

Experimental Design

STAT 430

Spring 2021 (1215)

TeX: *Cameron Roopnarine*

Instructor: *Nathaniel Stevens*

May 12, 2021

Contents

Contents	1
1 Week 1	2
1.1 Notation and Nomenclature	2
1.2 Experiments versus Observational Studies	4
1.3 QPDAC: A Strategy for Answering Questions with Data	5
1.4 Fundamental Principles of Experimental Design	6

Chapter 1

Week 1

1.1 Notation and Nomenclature

EXAMPLE 1.1.1: Experiment 1 — List View vs. Tile View

Suppose that **Nike**, the athletic apparel company, is experimenting with their mobile shopping interface, and they are interested in determining whether changing the user interface from *list view* to *tile view* will increase the proportion of customers that proceed to checkout.

EXAMPLE 1.1.2: Experiment 2 — Ad Themes

Suppose that **Nixon**, the watch and accessories brand, is experimenting with four different video ads that are to be shown on Instagram. The first has a surfing theme, the second has a rock climbing theme, the third has a camping theme, and the fourth has an urban professional theme. Interest lies in determining which of the four themes, on average, is watched the longest.

DEFINITION 1.1.3: Metric of interest

The **metric of interest** (MOI) is the statistic the experiment is meant to investigate.

REMARK 1.1.4

Typically, we want to optimize for the metric of interest; that is, we would like to either maximize or minimize it.

EXAMPLE 1.1.5: Metric of Interest

- Key performance indicators (KPIs): a statistic that quantifies something about a business.
 - Click-through rates (CTRs).
 - Bounce rate.
 - Average time on page.
 - 95th percentile page load time.
- *Nike Example*: checkout rate (COR).
- *Nixon Example*: average viewing duration (AVD).

DEFINITION 1.1.6: Response variable

The **response variable**, denoted y , is the variable of primary interest.

REMARK 1.1.7

The response variable is what needs to be measured in order for the MOI to be calculated.

EXAMPLE 1.1.8: Response Variable

- *Nike Example*: binary indicator indicating whether a customer checked out.
- *Nixon Example*: the continuous measurement of viewing duration for each user.

DEFINITION 1.1.9: Factor

The **factor**, denoted x , is the variable(s) of secondary interest.

Also known as: **covariates, explanatory variates, predictors, features, independent variables.**

REMARK 1.1.10

The factors are thought to influence the response (dependent) variable.

EXAMPLE 1.1.11: Factor

- *Nike Example*: the factor is the *visual layout*.
- *Nixon Example*: the factor is the *ad theme*.

DEFINITION 1.1.12: Experimental conditions

The **experimental conditions** are the unique combinations of levels of one or more factors.

Also known as: **treatments, variants, buckets.**

DEFINITION 1.1.13: Levels

The **levels** are the values that a factor takes on in an experiment.

EXAMPLE 1.1.14: Levels

- *Nike Example*: {tile view, list view}.
- *Nixon Example*: {surfing, rock climbing, camping, business}.

DEFINITION 1.1.15: Experimental units

The **experimental units** are what is assigned to the experimental conditions, and on which the response variable is measured.

EXAMPLE 1.1.16: Experimental Units

- *Nike Example*: Nike mobile customers.
- *Nixon Example*: Instagram users.

REMARK 1.1.17

Often, in online experiments, the unit is a user/customer (i.e., person), but it does not have to be.

EXAMPLE 1.1.18

Uber matching algorithm experiment.

1.2 Experiments versus Observational Studies

DEFINITION 1.2.1: Experiment

An **experiment** is composed of a collection of conditions defined by *purposeful changes* to one or more factors. Here, we intervene in the data collection.

- The goal is to identify and quantify the differences in response variable values across conditions.
- In determining whether a factor significantly influences a response, like whether a video ad's theme significantly influences its AVD, it is necessary to understand how experimental units' response when exposed to each of the corresponding conditions.
- However, it would be nice if we could observe how the *same* units behave in each of the experimental conditions, but we can't. We only observe their response in a single condition.
- **Counterfactual**: the hypothetical and unobservable value of a unit's response in a condition to which they were not assigned. We may think of this as an "alternate reality."

EXAMPLE 1.2.2

Nixon Example: the "camping" response variable for units assigned to the "surfing" condition.

- Because counterfactual outcomes cannot be observed, we require a **proxy**. Instead, we randomly assign *different units* to *different experimental conditions*, and we compare their responses.
- Ideally, the only difference between the units in each condition is the fact that they are in different conditions.
 - We want the units to be as homogenous as possible, this will help facilitate **causal inference** (establishing causal connections between variables).
 - This is typically guaranteed by *randomization*.
- The key here is that the factors are purposefully controlled in order to observe the resulting effect on the response. This facilitates causal conclusions.
- In an **observational study**, on the other hand, there is no measure of control in the data collection process. Instead, data is collected passively and the relationship between the response and factor(s) is observed organically.
- This hinders our ability to establish causal connections between the factor(s) and the response variables. However, sometimes we have no choice.

EXAMPLE 1.2.3: Unethical Experiments

- *Unethical Experiment 1*: In evaluating whether smoking lung cancer, it would be unethical to have a 'smoking' condition in which subjects are forced to smoke.
- *Unethical Experiment 2*: In dynamic pricing experiments, it would be unethical to show different users different prices for the same products. For example, surge pricing in Uber/Lyft.
- *Unethical Experiment 3*: In social contagion experiments, it would be unethical to show some network users consistently negative content and others consistently positive content. **But Facebook did this anyway.**
- *Unethical Experiment 4*: Mozilla conducted an investigation in which the company was interested in determining whether Firefox users that installed an ad blocker were more engaged with the browser. However, it would have been unethical to force users to install an ad blocker, and so they were forced to perform an observational study with *propensity score matching* instead.

	Advantages	Disadvantages
Experiment	causal inference is clean	experiments might be unethical, risky, or costly
Observational Study	no additional cost, risk, or ethical concerns	causal inference is muddy

1.3 QPDAC: A Strategy for Answering Questions with Data

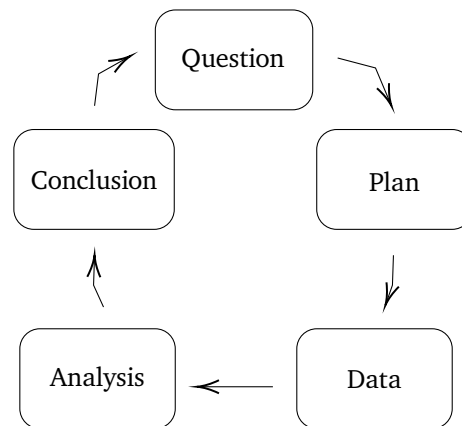


Figure 1.1: QPDAC Cycle

Question: Develop a clear statement of the question that needs to be answered.

- It is important that this is clear and concise and widely communicated, so all stakeholders are on the same page.
- The question should be quantifiable/measurable and typically stated in terms of the metric of interest.

EXAMPLE 1.3.1

- *Nike Example*: “which visual layout, tile view or list view, corresponds to the highest checkout rate?”
- *Nixon Example*: “which ad theme, camping, surfing, rock climbing, business, corresponds to the highest average viewing duration?”

Plan: In this stage, the experiment is designed and all pre-experimental questions should be answered.

- Choose the response variable. This should be dictated by the **Question** and the metric of interest.
- Choose the factor(s): brainstorm all factors that might influence the response and make decisions about whether and how they will be controlled in the experiment.
 - Design factors:** factors that we will manipulate in the experiment. The factors we’ve discussed in the Nike and Nixon examples are design factors.
 - Nuisance factors:** factors that we expect to influence the response, but whose effect we do not care to quantify. Instead, we try to eliminate their effects with *blocking*.
 - Allowed-to-vary factors:** factors that we *cannot* control and factors that we are unaware of

in an experiment.

- *Nixon Example*: users' age, gender, nationality.
- Choose the experimental units. These are what the response variable is measured on.
- Choose the sample size and sampling mechanism.
 - Sample size: how many units per experimental condition?
 - Sampling mechanism: how are they selected?

Data: In this stage, the data are collected according to the **Plan**. It is extremely important that this step be done correctly; the suitability and effectiveness of the analysis relies on the data being collected correctly. Computer scientists often use the phrase “garbage in, garbage out” to describe the phenomenon whereby poor quality input will always provide faulty output.

- A/A Test: units are assigned to one of two *identical* conditions.
 - We do this to ensure the assignment of units to conditions is truly random.
 - Two groups should be indistinguishable in terms of response distribution and other demographics.
 - If things aren't indistinguishable, there is a problem.
 - *Simple Ratio Mismatch Test*: check whether the observed sample ratios match what would be expected if assignment was truly done at random.
 - * Hypothesis test can be used to determine whether the proportion of units in each condition match what would have been expected under random assignment.

Analysis: In this stage, the **Data** are statistically analyzed to provide an objective answer to the **Question**.

- This is typically achieved by way of estimating parameters, fitting models, and carrying out statistical hypothesis tests. This is where we spend most of our time in the course.
- If the experiment was well-designed and the data was collected correctly, this step should be straightforward.

Conclusion: In this stage, the results of the **Analysis** are considered and one must draw conclusions about what has been learned.

- These conclusions should then be clearly communicated to all parties involved in — or impacted by — the experiment.
- Communicating “wins” and “loses” will help to foster the culture of experimentation.

1.4 Fundamental Principles of Experimental Design

DEFINITION 1.4.1: Randomization

Randomization refers both to the manner in which experimental units are *selected for inclusion* in the experiment and the manner in which they are *assigned to experimental conditions*.

REMARK 1.4.2

Typically, we don't include the entire target/study population.

Thus, we have two levels of randomization:

- The first level of randomization exists to ensure the sample of units included in the experiment is *representative of those that were not*.
 - Allows us to generalize conclusions beyond just the experimental units to units in the population not in the experiment.
- The second level of randomization exists to *balance* the effects of *extraneous variables* not under study (i.e., the allowed-to-vary factors).
 - Balancing the effects of allowed-to-vary factors makes our conditions homogenous and thus best mimics the counterfactual, thereby making causal inference easy.

DEFINITION 1.4.3: Replication

Replication refers to the existence of multiple response observations within each experimental condition and thus corresponds to the situation in which more than one unit is assigned to each condition.

- Assigning multiple units to each condition provides *assurance* that the observed results are genuine, and *not just due to chance*.
- For instance, consider the [Nike experiment](#) introduced previously. Suppose the CORs in the *list view* and *tile view* conditions were 0.5 and 1 respectively. This conclusion would be a lot more convincing if each condition had $n = 1000$ units as opposed to $n = 2$, where n is the sample size in *each* condition.
- How much replication is needed?
 - How big a sample size is needed?
 - Power analysis + sample size calculations will help answer this.

DEFINITION 1.4.4: Blocking

Blocking is the mechanism by which the nuisance factors are controlled for.

- To *eliminate* the influence of nuisance factors, we hold them fixed during the experiment.
- Thus, we run the experiment *at fixed levels of the nuisance factors*, i.e., within **blocks**.

EXAMPLE 1.4.5: GAP

Consider an email promotion experiment in which the primary goal is to test different variations of the *message in the subject* line with the goal of maximizing 'open rate.' However, suppose that it is known that the 'open rate' is also influenced by the time of the day and the day of the week that the email is sent.

We send all the emails at the same time of day and on the same day of week to control/eliminate the effect of time/day nuisance factor. By *blocking*, in this way, the nuisance factor can't confound our conclusions.