# Design and Analysis of Sample Surveys

Andrew Gelman
Department of Statistics and Department of Political Science
Columbia University

Class 14b: Review

- ▶ Goal: learning about the population
- ▶ Intermediate steps:
    - ▶ Sample to population (sampling)
    - ▶ Survey response to question of interest (measurement)

# Happiness and the Tea Party movement

- A Brooks *New York Times* op-ed:

  > *People at the extremes are happier than political moderates .... none, it seems, are happier than the Tea Partiers ...*
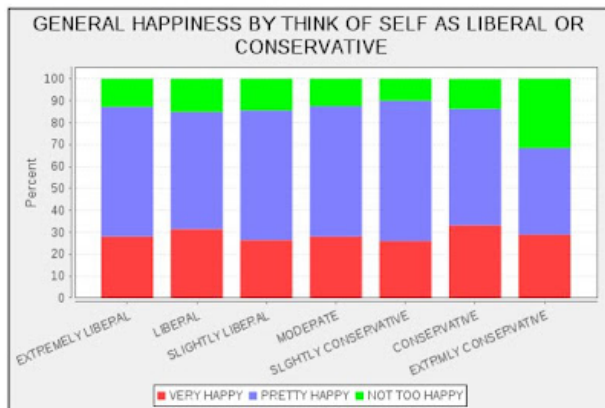
- But sociologist Jay Livingston writes:

  > *The GSS does not offer "bitter" or "Tea Party" as choices, but extreme conservatives are nearly three times as likely as others to be "not too happy."*
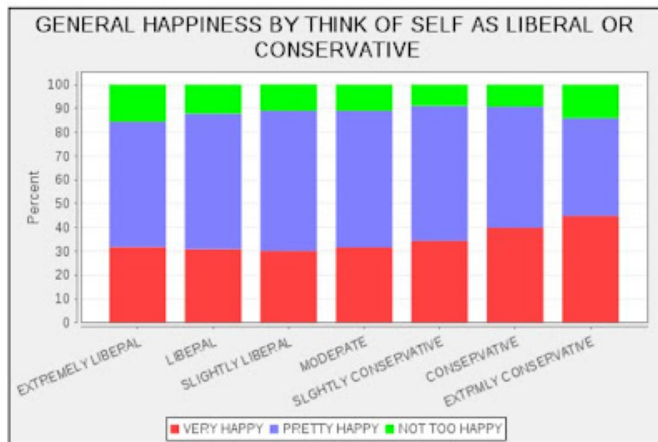
- Let's look at the data!

## Data from General Social Survey



**Chart for YEAR = 4(2009-2010)**

GENERAL HAPPINESS BY THINK OF SELF AS LIBERAL OR CONSERVATIVE

EXTREMELY LIBERAL, LIBERAL, SLIGHTLY LIBERAL, MODERATE, SLIGHTLY CONSERVATIVE, CONSERVATIVE, EXTPMLY CONSERVATIVE

VERY HAPPY  PRETTY HAPPY  NOT TOO HAPPY

▶ Is this just sampling variation?
  ▶ Sample size for "Extremely Conservative" here is 80
  ▶ Thus the standard error for that green bar on the right is approx $\sqrt{0.3 \cdot 0.7/80} = 0.05$

# How did Brooks get this wrong?



GENERAL HAPPINESS BY THINK OF SELF AS LIBERAL OR CONSERVATIVE

- ▶ Averaging over all the years, conservatives seem pretty happy!
- ▶ The importance of descriptive inference
  - ▶ Be careful about explaining patterns that aren't real!

# Class 1b: Statistical inference and linear regression

- $\sqrt{p(1-p)/n}$ or $\sigma/\sqrt{n}$
- $\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}$
- Sample size calculations
- $(y+2)/(n+4)$
- Weighted averages
- Living with uncertainty

# Linear regression

- ► Interpreting coefficients
- ► Building models
- ► The role of statistical significance

# Class 2a: Logistic regression

- ▶ Logistic curve
- ▶ Divide-by-4 rule
- ▶ Latent continuous variable
- ▶ Choice models

- The statistical significance filter
- When possible, study large effects
- Study effects in context

# Business-relevant examples . . .

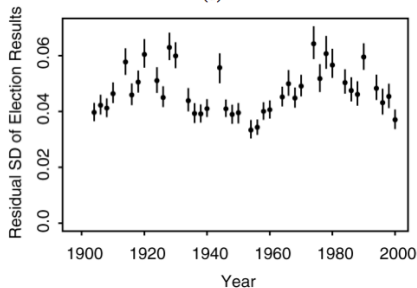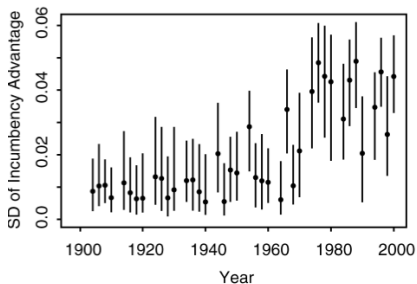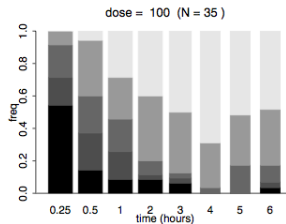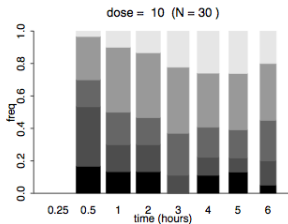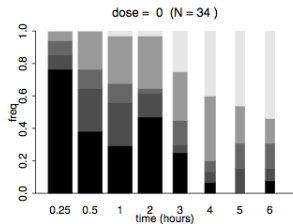# Example: incumbency advantage over time
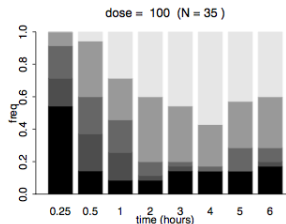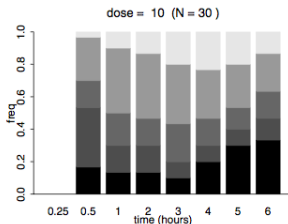
# Class 3a: Nonresponse and survey adjustment
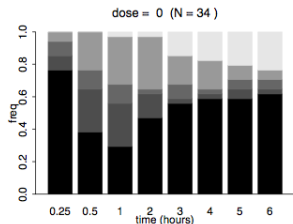
- Unit nonresponse
- Item nonresponse
- Intentional missing data
- Structural missing data
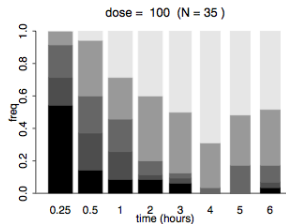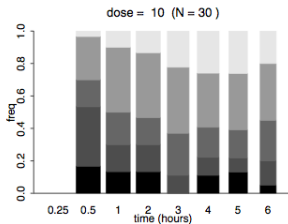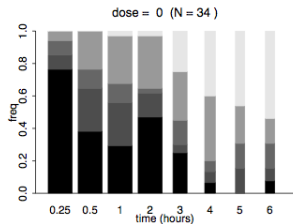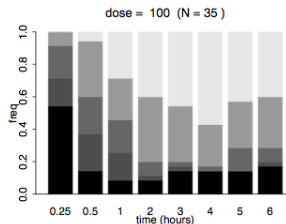
# Observed and completed data



Observed data display

Completed data display

# Class 3b: Adjusting for nonresponse

- ▶ Subsetting and imputation for item nonresponse
- ▶ Poststratification and weighting for unit nonresponse

- Where do weights come from?
- Poststratification identity: $\theta = \frac{\sum_j N_j \theta_j}{\sum_j N_j}$
- Challenges in weighting
- Challenges in poststratification
- Weighted regression using the "survey" package in R

# Example: CBS/New York Times pre-election polls

| id | org | y | state | edu | age | adults | weight |
|-----|--------|-----|-------|-----|-----|--------|--------|
| 6140 | cbsnyt | NA | 7 | 3 | 1 | 2 | 923 |
| 6141 | cbsnyt | 1 | 39 | 4 | 2 | 2 | 558 |
| 6142 | cbsnyt | 0 | 31 | 2 | 4 | 1 | 448 |
| 6143 | cbsnyt | 0 | 7 | 3 | 1 | 2 | 923 |
| 6144 | cbsnyt | 1 | 33 | 2 | 2 | 1 | 403 |

- The weight is listed as just another survey variable
- But they are actually constructed *after* the survey
- Weights $w_i = g(X_i, \theta)$

- ▶ Ratio estimation of a ratio
- ▶ Ratio estimation of a population average
- ▶ Regression estimation
- ▶ Robustness through model and design

# Regression estimation as a general framework

- Fit a regression, $y_i = a + bx_i + \text{error}$
- Regression estimate of $\overline{Y}$ is $\bar{y} + b(\overline{X} - \bar{x})$
- Special cases:
  - $b = 0$: unadjusted sample average
  - $b = \frac{\bar{y}}{\bar{x}}$: ratio estimation
  - $b = 1$: simple adjustment
- Regression estimation is valid for any $b$
  - Optimal for $b = \text{least-squares estimate}$

# Class 5a: Simple and stratified random sampling

- Sampling from a list
- Systematic sampling
- Stratified sampling
- Design and analysis

- ▶ Why do cluster sampling?
- ▶ Goal of equal-probability sampling
- ▶ Analysis of cluster data
- ▶ Design effects

- Sampling with equal probability at both stages
- Sampling clusters with probability proportional to size
- Adjusting for unequal sampling probabilities
- Design effects

# Class 6b: Inference for regression coefficients

- ► Option 1: weighted regression
- ► Option 2: unweighted regression, including in the model all variables that affect the probability of inclusion in the survey
- ► Discuss
- ► Population mean as a special case of a regression coefficient
- ► Population difference as a special case of a regression coefficient
- ► Practical limitations of weighting
- ► Practical limitations of modeling

- ▶ Our ideal procedure:
    - ▶ As easy to use as hierarchical regression
    - ▶ Population info included using poststratification
- ▶ Smooth transition from classical weighting
    - ▶ Equivalent weights
    - ▶ When different methods give different results, we can track it back to an interaction

# Class 7a: Survey interviewing

- Questions and answers in surveys
- Evaluating survey questions
- Survey interviewing
- Surveys vs. other sources of information (administrative data, economic activity, . . . )
- Conceptual or specification errors
- Sampling and nonsampling errors
- Errors of measurement, interviewers, question wording, . . .
- Errors in reporting

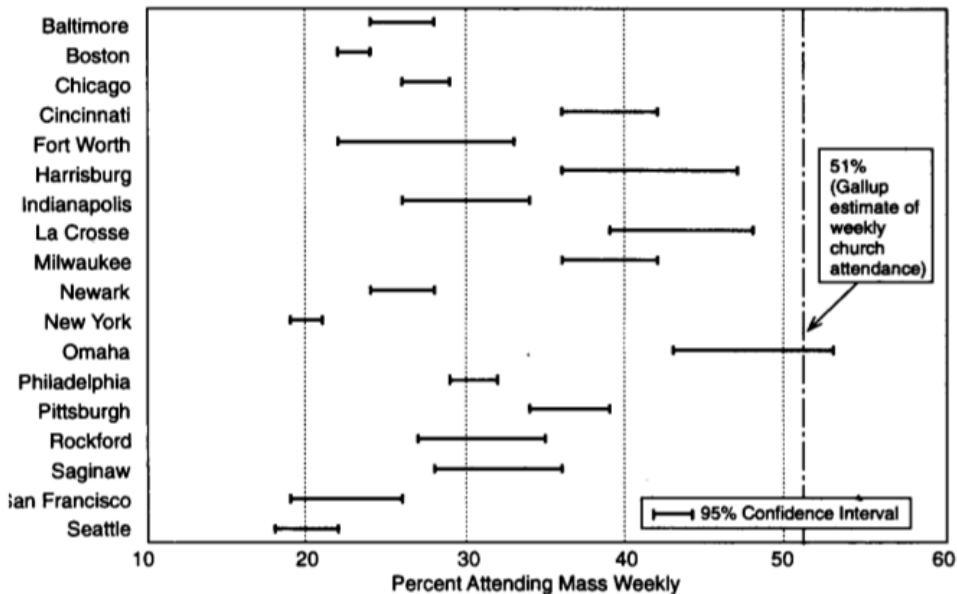- Difficulty of estimating small percentages
- Framing and question wording

# Estimating measurement effects

- You are conducting a survey and are concerned about the possible effects of the wording of one particular question. You decide to do one of two experiments:
    - (a) Within-subject design: Put the two different wordings on the same survey form (randomizing the order of the two questions) and compare responses to the two wordings.
    - (b) Between-subject design: Randomly give one wording to half the respondents and the other wording to the other half. Compare the average responses under the two wordings.
- Give a reason why you might prefer design (a).
- Give a reason why you might prefer design (b).
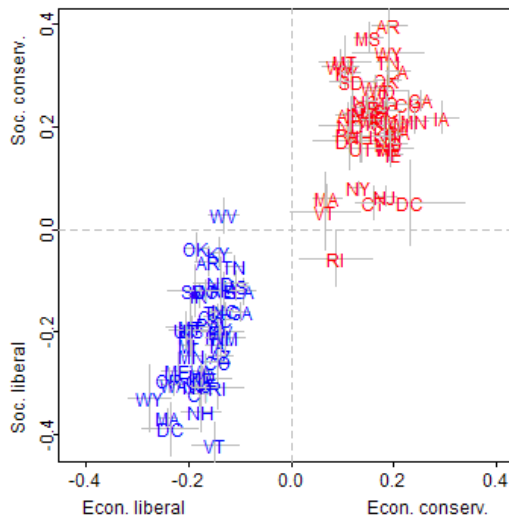- Computing standard errors for different designs

# Measurement: Church attendance

# Class 8a: Using surveys to answer questions in political science

- ▶ Political attitudes and behavior
- ▶ Comparing different groups

Average economic and social ideology scores
among Bush voters (red) and Gore voters (blue) in each state

# Class 8b: Conducting a survey in the real world

- Research goals
- Population, frame, and sampling design
- Constructing and testing the survey instrument
- Sampling and data collection
- Collection of auxiliary data
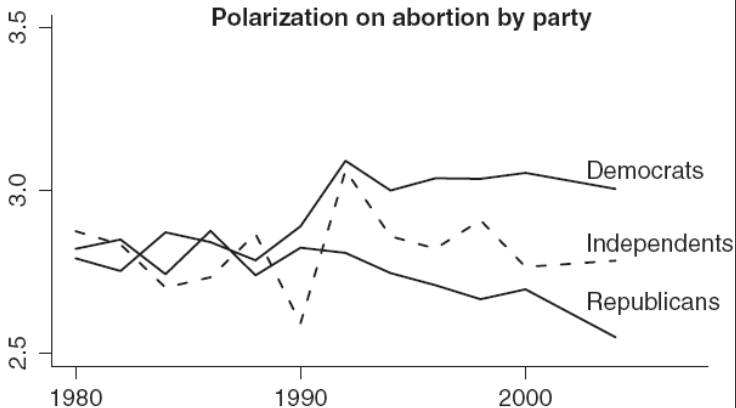- Data cleaning and manipulation
- Data analysis

# Class 9a: Voting

- Survey questions on demographics
- Political affiliation
- Issue attitudes

- ▶ Partisan polarization
- ▶ Elite and mass attitudes
- ▶ Uniform partisan swing

# Political polarization since 1990
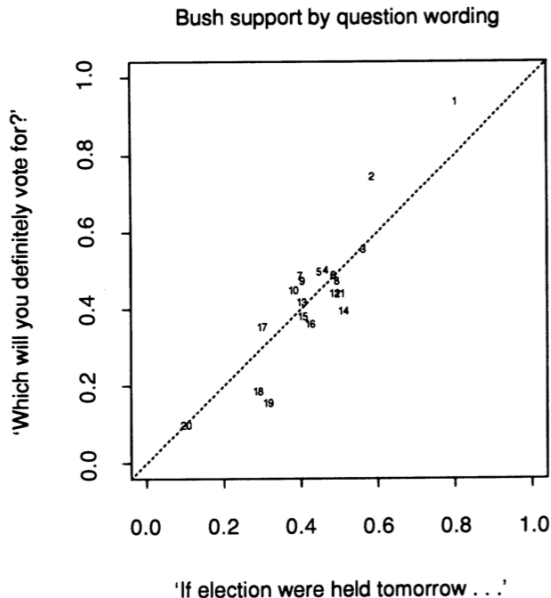


Polarization on abortion by party
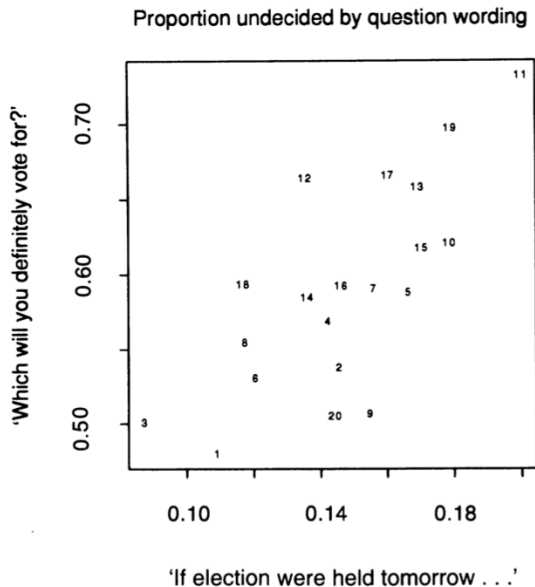
- ▶ Asking about voting
- ▶ Other forms of political participation
- ▶ Rationality of voting and responding to surveys
- ▶ Data sources

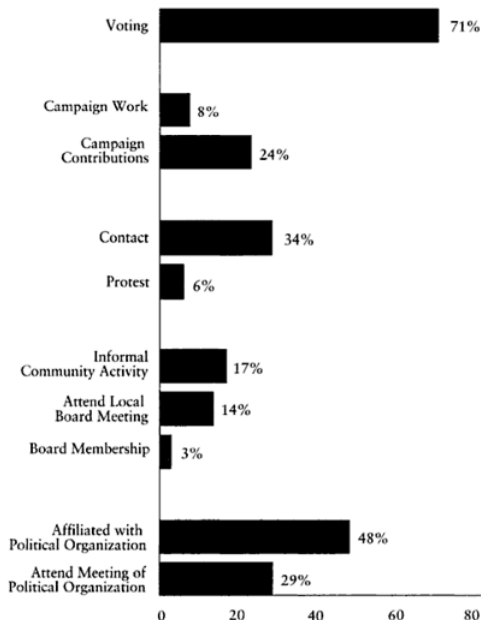# Question wording and vote intention



Bush support by question wording
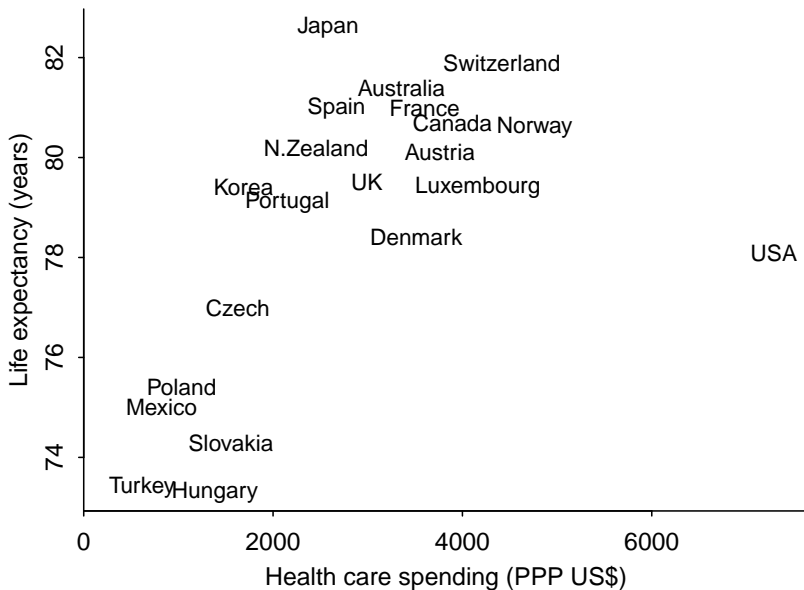
# Question wording and nonresponse



Proportion undecided by question wording

From Verba, Schlozman, Brady, *Voice and Equality* (1995)

- Manipulating data in R
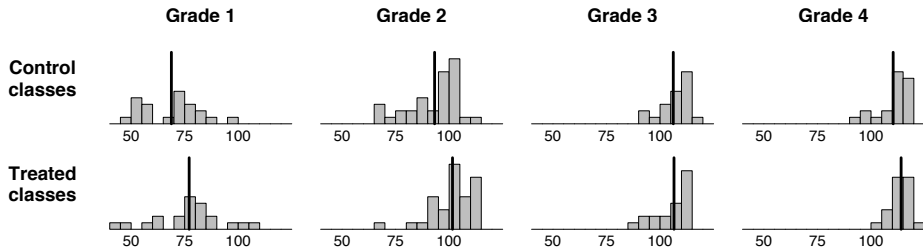- Looking carefully at the data
- Effective graphing

# Real simple
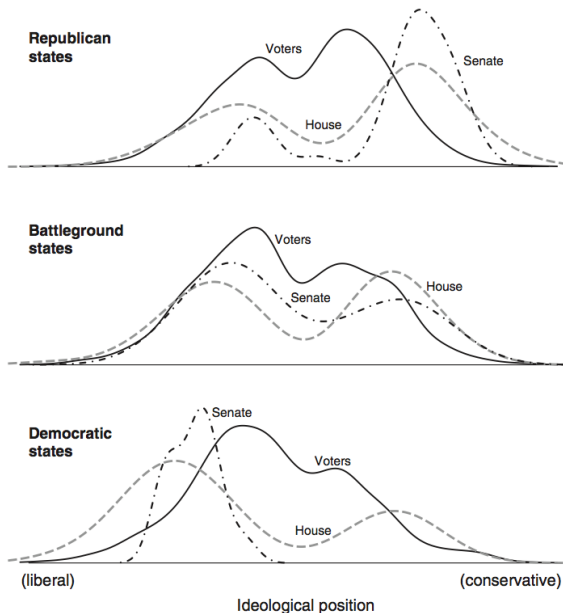
# Re-expression saves space and adds clarity

# Re-expression saves space and adds clarity

# Class 11a: Bayesian inference

- ▶ Calibration of probabilities
- ▶ Combining prior and data information
- ▶ Example: forecasting
- ▶ MRP

# Aligning voters with Congress

- Applications: ability testing, ranking
- Compare to ideal-point modeling
- Implications for education and for data collection more generally

- State-level opinions
- Comparisons to state policies
- Demographic breakdowns
- Displaying inferences

# Class 12b: Challenges in multilevel regression and poststratification

- Many factors, deep interactions
- Fitting and understanding models
- Adjusting for non-census variables
- Differential nonresponse within cells

- ▶ When does it matter?
- ▶ Validation of survey data
- ▶ Methods for increasing response rates

- Iraq mortality survey

- ▶ Fractal sampling
- ▶ Penumbra sampling
- ▶ Learning about networks
- ▶ Averaging over networks