

Questions

- websiteoptimizer
- Gmail search:
 - <https://www.greplin.com/>
 - <http://cloudmagic.com/>

Qualitative vs. Quantitative

- Qualitative: Develop understanding of human experience
- Quantitative: Objectively measure human performance

- Less about more vs. more about less

When are each appropriate?

Quantitative Evaluation

- Gather (performance) measurements
- Methods
 - User interaction collection (i.e., logging)
 - *Mouse clicks, keys pressed,...*
 - *Data collected during system use*
 - Google, Amazon
 - Controlled experiments
 - *Set forth a testable hypothesis*
 - *Manipulate one or more **independent** variable*
 - *Observe effect on one or more **dependent** variable*
 - *Can be reproduced by others*

Controlled experiment

- State a lucid, testable hypothesis
- Identify independent and dependent variables
- Design the experimental protocol
- Choose the user population
- Apply for human subjects protocol review (IRB)
- Run some pilot participants
- Fix the experimental protocol
- Run the experiment
- Perform statistical analysis
- Draw conclusions
- Communicate results

Experiment Design

- Is it reliable (repeatable)?
 - Will you get the same result if someone else repeats the experiment?
 - *Confounding variables?*
 - Does the experiment take into account variations between subjects?

- Is it valid?
 - Does the experiment reflect target use?
 - *Were users typical?*
 - *Were tasks typical?*
 - *Was the setting realistic?*
 - *Was the experience biased?*

Are results significant?

- Statistical significance
 - Is observed result is due to chance?
 - Type I errors are the most disruptive

Researcher's Decision	Actual Situation	
	NO effect	Effect
NO effect	Correct decision	Type II error
Effect	Type I error	Correct decision

- Design significance?
 - 3.00s versus 3.05s?

Are results significant?

- Statistical significance
 - Comparing to the null hypothesis: “There is no effect”
 - Type I errors are the most disruptive

Researcher’s Decision	Actual Situation: Null Hypothesis is	
	True	False
Fail to reject the null hypothesis	Correct decision	Type II error
Reject the null hypothesis	Type I error	Correct decision

- Design significance?
 - 3.00s versus 3.05s?

Example

- Examine value of animation during scrolling
[Klein & Bederson 2005]
<http://portal.acm.org/citation.cfm?doid=1056808.1057068>
 - Various reading tasks
 - For various document types
 - While scrolling with varying kinds of animated transitions

State a lucid, testable hypothesis

“Animated scrolling will speed up reading and decrease errors, especially for plain documents.”

Choose the variables

- Manipulate one or more *independent* variable (the thing you change)
 - Document type
 - Animation speed
- Observe effect on one or more *dependent* variable (the thing you measure)
 - Time to completion
 - Accuracy (i.e., error rate)

Design the experimental protocol

- Between or within subjects?
 - Between subjects: each subject runs one condition
 - *Need more subjects*
 - *Difference between subjects might introduce a bias*
 - + *No learning effects across conditions*
 - Within subjects: each subject runs several conditions
 - + *Need fewer subjects*
 - + *No bias across participants*
 - *Possible learning effect across conditions*
- Which tasks?
 - Must reflect the hypothesis
 - Must avoid bias
 - *Instructions, ordering...*
 - *In doubt, always favor the null hypothesis*

Design the experimental protocol

- Running Example:
 - Reading while scrolling

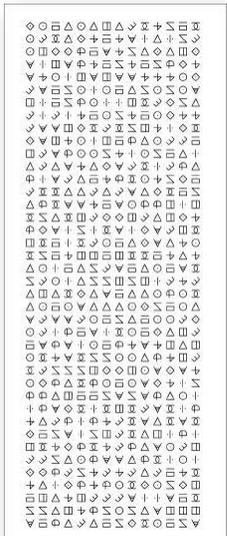
Two days before Christmas in 1947, Walter Brattain took a razor blade and sliced across a tiny piece of gold foil. He pressed the cut edges against a chunk of germanium metal, and connected wires to the foil. "I found that if I wiggled it just right," he said later, "I had an amplifier." He had created the world's first transistor. Brattain and his team at Bell Laboratories had been working toward that moment for months. They knew their invention would replace tubes—the fragile and unreliable glass tubes inside the radios and other electronic devices of the time. But what they could not have understood that winter day was how completely their invention would change society. Over the next 40 years, the transistor and its successor, the integrated circuit, swept away technology that had hardly changed since the invention of steam power. It was like a second industrial revolution. The first industrial revolution changed how the very things we made, replacing skilled workers with machines that could do the same jobs more quickly and cheaply. But the transistor did not just change factories or what they manufactured. It replaced many of them with an industry that created, stored, bought, and sold information. Compared with creating and selling real things, such as cars or cars, feeding information is perhaps a strange idea, but it is not a new one. Back in the mid-19th century, one of the first applications of the electric telegraph was in news distribution. News agencies such as Reuters and Associated Press collected news from their correspondents, and then used the electric telegraph to sell information to newspaper publishers. The transistor speeded up the collection, processing, and distribution of many different kinds of information—long before it replaced expensive glass tubes on the circuit boards of computers.

The Process of Invention During the Industrial Revolution

The Patent System
By the 18th century, the patent system, which had begun in Italy in 1471, was widespread and encouraging invention. Inventors disclosed the full details of their work in a patent application. On approval of the idea by the government patent office, an account of the invention was published, so protecting it from being copied for a fixed number of years—typically 16. Patents stimulated invention in three ways: they informed people of the latest technology, inventors had time to profit from their ideas before they were copied, and, by limiting the period of protection, the patent system ensured that good ideas would eventually benefit everybody, not just the inventor.

Warfare Fuels Invention
Warfare has always stimulated invention. The first war of the industrial age, the Crimean War (1853-56), provoked a surge of weapons patents—and it is not hard to see why. Weapons inventors needed to sell their idea to just one wealthy customer—the government. If successful, they could be sure of enormous orders. One inventor who grasped the advantages of military technology was American-born engineer William S. Mason, whose water-cooled machine gun is remembered long after his other inventions—including a better mousetrap—were forgotten. Skills learned making weapons may also affect nonmilitary technologies. Metallurgy techniques improved as foundries struggled to make stronger gun barrels. And it is no coincidence that the costly technology of interchangeable parts first appeared in the manufacture of rifles.

Mass-Production
Though costly, the factory system and interchangeable parts made it possible to manufacture large numbers of goods. In the mass-production of muskets, inspectors carried gauges like these from army to army to check that each dimension of each piece of a gun was correct—and interchangeable with all the other similar pieces.



Chose the user population

- Pick a well balanced sample
 - Novices, experts, average
 - Age group
 - Sex...
- Population group may be one of the independent variable
- Running example
 - Varied population, did not control

Run the experiment

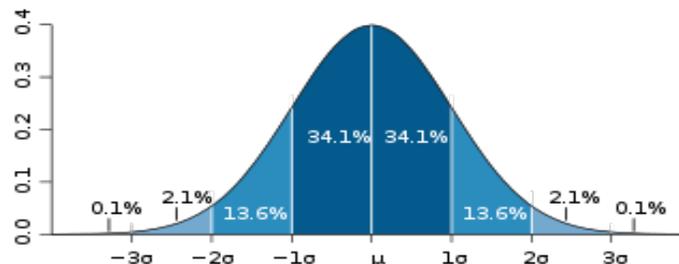
- Always run pilots first!
 - There are always unexpected problem!
 - When the experiment has started you cannot pick and choose
- Use a check-list so that all subjects follow the same steps
- IRB - Don't forget the consent form!
- Don't forget to debrief each subject

Run statistical analysis

- Properties of our result data
 - Mean, variance...
- How different data sets relate to each other
 - Are we sampling from similar or different distributions?
- Probability that our claims are correct
 - Statistical significance:
 - “The hypothesis that technique X is faster is accepted ($p < .05$)” means that there is a higher than 95% chance the hypothesis is true
 - Typical levels are .05 and .01 level
 - These levels are socially determined

Statistical tools I - Descriptive

- Mean
 - Average
- Median
 - Central value
- Standard Deviation
 - Measure of variation



Statistical tools I - Descriptive

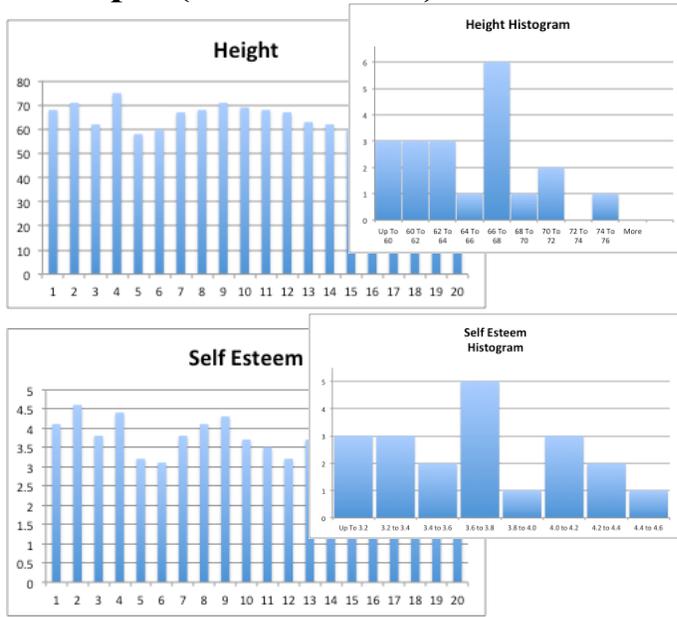
- Correlation
 - Measure the extent to which 2 concepts are related
 - Caveats
 - *Correlation does not imply cause and effect (hidden variable)*
 - Ice cream consumption and drowning
 - Third variable problem
 - Directionality problem
 - *Need a large enough group*
 - $r = 0$ (no correlation)
 - $r = 1$ (positively correlated)
 - $r = -1$ (negatively correlated)



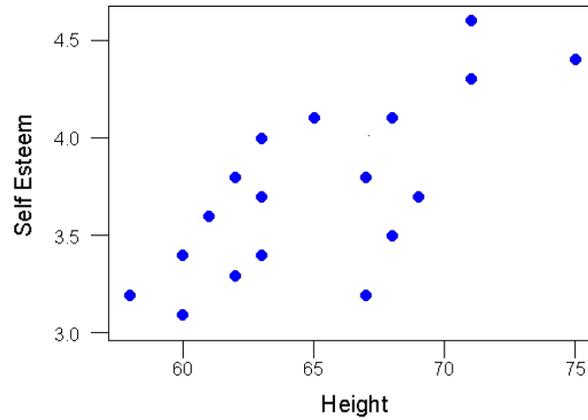
Example (random data)

Height	Self Esteem
68	4.1
71	4.6
62	3.8
75	4.4
58	3.2
60	3.1
67	3.8
68	4.1
71	4.3
69	3.7
68	3.5
67	3.2
63	3.7
62	3.3
60	3.4
63	4.0
65	4.1
67	3.8
63	3.4
61	3.6

Mean 65.4
Median 66.0
SD 4.4



Example (continued)

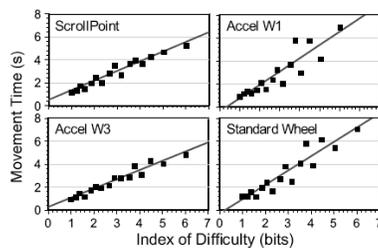


Correlation: $r = 0.73$

Statistical tools I - Descriptive

- Regression
 - Calculate the “best fit”

	R	R ²	Slope	Intercept (s)	IP (bps)
ScrollPoint	0.97	0.94	0.84	0.42	1.19
Accel W1	0.90	0.81	1.16	-0.51	0.86
Accel W3	0.97	0.95	0.80	0.18	1.25
Wheel Std	0.94	0.88	1.25	-0.42	0.80

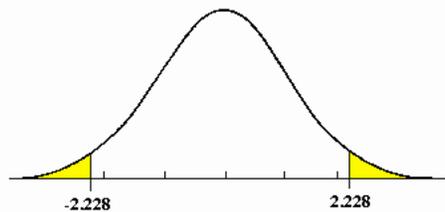


Statistical tools II - Analytical

- T-test
 - Compare the mean of 2 populations
 - *Null hypothesis: no difference between means*
 - *Can only examine a single independent variable*
 - Assumptions
 - *Samples are normally distributed*
 - Very robust in practice
 - *Population variances are equal*
 - Reasonably robust for differing variances
 - *Individual observations in samples are independent*
 - Very important

T-Tests

- Increase statistical power by:
 - Increasing number of participants
 - Decreasing variance
 - Running 1-tail instead of 2-tail test



Statistical tool II - Analytical

- ANOVA
 - Single factor analysis of variance
 - *Compare three or more means*
 - Analysis of variance
 - *Compare relationship between many factors*
 - Beginners type at the same speed on all keyboards,
 - Touch-typist type fastest on the qwerty
- Your protocol influences the kind of test you can use
 - If in doubt, consult with a statistician before starting the experiment!

Reporting Results

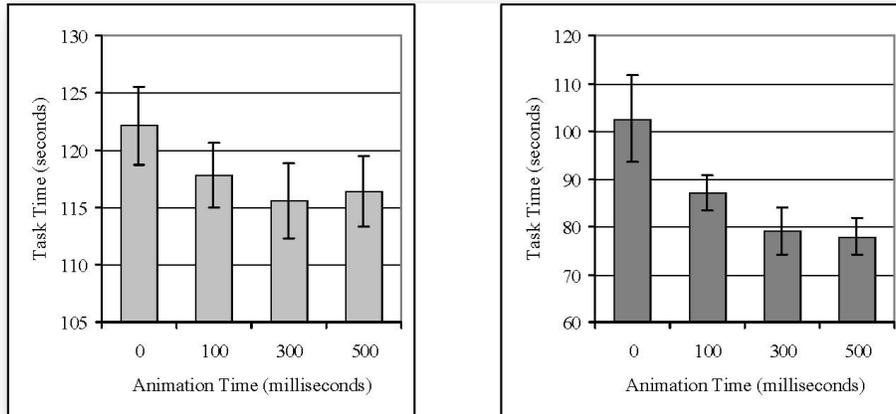
T-TEST

“There was a significant difference in the scores for sugar (M=4.2, SD=1.3) and no sugar (M=2.2, SD=0.84) conditions; $t(8)=2.89$, $p = 0.020$.”

ANOVA

“As one would expect, movement times increased as either W decreased or D increased (i.e., as the task got more difficult: for W, $F(2,25)=801$, $p<0.001$; and for D, $F(3,54)=1429$, $p<0.001$).”

Running example data



Error bars show +/- one standard deviation

Running Example Analysis

- Measured speed, error and RSD (relative subjective duration)
- Significant results for
 - reading time
 - reading error
 - counting time
- Animated scrolling
 - reduces reading errors by up to 54%
 - task time by up to 3%
 - task time by up to 24% for counting tasks
- 300ms seems to be best overall animation time for reading.
- Formatted documents had a higher base rate perf, and lower improvement – thus visual landmarks appear to be a good idea

In Class Experiment

- Compare mouse vs. keyboard
- Independent variable:
 - Interaction type (conditions: mouse, keyboard)
- Dependent variable:
 - Speed
- Experiment:
 - Bring up two GDoc windows with 10 lines of supplied text in one
 - Copy each line of text (one at a time) from one window to the other
 - Conditions:
 1. Select text with mouse and press ctrl-c, then ctrl-tab, ctrl-v
 2. Select text with mouse and choose edit->copy menu, then click on other tab and choose edit->paste
 - Manually measure entire task and divide by 10
 - Train on 3 lines
 - Enter data in shared spreadsheet.
 - Work in pairs with one person timing the other, then switch.