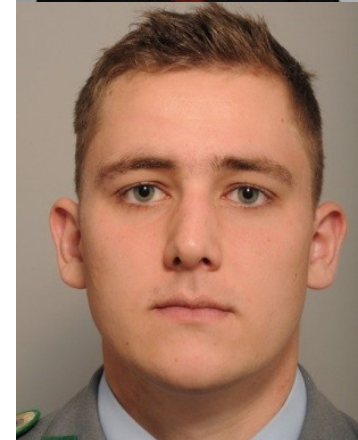


Uncertain Group t-Test

Timo von Oertzen, Tina Braun,
T. Alexander Bauer, Alexandro Folster

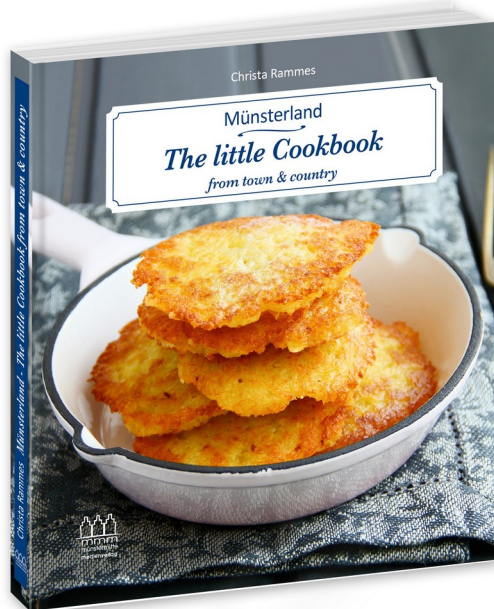
February 27th 2020



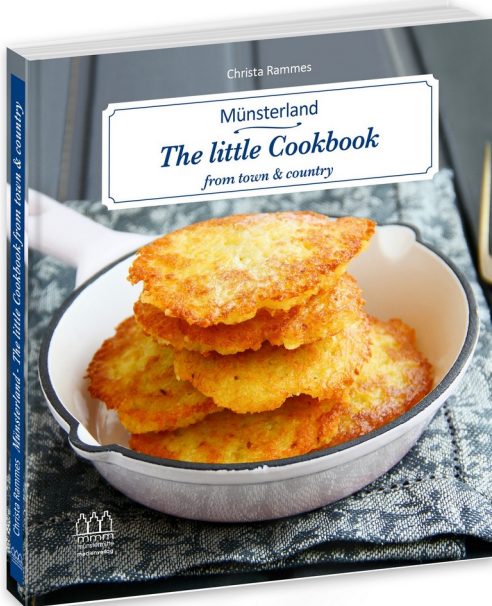
Do Aggressive People Earn Higher Salaries ?



Do Aggressive People Earn Higher Salaries ?



Do Aggressive People Earn Higher Salaries ?



Each sciene contains only as much truth as it contains mathematics.

**Roger Bacon, 1214 – 1292
(talking about Theology)**

What Do We Know ?

Committed crime	Male	Female
Total	13.20	4.32

What Do We Know ?

Committed crime	Male	Female
18 years old	3.93	1.58
21 years old	8.50	2.78
25 years old	12.04	3.70
50 years old	19.16	5.82

What Do We Know ?

Committed crime	Male	Female
18 years old	3.93	1.58
21 years old	8.50	2.78
25 years old	12.04	3.70
50 years old	19.16	5.82

- physical fitness
- neuroticism
- friends' rating
- ⋮

So I can get a probability for each of you. Does that help?

Uncertain Group t-Test: Setting

Every participant has:

x_i (dependent Variable)

p_i (probability to be in group 1)

Uncertain Group t-Test: Setting

Every participant has:

x_i (dependent Variable)

p_i (probability to be in group 1)

Assuming normality and variance homogeneity.

Uncertain Group t-Test: Setting

Every participant has:

x_i (dependent Variable)

p_i (probability to be in group 1)

Assuming normality and variance homogeneity.

We want:

estimator for the group
mean difference

Uncertain Group t-Test: Setting

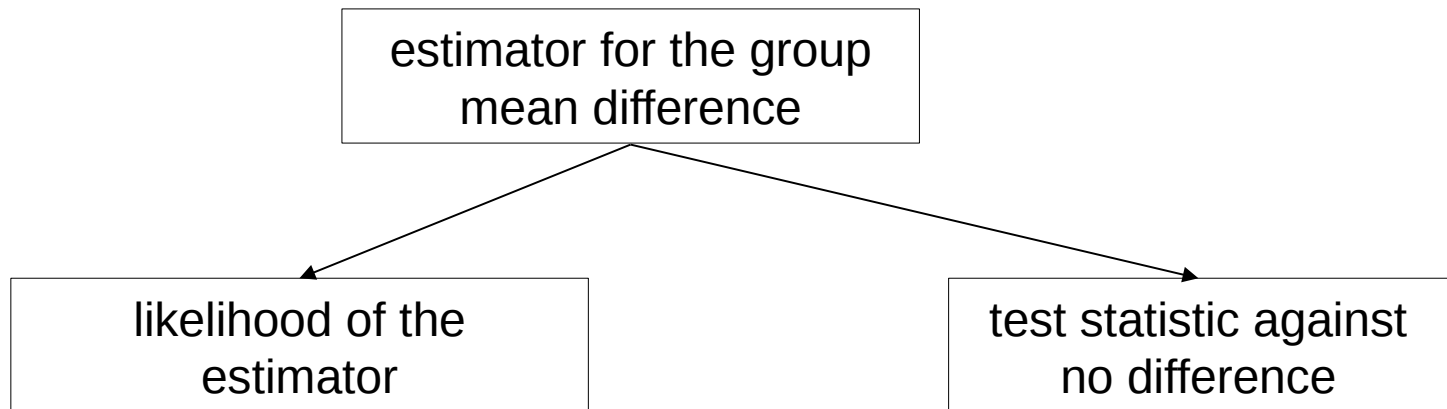
Every participant has:

x_i (dependent Variable)

p_i (probability to be in group 1)

Assuming normality and variance homogeneity.

We want:



Uncertain Group t-Test: Setting

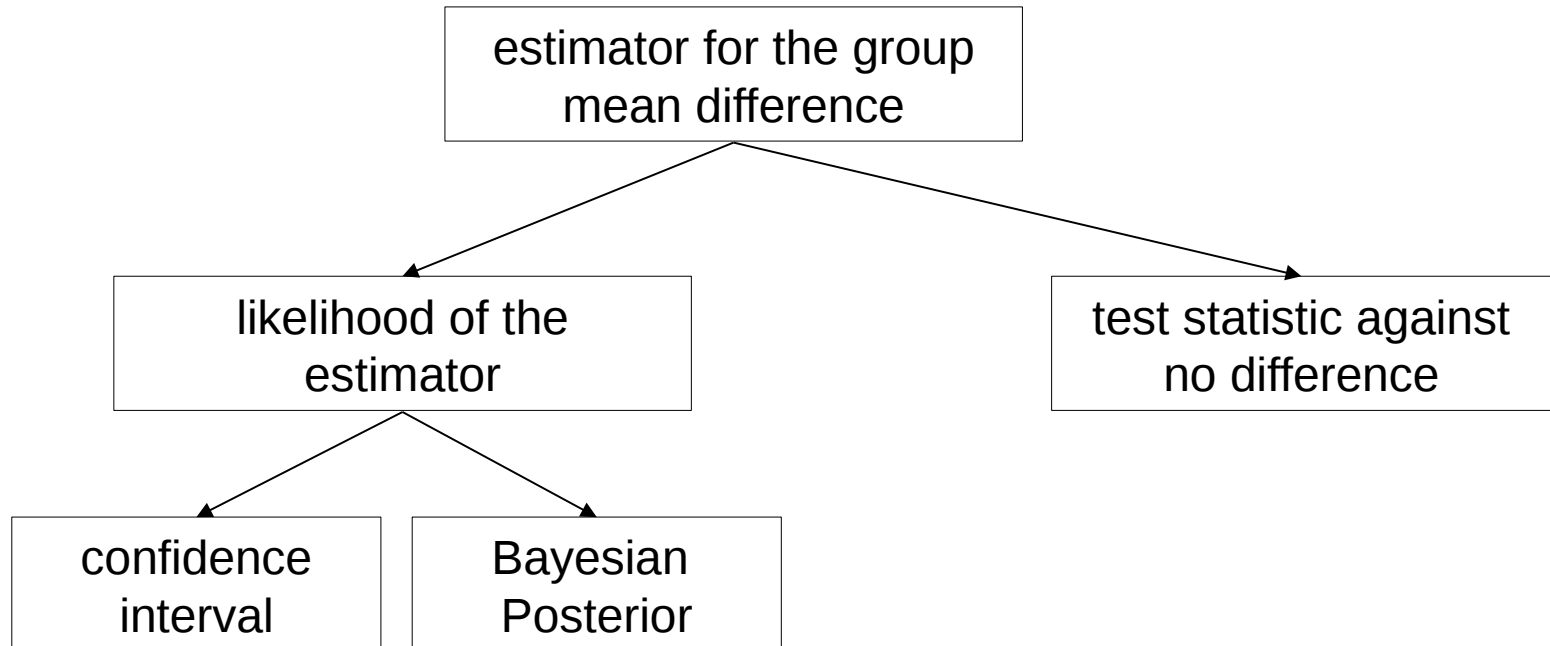
Every participant has:

x_i (dependent Variable)

p_i (probability to be in group 1)

Assuming normality and variance homogeneity.

We want:



Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i^{th} participant, and \bar{p} the average of the p_i .

$$z = \sum_{i=1}^N (p_i - \bar{p})x_i$$

Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i^{th} participant, and \bar{p} the average of the p_i .

$$z = \sum_{i=1}^N (p_i - \bar{p})x_i$$

then,

$$\mathbb{E}(z) = (\mu_1 - \mu_2)N\bar{p}$$

Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i^{th} participant, and \bar{p} the average of the p_i .

$$z = \sum_{i=1}^N (p_i - \bar{p})x_i$$

then,

$$\mathbb{E}(z) = (\mu_1 - \mu_2)N\mathbb{V}(p)$$

so

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p})x_i}{N\mathbb{V}(p)}$$

is an estimate for the group mean difference.

Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i^{th} participant, and \bar{p} the average of the p_i .

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p})x_i}{N\mathbb{V}(p)}$$

by the Central Limit Theorem, d is asymptotically normally distributed, with standard deviation

$$\frac{\hat{\sigma}}{\sqrt{N\mathbb{V}(p)}}$$

where $\hat{\sigma}$ is a $N/2$ -degree estimate of the standard deviation of the target variable.

Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i^{th} participant, and \bar{p} the average of the p_i .

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p})x_i}{\text{NV}(p)}$$

by the Central Limit Theorem, d is asymptotically normally distributed, with standard deviation

$$\frac{\hat{\sigma}}{\sqrt{\text{NV}(p)}}$$

where $\hat{\sigma}$ is a $N-2$ degree estimate of the standard deviation of the target variable. Therefore,

$$t = d \cdot \frac{\sqrt{\text{NV}(p)}}{\hat{\sigma}}$$

is t-distributed with $N-2$ degrees of freedom.

Group Mean Difference Estimator

Let x_i and p_i be the target value and probability for group 1 of the i th participant, and \bar{p} the average of the p_i .

$$t = d \cdot \frac{\sqrt{N\mathbb{V}(p)}}{\hat{\sigma}}$$

An unbiased estimate of $\hat{\sigma}$ is given by

$$\hat{\sigma}^2 = \frac{1}{N-2} \left(\left(\sum_{i=1}^N (x_i - \bar{x})^2 \right) - N\bar{p}(1-\bar{p})d^2 \right)$$

where \bar{x} is the average of the x_i .

Simulation

Data generation:

- p_i drawn from a uniform random distribution
- ,true' group drawn from a Bernoulli with parameter p_i
- x_i generated from a normal with mean dependent on group

Conditions:

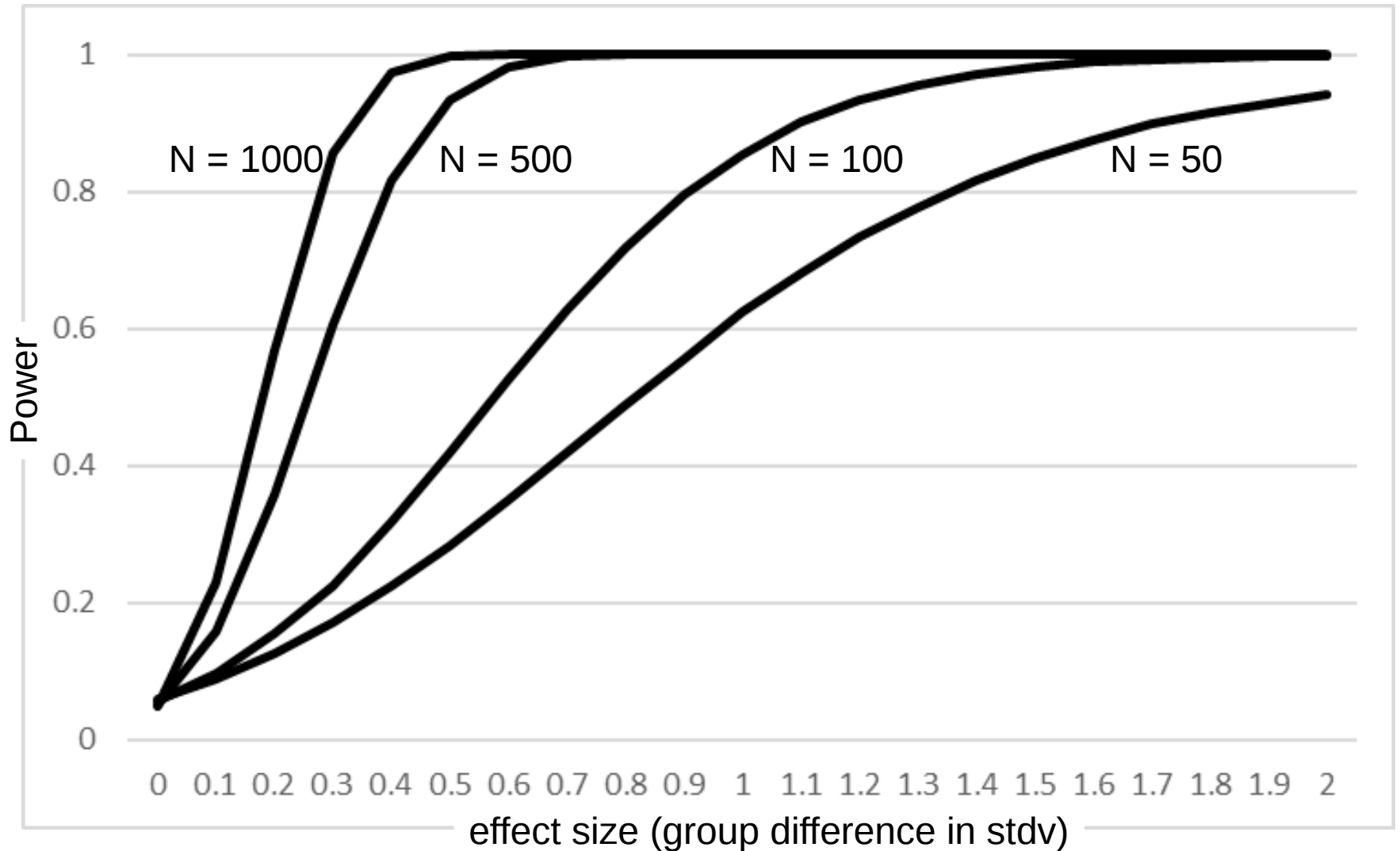
- N from {50,100,500,1000}
- (standardized) effect from 0 to 2

100,000 trials per condition

Simulation Results: α Error

N	α Error (95% confidence interval)
50	[0.0597 ; 0.0603]
100	[0.0555 ; 0.0561]
500	[0.0498 ; 0.0504]
1000	[0.0477 ; 0.0483]

Simulation Results: Power



Power Equivalence

Asymptotically, d is normally distributed and $\hat{\sigma}$ known. The standard error is

$$\frac{\hat{\sigma}}{\sqrt{N\mathbb{V}(p)}}$$

Assuming $p \sim B(\alpha, \alpha)$ we have

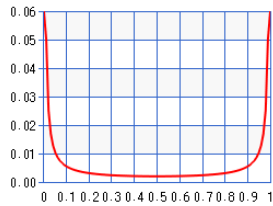
$$\mathbb{V}(p) = \frac{1}{4(2\alpha + 1)}$$

Power Equivalence

$$\mathbb{V}(p) = \frac{1}{4(2\alpha + 1)}$$

$$\text{stderr}(d) = \frac{\hat{\sigma}}{\sqrt{N\mathbb{V}(p)}}$$

$$\mathbb{V}(p)$$
$$\mathbb{V}(p) = \frac{1}{4}$$



Classical t-test
 $p_i \in \{0, 1\}$

α

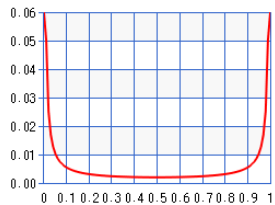
Power Equivalence

$$\mathbb{V}(p) = \frac{1}{4(2\alpha + 1)}$$

$$\text{stderr}(d) = \frac{\hat{\sigma}}{\sqrt{N\mathbb{V}(p)}}$$

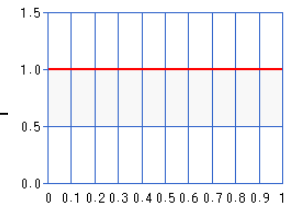
$$\mathbb{V}(p)$$

$$\mathbb{V}(p) = \frac{1}{4}$$



Classical t-test
 $p_i \in \{0, 1\}$

$$\mathbb{V}(p) = \frac{1}{12}$$



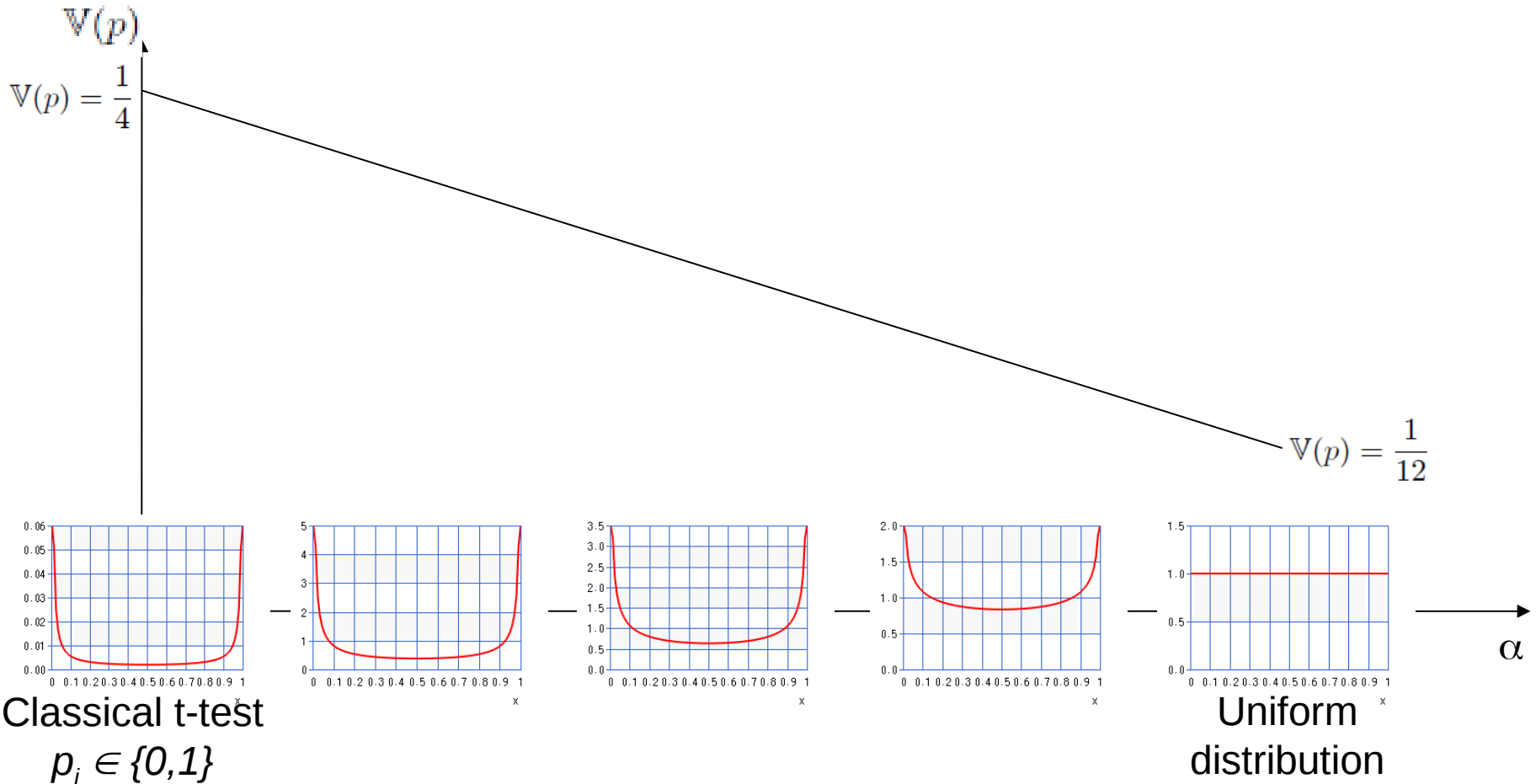
Uniform
distribution

→
 α

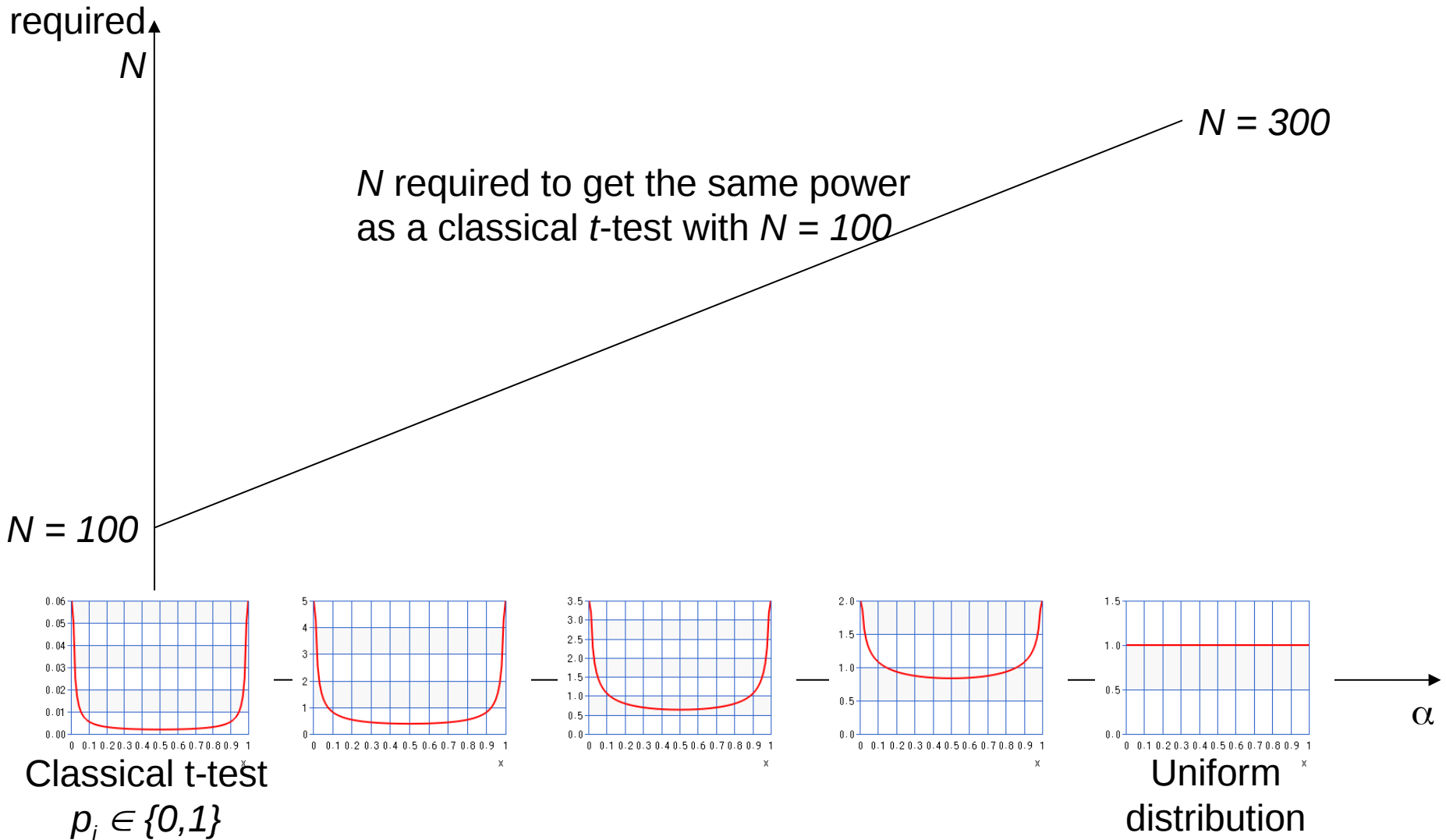
Power Equivalence

$$\mathbb{V}(p) = \frac{1}{4(2\alpha + 1)}$$

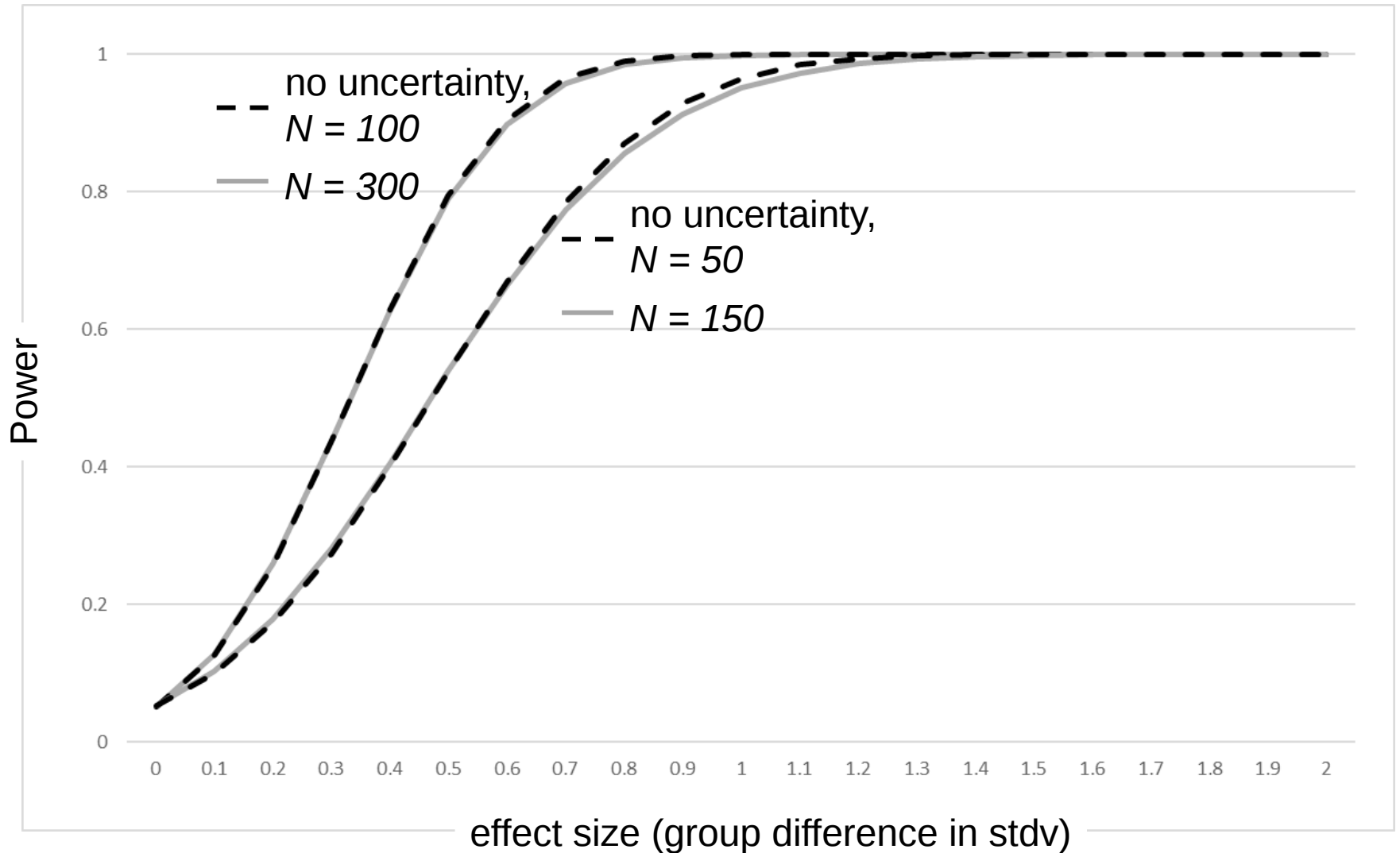
$$\text{stderr}(d) = \frac{\hat{\sigma}}{\sqrt{N\mathbb{V}(p)}}$$



Power Equivalence



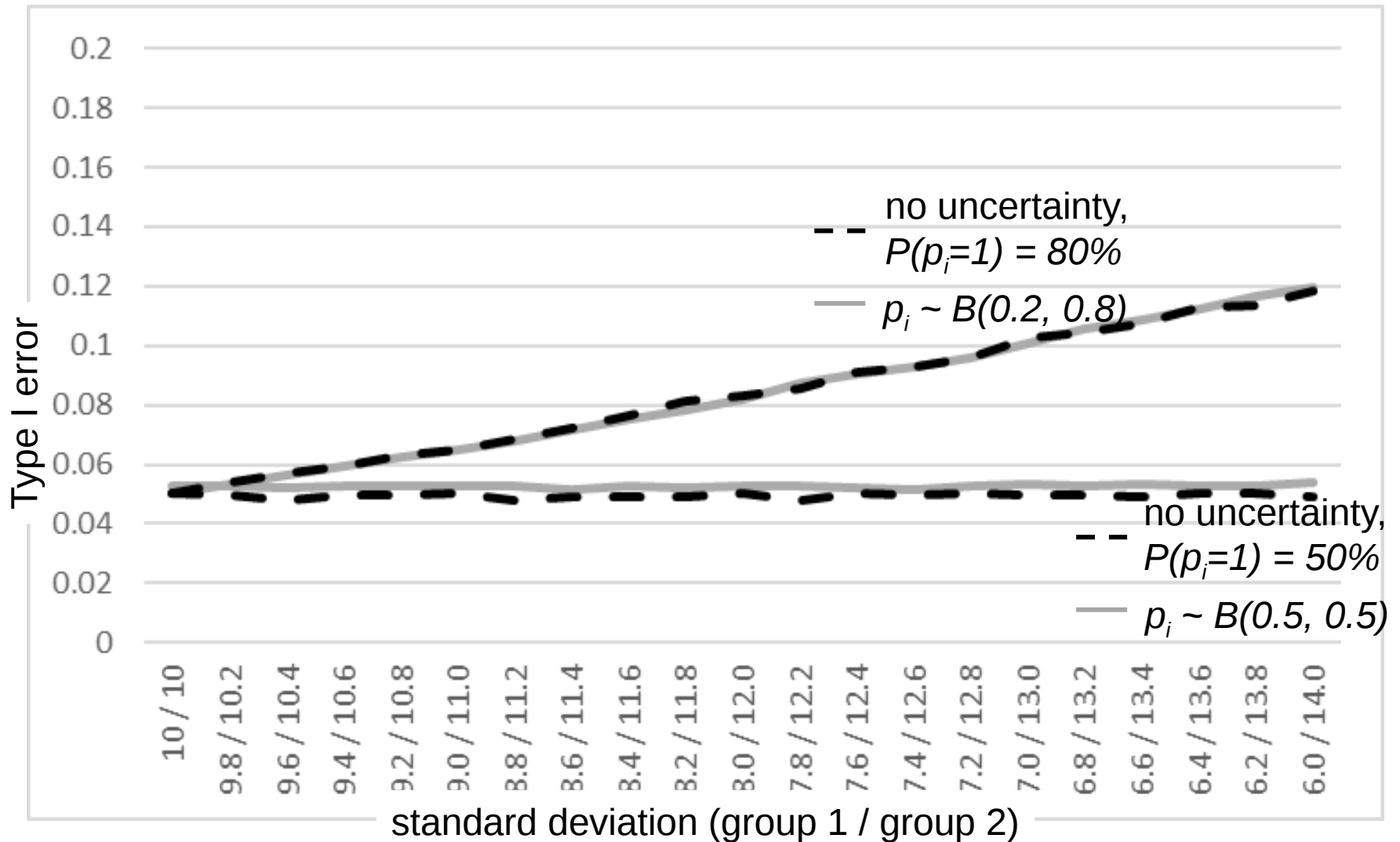
Power Equivalence Simulation



Assumption Violation Variance Homogeneity

If the standard deviation of group 1 and group 2 differ, we expect an α inflation both for the classical t -Test and the uncertain group t -Test.

Assumption Violation Variance Homogeneity

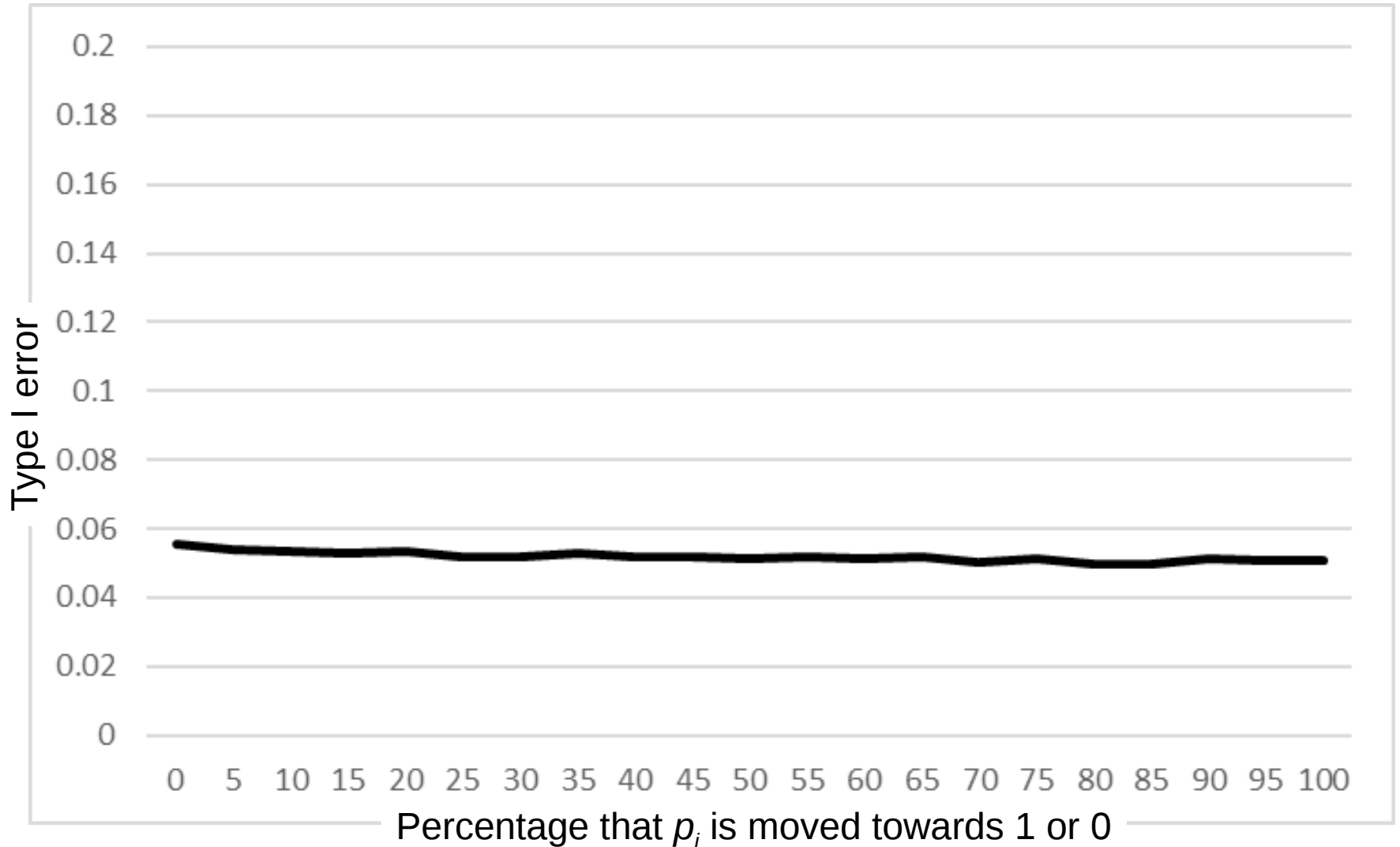


Assumption Violation

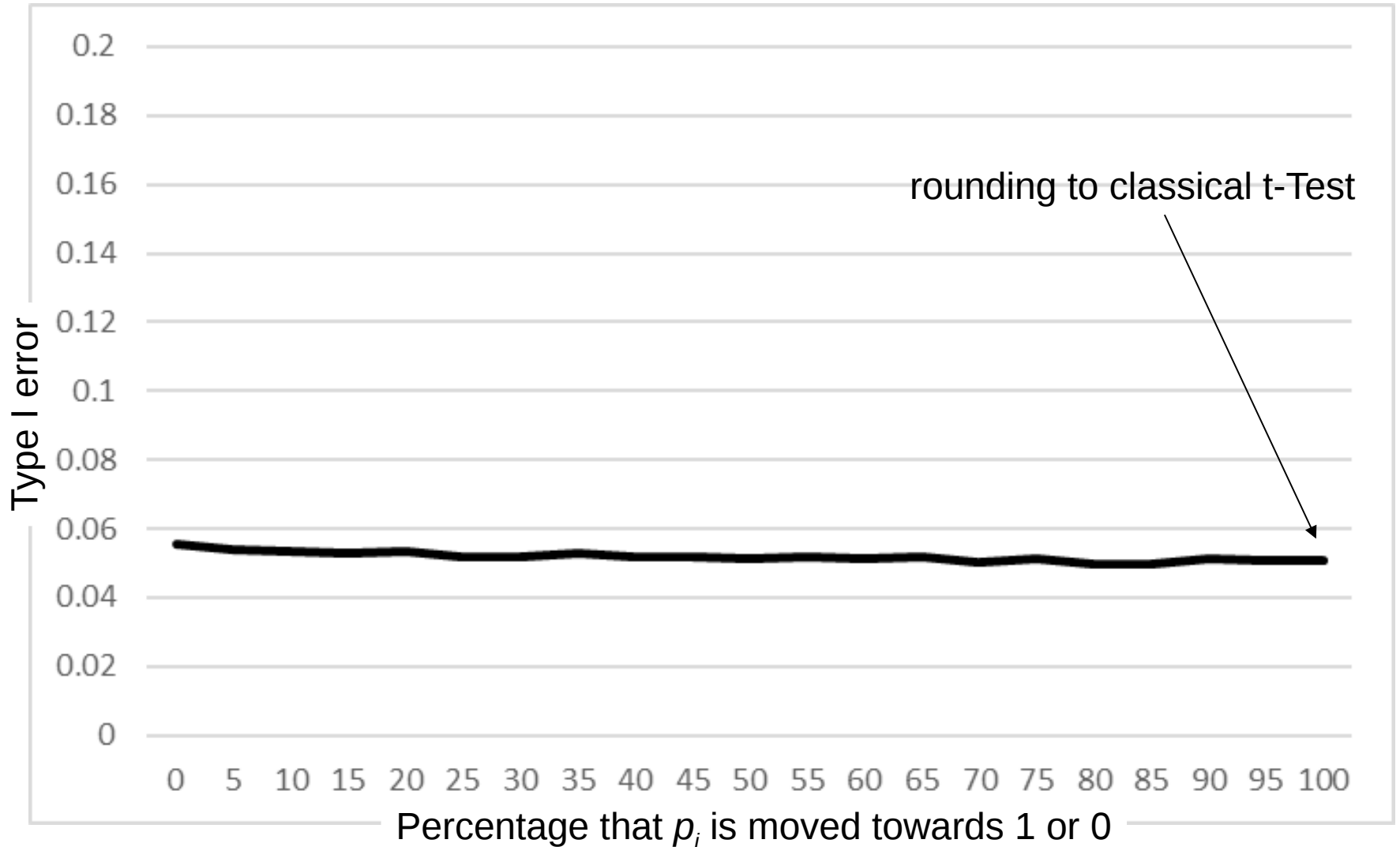
Exaggerated p_i

If we are too confident in the probability estimates, do we have an α inflation?

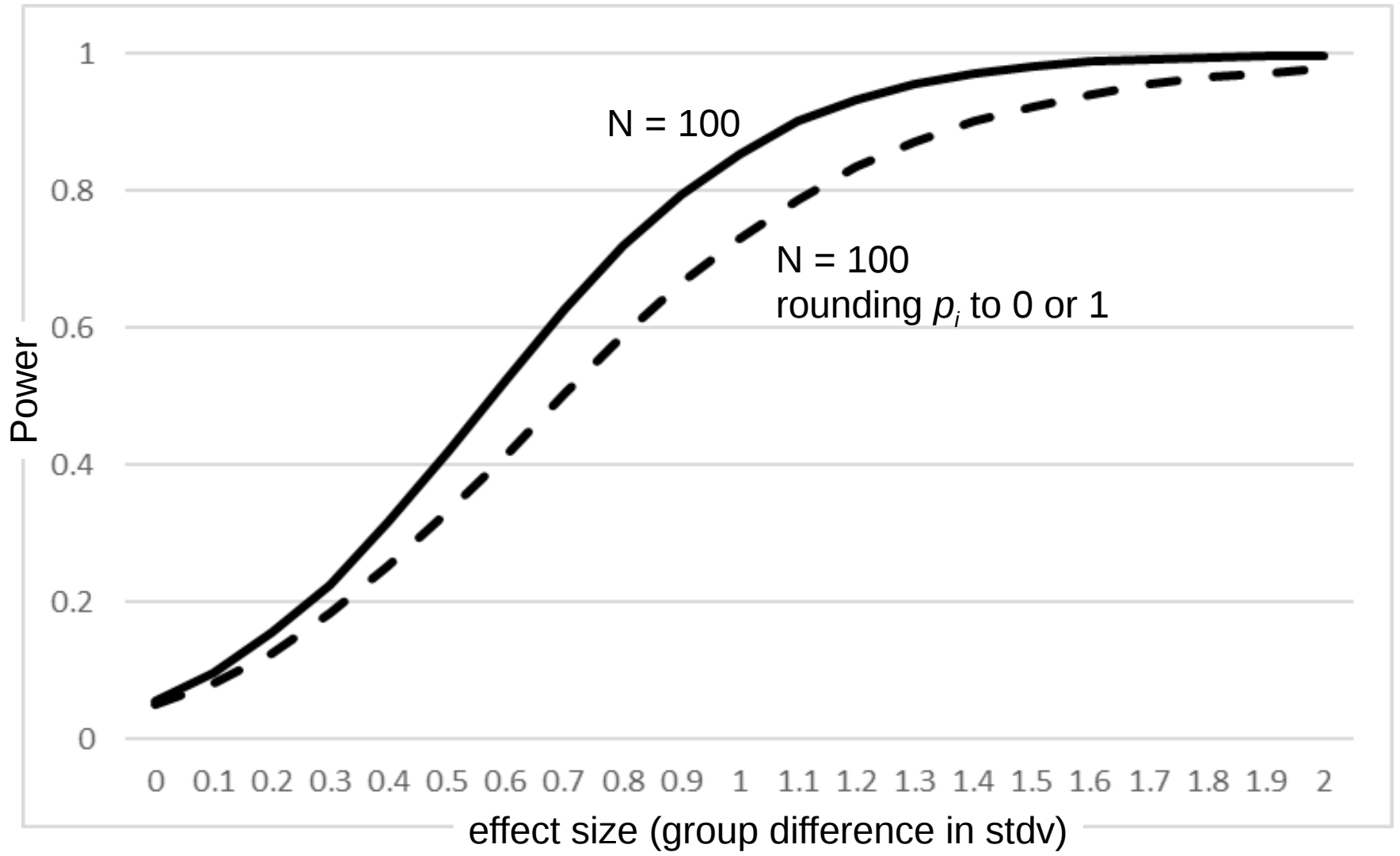
Assumption Violation Exaggerated p_i



Assumption Violation Exaggerated p_i



Simulation Results: Power in Comparison to standard t -Test



Summary

Fairly simple computations allow a mean comparison between two groups even if we don't know the group of any participant.

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p}) x_i}{N\mathbb{V}(p)}$$

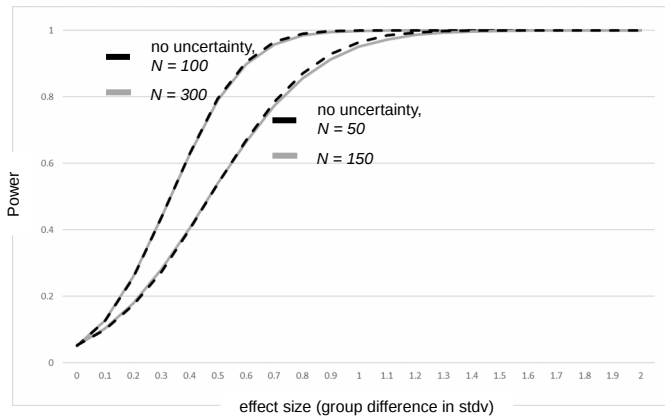
$$t = d \cdot \frac{\sqrt{N\mathbb{V}(p)}}{\hat{\sigma}}$$

Summary

Fairly simple computations allow a mean comparison between two groups even if we don't know the group of any participant.

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p}) x_i}{N\mathbb{V}(p)}$$

$$t = d \cdot \frac{\sqrt{N\mathbb{V}(p)}}{\hat{\sigma}}$$



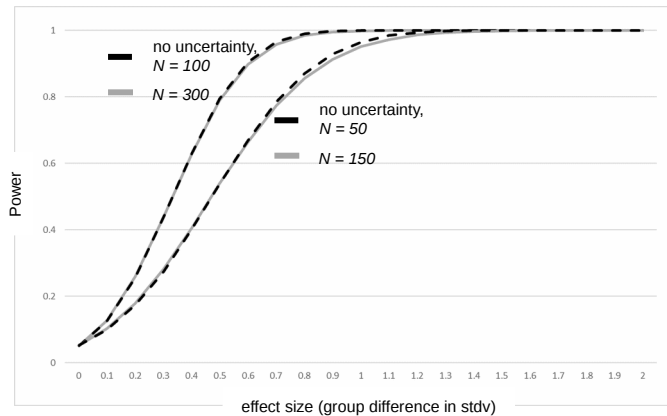
With uniformly distributed probabilities, we need three times as many participants to have the same power as if we know the groups.

Summary

Fairly simple computations allow a mean comparison between two groups even if we don't know the group of any participant.

$$d = \frac{\sum_{i=1}^N (p_i - \bar{p}) x_i}{N\mathbb{V}(p)}$$

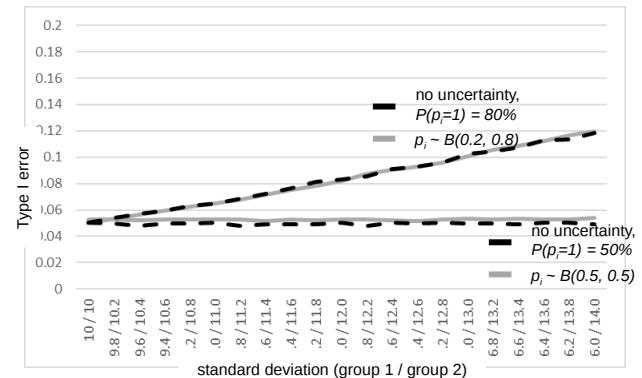
$$t = d \cdot \frac{\sqrt{N\mathbb{V}(p)}}{\hat{\sigma}}$$



With uniformly distributed probabilities, we need three times as many participants to have the same power as if we know the groups.

Lack of variance homogeneity is equally bad as with the standard t-Test.

Exaggeration p-values loses power, but not correctness.



Thank You !