# Special Topics in Biostatistics and Bioinformatics
# Week V
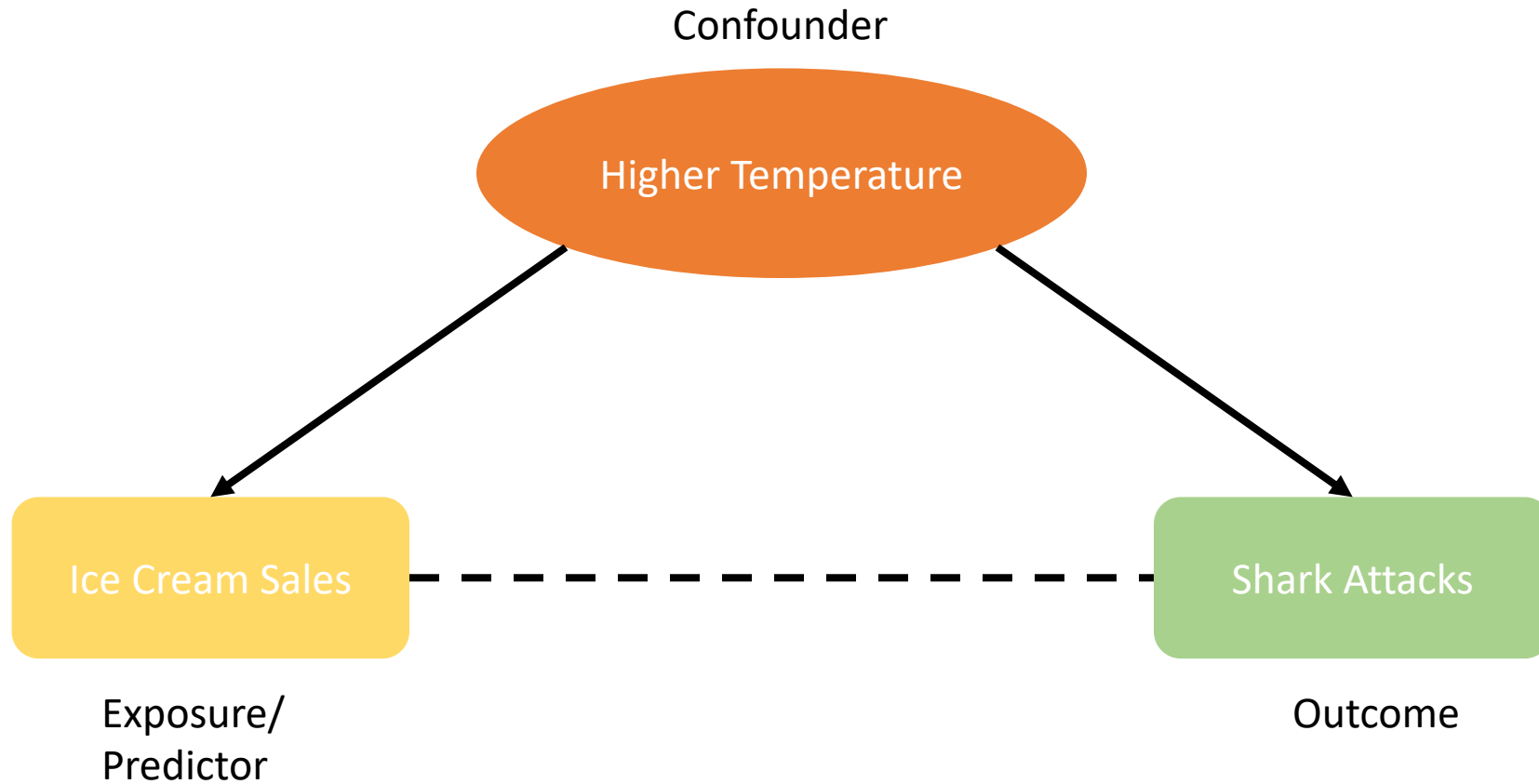
Ege Ülgen, M.D.

31 March 2022
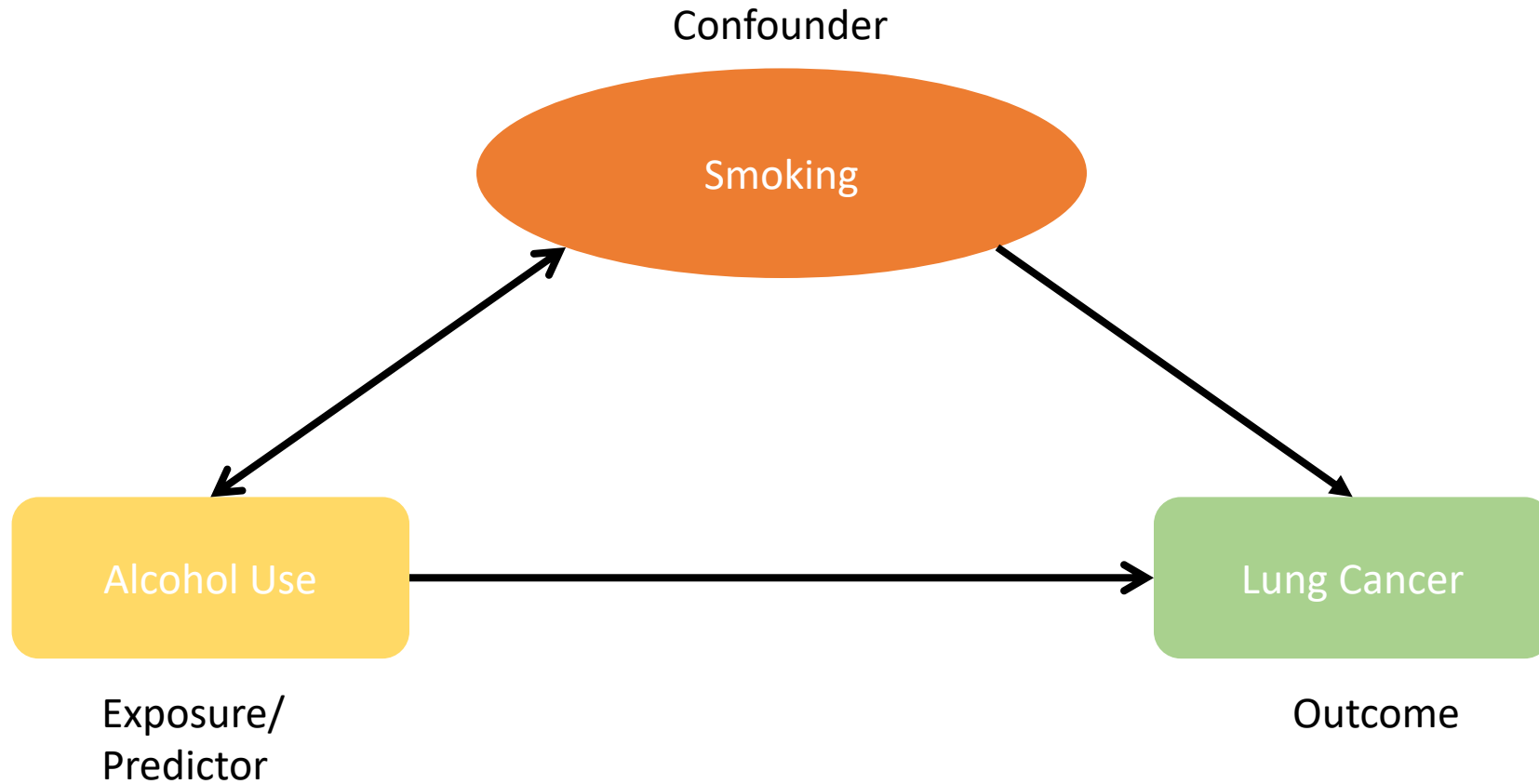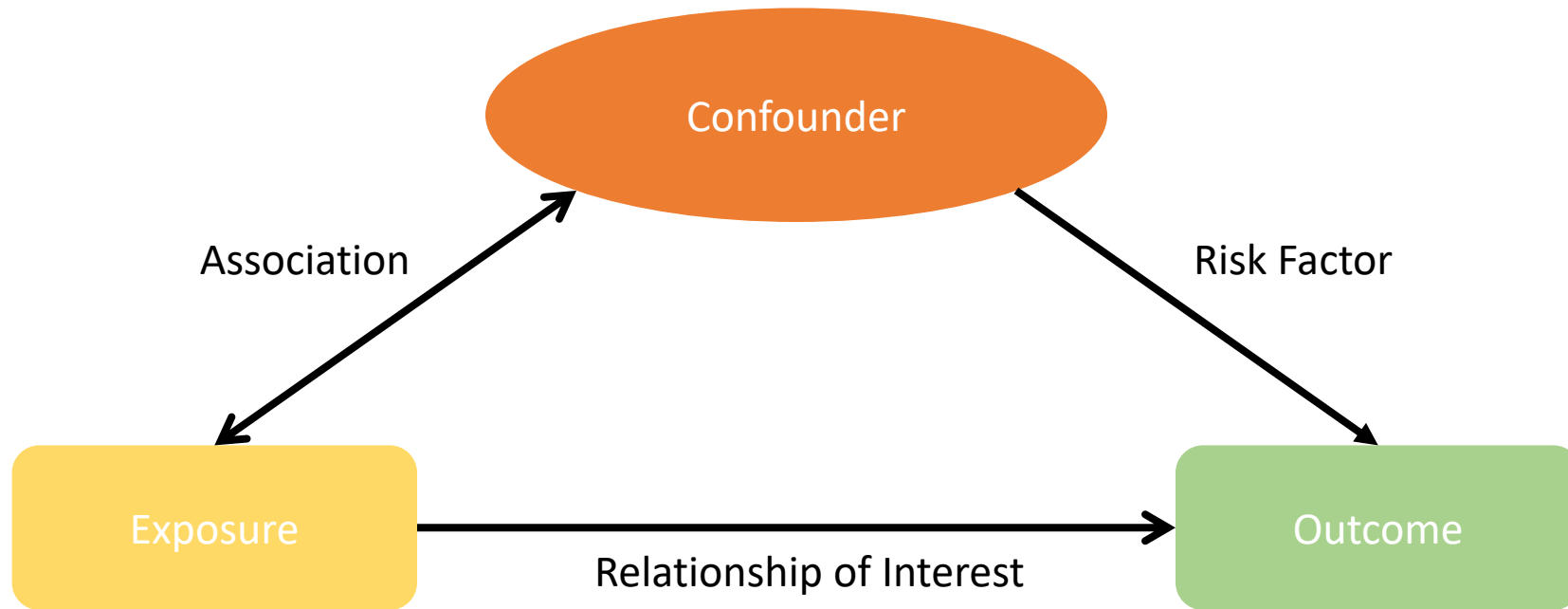
# What is confounding?

# What is confounding?

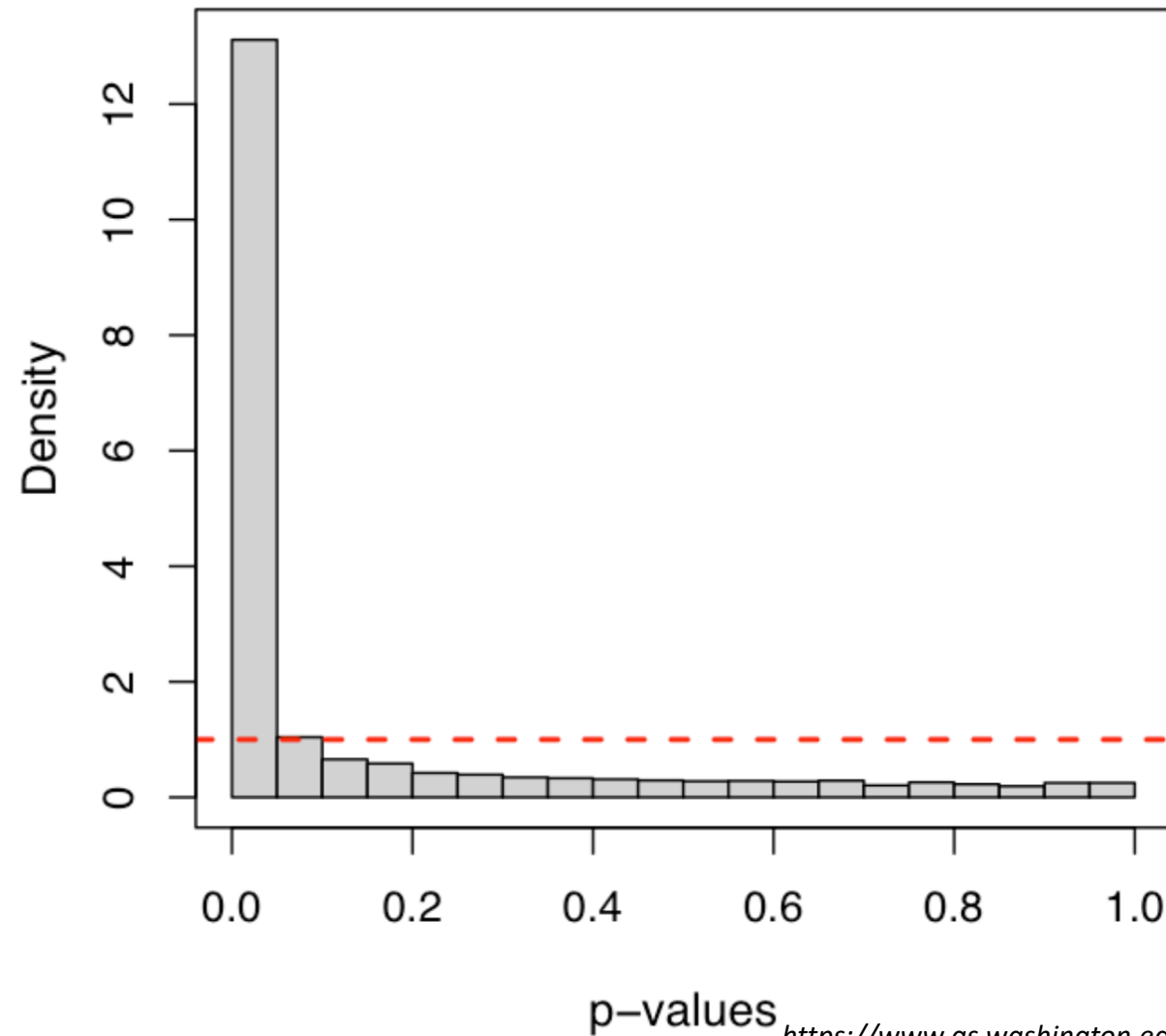# What is confounding?

# The Most Common Confounder

- Batch Effects

Common genetic variants account for differences in
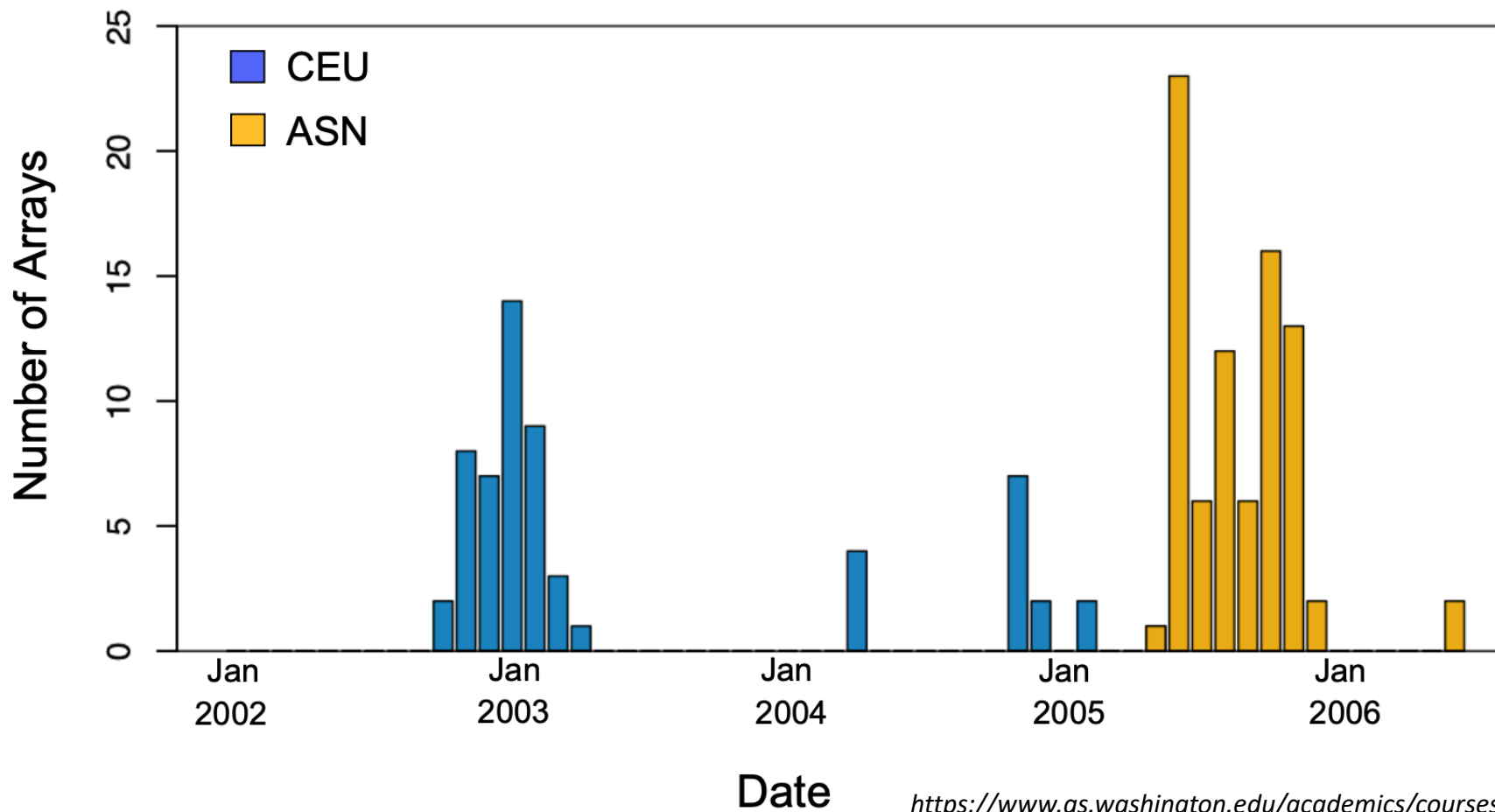gene expression among ethnic groups

Richard S Spielman[1], Laurel A Bastone[2], Joshua T Burdick[3], Michael Morley[3], Warren J Ewens[4] &
Vivian G Cheung[1,3,5]

- Compared gene expression levels between 60 CEU and 82 ASN HapMap individuals

- Tests of differential expression performed by **parametric t-tests** and adjustment for **multiple testing** through Sidak corrections

- Estimate ~26% of genes to be differentially expressed

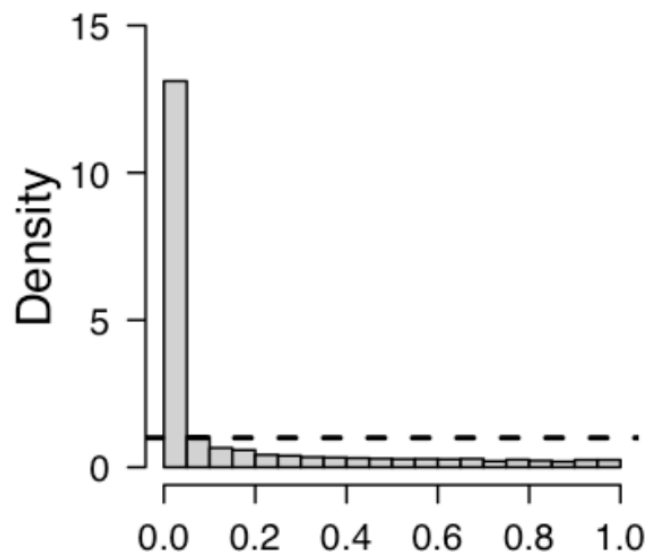# 78% of Genes Are Estimated To Be Differentially Expressed

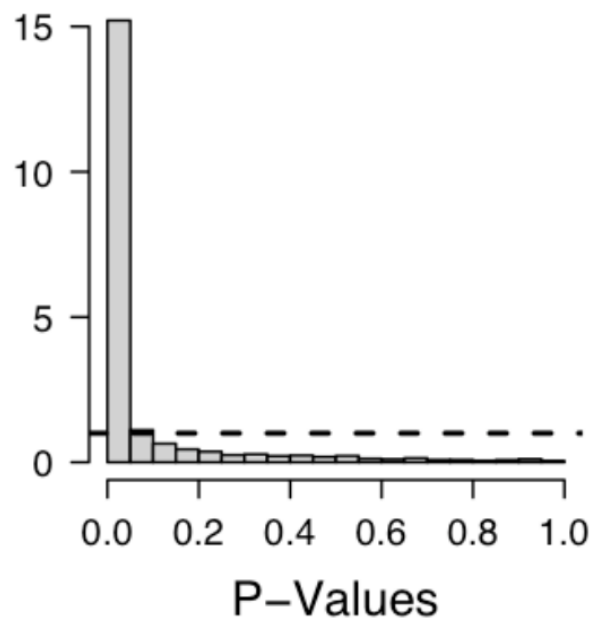# Population and Time of Processing Are Confounded

# Batch Effects Can Completely Account For Differential Expression
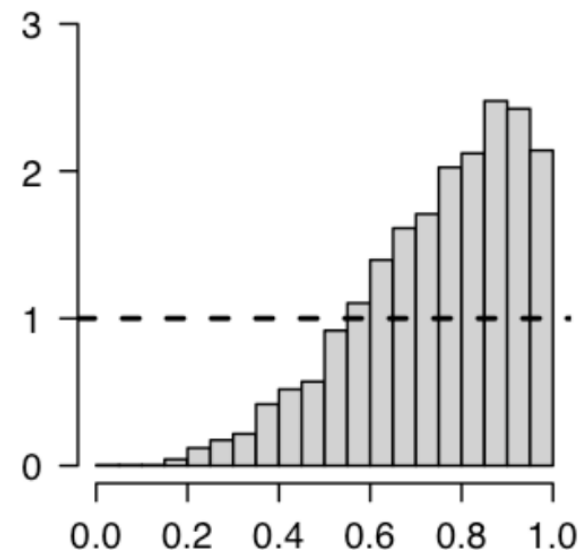
**Between Population**

**Between Years**

**Between Populations, Adjusting For Years**



**78%** of genes estimated to be differentially

**96%** of genes estimated to be differentially

**0%** of genes estimated to be differentially

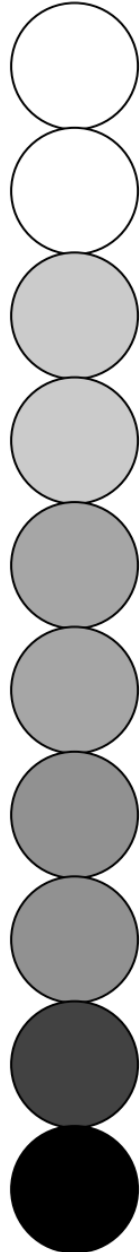Experimental units    Treatments

Confounding variable

Without randomization, the confounding variable differs among treatments
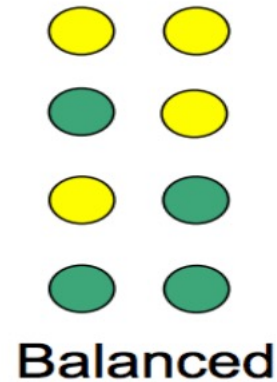
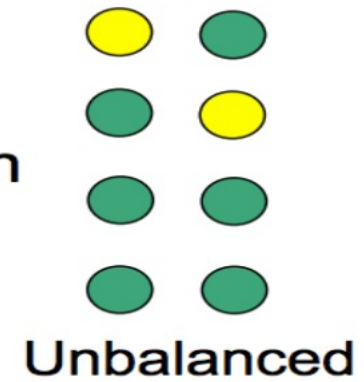Experimental units     Treatments

Confounding variable

With randomization, the confounding variable does not differ among treatments

# More good study design characteristics

- Balanced
- Replicated
- Has controls

Balanced

Better than

Unbalanced

# Batch Effects
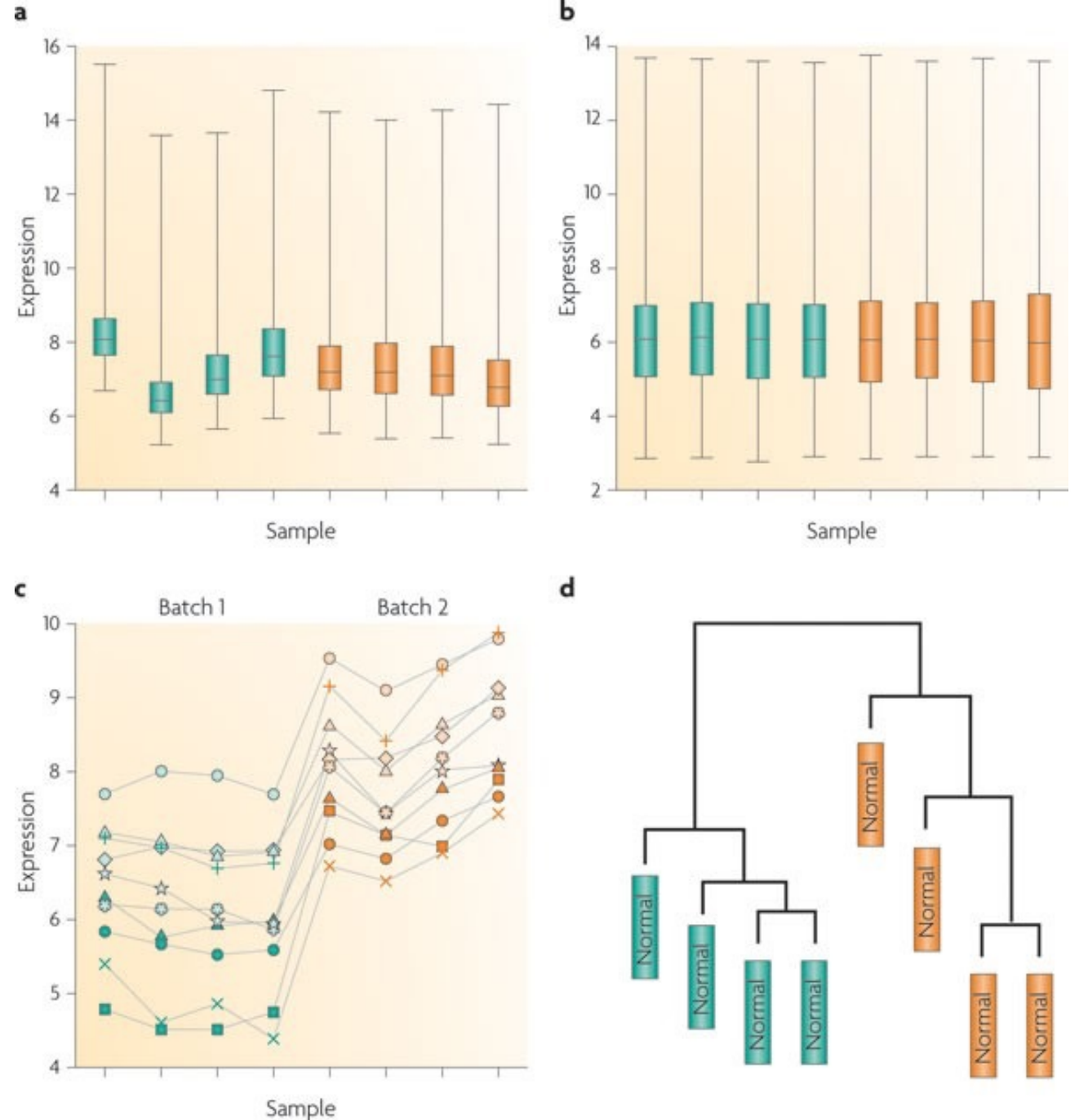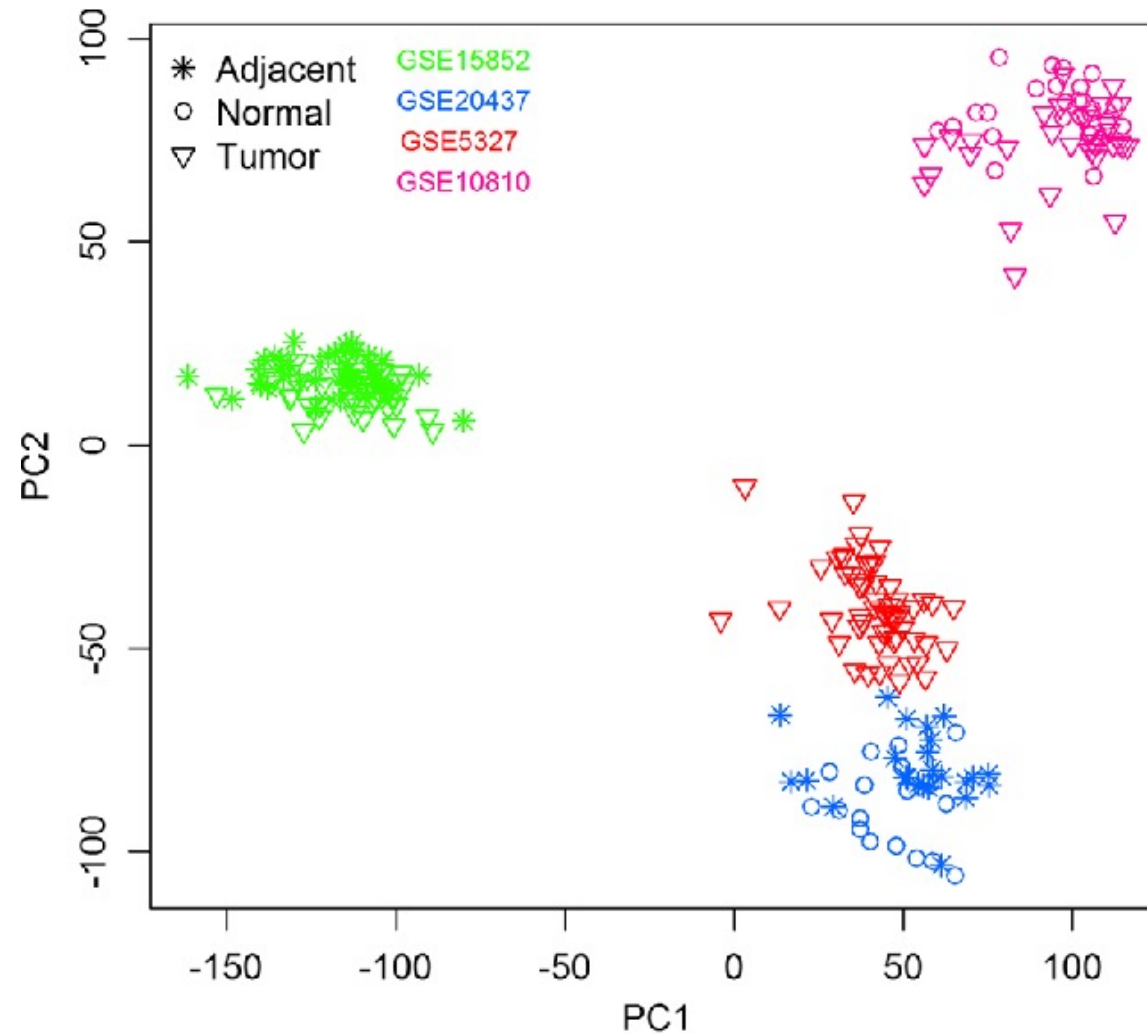
- Batch effects are **unwanted sources of variation** caused by different processing date, handling personnel, reagent lots, equipment/machines, etc.

- A common method for visualizing the existence of batch effects is PCA

- The first two principal components are plotted with each sample colored by the suspected batch, and **separation of colors** is taken as evidence of a batch effect

Green and orange represent two different processing dates

a. Box plot of raw gene expression data (log base 2)
b. Box plot of data processed with quantile normalization
c. Example of ten genes that are susceptible to batch effects even after normalization (Hundreds of genes show similar behavior but, for clarity, are not shown)
d. Clustering of samples after normalization. Note that the **samples perfectly cluster by processing date**

Leek JT, Scharpf RB, Bravo HC, Simcha D, Langmead B, Johnson WE, et al. Tackling the widespread and critical impact of batch effects in high-throughput data. Nature Reviews Genetics. 2010 Oct;11(10):733–9.



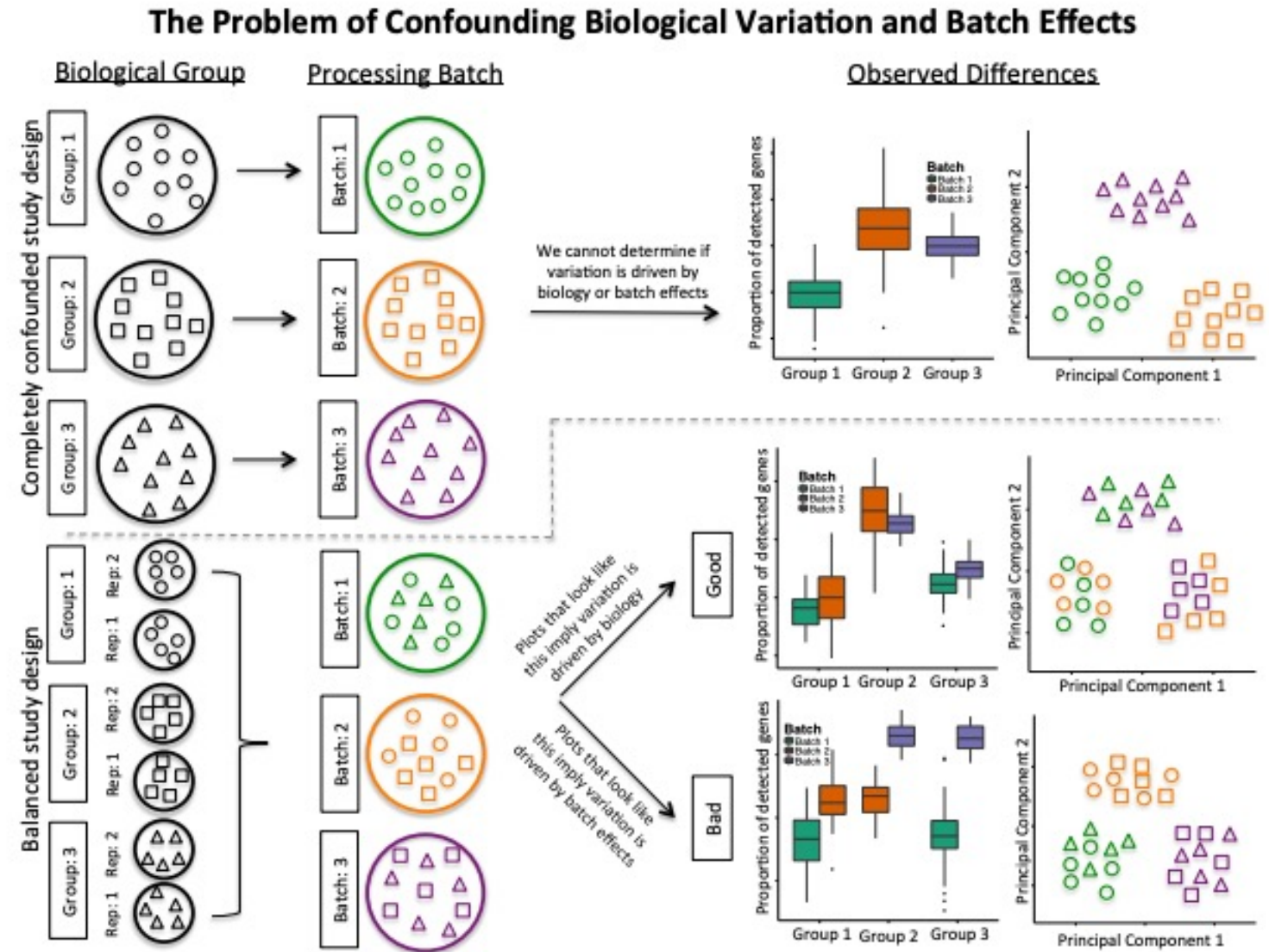Nature Reviews | Genetics

# Batch Effect Example

# Sources of "Batch" Effects

- External Factors (e.g., environment)
- Genetics/Epigenetics
- Technical Factors

# When can you remove batch effects?

- When they don't perfectly overlap with what you care about



**The Problem of Confounding Biological Variation and Batch Effects**
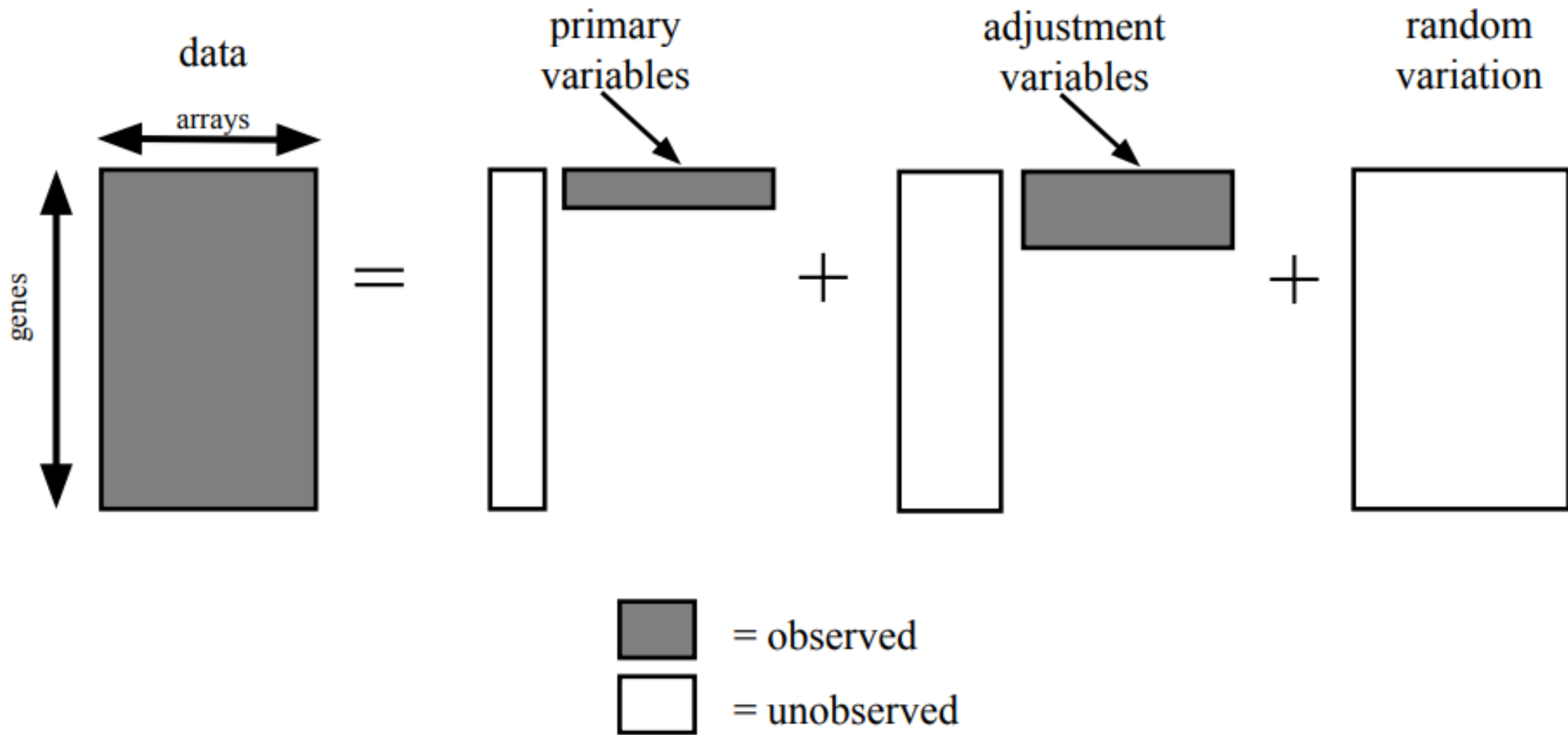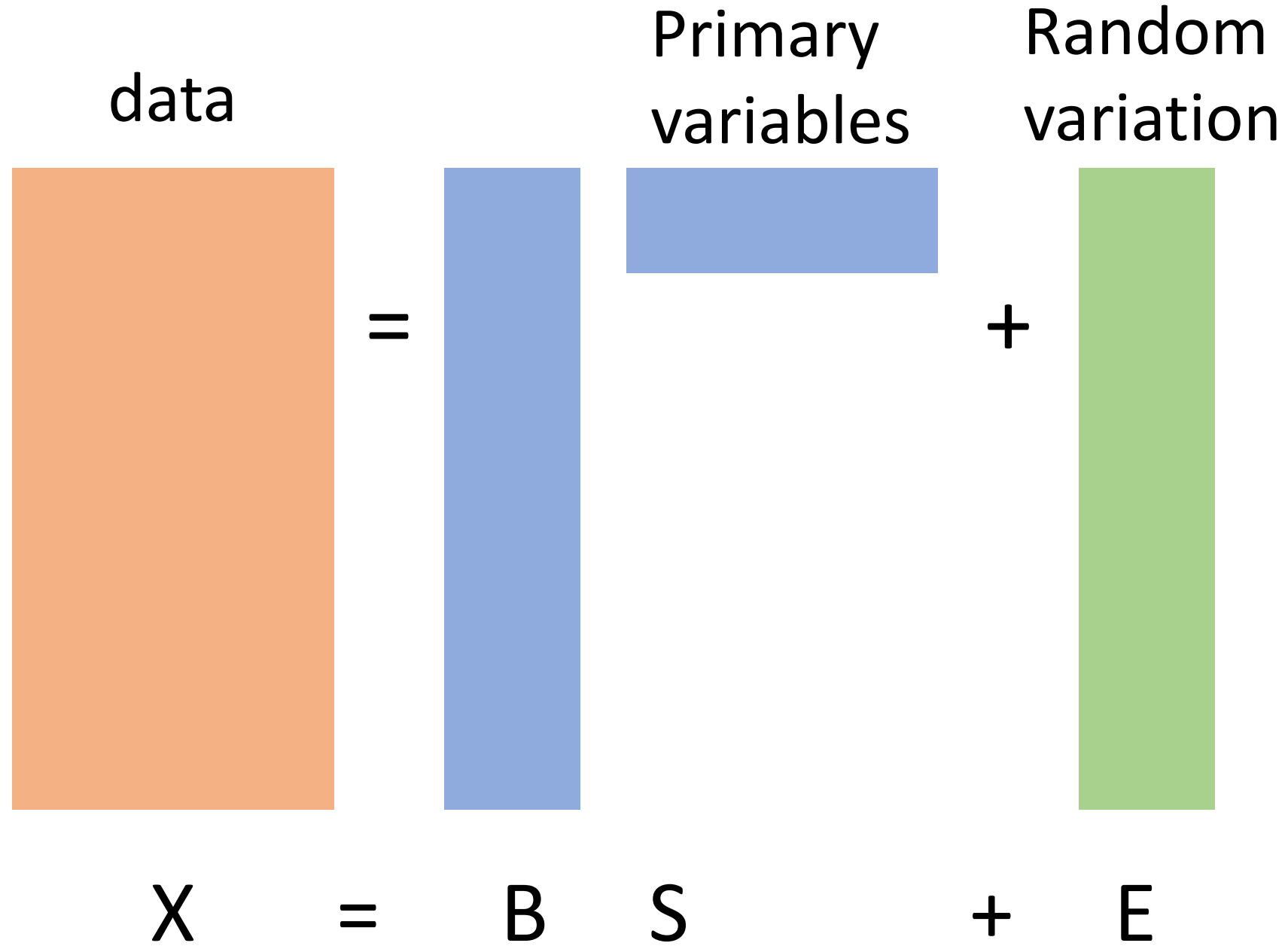
# When batch is known

$$Y = \beta_0 + \beta_1 P + \beta_2 B + \varepsilon$$

P: phenotype you're investigating

B: Batch

# When batch is unknown

data

Primary variables

Random variation

X  =  B  S  +  E

data   =   Primary variables   Dependent variation   Independent variation

$$X = BS + H + U$$

data    Primary
variables    Estimated
batch    Independent
variation

$$X = BS + \Gamma G + U$$

# Surrogate Variable Analysis

- https://www.pnas.org/content/105/48/18718.full
- https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.0030161