

# CATEGORICAL DATA AND SAMPLE SIZE

Meena Choi, Ting Huang, Olga Vitek

College of Science

College of Computer and Information Science

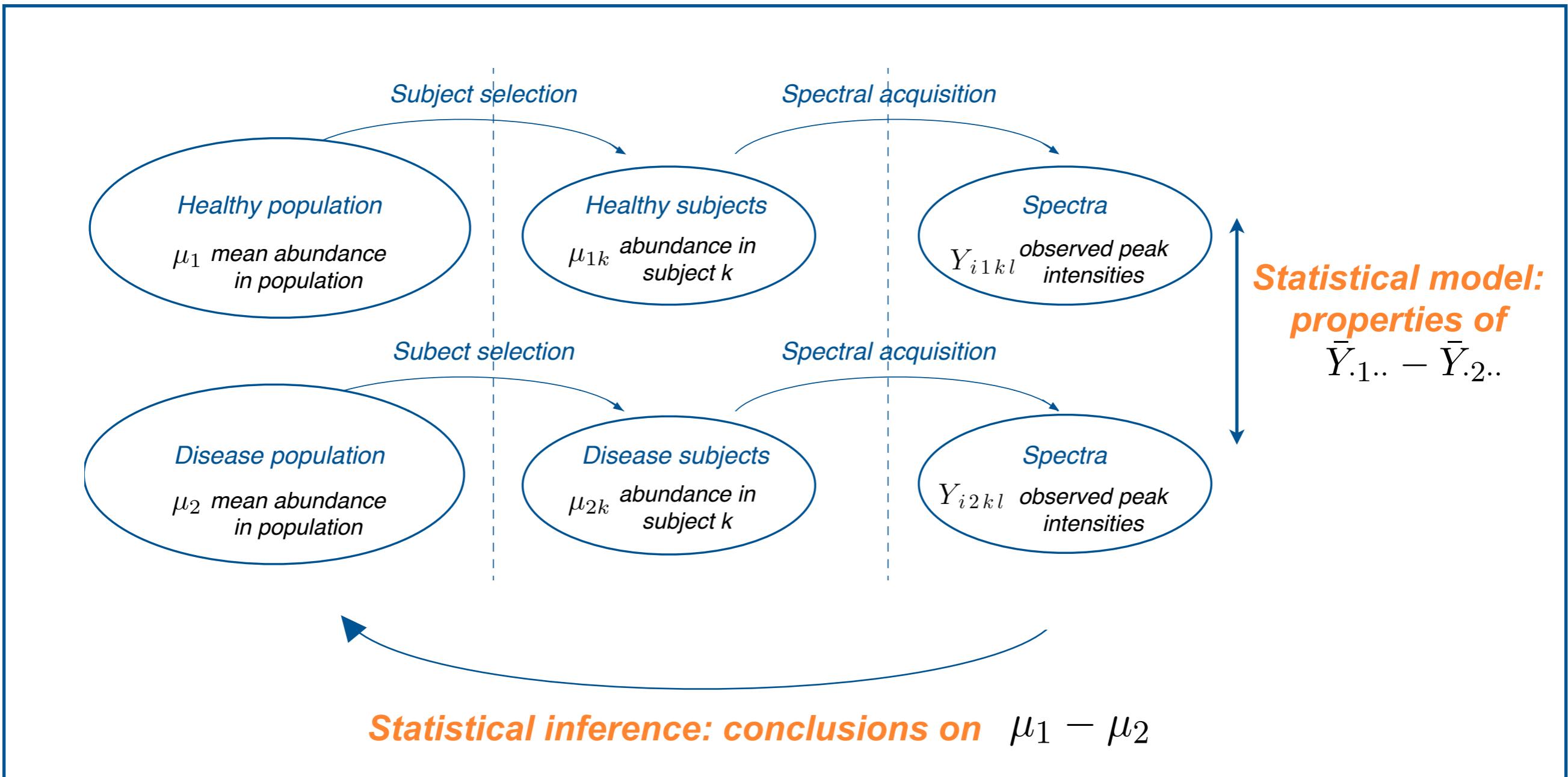


Northeastern University

# OUTLINE

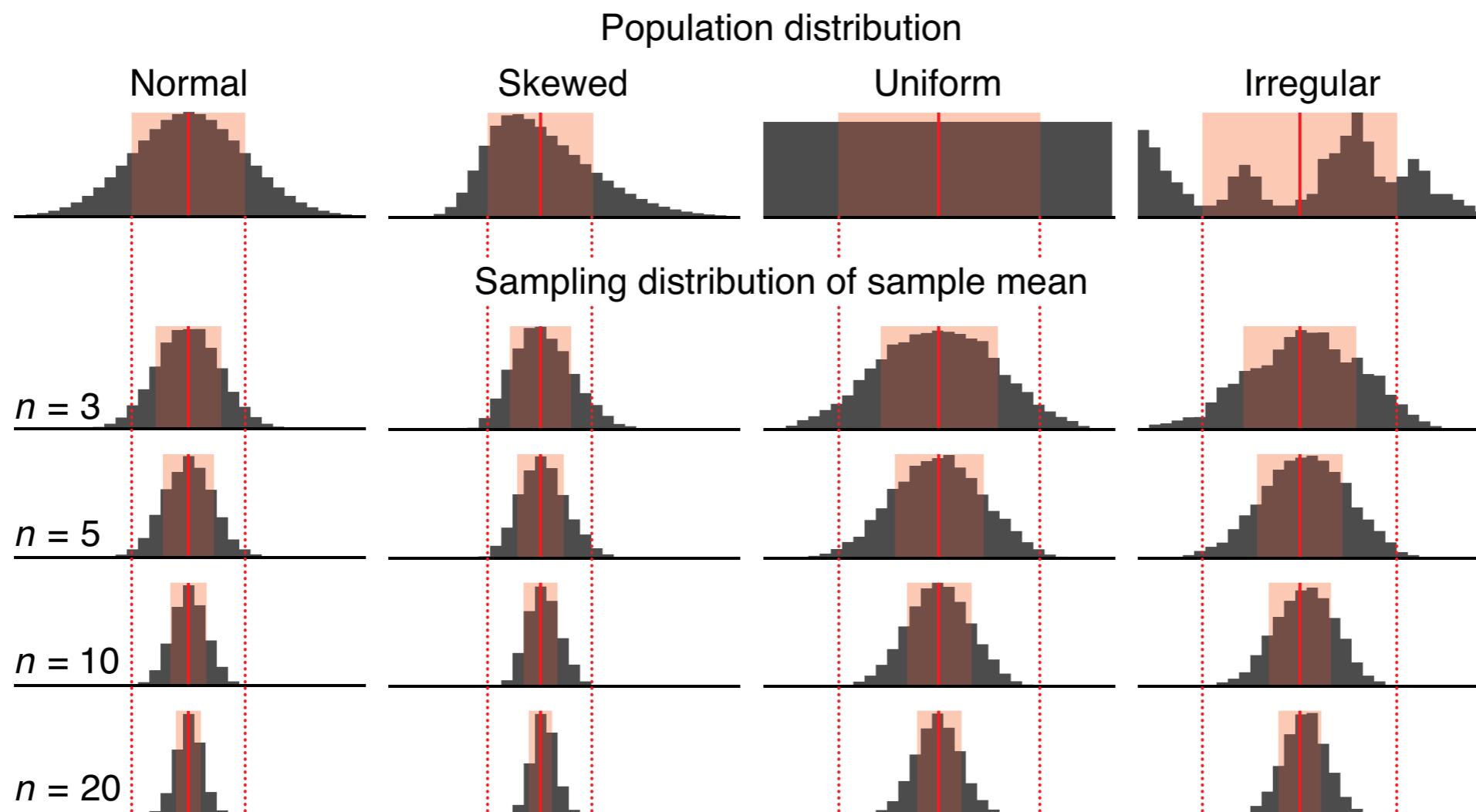
- Statistical analysis of count data
  - Comparison of proportions
  - 2-way tables: chi2 test, Fisher's exact test
- So how many replicates do I need?
  - Design of complex experiments

# RECALL THE GOAL OF GROUP COMPARISON



*What if measurements are binary?  
(E.g., up or down?)*

# RECALL THE CENTRAL LIMIT THEOREM



*Probability  
distribution  
of the data*

*Repeatedly  
selecting  $n$   
data points  
and  
calculating  
means*

*Averages of binary variables will approach  
Normal distribution as sample size increases*

# ANALYZING PROPORTIONS

## Hypothesis testing and confidence interval

*Population proportions*

$H_0: \text{'status quo', no change in abundance, } \pi_1 - \pi_2 = 0$

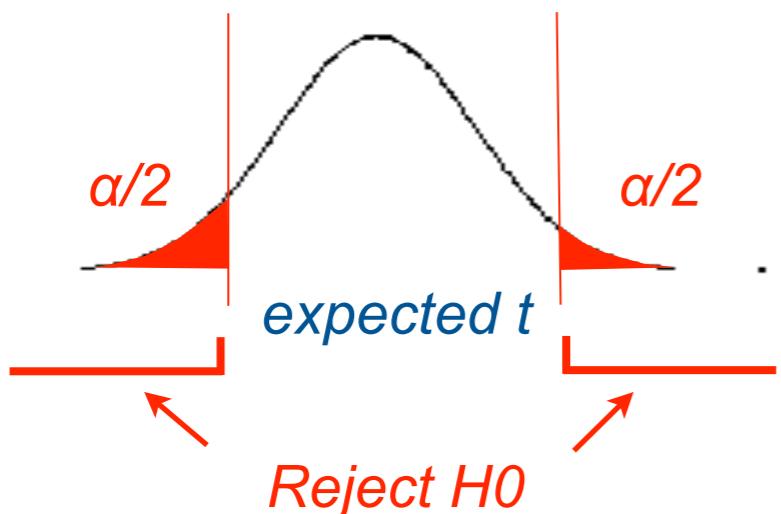
$H_a: \text{change in abundance, } \pi_1 - \pi_2 \neq 0$

observed  $t = \frac{\text{difference of group means}}{\text{estimate of variation}} = \frac{p_1 - p_2}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$

*Proportions in each group*

Distribution of the score if  $H_0$  is true  
 $\sim \text{Normal}(0, 1)$

$\alpha = \text{False Positive Rate}$



Confidence interval for difference of proportions

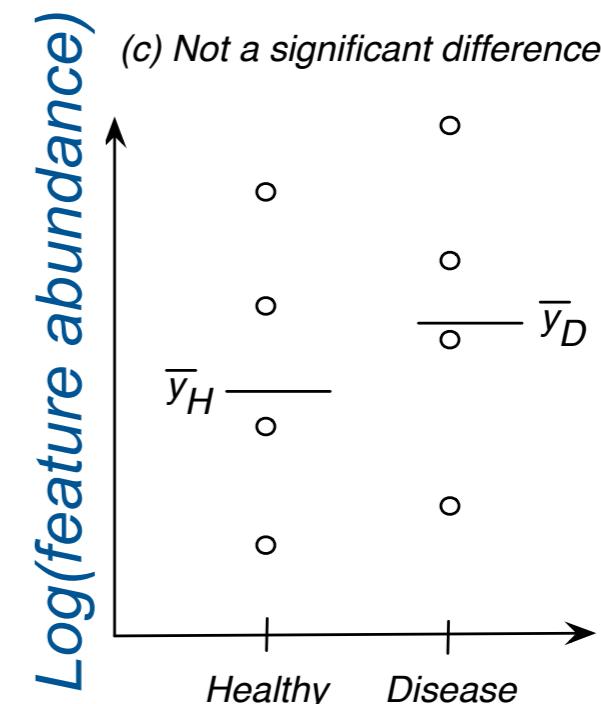
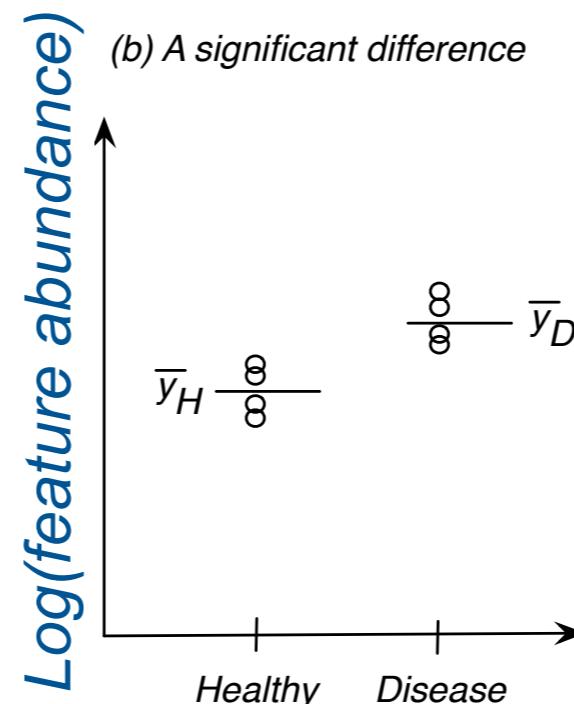
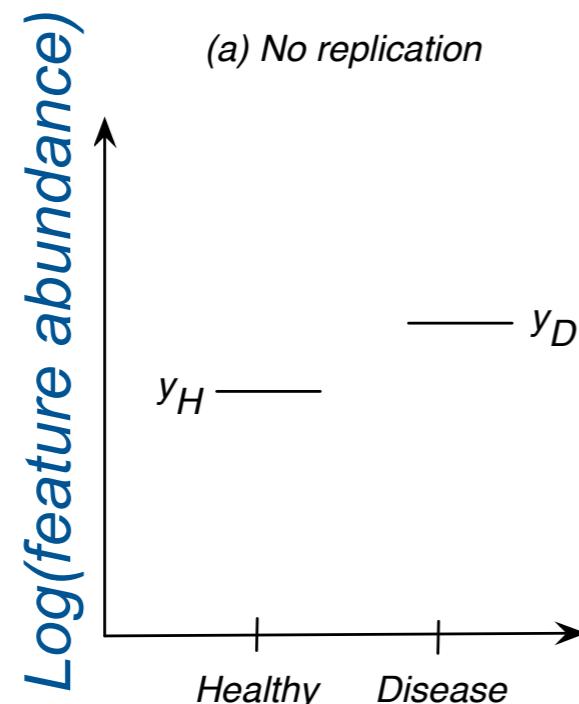
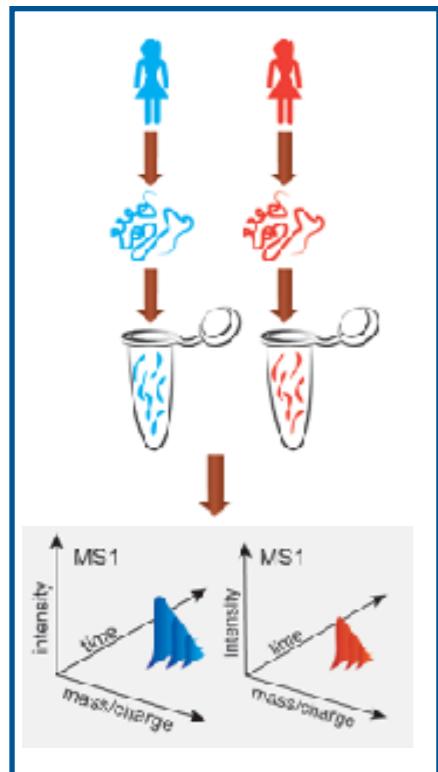
$$p_1 - p_2 \pm z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

# OUTLINE

- Statistical analysis of count data
  - Comparison of proportions
  - 2-way tables: chi<sup>2</sup> test, Fisher's exact test
- So how many replicates do I need?
  - Design of complex experiments

# PRINCIPLE I: REPLICATION

(1) carries out the inference and (2) minimizes inefficiencies

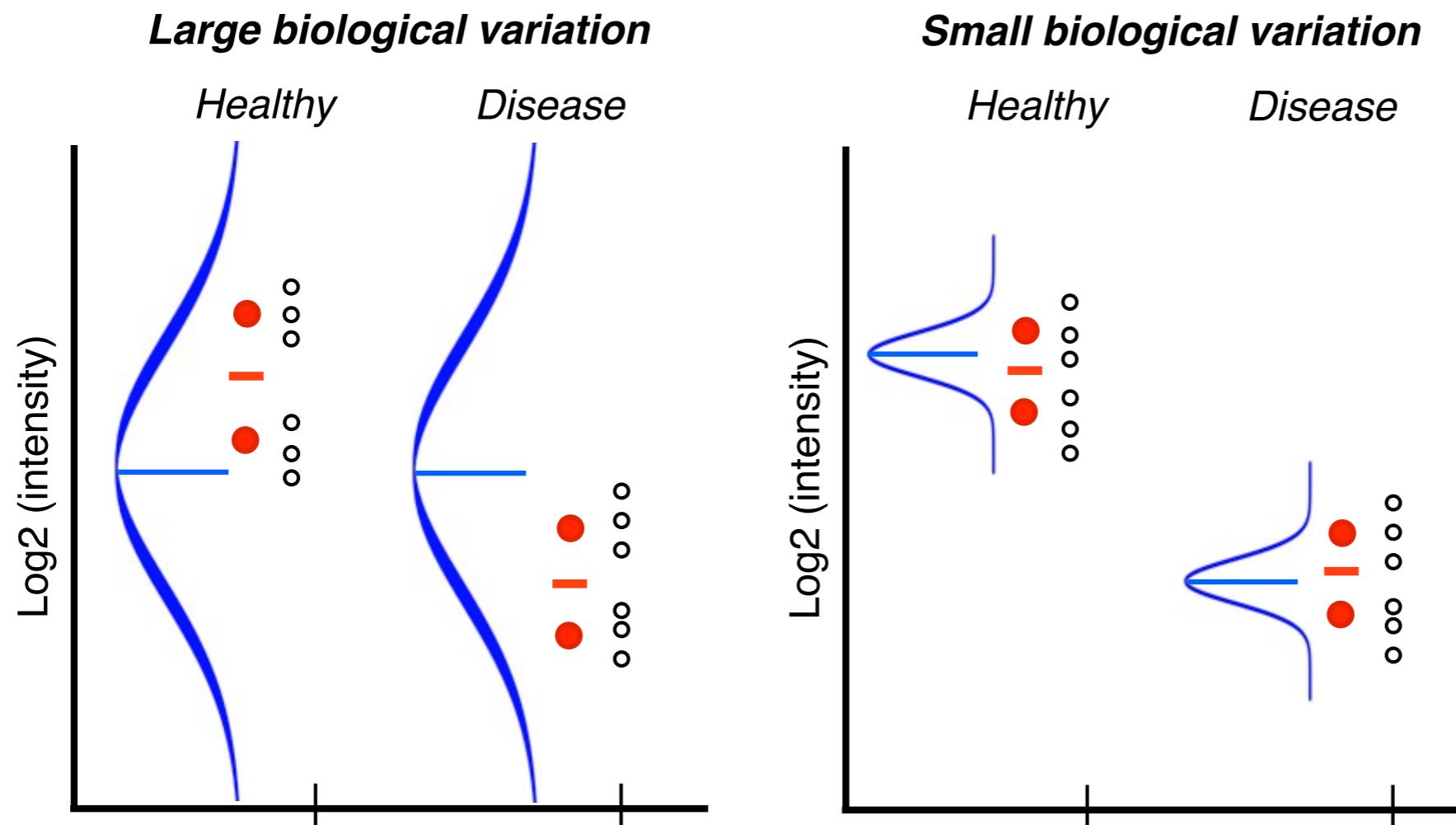


Two levels of randomness imply two types of replication:

- ◆ *Biological replicates*: selecting multiple subjects from the population
- ◆ *Technical replicates*: multiple runs per subject

# MULTI-LAYER DESIGN AND ANALYSIS

Multiple sources are responsible for variation in measurements



*When biological variance is large, more biological replicates are needed to accurately estimate the variance*

# BIOLOGICAL REPLICATION IS MOST IMPORTANT

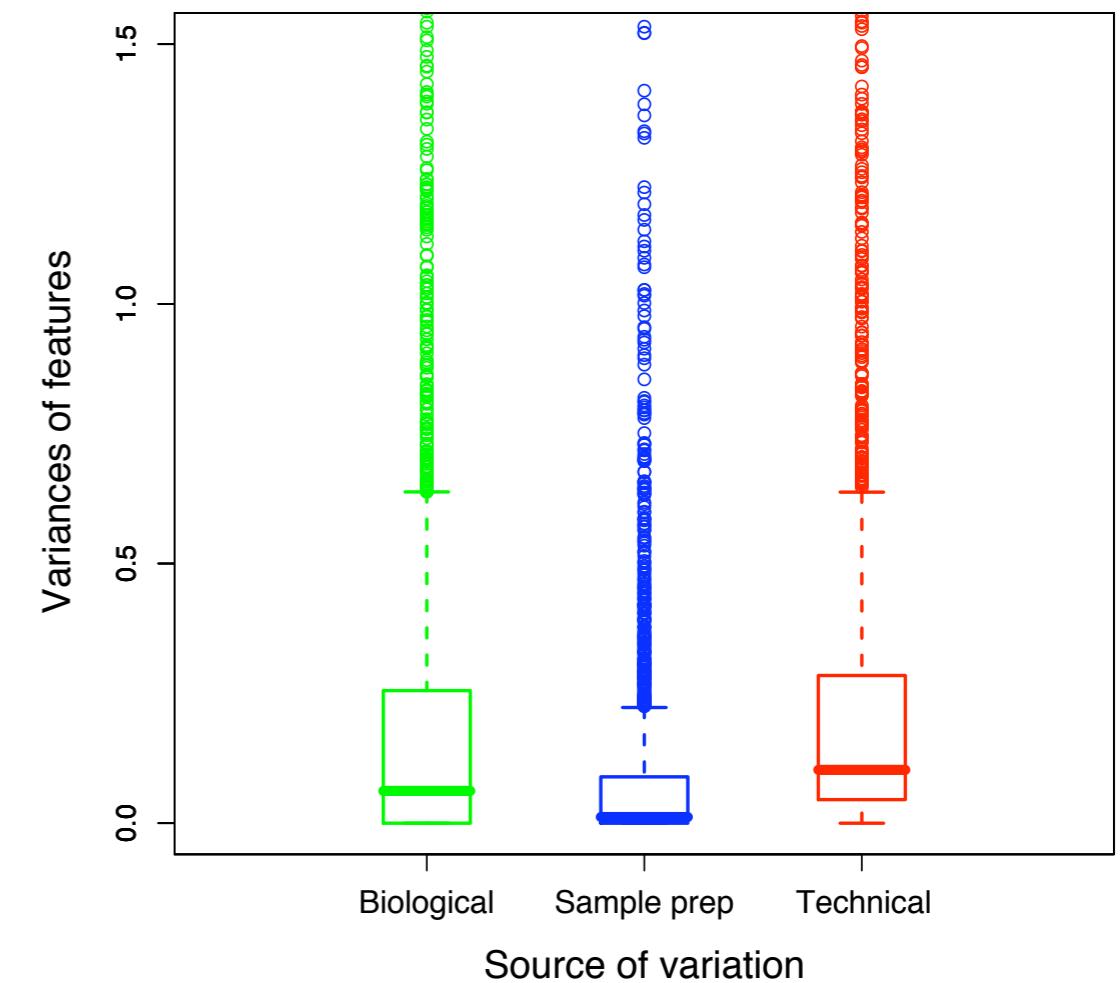
<b>Observed feature intensity</b> $y_{ijkl}$	<b>Systematic mean signal of disease group</b> <b>Group mean<sub>i</sub></b>	<b>Random deviation due to individual</b> <b>Indiv(Group)<sub>j(i)</sub></b> $\sim N(0, \sigma_{\text{Indiv}}^2)$	<b>Random deviation due to sample preparation</b> <b>Prep(Indiv)<sub>k(ij)</sub></b> $\sim N(0, \sigma_{\text{Prep}}^2)$	<b>Random deviation due to measurement error</b> <b>Error<sub>l(ijk)</sub></b> $\sim N(0, \sigma_{\text{Error}}^2)$
---------------------------------------------------------	-----------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------

$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Indiv}}^2}{I} + \frac{\sigma_{\text{Prep}}^2}{IJ} + \frac{\sigma_{\text{Error}}^2}{IJK} \right)$$

**I:** # individuals per disease group  
**J:** # sample preps  
**K:** # replicate runs

## A pilot experiment

- 2 healthy individuals, 2 with diabetes
- multiple sample preparations
- multiple LC-MS replicates



# BIOLOGICAL REPLICATION IS MOST IMPORTANT

<b>Observed feature intensity</b> $y_{ijkl}$	<b>Systematic mean signal of disease group</b> <b>Group mean<sub>i</sub></b>	<b>Random deviation due to individual</b> $\text{Indiv}(\text{Group})_{j(i)} \sim N(0, \sigma_{\text{Indiv}}^2)$	<b>Random deviation due to sample preparation</b> $\text{Prep}(\text{Indiv})_{k(ij)} \sim N(0, \sigma_{\text{Prep}}^2)$	<b>Random deviation due to measurement error</b> $\text{Error}_{l(ijk)} \sim N(0, \sigma_{\text{Error}}^2)$
---------------------------------------------------------	-----------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------

$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Indiv}}^2}{I} + \frac{\sigma_{\text{Prep}}^2}{IJ} + \frac{\sigma_{\text{Error}}^2}{IJK} \right)$$

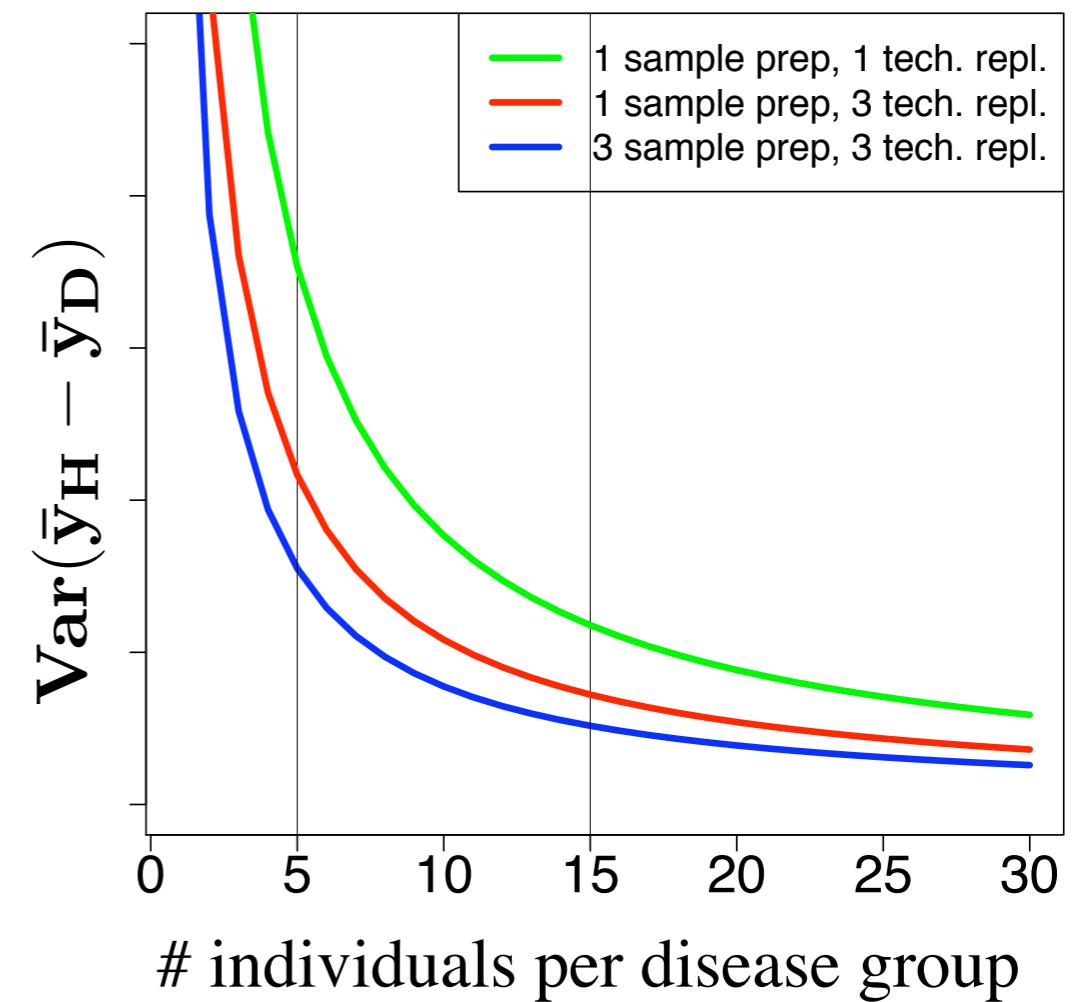
**I:** # individuals per disease group

**J:** # sample preps

**K:** # replicate runs

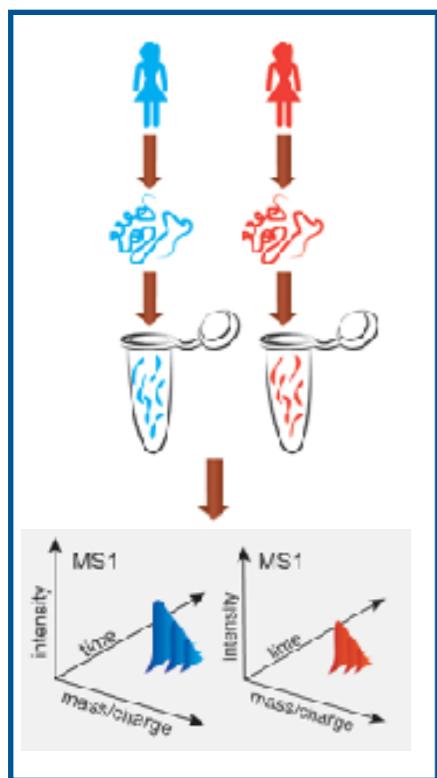
## Conclusion:

Maximize the number  
of biological replicates

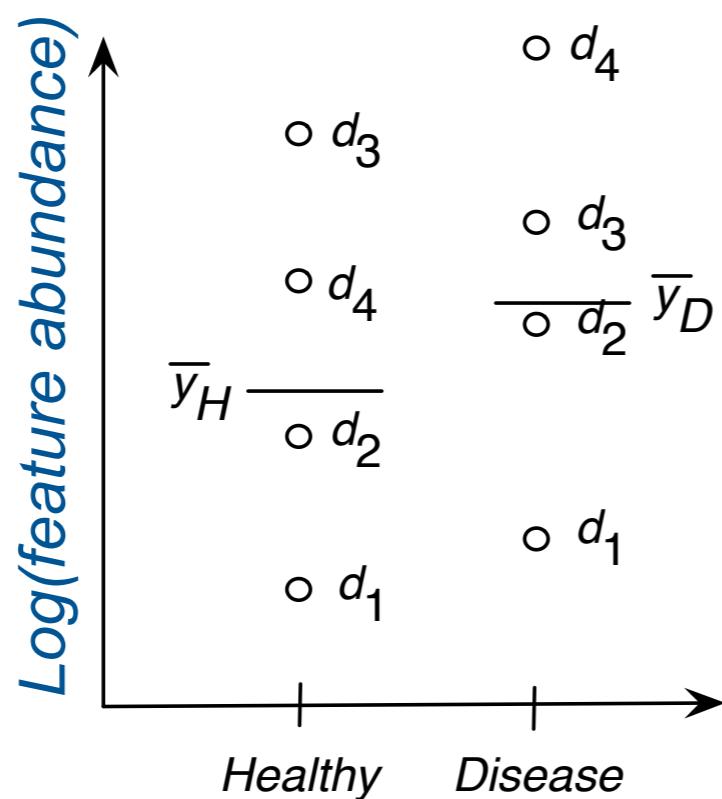


# PRINCIPLE 3: BLOCKING

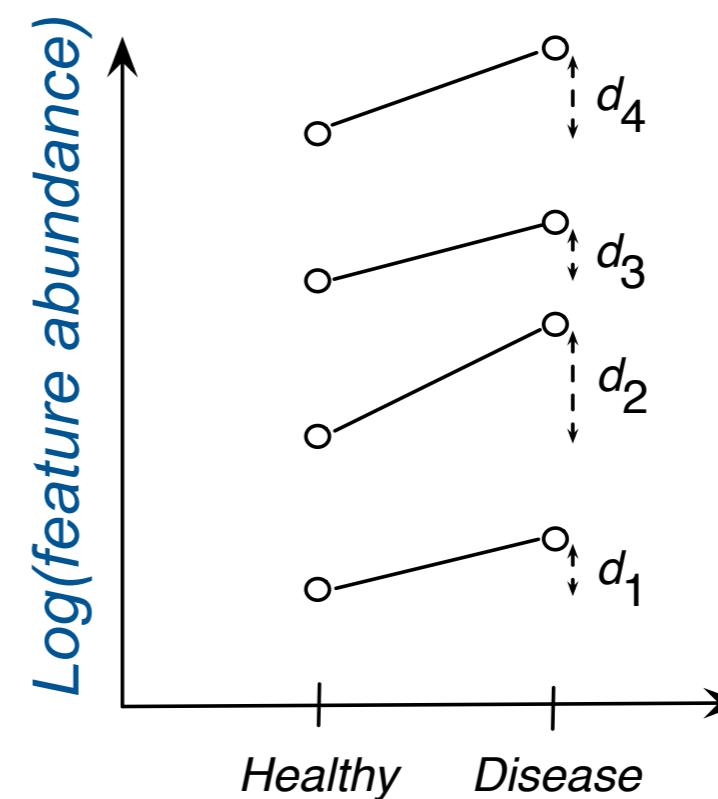
## Helps reduce both bias and inefficiency



(b) Complete randomization



(c) Day = block



Complete randomization  
= inflated variance

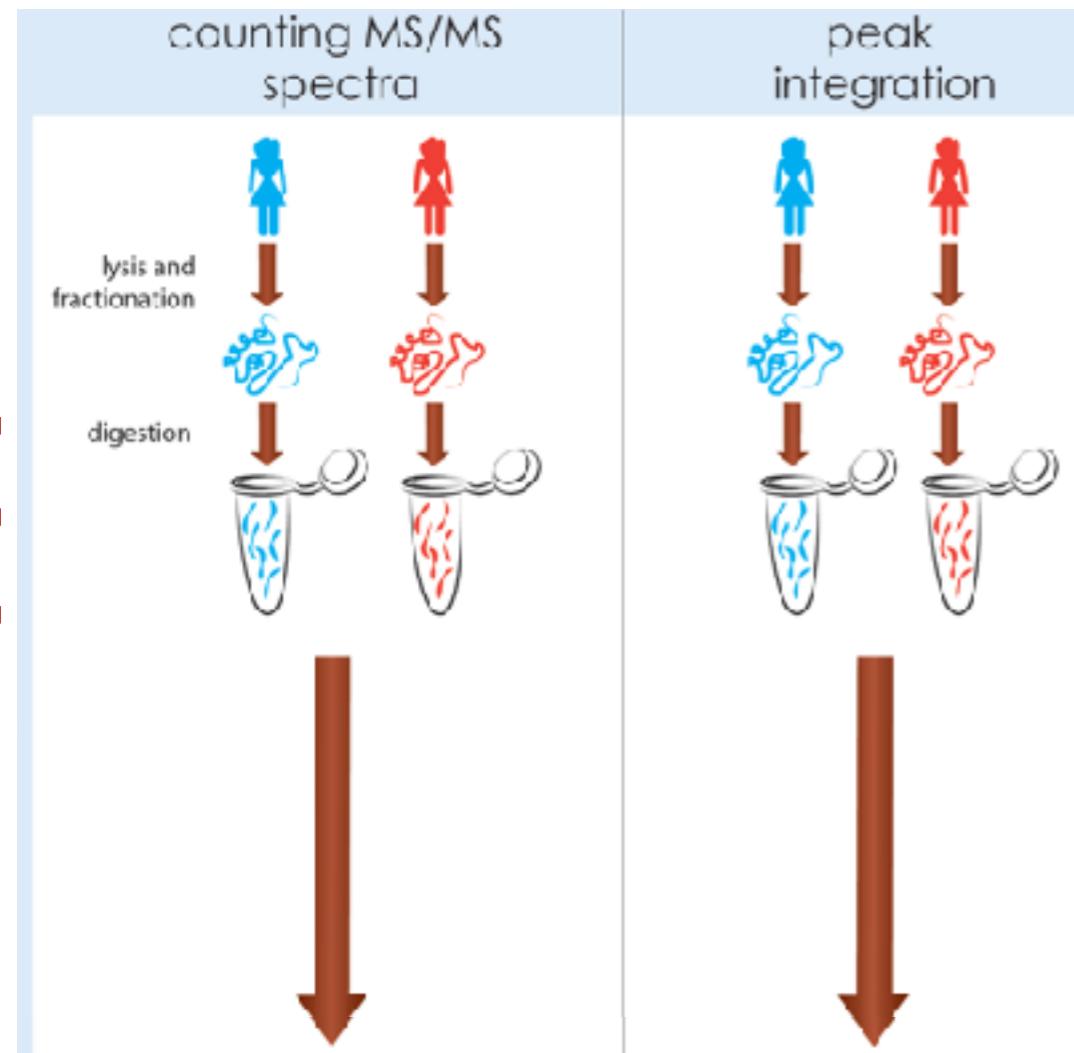
Block-randomization  
= restriction on randomization  
= systematic allocation

Two levels of randomness imply two types of blocks:

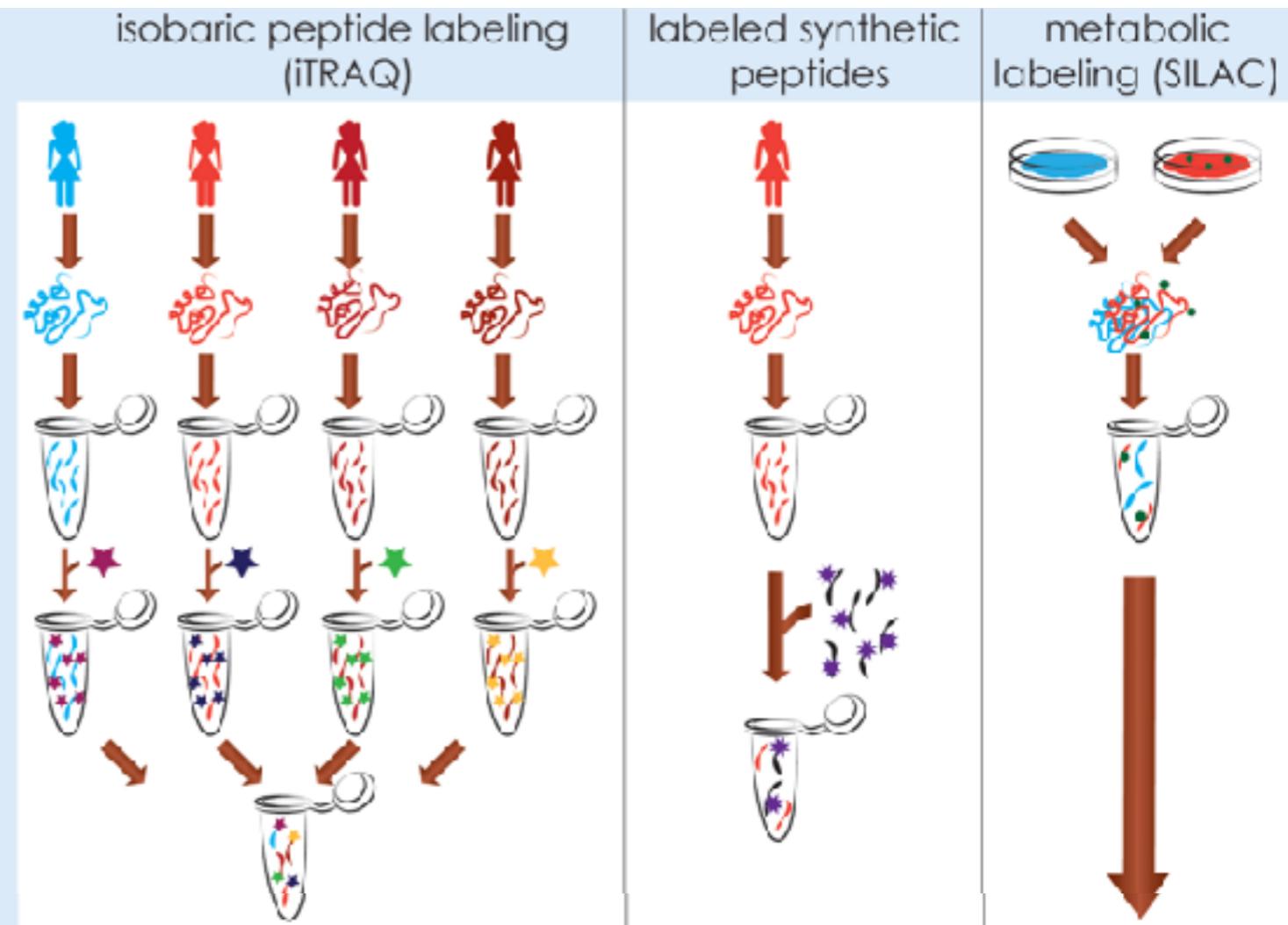
- ◆ *Biological replicates*: subjects having similar characteristics (e.g. age)
- ◆ *Technical replicates*: samples processed together (e.g. in a same day)

# LABELED QUANTITATIVE WORKFLOWS

## *Label-free*



## *Label-based*



# BLOCKING IS HELPFUL

## Especially when between-block variance is large

Observed feature intensity	=	Systematic mean signal of disease group	+	Random deviation due to block (e.g. plate or day)	+	Random deviation due to individual	+	Random deviation due to measurement error
$y_{ijkl}$	=	Group mean <sub>i</sub>	+	$\text{Block}_k \sim N(0, \sigma_{\text{Block}}^2)$	+	$\text{Indiv}(\text{Group})_{j(i)} \sim N(0, \sigma_{\text{Indiv}}^2)$	+	$\text{Error}_{l(ijk)} \sim N(0, \sigma_{\text{Error}}^2)$

A completely randomized design

I: # individuals per disease group

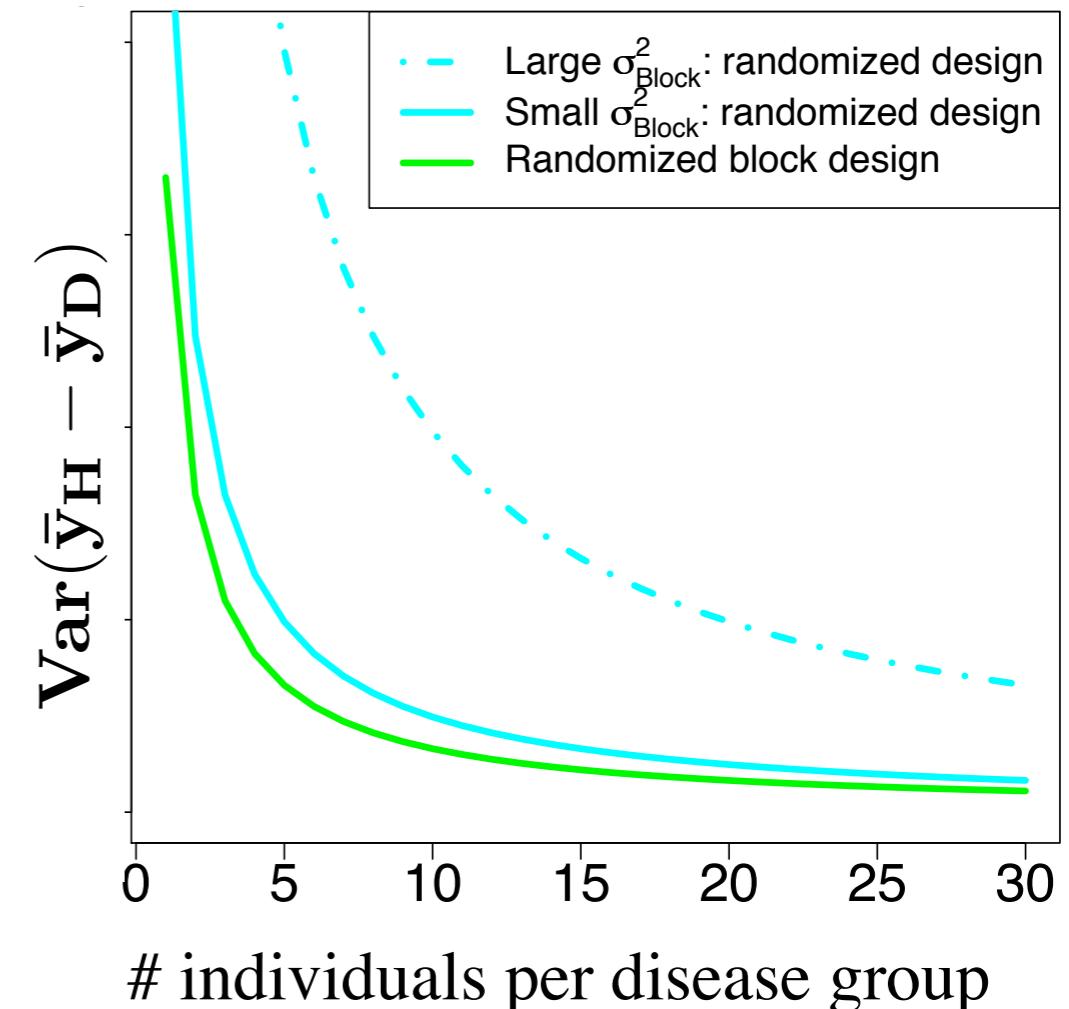
$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Block}}^2 + \sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$

A block-randomized design

$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$

**Conclusion:** Block-randomize

- if can not control a large source of variation
- if moderate sample size



# SO HOW MANY REPLICATES DO I NEED?

If we only have one feature:

Fix:  $\alpha$  - probability of a false positive discovery

$\beta$  - probability of a true positive discovery

$\Delta$  - anticipated fold change

$\sigma_{\text{Indiv}}^2$  and  $\sigma_{\text{Error}}^2$  - anticipated variability

Write:

$$\text{Var}(\bar{Y}_{1.} - \bar{Y}_{2.}) \leq \left( \frac{\Delta}{z_{1-\beta} + z_{1-\alpha/2}} \right)^2$$

where  $z_{1-\beta}$  and  $z_{1-\alpha/2}$  are quantiles of Normal distribution

$$\text{Var}(\bar{Y}_{1.} - \bar{Y}_{2.}) = s_1^2/n_1 + s_2^2/n_2 \leq \left( \frac{\Delta}{z_{1-\beta} + z_{1-\alpha/2}} \right)^2$$



solve for the number of individuals  $n_1$  and  $n_2$

# SO HOW MANY REPLICATES DO I NEED?

If we only have one feature:

Fix:  $\alpha$  - probability of a false positive discovery

$\beta$  - probability of a true positive discovery

$\Delta$  - anticipated fold change

$\sigma_{\text{Indiv}}^2$  and  $\sigma_{\text{Error}}^2$  - anticipated variability

Write:

$$\text{Var}(\bar{y}_H - \bar{y}_D) \leq \left( \frac{\Delta}{z_{1-\beta} + z_{1-\alpha/2}} \right)^2$$

where  $z_{1-\beta}$  and  $z_{1-\alpha/2}$  are quantiles of the Normal distribution

A completely randomized design

I: # individuals per disease group

$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Block}}^2 + \sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$

A block-randomized design

$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left( \frac{\sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$



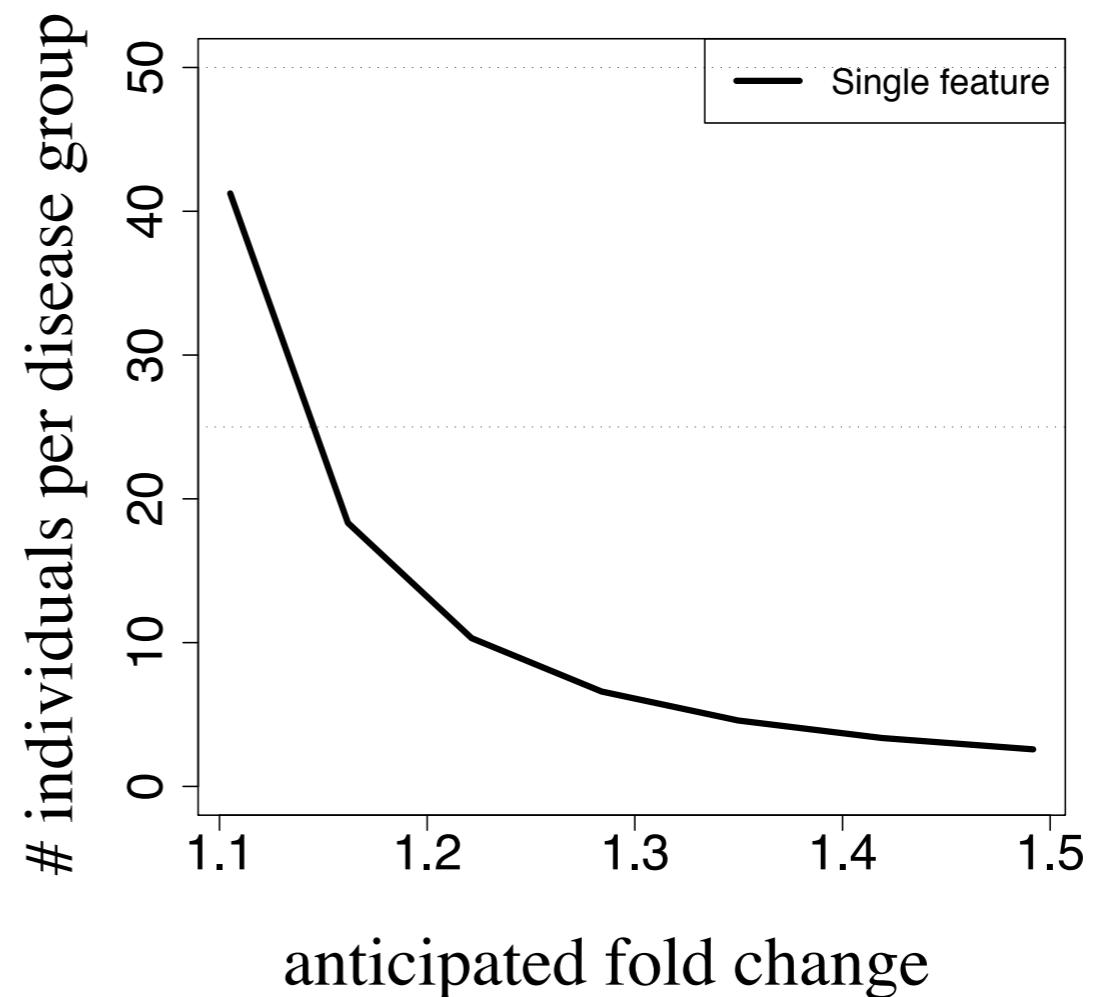
solve for the number of individuals I

# SO HOW MANY REPLICATES DO I NEED?

Example: pilot study with diabetes patients.

A block-randomized design

If we only had one feature:



## Conclusion:

The smaller the anticipated difference, the larger the sample size