

# Introduction to Experimental Design at CRUK-CI

Chandra Chilamakuri

James Hadfield

Jing Su

Rory Stark

[tinyurl.com/cruk-edesign](http://tinyurl.com/cruk-edesign)

# Agenda

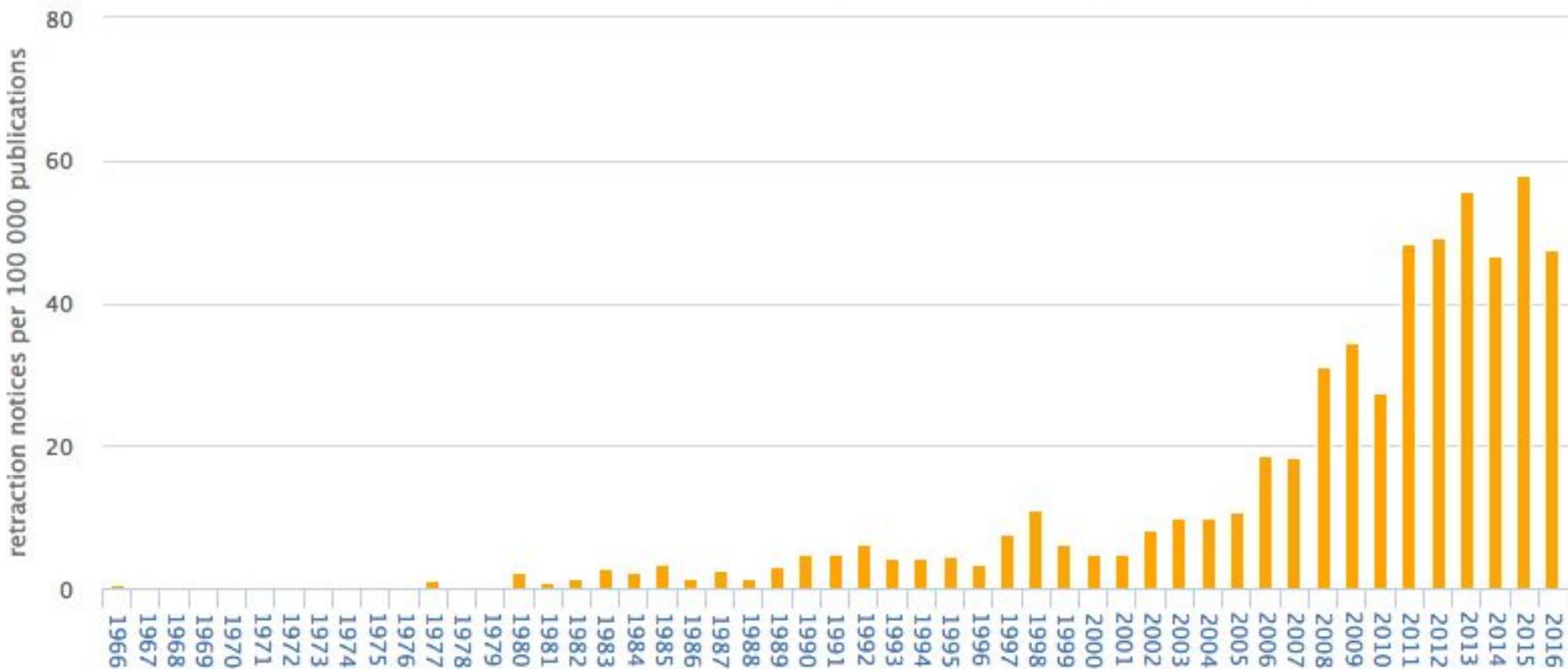
- Why perform experiments?
- Why think about experimental design?
- What makes for a well designed experiment?
- Aspects of experimental design
  - Experimental variables
  - Power: variance and replicates
  - Bias: confounding factors, randomisation, and controls
- Experimental design types
- Experimental design at CRUK-CI

# Why Perform Experiments?

# Why Think About Experimental Design?

# *Crises in Reproducible Research!!*

Retraction notices per 100 000 publications by year of Entrez record creation



# 47 of 53 high-profile cancer studies were not reproducible!

The screenshot shows the header of the Nature journal website with the title "nature International weekly journal of science". Below the header is a navigation bar with links: Home, News & Comment, Research, Careers & Jobs, Current Issue, Archive, Audio & Video, and a partially visible link starting with 'F'. Below the navigation bar is a breadcrumb trail: Archive > Volume 483 > Issue 7391 > Comment > Article. The main content area features the title "NATURE | COMMENT" and the article title "Drug development: Raise standards for preclinical cancer research" by C. Glenn Begley & Lee M. Ellis. To the right of the article title are three sharing icons: a person icon, an envelope icon, and a square icon.

## Drug development: Raise standards for preclinical cancer research

C. Glenn Begley & Lee M. Ellis

Affiliations | Corresponding author

*Nature* 483, 531–533 (29 March 2012) | doi:10.1038/483531a

Published online 28 March 2012

# Consequences of Poor Experimental Design...

- **Cost** of experimentation. We have a responsibility to CRUK donors!
- **Limited & Precious** material, esp. clinical samples.
- **Immortalization** of data sets in public databases and methods in the literature. Our bad science begets more bad science.
- **Ethical concerns** of experimentation: animals and clinical samples.

# A Well-Designed Experiment

Should have

- Clear objectives
- Focus and simplicity
- Sufficient power
- Randomised comparisons

And be

- Precise
- Unbiased
- Amenable to statistical analysis
- Reproducible

# Ronald A. Fisher(1890-1962)



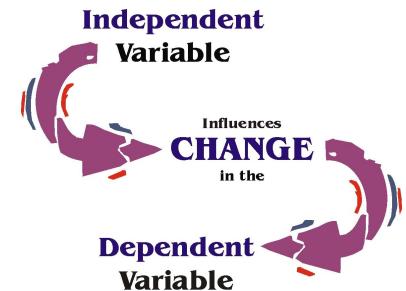
*“To consult the statistician **after** an experiment is finished is often merely to ask him to conduct a **post mortem** examination. He can perhaps say what the experiment died of.” (1938)*

# Aspects of Experimental Design

- Experimental Factors
- Power
  - Sources of Variance
  - Replicates
- Minimising Bias
  - Confounding factors
  - Randomisation wherever a decision is to be made
    - Controls for both measured and unmeasured factors
  - Controls

# Experimental Factors

- Factors: aspects of experiment that change and **influence the outcome** of the experiment
  - e.g. time, weight, drug, gender, ethnicity, country, plate, cage etc.
- Variable type depends on type of measurement:
  - Categorical (**nominal**), e.g. gender
  - Categorical with ordering (**ordinal**), e.g. tumour grade
  - **Discrete**, e.g. shoe size, number of cells
  - **Continuous**, e.g. body weight in kg, height in cm
- Independent and Dependent variables
  - Independent variable (IV): what you change
  - Dependent variable (DV): what changes due to IV
  - “**If (independent** variable), **then (dependent** variable)”

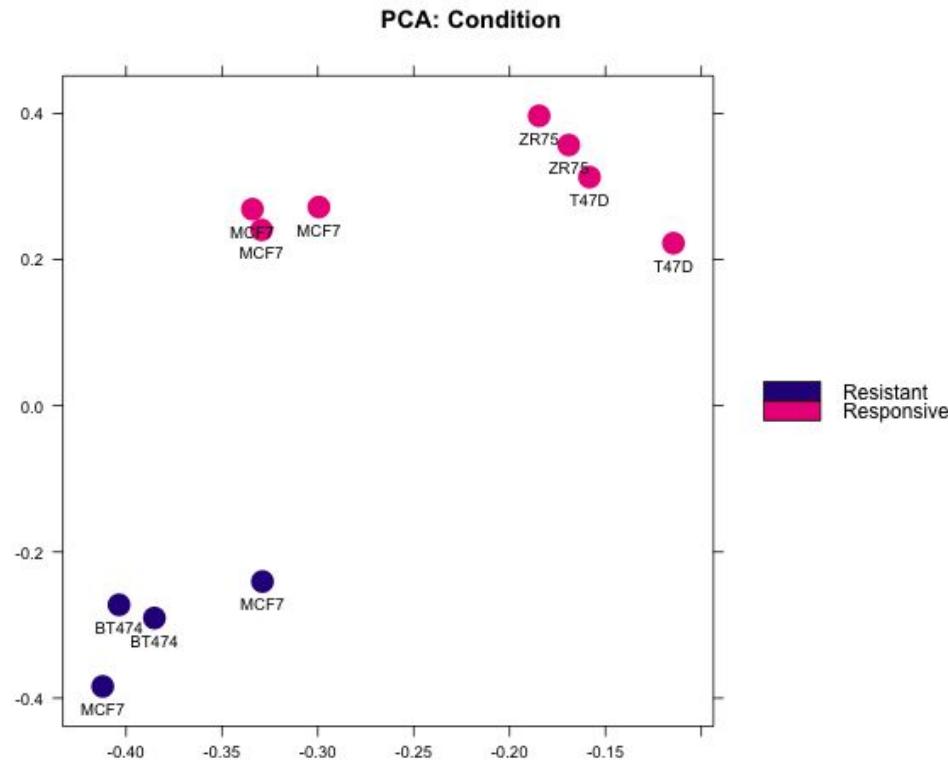


# Sources of Variation

- Biological “noise”
  - Biological processes are inherently stochastic
  - Single cells, cell populations, individuals, organs, species....
  - Timepoints, cell cycle, synchronized vs. unsynchronized
- Technical noise
  - Reagents, antibodies, temperatures, pollution
  - Platforms, runs, operators
- Consider in advance and control

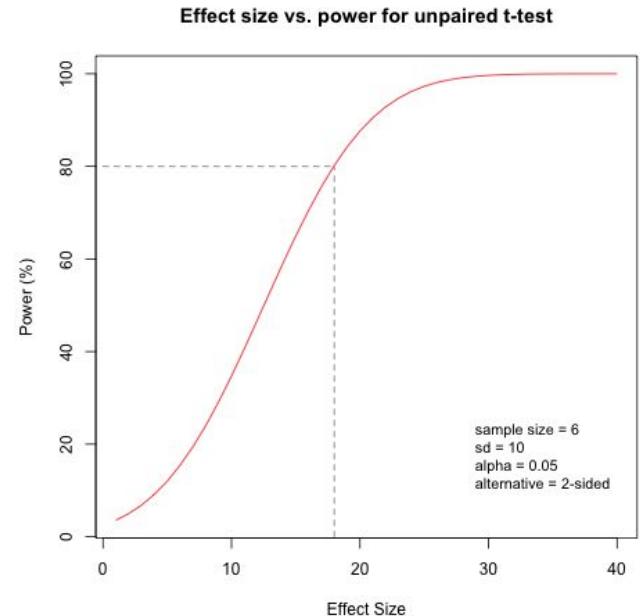
# Types of Replication

- Biological replication:
  - *In vivo*:
    - Patients
    - Mice
  - *In vitro*:
    - Different cell lines
    - Re-growing cells (passages)
- Technical replication:
  - Experimental protocol
  - Measurement platform (i.e. sequencer)



# How many samples?

- Why do you need replicates?
- Calculating appropriate sample sizes
  - Power calculations
  - Planning for precision
  - Resource equation
- Power: the **probability** of detecting an **effect** of a specified size if present.
  - Identify and control the **sources of variability**
    - Biological variability
    - Technical variability
  - Using **appropriate numbers** of samples (sample size/replicates)
  - Power calculations estimate sample size required to detect an effect *if degree of variability is known*
    - Depends on  $\delta$ , n, sd,  $\alpha$ ,  $H_A$
  - If adding samples increases variability, that alone won't add power!



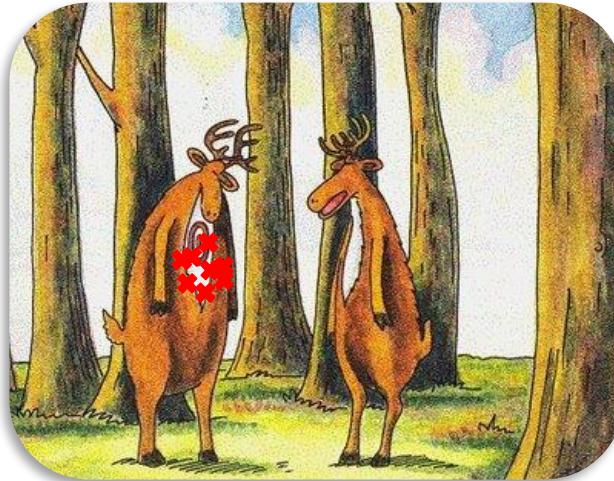
# Forms of Bias

Type of Bias	Description
Selection bias	Systematic differences between <b>baseline characteristics</b> or <b>treatment groups</b> that are being compared.
Performance bias	Systematic differences between groups in exposure to factors other than the interventions of interest (aka <b>confounding</b> or <b>extraneous factors</b> ).
Attrition bias	Systematic differences between groups due to samples being <b>withdrawn</b> from the study or <b>excluded</b> from the analyses.
Detection or Measurement bias	Systematic differences between groups in how outcomes are assessed or determined, e.g. <b>measurement errors</b> and inefficient use of data.
Reporting bias	Systematic differences between reported and unreported findings due to manipulation in the reporting of findings such as <b>selective or distorted reporting</b> , e.g. papers with more 'interesting results' are more likely to be submitted and accepted for publication.

# Precision, Accuracy & Bias

Precise

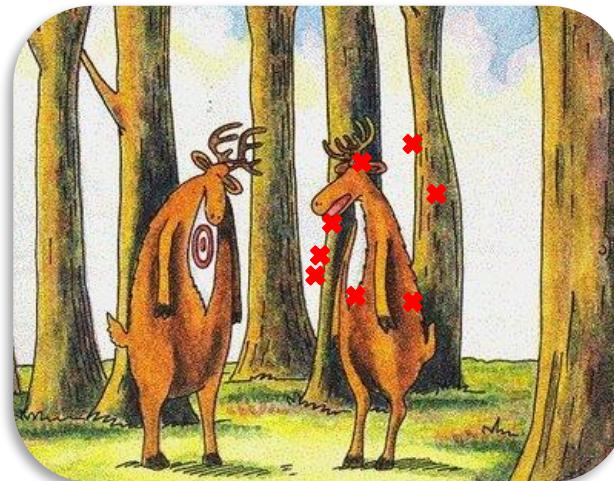
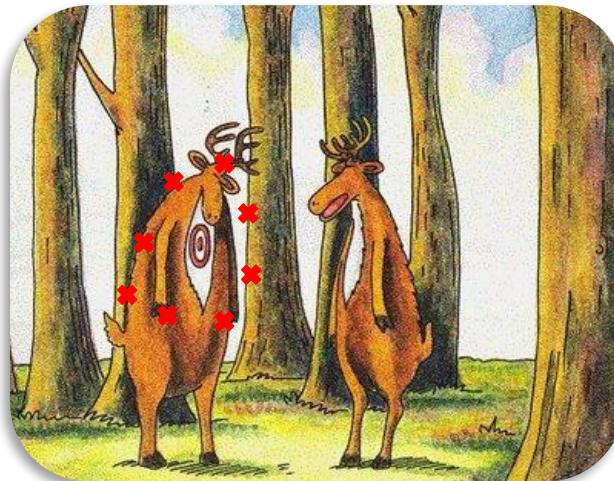
Accurate



Biased

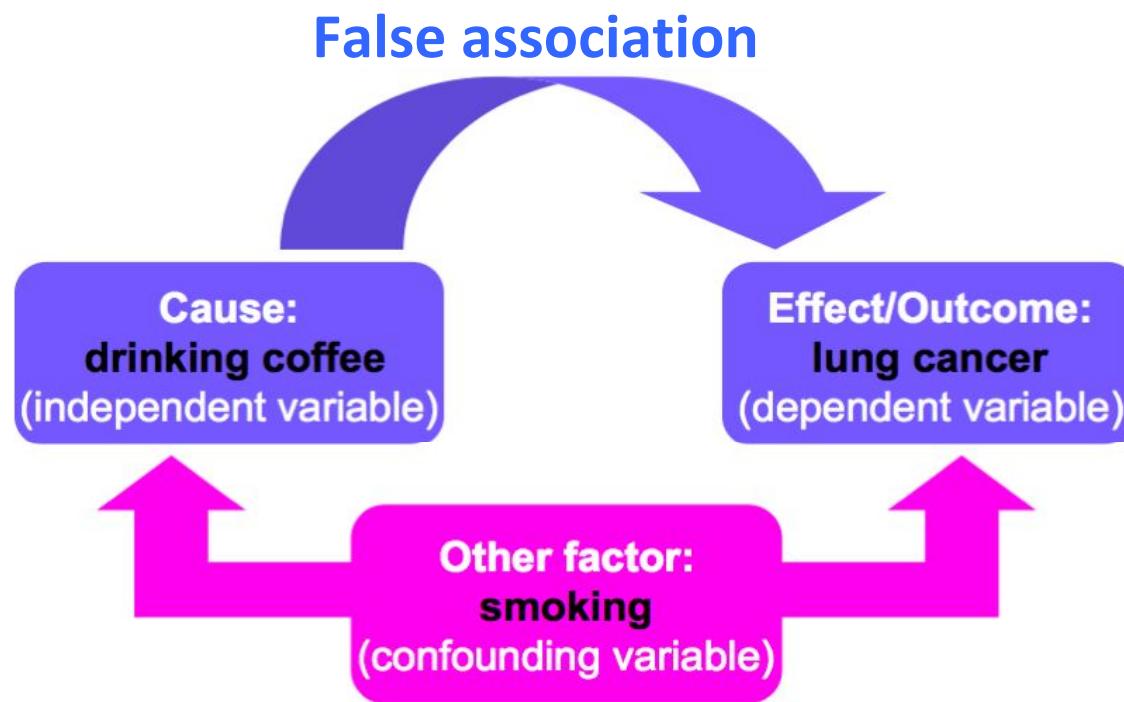


Imprecise



# Confounding Factors

- Also known as **extraneous, hidden, lurking or masking factors**, or the **third variable or mediator variable**.
- May mask an actual association or **falsely** demonstrate an apparent association between the independent & dependent variables.
- Hypothetical Example would be a study of coffee drinking and lung cancer.



# Confounding Factors

- Other examples:
  - Democrats were less satisfied with their sex lives than Republicans.  
(ABC poll report).
  - Slightly overweight people live longer than thin people  
(US Centre for Disease Control).
- **Inadequate management and monitoring** of confounding factors
  - one of the most common causes of researchers wrongly assuming that a correlation leads to a causality.
- If a study does not consider confounding factors,  
**don't believe it!**

### Genetic Signatures of Exceptional Longevity in Humans

Paola Sebastiani,<sup>1,\*</sup> Nadia Solovieff,<sup>1</sup> Annibale Puca,<sup>2</sup> Stephen W. Hartley,<sup>1</sup> Efthymia Melista,<sup>3</sup> Stacy Andersen,<sup>4</sup> Daniel A. Dworkis,<sup>3</sup> Jemma B. Wilk,<sup>5</sup> Richard H. Myers,<sup>5</sup> Martin H. Steinberg,<sup>6</sup> Monty Montano,<sup>3</sup> Clinton T. Baldwin,<sup>6,7</sup> Thomas T. Perls<sup>4,\*</sup>

<sup>1</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA. <sup>2</sup>IRCCS Multimedica, Milano, Italy; Istituto di Tecnologie Biomediche, Consiglio Nazionale delle Ricerche, Segrate, 20122, Italy. <sup>3</sup>Department of Medicine, Boston University School of Medicine, Boston, MA 02118, USA. <sup>4</sup>Section of Geriatrics, Department of Medicine, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>5</sup>Department of Neurology, Boston University School of Medicine, Boston, MA 02118, USA. <sup>6</sup>Departments of Medicine and Pediatrics, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>7</sup>Center for Human Genetics, Boston University School of Medicine, Boston, MA 02118, USA.

- GWAS study: 800 centenarians vs. controls
- Found 150 SNPs predicting centenarians with 77 % accuracy
- Problem: they used **different SNP chips** for centenarians and controls
- Retracted in 2011 following independent review and QC of data

# Technical Confounding Factors: Batch Effects



RNA Extraction

Day1, Plate 1

	1	2	3	4	5	6	7	8	9	10	11	12
A	O											
B												
C												
D												
E												
F												
G												
H												

Control

Day2, Plate 2

	1	2	3	4	5	6	7	8	9	10	11	12
A	O											
B												
C												
D												
E												
F												
G												
H												

Treatment 1

Day3, Plate 3

	1	2	3	4	5	6	7	8	9	10	11	12
A	O											
B												
C												
D												
E												
F												
G												
H												

Treatment 2

The difference between Control, Treatment 1  
and Treatment 2 is confounded by **day** and **plate**.

# Solutions

- Randomisation
  - Statistical analysis assume randomised comparisons
  - May not see issued caused by non-randomised comparisons
  - Make every decision random not arbitrary
- Blinding
  - Especially important where subjective measurements are taken
  - Every experiment should reach its potential degree of blinding

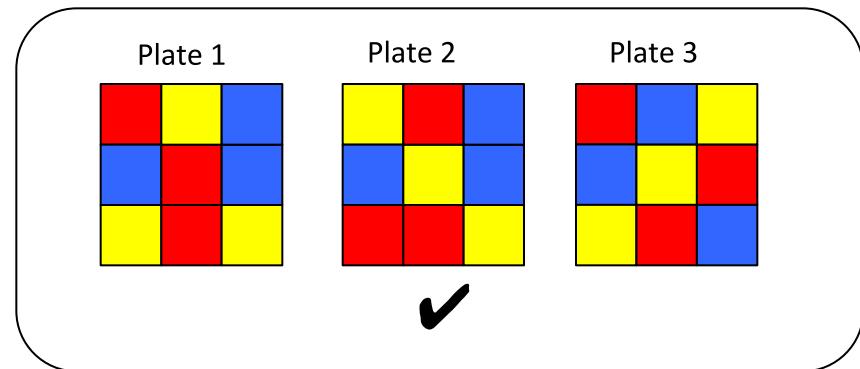
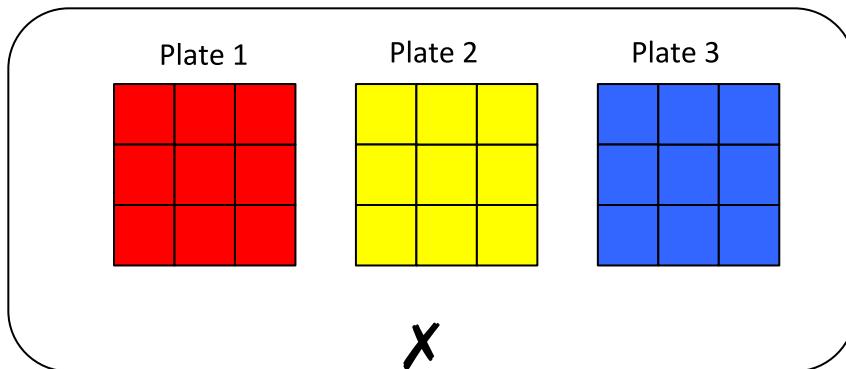
# Randomised Block Design

- **Blocking** is the arranging of *experimental units* in groups (blocks) that are similar to one another.

Control

Treatment 1

Treatment 2



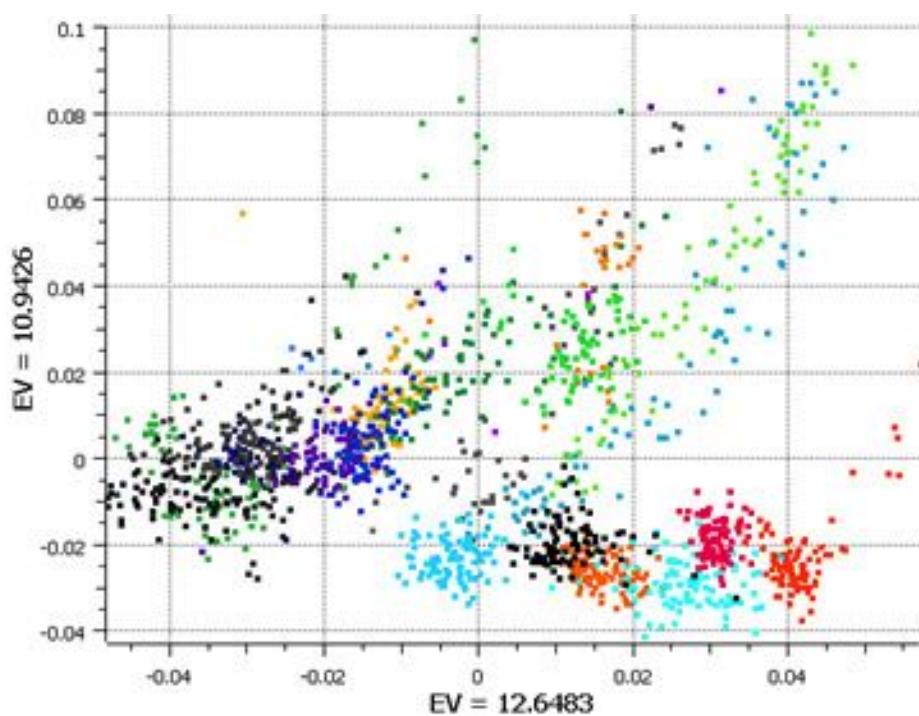
- RBD across plates so that each plate contains spatially randomised **equal proportions** of:
  - Control
  - Treatment 1
  - Treatment 2controlling plate effects.

# Randomised Block Design

**Good** design example: Alzheimer's study from GlaxoSmithKline

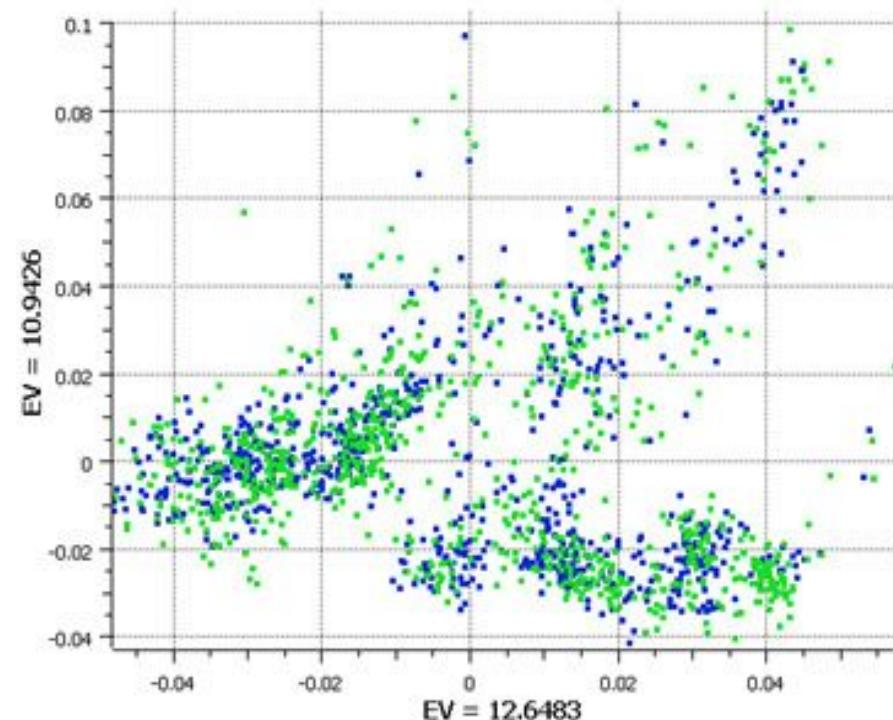
## Plate effects by *plate*

Left PCA plot show *large plate effects*.  
Each colour corresponds to a different plate



## Plate effects by *case/control*

Right PCA plot shows each plate cluster contains *equal proportions* of cases (blue) and controls (green).



# Experimental Controls

- Controlling errors
  - Type I: FP
    - Negative controls: should have minimal or no effect
  - Type II: FN
    - Positive controls: known effect
- Technical controls
  - Detect/correct technical biases
  - Normalise measurements (quantification)

# Examples of Experimental Controls

- Wild-type organism (knockouts)
- Inactive siRNA (silencing)
- Vehicle (treatments)
- Input: fragmented chromatin (ChIP)
- Spike-ins (quantification/normalisation)
- “Gold standard” datapoints
- Multi-level controls
  - e.g. contrast Vehicle/Input vs. Treatment/Input

# Design Issues: Sequencing Experiments

- Platforms (MiSeq, HiSeq, etc.)
- Library preps
- Multiplexing and pooling strategies
- Single-end vs paired end
- Sequencing depth
  - Coverage
  - Lanes
- Validation

# Types of Experimental Designs

- Full random: compute associations after
- Block designs: randomisation
- Matched: tumour/normal
- Factorial/multifactorial designs: GLM
- Time series
- Hierarchical designs

# CRI Experimental Design Meetings

- **Tuesday** 30 min slots (2:00-3:00pm) with Bioinformatics & Genomics Cores
- **Requirements:**
  - Email [CRIExperimentalDesign@cruk.cam.ac.uk](mailto:CRIExperimentalDesign@cruk.cam.ac.uk) to request meeting
  - Fill in Experimental Design Form and return 1 week prior to meeting
  - **Your attendance**
  - Provide **project background** (a few slides from you)
- **Discussion:**
  - Planning, time-scale, cost, aims, scope, questions
  - Choosing the correct technology
  - Technical issues e.g. what sequencing depth?
  - Sample collection and processing methods
  - Sample information (meta-data) collection
  - Randomisation, Blocking and Replication issues
  - Analyst?
  - Pilot study?
  - Effect size & Sample-size calculation?

# Experimental Design Guide

- <https://rawgit.com/bioinformatics-core-share-d-training/experimental-design/master/ExperimentalDesignManual.pdf>

[tinyurl.com/cruk-edesign](http://tinyurl.com/cruk-edesign)

# Practicals

1. Clinical: Identification of prognostic biomarkers in human prostate cancer patients (Rory)
2. Animal: Effects of mutant vs wildtype HHEX in liver and brain development (Jing)
3. ~~PROTAC and its effect on transcription factors in drug responsive and resistant breast cancer cell line (Chandu) (James)~~
4. ChIP-seq: Transcription factor binding divergence in mice ( James )