# *Repeatability* and *Reproducibility* of Categorical Measurements (Classification)

Yariv N. Marmor
Emil Bashkansky
Tamar Gadrich

BRAUDE College of Engineering, Karmiel

Leuven - ENBIS 2024
16/09/2024, 13:50

1

# The Purpose of the Presentation

To propose the model of estimating different components of classification precision, when it is provided by defined number of collaborators (classifiers) according to fixed and random model of their selection.

# Definition of *"categorical"* measurement

*"Classification of the analyzed property value of the objects under study (OUS) into one of K exclusive categories forming a comprehensive spectrum (scale) of the studied property will be considered as* <u>**categorical measurement**</u>*." (\*)*

Note 1: The results of classification are presented by so-called categorical data. In cases where the spectrum of possible values consists only of two categories such data are binary, and the appropriate activity is also often called *testing.*

Note 2: In this presentation categories are not ordered (*nominal* scale)

\* T. Gadrich, E. Bashkansky, (2016) " A Bayesian approach to evaluating uncertainty of inaccurate categorical measurements", *Measurement* 91, 186–193.

# Examples of *Classifiers* ( $K > 2$ )
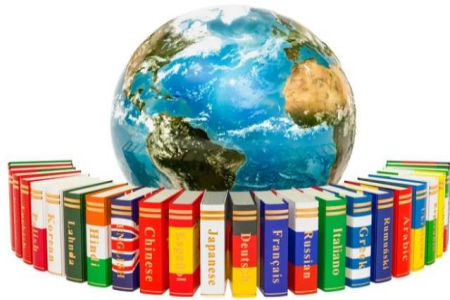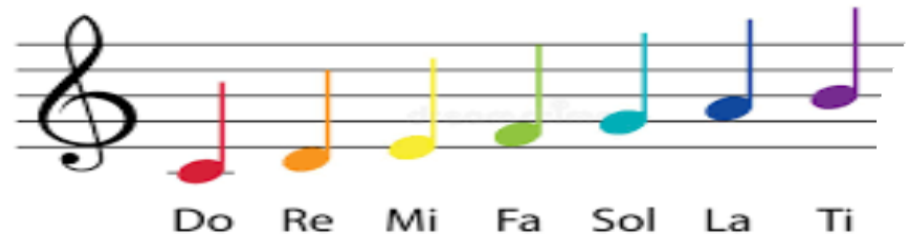


Coins sorting machine



Egg classification machine



Plastic color sorting machine



Google language detector



Musical tone recognition

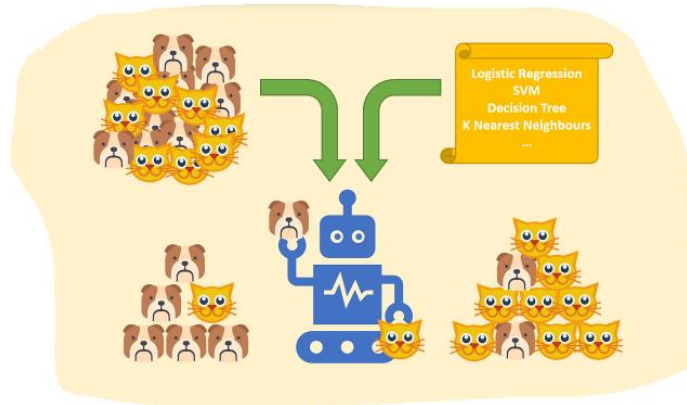# Examples of *Binary Classifiers* ($K = 2$)

**Pregnancy tester**

**Spam filtering**
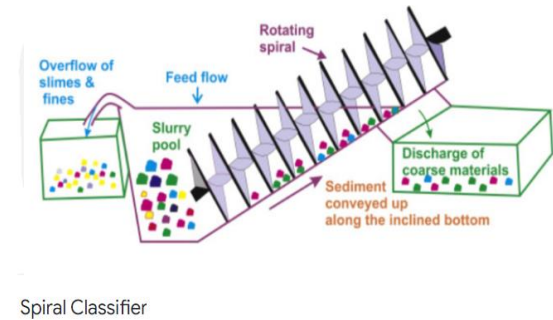
**Counting machine: bank notes are classified as forged or accepted**

**Covid - 19 tester**

**Classification algorithm**

**Spiral classifier**

5

# *Ability* of *Classifier* (General Case)

The conditional probabilities that an object will be

classified as category $k$, given that its actual/true category

is $i$ - $P_{k|i}$

$$1 \leq i, k \leq K; \qquad \sum_{k=1}^{K} P_{k|i} = 1$$

**Ideal classifier:** $every\ P_{k|i} = 0, except\ of\ P_{i|i}$

# *Classification (Confusion) Matrix and Repeatability for the General Case of K Categories*

$$P = \begin{pmatrix} p_{1|1} & p_{2|1} & \cdots & p_{k|1} & \cdots & p_{K|1} \\ p_{1|2} & p_{2|2} & \cdots & p_{k|2} & \cdots & p_{K|2} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ p_{1|i} & p_{2|i} & \cdots & p_{k|i} & \cdots & p_{K|i} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ p_{1|K} & p_{2|K} & \cdots & p_{k|K} & \cdots & p_{K|K} \end{pmatrix} \quad I = \begin{pmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

$$\textbf{\textit{Repeatability}} = \begin{bmatrix} Repeatability_1 \\ Repeatability_2 \\ \cdots \\ Repeatability_i \\ \cdots \\ Repeatability_K \end{bmatrix} \text{ - (for the \textbf{same} classifier)}$$

$$Repeatability_i = \frac{K}{K-1} \sum_{k=1}^{K} \left[ p_{k|i} \cdot (1 - p_{k|i}) \right] = \frac{K}{K-1} \left( 1 - \sum_{k=1}^{K} [p_{k|i}]^2 \right) = VAR_{within}$$

$$0 \leq VAR_{within} \leq 1$$

# The *Closeness of Agreement* between classifications obtained by <u>different</u> *Classifiers* participating in collaborative study (from here on - *Classifiers effect* )

$$Classifiers\ effect = \begin{bmatrix} Classifier\ effect_1 \\ Classifier\ effect_2 \\ \dots \\ Classifier\ effect_i \\ \dots \\ Classifier\ effect_K \end{bmatrix}$$

Where:

$$Classifiers\ effect_i = \frac{K}{K-1} \sum_{k=1}^{K} VAR(p_{k|i})$$

$VAR\ (p_{k|i})$ = classic variation of $p_{k|i}$ between collaborators/classifiers

According to **ISO 5725:2023 "Accuracy (trueness and precision of measurement methods and results)"** (differ from that in the Gauge R&R studies): $Reproducibility = Repeatability + Classifier\ effect$
Reproducibility describes the **total precision** of classification procedure (i.e., categorical measurements)

# $H_0$ - *Homogeneity Hypothesis : Acceptance and Rejecting ( for the fixed group of classifiers)*

$H_0$ : **all classifiers' counts are drawn from the same population characterized by vector ($p_1, p_2,\ldots, p_k ,\ldots p_K$).**

**If** *SP* **(segregation/separation power) statistics is less than some critical value,**

$H_0$ **is not rejected. Otherwise, see (*) for unbiased estimators of variation components, i.e.:** *repeatability, classifier's effect* **and** *total precision.*

*T. Gadrich, E. Bashkansky, I. Kuselman (2013) "Comparison of biased and unbiased estimators of variances of qualitative and semi-quantitative results of testing", *Accreditation and Quality Assurance*, **18** (2), 85-90

*T. Gadrich, E. Bashkansky, R. Zitikis, (2015) "Assessing variation: a unifying approach for all scales of measurement", *Quality and Quantity*, **49** (3) , 1145-1167
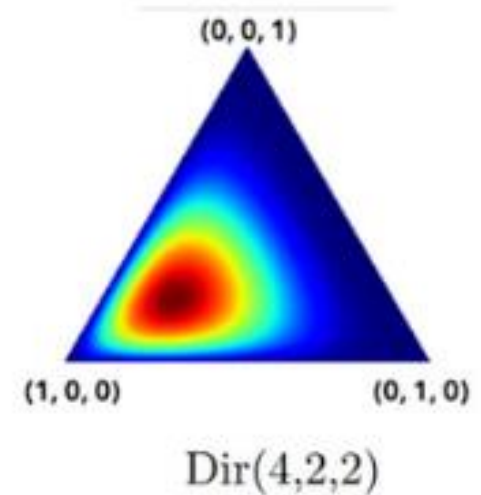
# Random *Dirichlet - Multinomial (DM)* Model - 1

❑ $L$ classifiers are randomly sampled from the population which classification abilities related to category $i$ are distributed according to the Dirichlet distribution:

$$f\left(p_{1|i}, p_{2|i}, \cdots, p_{K|i}\right) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^{K} p_{k|i}^{\alpha_{k|i}-1} \qquad \left(\sum_{k=1}^{K} p_{k|i} = 1\right)$$

where: $\boldsymbol{\alpha_i} = \left(\alpha_{1|i}, \alpha_{2|i}, \cdots, \alpha_{K|i}\right)$ - *parameters of the Dirichlet distribution characterizing the i-th category classification*;

$E\left(p_{k|i}\right) = \alpha_{k|i} / \alpha_{0|i}$ ;

$\alpha_{0|i} = \sum_{k=1}^{K} \alpha_{k|i}$ - *concentration parameter* (*)



(0, 0, 1)

(1, 0, 0)          (0, 1, 0)

Dir(4,2,2)

* Location is determined by repeatability, and dispersion ($\alpha_{0|i}$) is determined by degree of variation between classifiers.

10

# *Repeatability (within)* and *Classifiers (between-classifiers) effect*

$$Repeatability_i = \frac{K}{K-1}\sum_{k=1}^{K} E\left[p_{k|i} \cdot \left(1 - p_{k|i}\right)\right] = \frac{K}{K-1}\frac{\alpha_{0|i}}{\alpha_{0|i}+1}\left[1 - \sum_{k=1}^{K}\frac{\alpha_{k|i}^2}{\alpha_{0|i}^2}\right]$$

$$Classifier\ effect_i = \frac{K}{K-1}\sum_{k=1}^{K} VAR\left(p_{k|i}\right) = \frac{K}{K-1}\frac{1}{\alpha_{0|i}+1}\left[1 - \sum_{k=1}^{K}\frac{\alpha_{k|i}^2}{\alpha_{0|i}^2}\right] = \frac{Repeatability}{\alpha_{0|i}}$$

$$Reproducibility_i = \frac{K}{K-1}\left[1 - \sum_{k=1}^{K}\frac{\alpha_{k|i}^2}{\alpha_{0|i}^2}\right]$$

| $\alpha_{0|i} \to 0$ | $\alpha_{0|i} \to \infty$ |
|---|---|
| Repeatability = 0 | Classifiers effect = 0 |
| $Classifiers\ effect_i$ $= \frac{K}{K-1}\left[1 - \sum_{k=1}^{K}\frac{\alpha_{k|i}^2}{\alpha_{0|i}^2}\right]$ | $Repeatability_i = \frac{K}{K-1}\left[1 - \sum_{k=1}^{K}\frac{\alpha_{k|i}^2}{\alpha_{0|i}^2}\right]$ |

11

# Random *Dirichlet - Multinomial* (*DM*) Model - 2

❑   From an experiment consisting of **N repeated classifications** of each property that actually belongs to the category $i$

$$\left(n_{1|i}^{(l)}, n_{2|i}^{(l)}, \ldots, n_{K|i}^{(l)} \middle| \boldsymbol{p} = (p_{1|i}, p_{2|i}, , \ldots, p_{K|i})\right) \sim Multinomial(N, \boldsymbol{p} = (p_{1|i}, p_{2|i}, , \ldots, p_{K|i}))$$

where $\boldsymbol{p} = (p_{1|i}, p_{2|i}, , \ldots, p_{K|i}) \sim Dirichlet(\boldsymbol{\alpha_i})$

Estimator of classification matrix of the $l$-th classifier $(l = 1, 2, \ldots, L)$.

$$\hat{P}^{(l)} = \begin{pmatrix} \hat{p}_{1|1}^{(l)} & \hat{p}_{2|1}^{(l)} & \cdots & \hat{p}_{k|1}^{(l)} & \cdots & \hat{p}_{K|1}^{(l)} \\ \hat{p}_{1|2}^{(l)} & \hat{p}_{2|2}^{(l)} & \cdots & \hat{p}_{k|2}^{(l)} & \cdots & \hat{p}_{K|2}^{(l)} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \hat{p}_{1|i}^{(l)} & \hat{p}_{2|i}^{(l)} & \cdots & \hat{p}_{k|i}^{(l)} & \cdots & \hat{p}_{K|i}^{(l)} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \hat{p}_{1|K}^{(l)} & \hat{p}_{2|K}^{(l)} & \cdots & \hat{p}_{k|K}^{(l)} & \cdots & \hat{p}_{K|K}^{(l)} \end{pmatrix} \quad \text{where } \hat{p}_{k|i}^{(l)} = \frac{n_{k|i}^{(l)}}{N}$$

# *Unbiased Estimators*

**Unbiased estimator for <u>repeatability</u> (within) variation:**

$$\widehat{Repeatability}_{u(i)} = \frac{N}{N-1}\frac{1}{L}\sum_{l=1}^{L}\widehat{Repeatability}_i =$$

$$= \frac{N}{N-1}\frac{1}{L}\sum_{l=1}^{L}\frac{K}{(K-1)}\sum_{k=1}^{K}\hat{p}_{k|i}^{(l)}\cdot\left(1-\hat{p}_{k|i}^{(l)}\right)$$

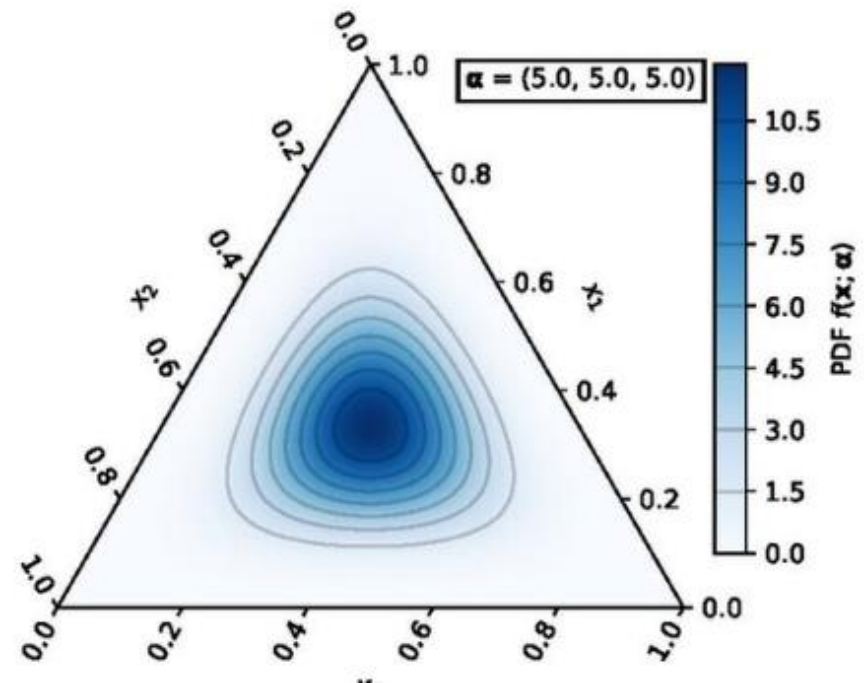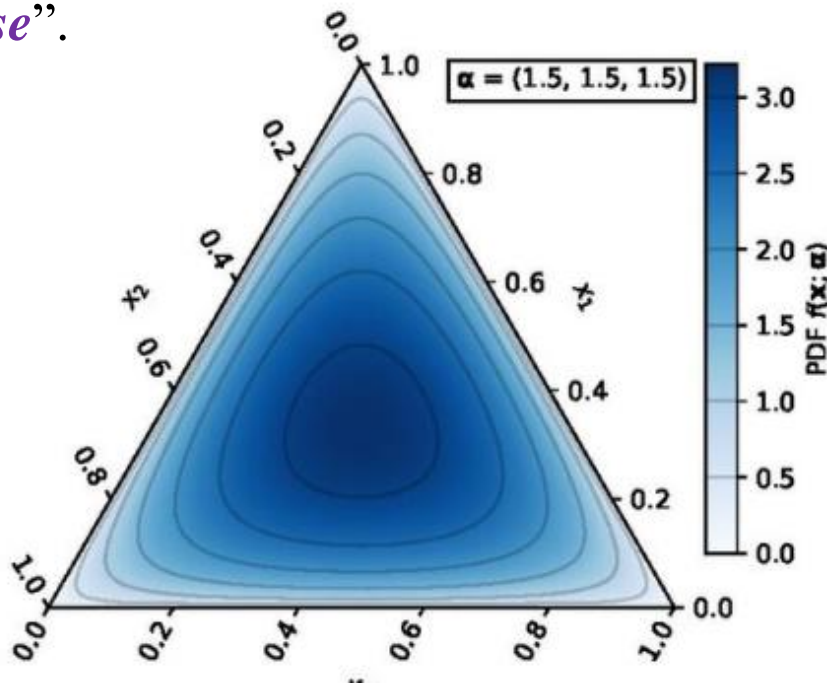**Unbiased estimator for <u>between classifiers</u> variation:**

$$\widehat{Classifier\ effect}_{u(i)} = \frac{K}{K-1}\sum_{k=1}^{K}\frac{1}{L-1}\sum_{l=1}^{L}\left(\hat{p}_{k|i}^{(l)}-\frac{1}{L}\sum_{l=1}^{L}\hat{p}_{k|i}^{(l)}\right)^2$$

$$-\frac{\widehat{Repeatability}_{u(i)}}{N}$$

**Unbiased estimator for <u>total</u> classification variation:**

$$\widehat{VAR_{TOTAL}} = \frac{K}{K-1}\sum_{k=1}^{K}\frac{1}{L-1}\sum_{l=1}^{L}\left(\hat{p}_{k|i}^{(l)}-\frac{1}{L}\sum_{l=1}^{L}\hat{p}_{k|i}^{(l)}\right)^2 +$$

$$+\frac{1}{L}\sum_{l=1}^{L}\frac{K}{K-1}\sum_{k=1}^{K}\hat{p}_{k|i}^{(l)}\cdot\left(1-\hat{p}_{k|i}^{(l)}\right)$$

# How to estimate vector $\boldsymbol{\alpha_i} = \left(\alpha_{1|i}, \alpha_{2|i}, \cdots, \alpha_{K|i}\right)$, characterizing *Population Distribution* of classification abilities?

It is considered that <u>the key parameter is</u> $\alpha_{0|i}$ , characterizing the so called "***Dirichlet noise***".



We propose to estimate $\alpha_{0|i}$ by help of Repeatability to *Classifiers effect* ratio:

$$\hat{\alpha}_{0|i} = \frac{\text{Repeatability}_{u(i)}}{Classifier\ effect_{u(i)}} \qquad \text{and then} \qquad \hat{\alpha}_{k|i} = \hat{\alpha}_{0|i} \frac{1}{L} \sum_{L=1}^{L} \hat{p}_{k|i}^{(l)}$$

# Simulation Study for $K = 3$

The population was $\boldsymbol{p} = (0.8200, 0.1667, 0.0133)$ (*)

Accordingly, we simulated for different $\boldsymbol{\alpha}$ and $\alpha_0$:

- $\boldsymbol{\alpha} = (2.46, 0.5, 0.04)$, i.e., $\alpha_0 = 3$

- $\boldsymbol{\alpha} = (8.2, 1.667, 0.133)$, i.e., $\alpha_0 = 10$

- $\boldsymbol{\alpha} = (16.4, 3.333, 0.267)$, i.e., $\alpha_0 = 20$

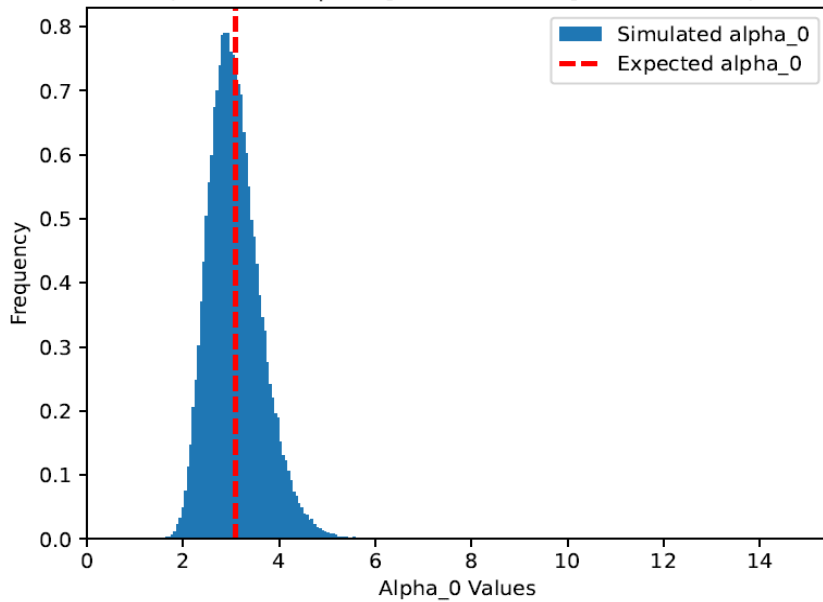- $\boldsymbol{\alpha} = (24.6, 5.0, 0.4)$, i.e., $\alpha_0 = 30$

While

- $L$ = 5, 10, 20, 30 or 100

- $N$ = 1000

* E. Bashkansky, S. Dror, R. Ravid, P. Grabov (2007) "Effectiveness of a Product Quality Classifier", *Quality Engineering*, vol. 19, issue 3, pp.235-244
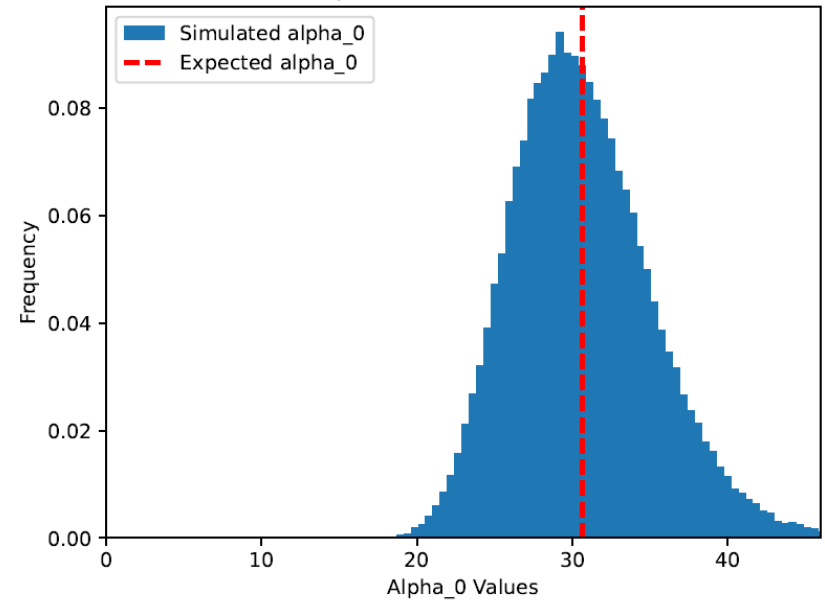
# Influence of $\alpha_0$ magnitude

$$\alpha_0 = 3$$

$$\alpha_0 = 30$$



Distribution of Alpha_0
(N=1000; alpha=[2.46 0.5 0.04]; K=3; L=100)

Distribution of Alpha_0
(N=1000; alpha=[24.6 5. 0.4]; K=3; L=100)

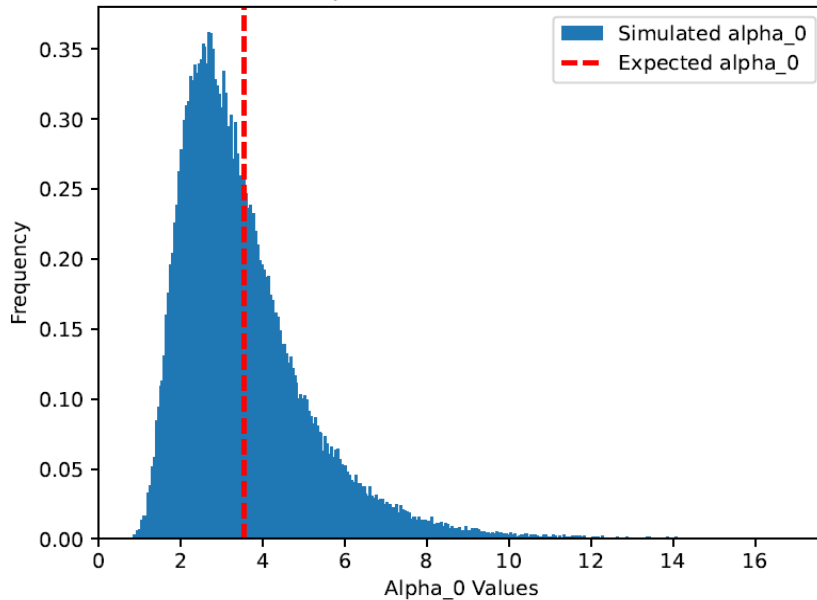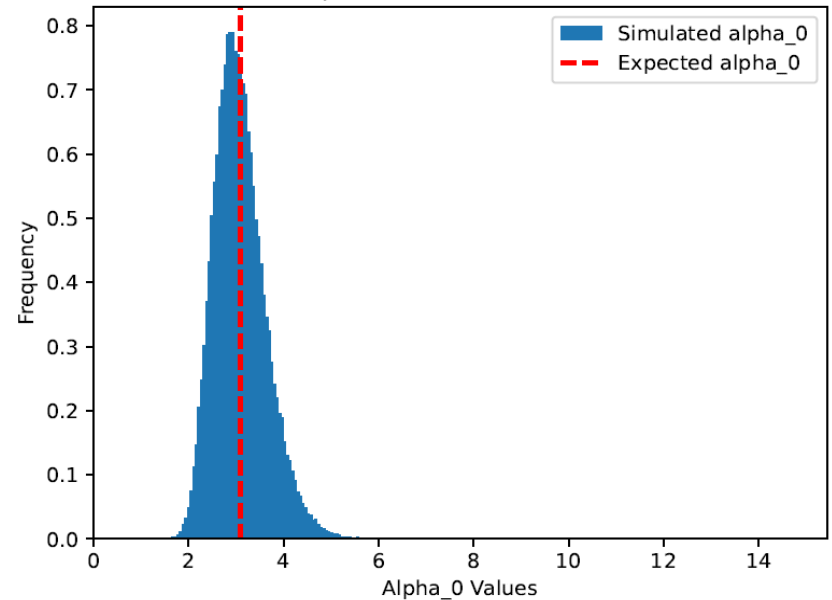# Influence of the number of *classifiers - L*

$L = 20$

$L = 100$



$\alpha_0 = 3$

# How the number of *classifiers* $L$ influences the *relative deviation* of $\alpha_0$ estimator



$L$ = 30 brings quite good results!

# *Summary*

1. Classification precision is particularly crucial in scenarios where the cost of false output is high, e.g. medical diagnosis, search engine results, product quality control etc.

2. A statistical model for analyzing classification's precision and its *intra* and *inter* components from collaborative studies is presented.

3.  Unbiased estimators of repeatability, classifiers' effect and total precision are  provided for both the fixed and the random model.

4. Using simulation, it was found that the greatest influence on the accuracy and effectiveness of the model estimations in a random model is exerted mainly by the number of classifiers.

# Thank You for Your Attention!



**E-mail:** *myariv@braude.ac.il*