# Improved Differentially Private Analysis of Variance
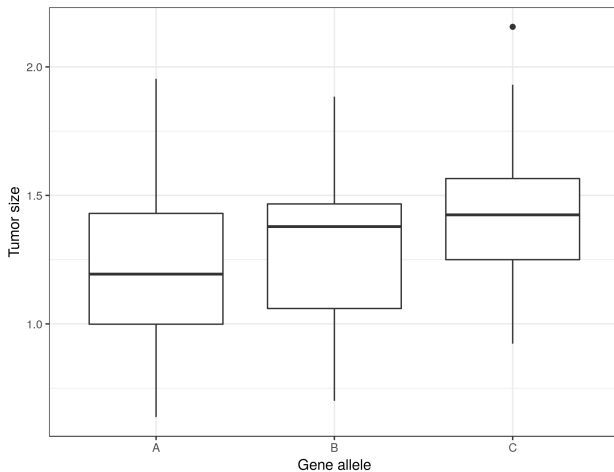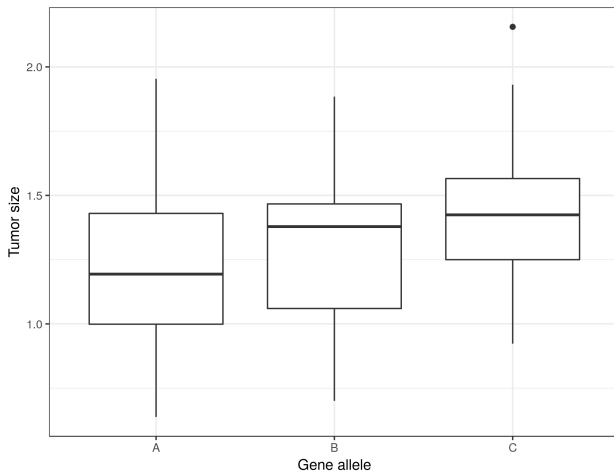
Marika Swanberg     Ira Globus-Harris     Iris Griffith
Anna Ritz     Andrew Bray     Adam Groce

Reed College

# Observed Data, $n = 30$

# Observed Data, $n = 30$



Are gene allele and tumor size dependent?

# Metric for dependency

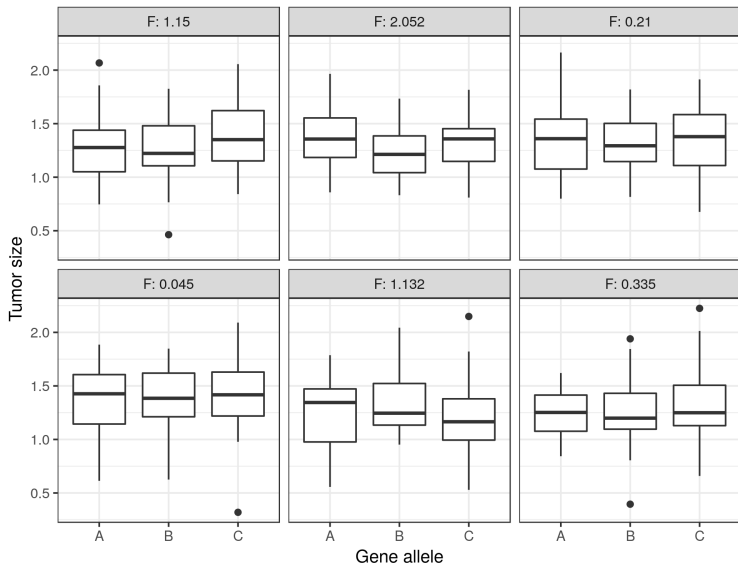$$F\text{-test} = \frac{\text{Variation between groups}}{\text{In-group variation}}$$

# Metric for dependency

$$F\text{-test} = \frac{\text{Variation between groups}}{\text{In-group variation}}$$

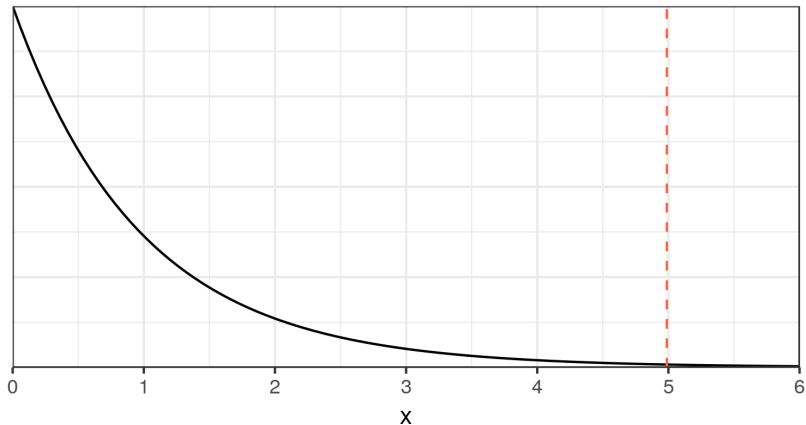$$SSA(D) = \sum_{j=1}^{k} n_j(\bar{y}_j - \bar{y})^2/(k-1)$$

$$SSE(D) = \sum_{i=1}^{n} (y_i - \bar{y}_{c_i})^2/(n-k)$$

# Simulate Random Data

## Reference Distribution of F



Now do we think our data supports independence of gene allele and tumor size?

# Why is $F$-test optimal?

# Why is *F*-test optimal?

High probability of indicating dependence when variables <u>are</u> dependent . . .

# Why is *F*-test optimal?

High probability of indicating dependence when variables <u>are</u> dependent . . .

. . . even when dataset is small.

# Why is *F*-test optimal?

High probability of indicating dependence when variables <u>are</u> dependent . . .

. . . even when dataset is small.

### Definition (Power)

The **power** of a hypothesis test is the probability it rejects $H_0$. It depends on the alternate distribution $H_A$ and $n$.

# Why is $F$-test optimal?

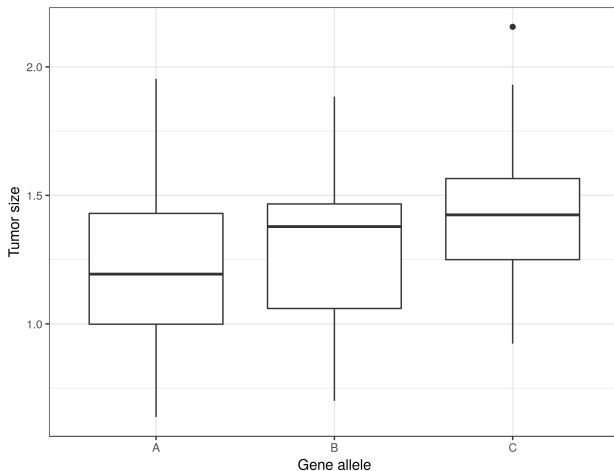High probability of indicating dependence when variables <u>are</u> dependent . . .

. . . even when dataset is small.

### Definition (Power)

The **power** of a hypothesis test is the probability it rejects $H_0$. It depends on the alternate distribution $H_A$ and $n$.

Goal of any test statistic is achieving high power[*].

# Observed Data, $n = 30$



What if we want to keep this data private?

# Differential privacy [DMNS06]

### Definition

Two databases are **neighboring** if they differ only in the data of one
individual.

# Differential privacy [DMNS06]

### Definition

Two databases are **neighboring** if they differ only in the data of one individual.

### Definition

A query $f$ is $\varepsilon$-**differentially private** if for all neighboring databases $D, D'$ and all output sets $S$

$$\Pr[f(D) \in S] \leq e^{\varepsilon} \Pr[f(D') \in S].$$

# Properties of differential privacy [DMNS06]

### Theorem (Post-processing)

*If f is $\varepsilon$-differentially private then for any (randomized) function g, then if $h(D) = g(f(D))$, h is also $\varepsilon$-differentially private.*

# Properties of differential privacy [DMNS06]

## Theorem (Post-processing)

*If f is $\varepsilon$-differentially private then for any (randomized) function g, then if $h(D) = g(f(D))$, h is also $\varepsilon$-differentially private.*

## Theorem (Composition)

*If f is $\varepsilon_1$-differentially private and g is $\varepsilon_2$-differentially private then if $h(D) = (g(D), f(D))$, h is $(\varepsilon_1 + \varepsilon_2)$-differentially private.*

# Laplace mechanism

### Definition (Sensitivity)

The sensitivity $\Delta f$ of a deterministic, real-valued function $f$ on databases is the maximum over all pairs of neighboring $D, D'$ of $|f(D) - f(D')|$.

# Laplace mechanism

### Definition (Sensitivity)

The sensitivity $\Delta f$ of a deterministic, real-valued function $f$ on databases is the maximum over all pairs of neighboring $D, D'$ of $|f(D) - f(D')|$.

### Theorem (Laplace Mechanism)

*Given any deterministic, real-valued function $f$ on databases, define $\widehat{f}$ as*

$$\widehat{f}(D) = f(D) + Y,$$

*where $Y \leftarrow \text{Lap}(\Delta f / \varepsilon)$. The Laplace mechanism is $\varepsilon$-differentially private.*

# Related Works

Other work on private hypothesis testing:

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]

## Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - Binomial data [AS18] (Proven optimal!)

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - Binomial data [AS18] (Proven optimal!)
  - Difference of two means [OHK15, DNLI18]

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - ▶ Binomial data [AS18] (Proven optimal!)
  - ▶ Difference of two means [OHK15, DNLI18]
  - ▶ Linear regression [BRMC17, Sheffet17]

## Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - Binomial data [AS18] (Proven optimal!)
  - Difference of two means [OHK15, DNLI18]
  - Linear regression [BRMC17, Sheffet17]

Earlier work is often missing:

## Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - Binomial data [AS18] (Proven optimal!)
  - Difference of two means [OHK15, DNLI18]
  - Linear regression [BRMC17, Sheffet17]

Earlier work is often missing:

- Rigorous p-value computations

# Related Works

Other work on private hypothesis testing:

- Asymptotic analysis [WZ10, Smith11, CKMSU19]
- Chi-squared test (difference of discrete distributions) [VS09, FSU11, JS13, USF13, WLK15, GLRV16, RK17]
- Other tests:
  - ▶ Binomial data [AS18] (Proven optimal!)
  - ▶ Difference of two means [OHK15, DNLI18]
  - ▶ Linear regression [BRMC17, Sheffet17]

Earlier work is often missing:

- Rigorous p-value computations
- Power analysis

# Private $F$-statistic [CBRG18]

Assume data is on the $[0, 1]$ interval.

# Private *F*-statistic [CBRG18]

Assume data is on the $[0, 1]$ interval.

**Theorem**

*SSE has sensitivity bounded by 7.*

# Private $F$-statistic [CBRG18]

Assume data is on the $[0, 1]$ interval.

**Theorem**

*SSE has sensitivity bounded by 7.*

$$\widehat{SSE}(D) = SSE(D) + \mathsf{Lap}(7/\varepsilon)$$

# Private $F$-statistic [CBRG18]

Assume data is on the $[0, 1]$ interval.

**Theorem**

*SSE has sensitivity bounded by 7.*

$$\widehat{SSE}(D) = SSE(D) + \mathsf{Lap}(7/\varepsilon)$$

**Theorem**

*SSA has sensitivity bounded by $9 + 5/n$.*

# Private $F$-statistic [CBRG18]

Assume data is on the $[0, 1]$ interval.

### Theorem
*SSE has sensitivity bounded by 7.*

$$\widehat{SSE}(D) = SSE(D) + \text{Lap}(7/\varepsilon)$$

### Theorem
*SSA has sensitivity bounded by $9 + 5/n$.*

$$\widehat{SSA}(D) = SSA(D) + \text{Lap}\left(\frac{9 + 5/n}{\varepsilon}\right)$$
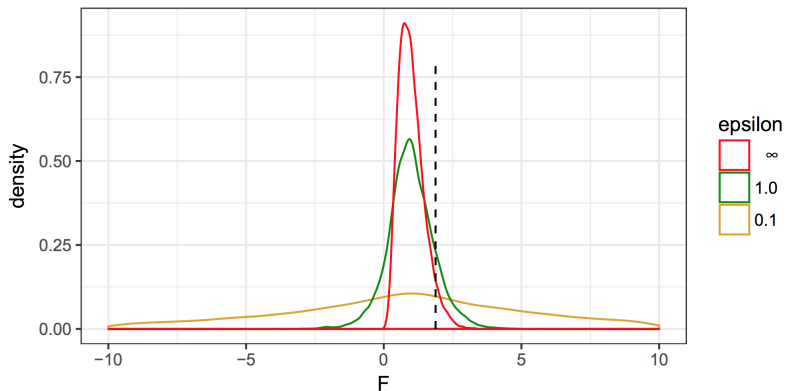
# Private ANOVA [CBRG18]

$$\widehat{F}(D) = \frac{\widehat{SSA}(D)/(k-1)}{\widehat{SSE}(D)/(n-k)}$$

# Private ANOVA [CBRG18]

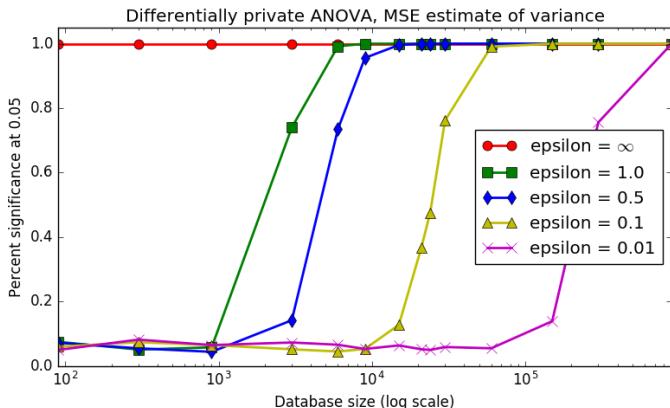$$\widehat{F}(D) = \frac{\widehat{SSA}(D)/(k-1)}{\widehat{SSE}(D)/(n-k)}$$

Problem: What is the reference distribution now?

# Private ANOVA [CBRG18]



Public reference distribution gives inaccurate *p*-values.

# Private ANOVA [CBRG18]



Differentially private ANOVA, MSE estimate of variance

## Definition (Power)

The **power** of a hypothesis test is the probability it rejects $H_0$. It depends on the alternate distribution $H_A$ and $n$.

# Improving ANOVA [SHGRGB19]

Are there other ways of measuring "dispersion" (analogous to variance)?

# Improving ANOVA [SHGRGB19]

Are there other ways of measuring "dispersion" (analogous to variance)?

$$(x_1 - x_2)^2, \quad |x_1 - x_2|, \text{ or maybe } |x_1 - x_2|?$$

$$SSA(D) = \sum_{j=1}^{k} n_j(\bar{y}_j - \bar{y})^2 \implies SA(D) = \sum_{j=1}^{k} n_j|\bar{y}_j - \bar{y}|$$

$$SSE(D) = \sum_{i=1}^{n}(y_i - \bar{y}_{c_i})^2 \implies SE(D) = \sum_{i=1}^{n} |y_i - \bar{y}_{c_i}|$$

$$F(D) = \frac{SSA(D)/(k-1)}{SSE(D)/(n-k)} \implies F_1(D) = \frac{SA(D)/(k-1)}{SE(D)/(n-k)}$$

# Improving ANOVA [SHGRGB19]

Are there other ways of measuring "dispersion" (analogous to variance)?

$$(x_1 - x_2)^2, \quad |x_1 - x_2|, \text{ or maybe } |x_1 - x_2|^?$$

$$SSA(D) = \sum_{j=1}^{k} n_j (\bar{y}_j - \bar{y})^2 \implies SA(D) = \sum_{j=1}^{k} n_j |\bar{y}_j - \bar{y}|$$

$$SSE(D) = \sum_{i=1}^{n} (y_i - \bar{y}_{c_i})^2 \implies SE(D) = \sum_{i=1}^{n} |y_i - \bar{y}_{c_i}|$$

$$F(D) = \frac{SSA(D)/(k-1)}{SSE(D)/(n-k)} \implies F_1(D) = \frac{SA(D)/(k-1)}{SE(D)/(n-k)}$$

The new $F_1$ statistic has:

# Improving ANOVA [SHGRGB19]

Are there other ways of measuring "dispersion" (analogous to variance)?

$$(x_1 - x_2)^2, \quad |x_1 - x_2|, \text{ or maybe } |x_1 - x_2|^?$$

$$SSA(D) = \sum_{j=1}^{k} n_j(\bar{y}_j - \bar{y})^2 \implies SA(D) = \sum_{j=1}^{k} n_j|\bar{y}_j - \bar{y}|$$

$$SSE(D) = \sum_{i=1}^{n}(y_i - \bar{y}_{c_i})^2 \implies SE(D) = \sum_{i=1}^{n} |y_i - \bar{y}_{c_i}|$$

$$F(D) = \frac{SSA(D)/(k-1)}{SSE(D)/(n-k)} \implies F_1(D) = \frac{SA(D)/(k-1)}{SE(D)/(n-k)}$$

The new $F_1$ statistic has:

- Lower sensitivity (3 for $SE$, 4 for $SA$)

# Improving ANOVA [SHGRGB19]

Are there other ways of measuring "dispersion" (analogous to variance)?

$$(x_1 - x_2)^2, \quad |x_1 - x_2|, \text{ or maybe } |x_1 - x_2|^?$$

$$SSA(D) = \sum_{j=1}^{k} n_j(\bar{y}_j - \bar{y})^2 \implies SA(D) = \sum_{j=1}^{k} n_j|\bar{y}_j - \bar{y}|$$

$$SSE(D) = \sum_{i=1}^{n}(y_i - \bar{y}_{c_i})^2 \implies SE(D) = \sum_{i=1}^{n} |y_i - \bar{y}_{c_i}|$$

$$F(D) = \frac{SSA(D)/(k-1)}{SSE(D)/(n-k)} \implies F_1(D) = \frac{SA(D)/(k-1)}{SE(D)/(n-k)}$$

The new $F_1$ statistic has:

- Lower sensitivity (3 for $SE$, 4 for $SA$)
- Much higher typical value

Making $F_1$ private

# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \mathsf{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \mathsf{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F}_1 = \frac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \mathsf{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \mathsf{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F_1} = \dfrac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

3. Private estimate of standard deviation:
   - Allocate some of epsilon budget?

# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \mathsf{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \mathsf{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F}_1 = \frac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

3. Private estimate of standard deviation:
   - Allocate some of epsilon budget? Makes power worse

# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \mathrm{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \mathrm{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F}_1 = \frac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

3. Private estimate of standard deviation:
   - Allocate some of epsilon budget? Makes power worse
   - Solution: derive an unbiased estimator for $\sigma$
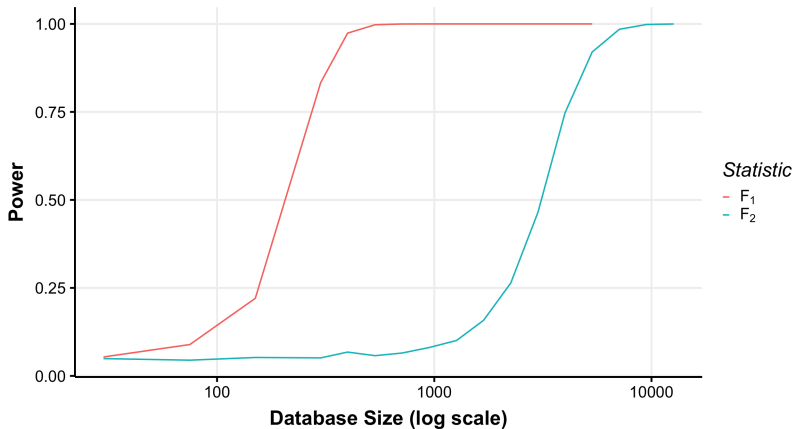
# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \text{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \text{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F_1} = \frac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

3. Private estimate of standard deviation:
   - Allocate some of epsilon budget? Makes power worse
   - Solution: derive an unbiased estimator for $\sigma$

$$\hat{\sigma} = \sqrt{\pi/2} \cdot \frac{\widehat{SE}}{(N-k)}$$

# Improving ANOVA [SHGRGB19]

Making $F_1$ private

1. Use Laplace mechanism on $SE$ and $SA$
   - $\widehat{SA} = F_1 = SA + \text{Lap}(4/\rho\varepsilon)$
   - $\widehat{SE} = F_1 = SA + \text{Lap}(3/(1-\rho)\varepsilon)$
   - $\widehat{F}_1 = \frac{\widehat{SA}/(k-1)}{\widehat{SE}/(n-k)}$

2. Use simulation for reference distribution
   - Problem: need standard deviation

3. Private estimate of standard deviation:
   - Allocate some of epsilon budget? Makes power worse
   - Solution: derive an unbiased estimator for $\sigma$

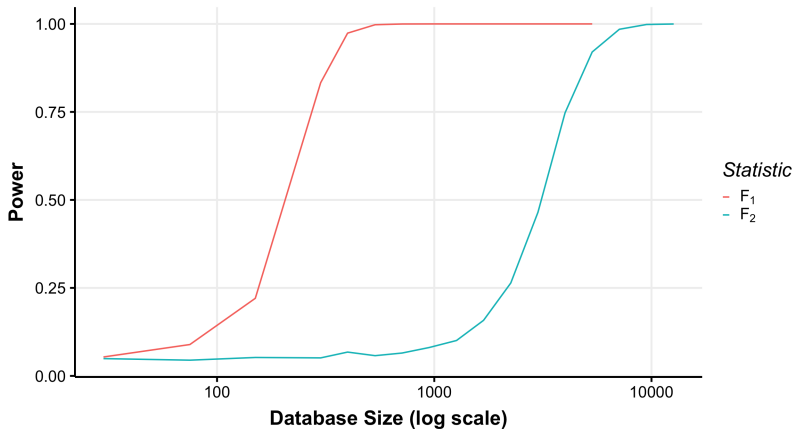$$\hat{\sigma} = \sqrt{\pi/2} \cdot \frac{\widehat{SE}}{(N-k)}$$

Note: empirically verified to have valid *p*-values.

Power comparison of $F_1$ and $F_2$ statistics

Power comparison of $F_1$ and $F_2$ statistics

Optimal public test $\;\not\!\!\Longrightarrow\;$ optimal private test [1].

[1] In figure: $\varepsilon = 1, \mu = [0.35, 0.5, 0.65],$ and $\sigma = 0.15$

# Further Optimization

# New Developments [CKSBG19]

Kruskal-Wallis test analogous to F-test
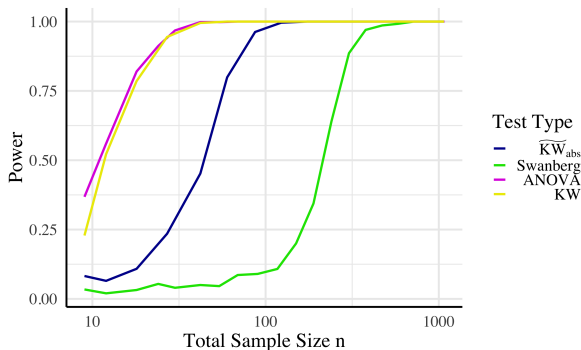
# New Developments [CKSBG19]

Kruskal-Wallis test analogous to F-test

Modified KW test, similar methods as [SHGRGB19]

# New Developments [CKSBG19]

Kruskal-Wallis test analogous to F-test

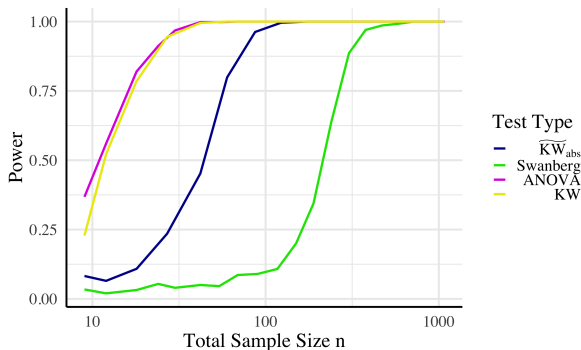Modified KW test, similar methods as [SHGRGB19]



For 80% power, need only 23% as much data as $F_1$ ([SHGRGB19]) ...

# New Developments [CKSBG19]

Kruskal-Wallis test analogous to F-test

Modified KW test, similar methods as [SHGRGB19]



For 80% power, need only 23% as much data as $F_1$ ([SHGRGB19]) ... and about 1-2% as much data as $F_2$ ([CBRG18])

# Thank you