# Thinking About Data

MMED

African Institute for the Mathematical Sciences

Muizenberg, South Africa

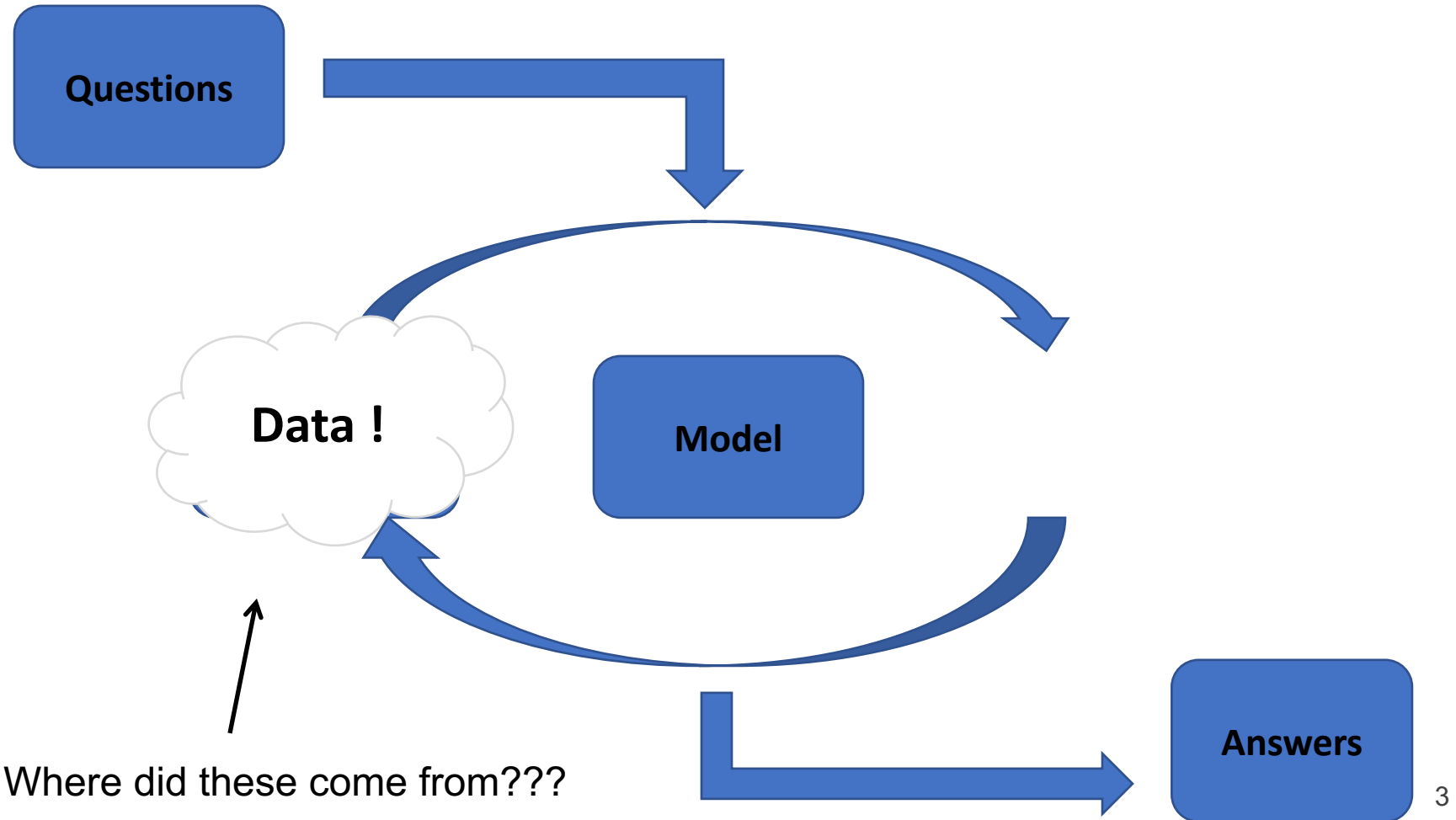May, 2019

Jim Scott, Ph.D, M.A., M.P.H.

# Goals:

- Recognize that data are better than anecdotes

- Understand that how the data are collected matters

- Be aware of variability – don't be fooled by it

- Consider confounding and other forms of bias

- Healthy skepticism

ICI3D

# Why do we care about data?

Questions

Data !

Model

Answers

Where did these come from???

3

ICI3D

"Data, data, data!", he cried impatiently. "I can't make bricks without clay"

- Sherlock Holmes

February 3, 2010

# Journal Retracts 1998 Paper Linking Autism to Vaccines

By GARDINER HARRIS

A prominent British medical journal on Tuesday retracted a 1998 research paper that set off a sharp decline in vaccinations in Britain after the paper's lead author suggested that vaccines could cause autism.

The retraction by The Lancet is part of a reassessment that has lasted for years of the scientific methods and financial conflicts of Dr. Andrew Wakefield, who contended that his research showed that the combined measles, mumps and rubella vaccine may be unsafe.

"The story became credible because it was published in The Lancet," Alison Singer, president of the Autism Science Foundation, said Tuesday. "It was in The Lancet, and we really rely on these medical journals."

5

# Consequences



Commons.wikimedia.org

## National MMR vaccination catch-up programme announced in response to increase in measles cases

A national catch-up programme to increase MMR vaccination uptake in children and teenagers is announced today by Public Health England, NHS England and the Department of Health.

Experts believe the rise in measles cases can be mostly attributed to the proportion of unprotected 10-16 year-olds who missed out on vaccination in the late 1990s and early 2000s when concern around the discredited link between autism and the vaccine was widespread. At this time measles had been eliminated in the UK, but coverage fell nationally to less than 80% in 2005, with even lower uptake in some parts of the country. After many years of low vaccination uptake, measles became re-established in 2007.

# How the Data are Collected Matters

- "Always do right.  This will gratify some people, and astonish the rest"

  - Mark Twain

- Beware:  All data are not created equal

*Source: Statistics* 3rd, ed. Pisani, Purves, Freedman

7

# 1936 Presidential Election

**Franklin D. Roosevelt**

Alf Landon

# 1936 Presidential Election

- 1936 Literary Digest Poll

- Literary Digest had predicted the winner of every US presidential election since 1916.

- In 1936, Literary Digest mailed questionnaires to 10 million people (25% of voters).

- 2.4 million people responded

  - Returned questionnaires:
    - **Landon: 1,293,668                    57%**
    - **FDR:        972,897                    43%**

Source:   http://historymatters.gmu.edu/d/5168/

# Results

- **Actual Result: Roosevelt 61%, Landon 37%.**

- **One of the biggest landslides in U.S. history**



| | | |
|---|---|---|
| Electoral vote | 523 | 8 |
| States carried | 46 | 2 |
| Popular vote | 27,752,648 | 16,681,862 |
| Percentage | 60.8% | 36.5% |

# What went wrong?

- How were the data collected?

  - Those who received the questionnaire were systematically different than those who didn't
    - 10 million sent out  (~25% of voters)
    - 2.3 million returned – sample of convenience
      - Not representative
  - Sampling frame:
    - Telephone books
    - Automobile registries

ICI3D

# Bias and Variability

- Sample size doesn't matter if the data collection scheme is flawed

Observed value = The TRUTH + **Bias** + *Chance Error*

ICI3D

# Another example?

- IBM created Watson
  - Uses data to make predictions – finds patterns in knowledge database

- Trained Watson to play Jeopardy

- Crushed all humans

- Triumph for A.I.
  - …..mostly…..

# U.S. Cities

Its largest airport is named for a World War II hero; its second largest, for a World War II battle

Watson:
What is…… Toronto(?!)

# Oops

# Big Data Limitations

# Variability

- "When the facts change, I change my mind. What do you do sir?"

  - John Maynard Keynes

- Variation is everywhere

Observed value = Truth + Bias + Random Error

# Time Series of four "stock prices"



Stock #1

Stock #2

# Time series of four "stock prices"



20

# Don't be fooled by randomness



**1000 Random Stocks**

# Confounding



Per Capita Expenditures on Road Maintenance(?)

ICI3D

# What is the conclusion?

## January 2009
## Percent of Planes Delayed from City of Origin

| Airport | Continental | | | United | | |
|---|---|---|---|---|---|---|
| | Late | Total | % | Late | Total | % |
| Newark | 957 | 3998 | *23.9* | 100 | 399 | 25.1 |
| LaGuardia | 62 | 356 | *17.4* | 113 | 573 | 19.7 |
| Pittsburg | 8 | 60 | *13.3* | 17 | 119 | 14.3 |
| Detroit | 16 | 145 | *11.0* | 16 | 139 | 11.5 |
| Totals | 1043 | 4559 | 22.9 | 246 | 1230 | *20.0* |

slide credit: Jeff Witmer, data source: www.bts.gov

ICI3D

# How about now?

```
. logistic delay continental

Logistic regression                              Number of obs   =        5789
                                                 LR chi2(1)      =        4.72
                                                 Prob > chi2     =      0.0298
Log likelihood = -3067.3063                      Pseudo R2       =      0.0008

-------------------------------------------------------------------------------
       delay | Odds Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
 continental |   1.186576   .0943644     2.15   0.031     1.015318     1.38672
-------------------------------------------------------------------------------
```

Unadjusted
OR = 1.19
95% CI = (1.02, 1.39)

```
. logistic delay continental laguardia newark pittsburg

Logistic regression                              Number of obs   =        5789
                                                 LR chi2(4)      =       46.30
                                                 Prob > chi2     =      0.0000
Log likelihood = -3046.5183                      Pseudo R2       =      0.0075

-------------------------------------------------------------------------------
       delay | Odds Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
 continental |   .9161694   .0859693    -0.93   0.351     .7622593    1.101156
   laguardia |   1.807797   .3722533     2.88   0.004     1.207464    2.706607
      newark |    2.58219   .5031265     4.87   0.000     1.762527    3.783036
    pittsburg |   1.259086    .360524     0.80   0.421     .7183302     2.20692
-------------------------------------------------------------------------------
```

Adjusted
OR = 0.92
95% CI = (0.76, 1.10)

24

# Obesity in the United States: 1995 - 2006



Sources: CDC (BRFSS),
Merriman, Curhan, & Ford Fitness & Wellness report, 2006

# What do the data say?

- Caution:

- This may require thoughtful consideration

ICI3D

# What is the relationship?

# What is the relationship?

# Median Global Temperature During the Past 50 Years

# Love statistics(?)

# Are you a believer?


Cat injuries by # of stories fallen

# Cat conundrum

As long as it experiences acceleration, the cat probably extends its limbs reflexly, but on reaching terminal velocity it may relax and extend the limbs more horizontally in flying-squirrel fashion, thus not only reducing the velocity of fall but also absorbing the impact over a greater area of its body. This may explain the paradoxical decrease of mortality and injury in cats that fall more than 100 feet.
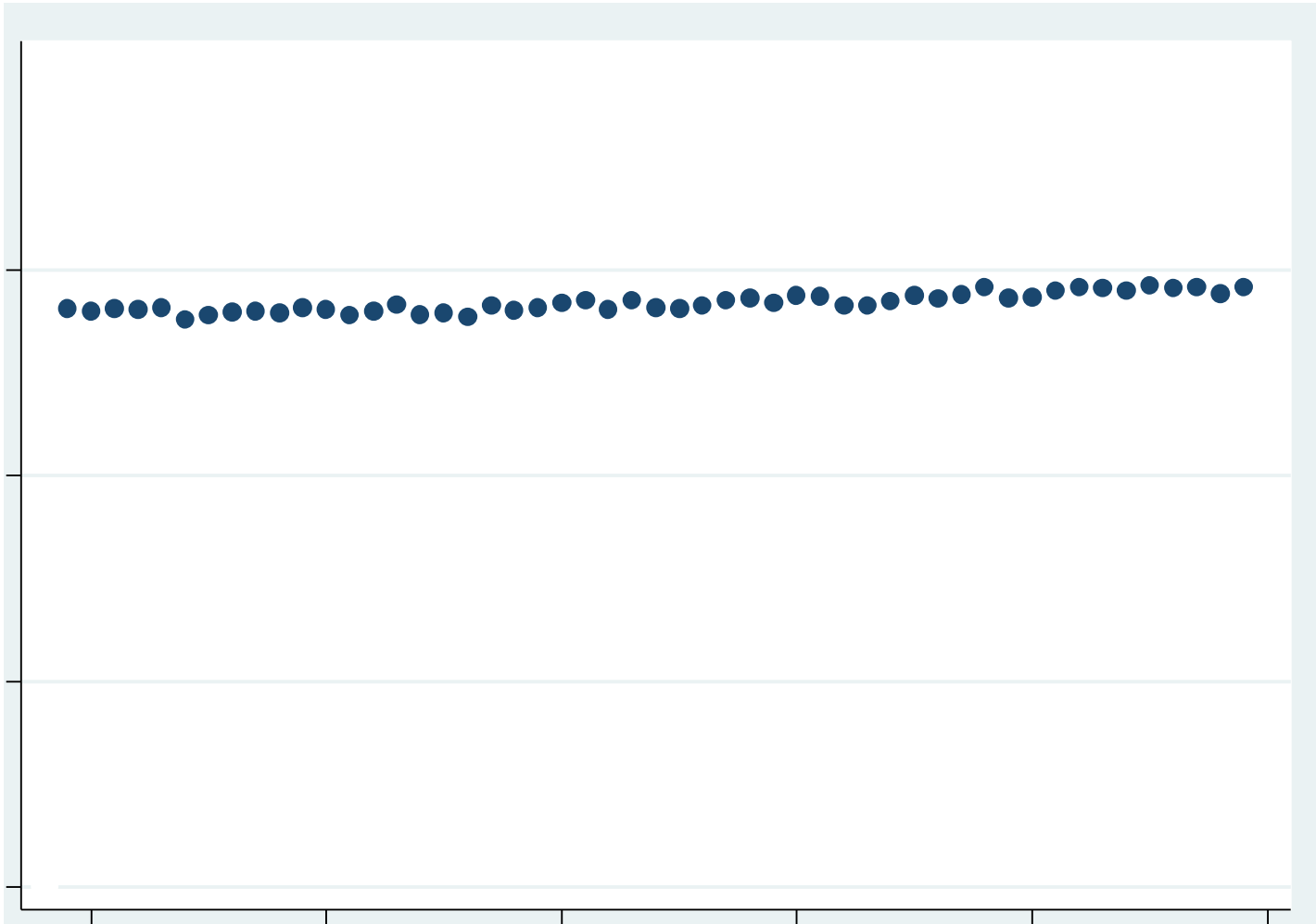
| Stories Fallen | # of cats |
|---|---|
| 1 | 0 |
| 2 | 8 |
| 3 | 14 |
| 4 | 27 |
| 5 | 34 |
| 6 | 21 |
| 7-8 | 9 |
| 9-32 | 13 |

Mortality rates for falling adult humans and cats (above), and number of total injuries and various types of injury per falling cat (below), as a function of number of stories fallen. (Based on the work by Waring and Demling and by Whitney and Mehlhoff.)

33

# What do you think?

A Randomized, Controlled Trial of the Effects of Remote, Intercessory Prayer on Outcomes in Patients Admitted to the Coronary Care Unit

**Table 3. Effects of Intercessory Prayer on Individual Components of the Mid America Heart Institute–Cardiac Care Unit (MAHI-CCU) Score***

| | No. (%) of Patients | | |
|---|---|---|---|
| MAHI-CCU Score Component | Usual Care Group (n = 524) | Prayer Group (n = 466) | P |
| Antianginal agents | 59 (11.3) | 47 (10.1) | .62 |
| Antibiotics | 82 (15.6) | 77 (16.5) | .77 |
| Unstable angina | 4 (0.8) | 1 (0.2) | .38 |
| Arterial monitor | 42 (8.0) | 32 (6.9) | .57 |
| Catheterization | 180 (34.4) | 162 (34.8) | .94 |
| Antiarrhythmics | 56 (10.7) | 50 (10.7) | .94 |
| Inotropes | 76 (14.5) | 69 (14.8) | .96 |
| Vasodilation | 78 (14.9) | 59 (12.7) | .36 |
| Diuretics | 112 (21.4) | 97 (20.8) | .89 |
| Pneumonia | 10 (1.9) | 12 (2.6) | .62 |
| Atrial fibrillation | 17 (3.2) | 12 (2.6) | .66 |
| Supraventricular tachycardia | 6 (1.1) | 2 (0.4) | .29 |
| Hypotension | 7 (1.3) | 8 (1.7) | .82 |
| Anemia/transfusion | 66 (12.6) | 50 (10.7) | .42 |
| Temporary pacer | 16 (3.0) | 13 (2.8) | .95 |
| Third-degree heart block | 1 (0.2) | 2 (0.4) | .60 |
| Readmit to cardiac care unit | 22 (4.2) | 25 (5.4) | .48 |
| Swan-Ganz catheter | 172 (32.8) | 123 (26.4) | .03 |
| Implanted cardiac defibrillator | 6 (1.1) | 10 (2.1) | .32 |
| Electrophysiology study | 15 (2.9) | 10 (2.1) | .61 |
| Radiofrequency ablation | 8 (1.5) | 2 (0.4) | .11 |
| Extension of infarct | 2 (0.4) | 0 (0.0) | .50 |
| Gastrointestinal bleed | 12 (2.3) | 5 (1.1) | .22 |
| Interventional coronary procedure | 155 (29.6) | 121 (26.0) | .21 |
|   PTCA alone | 69 (13.2) | 62 (13.3) | .95 |
|   PTCA with stent and/or rotablator | 86 (16.4) | 59 (12.7) | .10 |
| Permanent pacer | 21 (4.0) | 12 (2.6) | .28 |
| Congestive heart failure | 17 (3.2) | 19 (4.1) | .60 |
| Ventricular fibrillation/tachycardia | 12 (2.3) | 10 (2.1) | .95 |
| Intra-aortic balloon pump | 20 (3.8) | 12 (2.6) | .36 |
| Major surgery | 76 (14.5) | 51 (10.9) | .11 |
| Sepsis | 7 (1.3) | 7 (1.5) | .96 |
| Intubation/ventilation | 27 (5.2) | 26 (5.6) | .88 |
| Cardiac arrest | 6 (1.1) | 5 (1.1) | .84 |
| Death | 46 (8.8) | 42 (9.0) | .99 |

**Table 4. Effects of Intercessory Prayer on Mid America Heart Institute–Cardiac Care Unit (MAHI-CCU) Scores and Length of Stay in the CCU and in the Hospital***

| | Mean ± SEM | | | |
|---|---|---|---|---|
| | Usual Care Group (n = 52) | Prayer Group (n = 466) | Percentage Change | P |
| MAHI-CCU score | 7.13 ± 0.27 | 6.35 ± 0.26 | −11 | .04 |
| Unweighted MAHI-CCU score† | 3.00 ± 0.10 | 2.70 ± 0.10 | −10 | .04 |
| Length of CCU stay, d‡ | 1.23 ± 0.09 | 1.12 ± 0.08 | −9 | .28 |
| Length of hospital stay, d‡ | 5.97 ± 0.29 | 6.48 ± 0.54 | +9 | .41 |

34

ICI3D

# Skepticism vs Openness

- "It seems to me what is called for is an exquisite balance between two conflicting needs: the most skeptical scrutiny of all hypotheses that are served up to us and at the same time a great openness to new ideas …

- If you are only skeptical, then no new ideas make it through to you …

- On the other hand, if you are open to the point of gullibility and have not an ounce of skeptical sense in you, then you cannot distinguish the useful ideas from the worthless ones."

- Carl Sagan

# Summary

- We use data to help form and revise models

- Data collection should not be haphazard

- Data are not easily obtained – especially good data

- Think about how the data were collected and be sure to consider:
  - Bias and confounding
  - Variability
  - How the data are presented

- Be skeptical but open

ICI3D