

SAMPLE SIZE, ASSOCIATIONS AND CATEGORICAL DATA

Olga Vitek

College of Science
College of Computer and Information Science



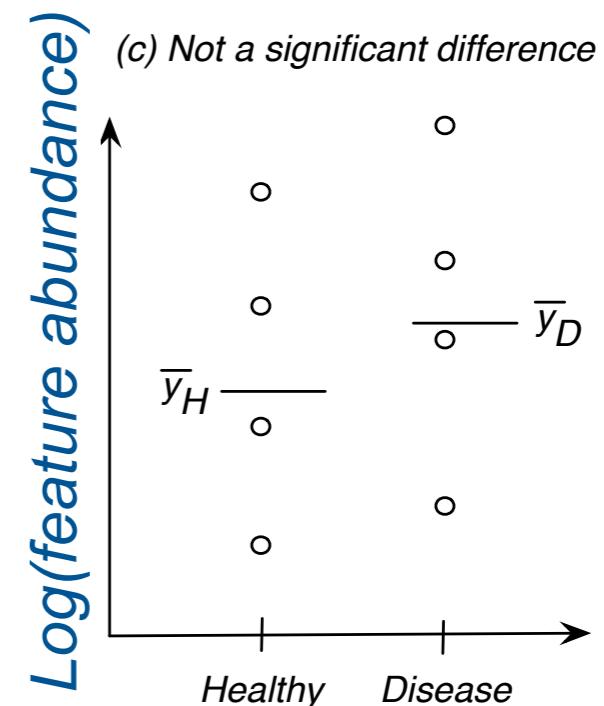
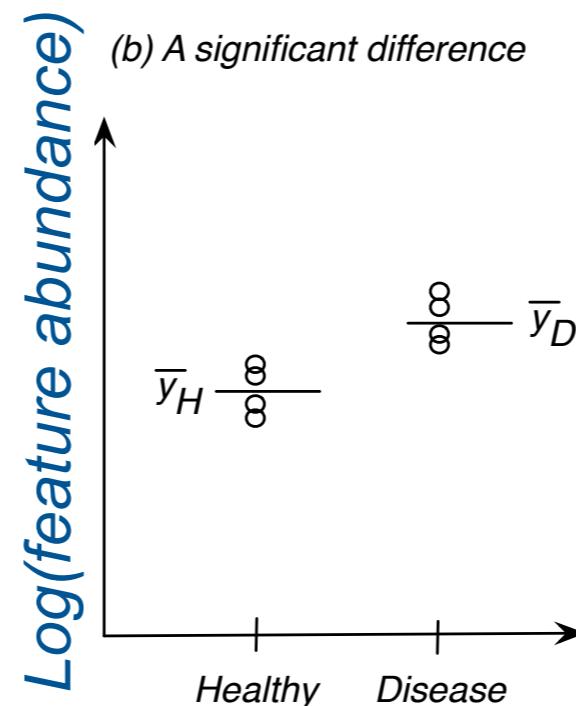
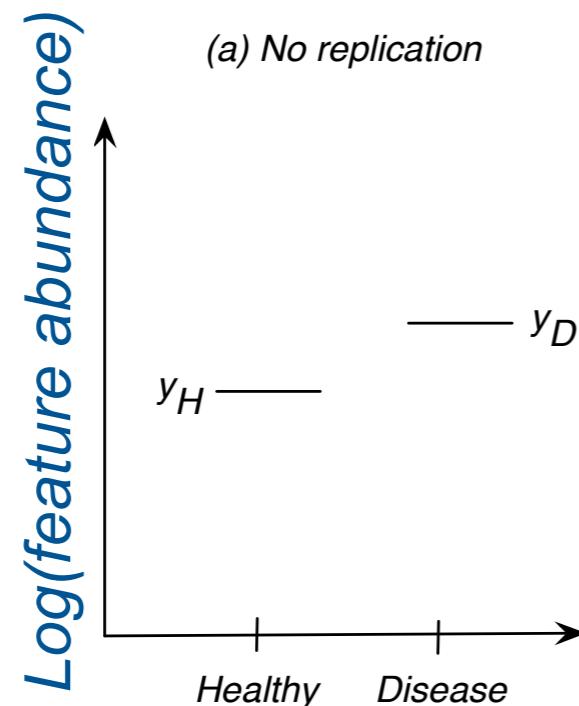
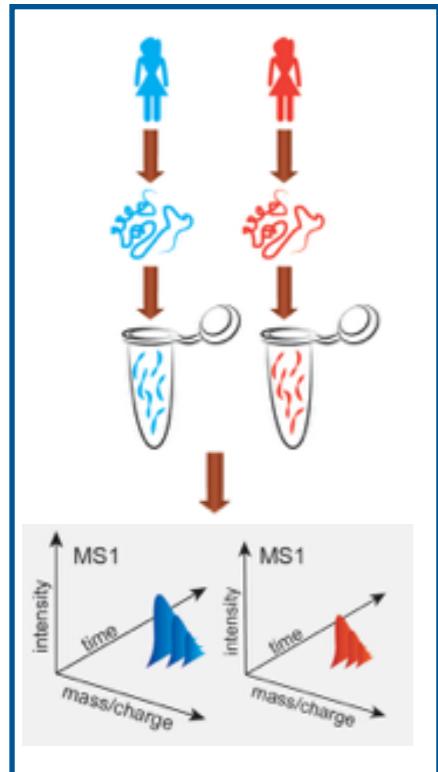
Northeastern University

OUTLINE

- So how many replicates do I need?
 - Design of complex experiments
- Associations between two measurements
 - Correlation and linear regression
- Statistical analysis of count data
 - Comparing proportions and spectral counts

PRINCIPLE I: REPLICATION

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

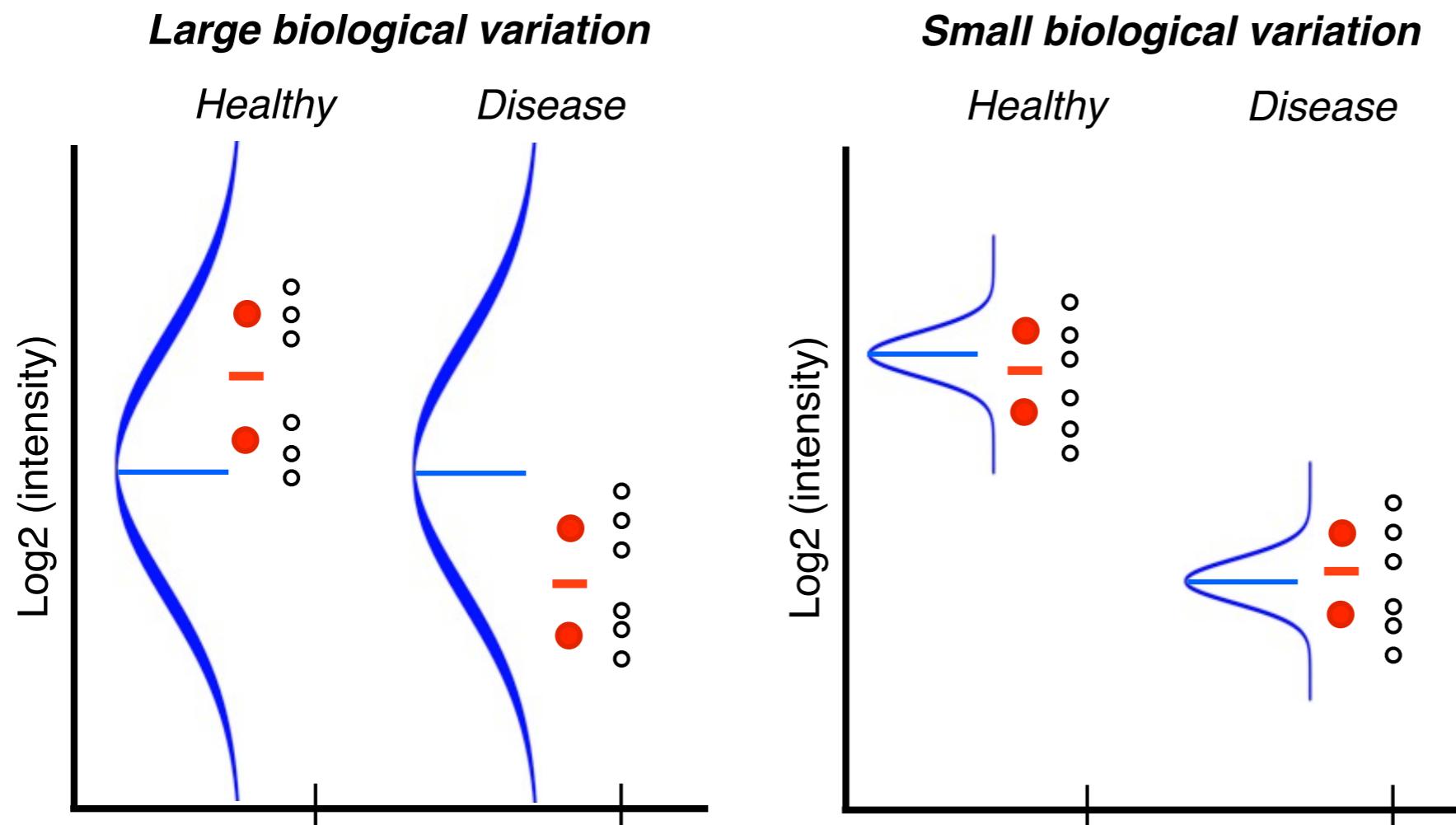


Two levels of randomness imply two types of replication:

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

MULTI-LAYER DESIGN AND ANALYSIS

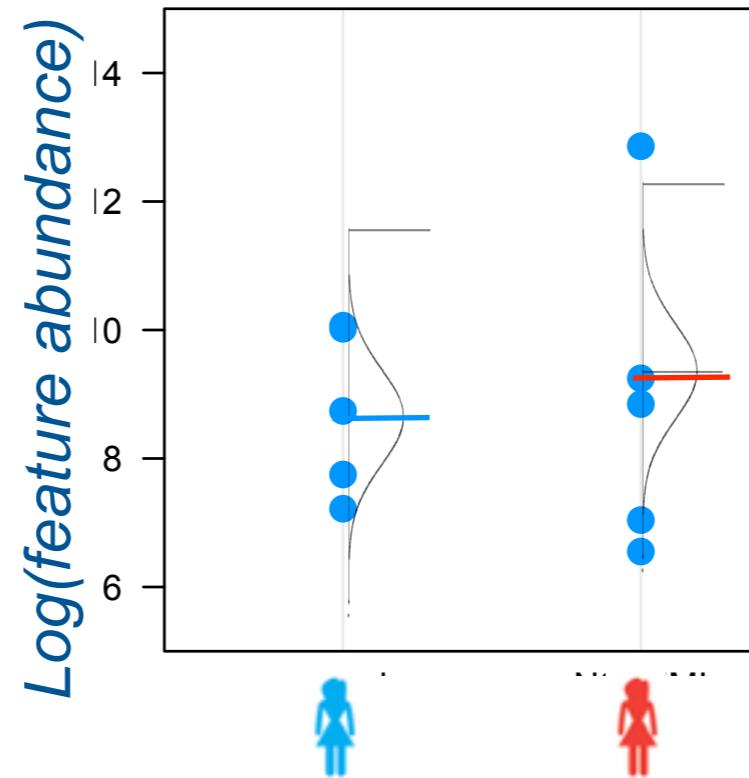
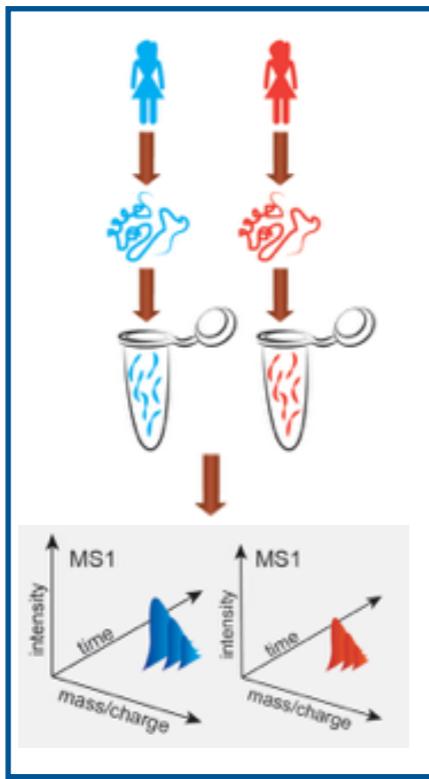
Multiple sources are responsible for variation in measurements



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

TWO-SAMPLE T-TEST

Simple example: label-free experiment, one feature/protein



H_0 : no change in abundance, $\mu_1 - \mu_2 = 0$

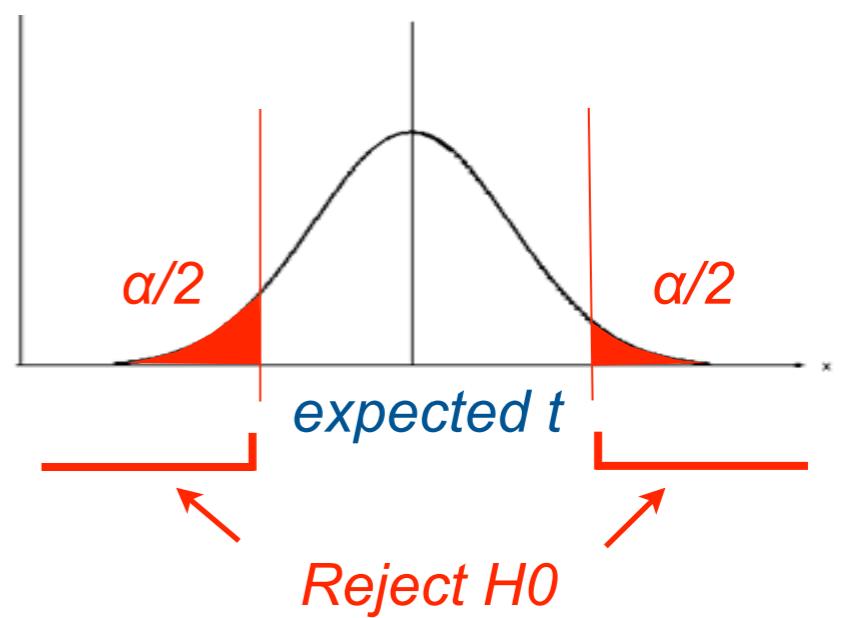
H_a : change in abundance, $\mu_1 - \mu_2 \neq 0$

$$t = \frac{\text{difference of group means}}{\text{estimate of variation}} = \frac{\bar{Y}_{1\cdot} - \bar{Y}_{2\cdot}}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

$$s_1^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_{1\cdot})^2$$

Distribution of the score if H_0 is true

α = False Positive Rate



SO HOW MANY REPLICATES DO I NEED?

One protein, two-group comparison

If we only have one feature:

Fix: α - probability of a false positive discovery

β - probability of a true positive discovery

Δ - anticipated fold change

σ_{Indiv}^2 and σ_{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?

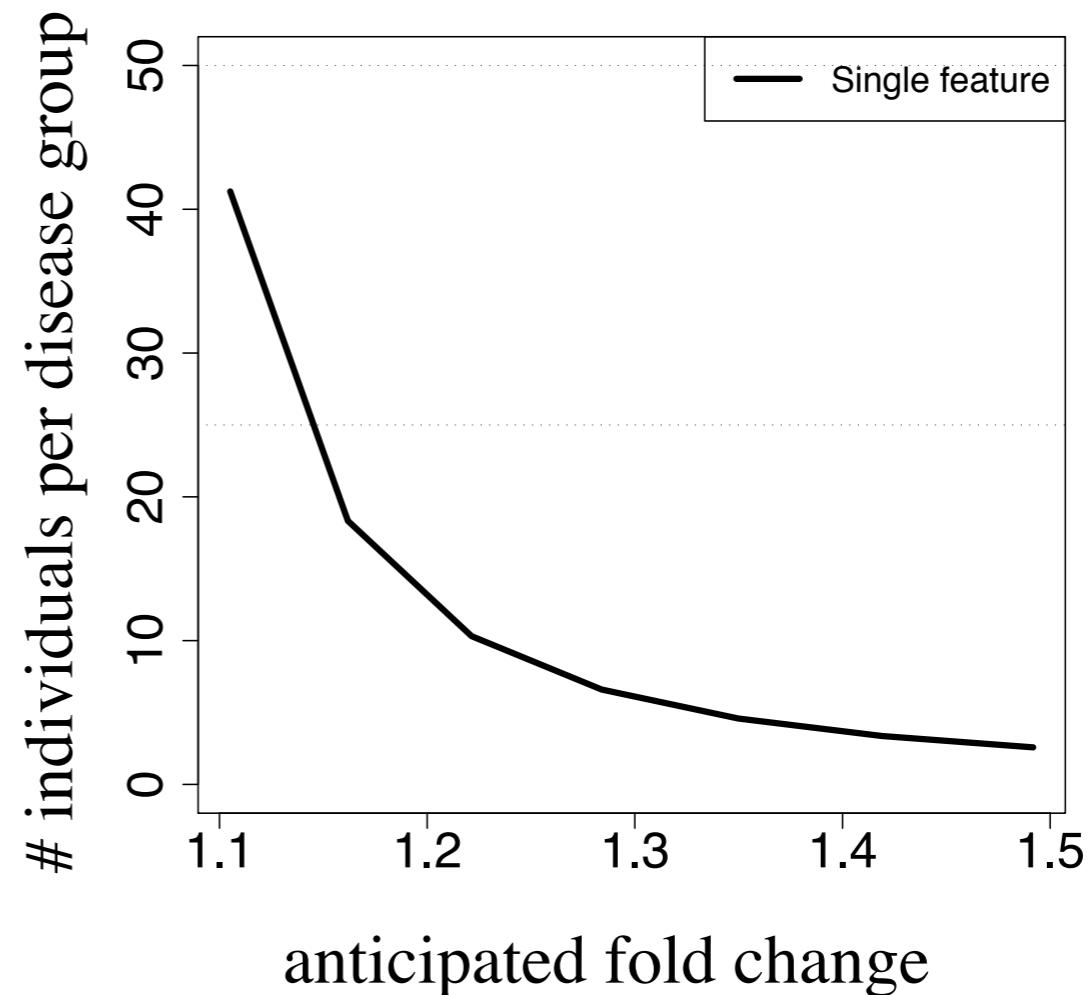
7

One protein, two-group comparison

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

SO HOW MANY REPLICATES DO I NEED?

Many proteins, two-group comparison

Would like to control the False Discovery Rate:

	# of features with no detected difference	# of features with detected difference	Total
# true non-diff. features	U	V	m_0
# true diff. features	T	S	$m_1 = m - m_0$
Total	$m - R$	R	m

$$q = E \left[\frac{V}{\max(R, 1)} \right] = \text{the “average” proportion of false positives}$$

This changes the sample size calculation:

Fix: q - the False Discovery Rate

m_0/m_1 - anticipated ratio of unchanging features

This defines

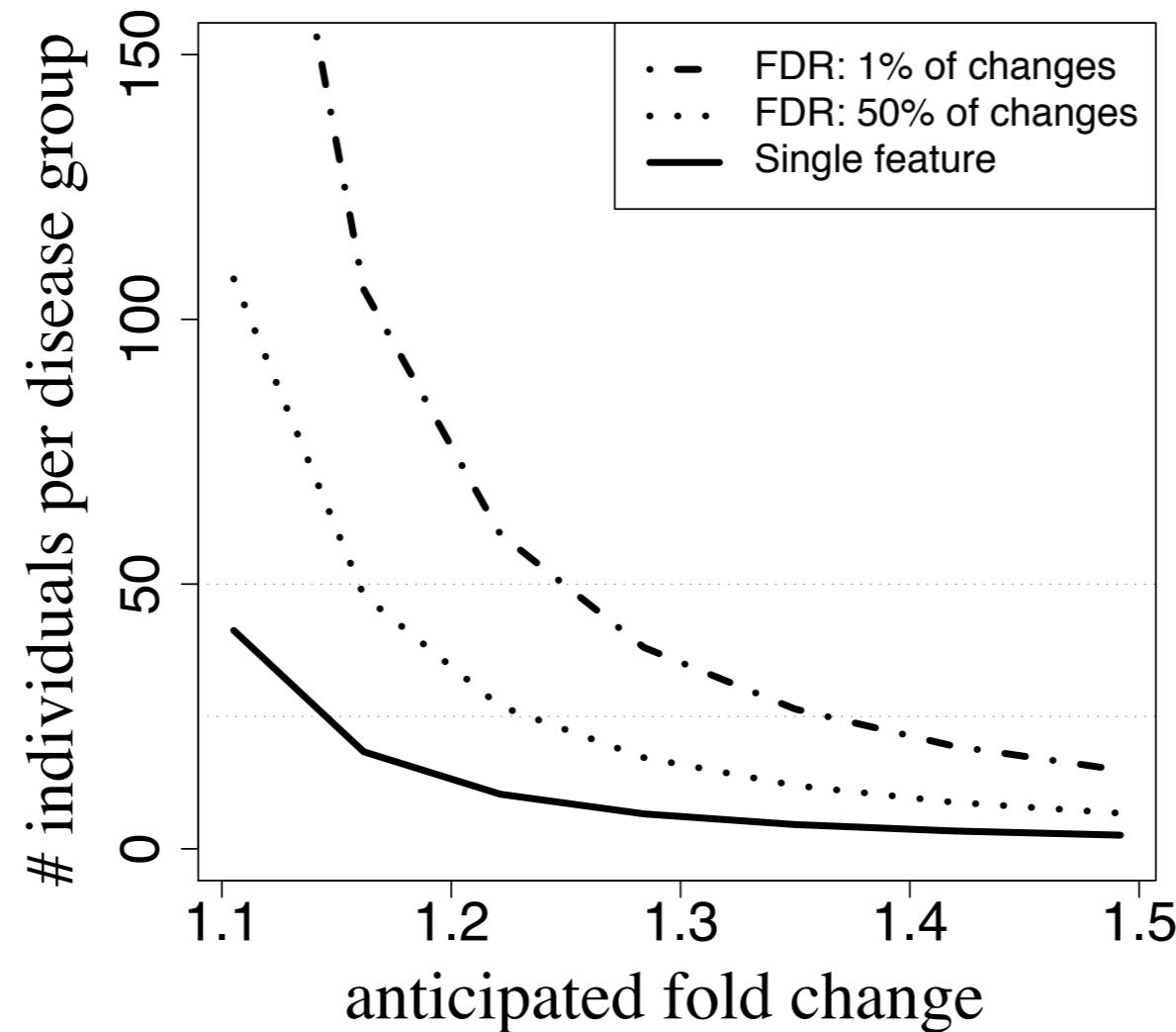
$$\alpha_{ave} \leq (1 - \beta)_{ave} \cdot q \frac{1}{1 + (1 - q) \cdot m_0/m_1},$$

- average probability of a false positive discovery

SO HOW MANY REPLICATES DO I NEED? ⁹

Many proteins, two-group comparison

Example: pilot study with diabetes patients.
A block-randomized design



Conclusion:

The fewer changes we expect, the larger the sample size

Oberg and
Vitek, *JPR*,
2009

COMPLEX DESIGNS: BIO AND TECH REPS

Biological replication is most important

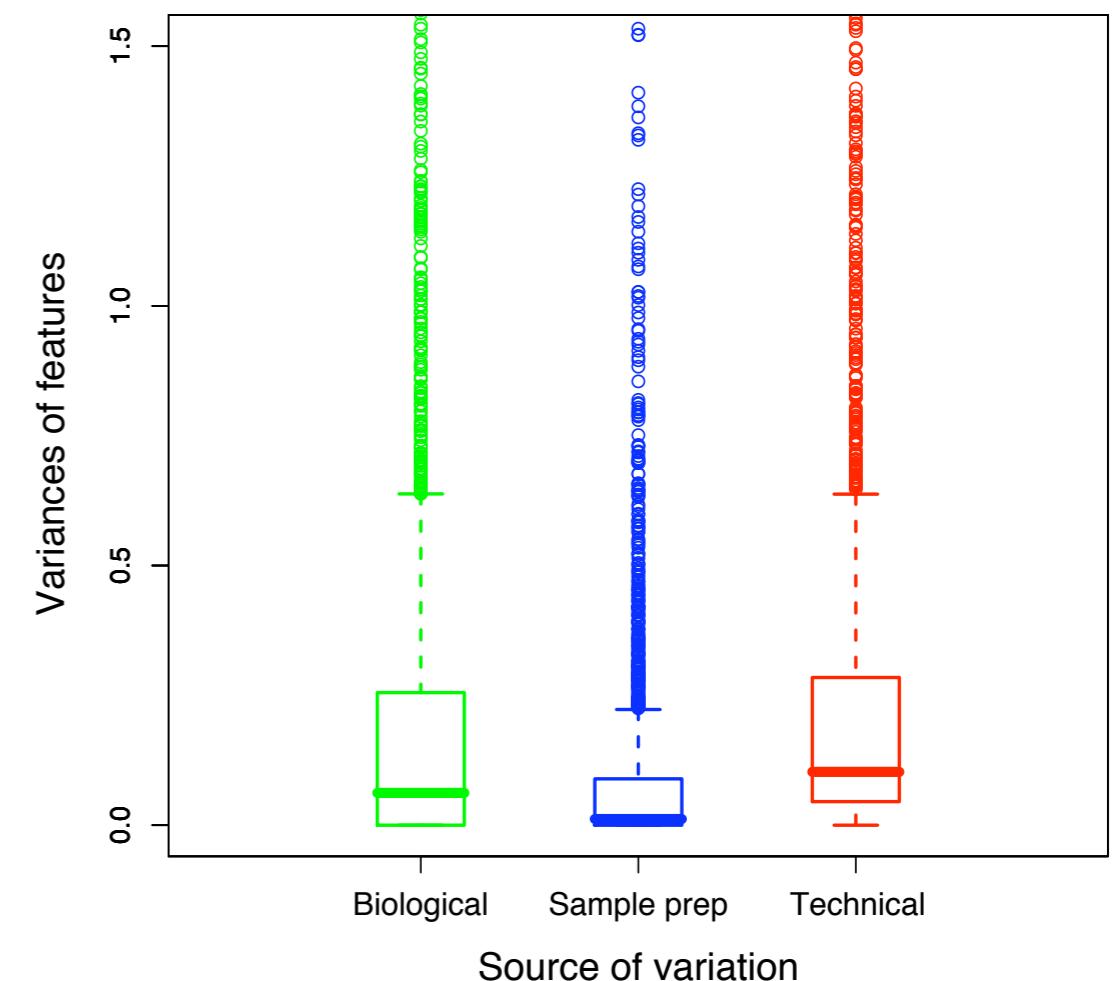
Observed feature intensity	=	Systematic mean signal of disease group	+	Random deviation due to individual	+	Random deviation due to sample preparation	+	Random deviation due to measurement error
y_{ijkl}	=	Group mean _i	+	$\text{Indiv}(\text{Group})_{j(i)} \sim N(0, \sigma_{\text{Indiv}}^2)$	+	$\text{Prep}(\text{Indiv})_{k(ij)}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

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

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

A pilot experiment

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



COMPLEX DESIGNS: BIO AND TECH REPS

Biological replication is most important

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

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

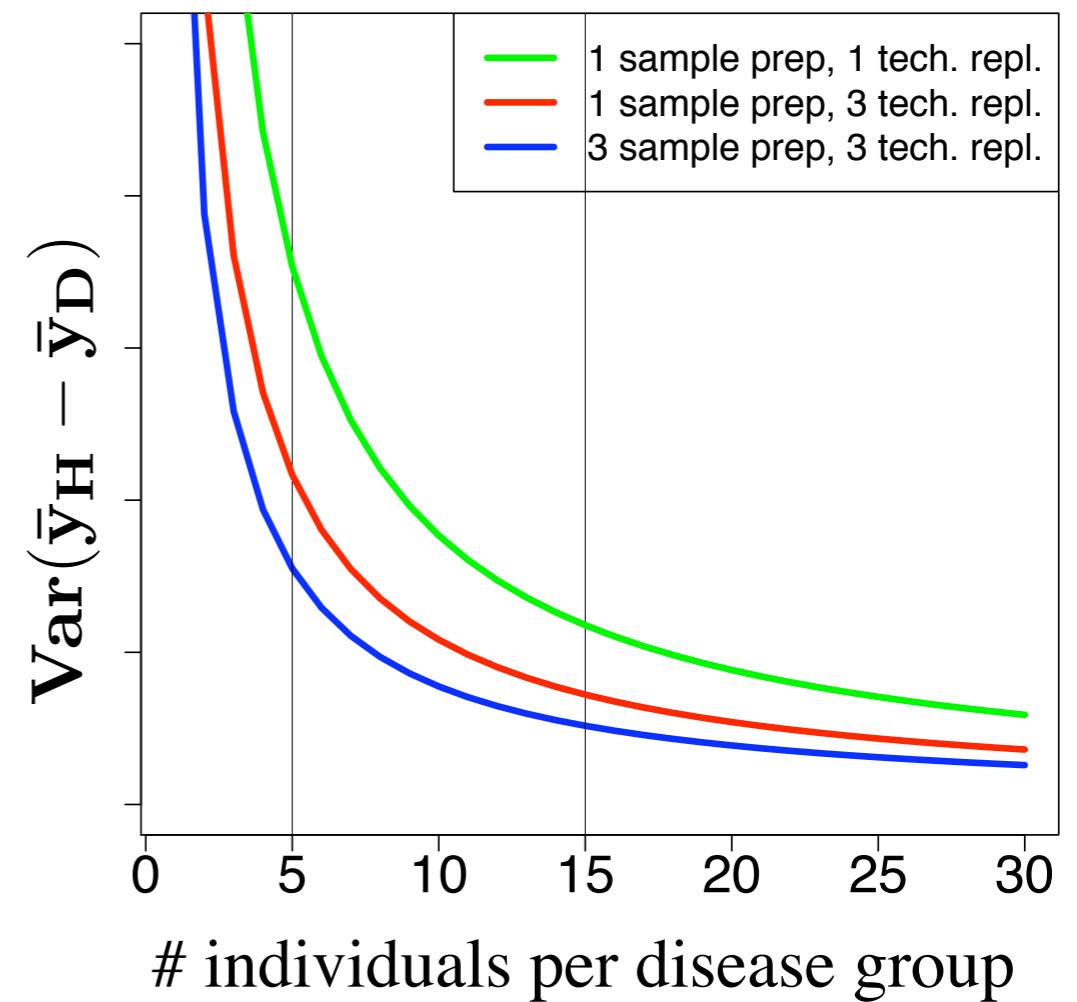
I: # individuals per disease group

J: # sample preps

K: # replicate runs

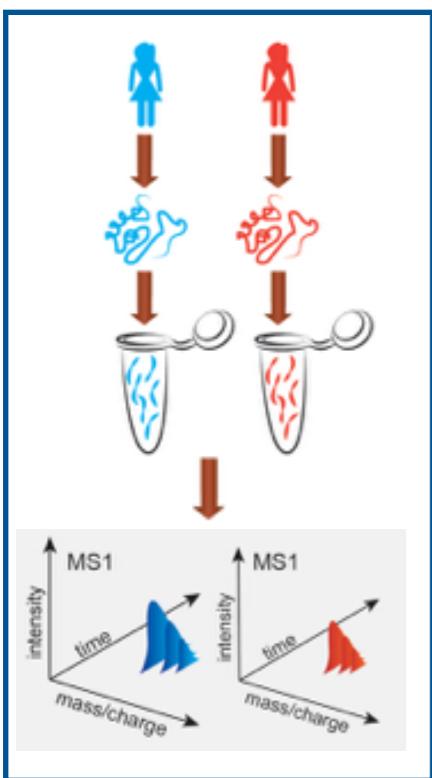
Conclusion:

Maximize the number
of biological replicates

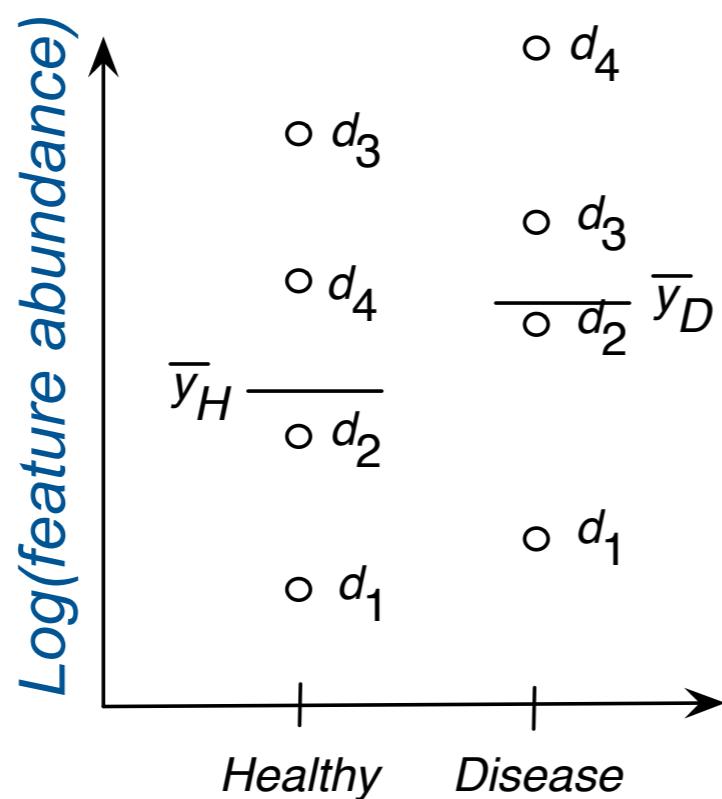


COMPLEX DESIGNS

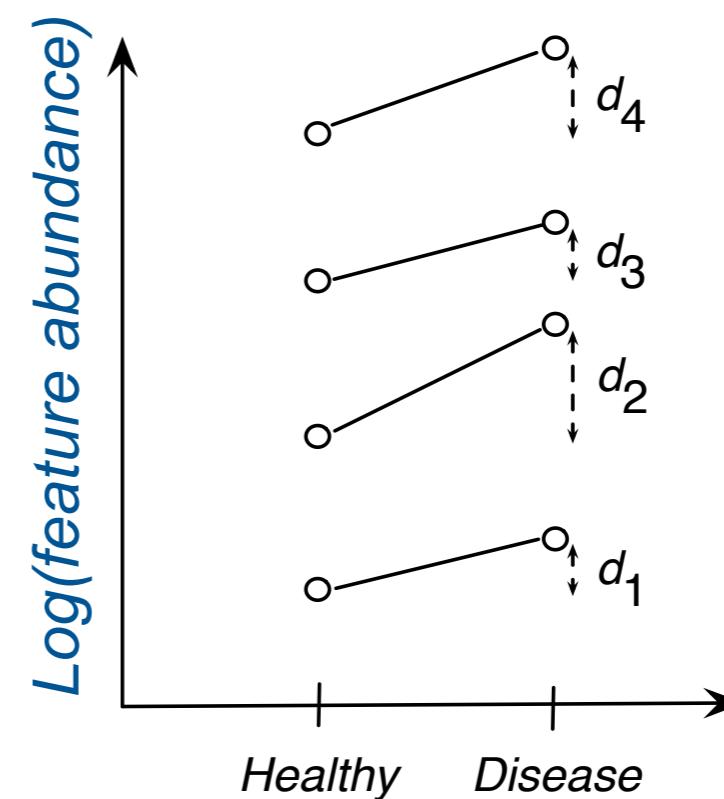
Blocking helps reduce the sample size



(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)

COMPLEX DESIGNS

Blocking helps most when between-block variance is large

Observed feature intensity y_{ijkl}	Systematic mean signal of disease group Group mean_i	Random deviation due to block (e.g. plate or day) $\sim N(0, \sigma_{\text{Block}}^2)$	Random deviation due to individual $\sim N(0, \sigma_{\text{Indiv}}^2)$	Random deviation due to measurement error $\sim N(0, \sigma_{\text{Error}}^2)$
---	---	--	---	--

A completely randomized design

I: # individuals per disease group

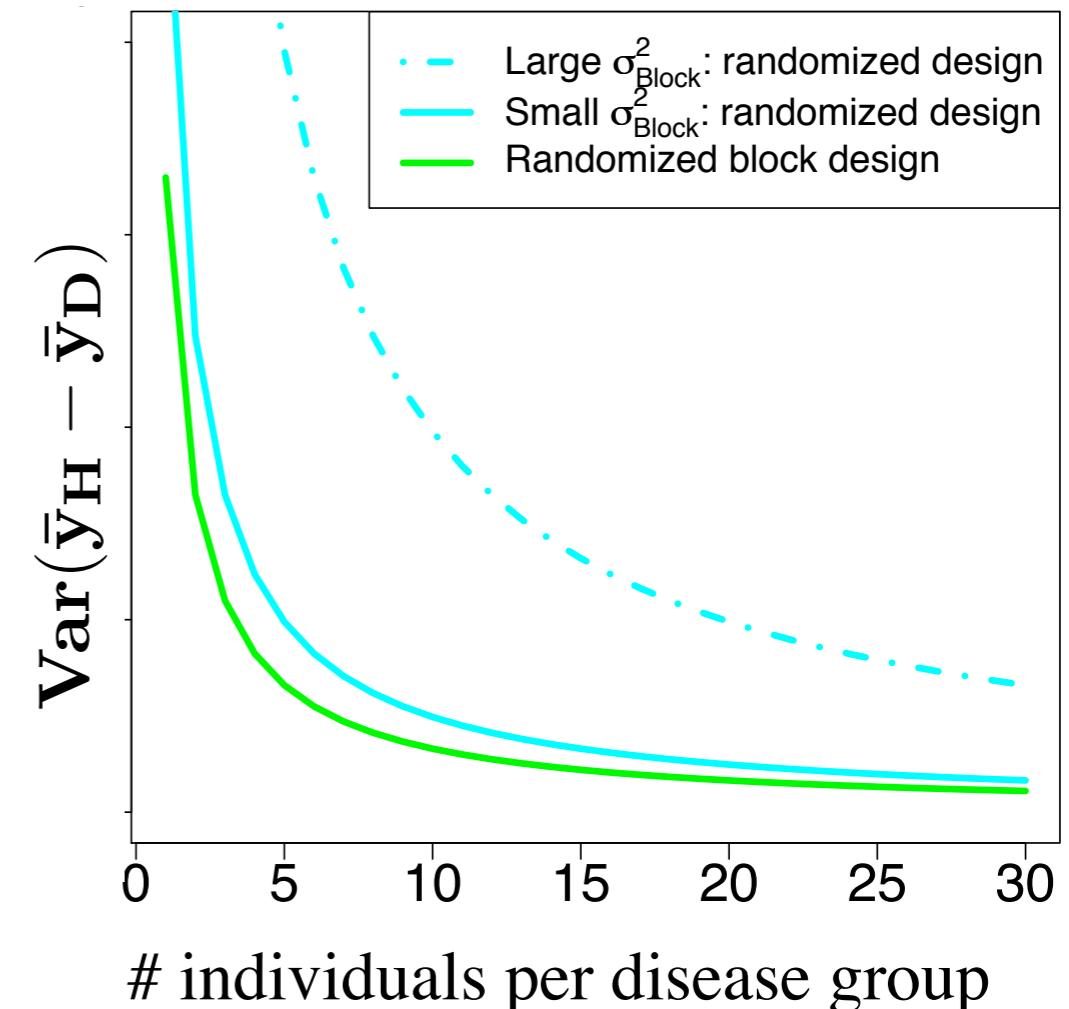
$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left(\frac{\sigma_{\text{Block}}^2 + \sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$

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?

One protein, blocking

If we only have one feature:

Fix: α - probability of a false positive discovery

β - probability of a true positive discovery

Δ - anticipated fold change

σ_{Indiv}^2 and σ_{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

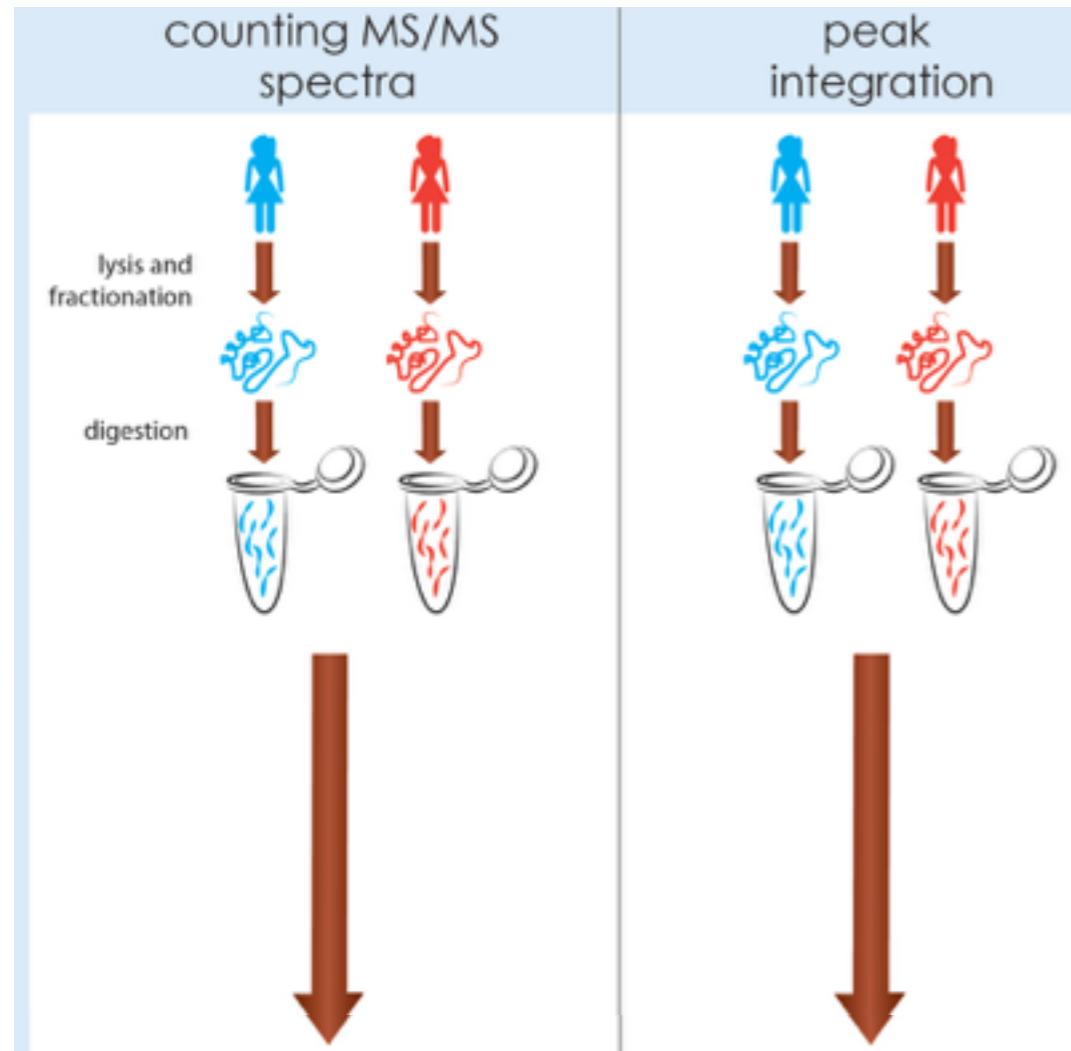
$$\text{Var}(\bar{y}_H - \bar{y}_D) = 2 \left(\frac{\sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2}{I} \right)$$



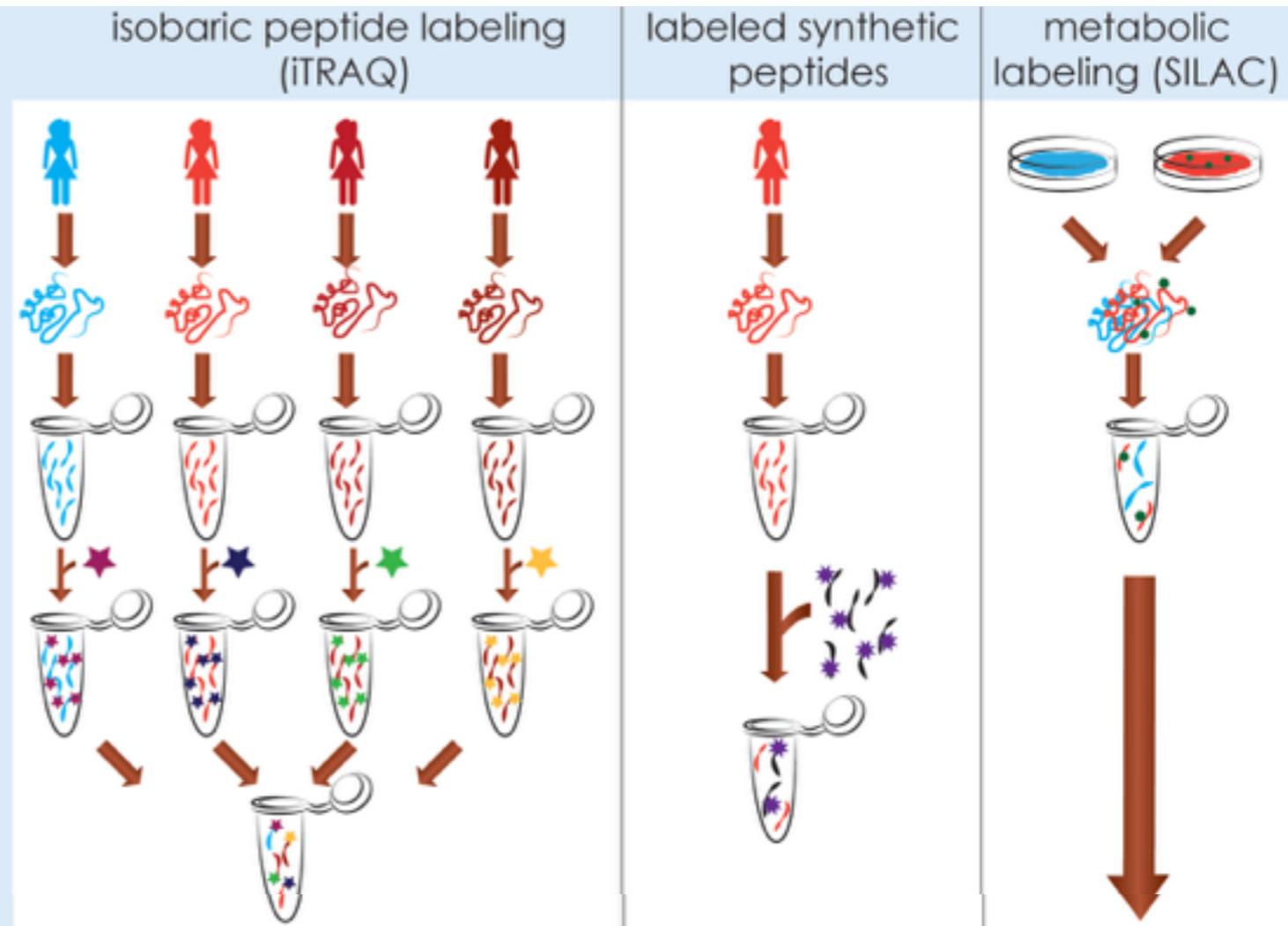
solve for the number of individuals I

COMPLEX DESIGNS: MULTIPLEXING

Label-free

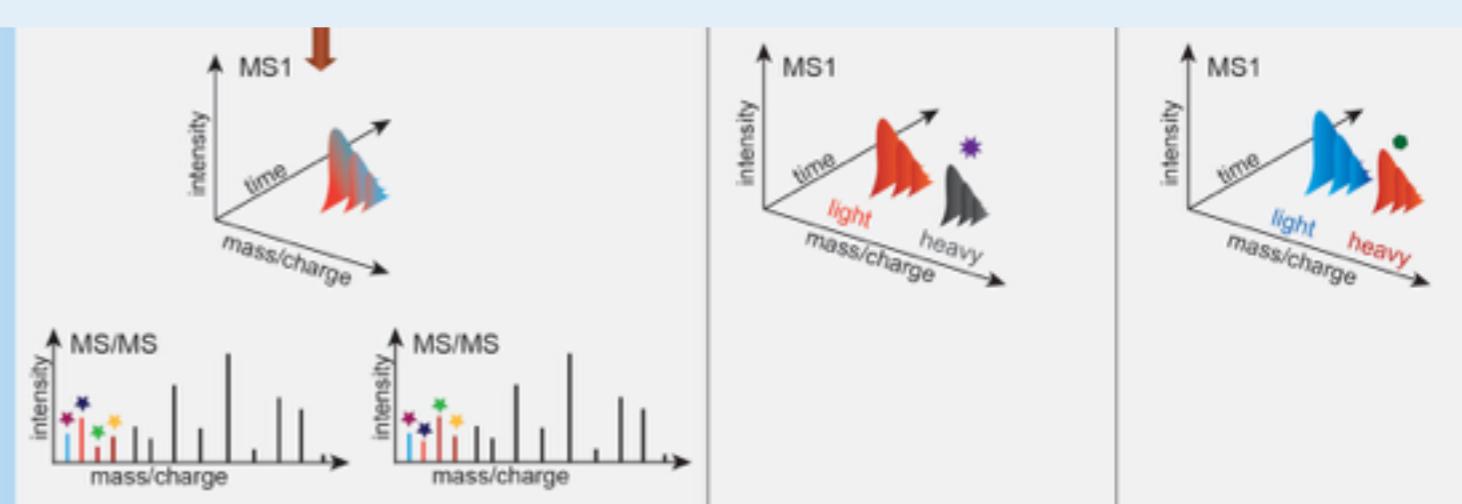


Label-based



Sample preparation

Global LC-MS/MS



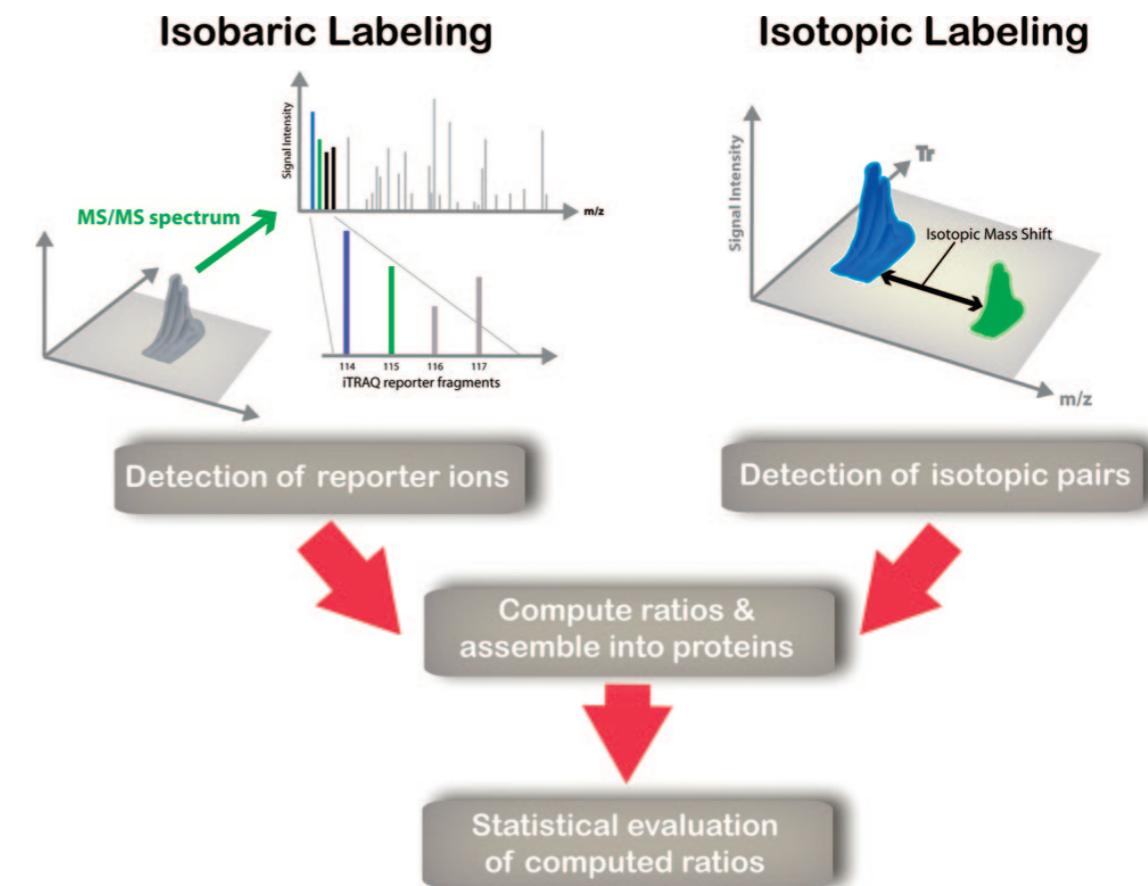
MULTIPLEXING: RUN = BLOCK

Balanced incomplete block design

Each sample appears with samples from other groups an equal number of times

Example: 5 groups, 4 labels, run=block

Disease group	Replicate set 1					...
	Block 1	Block 2	Block 3	Block 4	Block 5	
D_1	X	X	X	X		...
D_2	X	X	X		X	...
D_3	X	X		X	X	...
D_4	X		X	X	X	...
D_5		X	X	X	X	...



Example: 5 groups, 2 labels, run=block

Disease group	Block 1	Block 2	Block 3	Block 4	Block 5	Block 6	Block 7	Block 8	Block 9	Block 10	...
D_1	X_{L_1}	X_{L_2}	X_{L_1}	X_{L_2}							...
D_2	X_{L_2}				X_{L_1}	X_{L_2}	X_{L_1}				...
D_3			X_{L_1}			X_{L_2}			X_{L_1}	X_{L_2}	...
D_4				X_{L_2}			X_{L_1}		X_{L_2}		...
D_5					X_{L_1}		X_{L_2}		X_{L_1}	X_{L_2}	...

COMPLEX DESIGNS: MULTIPLEXING

Find allocations that minimize variance

A block-randomized design

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

Balanced incomplete block design

$$\text{Var}(\hat{D}_1 - \hat{D}_2) = 2 \frac{n_b}{n_g n_p n_s} (\sigma_{\text{Indiv}}^2 + \sigma_{\text{Error}}^2)$$

I: # individuals per disease group

n_b, n_g, n_p, n_s - counts of paired allocations

Conclusion:

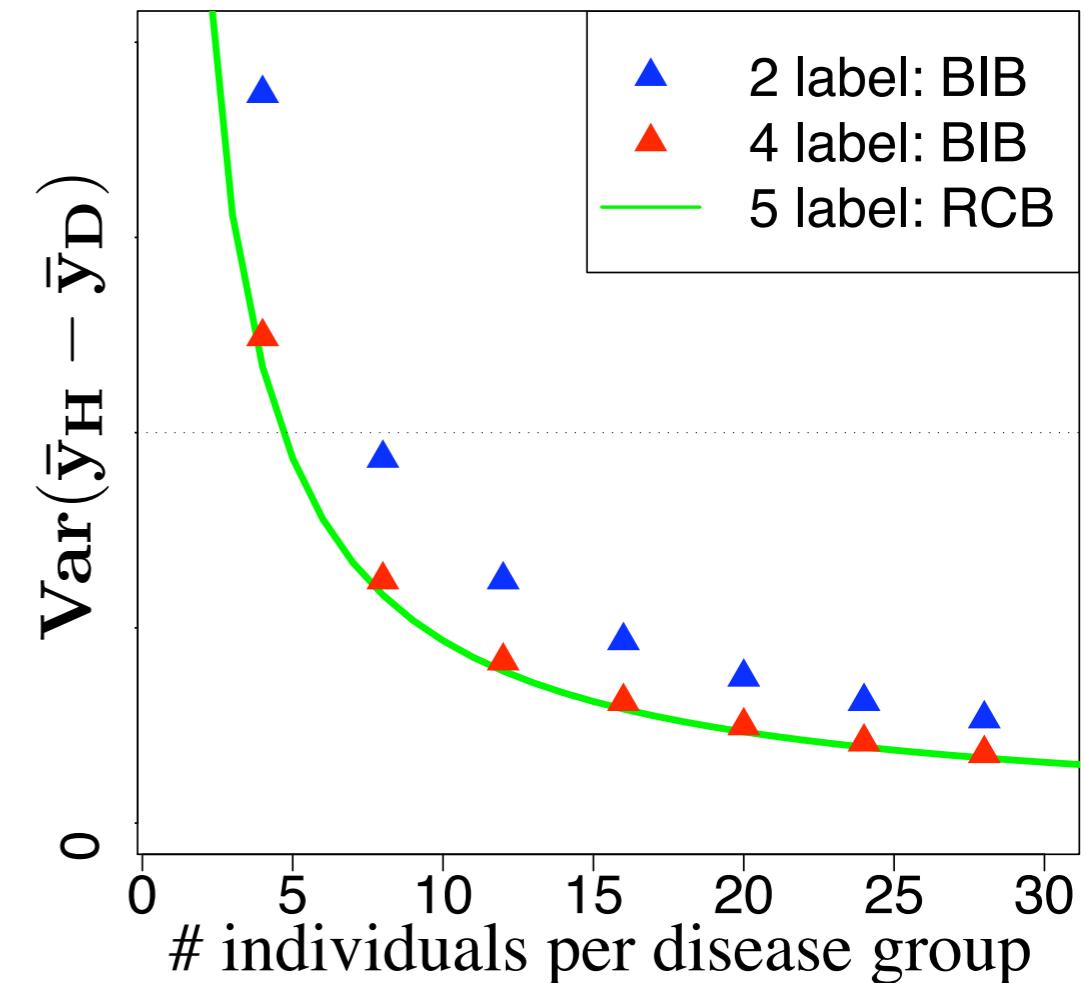
Limited sample size options

Larger blocks are more efficient

(assuming same variances in all workflows)

A pilot experiment

- 3 groups of cardiovascular disease
- iTRAQ workflow



COMPLEX DESIGN: POOLING

Helpful in theory; useful for a subset of research objectives

- Goal: pool biological specimens from a same group prior to the experiment
 - Insufficient biological material
 - Need to gain efficiency
- Strategy
 - Cannot pool all samples from a group
(will not be able to assess the biol. variation)
 - Can pool multiple disjoint subsets of subjects

COMPLEX DESIGN: POOLING

Helpful in theory; useful for a subset of research objectives

- If pooled samples are viewed as averages of biological material, then
 - For a CRD with 1 tech. rep.:

$$Var(\hat{\mu}_1 - \hat{\mu}_2) = 2 \left(\frac{\sigma_{Indiv}^2 + \sigma_{Prep}^2 + \sigma_{Error}^2}{I} \right).$$

- Pool of r specimens:

$$Var(\hat{\mu}_1 - \hat{\mu}_2) = 2 \left(\frac{\sigma_{Indiv}^2}{I r} + \frac{\sigma_{Prep}^2 + \sigma_{Error}^2}{I} \right).$$

COMPLEX DESIGN: POOLING

Helpful in theory; useful for a subset of research objectives

- Disadvantage:
 - Potential bias: pipetting errors cause unequal contributions of specimens to the pool
 - Difficult to detect outlying or contaminated specimens
 - Downstream analysis limited to testing (i.e. no prediction/classification for individual subjects)
- Required assumption: biological averaging
 - Model-based averaging on the log scale
 - Physical mixing on the original scale
 - Assume that the two scales are equivalent (generally not true)

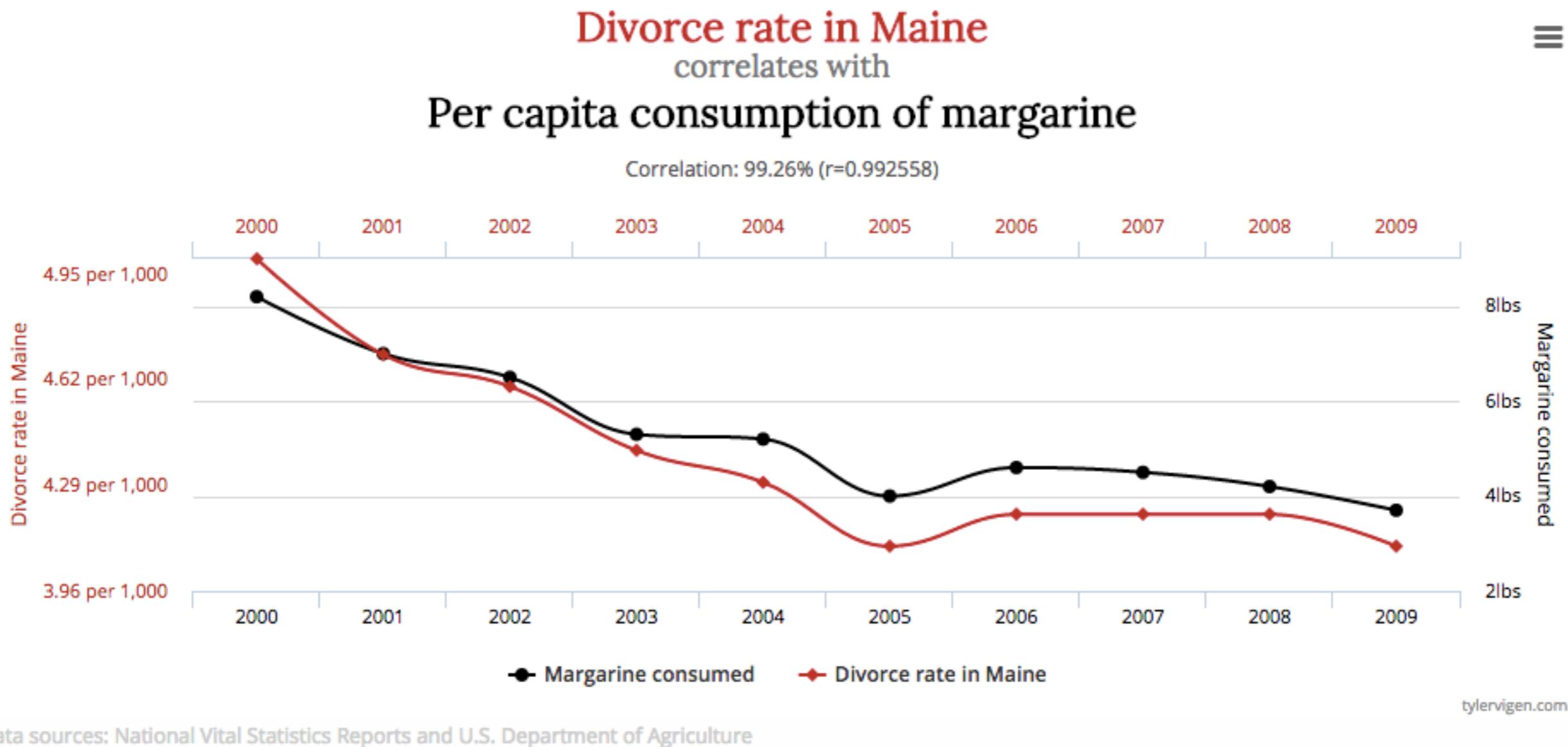
OUTLINE

- So how many replicates do I need?
 - Design of complex experiments
- Associations and calibration
 - Correlation and linear regression
- Statistical analysis of count data
 - Comparing proportions and spectral counts

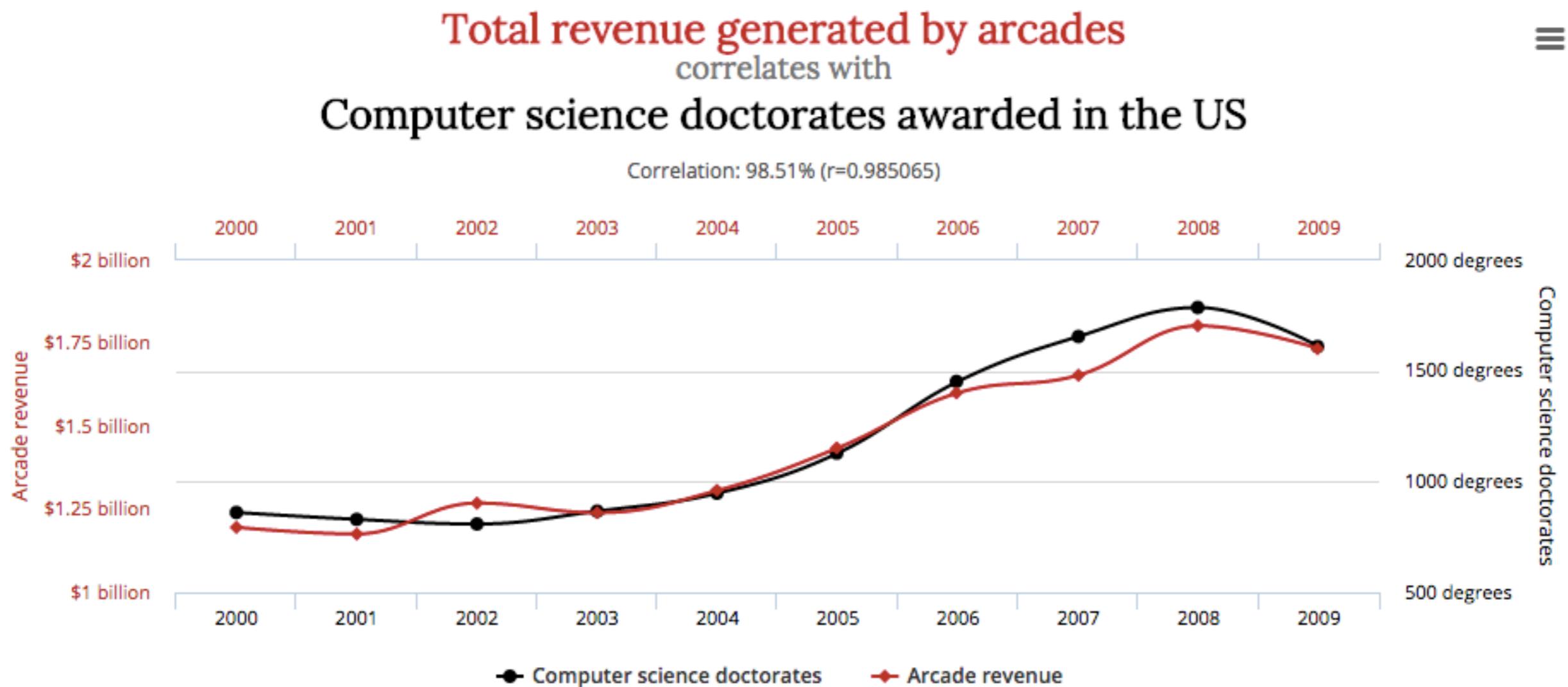
ASSOCIATIONS AND CALIBRATIONS

- Associations
 - do two proteins co-vary across samples?
- Calibrations
 - does protein abundance reflect the true abundance?

SPURIOUS ASSOCIATIONS ABOUND



SPURIOUS ASSOCIATIONS ABOUND



Data sources: U.S. Census Bureau and National Science Foundation

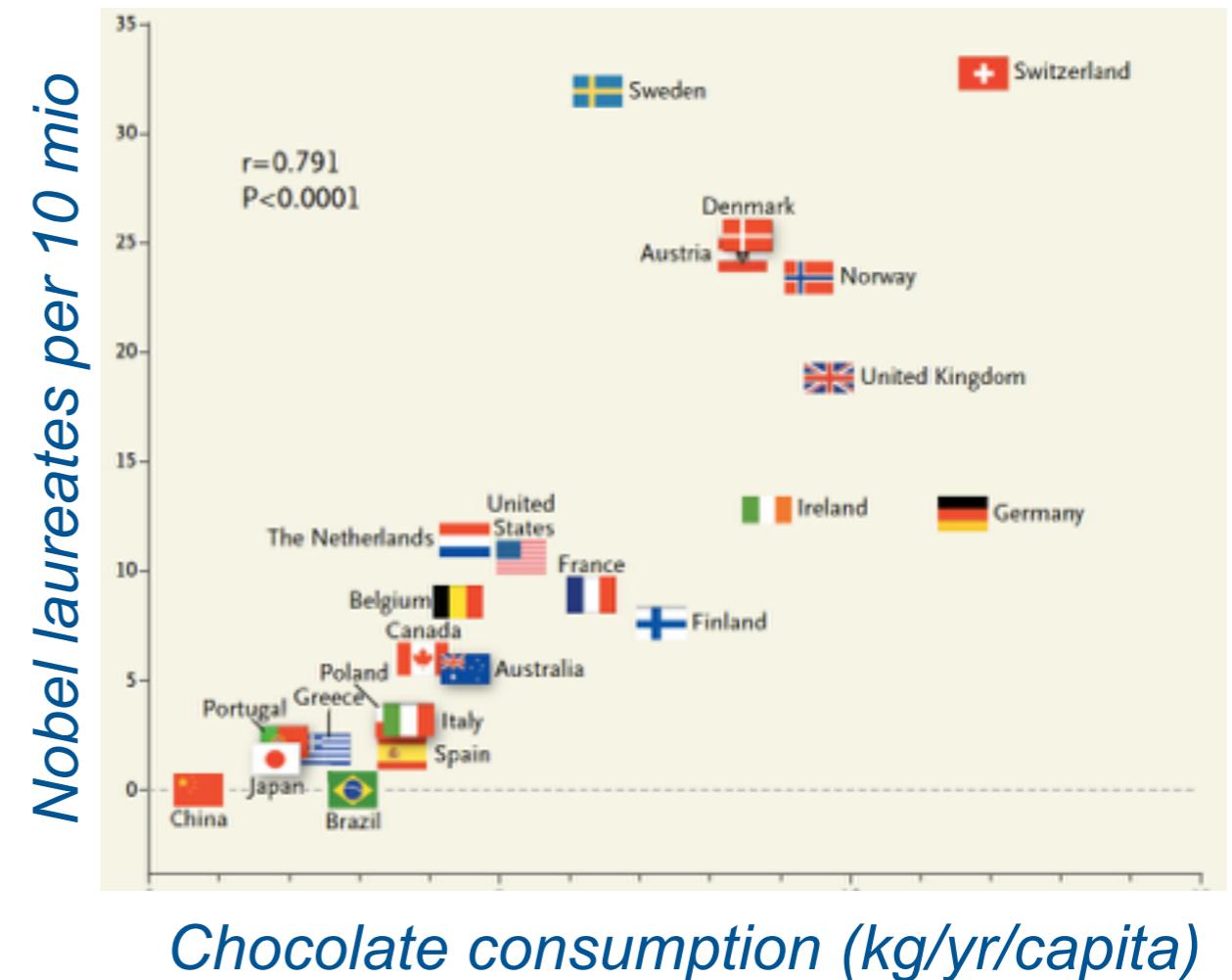
tylervigen.com

tylervigen.com/spurious-correlations

SPURIOUS ASSOCIATIONS ABOUND

Easy to dismiss when we understand the context

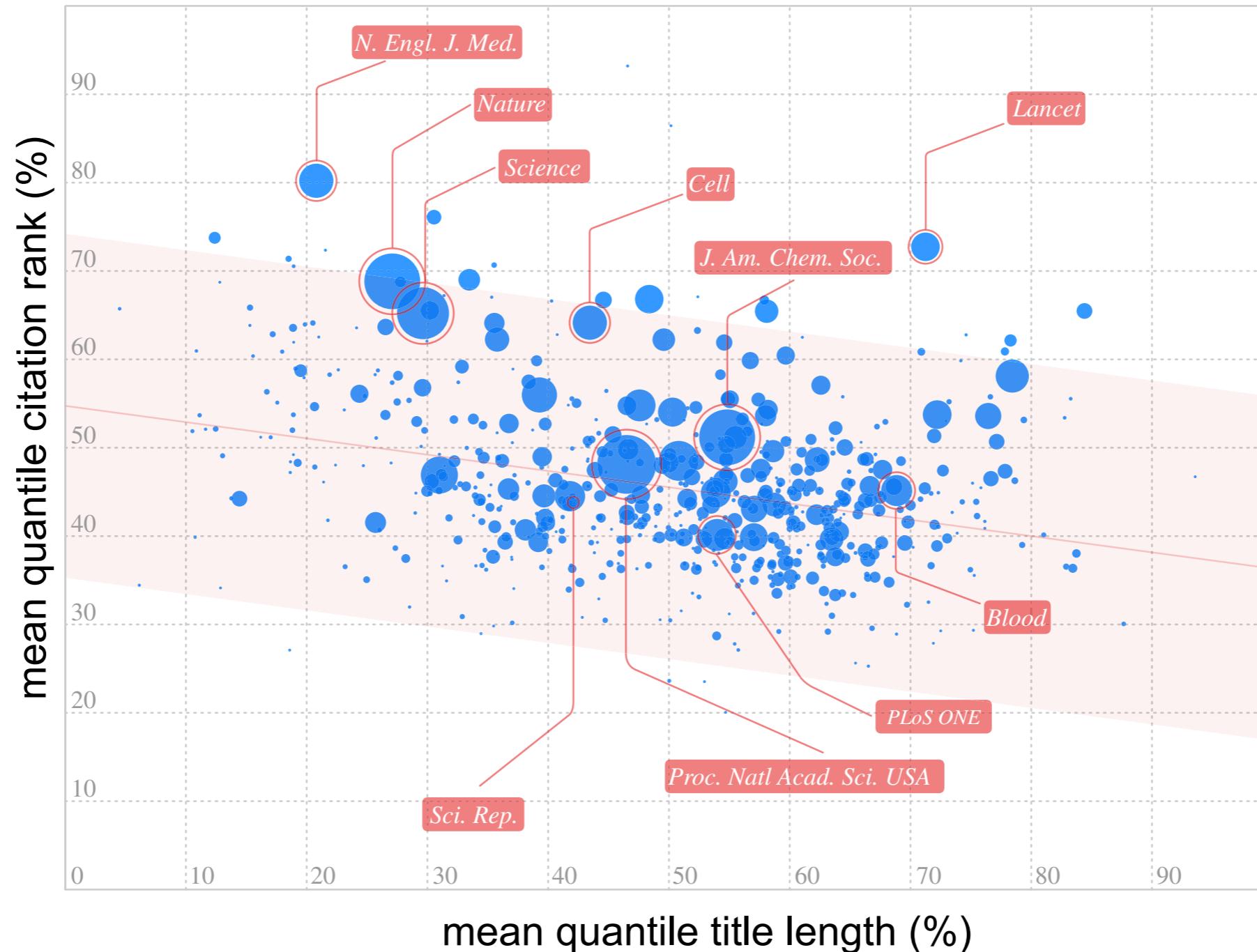
- Premier medical journal
 - *Nobel prize is related to cognitive ability*
 - *flavanols (organic molecules present in chocolate) are linked to cognitive ability*
- Technical flows
 - *Nobel prize winners between 1900-2011*
 - *Chocolate consumption after 2002*
 - *Countries with many Nobel prizes have a high Human Development Index and high per capita income*



New England Journal of Medicine, 367:1562 (2012)
A. Jogalekar, Scientific American, 2012

SPURIOUS ASSOCIATIONS ABOUND

Length of title negatively correlates with number of citations



A. Letchford et al., Royal Society Publishing, 2015

SPURIOUS ASSOCIATIONS ABOUND

Length of title negatively correlates with number of citations



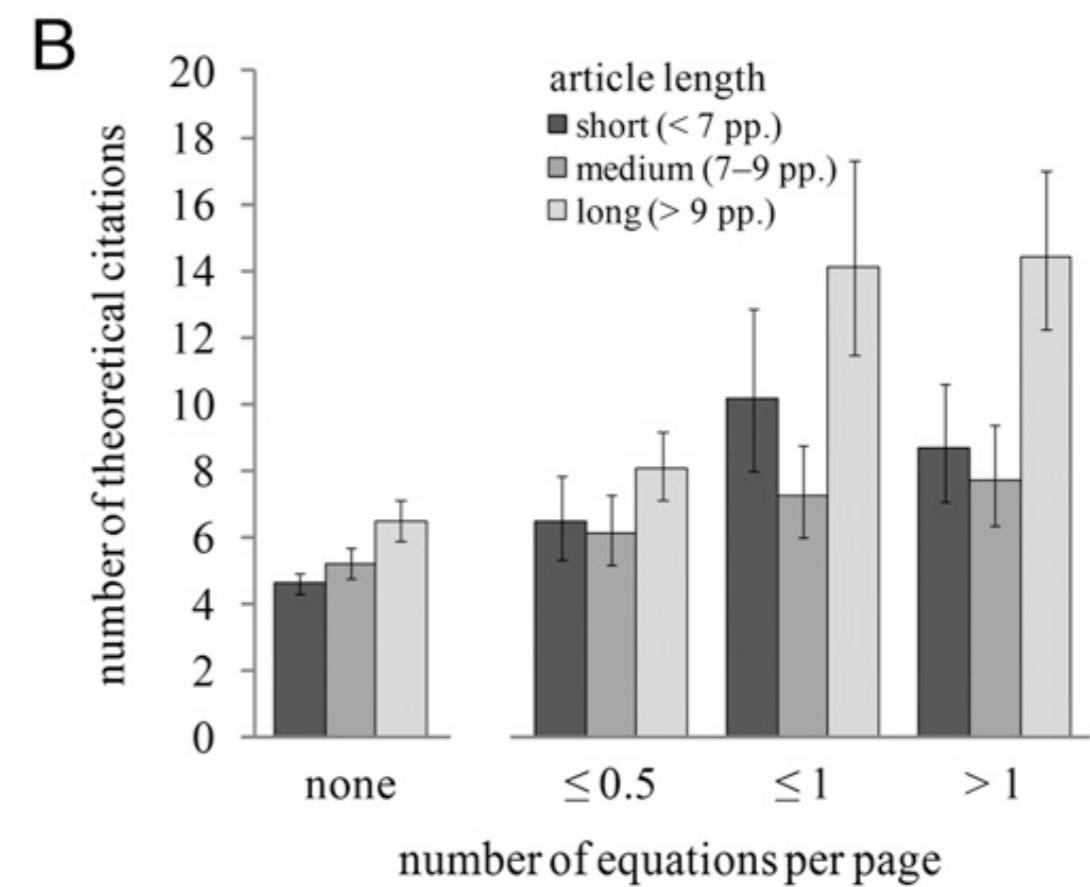
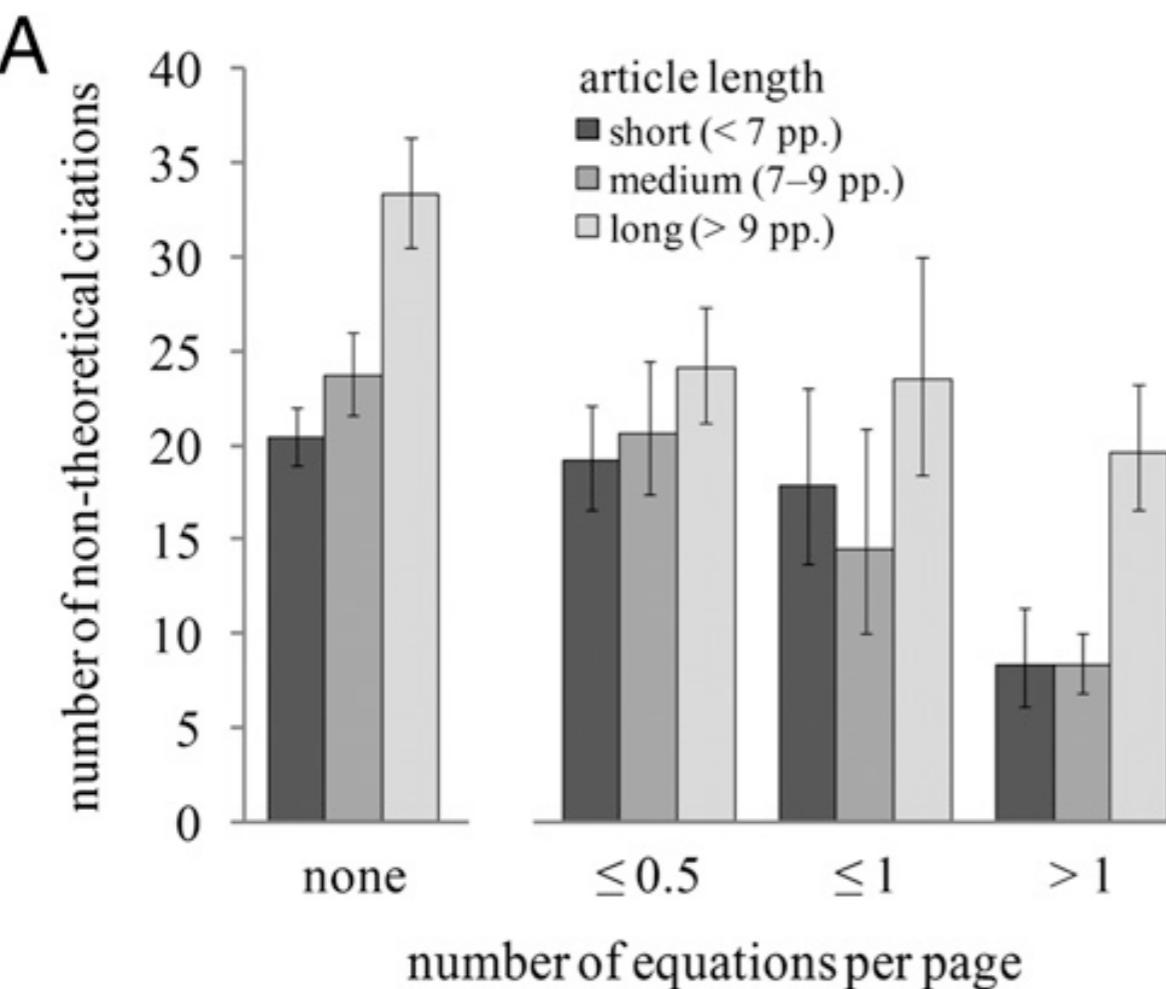
A. Letchford et al., Royal Society Publishing, 2015

Heavy use of equations impedes communication among biologists

Tim W. Fawcett¹ and Andrew D. Higginson

School of Biological Sciences, University of Bristol, Bristol BS8 1UG, United Kingdom

Edited[†] by Robert M. May, University of Oxford, Oxford, United Kingdom, and approved June 6, 2012 (received for review April 4, 2012)



SIMPLE LINEAR REGRESSION

Correlation between

x_1, \dots, x_n and y_1, \dots, y_n

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Linear regression for

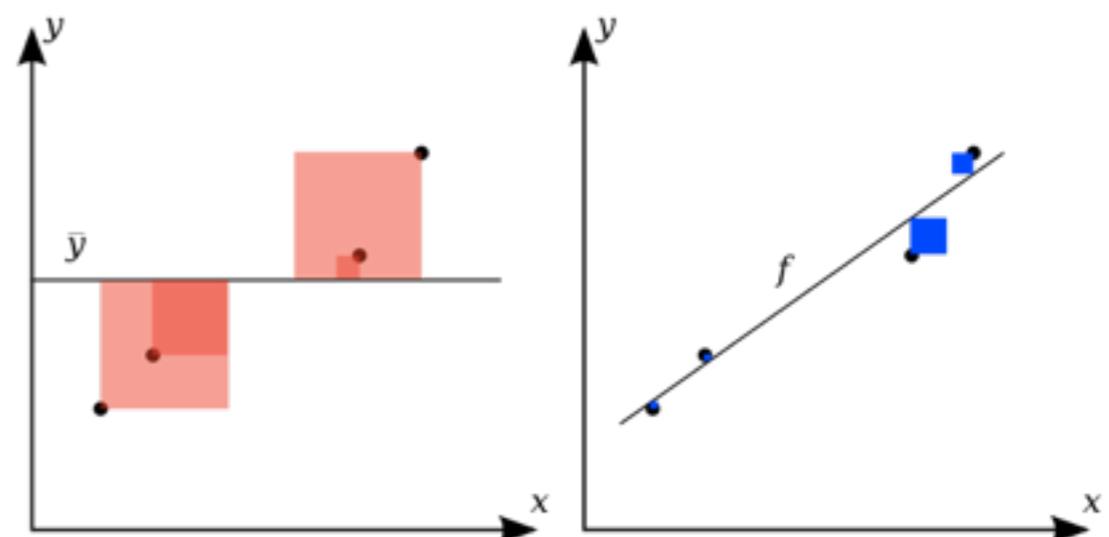
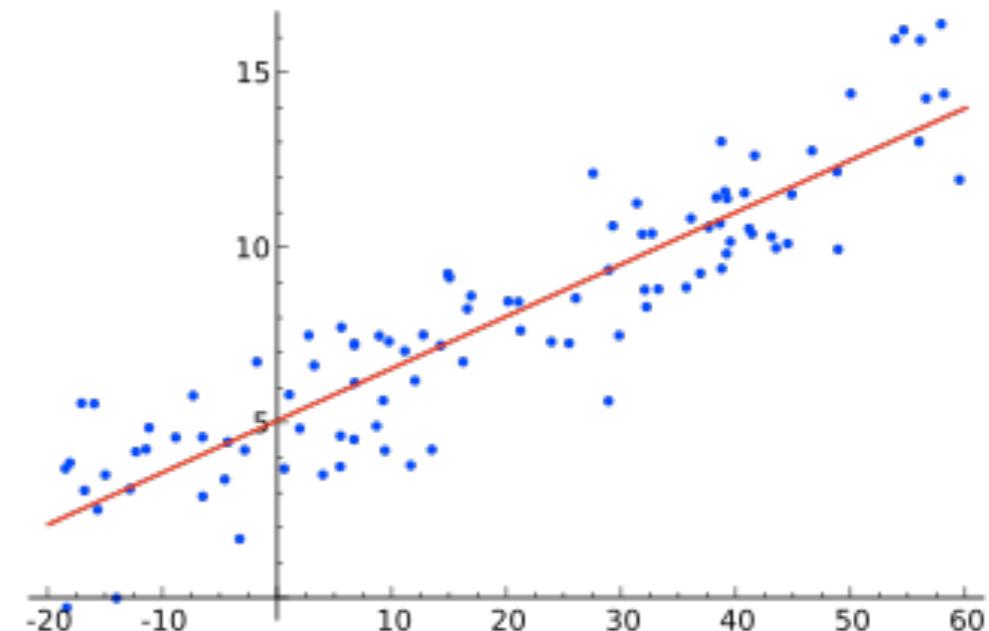
x_1, \dots, x_n and y_1, \dots, y_n

$$y_i = \alpha + \beta x_i + \varepsilon_i.$$

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\text{Cov}[x, y]}{\text{Var}[x]}$$

$$= r_{xy} \frac{s_y}{s_x},$$



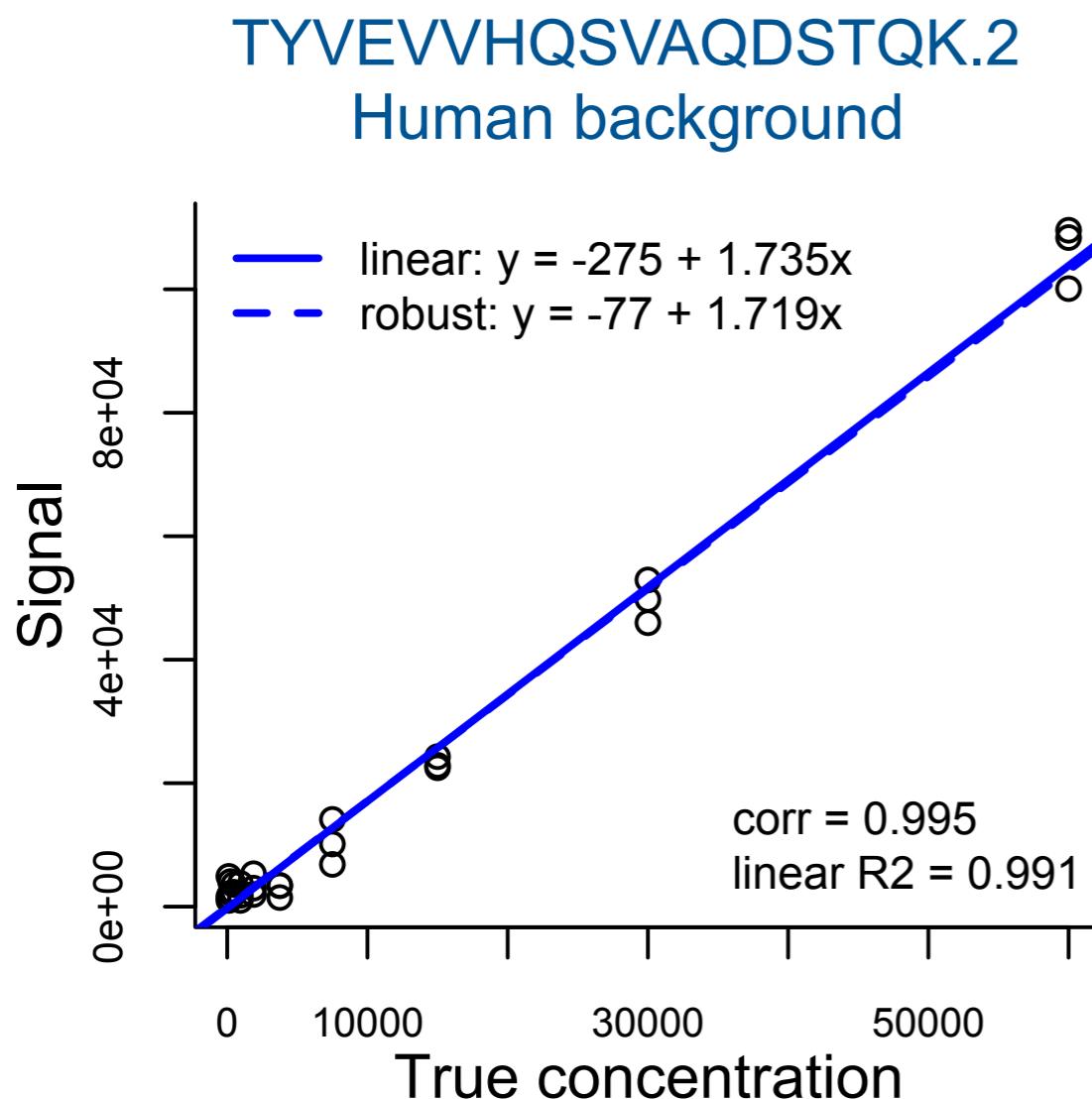
$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

https://en.wikipedia.org/wiki/Simple_linear_regression

https://en.wikipedia.org/wiki/Coefficient_of_determination

$$R^2 = r_{xy}^2$$

EXAMPLE: CALIBRATION



- Motivating example
 - ◆ OpenSWATH
(Röst *et al*, *Nature Biotech*, 2014)
 - ◆ 387 peptide ions
 - ◆ 10 log-spaced concentrations
 - ◆ water, yeast, human backgrounds
- Goal: performance evaluation
 - ◆ Detect known concentrations
 - ◆ Quantify background noise
 - ◆ Quantify the slope
 - ◆ Quantify LoD and LoQ

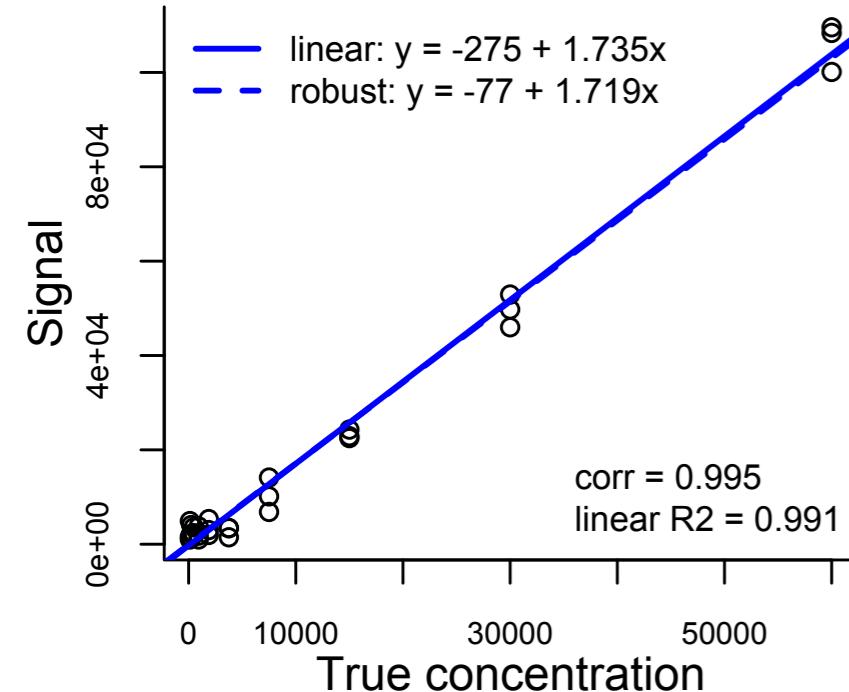
- From the graph
 - ◆ Linear relationship is theoretically plausible
 - ◆ High correlation, high R²
 - ◆ Ordinary least squares fit agrees with robust fit

→ False perception of a good agreement

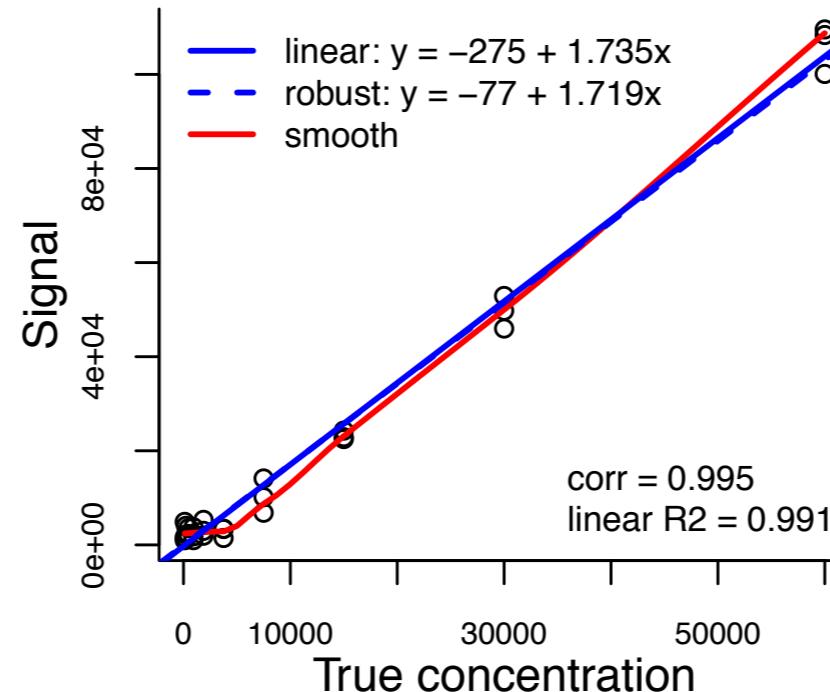
PROBLEM I

Linear fit may be uninterpretable, or wrong

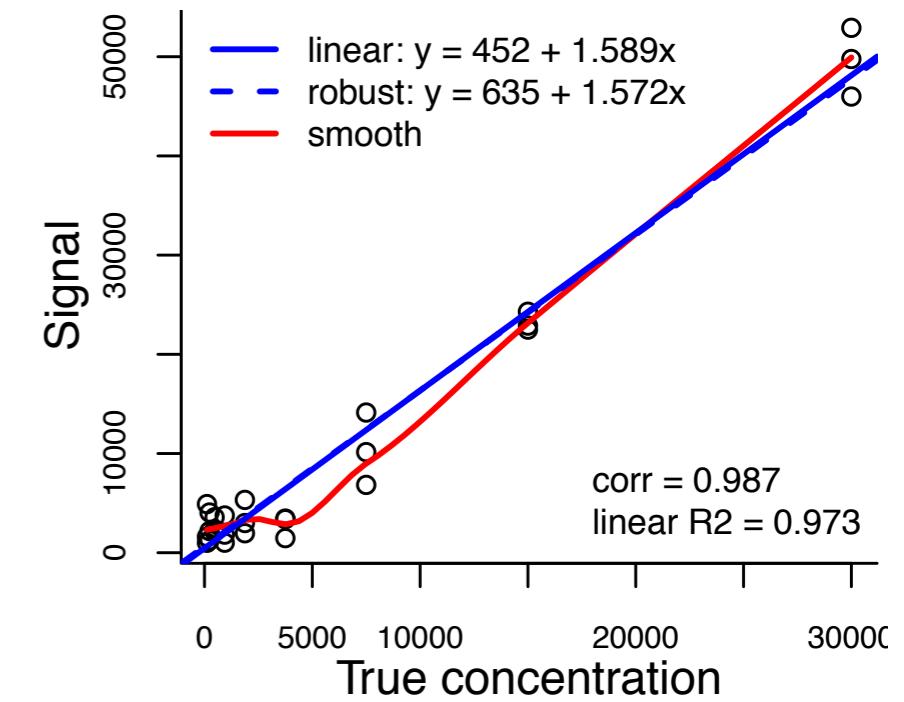
All concentrations



All concentrations,
smoothed



No top concentration,
smoothed



Negative intercept
(no interpretation)

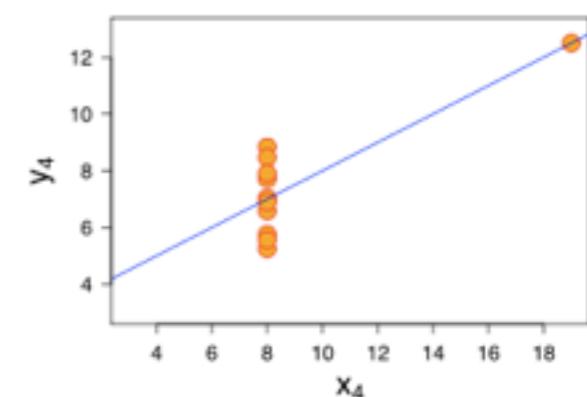
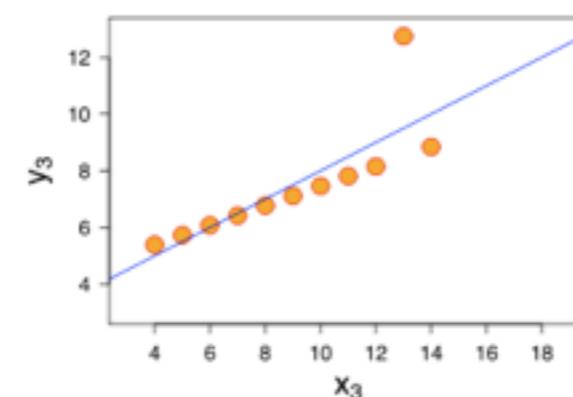
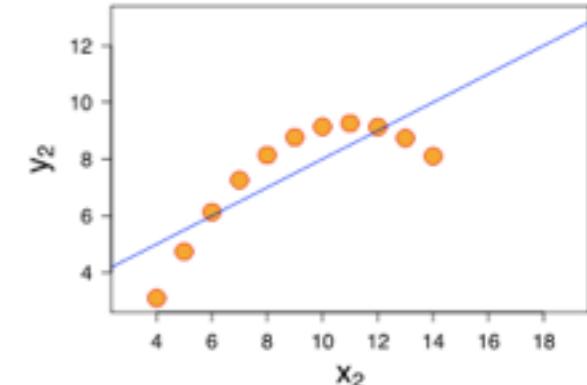
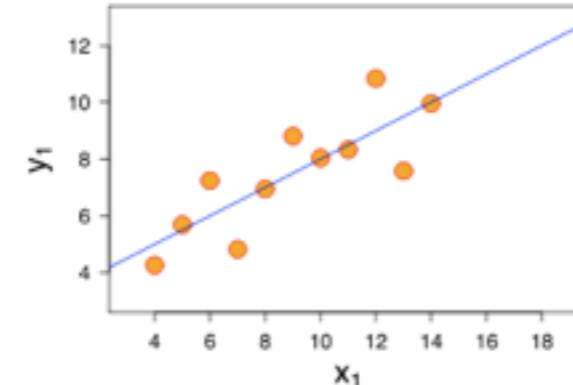
Perception of linearity depends
on the # of concentrations

Correlation, slope and R²
do not quantify linearity

PROBLEM I

Linear fit may be uninterpretable, or wrong

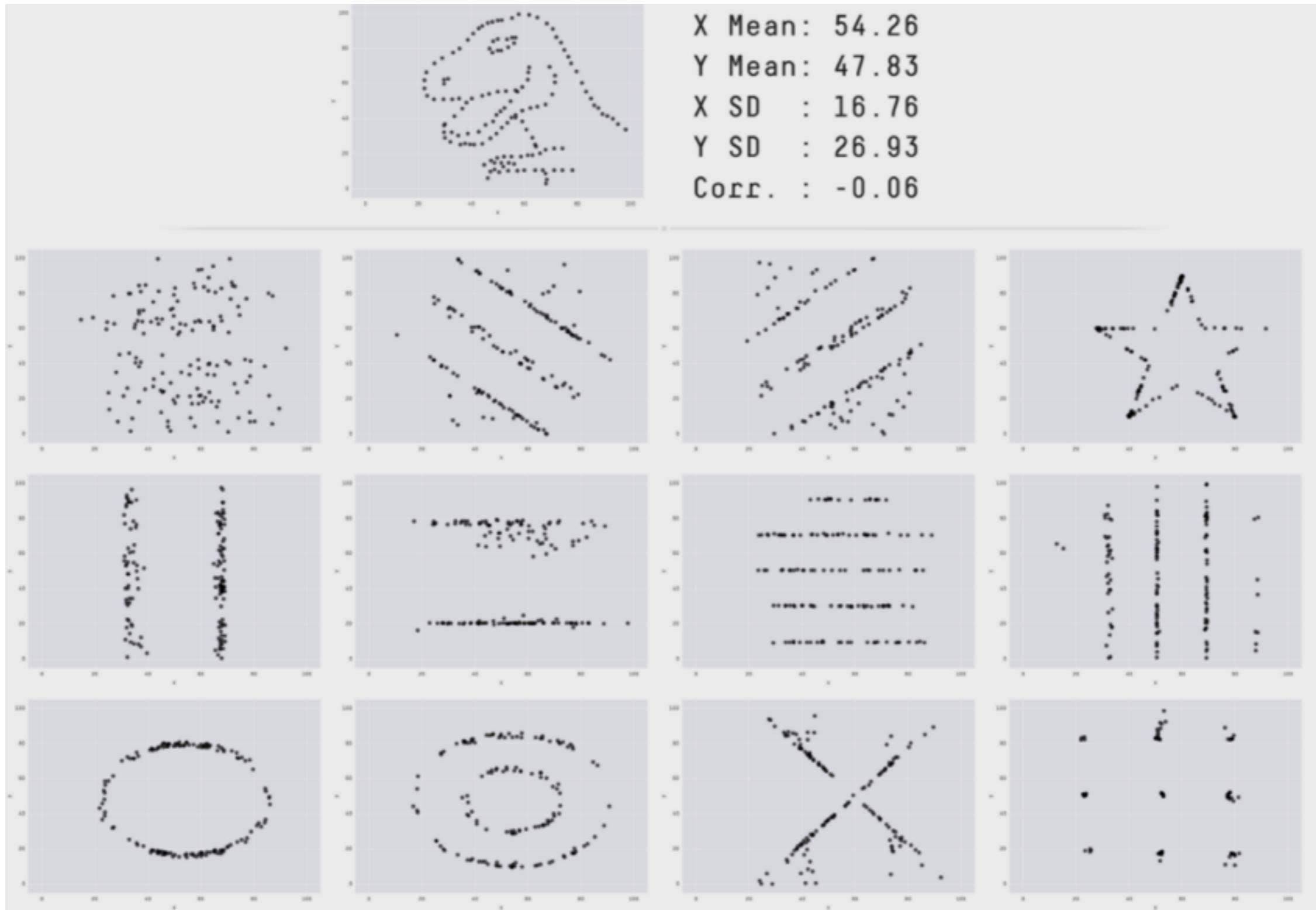
Property	Value
Mean of x	9
Sample variance of x	11
Mean of y	7.50
Sample variance of y	4.125
Correlation between x and y	0.816
Linear regression line	$y = 3.00 + 0.500x$



Anscombe quartet: all datasets have same means, variances, correlations and R^2

https://en.wikipedia.org/wiki/Anscombe%27s_quartet

DATASAURUS

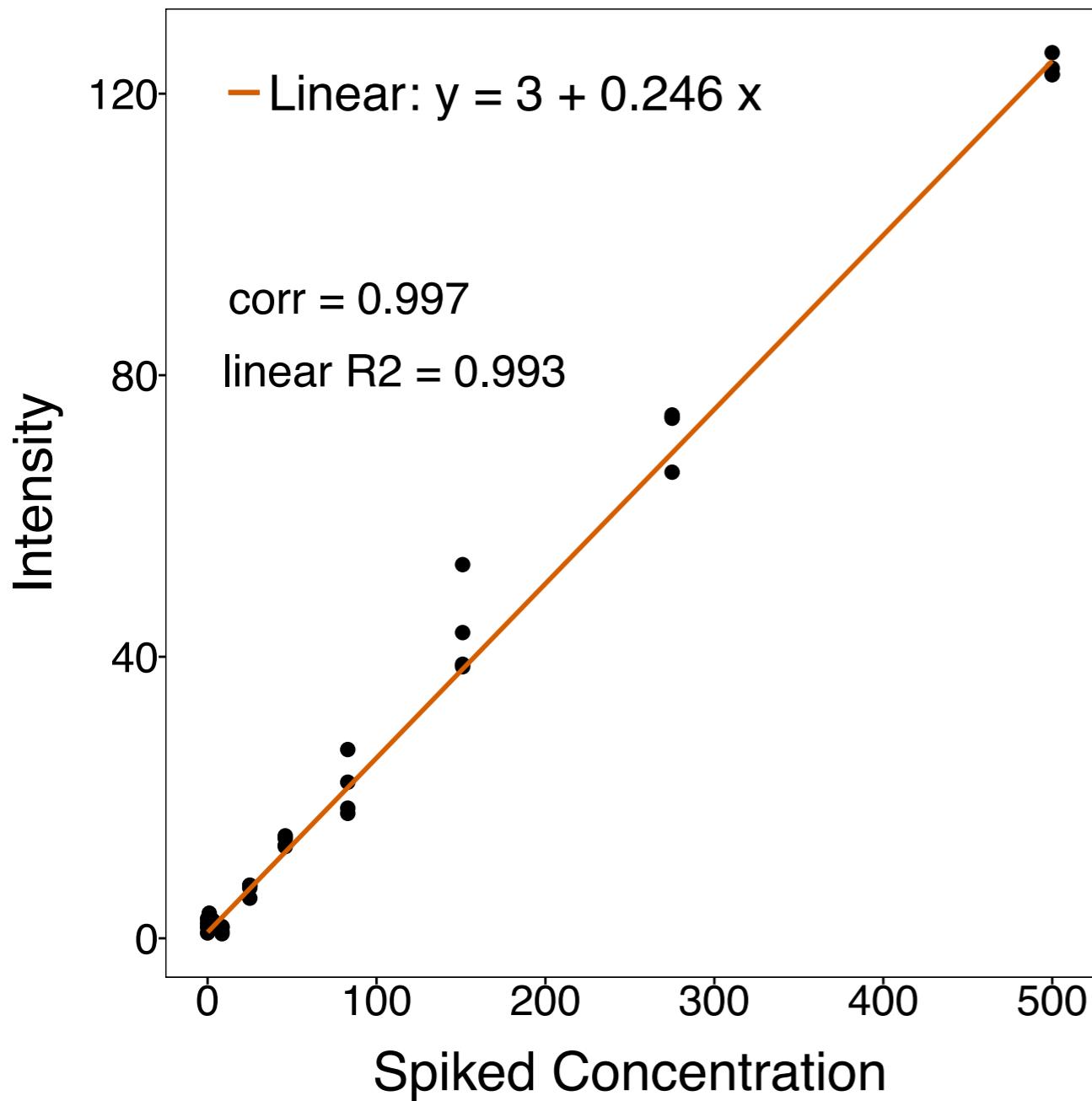


All datasets have same means, variances, correlations and R^2

<https://www.autodeskresearch.com/publications/samestats>

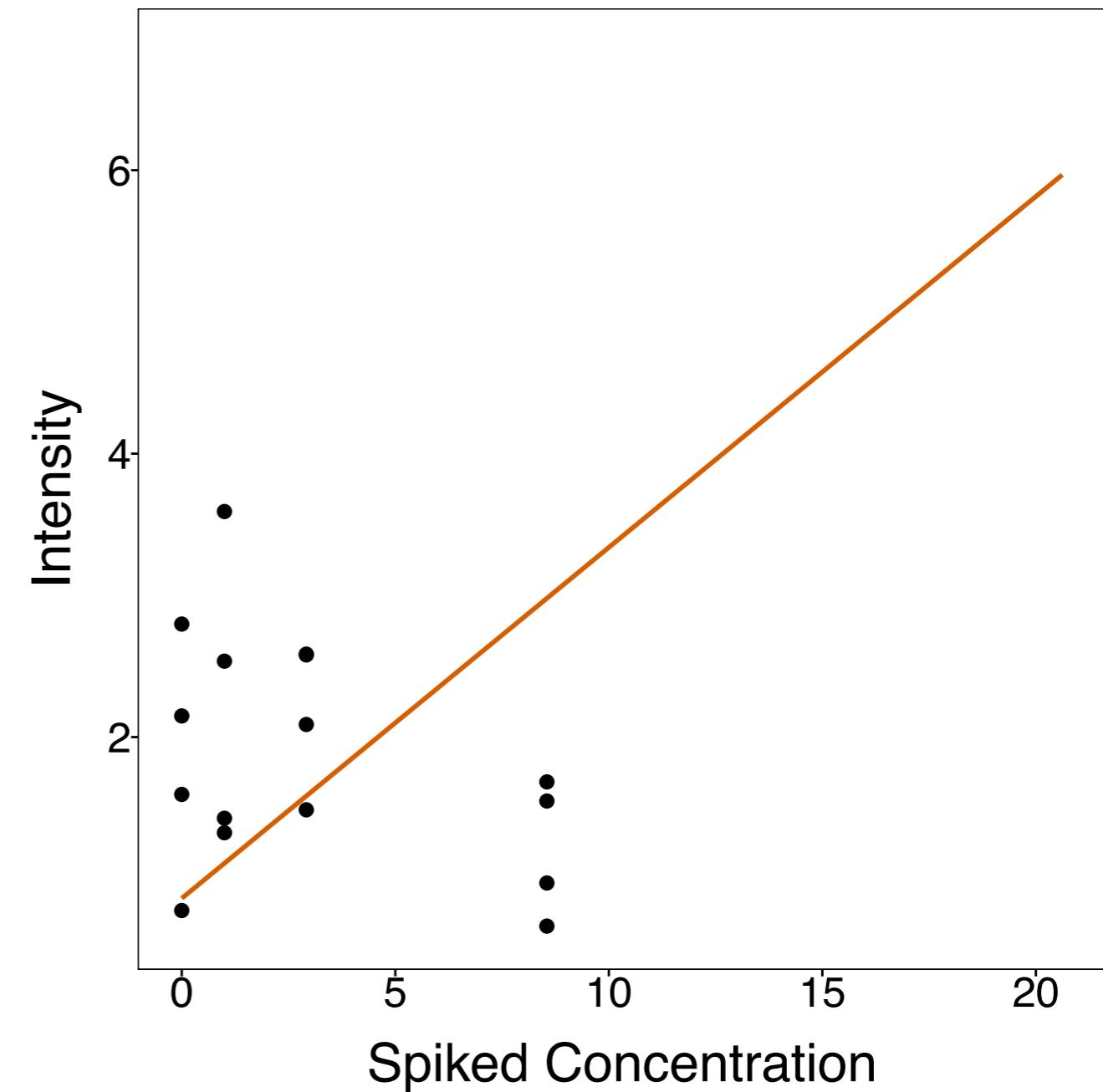
MORE EXAMPLES

High R² does not always mean good fit



Zoom out

Calibration experiment, CPTAC

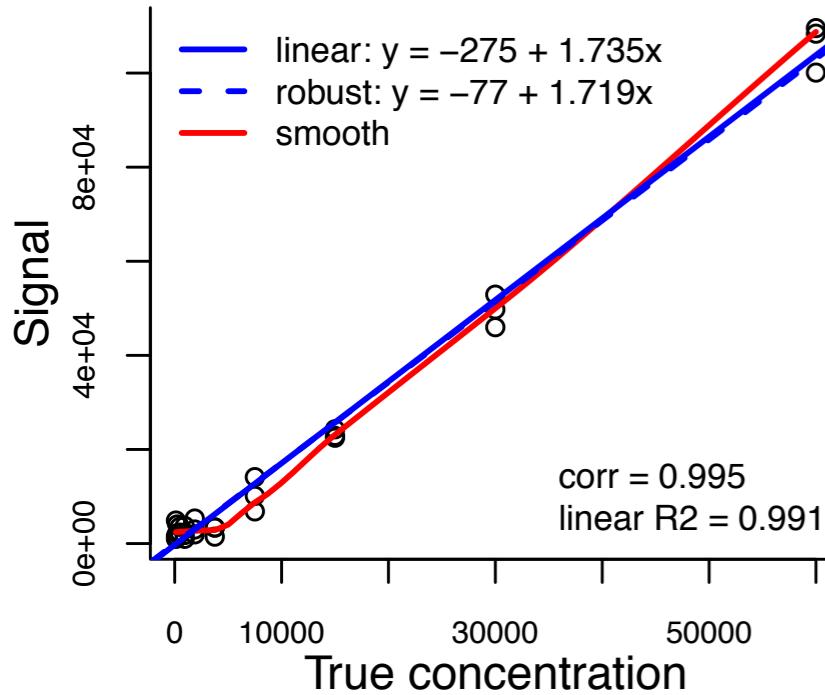


Zoom in

PROBLEM 2

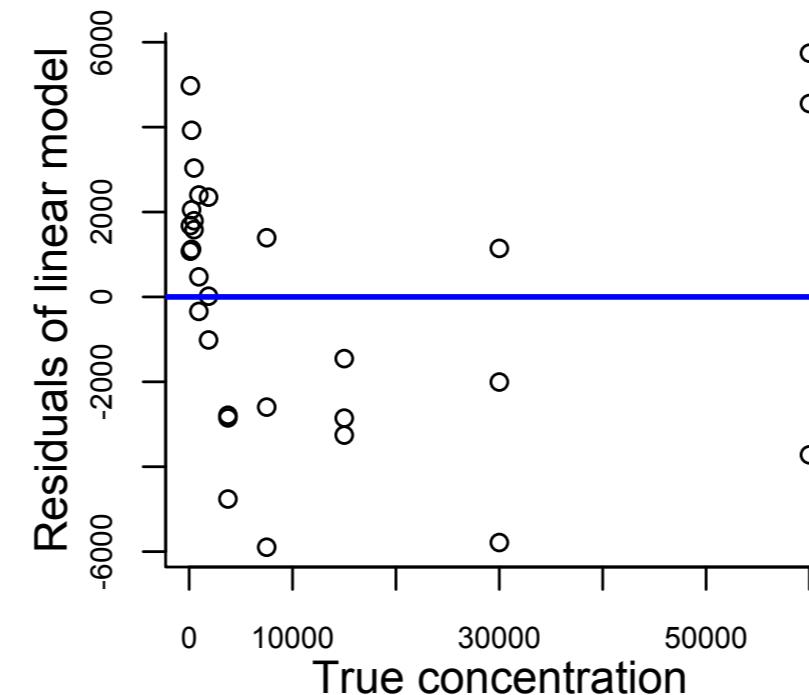
Unequal leverage, unequal variance

All concentrations



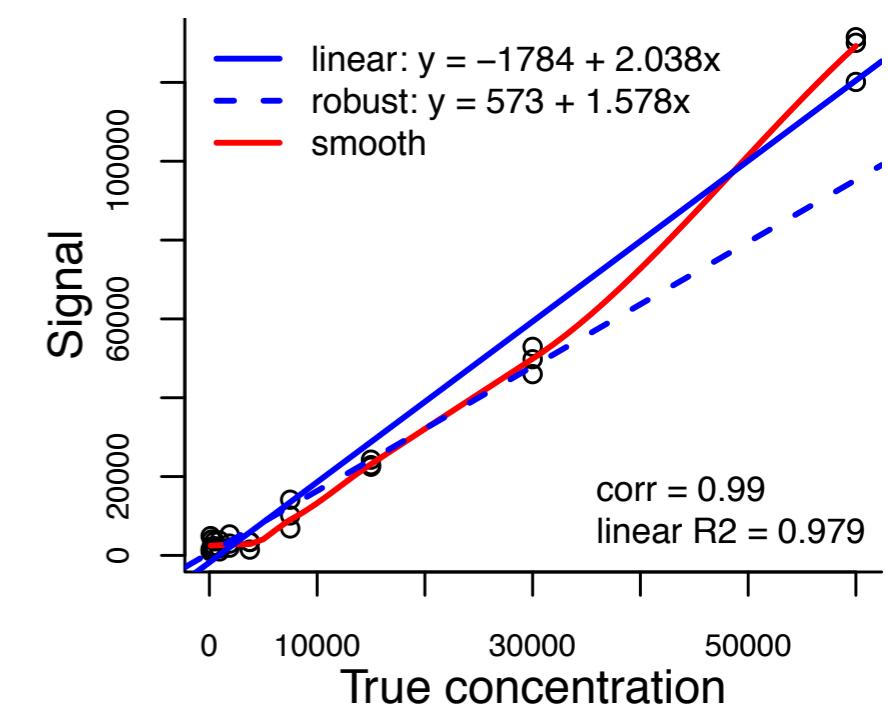
High concentrations are influential

Residuals of the linear model fit



High concentrations have more variance
Low concentrations have worse fit

Intensities at max concentration up by 20%



Moving the intensities produces a similarly good fit

The fit is unduly affected by highly variable values
Poor fit at low concentrations is unnoticed

CORRELATION

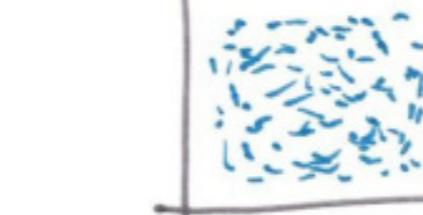
Correlation
Coefficient

Try our energy drink —
it's highly correlated with
performance!



athletic
performance

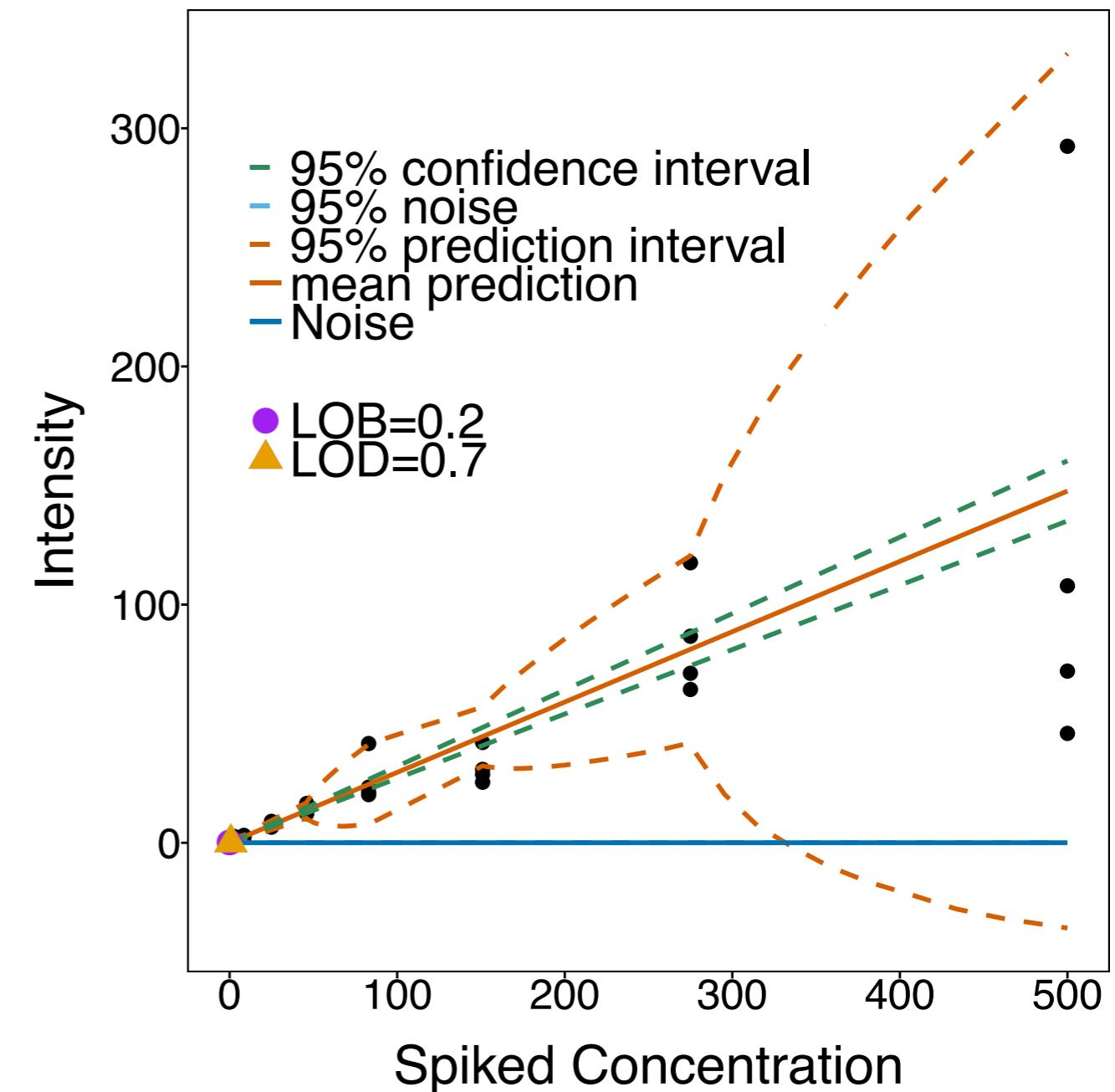
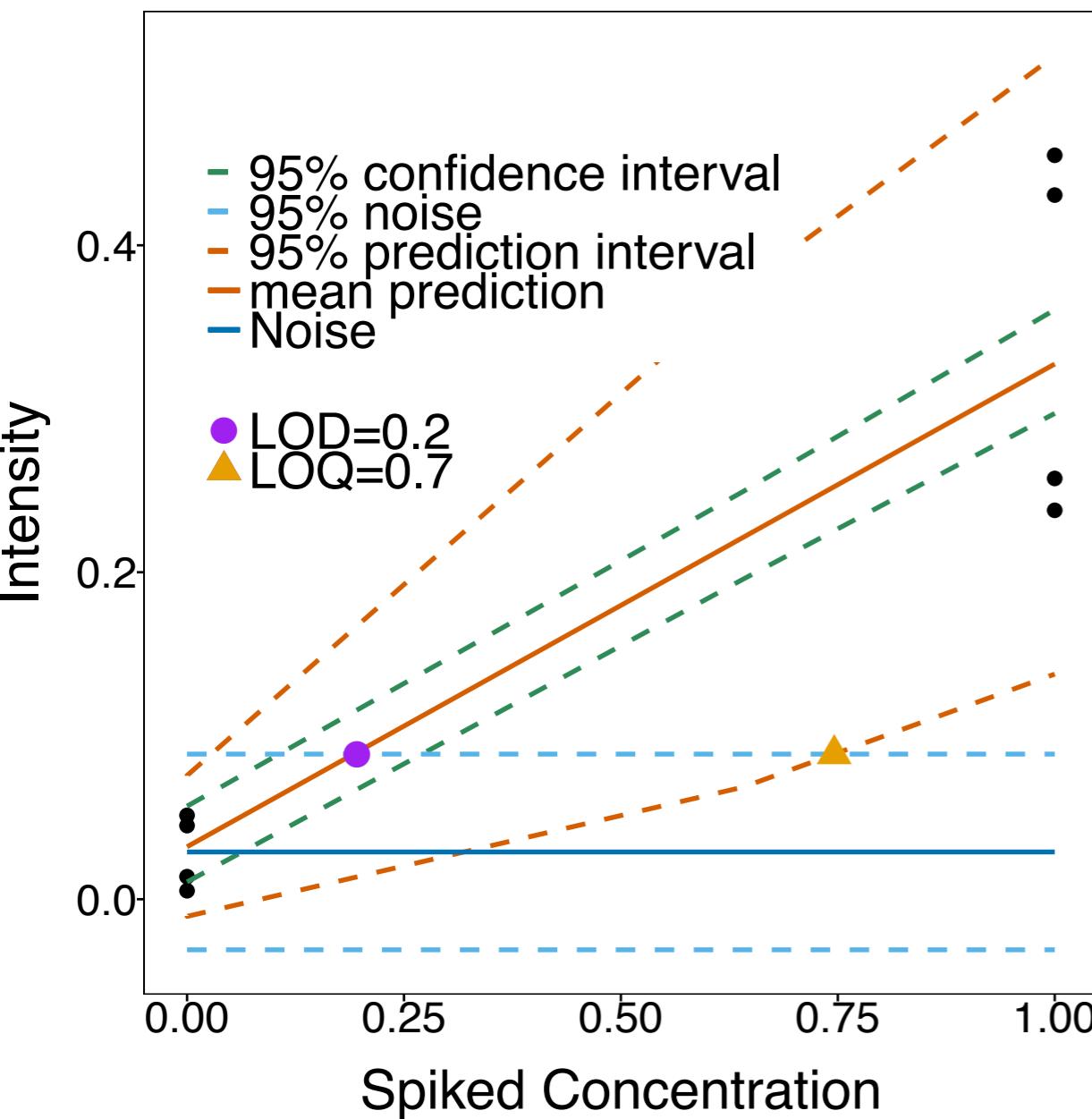
professional athletes we
paid to guzzle the stuff



amount of
drink consumed

CONFIDENCE VS PREDICTION

Prediction intervals should be used for determining LoD and LoQ



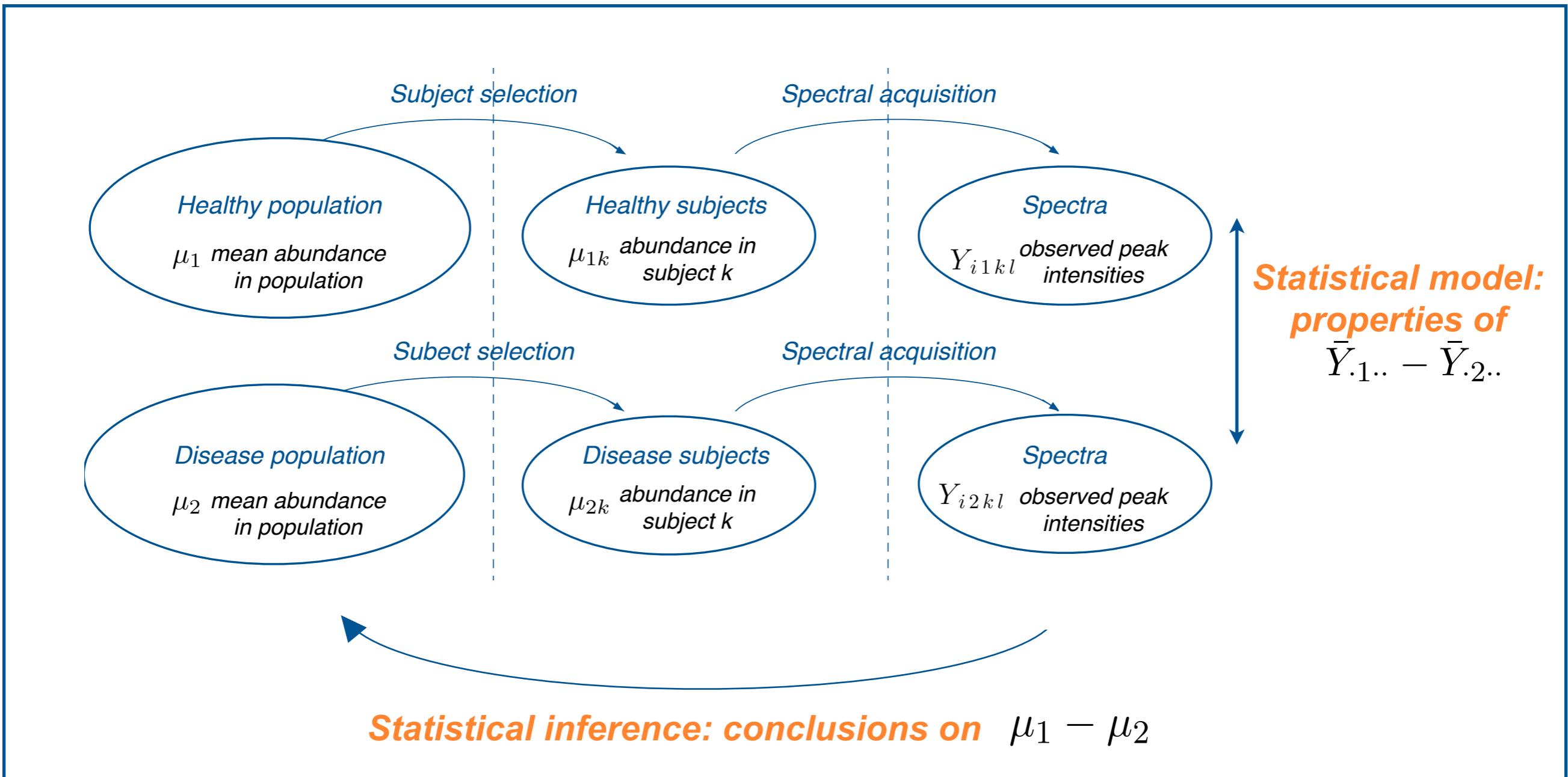
Confidence interval: $\left[\hat{y} - t_{\alpha/1}^{n-2} \cdot \sqrt{\text{var}\{\hat{y}\}}, \hat{y} + t_{\alpha/1}^{n-2} \cdot \sqrt{\text{var}\{\hat{y}\}} \right]$

Prediction interval: $\left[\hat{y} - t_{\alpha/1}^{n-2} \cdot \sqrt{\text{var}\{\hat{y}\} + s^2}, \hat{y} + t_{\alpha/1}^{n-2} \cdot \sqrt{\text{var}\{\hat{y}\} + s^2} \right]$

OUTLINE

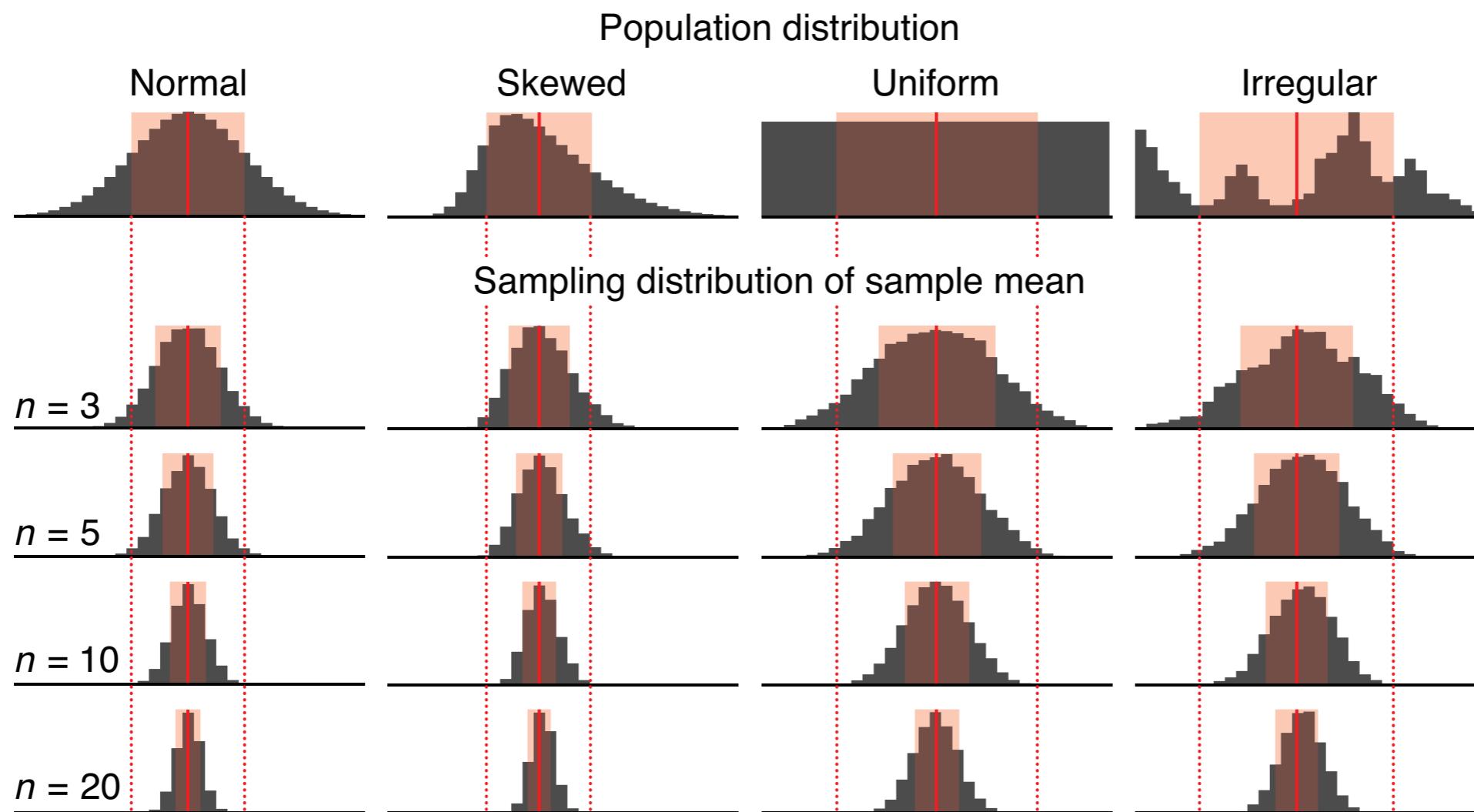
- So how many replicates do I need?
 - Design of complex experiments
- Associations between two measurements
 - Correlation and linear regression
- Statistical analysis of count data
 - Comparing proportions and spectral counts

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

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

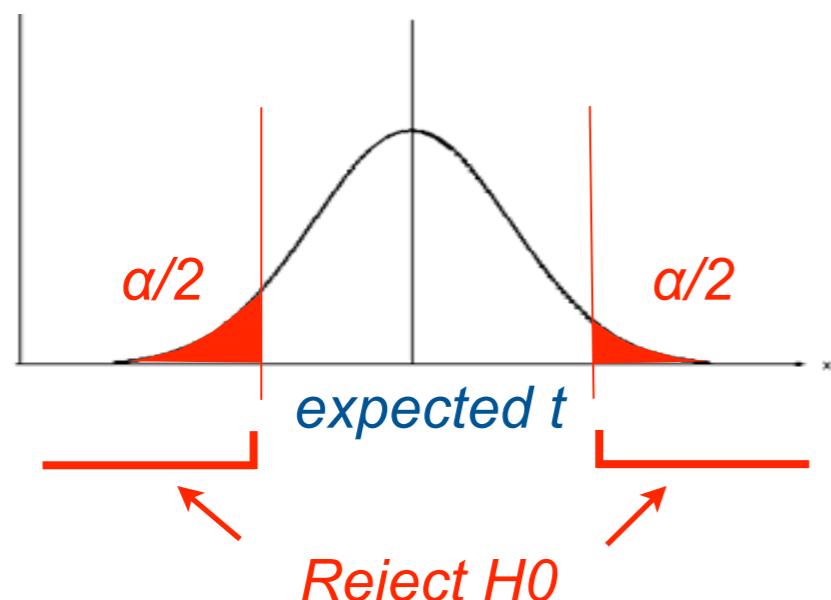
$H_a: \text{change in abundance, } \pi_1 - \pi_2 \neq 0$

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
 $\sim \text{Normal}(0, 1)$

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

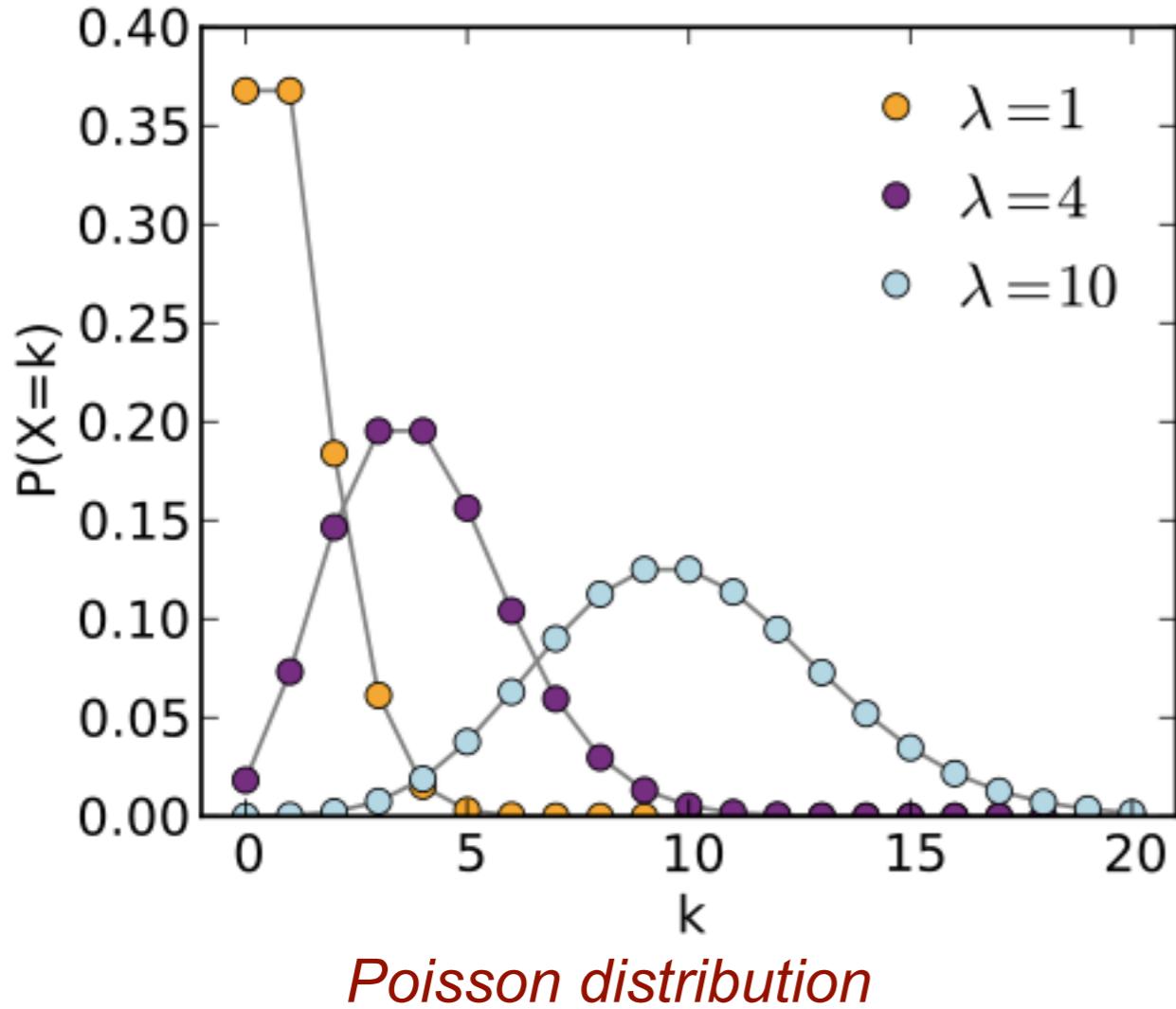


Confidence interval for difference of proportions

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

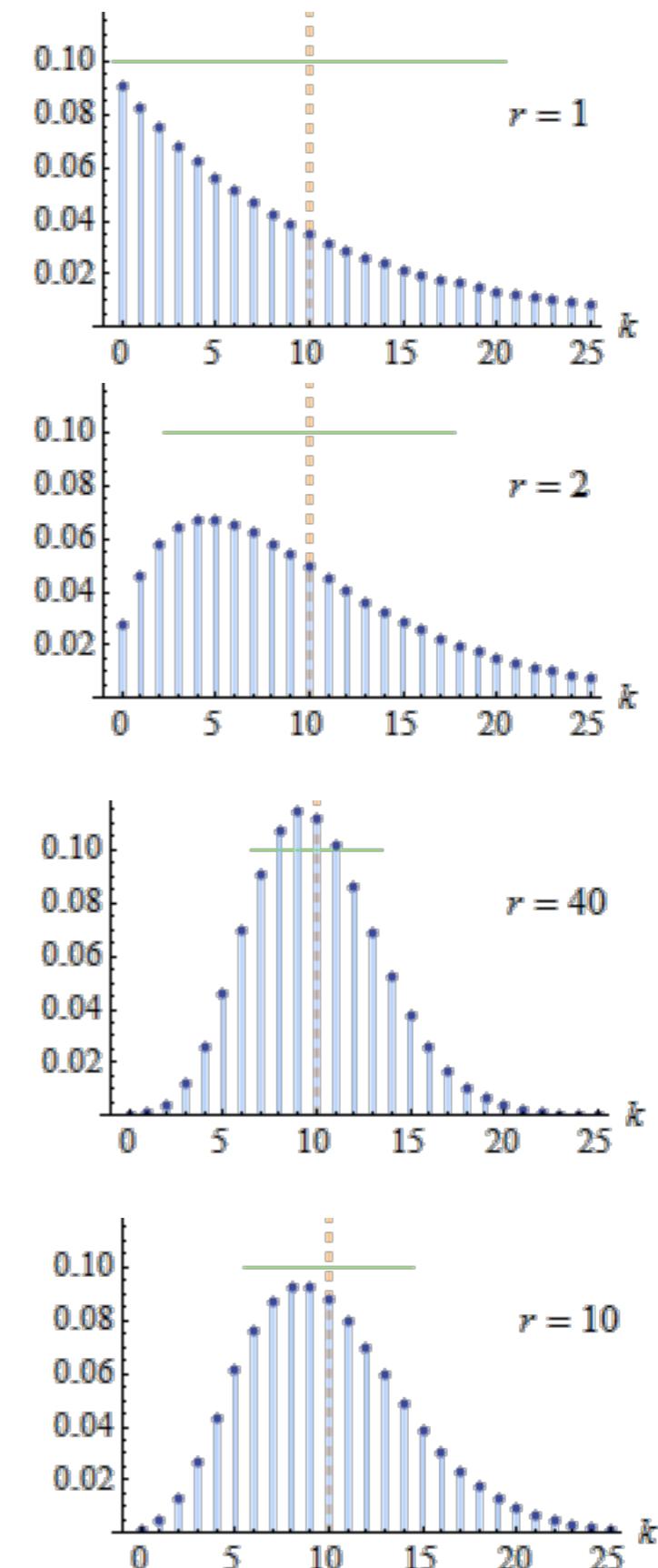
CATEGORICAL MEASUREMENTS

Comparing counts of MS/MS spectra



Poisson distribution

Negative Binomial distribution



https://en.wikipedia.org/wiki/Poisson_distribution

https://en.wikipedia.org/wiki/Negative_binomial_distribution

Northeastern University

THE COLLEGE OF SCIENCE AND
THE BARNETT INSTITUTE OF CHEMICAL AND BIOLOGICAL
ANALYSIS

PRESENTS

The Barry L. Karger Medal Lectures and Award

MAY 15-16, 2017

CURRY STUDENT CENTER
NORTHEASTERN UNIVERSITY
346 HUNTINGTON AVENUE
BOSTON, MASSACHUSETTS

WITH DISTINGUISHED LECTURER

DR. RUDOLF AEBERSOLD

*Head of Department of Biology
Institute of Molecular Systems Biology
ETH Zurich, Zurich, Switzerland*



SCHEDULE OF EVENTS

MONDAY, MAY 15, 2017

3:30 P.M. LECTURE RECEPTION

4:00 P.M. KARGER MEDAL PUBLIC LECTURE
BALLROOM, CURRY STUDENT CENTER

"The Proteome in Context"

TUESDAY, MAY 16, 2017

9:30 A.M. CONTINENTAL BREAKFAST

10:00 A.M. KARGER MEDAL TECHNICAL LECTURE
MCLEOD SUITES, CURRY STUDENT CENTER

*"SWATH-MS: Principles, current
state and new developments"*