# Sampling and Coverage

GOV 1010

# Pre-Election Polls

A Sample of Problems

# 1936 Presidential Election

Franklin Delano Roosevelt (D) is first-term incumbent

Elected in 1932 with 57.4% of vote

Alf Landon (R) , Kansas Governor, is challenger

# "People's Budget" vs Balanced Budget

# Pre Election Polls in 1936

# Literary Digest and Gallup

- Literary Digest
  - Weekly Newsmagazine
  - Founded in 1890
  - Polls since 1916
  - 136 million mailings 1916 – 1932
  - Predicted Roosevelt victory within 1% in 1932
  - "Uncanny accuracy"

- Methodology
  - Mail surveys
  - Large sample sizes
  - **Mailed 10 million surveys in 1936**
  - Auto owners / telephone directories

- George Gallup Ph.D.
  - Advertising Measurement specialist
  - "American Institute of Public Opinion"
  - New nationally syndicated newspaper column
  - Money-back guarantee (more accurate)

- Methodology
  - Face-to-face surveys with quota sampling
  - Also some mail polls

# GOOD YEAR POLL·O·METER

### For Recording
## The Literary Digest

**PRESIDENTIAL POLL·BROADCAST BY**
# THE GOODYEAR TIRE & RUBBER COMPANY, INC.
### OVER NBC BLUE NETWORK EACH MONDAY, WEDNESDAY & FRIDAY EVENINGS, SEPT. 2-NOV. 2

*R. Hillger*

## FIRST WEEK—SEPT. 2

| | Popular Vote | | | Electoral Vote | |
|------|------|------|------|------|------|
| Dem. | Rep. | Other | Dem. | Rep. | Other |
| 7,645 | 16,056 | 967 | | | |

## SECOND WEEK—SEPT. 9

| | Popular Vote | | | Electoral Vote | |
|------|------|------|------|------|------|
| Dem. | Rep. | Other | Dem. | Rep. | Other |
| 33,453 | 61,190 | 5,387 | 34 | 165 | 0 |

## THIRD WEEK—SEPT. 16

| | Popular Vote | | | Electoral Vote | |
|------|------|------|------|------|------|
| Dem. | Rep. | Other | Dem. | Rep. | Other |
| ...15 | 153,360 | 12,543 | 62 | 166 | 0 |

## FOURTH WEEK—SEPT. 23

| 1932 RESULTS | | Electoral Votes | STATES | FINAL 1936 POLL—L.D. | | | ACTUAL VOTES 1936 ELECTION | | |
|------|------|------|------|------|------|------|------|------|------|
| ROOSEVELT | HOOVER | | | Dem. | Rep. | Other | Dem. | Rep. | Other |
| 207,910 | 34,675 | 11 | Alabama | | | | | | |
| 79,264 | 36,104 | 3 | Arizona | | | | | | |
| 189,602 | 28,467 | 9 | Arkansas | | | | | | |
| 1,324,157 | 847,902 | 22 | California | | | | | | |
| 250,877 | 189,617 | 6 | Colorado | | | | | | |
| 281,632 | 288,420 | 8 | Connecticut | | | | | | |
| 54,319 | 57,073 | 3 | Delaware | | | | | | |
| 206,307 | 69,170 | 7 | Florida | | | | | | |
| 234,118 | 19,863 | 12 | Georgia | | | | | | |
| 109,479 | 71,312 | 4 | Idaho | | | | | | |
| 1,882,304 | 1,432,756 | 29 | Illinois | | | | | | |
| 862,054 | 677,184 | 14 | Indiana | | | | | | |
| 598,019 | 414,433 | 11 | Iowa | | | | | | |
| 424,204 | 349,498 | 9 | Kansas | | | | | | |
| 580,574 | 394,716 | 11 | Kentucky | | | | | | |
| 249,418 | 18,853 | 10 | Louisiana | | | | | | |
| 188,907 | 166,631 | 5 | Maine | | | | | | |

# The Literary Digest

NEW YORK — OCTOBER 31, 1936

## Topics of the day

# LANDON, 1,293,669; ROOSEVELT, 972,897

### Final Returns In The Digest's Poll of Ten Million Voters

Well, the great battle of the ballots in the Poll of ten million voters, scattered throughout the forty-eight States of the Union, is now finished, and in the table below we record the figures received up to the hour of going to press.

These figures are exactly as received from more than one in every five voters polled in our country—they are neither weighted, adjusted nor interpreted.

Never before in an experience covering more than a quarter of a century in taking polls have we received so many different varieties of criticism—praise from many; condemnation from many others—and yet it has been just of the same type that has come to us every time a Poll has been taken in all these years.

A telegram from a newspaper in California asks: "Is it true that Mr. Hearst has purchased The Literary Digest?" A telephone message only the day before these lines were written: "Has the Republican National Committee purchased The Literary Digest?" And all types and varieties, including: "Have the Jews purchased The Literary Digest?" "Is the Pope of Rome a stockholder of The Literary Digest?" And so it goes—all equally absurd and amusing. We could add more to this list, and yet all of these questions in recent days are but repetitions of what we have been experiencing all down the years from the very first Poll.

**Problem:** Now, are the figures in this Poll correct? In answer to this question we will simply refer to a telegram we sent to a young man in Massachusetts the other day in answer to his challenge to us to wager $100,000 on the accuracy of our Poll. We wired him as follows:

"For nearly a quarter century, we have been taking Polls of the voters in the forty-eight States, and especially in Presidential years, and we have always merely mailed the ballots, counted and recorded those returned, and let the people of the Nation draw their conclusions as to our accuracy. So far, we...

---

# Institute Forecasts the Re-election of Franklin D. Roosevelt, Gives Him 54% of Popular Vote, Minimum of 315 Electors

### Major Party Forecast Is 55.7; New York in F.D.R. "Sure" Column

### Election Forecast

# Literary Digest versus Gallup Forecasts

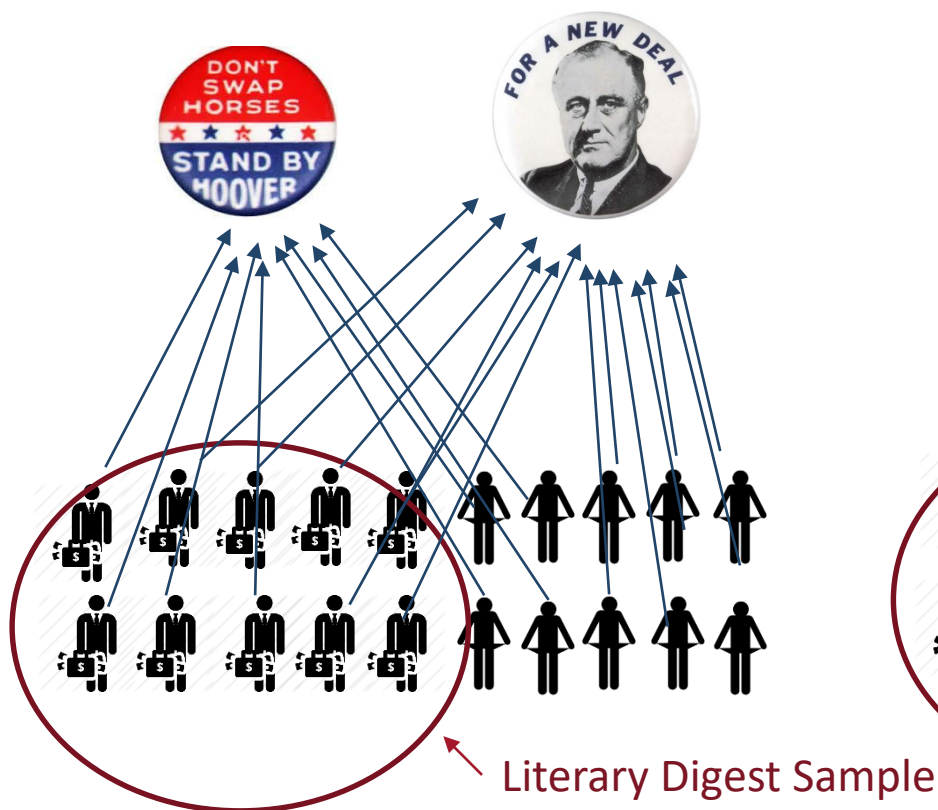|  | Literary Digest | Gallup | Election |
|---|---|---|---|
| Roosevelt | 43% | 56% | 62% |
| Landon | 57% | 43% | 37% |
|  |  |  |  |
| Roosevelt Electoral Votes | 161 | 315+ | 523 |
|  |  |  |  |
| Sample Size | 2,376,523 | 40,000 MAX |  |

# Discussion

# Differences Between 1932 and 1936



1932
Few Differences on Vote by Income

1936
Large Differences on Vote by Income

Literary Digest Sample

Literary Digest Sample

Inference in Scientific Samples

SAMPLING

# Populations in Surveys

## Population of Inference

- The general set of persons to whom one wishes to generalize results.
- This population may be infinite

## Target Population

- The inferential population as operationalized by the researcher.
- Bounded by time
- Observable (i.e. can be reached)

## Frame Population

- Population that could be measured by the sample frame
- Generally, all members of the frame
- Or... all members who could be enumerated by the sample frame

## Survey Population:

- The set of people who can be reached through your sample frame as implemented in the      survey

Inferential Population

*Social Science Theory*

Target Population

*Coverage Error*

Sample Frame

*Sample Error*

*"Margin of Error"*

Sample Records

*Nonresponse*

Respondents

All Americans

↓

*Social Science Theory*

All American Adults Home at Some Point Aug 19 – Aug 20, 2018

↓

*Coverage Error*

Americans with Domestic Telephone Service

↓

*Sample Error*

Approximately 5,000 (valid) Telephone Numbers Sampled

↓

*Nonresponse*

1,001 Americans who Responded to Survey

# Sample Frames

- **List, or**

- **Set of procedures**

- Sometimes requires two or more stages of selection

- Designed to cover target population

# Sampling Harvard Graduates

Sampling People at Burning Man Festival (Black Rock City)

# Ideal Sample Frame

- Simple list
- Available and accessible
- All Members of target population are on list
- All members of list are eligible respondents
- Contact information available for all elements of frame

# Example of Typical Sample Frame

Population | Sample
--- | ---
*Ed Ash* ⟶ | *617-555-1234*
*Peggy Birch* ⟶ | *617-555-1235*
*Tony Birch* | 
*Marie Chestnut* → | *617-555-1236*
 | *617-555-1237*
*Philip Elm* | *(No telephone)*
*(Not Assigned)* | *617-555-1238*

# Relationship Between Target Population and Sample Frame

| Relationship | Name | Problems |
|---|---|---|
| One-to-one | Perfect | None |
| One-to-None | Undercoverage | Bias |
| None-to-One | Low Incidence | Cost-Effectiveness |
| One to Many | Multiplicity | Probabilities |
| Many to One | Clustering | Probabilities |

# Example of Sample Frame

Population      Sample

*Ed Ash*     ⟶     *617-555-1234*

*Peggy Bolton* ⟶ *617-555-1235*      *Clustering*

*Tony Bolton*

*David Chandler* → *617-555-1236*      *Multiplicity*

             *617-555-1237*

*Philip Elm*

*(Not Assigned)*      *617-555-1238*      *Empty Record*

Target Population

Undercovergae

Coverage

Overcoverage (Empty Units)

Sample Frame

Out of Range

Target Population

Undercoverage

*Multiplicities*

Coverage

Undercoverage

Out of Range

Sample Frame B

Sample Frame A

## Measuring the Effect of Coverage Error on Population Estimates:

$$Y = \frac{N_c}{N} Y_c + \frac{N_{nc}}{N} Y_{nc}$$

$Y$=The value of the statistic in the target population

$N_c$=Number in the target population covered by the frame population

$N$=Total number in the target population

$Y_c$=Value of the statistic for those covered by the frame population

$Y_{nc}$=Value of the statistic for those not covered by the frame population

In Words:

| The Population Value of a Statistic | = | The proportion of the population included in your frame | X | The value of that statistic for those people | + | The proportion of the population *Not Included* in the frame | X | The value of the statistic for the people not included in the frame |
|---|---|---|---|---|---|---|---|---|

**Thus the impact of coverage error is based on two things:**

➤ The percent of the total population excluded from the sample frame

➤ The difference on the statistic of interest between those included in the frame and those excluded

# Differences Between 1932 and 1936



High Noncoverage
Low Differences Between Covered and Non-covered

High Noncoverage
High Differences Between Covered and Non-Covered

Uncanny Accuracy

Bankruptcy [1938]

Literary Digest Sample Frame

Literary Digest Sample Frame

# Considerations in Design

- What frames might be available for population?
  - Lists
  - Sets of procedures
- What is relationship of unit of frame to population?
- What is coverage of population in potential sample frames?
- What is incidence of respondents in potential sample frames?

| General Populations | Special Populations |
|---|---|
| • Broad populations of residents | • Narrow definition |
| • No list available | • Lists may be available |
| • No easy way to target | • Targeted frames may be feasible |

General versus Special Populations

## ➤ Levels of analysis

*Sometimes the conceptual population to which we infer our data doesn't match the survey population from which we collect information.*

**Contents:**
*One survey may gather information about different things*

**Example:** A survey of fast-food customers may yield information about different visits to fast food restaurants, people who eat fast food, households who eat fast food

**Units**
*One survey may gather information about a different conceptual unit than the person*

**Example:** A survey of households may interview a head of household or other household member

**Example:** A survey of land use may interview landowners to infer to acreage.

**Example:** A survey to estimate the number of job applicants who receive pre-employment drug tests might interview human resource officers at businesses.

**Time-Frames**
*A survey may gather information about more than one time period*

**Example:** A survey of investors might look at current and past ownership of investment instruments

# Examples of Special Populations

- Usage – customers, visitors, participants, etc.
- Occupations
    - Journalists
    - Firefighters
    - People wo work three jobs
    - Jazz Musicians
    - People who barter at flea markets
- Employees
- Companies

# Chocolate Activity

- Stick hand in bag and mix
- Select **ten** chocolates randomly
- Count number of Green Wrapper chocolates (out of ten)
- Put chocolates back in bag
- Enter Data in Qualtrics Survey

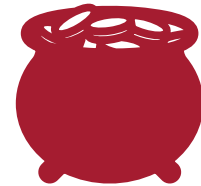# Scientific and Non-Scientific Samples

# Scientific Samples

- ***Based on Probability Theory***
- ***Allow Inference to Sample Frame***
- ***Sample Variance and Error Can Be Calculated***
  - Sample Records Are Drawn From a Well-Specified Frame
  - Sample Records Are Drawn According to Well-Specified Procedures With Known Properties
  - Each Sample Record Has a Known Non-Zero Probability of Selection
  - Data are Adjusted (Weighted) As Required To Reflect Sample Design
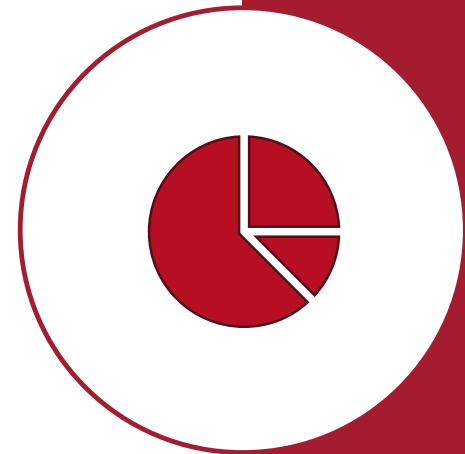
# Non-Probability Samples

- Availability Samples
  - Convenience Samples
  - Volunteer Cases

- Purposive Cases
  - Typical Cases
  - Critical Cases

- Respondent Driven Samples

- Quota Samples

# How Could We Sample Jazz Musicians?

# Respondent Driven Sampling (RDS)

- Useful when population definition may be rich or complex

- Useful when population may be rare

- Most useful if rare populations are part of networks


- Select set of seeds

- Give coupons to respondents to recruit other members of population

- Continue process multiple times


- Examples: Prostitutes, IV Drug Users, Illegal immigrants, Jazz Musicians,

# Scientific Samples

- *Based on Probability Theory*
- *Allow Inference to Sample Frame*
- *Sample Variance and Error Can Be Calculated*
  - Sample Records Are Drawn From a Well-Specified Frame
  - Sample Records Are Drawn According to Well-Specified Procedures With Known Properties
  - Each Sample Record Has a Known Non-Zero Probability of Selection
  - Data are Adjusted (Weighted) As Required To Reflect Sample Design

# Green Chocolate Example

- Population is Chocolates in Bag

- Statistic is whether chocolate has Green or Orange wrapper

- Chocolates are randomly mixed in bag

- Each chocolate has an equal probability of being selected

- In samples of ten, number of green chocolates is estimate of percentage of population that is green

# Green Chocolate Example

- Population: 142 total chocolates
  - 52 Green wrapped chocolates (36.62%)
  - 90 Orange wrapped chocolates (63.38%)

  - Probability of selection:

  - In samples of ten:
    - 10/142 ≈ .07

  - Each chocolate had an equal probability of being selected

# Two Key Statistical Elements Found in Any Population

## Central Tendency
### (Mean)

## Dispersion
### (Variance or Standard Deviation)

# Why Randomize?

- Statistical Theory is based on randomization
- If a sample is randomized, errors are randomly distributed
- In the long run, errors or biases cancel each other out
- If these biases cancel each other out in the long run, then in the long run, the sample mean equals the population mean, and the sample variance equals the population variance
- Randomization works across all biases and errors, including those that we don't think about or know about.

# Simple Random Samples (SRS)

- All population members have an equal chance of being selected

- Statistics are easy to calculate

- An **Equal Probability Selection Method (EPSEM)** sample

- Most statistics assume Sampling with Replacement

- In practice **Sampling Without Replacement (SWOR)**is most practical

**How to Pull a Simple Random Sample From a Complete Frame:**

Determine the size of your sample frame (N)

Determine the desired number of sample records you need (n)

Calculate your sampling fraction (k): $\left(\dfrac{n}{N}\right)$

Generate a Random Number for each frame element
*It's easiest to calculate a random number between 0 and 1, and carry it out to many decimal places*
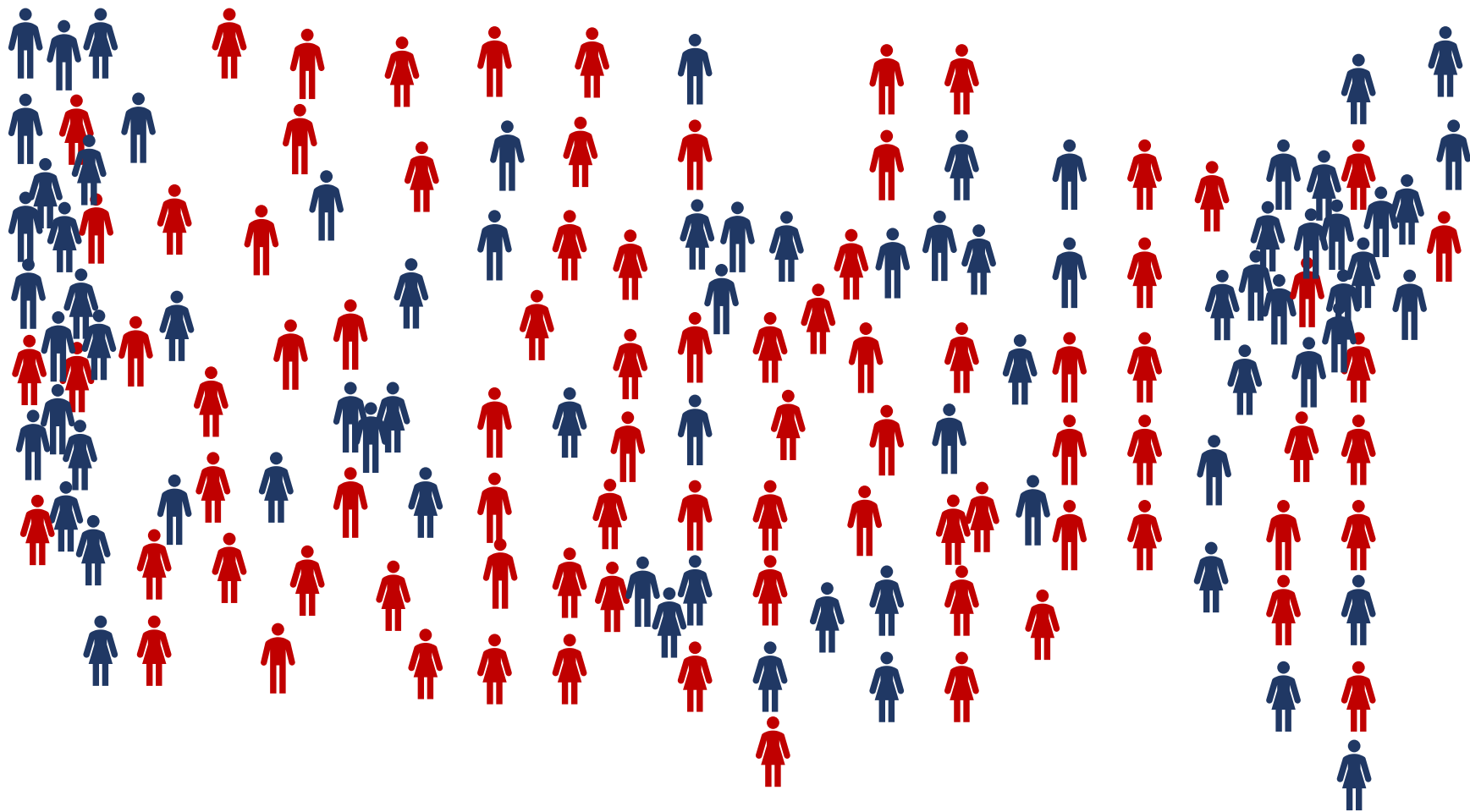
If the random number is Less Than or Equal to the Sampling Fraction, Include the element in your sample. Otherwise, exclude it.

**How to Pull a Systematic Random Sample From a Complete Frame:**

➢ Determine the total number of records in your frame (List) (N)

➢ Determine the desired number of sample records you need (n)

➢ Calculate your sampling fraction (*k*): $\left(\dfrac{n}{N}\right)$

➢ Generate a Random Number between 1 and k

➢ Count until you reach this random number and select this record

➢ Count from this randomly selected record and select every *k*th record until you reach the end of the list
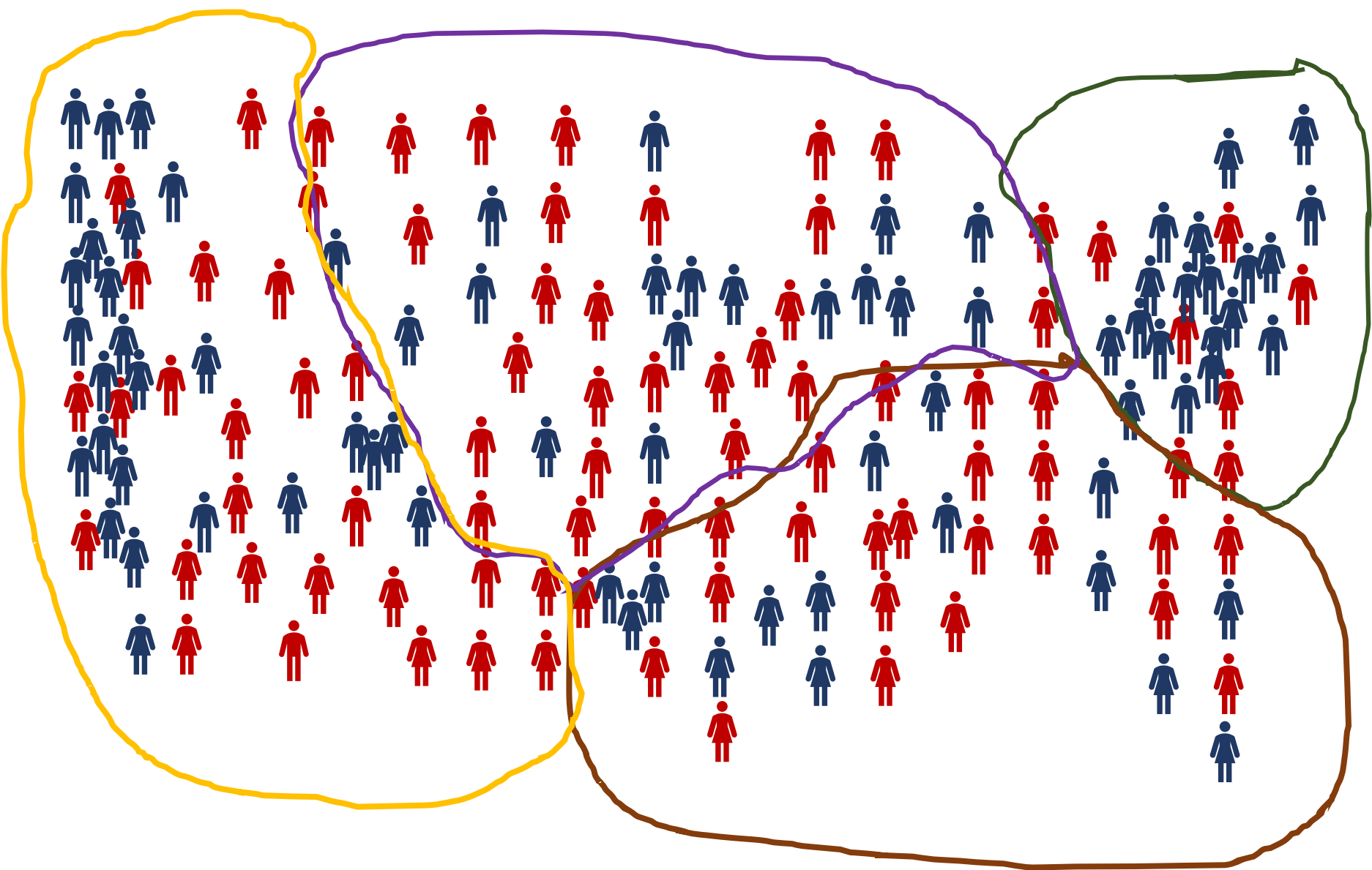
**Stratification:**

➢ Sample is divided into pre-designated units or partitions

➢ Strata should be homogeneous (heterogeneity between strata)

➢ Sample is drawn separately from each stratum

**Benefits of Stratification:**

- Insures that stratification elements are represented proportionally in overall survey data

    - Protects against error or bias *in the short run*

- Reduces overall sample variance

    - Increases precision of estimates

- Can allow more precise estimates of sub-groups if disproportionate stratification is used

## Example of Proportionate Stratification:

**Research Objective:** Conduct a survey to assess attitudes toward the electoral college

**Research Method:** Telephone survey of 1,000 telephone households in United States during November 2000

**Sample Frame:** List of all possible residential telephone numbers in US

**Effective Sampling Fraction:** You need to draw 6 telephone numbers to reach 1,000 households

Option 1: Draw a Simple Random Sample of 6,000 telephone numbers from your frame

Option 2:

➤ Divide the US into four (4) regions

➤ Determine the proportion of all households that are located in each region

➤ Draw a separate simple random sample within each of these four regions with an n that is proportionate to that regions size

## Regional Sample Stratification For US National Survey

| Stratum (Census Region) | Population | Percent | Sample Records | Estimated Interviews |
|---|---|---|---|---|
| Northeast | 39,418,789 | 19.97% | 1,198 | 200 |
| North Central | 46,076,032 | 23.35% | 1,401 | 233 |
| South | 69,894,352 | 35.42% | 2,125 | 354 |
| West | 41,954,834 | 21.26% | 1,276 | 213 |
| *Total* | *197,344,007* | *100%* | *6,000* | *1,000* |

**Sample Design For Connecticut Child Care Facilities**
**Proportional to Region and Type of Facility**
**(n=1,200)**

| Facility Type | Region | N | Percent | Sample if Proportional |
|---|---|---|---|---|
| Child Care Center | Southwest | 351 | 5.14% | 62 |
| | South Central | 308 | 4.51% | 54 |
| | Eastern | 290 | 4.25% | 51 |
| | North Central | 530 | 7.77% | 93 |
| | Northwest | 288 | 4.22% | 51 |
| Family Home | Southwest | 647 | 9.48% | 114 |
| | South Central | 951 | 13.94% | 167 |
| | Eastern | 1,019 | 14.93% | 179 |
| | North Central | 1,744 | 25.56% | 307 |
| | Northwest | 695 | 10.19% | 122 |
| *Total:* | | *6,823* | *100.00%* | *1200* |

## Sample Design And Weights For Connecticut Child Care Facilities
## Disproportionate Sample Design
## (n=1,200)

| Facility Type | Region | N | Population Percent | Expected | Observed | Sample Weight |
|---|---|---|---|---|---|---|
| Child Care Center | Southwest | 351 | 5.14% | 62 | 120 | 0.52 |
| | South Central | 308 | 4.51% | 54 | 120 | 0.45 |
| | Eastern | 290 | 4.25% | 51 | 120 | 0.43 |
| | North Central | 530 | 7.77% | 93 | 120 | 0.78 |
| | Northwest | 288 | 4.22% | 51 | 120 | 0.43 |
| Family Home | Southwest | 647 | 9.48% | 114 | 120 | 0.95 |
| | South Central | 951 | 13.94% | 167 | 120 | 1.39 |
| | Eastern | 1,019 | 14.93% | 179 | 120 | 1.49 |
| | North Central | 1,744 | 25.56% | 307 | 120 | 2.56 |
| | Northwest | 695 | 10.19% | 122 | 120 | 1.02 |
| *Total:* | | *6,823* | *100.00%* | *1,200* | *1,200* | |

# Considerations in Design

- What frames might be available for population?
  - Lists
  - Sets of procedures
- What is relationship of unit of frame to population?
- What is coverage of population in potential sample frames?
- What is incidence of respondents in potential sample frames?

# Survey Sampling -- Summary

Population

↓

Coverage Error

Sample Frame

↓

Nonresponse Error

Respondents

# Population Specification

## Conceptualize Inferential Population

- ➢ Begin by considering overall analytic goals of survey

- ➢ Determine extent and type of inference

- ➢ Fully elaborate and consider the nature of the inferential population, including different types of unusual cases and variants

## Operationalize Inferential Population in Target Population

- ➢ Develop specific definition

- ➢ Specify selection criteria and rules in detail

## Specify Frame Population in Relation to Target Population

- ➢ Specifically list details of frame population that meet and do not meet criteria of target population

- ➢ Be very specific about details, including sources and dates of databases or lists

- ➢ Specify procedures for dealing with potential incongruities between frame and target population

**Systematic Sampling:**

➢ Useful if full population list is only available electronically

    o Data entry (or scanning) of "hard-copy" list represent alternatives

➢ Sample records are systematically taken from a list so that every "kth" record is taken, and the list is sampled from beginning to end (*sometimes called "nth-ing"*

➢ Random start point should be used

➢ If list is ordered in a periodic manner serious bias can occur

➢ If list is ordered (or electronic list is sorted appropriately) systematic sampling represents a method of implicitly stratifying a sample

## How to Pull a Systematic Random Sample From a Complete Frame:

➢ Determine the total number of records in your frame (List) (N)

➢ Determine the desired number of sample records you need (n)

➢ Calculate your sampling fraction (*k*): $\left(\dfrac{n}{N}\right)$

➢ Generate a Random Number between 1 and k

➢ Count until you reach this random number and select this record

➢ Count from this randomly selected record and select every *k*th record until you reach the end of the list

**Implicit Stratification:**

- ➤ Utilizes Relevant Information About Frame to Order Frame

- ➤ Provides Many Benefits of Explicit Stratification

- ➤ Can Incorporate More Information Than Explicit Stratification

- ➤ Typically Used in Conjunction with Explicit Stratification

**How to Implicitly Stratify a Sample:**

- ➤ Sort Frame On Key Variables

- ➤ Take Systematic Sample or Systematic Random Sample From Sorted Frame

**Cluster Sampling:**

➢ Use when full population enumeration is not possible

➢ Use for cost efficiencies

➢ Use when physically required by research design

➢ Use when clustering of population is of analytic interest

**Overview of Clustered Sample:**

➢ Enumerate initial or Primary Sampling Units (PSU's)

➢ Select Sample (may be stratified) of PSU's

➢ If necessary, select further clusters below PSU stage

➢ If necessary, enumerate further clusters or elements in sufficient detail to calculate probabilities of selection

➢ Select n sample records from final or ultimate cluster

**How Many Clusters?  How many units?**

➤ The fewer clusters, the more economical the sample

➤ The more clusters, the more precise the overall sample

➤ The more units within each cluster, the more precise the estimates within that cluster

➤ In general population surveys, five units within each cluster is the norm

**Probability Proportionate to Size Sampling:**

➢ Method of multistage cluster sampling

➢ Results in an EPSEM sample

➢ Often called a "self-weighting sample"

➢ Typically results in better population coverage than cluster sampling with PSU's selected with equal probability

**Method:**

> ➢ Select PSU's with a probability proportionate to their overall size

> ➢ Select equal number of elements from each PSU

**Probability of Selection:**

$$\boxed{\text{Element Probability}} = \boxed{\text{Number of Clusters Selected}} \quad X \quad \boxed{\dfrac{\text{Cluster Size}}{\text{Population Size}}} \quad X \quad \boxed{\dfrac{\text{Elements Selected Per Cluster}}{\text{Cluster Size}}}$$

## PPS Example:
### EPSEM Exit Poll With PPS Design

|  | N | Prob. | n | f | Total Prob. |
|---|---|---|---|---|---|
| Precinct 1 | 500 | .13 | 100 | 1/5 | .03 |
| Precinct 2 | 1,000 | .25 | 100 | 1/10 | .03 |
| Precinct 3 | 200 | .05 | 100 | 1/2 | .03 |
| Precinct 4 | 800 | .20 | 100 | 1/8 | .03 |
| Precinct 5 | 500 | .13 | 100 | 1/5 | .03 |
| Precinct 6 | 1,000 | .25 | 100 | 1/10 | .03 |
| Total: | 4,000 |  |  |  |  |

**Example:**
*EPSEM Exit Poll With PSU's Selected With Equal Probability*

|  | N | Prob. | f | n | Total Prob. |
|---|---|---|---|---|---|
| Precinct 1 | 500 | .17 | 1/7 | 75 | .03 |
| Precinct 2 | 1,000 | .17 | 1/7 | 150 | .03 |
| Precinct 3 | 200 | .17 | 1/7 | 30 | .03 |
| Precinct 4 | 800 | .17 | 1/7 | 120 | .03 |
| Precinct 5 | 500 | .17 | 1/7 | 75 | .03 |
| Precinct 6 | 1,000 | .17 | 1/7 | 150 | .03 |
| Total: | 4,000 | | | | |

Sampling People who Go to Baseball Games at Fenway Park

# A Taxonomy of Probability Selection Methods

| | | |
|---|---|---|
| I. | **Equal Probabilities of Selection (EPSEM):** (a) Equal Probabilities at all stages of sample design (b) Equal overall probabilities obtained through compensating unequal probabilities at several stages | **Unequal Probabilities for different stages;** ordinarily compensated with inverse weights (a) Caused by irregularities in selection frames and procedures (b) disproportionate allocation designed for optimum allocation |
| II. | **Element Sampling:** Single stage, sampling unit contains only one element | **Cluster Sampling:** Sampling units are clusters of elements (a) One-stage cluster sampling (b) Subsampling or multistage sampling (c) Equal clusters (d) unequal clusters |
| III. | **Unstratified Selection:** Sampling units selected from entire population | **Stratified Sampling:** separated selections from partitions, or strata, of population |
| IV. | **Random Selection** of individual sampling units from entire stratum or population | **Systematic Selection** or sampling units with selection interval applied to list |
| V. | **One-Phase sampling:** Final sample selected directly from entire population | **Two-Phase (or double) sampling:** final sample selected from first-phase sample, which obtains information for stratification or estimation |

**Source:** Leslie Kish; *Survey Sampling*, New York: John Wiley & Sons, 1965