





Baselines

Assume that you have three independendent measurements X_1, X_2, X_3 , with $var(X_i) = \sigma^2$.

 \longrightarrow It follows that $var(\bar{X}) = \sigma^2/3$.

Assume that you have to subtract a baseline *B* from each measurement. What is $var(X_i - B)$?

If B is just a constant, then

 \longrightarrow var $(X_i - B) =$ var $(X_i) = \sigma^2$.

If *B* is a measurement with $var(B) = \sigma_B^2$, then $\rightarrow var(X_i - B) = var(X_i) + var(B) = \sigma^2 + \sigma_B^2$.

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

[140.615]

Baselines

What is the variance of the average of those values, i.e. what is the variance of $\sum_{i} (X_i - B)/3$..?

$$\operatorname{var}\left(\sum_{i} (X_{i} - B)/3\right) = \operatorname{var}\left(\frac{1}{3} (X_{1} + X_{2} + X_{3} - 3 \times B)\right)$$
$$= \operatorname{var}\left(\bar{X} - B\right) = \operatorname{var}(\bar{X}) + \operatorname{var}(B)$$
$$= \frac{\sigma^{2}}{3} + \sigma^{2}_{B}.$$

In other words: if you use the same baseline in all experiments, you will never be able to get rid of the baseline error, regardless how many replicates you use!

<text><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block><equation-block>

t-tests

Suppose that

• $X_1, X_2, ..., X_n$ are iid Normal(mean= μ_A , SD= σ), and • $Y_1, Y_2, ..., Y_m$ are iid Normal(mean= μ_B , SD= σ).

Then

$$\longrightarrow \mathsf{E}(\overline{X} - \overline{Y}) = \mathsf{E}(\overline{X}) - \mathsf{E}(\overline{Y}) = \mu_{\mathsf{A}} - \mu_{\mathsf{B}}$$

$$\longrightarrow SD(\overline{X} - \overline{Y}) = \sqrt{SD(\overline{X})^2 + SD(\overline{Y})^2} = \sqrt{\left(\frac{\sigma}{\sqrt{n}}\right)^2 + \left(\frac{\sigma}{\sqrt{m}}\right)^2} = \sigma\sqrt{\frac{1}{n} + \frac{1}{m}}$$

Note: If n = m, then $SD(\overline{X} - \overline{Y}) = \sigma\sqrt{2/n}$.

t-tests

$$\widehat{SD}(\overline{X} - \overline{Y}) = \widehat{\sigma}_{\text{pooled}} \sqrt{\frac{1}{n} + \frac{1}{m}}$$
$$= \sqrt{\left[\frac{S_{A}^{2}(n-1) + S_{B}^{2}(m-1)}{n+m-2}\right] \cdot \left[\frac{1}{n} + \frac{1}{m}\right]}$$

In the case n = m,

$$\widehat{\mathrm{SD}}(\overline{X} - \overline{Y}) = \sqrt{\frac{S_{\mathrm{A}}^2 + S_{\mathrm{B}}^2}{n}}$$

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

t-tests

If the null hypothesis is true and there there are no differences in group means, then



Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017





Hypothesis testing

We consider two hypotheses:

Null hypothesis, H₀: $\mu = 0$ Alternative hypothesis, H_a: $\mu \neq 0$

Type I error: Reject H₀ when it is true (false positive)

Type II error: Fail to reject H₀ when it is false (false negative)

We set things up so that a Type I error is a worse error (and so that we are seeking to prove the alternative hypothesis). We want to control the rate (the significance level, α) of such errors.

Two-sample t-test



ARTICLE

PMID 17357068

Overcoming the Winner's Curse: Estimating Penetrance Parameters from Case-Control Data

Sebastian Zöllner and Jonathan K. Pritchard

Genomewide association studies are now a widely used approach in the search for loci that affect complex traits. After detection of significant association, estimates of penetrance and allele-frequency parameters for the associated variant indicate the importance of that variant and facilitate the planning of replication studies. However, when these estimates are based on the original data used to detect the variant, the results are affected by an ascertainment bias known as the "winner's curse." The actual genetic effect is typically smaller than its estimate. This overestimation of the genetic effect may cause replication studies to fail because the necessary sample size is underestimated. Here, we present an approach that corrects for the ascertainment bias and generates an estimate of the frequency of a variant and its penetrance parameters. The method produces a point estimate and confidence region for the parameter estimates. We study the performance of this method using simulated data sets and show that it is possible to greatly reduce the bias in the parameter estimates, even when the original association study had low power. The uncertainty of the estimate decreases with increasing sample size, independent of the power of the original test for association. Finally, we show that application of the method to case-control data can improve the design of replication studies considerably.

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

Wilcoxon rank-sum test

Rank the X's and Y's from smallest to largest (1, 2, ..., n+m)

R = sum of ranks for X's

rank

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

46.8

49.7

51.9

74.1

75.1

84.5

90.0

95.1

101.5

Х

35.0

38.2

43.3

50.0

57.1

61.2

(Also known as the Mann-Whitney Test)

 $\mathsf{R} = 1 + 2 + 3 + 6 + 8 + 9 = 29$

P-value = 0.026

 \rightarrow use wilcox.test()

Note: The distribution of R (given that X's and Y's have the same dist'n) is calculated numerically

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017



This leaves the relationship between the features (such as genes) unchanged.

[140.688]



Another reason to like balanced designs

\downarrow	σ_1^2/σ_2^2						
n_{1}/n_{2}	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{2}$	1	2	4	8
$\frac{1}{2}$ 1 2 4 8	0.011 0.050 0.133 0.237 0.331	0.016 0.050 0.110 0.179 0.237	0.028 0.050 0.080 0.110 0.133	0.050 0.050 0.050 0.050 0.050	0.080 0.050 0.028 0.016 0.011	0.110 0.050 0.016 0.004 0.001	0.133 0.050 0.011 0.001 0.000

-



Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

Paired data in ten observations. Between subject variability (SD) ten times larger than error.