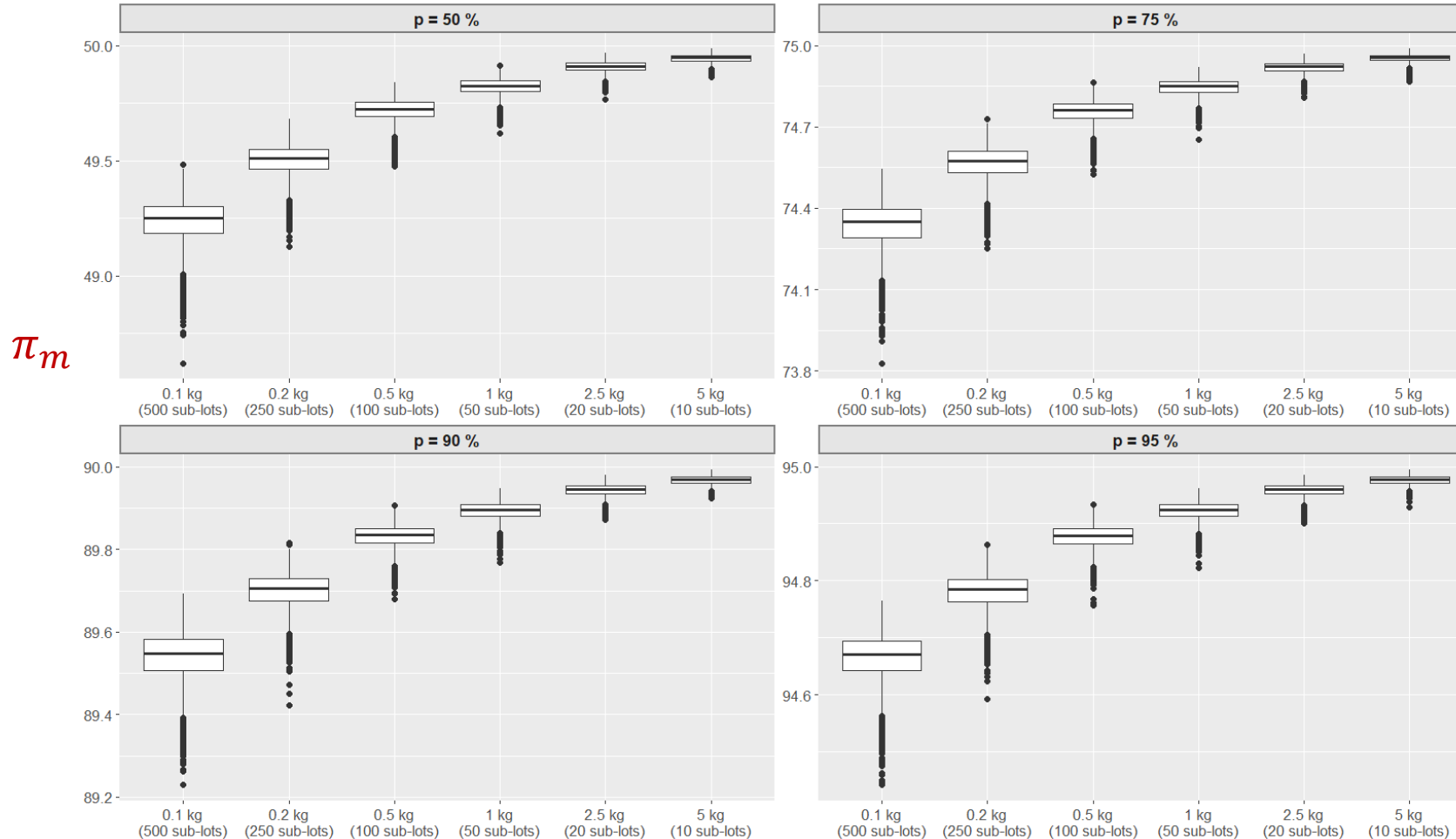
# of seeds to germinate in each sub-lot	Sub-lot 1	Sub-lot 2	Sub-lot 3	Sub-lot 4
	360,023	360,007	359,968	360,002
	359,270	312,987	368,976	398,767
	360,030	292,345	387,969	399,656
	...			

Distribution of the number  
of seeds to germinate  
in the sub-lots:  
**multivariate hypergeometric**  
distribution

## 2. Number of sub-lots determination – Some details



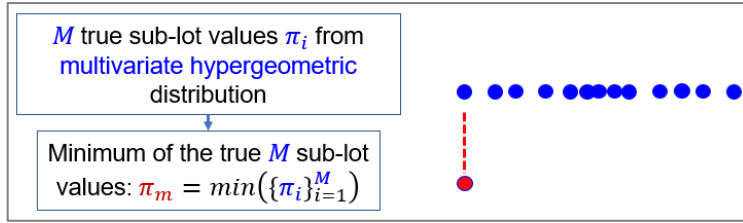
Lot weight = 50 kg



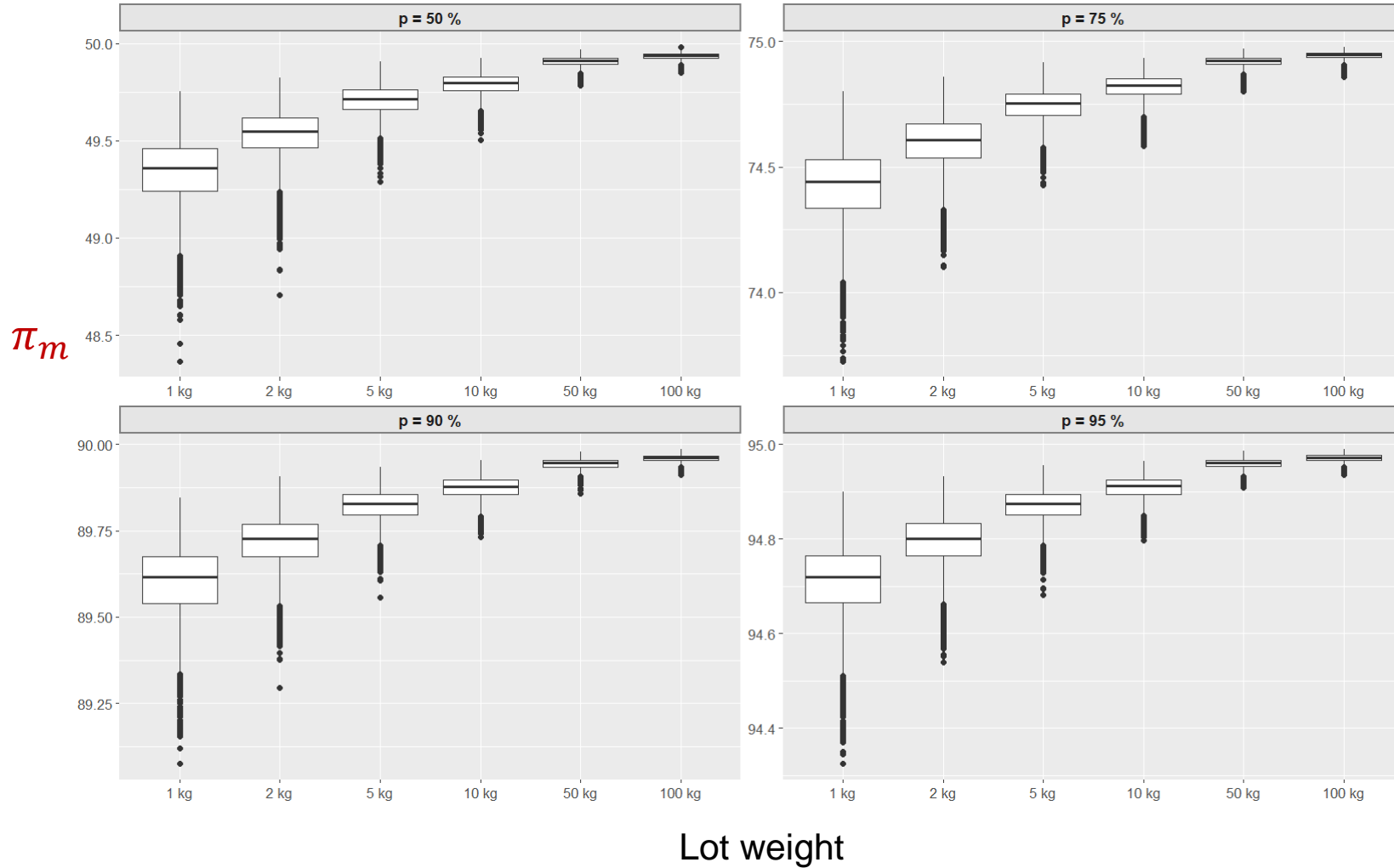
The larger the number of sub-lots  $M$ , the lower the minimum values  $\pi_m$

Sub-lots weight/Number of sub-lots ( $M = \frac{\text{Lot weight}}{\text{Sub-lot weight}}$ )

## 2. Number of sub-lots determination – Some details

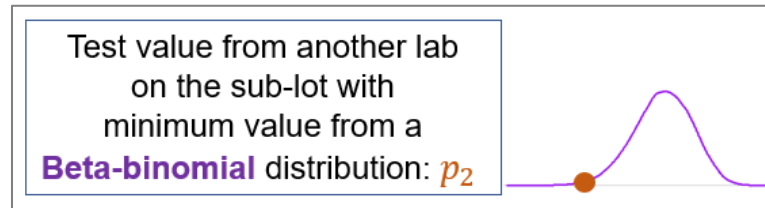


Number of sub-lots = 20



The lower the lot weight,  
the lower the  
minimum values  $\pi_m$

## 2. Number of sub-lots determination – Some details



### Refresher:

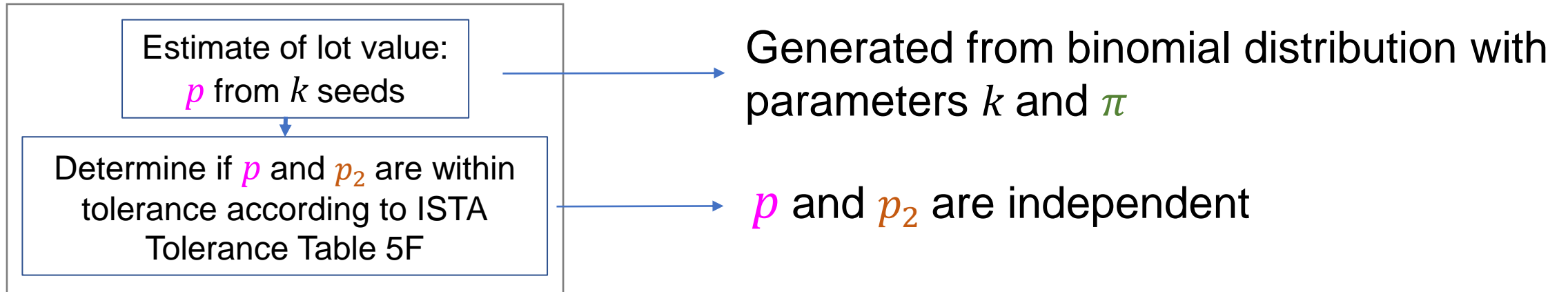
- The result from a different lab is not from a **Binomial**( $k, \pi_m$ ) but from a distribution with a variance larger than the binomial variance
- The over-dispersion has been quantified by Miles (1963) and is taken into account in tolerance tables for comparing different laboratories

$$\sqrt{\frac{\text{Over\_dispersed\_variance}}{\text{Binomial\_variance}}} = f = 2.38 - 0.8321\pi_m$$

A model for generating over-dispersed binomial data is the **Beta-binomial** model with parameters  $k$ ,  $\alpha$  and  $\beta$ :

$$\alpha = \pi_m \left( \frac{k-1}{f^2-1} - 1 \right) \quad \beta = \alpha \left( \frac{1}{\pi_m} - 1 \right)$$

## 2. Number of sub-lots determination – Some details



## 2. Number of sub-lots determination – Results

	$W_{Lot} = 1.5 \text{ kg}, W_{Sub} = 0.1 \text{ kg} (M = 15)$	$W_{Lot} = 50 \text{ kg}, W_{Sub} = 1 \text{ kg} (M = 50)$	$W_{Lot} = 50 \text{ kg}, W_{Sub} = 0.1 \text{ kg} (M = 500)$
$\pi$ (%)	Prob( $p$ and $p_2$ are within Tol)	Prob( $p$ and $p_2$ are within Tol)	Prob( $p$ and $p_2$ are within Tol)
50	0.9865	0.9857	0.9887
55	0.9858	0.9854	0.9836
60	0.9870	0.9871	0.9867
65	0.9852	0.9866	0.9854
70	0.9850	0.9853	0.9860
75	0.9840	0.9845	0.9842
80	0.9848	0.9858	0.9832
85	0.9861	0.9851	0.9849
90	0.9861	0.9892	0.9887
95	0.9893	0.9894	0.9885
99	0.9951	0.9959	0.9958

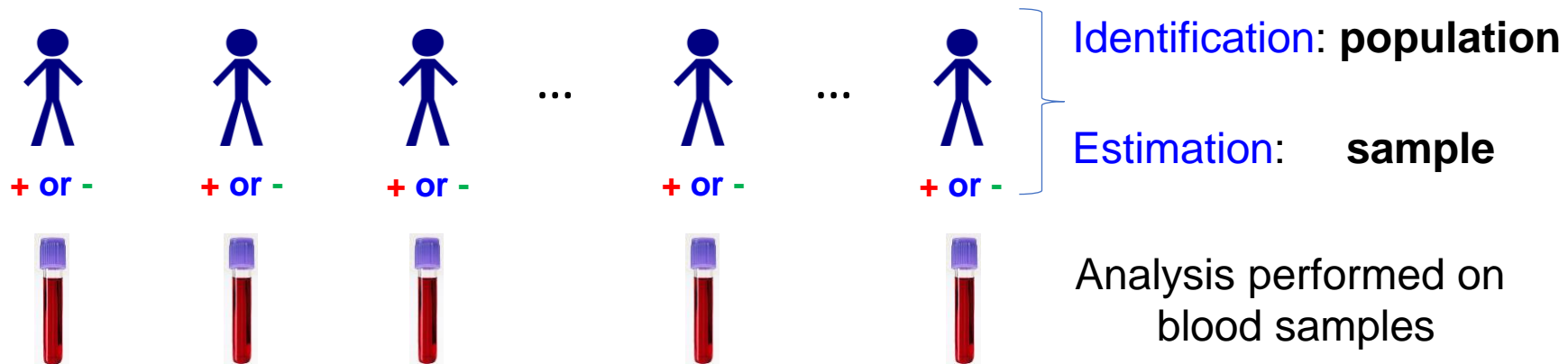
- For two extreme lot sizes (1.5 kg and 50 kg) and different number of sub-lots (15, 50 and 500), all the probabilities are very high (**above 0.98**)
- Evidence that given that the original lot is **homogeneous**, there is **no limit** in the number of sub-lots that can be elaborated from it

3. Group testing: number of groups to ensure that estimation is possible

### 3. What is group testing

- Suppose people are tested for a disease
- Who has the disease? → **identification**  
 What is the **prevalence** of the disease? → **estimation**  
 (i.e. what is the proportion of people with the disease?)

- One solution: individual testing

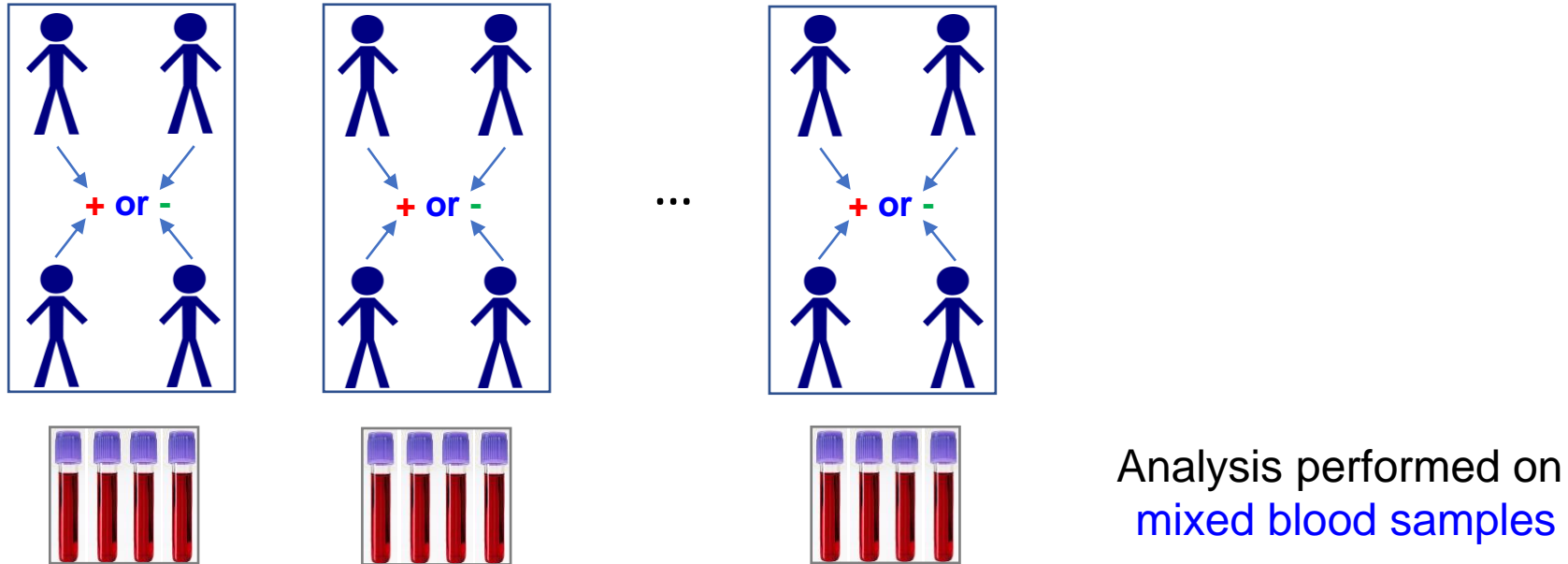


- **Problem:** can be expensive
- **Group testing:** **cost savings**



### 3. What is group testing

- **Group testing:**



- **Identification:** if the group is **positive**, the individuals making up the group are retested to determine which of the members have the disease.
- **Estimation:**  $\hat{p}$



### 3. What is group testing

- **Identification:** original development of group testing by Robert Dorfman in 1943:

The Detection of Defective Members of Large Populations

Author(s): Robert Dorfman

Source: *The Annals of Mathematical Statistics*, Vol. 14, No. 4 (Dec., 1943), pp. 436-440

The inspection of the individual members of a large population is an expensive and tedious process. Often in testing the results of manufacture the work can be reduced greatly by examining only a sample of the population and rejecting the whole if the proportion of defectives in the sample is unduly large. In many inspections, however, the objective is to eliminate all the defective members of the population. This situation arises in manufacturing processes where the defect being tested for can result in disastrous failures. It also arises in certain inspections of human populations. Where the objective is to weed out individual defective units, a sample inspection will clearly not suffice. It will be shown in this paper that a different statistical approach can, under certain conditions, yield significant savings in effort and expense when a complete elimination of defective units is desired.

...

The method will be described by showing its application to a large-scale project on which the United States Public Health Service and the Selective Service System are now engaged. The object of the program is to weed out all syphilitic men called up for induction. Under this program each prospective inductee is subjected to a "Wasserman-type" blood test. The test may be divided conveniently into two parts:

...

### 3. Group testing estimation

- $p$  : proportion of individuals in the population with the attribute
- $n$  : number of groups
- $m$  : number of individuals per group
- $X$  : number of **positive** groups out of  $m$

$$\hat{p} = 1 - \left(1 - \frac{x}{n}\right)^{\frac{1}{m}}$$

When all the groups are **positive**,  $\hat{p} = 1$  and the result is not considered.

<p><math>Y_i</math> : random variable = 1 if group <math>i</math> is <b>positive</b> = 0 if group <math>i</math> is <b>negative</b></p> <p>→ <math>Y_i</math> has a <b>Bernoulli distribution</b> with pmf :</p> $P(Y_i = y) = \begin{cases} (1-p)^m & \text{if } y = 0 \\ 1 - (1-p)^m & \text{if } y = 1 \end{cases}$ <p><small><math>P(\text{individual 1 neg} \ \&amp; \ \text{individual 2 neg} \ \&amp; \ \dots \ \&amp; \ \text{individual } m \ \text{neg})</math> = <math>P(\text{individual 1 neg}) \times P(\text{individual 2 neg}) \times \dots \times P(\text{individual } m \ \text{neg})</math> = <math>(1-p) \times (1-p) \times \dots \times (1-p)</math></small></p> <p>→ <math>X = \sum_{i=1}^n Y_i</math> has a <b>binomial distribution</b> <math>B(n, 1 - (1-p)^m)</math></p>	<p><b>Maximum Likelihood Estimation:</b></p> <p><b>Functional invariance</b> of MLEs property: If <math>\hat{\theta}</math> is the MLE for <math>\theta</math>, and if <math>g(\theta)</math> is any transformation of <math>\theta</math>, then the MLE for <math>\alpha = g(\theta)</math> is : <math>\hat{\alpha} = g(\hat{\theta})</math>.</p> <p>Let <math>\theta = 1 - (1-p)^m</math>.</p> <p>Rearranging: <math>p = 1 - (1-\theta)^{\frac{1}{m}} = g(\theta)</math></p> <p>The random variable <math>X</math> giving the number of <b>positive</b> groups <math>x</math> out of <math>n</math> has a <b>binomial distribution</b> <math>B(n, \theta)</math>.</p> <p>Then : <math>\hat{\theta} = \frac{x}{n}</math> is the MLE for <math>\theta</math>.</p> <p>→ <math>\hat{p} = g(\hat{\theta}) = 1 - \left(1 - \frac{x}{n}\right)^{\frac{1}{m}}</math></p>
---	--



### 3. Group testing estimation

- Taking into account **assay errors**:

$\lambda$  : **false negative** (group tests **neg** when at least 1 individual is **pos**) rate = 1 – **sensitivity**

$\delta$  : **false positive** (group tests **pos** when all individuals are **neg**) rate = 1 – **specificity**

$$\hat{p} = 1 - \left( 1 - \frac{\frac{x}{n} - \delta}{1 - \lambda - \delta} \right)^{1/m}$$

with  $x$  being the number of groups out of  $n$  testing **positive**

See:

**Statistical considerations in seed purity testing for transgenic traits**

*Seed Science Research* (2001) 11, 101–119

Kirk M. Remund<sup>1\*</sup>, Doris A. Dixon<sup>1</sup>, Deanne L. Wright<sup>2</sup> and Larry R. Holden<sup>3</sup>

### 3. Group testing estimation example

Seed lot



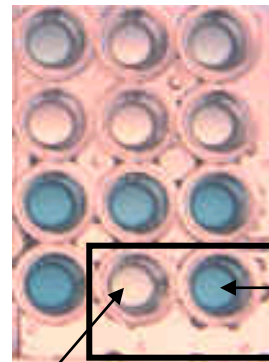
**10** groups of **150** seeds are taken from the lot using an appropriate sampling procedure (e.g. ISTA)



Seeds are ground into flour



Each sample is tested for **presence/absence** of GM seeds



Negative control

Positive control

**6 groups are negative**

**4 groups are positive**

### 3. Group testing estimation example

**6 groups of  
150 seeds  
are negative**

**4 groups of  
150 seeds  
are positive**

Point estimate  
of  $p$ , the true  
proportion of  
GM seeds in  
the lot

$$\hat{p} = 1 - \left(1 - \frac{4}{10}\right)^{\frac{1}{150}}$$
$$= 0.34\%$$

### 3. 2023 project around group testing

- $p$  : proportion of individuals in the population with the attribute
- $n$  : number of groups
- $m$  : number of individuals per group
- $x$  : number of **positive** groups out of  $m$

$$\hat{p} = 1 - \left(1 - \frac{x}{n}\right)^{\frac{1}{m}}$$

When all the groups are **positive**,  
**estimation is not possible!**



The probability that all groups are positive could help to ensure that the testing plan  $(n, m)$  will ensure estimation of  $p$



Probability that all groups are positive: 2 cases

1. Infinite population size (e.g. seed lot)
2. Finite population size (e.g. sample distributed for a Proficiency Test)





# 1. Infinite population size: **easy**

$k$  groups of  $m$  balls are sampled from a population of balls with a proportion  $\pi$  of white balls. The random variable  $Y_i$  “number of white balls in group  $i$ ” has a binomial distribution with parameters  $m$  and  $\pi$ :

$$P(Y_i = n_i) = \binom{m}{n_i} \pi^{n_i} (1 - \pi)^{m - n_i} .$$

Let  $A_i$  be the event “the  $i^{\text{th}}$  group has at least one white ball”. Then, the probability that the 1<sup>st</sup> group is positive is:

$$P(A_1) = 1 - \binom{m}{0} \pi^0 (1 - \pi)^m = 1 - (1 - \pi)^m .$$

The probability that the 1<sup>st</sup> and the 2<sup>nd</sup> groups have at least one white ball is:

$$P(A_1 \cap A_2) = (1 - (1 - \pi)^m)(1 - (1 - \pi)^m) = (1 - (1 - \pi)^m)^2 .$$

The probability that all the groups have at least one white balls (i.e. that all the groups are positive) is:

$$P\left(\bigcap_{i=1}^k A_i\right) = (1 - (1 - \pi)^m)^k$$

## 2. Finite population size: less easy

$n_1$  white balls and  $n_2$  black balls ( $n_1 + n_2 = n$ ) are placed into  $k$  bins of maximum capacity  $m$ ;  $km = n$ . Let  $X$  be the random variable “number of bins without any white balls”. The random variable  $Y$  “number of white balls in a sample of  $m$  balls” has a **hypergeometric distribution** with parameters  $n$ ,  $n_1$  and  $m$ .

$$P(Y = w) = \frac{\binom{n_1}{w} \binom{km-n_1}{m-w}}{\binom{km}{m}}.$$

Let  $A_1$  be the event “the 1<sup>st</sup> sample has no white ball”. The probability that the 1<sup>st</sup> sample has no white ball is:

$$P(A_1) = \frac{\binom{km-n_1}{m}}{\binom{km}{m}} = \frac{\binom{km-m}{n_1}}{\binom{km}{n_1}}.$$

The probability that the 1<sup>st</sup> and the 2<sup>nd</sup> samples have no white ball is:

$$P(A_1 \cap A_2) = \frac{\binom{km-n_1}{m}}{\binom{km}{m}} \times \frac{\binom{km-n_1-m}{m}}{\binom{km-m}{m}} = \frac{\binom{km-2m}{n_1}}{\binom{km}{n_1}}.$$

The probability that the first  $s$  samples ( $s < k$ ) have no white balls is:

$$P(\cap_{i=1}^s A_i) = \frac{\binom{km-sm}{n_1}}{\binom{km}{n_1}}.$$

The probability that any  $s$  particular bins have no white balls is:

$$P\left(\bigcap_{\substack{i \in I \\ I \subset \{1,2,\dots,k\} \\ \text{Card}(I)=s}} A_i\right) = \binom{k}{s} \frac{\binom{km-sm}{n_1}}{\binom{km}{n_1}} = S_s$$

(there are  $\binom{k}{s}$  possible combinations for  $s$  (out of  $k$ ) bins without white balls).

Probability that at least one bin has no white ball:

$$\begin{aligned} P\left(\bigcup_{i=1}^k A_i\right) &= \sum_{i=1}^k (-1)^{i+1} S_i \quad (\text{principle of inclusion-exclusion for probability}) \\ &= \sum_{i=1}^k (-1)^{i+1} \binom{k}{i} \frac{\binom{km-im}{n_1}}{\binom{km}{n_1}} \\ &= \frac{1}{\binom{km}{n_1}} \sum_{i=1}^k (-1)^{i+1} \binom{k}{i} \binom{m(k-i)}{n_1}. \end{aligned}$$

The probability of having no bin without any white balls is:

$$P(X = 0) = 1 - \frac{1}{\binom{km}{n_1}} \sum_{i=1}^k (-1)^{i+1} \binom{k}{i} \binom{m(k-i)}{n_1} = \frac{1}{\binom{km}{n_1}} \sum_{i=0}^k (-1)^i \binom{k}{i} \binom{m(k-i)}{n_1}$$

And therefore the probability that all the groups are positive is:

$$1 - \frac{1}{\binom{km}{n_1}} \sum_{i=0}^k (-1)^i \binom{k}{i} \binom{m(k-i)}{n_1}$$

### 3. 2023 project around group testing



An Excel calculator has been developed with an implementation of these computations as well as the computation of the expected number of positive groups:

	A	B
1	<b>Group testing: on the number of positive groups</b>	
2	<b>Hypothesis: infinite population</b>	
3		
4	<b>Number of groups</b>	<b>10</b>
5	<b>Number of units per group</b>	<b>300</b>
6	<b>True characteristic content (%)</b>	<b>0.50%</b>
7		
8	<b>Probability that all groups are positive</b>	<b>8.09%</b>
9	<b>Expected number of positive groups</b>	<b>7.8</b>
10		
11	<b>Change any value in a yellow cell</b>	
12		
13		
14		
15		
16		

	A	B
1	<b>Group testing: on the number of positive groups</b>	
2	<b>Hypothesis: finite population (size = 3000 units)</b>	
3		
4	<b>Number of groups</b>	<b>10</b>
5	<b>Number of units per group</b>	<b>300</b>
6	<b>True characteristic (%)</b>	<b>0.50%</b>
7		
8	<b>Probability that all groups are positive</b>	<b>4.66%</b>
9	<b>Expected number of positive groups</b>	<b>7.9</b>
10		
11	<b>Change any value in a yellow cell</b>	
12		
13		
14		
15		
16		

Not yet published



## 4. Opportunities

## 4. Group testing estimator

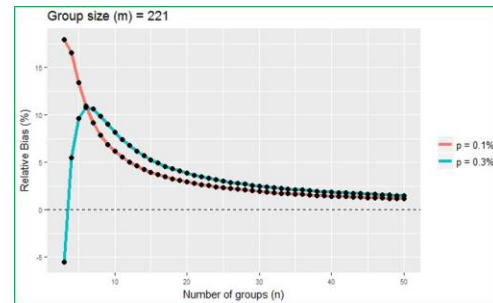
➔ Revisiting group testing estimator properties:  $\text{Var}[\hat{p}] = E[(\hat{p} - E[\hat{p}])^2]$  ,

$$\text{Bias}[\hat{p}] = E[\hat{p}] - p \text{ ,}$$

$$\text{MSE}[\hat{p}] = E[(\hat{p} - p)^2]$$

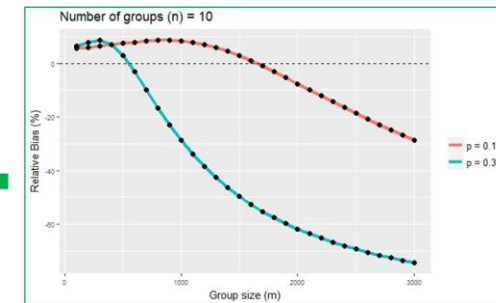
$$\text{Bias}[\hat{p}] = \left[ \sum_{x=0}^{n-1} \left( 1 - \left( 1 - \frac{x}{n} \right)^{1/m} \right) \binom{n}{x} \frac{(1 - (1-p)^m)^x (1-p)^{m(n-x)}}{1 - (1 - (1-p)^m)^n} \right] - p$$

➔ New insights for number of groups and group sizes recommendations



Recommendations on the number of groups to avoid **under-estimation**

Recommendations on group size to avoid **under-estimation**



## 4. Group testing estimator

➔ Estimator of the proportion of the number of white balls from the observed number of empty bins for the **hypergeometric group testing problem**

➔ Needs to solve the equation for  $n_1$  :

$$d = k \left( 1 - \frac{\binom{m(k-1)}{n_1}}{\binom{km}{n_1}} \right)$$

Difficult

where  $d$  is the observed number of positive bins.

**Table 2.1.** Minimum sampling intensity for seed lots in containers holding up to and including 100 kg seed

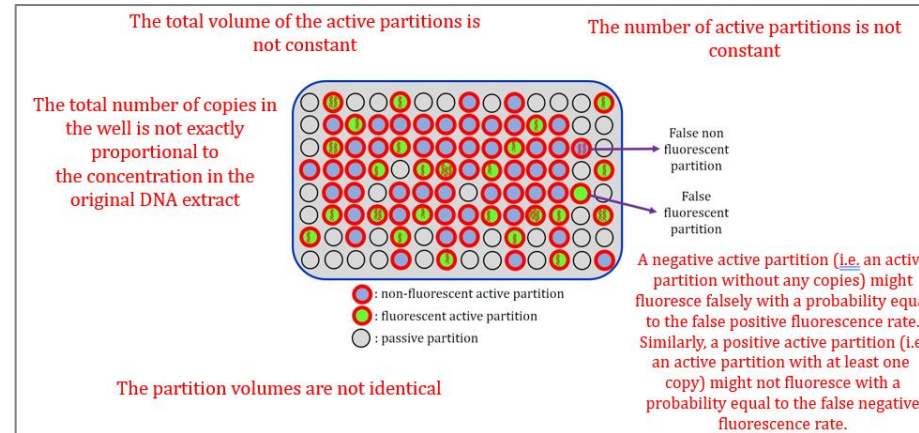
Number of containers	Minimum number of primary samples to be taken
1–4	3 primary samples from each container
5–8	2 primary samples from each container
9–15	1 primary sample from each container
16–30	15 primary samples, one each from 15 different containers
31–59	20 primary samples, one each from 20 different containers
60 or more	30 primary samples, one each from 30 different containers

These numbers have been elaborated over the years, using the results from sampling experiments and results from simulation studies.

➔ Can we fine-tune these numbers using sampling theory? (e.g. taking into account the size of the primary samples?)

➔ Use of theoretical results on two-stage sampling?

- dPCR modeling



- Method validation:
  - Revising *ISTAgermMV* R package
  - Reviewing needs in terms of number of labs, number of lots,...
- ...



## 4. Suggestions ?





Thank you!



Follow us on social media: