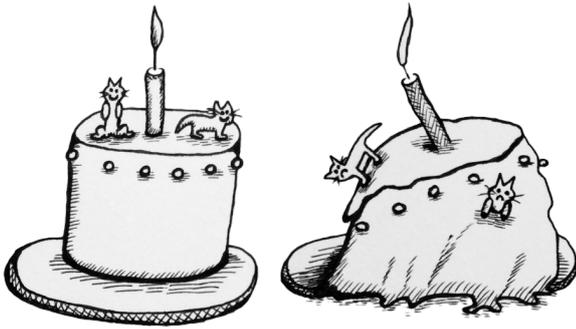


Controlling for Batch Effects

What are Batch Effects?

When cooking, you sometimes follow the same recipe, but get a complete different outcome - a batch effect.



In biological experiments, “the batch effect represents the systematic technical differences when samples are processed and measured in different batches and which are unrelated to any biological variation” (Leek, 2012)

Batch effects can arise from multiple sources such as; different reagents, experiment times, personnel, and instruments to name a few. It is important to note that batch effects are almost impossible to avoid.

Why is it a Problem?

Batch effects tend to lead to increased variability and decreased power when finding true positive biological signals.

When conducting gene expression studies, it is common to find that the greatest cause of differential expression is due to batch effects rather than biological groups (Leek, 2012).

Luckily, a well-designed high-throughput technology experiment will provide enough data to both detect and control for any batch effects.

Sources of batch effect

Batch effect is typically introduced by changes in what are through of as a benign change in any of the multiple steps that comprise your experiment. Examples include:

- Processing on different days
- Processing by different people
- Processing in different labs
- Differences in laboratory temperature
- Different batches of reagents
- Inconsistent sample handling (e.g. differences in storage and shipment or samples left on the bench, or off ice for too long)
- Processing and/or quantification using different instruments.

What Should You do if Batch Effects are Discovered in Your Data?

There are a number of technical approaches that can be applied at the analysis stage to control for batch effect, but the best thing you can do is design your experiment with due care, eliminating as many sources of variability as possible.

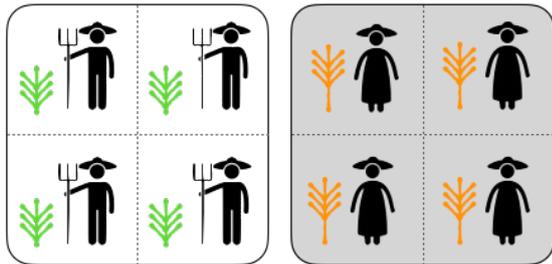
At the analysis stage, methods using some form of linear model can often include the batch variables as factors in the design (e.g. limma, DeSeq). Other approaches use Bayesian methods to attempt to control and remove batch effects (e.g. ComBat).

It should be noted that whilst batch effect algorithms can be used to remove batch effects they must be used carefully as there can be adverse effects to using these incorrectly.

An Example of a Badly Designed Experiment

When designing a high-throughput experiment, batches should be distributed fairly across biological groups.

For example, say you are wanting to compare two plant samples and you have two fields. The fields are also farmed with different personnel that are using different tools to farm the crops.



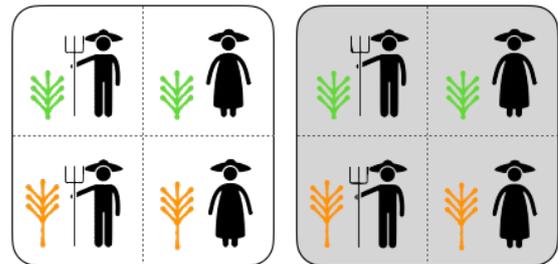
Each of the fields displayed show that they contain independent samples. They are also farmed by different personnel and tools.

This is bad experiment design as the real effect on each of the crops could be due to other factors such as personnel, tools, environmental factors etc. Designing an experiment in this way will almost always lead to confounded effects which can **not** be separated after data has been collected.

An Example of a Well Designed Experiment

A well designed experiment (as stated above) should distribute batches fairly across all biological groups to lessen the effects or chances of confounding.

An example of a well designed experiment, using the same field experiment as before.



This time, each of the fields contains both samples, with both personnel and both types of tools used. This is a well designed experiment as any real effect will most likely be due to biological effects rather than batch effects.

You must note, designing an experiment like this will not eradicate batch effects completely as they are almost impossible to avoid. However, designing an experiment like this will minimise the chances of confounding due to batch effects. Thus, leading to more meaningful conclusions.