# 1. Introduction/Overview

Readings:        Benjamin et al (2017)
                 Goodman (1999)

R:               Lab 0 and Lab 1

## Overview

- Statistics as a discipline arose from the need to use data to answer scientific questions in the face of uncertainty

- Statistical concepts are at the heart of scientific inquiry in the health sciences

- Your mastery of fundamental concepts will facilitate:
    o a better understanding of published research
    o a better understanding of how to structure effective scientific research
    o interpretation and presentation of results
    o collaboration with other biostatisticians and scientific investigators
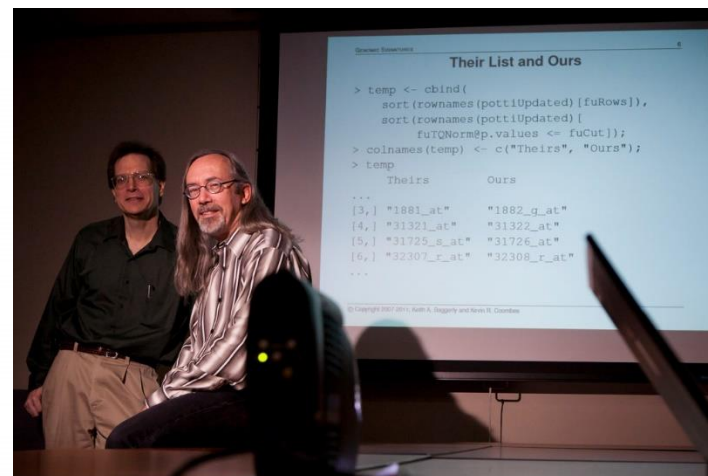
**Compelling example that underscores the need for careful design, data handling, and analysis:**

## DERIVING CHEMOSENSITIVITY FROM CELL LINES: FORENSIC BIOINFORMATICS AND REPRODUCIBLE RESEARCH IN HIGH-THROUGHPUT BIOLOGY

By Keith A. Baggerly[1] and Kevin R. Coombes[2]

University of Texas

*Source: New York Times*

# THE CANCER LETTER

Inside information on cancer research and drug development

*Nov. 19, 2010*

**JCO Retracts Key Duke Genomics Paper; Duke Shuts Down Three Phase II Trials;**

**Anil Potti Resigns From Duke University**

*Oct. 18, 2013*

**NCI Sets Rules For Omics Studies**

**Exit Joseph Nevins: Duke's Genomics Luminary Quietly Leaves**

# Algorithms vs. Inference
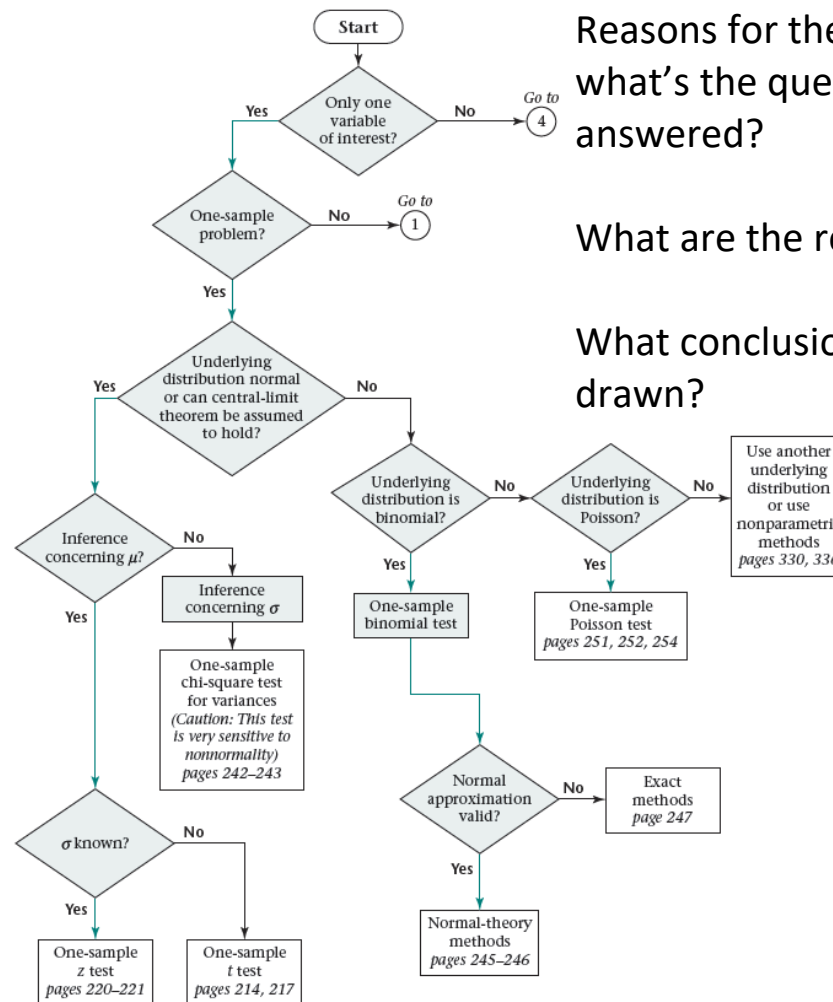
## Types of algorithms
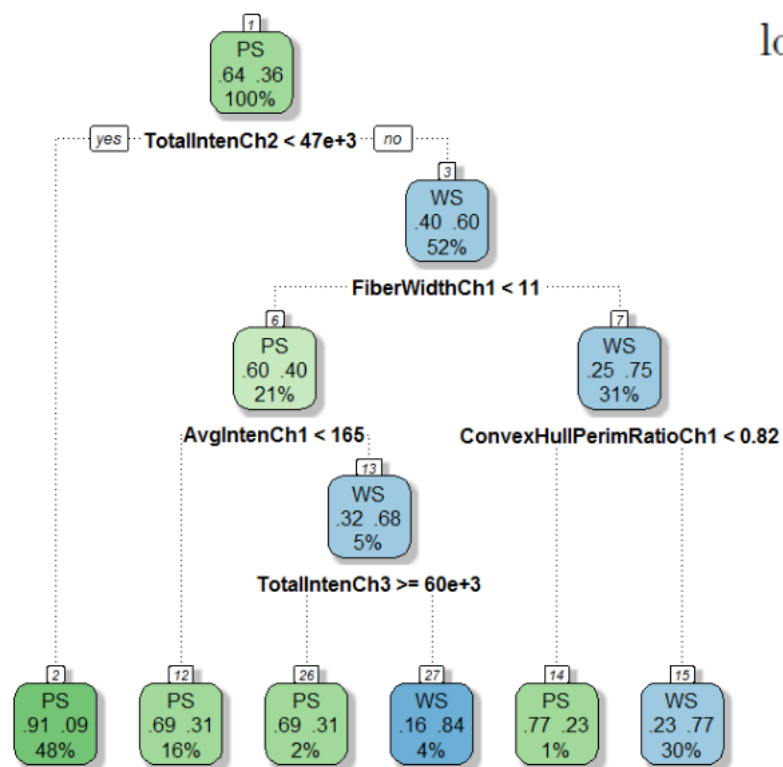


## Inference

**Inference**
Reasons for the algorithm – what's the question being answered?

What are the results?

What conclusions can be drawn?

# Algorithms vs. Models



Rattle 2013-Jun-19 15:43:30 Joe.Rickert

$$logit(\mathbb{E}[Y_i|x_{1,i}, ..., x_{m,i}]) = logit(p_i)$$

$$= log\left(\frac{p_i}{1 - p_i}\right)$$

$$= \beta_0 + \beta_1 x_{1,i} + \cdots + \beta_m x_{m,i}$$

$$logit(\mathbb{E}[Y_i|\mathbf{X}_i]) = logit(p_i)$$

$$= log\left(\frac{p_i}{1 - p_i}\right)$$

$$= \beta\mathbf{X}_i$$

# Randomness

Randomness is fundamental to statistical inference

Control of randomness is key to experimentation and the scientific method

# Random Number Generation

True Random Number Generator (TRNG) vs. Pseudo Random Number Generator (PRNG)

| Characteristic | PRNG | TRNG |
|---|---|---|
| Efficiency | Excellent | Poor |
| Determinism | Deterministic | Nondeterministic |
| Periodicity | Periodic | Aperiodic |

| Application | Most Suitable Generator |
|---|---|
| Lotteries and Draws | TRNG |
| Games and Gambling | TRNG |
| Random Sampling (e.g., drug screening) | TRNG |
| Simulation and Modelling | PRNG |
| Security (e.g., generation of data encryption keys) | TRNG |
| The Arts | Varies |

Source: www.random.org/randomness

# Applications of random number generation

## A.    Research Design

Randomization/random allocation, random sampling

1. **Random Digits:**  (Combinations of ) 0, 1, …, 9 occur with the same relative frequency (uniformly distributed).
   a. Digits occur independently of occurrence of other digits
   b. Use these sequences for random selection or allocation
   c. For example, a random number table

2. **Random Selection (Sampling):**  Selecting random portion of large population (e.g., select 10 units randomly from 1000)
   a. Tools in R include the "runif" and "sample" functions
   b. Simple random sampling, cluster sampling, stratified sampling

3. **Random Allocation:**  Assigning treatments randomly to individual units or groups of units
   a. Tools in R include the "blockTools" package

# Applications of random number generation

**B.     Random Sampling from Theoretical Distributions**

Also known as *simulation*

**Overview**

- Simulation is a fundamental and powerful tool in statistical practice

- Simulation for understanding

- Simulation for experimentation

- https://www4.stat.ncsu.edu/~davidian/st810a/simulation_handout.pdf

## IRReproducibility

**Reliance on p-values** –
- long-term trends leading up to the data + evidence provided by data; conflation of these has resulted in current approach which generally satisfies neither (Goodman, 1999)
- $< 0.05$ actually provides weak evidence against $H_0$ (Benjamin et al. 2017)

**Multiplicity** – leads to *selection* of results and this has an *effect* on observed ability to replicate results (see 2005 *JAMA* paper by Ioannidis on course website)
- Multiple variables, endpoints, time points, subgroups, comparisons
- Multiple hypothesis testing, multiple looks at the data
- Multiple models and adjustments
- Fishing expeditions, mountains of output without *a priori* thought and justification - *…what exactly <u>was</u> my (their) research hypothesis or question?…p-hacking …*

**Publication bias** – file drawer problem, publish or perish, sensationalism

**Cognitive biases** – preconceived notions about what effects are real, what effects could be real, what effects are likely not to be real

# Reproducible Research

**Reproducible research** is the idea that data analyses, and more generally, scientific claims, are published with their data and software code so that others may verify the findings and build upon them.

Reproducible Research | Coursera
https://www.coursera.org/learn/reproducible-research

# Open Science Paradigm

**Open Access** - Publications

    Public access – free to anyone with an internet connection

    Free: Use, reuse, remix

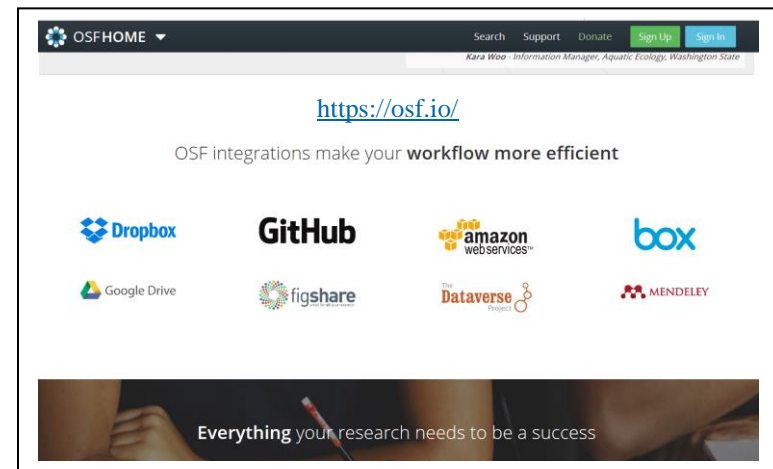    Pre-study: Study registration

    Post-study: Preprints

**Open Data**

    Known Provenance

    Confidentiality assured

    Portability (interoperability) built in

    Excel-free, reproducible data manipulation/management best practices applied

**Open Code**

    Version control used to track changes

    Collaborative model – team science

    Crowd sourcing – best solutions openly available
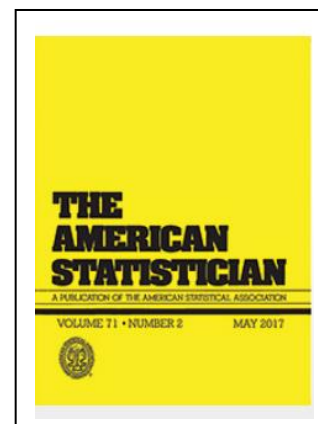
## Miscellaneous

Throughout semester we'll:

### Read

The statistical methods literature



**Practice** Reproducible Research Principles
"Lite" version …move towards compliance

### Appreciate

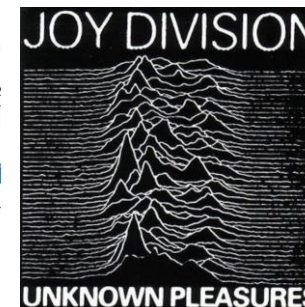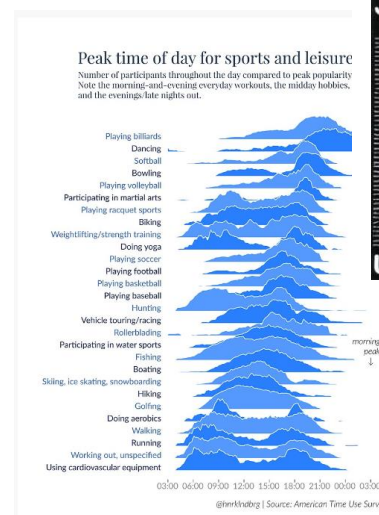("Peak") science writing – WIRED, Quanta, your favorites

**Be on the lookout** for cool, innovative graphics

ggjoy – R package
https://cran.r-project.org/web/packages/ggjoy/vignettes/introduction.html





R Gallery:
https://cran.r-project.org/web/packages/ggjoy/vignettes/gallery.html

### Share

Resources …using Canvas