Richard W. Hamming



Learning to Learn

The Art of Doing Science and Engineering

Session 27: Unreliable Data

Bad data



The data you get is not as reliable as advertised

Simulation data is normally assumed to be accurate, but we cannot depend on data the way we expect

First, some examples of bad physical data...

Life Testing



Life testing - how reliable is equipment

Example: testing reliability of vacuum tubes

- Wanted 20 year life expected failure cost = \$1M
- Tubes available for testing 18 months before use

Why believe testing equipment is as reliable as what you are testing?

Life Testing (cont')



Tricks in Life Testing

- Raise temp 17 degrees C and item ages twice as fast
- Increase voltage and stress resistance breakdown
- Increase frequency and expose failures earlier

Hard to adequately test highly reliable equipment in a short time period

- Too anxious to get latest technology into the field
- Not enough time to test it right, but plenty of time to fix

Accurate Measurements



Never assume people, equipment can't make measurement mistakes

Take apart test equipment and recalibrate it

Never trust the manufacturer on advertised accuracy

Just because a machine records data doesn't make it correct

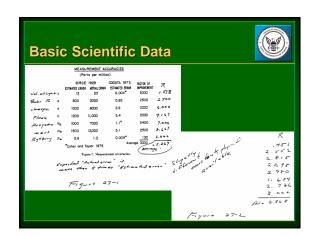
Never process data before checking for errors first

· it wastes more time to redo work

Test Studies



Never trust a pilot program (i.e. a test study) to accurately predict the success of the main study







Hamming examines increased accuracy of National Bureau of Standards physical constants

 Assuming new constants are correct, averaged 5 sigma for previous estimated error

Any number with a probable error is typically much smaller than it should be

• Scientists fine tune equipment to minimize variance, not error

90% of the time the next independent measurement will fall outside the previous 90% confidence limits

Fitting Data to a Model



There can be errors in both the model and the data

- Initially, the model is a guess and the data is bad
- As the data & measurements improve, the model is rarely reexamined for error
- Thus, as measurement accuracy increases, the error in the model increases

Estimates combined with reliable measurements are taken at the reliability of the best measurement

"A chain is only as strong as its weakest link"

Economic Definitions



An economic definition today is not necessarily what it was yesterday.

• e.g. gold flow, income, inventory figures, poverty rate

Hard to draw conclusions from a time series when the definition of what is being measured constantly changes.

Often a changing definition is used to satisfy the needs of a desired policy (e.g. poverty rate) or outcome (e.g. unstated inventory to avoid taxes).

Economic Predictors and Indices



... change very slowly & lag behind reality

- Poorly represents current state of affairs
- Example: manufacturing vs. processing information

Worthy argument: it's better to have an index that is comparable and wrong, versus gaining the right index that is incomparable

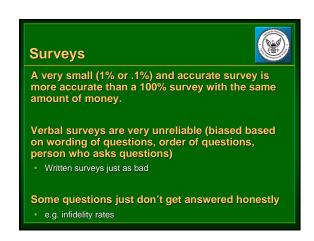
Economists

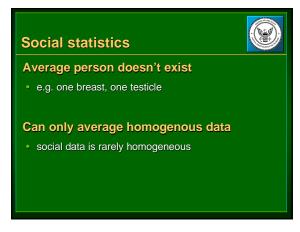


Refuse to admit a data problem exists

Economists & governments typically skew figures to enact a particular policy

(e.g. economic aid, gold holdings)





Management use of data Surveys are answered quicker by those who think they will rate well than those who don't • Presence of management will change data, surveys, observation, etc. • Subordinates will skew data for the boss so as not to break bad news. Keep all this in mind as you move through the ranks, and realize that the data that reaches you is often unreliable