

# Let's Agree to Disagree: Fixing Agreement Measures for Crowdsourcing

**Kevin Roitero** (and Stefano Mizzaro)



Funded by the Horizon 2020 Framework Programme  
of the European Union

This publication is based upon work from COST Action  
CA16105, supported by COST (European Cooperation in  
Science and Technology).

Lisbon, 15th March 2019

# but the real title should have been...

**Kevin Roitero** (and Stefano Mizzaro)



Funded by the Horizon 2020 Framework Programme  
of the European Union

This publication is based upon work from COST Action  
CA16105, supported by COST (European Cooperation in  
Science and Technology).

Lisbon, 15th March 2019

# The Elephant in the Room

**Kevin Roitero** (and Stefano Mizzaro)



Funded by the Horizon 2020 Framework Programme  
of the European Union

This publication is based upon work from COST Action  
CA16105, supported by COST (European Cooperation in  
Science and Technology).

Lisbon, 15th March 2019

# This Talk is Based on the Following Paper

## **Let's Agree to Disagree: Fixing Agreement Measures for Crowdsourcing.**

Alessandro Checco, Kevin Roitero, Eddy Maddalena, Gianluca Demartini and Stefano Mizzaro.  
Proceedings of the The fifth AAI Conference on Human Computation and Crowdsourcing, AAI HCOMP 2017.  
Quebec City, Canada. October 24-26 2017.

url: <https://aaai.org/ocs/index.php/HCOMP/HCOMP17/paper/viewFile/15927/15258>

# Setting

- micro-task crowdsourcing
- many workers do the same task
- agreement among workers can / should be leveraged
- leveraging agreement can be useful for:
  - estimating the reliability of collected data
  - understanding behavior of the workers

# Agreement Formalization

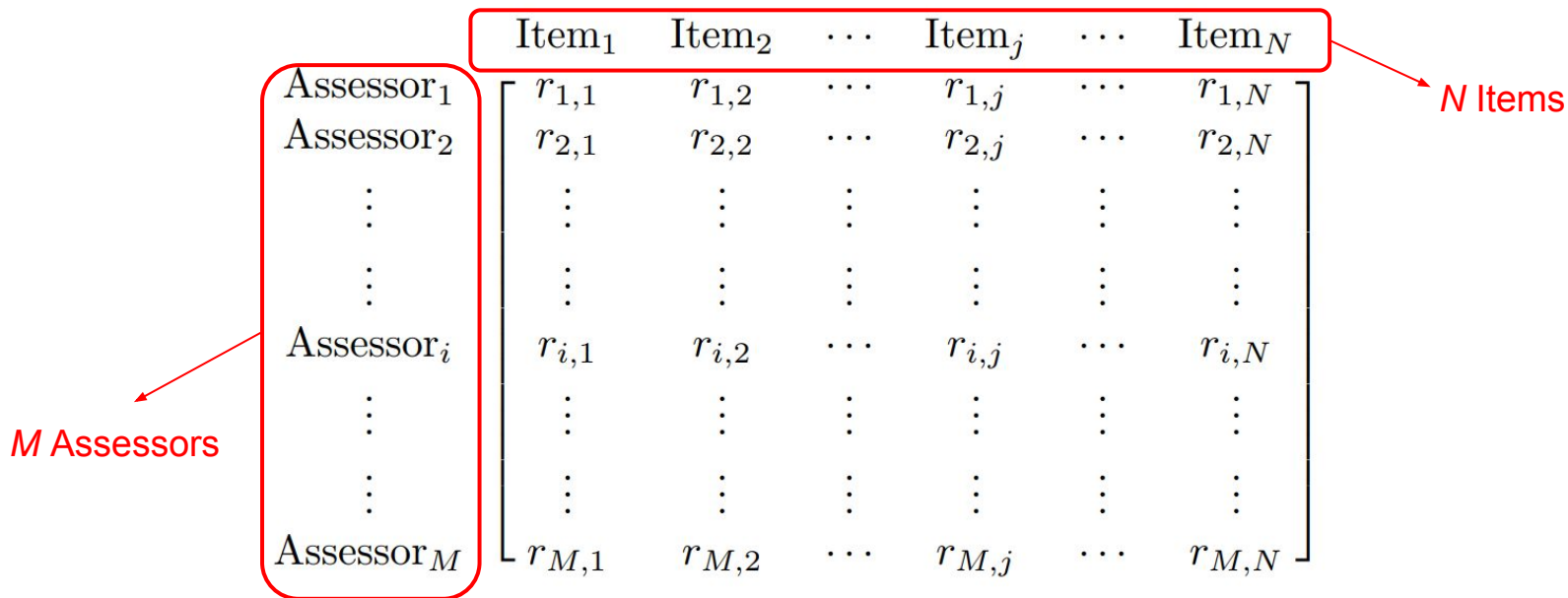
	Item <sub>1</sub>	Item <sub>2</sub>	...	Item <sub>j</sub>	...	Item <sub>N</sub>
Assessor <sub>1</sub>	$r_{1,1}$	$r_{1,2}$	...	$r_{1,j}$	...	$r_{1,N}$
Assessor <sub>2</sub>	$r_{2,1}$	$r_{2,2}$	...	$r_{2,j}$	...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>i</sub>	$r_{i,1}$	$r_{i,2}$	...	$r_{i,j}$	...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>M</sub>	$r_{M,1}$	$r_{M,2}$	...	$r_{M,j}$	...	$r_{M,N}$

# Agreement Formalization

	Item <sub>1</sub>	Item <sub>2</sub>	...	Item <sub>j</sub>	...	Item <sub>N</sub>
Assessor <sub>1</sub>	$r_{1,1}$	$r_{1,2}$	...	$r_{1,j}$	...	$r_{1,N}$
Assessor <sub>2</sub>	$r_{2,1}$	$r_{2,2}$	...	$r_{2,j}$	...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>i</sub>	$r_{i,1}$	$r_{i,2}$	...	$r_{i,j}$	...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>M</sub>	$r_{M,1}$	$r_{M,2}$	...	$r_{M,j}$	...	$r_{M,N}$

*N* Items

# Agreement Formalization



The diagram shows a matrix representing agreement between  $M$  assessors and  $N$  items. The rows are labeled with assessor names and the columns with item names. A red box highlights the row labels, with an arrow pointing to the text " $M$  Assessors". Another red box highlights the column labels, with an arrow pointing to the text " $N$  Items".

	Item <sub>1</sub>	Item <sub>2</sub>	...	Item <sub><math>j</math></sub>	...	Item <sub><math>N</math></sub>
Assessor <sub>1</sub>	$r_{1,1}$	$r_{1,2}$	...	$r_{1,j}$	...	$r_{1,N}$
Assessor <sub>2</sub>	$r_{2,1}$	$r_{2,2}$	...	$r_{2,j}$	...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub><math>i</math></sub>	$r_{i,1}$	$r_{i,2}$	...	$r_{i,j}$	...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub><math>M</math></sub>	$r_{M,1}$	$r_{M,2}$	...	$r_{M,j}$	...	$r_{M,N}$



# Agreement Formalization

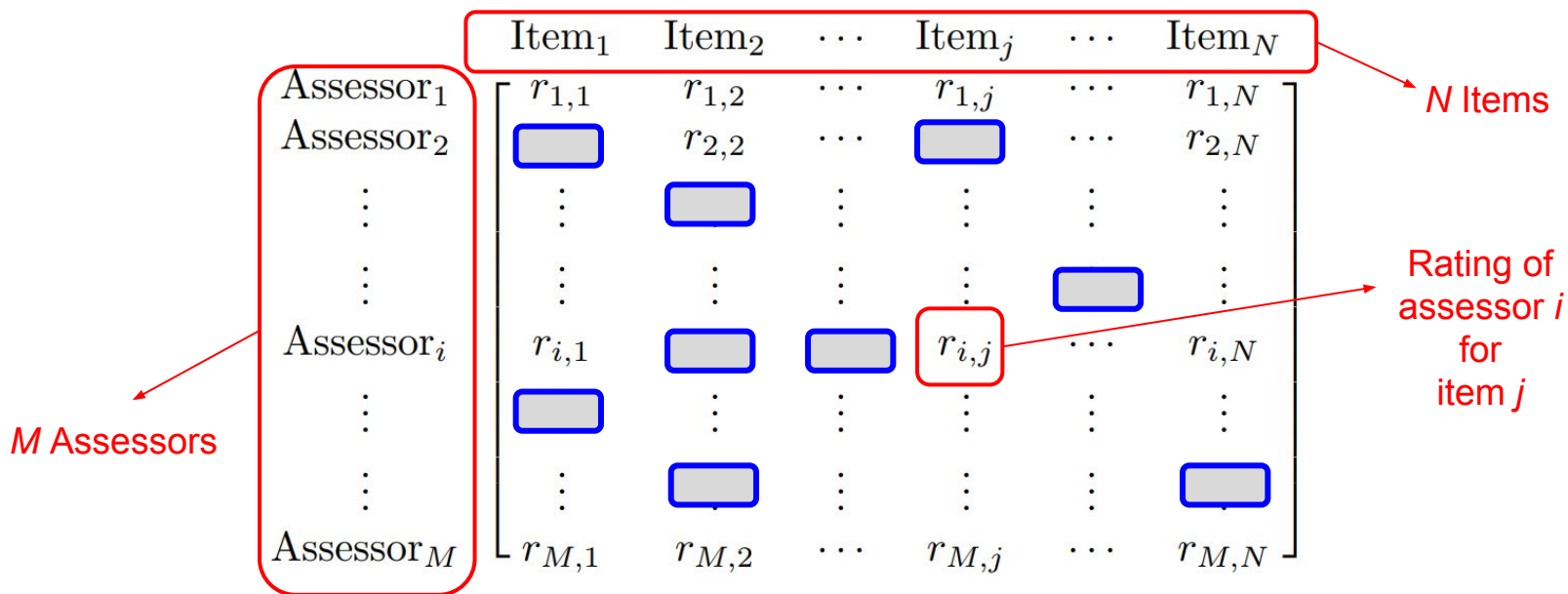
	Item <sub>1</sub>	Item <sub>2</sub>	...	Item <sub>j</sub>	...	Item <sub>N</sub>
Assessor <sub>1</sub>	$r_{1,1}$	$r_{1,2}$	...	$r_{1,j}$	...	$r_{1,N}$
Assessor <sub>2</sub>	$r_{2,1}$	$r_{2,2}$	...	$r_{2,j}$	...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>i</sub>	$r_{i,1}$	$r_{i,2}$	...	$r_{i,j}$	...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor <sub>M</sub>	$r_{M,1}$	$r_{M,2}$	...	$r_{M,j}$	...	$r_{M,N}$

$M$  Assessors

$N$  Items

Rating of assessor  $i$  for item  $j$

# Agreement Formalization



This matrix is often **very** sparse in crowdsourcing

# There are Several Agreement Measures

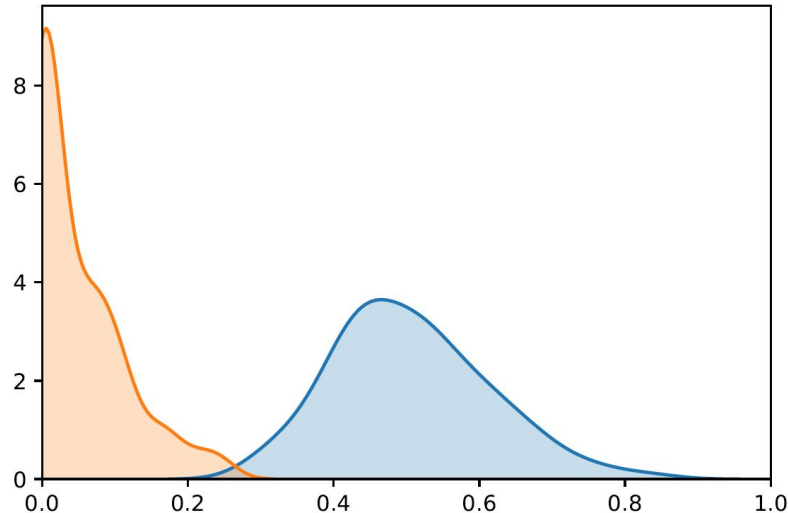
- Percentage Agreement (PA)
- Scott's  $\pi$
- Cohen's  $\kappa$
- Intraclass Correlation Coefficient (ICC)
- Fleiss  $\kappa$
- Krippendorff's Alpha

# Current Agreement Measures Are Inadequate

- measures often borrowed from other scenarios with **different assumptions** (which usually do not hold for crowdsourcing):
  - one assessor rates all items
  - all assessors rate all items
  - limited and fixed (= known) number of assessors
- measures are often designed for estimating **data reliability**, not **agreement**
  - **reliability**: the capacity of any measurement tool to differentiate between respondents when measured twice under the same conditions. [Berchtold]
  - **agreement**: the capacity of any other measurement tool applied twice on the same respondents under the same conditions to provide strictly identical results. [Berchtold]
  - reliability can be considered as a necessary but not sufficient condition to demonstrate agreement. [Berchtold]

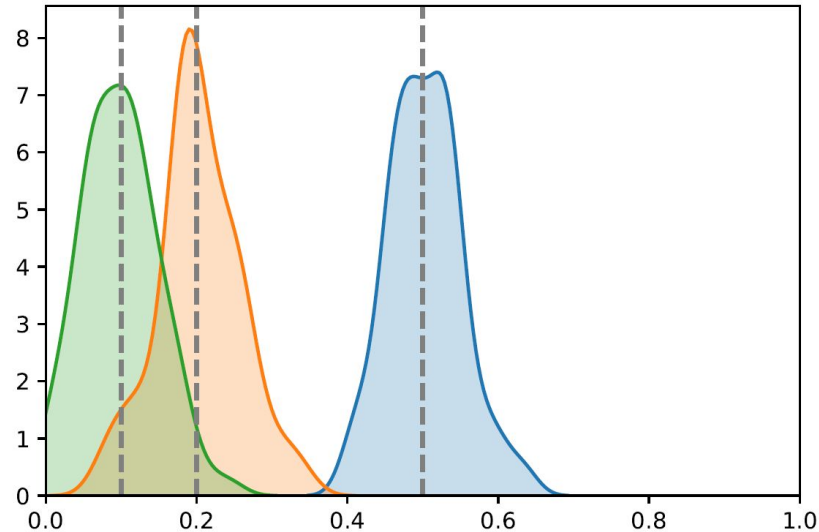
# Problems

- there is more variability of judgments in the centre of the scale w.r.t. scale boundaries.  
→ can lead to over-estimate agreement close to scale boundaries.



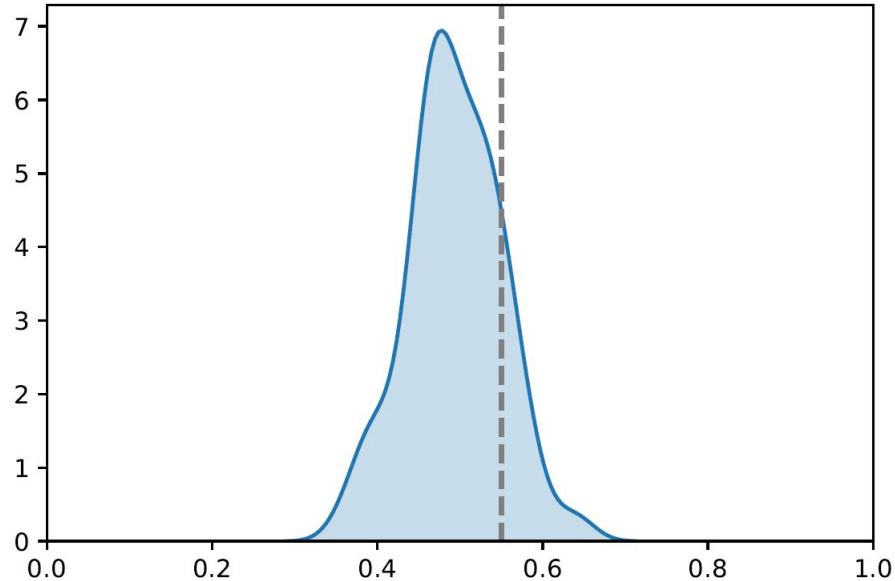
# Problems

- the concentration point can be different for different items  
→ can lead to over/under-estimate agreement



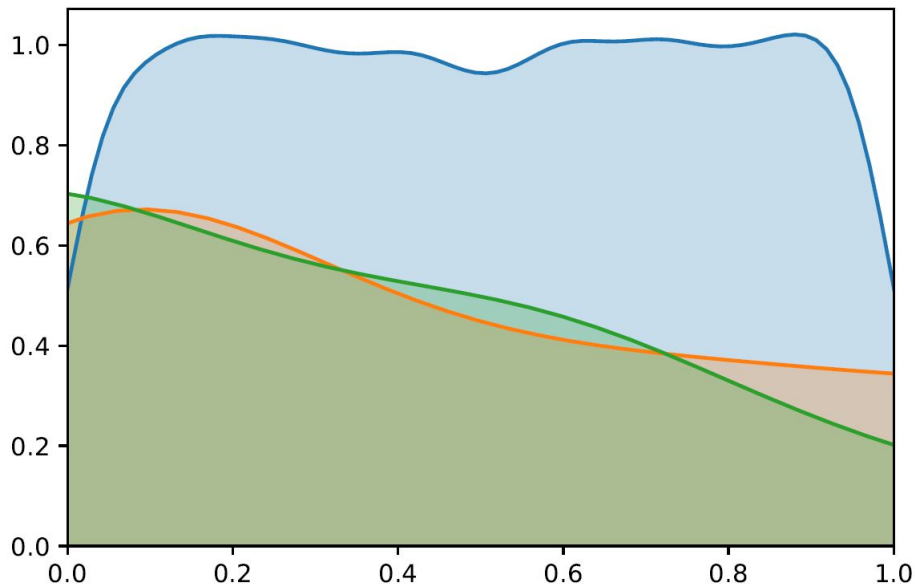
# Problems

- additional information is often not considered (e.g., gold questions)



# Problems

- different ideas of “**agreement by chance**” definition
- correction by chance assumptions are often violated in crowdsourcing setting





# Real Problems with State-of-the-Art Measures

- Percentage Agreement (PA)
  - does not consider agreement by chance
  - works only with nominal data
  - depends on the scale granularity (can not compare different scales)
- Scott's  $\pi$  and Cohen's  $\kappa$ 
  - work only with two assessors
  - work only with nominal data
- Intraclass Correlation Coefficient (ICC)
  - assessor have same marginal probability of an answer (not true in crowdsourcing)
  - equivalent to weighted Cohen's  $\kappa$
- Fleiss  $\kappa$ 
  - Generalizes  $\kappa$  to multiple assessors (i.e., shares the same issues)

# Real Problems with State-of-the-Art Measures

- Krippendorff's Alpha: an attempt to generalize previous metrics
  - **Random guessing can have high agreement**
  - **Random guessing may have more agreement than honest coding**
  - High agreement, low reliability
  - Zero change in percentage agreement causing radical drop in reliability.
  - **Eliminating disagreements does not improve agreement**
  - Honest work as bad as coin flipping.
  - Two datasets: same quality, same agreement; but higher reliability in one.
  - punishing larger sample and replicability (i.e., data quantity dependent)
  - **“reverse answer” problem**  $([1, 0, 0, 0, 1] \neq [0, 1, 1, 1, 0])$

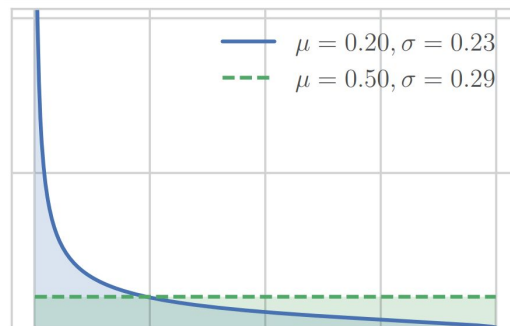
(a complete overview and all the mathematical details are available in our paper)

# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:

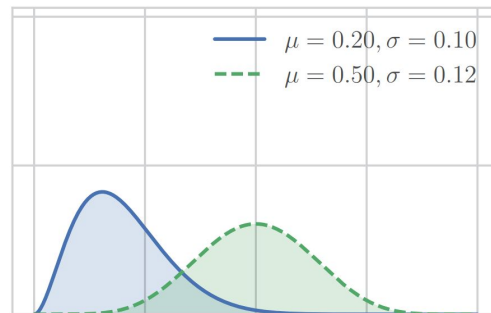
# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
  - random judgments  $\rightarrow$  flat distribution



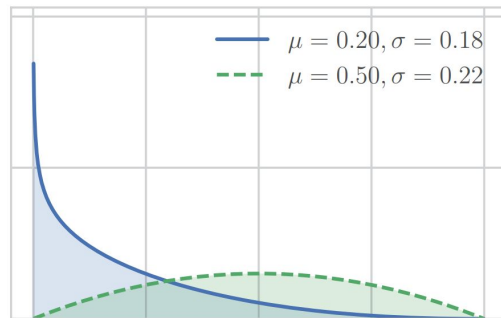
# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
  - random judgments → flat distribution
  - agreement → bell-shaped distribution



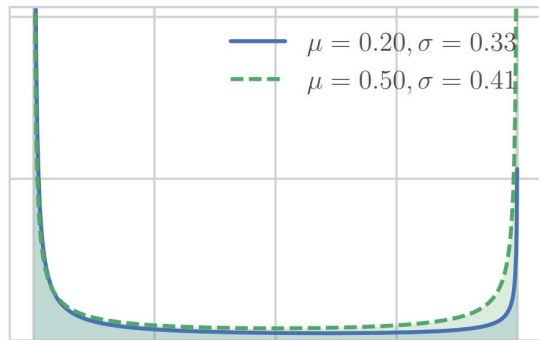
# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
  - random judgments → flat distribution
  - agreement → bell-shaped distribution
  - agreement around scale boundaries → J-distribution



# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
  - random judgments → flat distribution
  - agreement → bell-shaped distribution
  - agreement around scale boundaries → J-distribution
  - disagreement → U shaped distribution



# Our Measure: $\Phi$

- *agreement* definition as a key point:
  - “agreement is the amount of concentration around a data value”
- if we do not observe agreement (i.e., concentration around a point), we have disagreement, treated as negative agreement in our measure
- in practice:
  - first, we fit a distribution over the histogram of the ratings
  - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
  - random judgments → flat distribution
  - agreement → bell-shaped distribution
  - agreement around scale boundaries → J-distribution
  - disagreement → U shaped distribution
- we should have a minimal number of parameters, to avoid overfitting



## Our Measure: $\Phi$

- we use a Beta distribution to model our scenario:  $B(a, b)$
- we re-parametrize the distribution in terms of the mean value  $\mu$  and the precision  $p$  as  $\mu = \frac{a}{a+b}$ ;  $p = a + b$
- now, we can treat separately mean and dispersion
- we can have a metric that is agnostic of the mean value
- then, we transform to have values in the  $[-1, +1]$  range:


$$\Phi = 1 - 2^{\frac{-p \log 2}{2}}$$

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :


$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}}$$
$$\prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a common dispersion, given the observed data



# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :


$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}}$$
$$\prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

**probability** of observing the mean values, with a common dispersion, given the observed data

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}}$$
$$\prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the **mean values**, with a common dispersion, given the observed data

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}}$$
$$\prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a **common dispersion**, given the observed data

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}} \prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$


probability of observing the mean values, with a common dispersion, given the **observed data**

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}} \prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a common dispersion, given the observed data



- Then, we estimate  $\Phi$  using


$$\hat{\Phi} = \arg \max_{\Phi} P(\vec{\mu}, \Phi | X).$$

# Our Measure: $\Phi$

- we use Bayesian inference to compute  $\Phi$ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}} \prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a common dispersion, given the observed data



- Then, we estimate  $\Phi$  using

$$\hat{\Phi} = \arg \max_{\Phi} P(\vec{\mu}, \Phi | X).$$

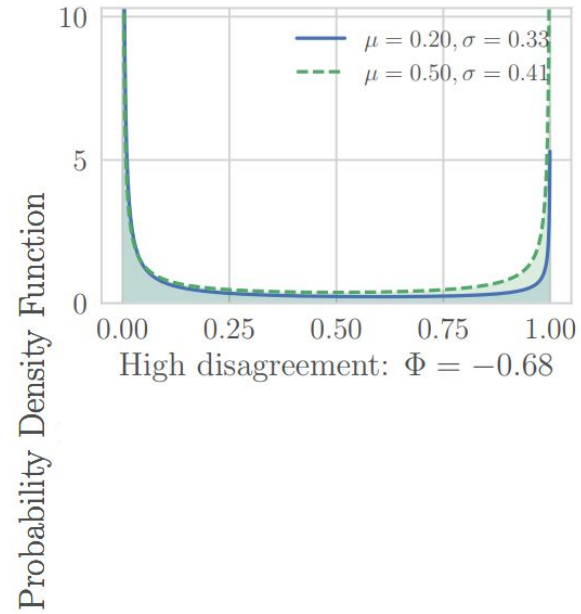
the formula can change to incorporate custom ground truth



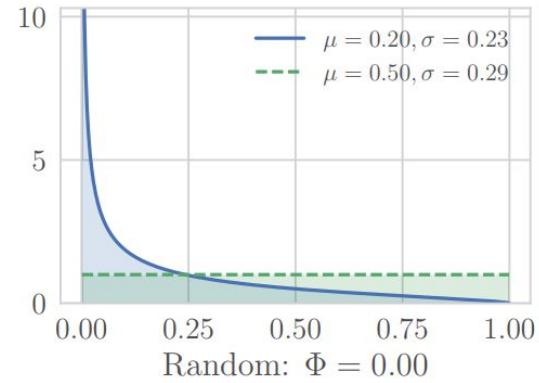
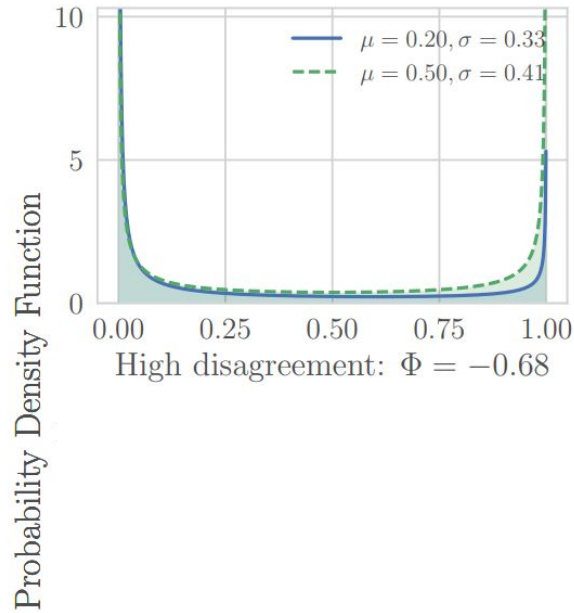
## $\Phi$ Interpretation

- **High Disagreement.** When  $\Phi < 0$ , there is no central tendency value but rather a tendency to exclude a central area (polarized behavior)
- **Random.** When  $\Phi=0$ , the behavior is equivalent with a unbounded uniform process censored on the scale
- **Weak Agreement.** When  $0 < \Phi \leq 0.5$ , the distribution has no inflection point, but there is a unique central tendency or a dispersion that is smaller than a uniform process
- **High Agreement.** When  $\Phi > 0.5$ , the distribution is bell shaped with two inflection points, more narrow around the mean as  $\Phi$  grows

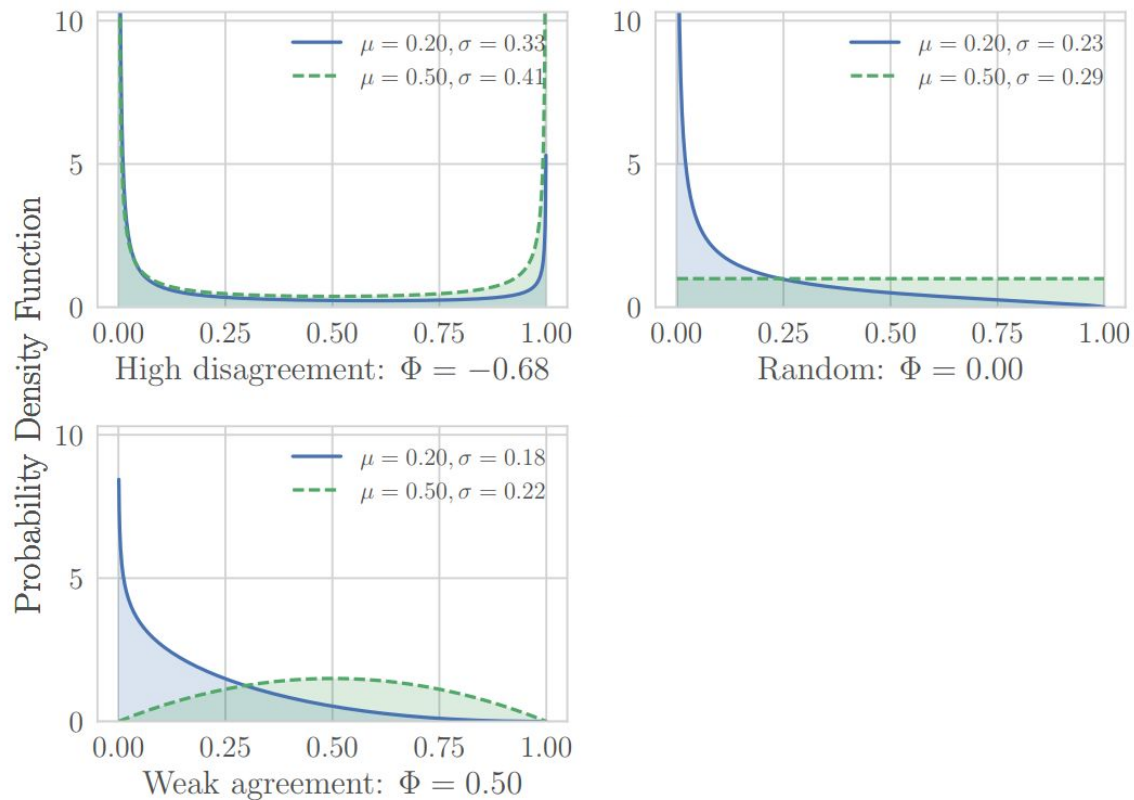
# Examples of $\Phi$ Shapes



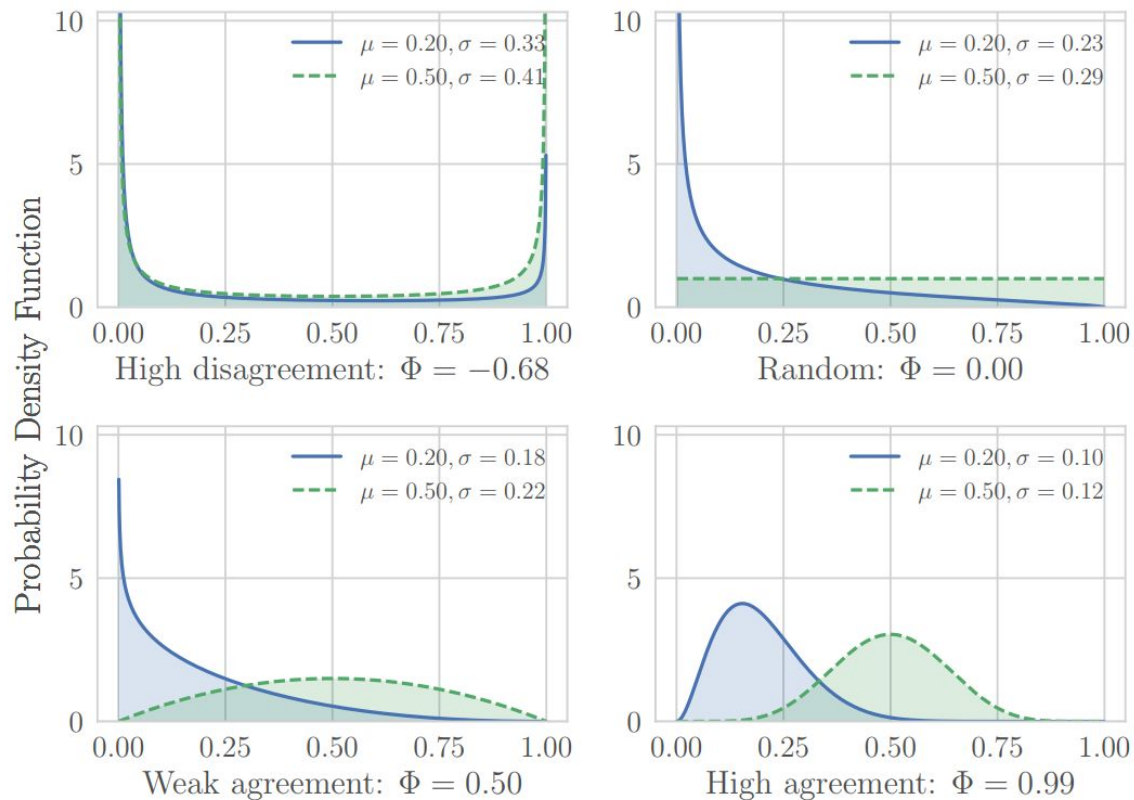
# Examples of $\Phi$ Shapes



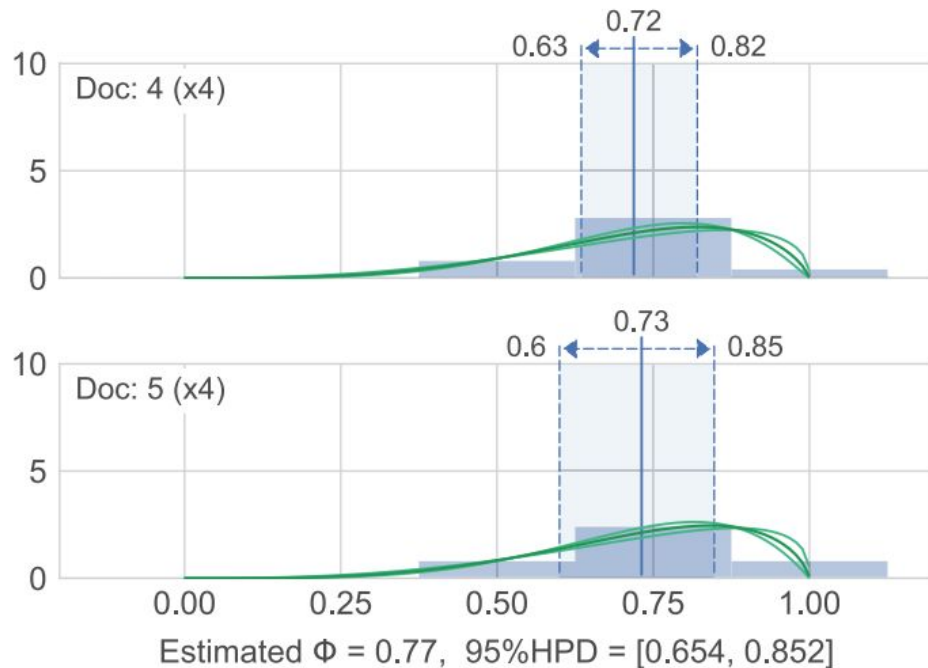
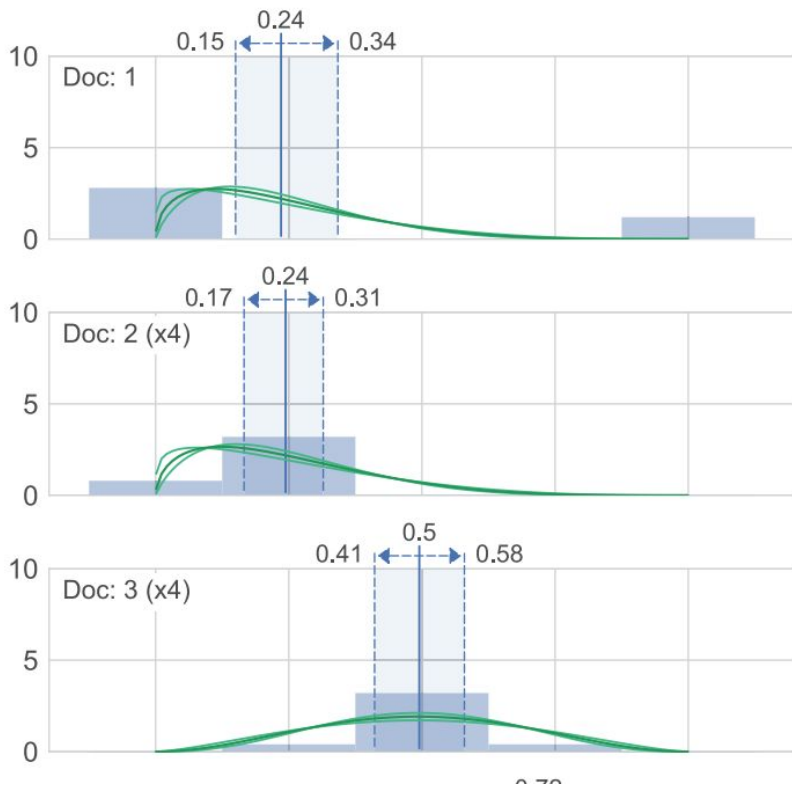
# Examples of $\Phi$ Shapes



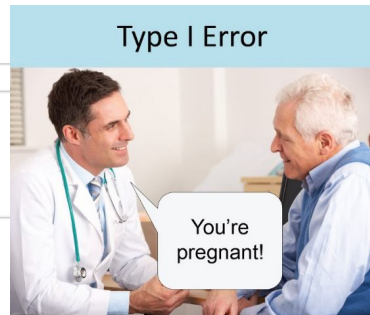
# Examples of $\Phi$ Shapes



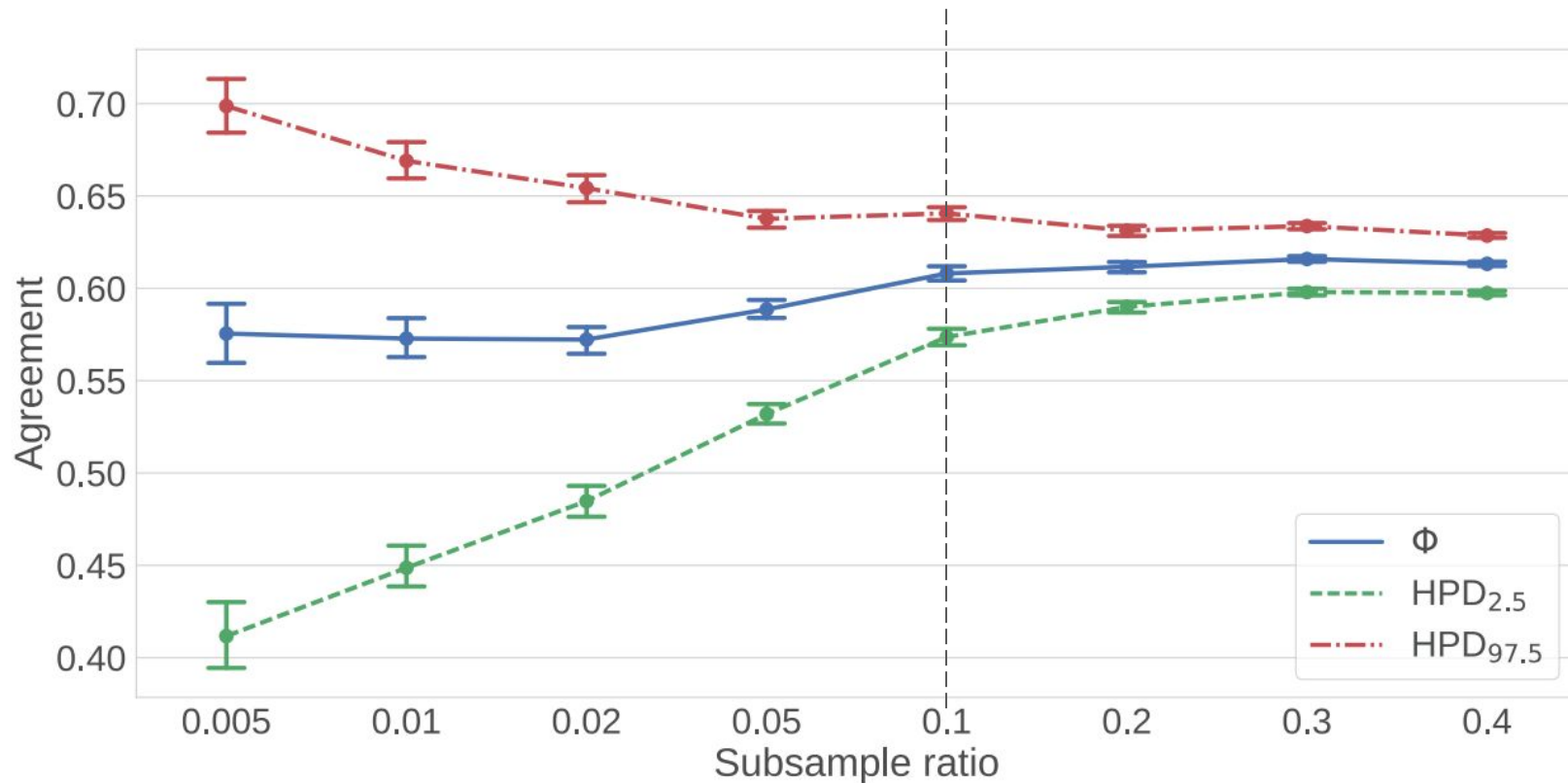
# $\Phi$ in Action on Real Data



# Robustness of $\Phi$



# Confidence Interval, Robustness of $\Phi$





## Done / Ongoing / Future Developments

- Incorporate agreement in metrics used for evaluation
- Incorporate agreement into aggregation methods
- Extend / fine tune  $\Phi$  for different scales (categorical, ratio, etc.)
- Deal with bias / reputation: different weights for different items / assessors

## Take Home Messages

- $\Phi$  is a new agreement measure
- $\Phi$  has a set of nice properties that makes it suitable for different (crowdsourcing) scenarios
- $\Phi$  can be customized and adapted to different situations

## ⊕ Properties Summary

- We have a **confidence interval** for the measure
- If we have **prior knowledge** on the domain (e.g., gold question), we could use that in the computation of the metric  
(by adding a set of priors to the model)
- We can deal with items having **different concentration points**

# Resources

- Paper: <https://www.aaai.org/ocs/index.php/HCOMP/HCOMP17/paper/viewFile/15927/15258>
- “Follow the Crowd” article: <https://blog.humancomputation.com/?p=9756>
- Python-library (pip) and GitHub Repository: <https://pypi.org/project/agreement-phi/>
- Live Demo / Online Tool: <http://agreement-measure.sheffield.ac.uk/>

## Comparison between inter-rater agreement measures

Welcome,

This demo compares the agreement between evaluators in a dataset, computed with Common Agreement Phi, Krippendorff's alpha, and Percent Agreement. Each row represents a different item, and each column represents a different assessor.

Note that this demo has a limited allocation of memory and CPU time. For extensive analyses please download the source code and run it locally.

You can input your own file or generate a random example:

### Random example

The following table contains an example randomly created:

1	3	3	2	2
3	3	5	5	1
4	2	2	5	4
4	2	3	1	4

Compute

### Upload your own file

You can also upload a csv numeric file and compute the agreement on it. Each row represents a different item, and each column represents a different assessor.

The file has to conform to the following:

- CSV format;
- no header;
- extremes of the scale appearing at least once in the file;

(You can access a version of Phi without such limitations [here](#))

Browse

Compute

### Results

Phi is: 0.273 with 95% HPD: [-0.327,0.594]  
Krippendorff's Alpha is: -0.0528  
Percent Agreement is: 0.133

Online tool

## agreement-phi 0.3.0

✓ Latest version

`pip install agreement-phi`

Last released: Aug 7, 2018

Inter-rater agreement Phi, as an alternative to Krippendorff's alpha, as described in <https://github.com/AlessandroCecco/agreement-phi>

### Navigation

Project description

Release history

Download files

### Project links

Homepage

Bug Reports

Source

### Project description

#### Agreement measure Phi

Source code for inter-rater agreement measure Phi. Live demo here: <http://agreement-measure.sheffield.ac.uk>

#### Requirements

python 3+, pymc3 3.3+. See requirements files for tested working versions on linux and osx.

#### Installation - with pip

Simply run `pip install agreement_phi`. This will provide a module and a command line executable called `run_phi`.