# Doing More With Less:

# From Sample to Population

Jenny Zhang

Allstate Insurance Company

May, 2015

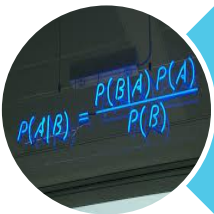# Table of Contents



Calculating Sample Size
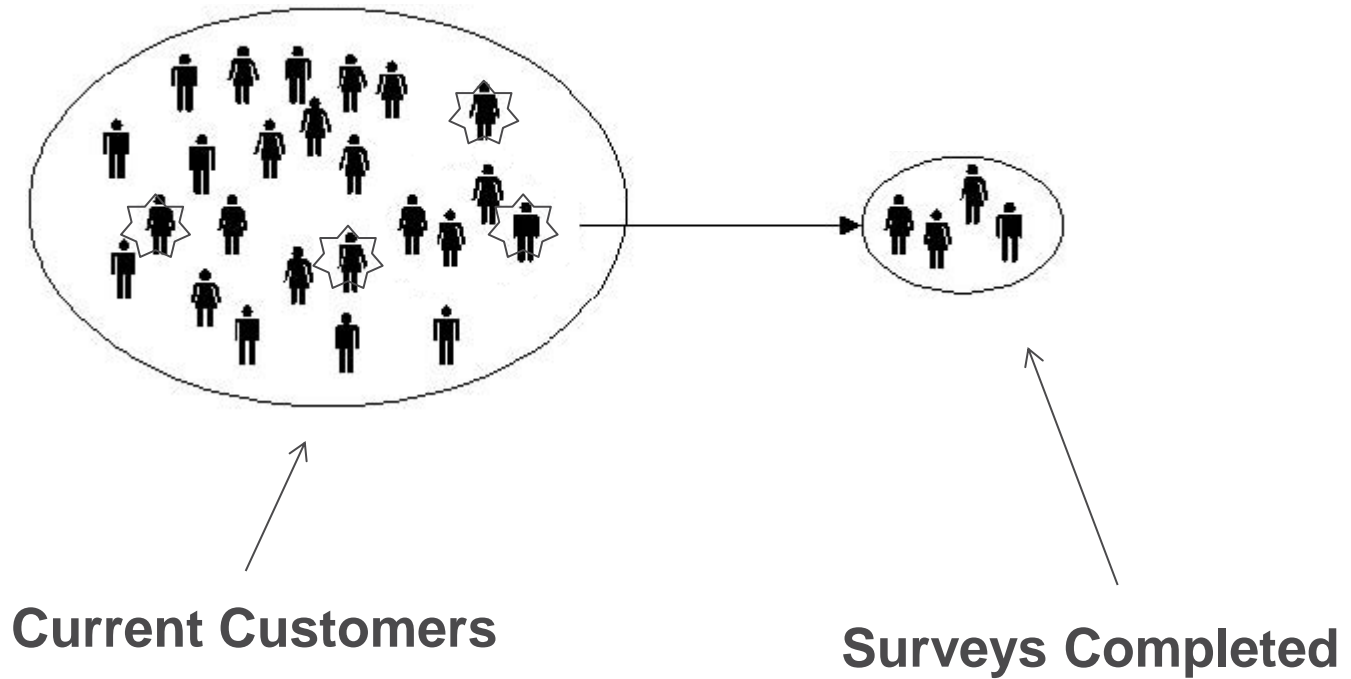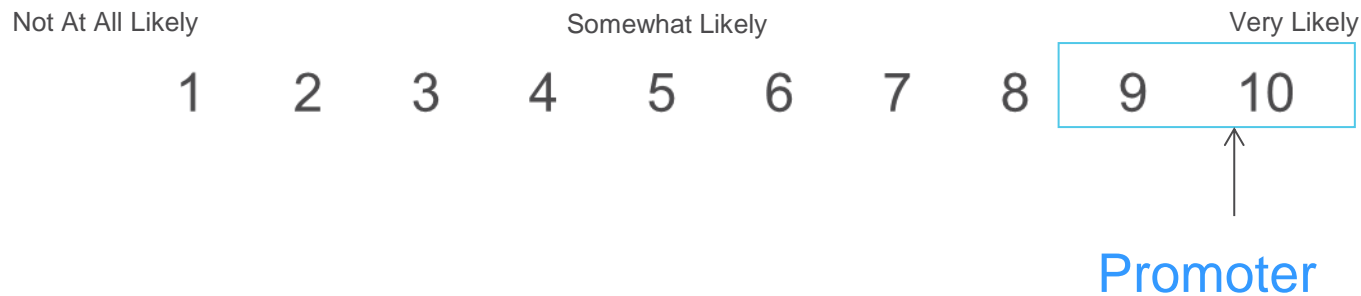


Bootstrapping



Downsampling



Bayesian Analysis

# Population and Sample



**Current Customers**

**Surveys Completed**

# Measure Proportion of Promoter

What is the likelihood that you would recommend insurance from XXX
company to a friend or colleague?

Not At All Likely                    Somewhat Likely                                    Very Likely

1    2    3    4    5    6    7    8    9    10

Promoter

**We want to know the proportion of Promoters $p$.**

# Hypothesis Testing

**Hypothesis testing** is a formal process to determine whether to reject a null hypothesis or not, based on the sample data.

- $H_0$ is null hypothesis
- $H_a$ is alternative hypothesis

In the statistical framework, we either reject a null hypothesis or fail to reject a null hypothesis.

# Type I and Type II Error

## Decision

|  | Reject $H_0$ | Fail to reject $H_0$ |
|---|---|---|
| **Truth** $H_0$ | **Type I Error** | Correct Decision |
| $H_a$ | Correct Decision | **Type II Error** |

# What is $z_\alpha$

A z-score indicates how many standard deviations an element is from the mean.

$f$

5%

$z = 0$

$z_\alpha = 1.645$

# Confidence Interval

# Traditional Sample Size

The point estimate of $p$ is $\hat{p} = \frac{y}{n}$ and an approximate $1 - \alpha$ confidence interval for $p$ is

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Here $y \sim binomial(n, p)$.

$\alpha$ is the Type I error.

$\varepsilon = z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ is the maximum error of the point estimate $\hat{p} = \frac{y}{n}$.

# Traditional Sample Size

$$\varepsilon = z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Assume $p$ is about to equal to $p^*$, then

$$Sample\ Size\ (SS) = n = \frac{z^2{}_{\alpha/2}\, p^*(1 - p^*)}{\varepsilon^2}$$

$$SS = \frac{1.645^2 * 0.82(1 - 0.82)}{0.10^2} \approx 39.94$$

From 40 sample surveys, we are 90% confident that the true $p$ belongs to the interval $\hat{p} \pm 0.10$.

| Confidence Level | Z score |
| --- | --- |
| 90% | 1.645 |

Test for Equality

How many surveys per agent are credible in order to reflect a statistically significant change on the Promoter Score?

# Steps

1. Set up hypotheses: null and alternative

2. Identify the distributions and key variables

3. Select desirable confidence level and statistical power

4. Calculate sample size needed

# Hypothesis on Promoter Score

$$H_0 : p_1 - p_2 = 0 \; vs \; H_a : p_1 - p_2 \neq 0$$

p1 is the current proportion of promoters among all customers for an agent.

p2 is the proportion of promoters for an agent in the next month.

$$\hat{p} = \frac{y}{n} \; and \; y \sim binomal(n, p)$$

# Derivation

We reject the null hypothesis if

$$\left| \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{(\hat{p}_1(1 - \hat{p}_1) + \hat{p}_2(1 - \hat{p}_2))/n}} \right| > z_{\alpha/2}$$

$$\frac{|p_1 - p_2|}{\sqrt{(\hat{p}_1(1 - \hat{p}_1) + \hat{p}_2(1 - \hat{p}_2))/n}} - z_{\alpha/2} = z_{\beta}$$

$$n_1 = n_2 = \frac{(z_{\alpha/2} + z_{\beta})^2 [p_1(1 - p_1) + p_2(1 - p_2)]}{(p_1 - p_2)^2}$$

# Sample Size Needed

$$n_1 = n_2 = \frac{(z_{\alpha/2} + z_\beta)^2 [p_1(1-p_1) + p_2(1-p_2)]}{(p_1 - p_2)^2}$$

$z_{\alpha/2}$ is z statistic given Type I error $\alpha$

$z_\beta$ is z statistic given Type II error $\beta$

# Illustrative Example

$$H_0 : p_1 - p_2 = 0 \: vs \: H_a : p_1 - p_2 \neq 0$$

$$n_1 = n_2 = \frac{(z_{\alpha/2} + z_{\beta})^2 [p_1(1-p_1) + p_2(1-p_2)]}{(p_1 - p_2)^2}$$

$$n_1 = n_2 = \frac{(1.645 + 0.84)^2 [0.80(1-0.80) + 0.86(1-0.86)]}{(0.80 - 0.86)^2} \approx 480.98$$

**To conclude a significant change, more samples are needed.**

# Aggregate the Results

We want to aggregate the results to a regional level.

We can use weighted average by number of customers to scale it back in order to minimize the sample size needed.

| | Promoter Proportion | Agency Satisfaction |
|---|---|---|
| Big Sample | 59.3% | 70.2% |
| Weighted Avg w/ 40 surveys per agent | 59.0% | 70.5% |

# Sample Size Calculation

Survey Analysis

Hypothesis Testing

Collect right amount of sample

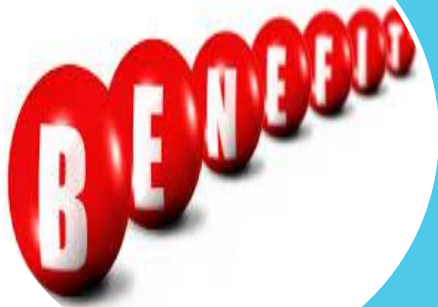# Bootstrapping



Random sampling with replacement.

# Illustrative Example Cont.

## - Customer Experience Survey

# Bootstrap

When Sample is not Sufficient

Establish Error Bounds

Statistical Inference with a Small Sample

# Downsampling

Down Sampling zero-claims can reduce the time of model convergence during the model development phase.

The standard deviation for parameter estimator will likely increase with the small sample. It won't necessarily reverse modeling decisions.

Modeler can periodically run the model with the entire training data set to verify the model structure.

Once the mode structure is finalized, the entire data set should be used for the final model.
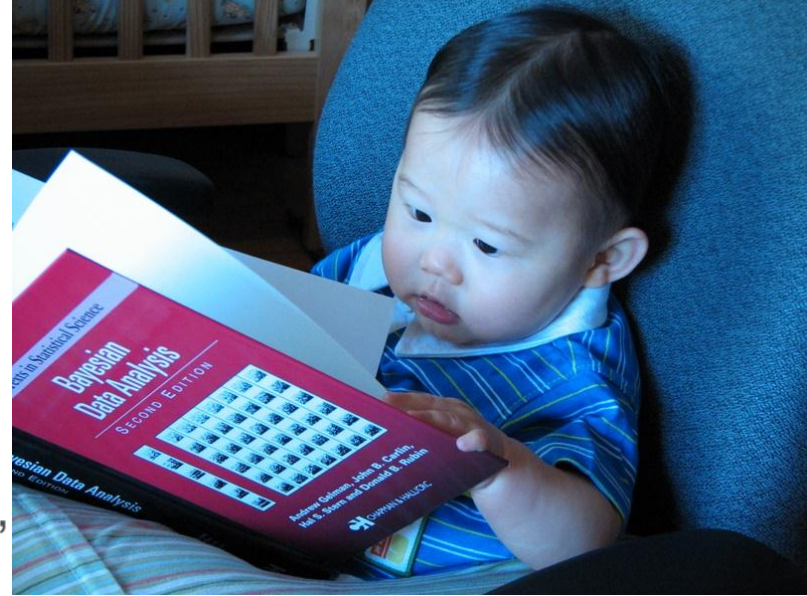
# Downsampling

During Model Development

Save Time

# Bayesian Analysis



$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Posterior is proportional to "likelihood*prior"

We use known knowledge as input for the prior distribution, and use data collected to calculate the likelihood. The posterior will be updated as more and more data is collected.

As a result, the sample size needed will be reduced.

# Bayesian Analysis

Need more flexibility

Less Sample Needed

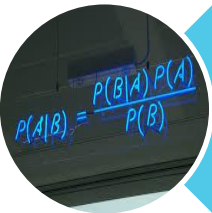Allow distributions for parameters

# Summary

Calculate Sample Size

Bootstrapping

Downsampling

Bayesian Analysis

# Questions